

Lecture da
LE SCIENZE

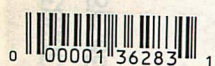
edizione italiana di
**SCIENTIFIC
AMERICAN**

VERITA' E DIMOSTRAZIONE

Questioni di matematica

presentazione di Carlo Ciliberto

100	99	98	97	96	95	94	93	92	91
65	64	63	62	61	60	59	58	57	90
66	37	36	35	34	33	32	31	56	89
67	38	17	16	15	14	13	30	55	88
68	39	18	5	4	3	12	29	54	87
69	40	19	6	1	2	11	28	53	86
70	41	20	7	8	9	10	27	52	85
71	42	21	22	23	24	25	26	51	84
72	43	44	45	46	47	48	49	50	83
73	74	75	76	77	78	79	80	81	82



K 136283
D 136258

DEPO 00566

SCIENZE

Sezione n. 3

V-DEPO 566

VERITÀ E DIMOSTRAZIONE

Questioni di matematica

136283
136258

Redazione del volume: Maurizio Negri

La copertina

L'illustrazione della copertina, tratta dalla copertina del numero di marzo del 1964 di «Scientific American» disegnata da Joan Starwood, rappresenta un curioso comportamento dei numeri primi, cioè dei numeri interi che sono divisibili solo per se stessi e l'unità, messo in luce per la prima volta da Stanislaw M. Ulam del Los Alamos Scientific Laboratory. Ulam ha scoperto che se si scrivono i numeri naturali su carta quadrettata disponendoli secondo una spirale, i numeri primi tendono a disporsi lungo linee diagonali. Sulla copertina la spirale è indicata con una linea nera marcata, i numeri primi sono in rosso e le linee diagonali sono in verde.

ISBN 88-7004-032-1

Copyright © 1964, 1965, by Scientific American Inc.; 1968, 1969, 1970, 1971, 1972, 1973, 1974, 1976, 1977, by Le Scienze S.p.A.

Lecture da
LE SCIENZE

VERITÀ E DIMOSTRAZIONE
Questioni di matematica

presentazione di
Carlo Ciliberto



LE SCIENZE S.p.A. editore
Milano, 1978

A 510 ER

mai 510.1 VER

È noto che la matematica ha svolto fin dalla antichità un ruolo attivo e di produttivo sviluppo sul piano culturale nella vita sociale dell'uomo, basti per questo ricordare l'importanza avuta dalla scuola matematica greca con Pitagora, Euclide, Archimede il quale peraltro occupa un posto particolare nella matematica per geniali scoperte da cui conseguirono importanti applicazioni. La matematica, che spesso, a torto, viene considerata come una scienza arida, noiosa, astratta, limitata e di dominio di una ristretta cerchia di cultori, è il permanente supporto di ogni scienza e di ogni attività pratica, concreta dell'uomo. Nel suo significato culturale, come pure nel suo uso ai fini applicativi, la matematica nelle sue varie estrinsecazioni e articolazioni ha legami con tutte le altre scienze, naturali e non, e queste sono legate alla matematica, che è il loro linguaggio e il loro strumento, spesso essenziale. Non è pertanto pensabile una separazione della matematica dal sostanziale sviluppo della vita sociale dell'uomo. Invero grandi conquiste, conseguite dall'uomo, in svariati importanti campi, sono basate proprio su notevoli risultati acquisiti facendo uso di raffinati strumenti matematici. Il ruolo della matematica si è poi ampliato nell'ultimo secolo in relazione al notevole sviluppo industriale che si è accompagnato a conquiste sociali che hanno comportato l'uso sempre più spinto e avanzato di applicazioni della matematica in svariati campi, dalla ingegneria e dalla fisica alla biomedicina, dalla agricoltura allo studio dei fenomeni geodinamici e idrogeologici, dalla statistica alle scienze umane, alla sociologia, ecc., tanto che addirittura nell'ultimo ventennio ha preso corpo una sua branca specifica, l'informatica, che costituisce una disciplina attraverso la quale lo strumento matematico viene posto al servizio dell'informazione nella sua più ampia accezione. In particolare è ben noto lo sviluppo acquisito negli ultimi anni dalla scienza dei calcolatori, che fonda le sue radici su strumenti matematici. D'altra parte è anche da sottolineare l'importanza del ruolo culturale della matematica: non è accettabile infatti che questa scienza sia confinata nel ghetto delle discipline prettamente strumentali, assegnandole il livello di una tecnica, spesso anche molto raffinata, che non si può non diffondere perché è molto importante per le applicazioni, ma che non ha nulla da dire sulla formazione dell'uomo. La matematica invece ha un suo posto insostituibile in questa formazione perché educa, attraverso il suo stretto rapporto con la logica, all'analisi critica dei concetti, all'astrazione, alla deduzione rigorosa e anche all'umiltà intellettuale. Sono i due profili culturali della matematica che vanno però riguardati globalmente, nel senso che attraverso una cultura matematica impostata sul rigore e l'astrazione è possibile pervenire a dare corrette interpretazioni e spiegazioni dei fenomeni naturali e creare così quegli strumenti applicativi occorrenti nella realtà quotidiana. A tale proposito va ricordato che spesso le più grandi scoperte hanno trovato conferma, spiegazione e sistemazione proprio attraverso lo studio matematico dei fenomeni a esse legati, e ne sono così derivate teorie che a loro volta hanno portato a nuove scoperte; basti per questo andare col pensiero all'astronomia e all'astronautica i cui risultati affondano le radici in notevoli teorie matematiche. D'altra parte spesso teorie matematiche create in maniera puramente astratta, a seguito di verifiche e riscontri, sono risultate essere la base concreta di rilevanti risultati pratici. È questo il caso dei profondi studi matematici in fluidodinamica che hanno portato poi a realizzazioni di alta tecnica, nonché alla soluzione del problema del superamento del muro del suono, determinando così una autentica rivoluzione nel campo aerospaziale.

La rivista «Le Scienze» nella sua ampia concezione culturale non poteva non cogliere l'importanza della diffusione delle conoscenze nel campo della matematica e pertanto sapientemente fin dalla sua nascita ha dato spazio ad argomenti matematici di alto livello culturale. Su un vasto piano internazionale, sono stati enucleati i temi di più vasto interesse corrente, presentando dei notevoli squarci di cultura matematica, di grande validità. Proprio per le motivazioni esposte, gli articoli matematici di «Le Scienze» non potevano certamente rimanere esclusi dall'iniziativa di essere raccolti almeno in parte in un volume della serie delle «Lecture», e in tale direzione si è giustamente mosso l'editore.

La scelta è caduta su una serie di saggi che, nell'arco degli ultimi anni, riguardano

PRESENTAZIONE

gli aspetti principali dello sviluppo del pensiero matematico. Si tratta di una selezione che tocca le maggiori aree culturali e le fasi più salienti, mettendo in evidenza come le ricerche fondamentali e quelle più direttamente applicative siano strettamente intrecciate tanto da rendere spesso oltremodo incerto e opinabile il confine tra le une e le altre. A tale selezione è stato aggiunto un paio di interessanti articoli pubblicati sull'edizione originale americana. Gli articoli raccolti nel presente volume sono stati ripartiti in tre parti, riguardanti rispettivamente i fondamenti della matematica e la teoria degli insiemi, la matematica e la logica, la matematica nelle applicazioni.

La prima parte si apre molto opportunamente con l'esposizione, tratta dalla rubrica «Giochi matematici», di M. Gardner, di una nuova e divertente variante dei noti paradossi di Zenone, dovuta ad A.K. Austin. Segue un articolo di W.V. Quine sui fondamenti della matematica. In esso, partendo da classici esempi, fonti di difficoltà concettuali, viene descritto il processo di riduzione di alcune nozioni a altre, e cioè il processo di riduzione dell'insieme dei concetti matematici con una risistemazione generale che elimina le difficoltà originarie. Il discorso viene poi ampliato al caso delle leggi e alla teoria della dimostrazione, originata dagli studi di Gödel. Segue poi un gruppo di tre articoli dedicati a questioni di teoria degli insiemi trattate con teorie non standard. P.J. Cohen e R. Hersh, servendosi di un'analogia con la geometria non euclidea, illustrano la non dimostrabilità (stabilita dallo stesso Cohen) dell'ipotesi del continuo formulata da Cantor alla fine del secolo scorso. L'articolo di M. Davis e R. Hersh è dedicato alla teoria dell'analisi non standard, un nuovo ramo della matematica di recente scoperto dal logico A. Robinson, che consente, utilizzando il linguaggio formale (peraltro legame fra logica e teoria dei calcolatori), di riportare in vita metodi infinitesimali nel calcolo differenziale e integrale, adoperati fin dall'antichità, ma spesso con dubbi. Nell'ordine di idee dell'adozione di teorie e tecniche non standard si inquadra l'articolo di G. Lolli sulla possibilità di costruzione di un nuovo modello del sistema dei numeri reali, particolarmente interessante per i suoi legami con la nozione di probabilità, e che è profondamente diverso da quello intuitivo costituito dalla retta euclidea. Un articolo di B. di Finetti è dedicato a una suggestiva e singolare illustrazione degli stretti legami che uniscono i tre numeri e , i , π , fondamentali in matematica, con una descrizione che non richiede alcuna conoscenza propedeutica, ma che proprio da se stessa fa scaturire le nozioni occorrenti, legandole anche a interessanti interpretazioni in fisica e in economia.

I rapporti fra matematica e logica sono così intrinseci da costituirne una sola disciplina, talché si può ormai considerare la logica come fondamento essenziale della matematica. In proposito va sottolineato che l'astrattezza e la sistemazione rigorosa di taluni rami della matematica è resa possibile dalla presenza di tecniche logiche, che, proprio per il loro astratto rigore formale, conferiscono naturalezza e sicurezza al procedere di ragionamenti per la loro natura poco provvisti di supporto intuitivo, sicché si può dire che la logica costituisce l'ossatura della matematica. La seconda parte degli articoli, che è dedicata a tali rapporti, si apre con una suggestiva e altamente divulgativa esposizione di M. Gardner di questioni di logica, riguardanti l'algebra di Boole, i diagrammi di Venn e il calcolo proposizionale. Segue un articolo di A. Tarski dedicato all'esame dei concetti di verità e di dimostrazione, diversi ma tra loro connessi, e allo studio del problema della coincidenza o meno dell'insieme delle proposizioni formalmente dimostrabili con quello delle proposizioni vere. (Per la sua suggestività, il titolo dell'articolo è stato assunto come titolo del presente volume.) Il periodo di transizione dalla logica tradizionale a quella contemporanea, propugnata da B. Russell, costituì una tappa importante nella logica e W.W. Bartley III in un suo articolo mette in luce il fatto che esso ha avuto uno dei più interessanti innovatori tecnici in C.L. Dodgson (noto come L. Carroll). L'evoluzione della matematica è legata strettamente a quella dei processi nell'astrazione; in questo quadro acquista enorme rilevanza una recente conquista del pensiero matematico, la «teoria delle categorie» che, come è illustrato nell'articolo di L. Lombardo Radice, costituisce un nuovo livello di astrazione, collocandosi a un più elevato grado di generalità rispetto alle astrazioni di secondo grado che hanno caratterizzato la ma-

tematica negli ultimi decenni del XIX secolo e nei primi del XX. Nel quadro del rapporto fra logica e matematica si pongono i nessi fra la teoria della probabilità e la moderna logica induttiva, e nell'articolo di D. Costantini e M. Mondadori viene posto in luce come l'evoluzione del concetto di probabilità e la «crisi dei fondamenti» hanno legato strettamente tali teorie. Uno dei problemi che hanno occupato per secoli i ricercatori è stato quello della conferma di un'ipotesi scientifica mediante l'osservazione e l'esperimento: W.C. Salmon nel suo articolo pone in rilievo come la logica di questo procedimento è ancor lungi dall'essere compresa, poiché allo stadio attuale di sviluppo gli studi in proposito «hanno prodotto più paradossi ed enigmi che non soluzioni convincenti o largamente accettate di problemi fondamentali». Nell'articolo di H. DeLong viene posto in evidenza il fatto che l'aver dimostrato che ci sono problemi che è impossibile risolvere è una delle grandi conquiste della matematica: si tratta dei cosiddetti teoremi limitativi, che sembrano esprimere delle limitazioni delle facoltà umane, fra cui quello di incompletezza di K. Gödel che viene discusso nell'articolo.

Non è inutile sottolineare ancora la rilevanza delle teorie matematiche e degli strumenti algoritmici matematici ai fini delle applicazioni. È la matematica che si cala nella realtà come scienza concreta. A.M. Turing, vissuto nella prima metà di questo secolo, fu uno dei primi studiosi di calcolatori, e fra l'altro ideò una macchina, nota con il suo nome, che, come dimostrò, può essere programmata per realizzare qualsiasi operazione alla portata del più potente elaboratore elettronico. M. Gardner, nell'articolo che apre questa parte, dopo una descrizione di tale congegno idealizzato, tratta la questione sollevata dallo stesso Turing: può una macchina pensare?, che conduce a profonde e dialettiche controversie filosofiche. Un altro legame fra teoria e applicazione è mostrato nell'articolo di H. Wang nel quale si pone in luce che la questione di stabilire con ragionamenti logici la vincita o meno in un gioco è analoga al fatto che un problema si possa o meno risolvere mediante calcolatore e si evidenzia così lo stretto legame intercorrente tra giochi, logica e calcolatori. Un grosso problema con riflessi in campo applicativo, posto da Hilbert e per lungo tempo non risolto, riguarda la possibilità di descrivere una procedura meccanica grazie alla quale sia possibile saggiare ogni equazione diofantea per decidere se possiede soluzione. Y. Matyasevich ha mostrato che a tale quesito si risponde negativamente, cioè non esiste un tale algoritmo. Si tratta di un importante contributo alla conoscenza delle proprietà dei numeri illustrato da M. Davis e R. Hersh nel loro articolo. È noto che il punto focale della scienza dei calcolatori è lo studio degli algoritmi che sono processi, procedure, metodi graduati atti a fornire uno specifico *output* come risposta a uno specifico *input*. Nell'articolo di D.E. Knuth viene data una lucida esposizione di modelli algoritmici, sottolineando il fatto «che gli algoritmi si occupano innanzitutto di manipolazioni di simboli, che non rappresentano necessariamente numeri»: questo concetto è alla base del linguaggio dei calcolatori. La serie degli articoli si chiude con un articolo di E.C. Zeeman sulla recente teoria delle catastrofi. Con metodi topologici di recente elaborazione, è possibile descrivere i fenomeni discontinui e divergenti. Proprio per la sua vasta generalità e per il fatto che può essere applicato efficacemente in situazioni in cui forze e motivazioni che mutano gradualmente portano a cambiamenti comportamentali improvvisi, il metodo è stato chiamato «teoria delle catastrofi». Molti esempi di catastrofi sono offerti da fenomeni fisici, ma applicazioni rilevanti della teoria possono essere effettuate in biologia e in scienze sociali, dove sono molto diffusi fenomeni discontinui e divergenti, e in cui altre tecniche matematiche si sono rivelate insufficienti.

La panoramica del presente volume, seppure ampia, ovviamente non può ritenersi esaustiva, ma vuole soltanto offrire uno squarcio dei più rilevanti problemi fondamentali e applicativi nel campo matematico. È augurabile che anche questo volume possa trovare, come gli altri, favorevole accoglienza oltre che fra gli addetti ai lavori, pure fra coloro che sono desiderosi di approfondire la loro cultura e le proprie conoscenze sui più rilevanti problemi scientifici del nostro tempo.

C.C.

INDICE

C. CILIBERTO Presentazione, 7

I I FONDEMENTI DELLA MATEMATICA

- M. GARDNER Sui paradossi di Zenone, 13
W. V. QUINE I fondamenti della matematica, 16
P. J. COHEN e R. HERSH La teoria non cantoriana degli insiemi, 25
M. DAVIS e R. HERSH L'analisi non-standard, 34
G. LOLLJ Nuovi modelli del sistema dei numeri reali, 42
B. DE FINETTI Tre personaggi della matematica, 50

II MATEMATICA E LOGICA

- M. GARDNER Algebra di Boole, diagrammi di Venn e calcolo proposizionale, 66
A. TARSKI Verità e dimostrazione, 70
W.W. BARTLY III Un libro di logica smarrito di Lewis Carroll, 80
L. LOMBARDO-RADICE Un nuovo livello di astrazione: la teoria delle categorie, 89
D. COSTANTINI e M. MONDADORI Induzione e probabilità, 96
W.C. SALMON I problemi della conferma, 104
H. DELONG Problemi non risolti dell'aritmetica, 113

III MATEMATICA E REALTÀ

- M. GARDNER La macchina di Turing e la questione da essa sollevata:
 può una macchina pensare? 126
H. WANG Giochi, logica e calcolatori, 130
M. DAVIS e R. HERSH Il X problema di Hilbert, 138
D.E. KNUTH Gli algoritmi, 147
E.C. ZEEMAN La teoria delle catastrofi, 158

Note biografiche e bibliografiche, 175
Indice analitico

I

I FONDAMENTI DELLA MATEMATICA

Sui paradossi di Zenone

di Martin Gardner

Una nuova e divertente variante dei paradossi di Zenone è stata presentata da A. K. Austin dell'Università di Sheffield nel « Mathematics Magazine » del gennaio 1971.

Un ragazzo, una ragazza e un cane si trovano nello stesso punto su una strada diritta. Il ragazzo e la ragazza camminano in avanti — il ragazzo a una velocità di quattro chilometri all'ora e la ragazza a tre chilometri all'ora. Mentre essi avanzano il cane corre avanti e indietro da uno all'altra a una velocità di 10 chilometri all'ora. Supponiamo che ogni inversione di direzione avvenga in un tempo nullo. Dopo un'ora dove sarà il cane e verso chi sarà rivolto?

Risposta: « Il cane può essere in un punto qualsiasi tra il ragazzo e la ragazza, e rivolto verso qualunque parte ». Dimostrazione: « Trascorsa un'ora, ponete il cane in un punto qualunque fra il ragazzo e la ragazza rivolto verso una direzione qualunque. Invertendo la successione temporale del movimento, tutti e tre, cane, ragazzo e ragazza, torneranno allo stesso istante nel punto di partenza ».

Il paradosso di Austin, come ha subito notato Wesley C. Salmon dell'Università dell'Indiana, si può presentare sotto altre innumerevoli forme; una delle

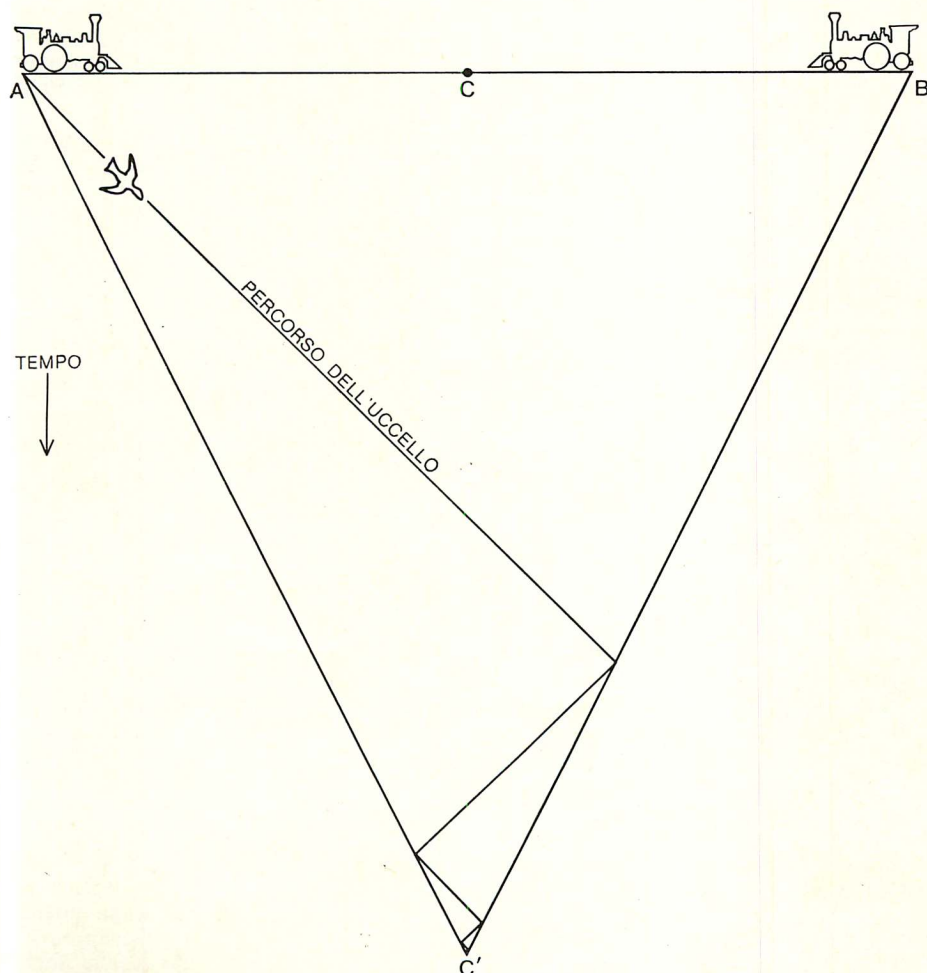


Grafico spazio-tempo del percorso dell'uccello tra le due locomotive in movimento.

più semplici è l'inversione temporale del noto enigma delle due locomotive e dell'uccello. Due locomotive distanti 30 chilometri partendo rispettivamente da *A* e da *B* si muovono sullo stesso binario una contro l'altra alla velocità di 15 chilometri all'ora fino a che si scontrano in *C*. Un uccello partendo da *A*, vola avanti e indietro tra le due locomotive a 60 chilometri all'ora fino allo scontro. Qual è la lunghezza del percorso dell'uccello? In questo caso non occorre ricorrere alla somma di una serie di infiniti termini; dato che l'uccello vola per un'ora, il percorso dovrà essere di 60 chilometri. Se si inverte la successione temporale degli eventi imponendo che l'uccello finisca in *A*, si determina un unico percorso a zig-zag che l'uccello può percorrere in entrambe le direzioni.

Supponiamo però di non aver stabilito dove debba trovarsi l'uccello dopo che le locomotive stanno ritornando verso i punti *A* e *B*. Senza questa informazione non è possibile determinare un percorso unico; infatti l'uccello ora può intraprendere un numero infinito di percorsi e la sola cosa che possiamo dire è che esso alla fine si troverà in un punto qualsiasi tra *A* e *B*.

E' però realmente possibile fare questa affermazione? Molti matematici la contestano poiché nella versione con inversione temporale degli eventi emerge una singolarità che rende contraddittorie le condizioni iniziali. Un matematico potrebbe osservare: « In analisi non esiste una giustificazione valida in generale per invertire l'operatore limite ». Quando le locomotive si muovono una verso l'altra, è soltanto la posizione dell'uccello che converge. « Il vettore velocità diverge, sicché (come nel problema di Austin) si ripresenta la stessa difficoltà di invertire in modo univoco il processo di limite. Lo sviluppo delle regole ormai sancite del calcolo differenziale è dovuto proprio al fatto che tali regole, se si seguono correttamente, evitano contraddizioni ».

Può essere di aiuto tracciare un grafico spazio-temporale del percorso dell'uccello da *A* a *C'* (si veda la figura della pagina precedente). Naturalmente non possiamo disegnare il percorso dell'uccello proprio fino in *C'* poiché gli zig-zag sono infiniti; possiamo però certamente supporre che una linea ideale esista. Se questa linea può scendere da *A* a *C'* è chiaro che non è possibile fare nessuna obiezione logica all'affermazione che essa può salire da *C'* ad *A*. Se non si specifica la destinazione finale dell'uccello, vi è una infinità più che numerabile di diagrammi di questo tipo che partono da *C'* e terminano in un punto qualsiasi del tratto tra *A* e *B*. E' vero che non è possibile risolvere col calcolo problemi analoghi a quello di Austin, se « risolvere » significa fissare proprio la posizione finale del cane, ma una « soluzione » del problema di Austin è precisamente quella che mostra che ciò non è possibile. Poiché non si dice come parta il cane, esso può partire in qualsiasi modo gli piaccia purché rimanga sempre tra il ragazzo e la ragazza e di conseguenza il suo percorso può finire in un punto qualsiasi tra la ragazza e il ragazzo.

Il commento di Salmon sul problema di Austin è stato il seguente. « Quasi tutti hanno udito la vecchia storia dell'uccello che vola avanti e indietro tra due locomotive che si avvicinano (quella che abbiamo raccontato prima). Per rispettare la prospettiva storica, supponiamo che Achille stia inseguendo la tartaruga mentre una mosca troiana vola avanti e indietro fra di loro. Dato un insieme di velocità e di distanze, e la certezza che Achille raggiungerà la tartaruga in un punto e a un istante ben determinato (si veda il mio libro *Zeno's Paradoxes*, Bobbs-Merrill, 1970), possiamo facilmente capire quanto sarà lungo il percorso della mosca. Fin a questo punto non abbiamo nuovi paradossi "zenoniani"... Il problema di Austin è l'inverso temporale del problema del treno e dell'uccello. »

« Sempre per rispettare la prospettiva storica, ritorniamo ad Achille e alla tartaruga. Achille, malgrado lo svantaggio iniziale che gli è imposto per tradizione, raggiunge la tartaruga e per riparare all'ingiustizia impostagli per lungo tempo da Zenone continua a correre aumentando costantemente il suo vantaggio sulla tartaruga. Consideriamo ora la mosca troiana che continua a tentar di volare avanti e indietro tra i due corridori anche dopo che il più veloce ha raggiunto il più lento. Quando Achille e la tartaruga sono alla pari, la mosca si trova nella stessa identica posizione del cane nel paradosso di Austin. »

« Supponiamo per esempio che la tartaruga si muova alla velocità di un chilometro all'ora, Achille a cinque chilometri all'ora (sta correndo dal 500 a.C. perciò non è più veloce come una volta) e la mosca a 10 chilometri all'ora. Tutti arrivano senza difficoltà al comune punto di incontro. Ma come possono proseguire? Se i tre partono contemporaneamente dal punto in comune, la mosca o supera subito ambedue oppure si sposta dietro a entrambi; in ogni caso si viola la condizione iniziale che la mosca stia sempre nell'intervallo tra i due (inclusi gli estremi). Sembrerebbe di poter dire che in ogni intervallo di tempo $\epsilon > 0$, piccolo a piacere, la tartaruga percorre la distanza di 1ϵ , Achille di 5ϵ e la mosca di 10ϵ . Quindi, per quanto piccolo si prenda l'intervallo di tempo, la mosca, dopo l'incontro, non si trova più nel tratto tra la tartaruga e Achille. Perciò anche se abbiamo mostrato come Achille possa compiere la "superimpresa" di raggiungere la tartaruga e a sua volta la tartaruga quella di iniziare il suo moto, si vede che la mosca deve affrontare la "superimpresa" di continuare a volare avanti e indietro tra Achille e la tartaruga anche dopo che questa è stata raggiunta: la mosca deve riuscire a non superare Achille. »

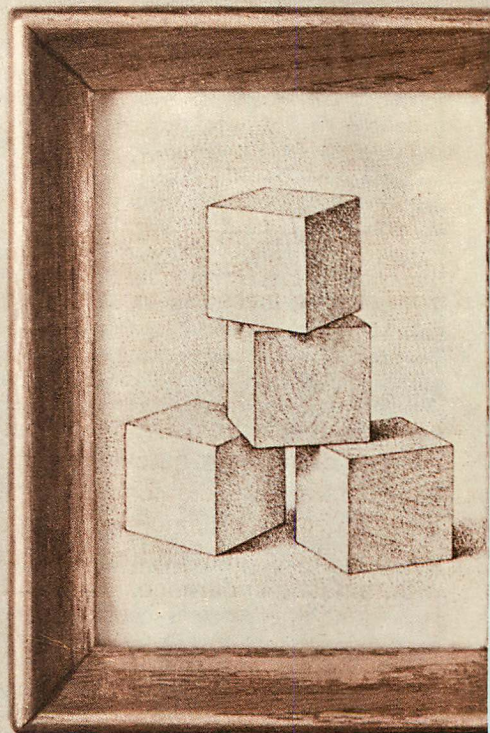
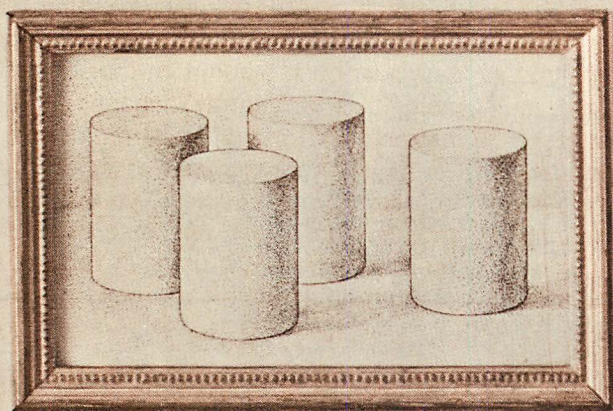
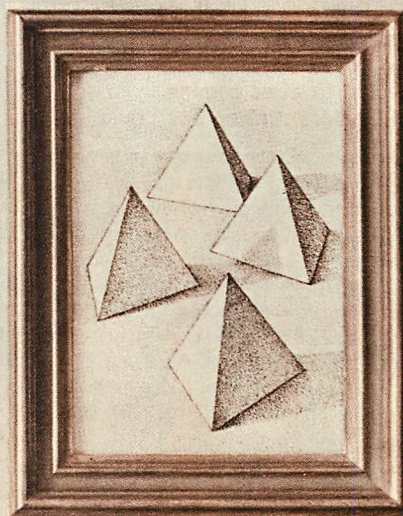
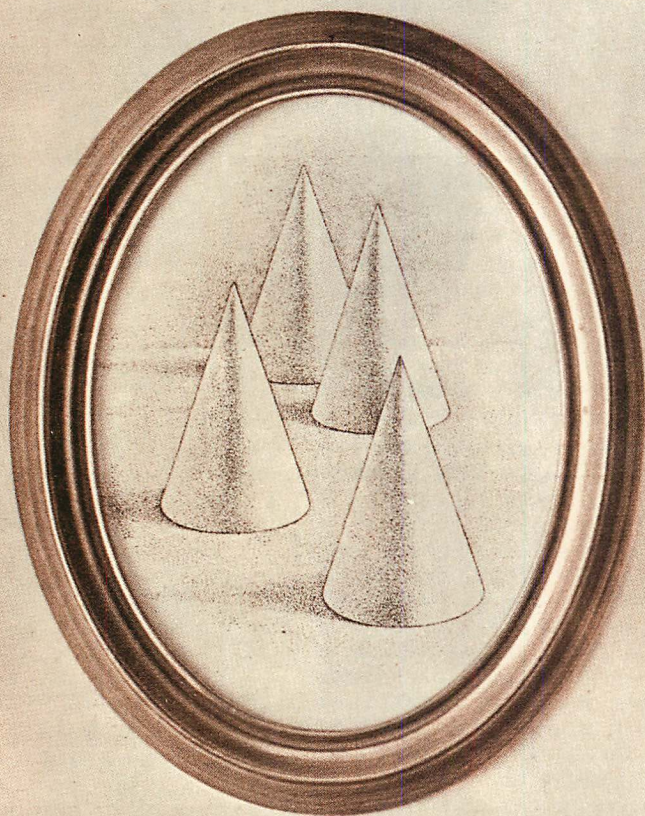
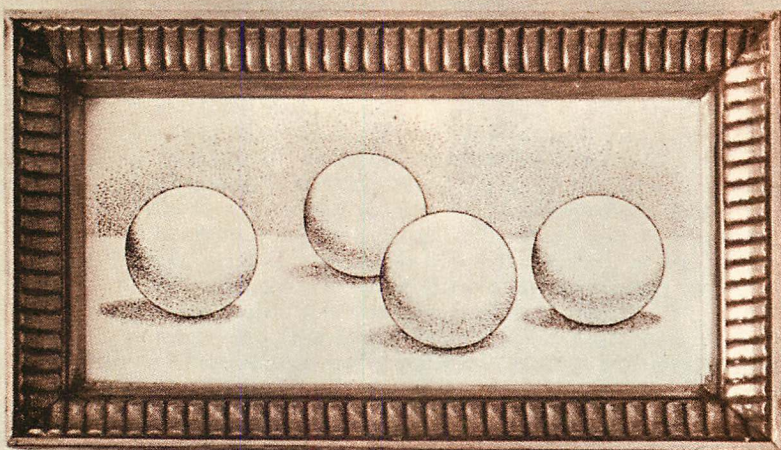
« A me sembra che il problema presenti una difficoltà evidente analoga a

quella sottolineata da Zenone nel suo paradosso della dicotomia regressiva. Non c'è dubbio che se *la mosca vola costantemente in una direzione senza tornare indietro*, distanzierà sia Achille che la tartaruga, anche in un tempo e arbitrariamente piccolo. Questo tuttavia non rende impossibile il moto della mosca poiché per quanto piccolo si scelga l'intervallo di tempo, la mosca in questo intervallo avrà già invertito la sua direzione di moto (un numero infinito di volte, sicché sarà completamente stordita). Questo significa semplicemente che non esiste un intervallo *iniziale* non nullo durante il quale la mosca vola dritta senza invertire la sua direzione; ne segue perciò che la mosca non può lasciare subito l'intervallo tra la tartaruga e Achille. In effetti possiamo vedere esattamente come i rapidi capovolgimenti permettano alla mosca di rimanere sempre tra la tartaruga e Achille se esaminiamo l'inversione temporale del moto della mosca quando essa si avvicina al punto d'incontro. Il fatto che la mosca non possa compiere un percorso *iniziale* dritto non nullo è analogo al fatto che la tartaruga lasciando il suo punto di partenza non può attraversare un qualunque segmento iniziale non nullo del suo percorso. In entrambi i casi la mancanza di un opportuno segmento iniziale non è un grave ostacolo. »

« La letteratura più recente sui paradossi di Zenone riporta molte discussioni sulle "macchine infinite". Si tratta di dispositivi ideali che dovrebbero eseguire una serie infinita di compiti; sono stati introdotti nella discussione a causa delle difficoltà che sembra si incontrino nel portare a termine una serie infinita di compiti (una "superimpresa"). La soluzione dei problemi relativi alle macchine infinite è molto simile alla soluzione del paradosso di Zenone della dicotomia espresso in forma progressiva. Esattamente le stesse considerazioni si possono applicare al moto della mosca troiana fino al momento in cui Achille raggiunge la tartaruga, quest'ultimo istante incluso. Non so se sia stato presentato esplicitamente il tipo di macchina infinita che sarebbe l'analogo del paradosso di Zenone della dicotomia in forma regressiva; una macchina la cui difficoltà consisterebbe nell'iniziare una serie infinita di prove a differenza della macchina infinita usuale la cui difficoltà consiste nel finire la sua serie di prove. La nostra mosca troiana, nel suo moto dal punto d'incontro di Achille con la tartaruga per tutta la parte successiva della corsa in cui Achille è davanti alla tartaruga, si comporta proprio come una macchina infinita (come pure il cane di Austin), potremmo dire una macchina infinita regressiva. Così come esiste uno stretto parallelismo tra il modo di trattare una macchina infinita e la soluzione del paradosso di Zenone della dicotomia in forma progressiva, dovrà pure esistere un parallelismo tra il modo di trattare la mosca troiana nella seconda parte del suo volo e la soluzione del paradosso in forma regressiva. »

« Un altro problema sul movimento della mosca merita particolare attenzione, e precisamente: qual è lo stato di moto della mosca nell'istante esatto dell'incontro? La posizione della mosca è ben determinata: essa coincide con la posizione di Achille e della tartaruga. Matematicamente si può descrivere la posizione della mosca con una funzione continua del tempo che all'istante assegnato passa per il punto di incontro. Viceversa la funzione velocità della mosca è discontinua. Il suo valore è $+10$ quando la mosca si muove in avanti, -10 quando essa si muove indietro e (possiamo dire) è zero quando la mosca si incontra con Achille o con la tartaruga (o con entrambi). Perciò nell'istante in cui tutti e tre si incontrano possiamo appropriatamente assegnare alla velocità della mosca il valore zero. E' ovvio che la funzione velocità in vicinanza del punto d'incontro comune ha da entrambe le parti infiniti punti di discontinuità. A ogni discontinuità *finita* della funzione velocità corrisponde un punto di discontinuità *infinita* per l'accelerazione poiché per cambiare istantaneamente la velocità da $+10$ a -10 e viceversa occorre un'accelerazione infinita. Inoltre, come si vede dal problema di Austin e dalla sua soluzione, lo stato di moto della mosca (o del cane) nel punto d'incontro non determina univocamente il tipo di moto che si avrà dopo. In altre parole, pur avendo mostrato come sia possibile ("possibile" in un certo senso) per la mosca continuare il suo movimento attraverso e oltre il punto d'incontro, ci sono però infiniti modi distinti di eseguire questo moto, tutti consistenti con le condizioni imposte dal problema. Dire che ci sono vari modi alternativi di eseguire un compito non significa tuttavia dimostrare che il compito è impossibile. »

« I paradossi di Achille e della dicotomia, nella formulazione usuale, comportano un numero finito di discontinuità del tipo che si è appena menzionato: si fa l'ipotesi che Achille e la tartaruga, al punto di partenza, accelerino in modo da raggiungere istantaneamente le loro rispettive velocità medie e che alla fine decelerino istantaneamente fino a una velocità nulla. Analogamente moltissime delle "macchine infinite" (per esempio le macchine trasferitrici di Black e la lampada di Thomson) implicano un numero infinito di tali discontinuità che si accumulano attorno a un certo istante finale (si veda *Zeno Paradoxes* alle pagine 204-244). Facendo uso di una funzione matematica introdotta da Richard Friedberg, Adolf Grünbaum ha mostrato come si possano modificare tali movimenti in modo da eliminare tutte le discontinuità pur raggiungendo il risultato complessivo desiderato. Si potrebbe pensare di applicare tale metodo al problema della mosca troiana (o al terzetto di Austin ragazzo-ragazza-cane) al fine di ottenere una descrizione del moto a cui non si possa muovere nessuna obiezione. »



I fondamenti della matematica

Quando si afferma una nuova idea in matematica, da essa sorge una sovrastruttura. Quindi l'idea originale può rivelarsi fonte di difficoltà, che è conveniente eliminare per non distruggere la sovrastruttura

di W.V. Quine

Inconfutabilità, il tuo nome è matematica. Gli studiosi di scienze naturali accettino pure l'evidenza empirica: il matematico chiede dimostrazioni. I criteri di scientificità devono essere diventati ben rigorosi se c'è qualcuno che si preoccupa dei fondamenti della matematica. Dove si potrebbero trovare fondamenti certi come ciò che in questo caso si desidera fondare?

L'interesse per i fondamenti della matematica è stato spesso espresso in situazioni d'emergenza, quando le idee fondamentali iniziano a sembrare instabili e i matematici sono obbligati a esaminarle. L'idea di infinitesimo fu soggetta a un tale esame molto tempo dopo che Isaac Newton e Gottfried Wilhelm von Leibniz svilupparono il calcolo infinitesimale. Questo concetto di quantità frazionaria infinitamente vicina allo zero, tuttavia differente dallo zero, forniva una fondazione allo studio dei tassi di variazione, l'argomento del calcolo differenziale.

Si consideri una macchina che da ferma accelera fino a una velocità di 90 miglia all'ora. Nell'istante in cui l'ago del tachimetro indica 60 sta procedendo alla velocità di un miglio al minuto, negli istanti precedenti la sua velocità era minore, in quelli successivi sarà maggiore. La velocità istantanea di un miglio al minuto non consiste nel compiere un miglio al minuto, perché non è mante-

nuta per un minuto. Di fatto la macchina non la mantiene per nessun tempo. La distanza coperta in ciascun istante è zero e la caratterizzazione «nessun miglio all'istante» elimina la distinzione tra una velocità e l'altra. Così i fondatori del calcolo infinitesimale assunsero l'esistenza di infinitesimi di poco diversi da zero e distinti fra loro. (Le frazioni sempre più piccole ci sono familiari, $1/8$, $1/16$ e così via, ma esse non sono infinitesimi; per definizione un infinitesimo sta in 1 non 16 volte ma un numero infinito di volte.)

Andare alla velocità di un miglio al minuto significa superare la metà di quella distanza infinitesimale in quel tempo infinitesimale. L'assurdità di questo approccio era ovvia, ma il calcolo risultante aveva reso possibile un ragionamento matematico sui tassi di variazione. Sorse così un problema caratteristico riguardante i fondamenti della matematica: come sbarazzarsi dell'infinitesimo sostituendolo con idee più chiare salvando nello stesso tempo la sovrastruttura utile.

Nel XIX secolo, Augustin Cauchy e i suoi seguaci risolsero il problema. Si considerino intervalli di tempo sempre più corti in modo che ciascuno di essi includa l'istante dato. Se su ciascun intervallo si scrive la distanza che la macchina ha percorso in esso, ogni rapporto distanza-tempo sarà vicino a un miglio al minuto se l'intervallo di tempo è corto. Fissato comunque un intervallo di approssimazione, esiste un intervallo di tempo tale che, per tutti gli intervalli all'interno di esso, i rapporti distanza-tempo si avvicineranno a un miglio al minuto con l'approssimazione stipulata. Una successione di rapporti distanza-tempo di questo tipo, determinati su intervalli sempre più piccoli, tende a un limite (che può essere calcolato con la tecnica conosciuta come differenziazione). La nozione di limite concerne distanze brevi ma non infinite-

simali e può essere usata per definire cosa significa viaggiare alla velocità di un miglio al minuto in un dato istante.

L'infinitesimo non è il solo concetto matematico che deve essere legittimato o eliminato. Si pensi ai numeri immaginari (le radici quadrate dei numeri negativi), frutto di un'idea sorta nel XVI secolo. Si elevi al quadrato un numero reale qualsiasi, negativo o positivo. Il risultato è positivo. Che cosa sono allora le radici quadrate dei numeri negativi? Come nel calcolo differenziale, anche in questo caso un esame dei fondamenti deve essere fatto con un occhio volto alla conservazione della sovrastruttura.

La radice quadrata di -1 è l'unità immaginaria i . Gli altri numeri immaginari sono multipli ottenuti moltiplicando i per i numeri reali. Il numero immaginario corrispondente al numero reale 3 è il numero immaginario $3i$. Il numero immaginario corrispondente al numero reale $1/2$ è il numero immaginario $1/2i$. Quello corrispondente al numero reale π è il numero immaginario πi . I numeri immaginari, costituiti in questo modo, si combinano con i numeri reali per mezzo dell'addizione, ottenendo in questo modo $3 + i$, $\pi + 2i$ e così via: questi numeri sono conosciuti come numeri complessi e si rivelano particolarmente utili. Qualsiasi numero complesso è una comoda codificazione di due numeri reali x e y , ciascuno dei quali può essere riottenuto in maniera univoca a richiesta. Questa corrispondenza può essere rappresentata su un piano definito da un asse reale x e da un asse immaginario y (si veda l'illustrazione della pagina seguente in basso). Un esame retrospettivo mostra che il mistero dei numeri immaginari avrebbe potuto essere evitato poiché il ruolo dei numeri complessi può essere giocato da «coppie ordinate» di numeri reali.

L'idea di coppia ordinata è utile in

Nella pagina a fronte il numero quattro è rappresentato come la classe contenente tutte le classi di quattro elementi. La cornice esterna non è chiusa a destra perché gli esempi mostrati non esauriscono tutti gli elementi della classe. Il titolo adeguato per il quadro completo sarebbe «4». Dire che una classe di tetraedri, di coni, di sfere o di cilindri ha quattro elementi significa dire che essa appartiene a «4». Questa concezione del numero è solo una tra le molte che sono state proposte dagli studiosi di fondamenti della matematica.

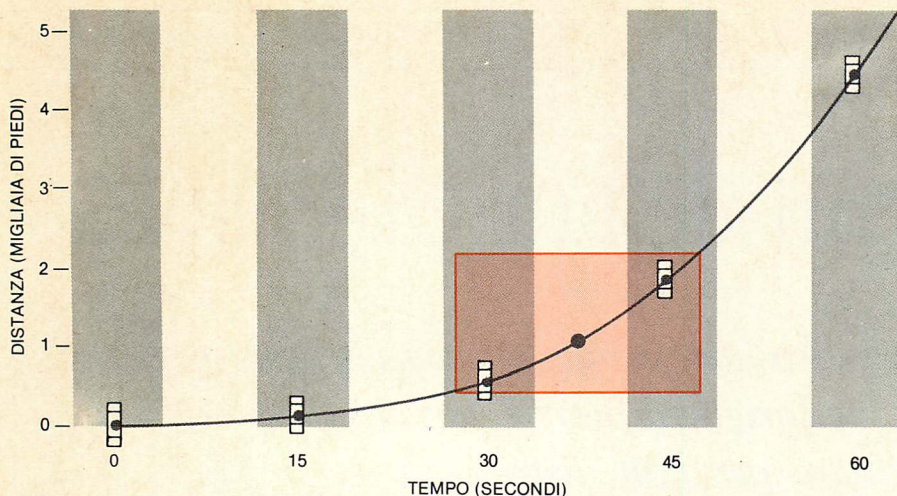
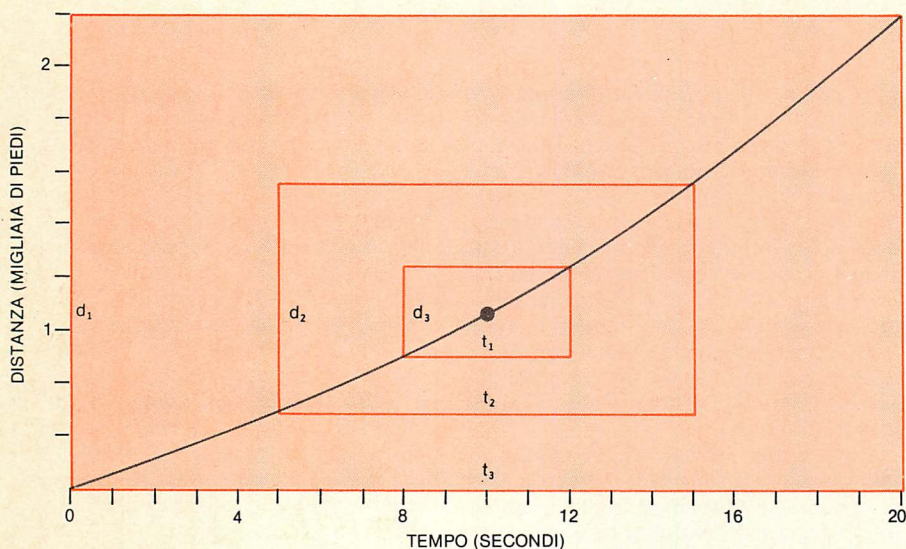
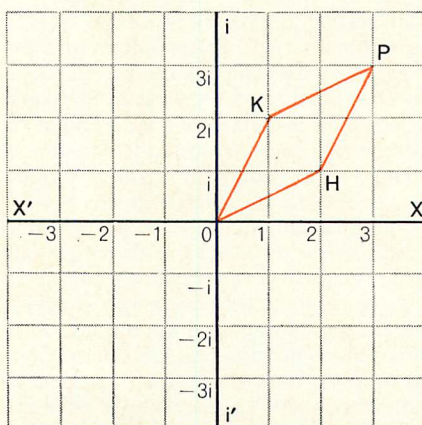
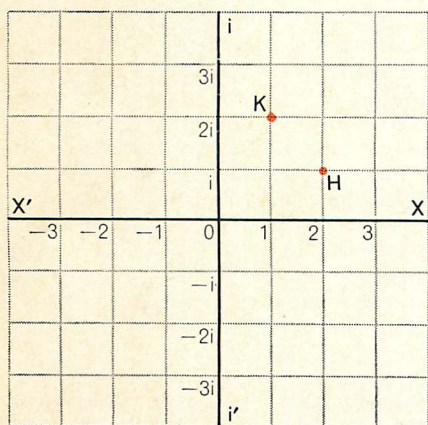


Illustrazione dell'accelerazione da fermo. Le cinque strisce verticali rappresentano il medesimo tratto di strada lungo un miglio, considerato a intervalli di 15 secondi. Nei primi 15 secondi l'auto (in colore) avanza di circa 70 piedi. Nel quarto intervallo avanza 32 volte più lontano.



$$\frac{d_1}{t_1} = \frac{1799,5}{20} = 90 \text{ piedi/secondo} \quad \frac{d_2}{t_2} = \frac{884,1}{10} = 88,4 \text{ piedi/secondo} \quad \frac{d_3}{t_3} = \frac{351,9}{4} = 88 \text{ piedi/secondo}$$

L'istante in cui l'auto raggiunge la velocità di un miglio al minuto è un punto sulla curva di accelerazione. (Questa sezione di curva corrisponde al rettangolo scuro della figura in alto.) I rapporti distanza-tempo in intervalli sempre più brevi tendono a un limite finito e determinabile.



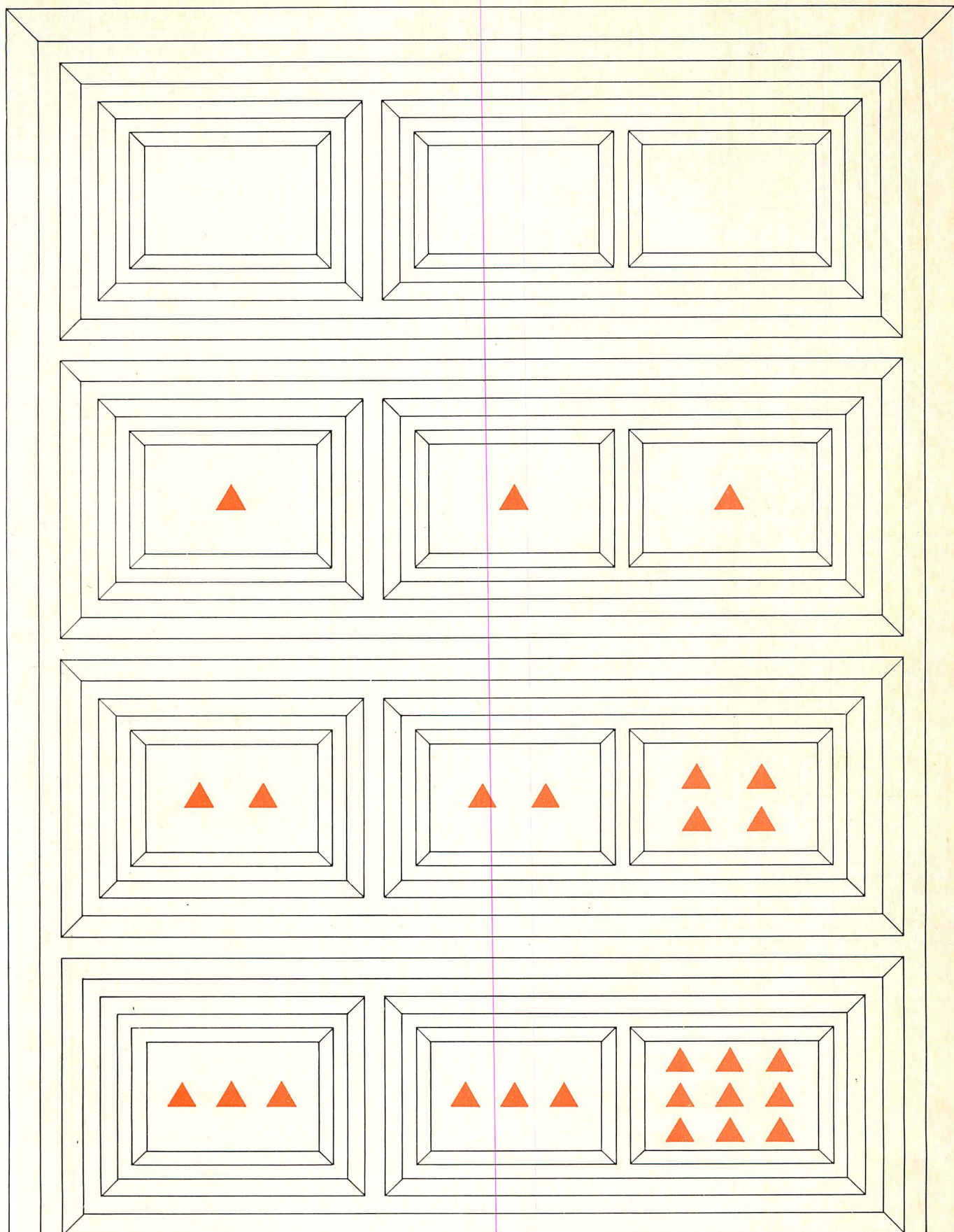
Il numero complesso $x + yi$ può essere rappresentato su un piano definito da un asse orizzontale reale e da un asse verticale immaginario. A sinistra vediamo i punti $K(1 + 2i)$ e $H(2 + i)$. Se questi punti e l'origine sono vertici di un parallelogramma, il quarto vertice (P) è la loro somma.

molte altre situazioni critiche della matematica. Il suo uso è sempre lo stesso: è un modo per considerare due cose come una sola senza per questo confonderle. Comunemente la coppia ordinata di x e di y , si tratti di numeri o di altre cose, come ad esempio padri e figli, è denotata (x, y) . Non ho detto che cosa sia una coppia ordinata e di solito tale problema è evitato: ciò che è importante è come si comporta. L'unica proprietà che importa è che se (x, y) è (z, w) , allora x è z e y è w .

Ho detto che in linea di principio il mito delle radici immaginarie potrebbe essere evitato. Tuttavia ha un valore: semplifica molto le leggi dell'algebra, un vantaggio che può essere conservato pur eliminando i numeri immaginari e quelli complessi. Questo si ottiene con una manovra che è comune negli studi fondamentali: *definendo* i numeri complessi come semplici coppie di numeri reali e quindi ridefinendo le solite operazioni algebriche di somma, moltiplicazione e potenza in modo tale che queste operazioni abbiano senso quando vengano applicate alle coppie ordinate. In questo modo si possono ideare definizioni che forniscono un'algebra delle coppie ordinate che sia formalmente indistinguibile dall'algebra dei numeri complessi. Si potrebbe dire che i numeri complessi sono stati spiegati in termini di coppie ordinate, ma si potrebbe anche dire che sono stati eliminati in favore delle coppie ordinate.

Invece di dire semplicemente ciò che le coppie ordinate fanno, si potrebbe proseguire tentando di determinare che cosa sono. Questo problema non ha l'urgenza dei problemi sugli infinitesimi o sui numeri immaginari e ha un sapore più filosofico. Qualsiasi risposta, per quanto artificiosa, servirà purché confermi la legge delle coppie: se (x, y) è (z, w) , allora x è z e y è w . Oggi la versione abitualmente adottata appartiene a Norbert Wiener e Casimir Kuratowski. Tale versione non identifica la coppia ordinata (x, y) semplicemente con la classe i cui elementi sono x e y , dove si potrebbe confondere (x, y) con (y, x) . Identifica (x, y) con la classe di due classi delle quali una è la classe il cui solo elemento è x , l'altra la classe i cui elementi sono x e y . Si può dire che si è data una spiegazione delle coppie ordinate in termini di classi di classi di due elementi, o che si sono eliminate le coppie ordinate in favore di queste classi di classi di due elementi. La differenza è puramente verbale, ma la prima descrizione ha il vantaggio di conservare la notazione « (x, y) » e la parola «coppia».

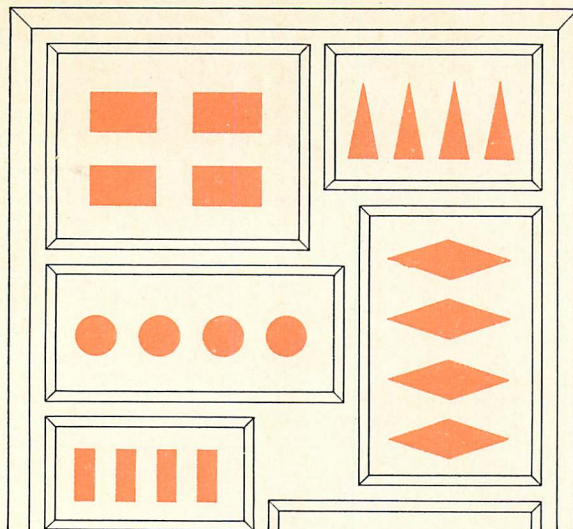
I problemi filosofici sembrano una via di mezzo fra le proteste indignate del senso comune offeso «Che cos'è un infinitesimo? Che cos'è la radice quadrata di un numero negativo?», e le domande petulantanti di un bambino annoiato in un sabato piovoso. Un problema filosofico più profondo di quello riguardante le coppie ordinate è: «Che cosa sono i nu-



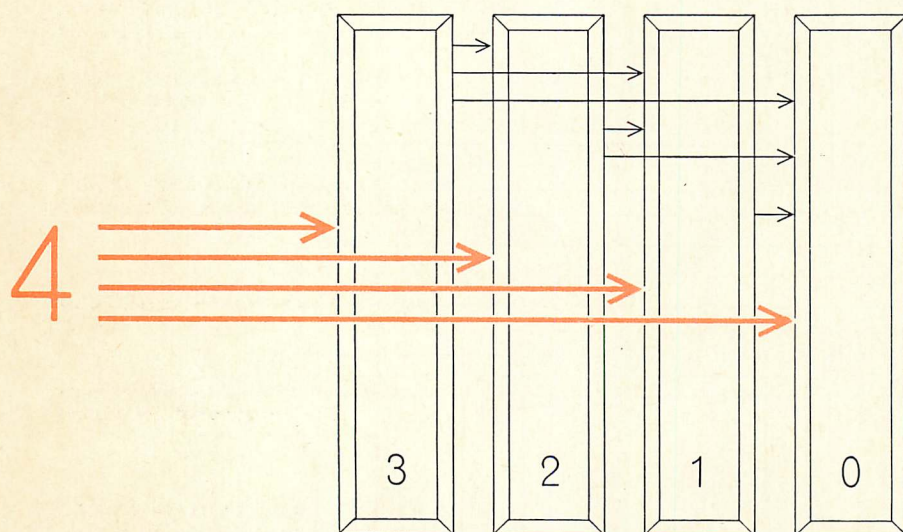
La funzione radice quadrata può essere rappresentata per mezzo di una cornice, aperta in basso, contenente le coppie ordinate (a, b) per le quali vale la relazione « b è uguale ad a volte a ». Ciascuna coppia ordinata viene identificata con una classe di due elementi. Il primo elemento è dato dalla classe che si trova a sinistra in ogni cornice, il cui unico elemento è il primo della coppia. Il secondo elemento è la

classe che si trova a sinistra contenente i due elementi della coppia. Usando questa convenzione, la coppia in basso $(3, 9)$ non può essere confusa con $(9, 3)$, coppia che appartiene alla funzione elevamento al quadrato. L'introduzione di questa rappresentazione astratta delle coppie ordinate, che è usata molto frequentemente nella matematica moderna, è dovuta all'opera di Norbert Wiener e di Casimir Kuratowski.

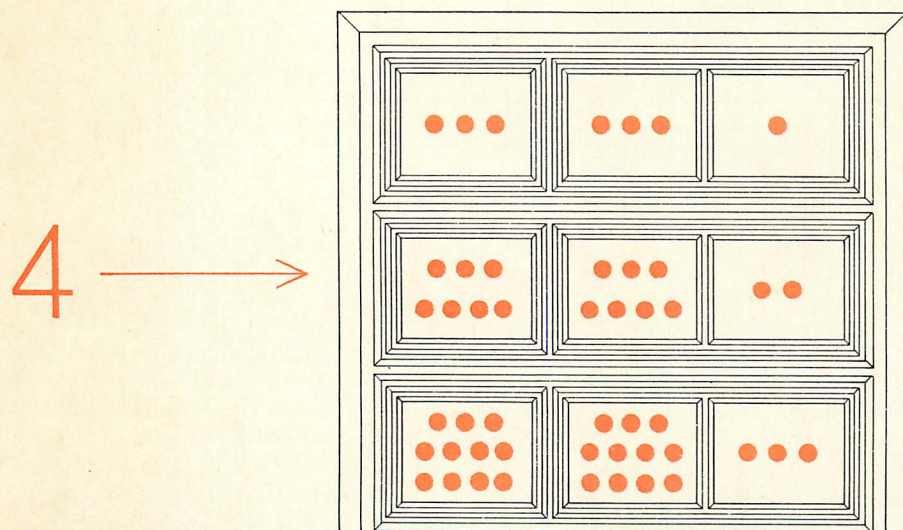
4 →



Il numero naturale quattro viene considerato come la classe di tutte le classi di quattro elementi, quindi la cornice esterna non è chiusa. Seguendo questa definizione di numero, suggerita per la prima volta alla fine dell'ottocento da Gottlob Frege, il numero 1 può essere considerato come la classe di tutte le classi che appartengono a 0 quando siano private di un elemento.



La definizione di John von Neumann del numero naturale 4 sottolinea l'esistenza di una corrispondenza biunivoca (indicata dalle frecce) con una classe di quattro elementi. Tre elementi costituenti questa classe sono definiti per mezzo di una corrispondenza mentre, l'altro è la classe vuota.



Il numero reale 4 è rappresentato come la classe, illimitata in basso (cornice esterna), di tutte le coppie ordinate (a, b) in cui a è meno di quattro volte b . Questa è la relazione che 3 ha con 1, 7 con 2, 11 con 3. Questa convenzione si applica sia ai numeri frazionari sia a quelli negativi.

meri?» Consideriamo questo problema prima in relazione ai numeri naturali che comprendono gli interi positivi e lo zero.

I numerali sono nomi per, ossia denotano, i numeri. Il simbolo «12» denota il numero 12. Possiamo riformulare il problema nel modo seguente: Che cosa denotano i numerali? Che cos'è il numero 12? Vi erano 12 apostoli ci sono 12 uova in un cestino e 12 mesi in un anno. Ma 12 non è solo una proprietà delle dozzine di uova, dei mesi e degli Apostoli, è la proprietà comune alla classe comprendente una dozzina di uova, alla classe comprendente una dozzina di mesi, e a quella comprendente una dozzina di Apostoli.

In matematica gran parte della chiarezza è dovuta alla tendenza a parlare di classi piuttosto che di proprietà. Qualunque cosa possa ottenersi riferendosi a una proprietà, può generalmente ottenersi altrettanto bene anche riferendosi alla classe di tutte le cose che godono di tale proprietà. Si guadagna in chiarezza perché, parlando di classi, si può stabilire più chiaramente un criterio di identità e di differenza; esso consiste semplicemente nell'avere o nel non avere il medesimo numero di elementi.

In particolare, quindi, è meglio spiegare il numero 12 non come la proprietà di essere una dozzina, ma come la classe di tutte le dozzine, la classe di tutte le classi di 12 elementi. Ciascun numero naturale n diventa la classe di tutte le classi di n elementi. La circolarità che consiste nell'usare n per definire n può essere evitata definendo ogni numero nei termini del suo predecessore. Una volta che si è ottenuto il numero 5, per esempio, si può definire il numero 6 come la classe di tutte le classi che, private di un elemento, appartengono a 5. Iniziando dal principio si può definire lo zero come la classe il cui solo elemento è la classe vuota; poi l'uno come la classe delle classi che, private di un elemento, appartengono a zero, poi il due come la classe di quelle classi che private di un elemento, appartengono a uno, e così via.

Proprio come ogni definizione di coppia ordinata serve allo stesso scopo se soddisfa la legge delle coppie, così ogni definizione di numero naturale servirà allo scopo se soddisfa questa legge: esiste un primo numero e un'operazione di successore che produce qualcosa di nuovo ogni volta. La precedente definizione di numero, presentata da Gottlob Frege nel 1884, soddisfa al requisito (si veda l'illustrazione in alto in questa pagina). Ciò vale anche per altre definizioni. La definizione di John von Neumann identifica ciascun numero con la classe di tutti i numeri che lo precedono (si veda l'illustrazione al centro in questa pagina). In questo sistema 0 è la classe vuota; 1 la classe il cui solo elemento è zero e 2 la classe formata da 0 e 1. Dove Frege direbbe che una classe di n elementi appartiene a n , von Neumann dice che una classe di n elementi è quella i cui n elementi possono mettersi in corrispondenza biunivoca con gli elementi di n .

Se si considerano i numeri nell'una e nell'altra di queste definizioni, o in qualche altra ancora, un passo successivo è quello di definire le operazioni aritmetiche. L'idea che sta dietro l'addizione è evidente: $m + n$ è il numero degli elementi di una classe se parte di essa ha m elementi e il resto ne ha n . Il prodotto $m \times n$ è il numero degli elementi di una classe costituita da m parti di n elementi.

Non si è ancora data una spiegazione dei numeri negativi, di quelli frazionari e di tutti i numeri irrazionali come $\sqrt{2}$ e π ; in breve di tutti i numeri reali eccetto i numeri naturali. Qui di nuovo qualsiasi definizione andrà bene se è in accordo con certe richieste. Da un punto di vista generale dotato di grande unitarietà è possibile definire un numero reale come una certa relazione tra nume-

ri naturali: di fatto una relazione basata su un paragone di grandezza. Per esempio, il numero reale $1/2$ si identifica con la relazione che 1 ha con ciascun intero a partire da 3 in poi, e che 2 ha con ciascun intero a partire da 5 in poi, e così via. Similmente, ogni numero reale positivo x s'identifica con la relazione «essere grande meno di x volte». Il numero reale $1/\pi$, per esempio, è conside-

PROPOSIZIONE RIGUARDANTE n	ESPRESSIONI DA DEFINIRE	DEFINIZIONI DELLE ESPRESSIONI
n è un numero primo	numero primo	n è un numero naturale e, presi comunque due numeri naturali h e k , se n è $h \times k$, allora h o k è 1.
n è un numero naturale e, presi comunque due numeri naturali h e k , se $n = h \times k$, allora h o k è 1.	n è $h \times k$	Una classe di n elementi è suddivisa in h parti aventi ognuna k elementi.
n è un numero naturale e presi comunque due numeri naturali h e k , se una classe di n elementi è suddivisa in h parti aventi k elementi ognuna, allora h o k è 1.	x è suddiviso in h parti aventi k elementi ognuna	Esiste una classe y di h elementi tale che ciascun elemento di y ha k elementi e nessun elemento di y ha elementi comuni con qualche altro elemento di y e tutti e soli gli elementi di y sono elementi di x .
n è un numero naturale e presi comunque due numeri naturali h e k , se per ogni elemento x di n esiste un elemento y di h tale che tutti gli elementi di y sono elementi di k e nessun elemento di y ha elementi in comune con qualche altro elemento di y e tutti e solo gli elementi degli elementi di y sono elementi di x , allora h o k è 1.	n è un numero naturale 0 successore 1	n elemento di ogni classe z tale che 0 è elemento di z e tutti i successori di elementi di z sono elementi di z . 0 è la classe il cui solo elemento è la classe senza elementi. Il successore di un m qualsiasi è la classe di tutte le classi che, private di un elemento, appartengono a m . 1 è la classe di tutte le classi che, private di un elemento, diventano la classe senza elementi.
n è un elemento di ogni classe z tale che la classe il cui solo elemento è la classe senza elementi è un elemento di z e, per ogni elemento m di z la classe di tutte le classi che, private di un elemento, appartengono a m , è un elemento di z e, per tutti gli h e k che sono elementi di ogni classe z tale che la classe, il cui solo elemento è la classe senza elementi, è un elemento di z e, per ogni elemento m di z la classe di tutte le classi che, private di un elemento, appartengono a m , è elemento di z , se per ogni elemento x di n esiste un elemento y di h tale che tutti gli elementi di y sono elementi di h e nessun elemento di y ha elementi comuni con qualche altro elemento di y e tutti e soli gli elementi degli elementi di y sono elementi di x , allora h o k è la classe di tutte le classi che, private di un elemento, sono la classe senza elementi.		

La proposizione semplificata è in basso a sinistra in questo schema. La definizione della proposizione « n è un numero primo» produce un vocabolario di espressioni (colonna centrale) che a loro volta devono essere definite (colonna a destra). La proposizione finale ha a che fare

solo con l'appartenenza a classi. Essa potrebbe risciversi usando solo le locuzioni «e», «non», «è un elemento di» e il quantificatore universale «ogni x tale che ... x ...», se non si hanno problemi di brevità. La definizione di numero usata è stata sviluppata da G. Frege.

rato come la relazione che 1 ha con tutti i numeri interi da 4 in poi, che 2 ha con tutti i numeri interi da 7 in poi, che 3 ha con tutti i numeri interi da 10 in poi, e così via. Per quanto riguarda i numeri reali negativi, questi sono considerati come le relazioni converse; poichè $1/2$ è la relazione «essere grande meno della metà», $-1/2$ è la relazione «essere grande più di due volte».

A rischio di aumentare lo strepito, come fa il bambino nel sabato piovoso, si potrebbe chiedere: che cos'è una relazione? Come è stato suggerito nella discussione sulle coppie ordinate, si può identificare una relazione con la classe di tutte le coppie ordinate (a, b) tali che a sta nella relazione in questione con b . In

questo modo i numeri reali sono in ultima analisi identificati con classi, come è accaduto prima per i numeri naturali. Il numero $1/2$ diventa la classe di tutte le coppie ordinate (a, b) tali che a è grande meno della metà di b . Mentre la classe delle coppie ordinate identificate con $1/2$ non può contenere la coppia $(2,4)$ senza violare la relazione «meno di» (in tal caso infatti a sarebbe grande esattamente x volte b), potrebbe invece contenere la coppia ordinata $(20,41)$ o qualsiasi coppia ordinata in cui a sia vicino quanto si vuole all'essere grande x volte b .

Ogni numero reale, dunque, corrisponde a una classe distinta di coppie ordinate. Questa distinzione può essere dimostrata per l'esistenza di un numero razio-

nale tra due punti qualsiasi sulla retta dei numeri reali. Cioè, se x e y sono numeri reali diversi (per esempio x è minore di y), allora esiste un numero razionale a/b (dove a e b sono interi) tale che a/b è minore di y ma maggiore di x . Allora la coppia ordinata (a,b) farà parte della classe che corrisponde a y ma non di quella che corrisponde a x , e in questo modo vengono distinte le due classi.

La descrizione del numero reale positivo x nei termini della relazione «grande meno di x volte» presenta lo stesso tipo di circolarità che si notò nella descrizione di n come la classe di tutte le classi di n elementi. Ma in questo caso, come in quello precedente, la descrizione ci aiuta a capire quali oggetti debbano essere i numeri. La ragione per cui serve è che l'uso circolare di « n » o di « x » all'interno della descrizione dipende dal contesto, che è quello del senso comune. La circolarità può essere effettivamente eliminata in entrambi i casi per mezzo di definizioni accurate e complesse.

Si deve adottare una definizione di numero naturale prima di costruire i numeri reali perché abbiamo considerato questi ultimi come relazioni definite sui numeri naturali. I numeri naturali devono quindi essere considerati distinti, almeno in linea di principio, dai numeri reali corrispondenti. Il numero reale 5, per esempio, è considerato la classe di tutte le coppie ordinate (a,b) di numeri naturali tali che a è grande meno di cinque volte b .

Il concetto di funzione in matematica è importante almeno quanto quello di numero, tuttavia il problema filosofico di che cosa sia una funzione può essere risolto in modo più spiccio, rispetto a quello riguardante i numeri. Una funzione si può identificare con la relazione che esiste tra gli argomenti e i valori della funzione. La funzione «radice quadrata di» può essere considerata come la relazione che sussiste tra la radice e il quadrato: la relazione che 0 ha con 0, che 1 ha con 1, che 2 ha con 4, che $2/3$ ha con $4/9$, e così via. In questo modo la funzione radice quadrata diventa la classe di tutte le coppie ordinate $(0,0)$, $(1,1)$, $(2,4)$, $(2/3,4/9)$, in generale (x,x^2) . L'illustrazione di pagina 19 rappresenta questa funzione come una classe.

Gli esempi di ricerca nel campo dei fondamenti della matematica che abbiamo visto prendevano l'avvio dal presentarsi di certe difficoltà per terminare in una risistemazione generale. Si tratta di un processo di riduzione di alcune nozioni a altre e perciò di riduzione dell'insieme dei concetti matematici fondamentali. Appliciamo questa tecnica nella riduzione di una nozione familiare, cioè la definizione di numero primo, nelle sue componenti elementari, scrivendo dettagliatamente tutte le definizioni successive e prendendo nota di tutti gli strumenti logici e matematici che sono usati in esse. Ciascuna definizione deve spiegare come eliminare qualche locuzione, o come eliminare gli enunciati in cui compaiono tali locuzioni, parafrasandola in

NUMERO	PASSO	GIUSTIFICAZIONE
1	$x = x - (y - y)$,	ASSIOMA
2	$x - (y - z) = z - (y - x)$	ASSIOMA

(Da questi assiomi si derivano i teoremi sostituendo un termine qualsiasi in tutte le occorrenze di una variabile qualsiasi oppure, data un'equazione, rimpiazzando ovunque il suo primo membro al posto del secondo membro).

3	$z = z - (y - y)$	PASSO 1
4	$z = z - (x - x)$	PASSO 3
5	$y - y = (y - y) - (x - x)$	PASSO 4
6	$x - (x - z) = z - (x - x)$	PASSO 2
7	$x - (x - (y - y)) = (y - y) - (x - x)$	PASSO 6
8	$x - x = (y - y) - (x - x)$	PASSI 1,7
9	$x - x = y - y$	PASSI 5,8

(L'espressione « $x + y$ » può definirsi come un'abbreviazione per « $x - ((y - y) - y)$ ». Le leggi dell'addizione diventano così semplici abbreviazioni delle leggi della sottrazione. Quindi la legge « $x + y = y + x$ » si dimostra nel modo seguente.)

10	$x - ((x - x) - y) = y - ((x - x) - x)$	PASSO 2
11	$x - ((y - y) - y) = y - ((x - x) - x)$	PASSI 9,10
12	$x + y = y + x$	PASSO 11, DEFINIZIONE

La procedura di dimostrazione per l'addizione e la sottrazione in aritmetica inizia con equazioni (Passi 1 e 2) prese come assiomi. Poi vengono definiti i passi per derivare i teoremi per mezzo di sostituzioni. I passi 3-9 seguono dai passi precedenti elencati nella colonna a destra. L'addizione viene definita nei termini della sottrazione prima che le sue leggi siano dimostrate (Passi 10-12).

un vocabolario sempre più ristretto che alla fine si riduce a termini elementari. Iniziamo trasformando l'enunciato

n è un numero primo

in

n è un numero naturale e, presi comunque due numeri naturali h e k, se n è $h \times k$, allora h o k è 1.

Il primo passo ha eliminato la locuzione «numero primo» dal vocabolario, ma nel vocabolario restante è rimasta la locuzione «n è un numero naturale», la notazione per la moltiplicazione di h per k e la notazione per 1. Sappiamo come eliminare la notazione per la moltiplicazione scrivendo, al posto di $n = h \times k$, che

una classe di n elementi è ripartita in h parti aventi ciascuna k elementi.

Per sostituire la nozione di «parte» di una classe con il concetto più semplice di appartenenza, questa espressione può essere trasformata nella seguente:

per ogni classe x con n elementi esiste una classe y con h elementi tale che ciascun elemento di y ha k elementi, nessun elemento di y ha elementi comuni con qualche altro elemento di y e tutti e soli gli elementi degli elementi di y sono elementi di x.

La scomodità aumenta velocemente, ma il vocabolario è stato ridotto in termini di appartenenza a classi. Ora si può eliminare «x ha n elementi» ed espressioni analoghe. Se si usa la definizione di Frege dei numeri naturali, l'enunciato menzionato sopra diventa «x è un elemento di n». La frase da cui eravamo partiti viene ora analizzata nel modo seguente:

n è un numero naturale e presi comunque due numeri naturali h e k, se per ogni elemento x di n esiste un elemento y di h tale che tutti gli elementi di y sono elementi di k e nessun elemento di y ha elementi in comune con qualche altro elemento di y e tutti e soli gli elementi degli elementi di y sono elementi di x, allora h o k è 1.

La diminuzione di chiarezza è meno importante della riduzione del vocabolario. Volendo maggior chiarezza le locuzioni eliminate possono essere reintrodotte come definizioni abbreviate.

La prossima espressione da eliminare è «è un numero naturale» (predicato di n, h e k). Dire che n è un numero naturale è come dire che n è 0 oppure un successore di 0 o un successore di quest'ultimo e così via. Frege mostrò come evitare l'idea del «così via» definendo un numero naturale come

un elemento di ogni classe z tale che 0 è elemento di z e tutti i successori di elementi di z sono elementi di z.

Frege considerò lo «0» come la classe il cui solo elemento è la classe senza elementi e il successore di un m qualsiasi come la classe di tutte quelle classi che, private di un elemento, appartengono a m. Se si eliminano le espressioni «0» e «successore» nel riscrivere la precedente versione di «è un numero naturale», e quindi si usa tale risultato per riscrivere la proposizione originaria, si finisce con una lunga storia raccontata con un vocabolario ridotto (si veda l'illustrazione di pagina 21). Anche il numero «1», alla fine della definizione intermedia, viene eliminato, perché 1 è successore di 0. Il vocabolario che rimane è formato da termini che riguardano l'appartenenza a classi e poco altro: un assortimento di particelle logiche elementari come «è», «e», «o», «se...allora», «ogni», «tutti» e simili.

Con un passo ulteriore tutte queste particelle possono essere ridotte a locuzioni fondamentali. Una di queste è «e», inteso come connettivo proposizionale. Un'altra è «non». Una terza è costituita dal quantificatore universale «ogni x tale che ...x ...» (dove ...x... è un'espressione solitamente contenente la variabile x). Il prefisso «ogni x è tale che» viene espresso in modo compatto col simbolo (x). Infine vi è il verbo « \in » che significa «è un elemento di». Questa lista dovrebbe includere anche le parentesi usate per raggruppare le proposizioni. La seguente breve proposizione illustra questo tipo di notazione:

$(x) \text{ non } (y) \text{ non } (x \in y \text{ e non } y \in x)$

Ciò significa: «Ogni cosa è elemento di qualcosa che non è suo elemento».

Ogni proposizione esprimibile nella notazione della matematica pura, sia nell'aritmetica sia nell'analisi (o in altre teorie), può essere parafrasata in questo vocabolario esiguo, anche se non con la stessa brevità di prima. L'espressione a cui siamo arrivati analizzando «n è un numero primo» è ancora chiara in confronto a quella che si otterrebbe limitando il vocabolario alle cinque locuzioni fondamentali. Queste cinque locuzioni non sono raccomandate come una *lingua franca* della matematica, né come un mezzo pratico di computazione. Ma è un fatto interessante da un punto di vista teorico che un così gran numero di idee matematiche possa essere generato a partire da una base così ridotta, e da questa base in particolare.

Quattro delle cinque locuzioni fondamentali appartengono alla logica. Una è caratteristica della teoria degli insiemi: « \in ». Si potrebbe dire che tutte e cinque appartengono alla teoria degli insiemi, dato che le locuzioni logiche sono contenute in ogni teoria e quindi anche nella teoria degli insiemi.

Sembrerebbe quindi che tutta la matematica possa essere riscritta nel vocabolario della teoria degli insiemi. Ogni problema matematico potrebbe essere trasformato in un problema della teoria

degli insiemi, quindi o questo fatto è promettente per la soluzione dei problemi della matematica oppure anche la teoria degli insiemi è essenzialmente problematica quanto la matematica classica.

Si tratta del secondo caso. E l'aspetto peggiore della teoria degli insiemi non è quello di poter formulare in essa proposizioni la cui verità è difficile da dimostrare, ma che in essa si possono formulare proposizioni di cui è anche troppo facile dimostrare contemporaneamente la verità e la falsità. Una di queste è la proposizione

$\text{non } (y) \text{ non } (x) [\text{non } (x \in y \text{ e } x \in x) \text{ e non } (\text{non } x \in y \text{ e non } x \in x)]$.

Parzialmente riscritta, con un occhio volto alla normale comunicazione umana, suona così:

Esiste almeno un y tale che (x) (x \in y se e solo se non x \in x).

Sembra una proposizione vera: basta considerare come y la classe di tutte le cose x tali che x non è elemento di se stesso. Tuttavia tale proposizione è anche falsa: considerato un y come il precedente si potrebbe prendere in particolare x uguale a y e concludere, in modo contraddittorio, che $y \in y$ se e solo se non $y \in y$.

Questo paradosso, scoperto da Bertrand Russell nel 1901, è il più semplice fra i paradossi della teoria degli insiemi. La morale che si può trarre da essi è che, data una condizione necessaria e sufficiente di appartenenza a una classe, questa condizione non garantisce l'esistenza di tale classe. Il paradosso di Russell mostra in particolare che non esiste la classe di tutte le cose che non sono elementi di se stesse. Di conseguenza il grande compito della teoria degli insiemi consiste nel decidere quali classi esistano. Non si conosce alcuna risposta naturale e indiscutibile; quella che sembrava la risposta più naturale, quella che esistesse una classe per ogni condizione di appartenenza, è insostenibile.

A partire dal 1901 si è verificata una proliferazione di teorie degli insiemi, ma nessuna di esse si rivelò migliore delle altre. Anche la non contraddittorietà diventa problematica, perché non si può più aver fiducia nel senso comune per quanto riguarda la plausibilità delle proposizioni. Nella teoria degli insiemi il senso comune è stato screditato dai paradossi. Come fondamento della matematica la teoria degli insiemi è molto meno solida di ciò che è fondato su essa.

Chiaramente non si deve considerare la fondazione della matematica sulla teoria degli insiemi come un modo per dissipare i timori riguardanti la solidità della matematica classica. Ciò che si sta cercando nel valutare i diversi progetti per una teoria degli insiemi è uno schema che riproduca nella sovrastruttura le leggi accettate della matematica classica. Ci troviamo a considerare la teoria degli insiemi come un vocabolario convenient-

temente ristretto in cui formulare un sistema generale di assiomi per la matematica classica, qualunque cosa possano essere gli insiemi.

Un tale programma di assiomatizzazione non può mai essere completato. Non c'è nessuna speranza di trovare una procedura di dimostrazione abbastanza forte da permettere di ottenere tutte le verità della matematica classica, o anche solo quelle dell'aritmetica, e da escludere tutte le proposizioni false. Questo importante risultato fu dimostrato da Gödel nel 1931.

La procedura di dimostrazione per la addizione e la sottrazione mostrata nell'illustrazione di pagina 22 è *completa*: ogni verità che può essere espressa nella notazione può essere dimostrata con tale procedura. Questa notazione, tuttavia, copre solo pochi aspetti dell'aritmetica, tralasciando la moltiplicazione e gli operatori logici. Se la notazione venisse ampliata per far fronte a questi ulteriori scopi, allora nessuna procedura di dimostrazione potrebbe far ottenere tutte le verità esprimibili, evitando nello stesso tempo le proposizioni false. Questo si verifica anche se si limitano i valori delle variabili ai numeri naturali.

Nella notazione della cosiddetta teoria elementare dei numeri possiamo esprimere, per esempio, la seguente proposizione vera:

(x) (y) *non* (z) [*non* ($x = y + z$) e *non* ($y = x + z$)].

Questo equivale a dire che, presi comunque due numeri naturali x e y , o $x = y + z$ oppure $y = x + z$, per qualche numero naturale z . Gödel dimostrò che, data una procedura di dimostrazione, è possibile costruire una proposizione, espressa in questa scarna notazione, che è falsa se può essere dimostrata in tale procedura, e vera se non può essere dimostrata. Quindi Gödel concluse che la procedura data o non è valida, in quanto permette di dimostrare proposizioni false, oppure è incompleta, in quanto non permette di dimostrare proposizioni vere della teoria elementare dei numeri.

La scoperta di Gödel fu uno shock. Infatti si supponeva che l'essenza delle verità matematiche fosse la dimostrabilità. Le cose tuttavia stavano diversamente. Ogni proposizione esprimibile in questa notazione scarna e perspicua della teoria elementare di numeri è dotata di significato, per ogni proposizione vale

che è vera o falsa e che essa o la sua negazione sono derivabili; tuttavia la verità non assicura la dimostrabilità. La differenza tra la verità matematica e quella delle scienze naturali è forse meno netta di quanto si pensi.

Lo studio sui fondamenti della matematica può aver a che fare con concetti e con leggi. In buona parte di questo articolo abbiamo lavorato su concetti, per esempio tramite la riduzione dei concetti per mezzo della definizione di alcuni di essi nei termini di altri.

Ma la scoperta di Gödel si riferisce allo studio delle leggi e alla loro codifica in assiomi e regole di dimostrazione. La importanza di questo tipo di lavoro non fu sminuita dalla conclusione che non si possono ottenere sistemi completi per importanti rami della matematica; se ne possono ottenere di incompleti che sono illuminanti per varie ragioni.

Di fatto il risultato di Gödel ha stimolato fortemente gli studi fondazionali dedicati alle leggi. Le notevoli tecniche usate nella dimostrazione di Gödel hanno dato origine a una importante branca della matematica: la teoria della dimostrazione. Ecco ad esempio un caso in cui la fondazione ha dato origine alla sovrastruttura.

La teoria non cantoriana degli insiemi

Nel 1963 è stato dimostrato che una celebre ipotesi avanzata da Georg Cantor non può essere dimostrata. Questo risultato viene illustrato mediante un'analogia con la geometria non euclidea

di Paul J. Cohen e Reuben Hersh

La teoria astratta degli insiemi si trova attualmente in uno stadio di evoluzione che da molti punti di vista è analogo alla rivoluzione verificatasi in geometria nel XIX secolo. Come in ogni rivoluzione, sia essa politica o scientifica, è difficile, per coloro che vi partecipano o ne sono testimoni, prevederne le estreme conseguenze, salvo forse affermare che esse saranno profonde. Una cosa che si può fare è tentare di usare il passato come guida per il futuro: non si tratta certo di una guida degna di affidamento, ma è sempre meglio che niente.

In questo articolo ci proponiamo di servirci della ormai vecchia storia della geometria non euclidea per illuminare la storia, oggi in pieno sviluppo, della teoria non standard degli insiemi.

Quella di insieme è una delle idee più semplici e primitive della matematica, tanto semplice che ai nostri giorni può essere introdotta a livelli scolastici elementari. Indubbiamente, proprio per questa ragione, la sua natura di concetto fondamentale di tutta la matematica non è stata messa in luce che attorno al 1880: solo in quel periodo Georg Can-

tor fece la prima scoperta non banale nell'ambito della teoria degli insiemi.

Per descrivere la sua scoperta dobbiamo dapprima spiegare cosa intendiamo per insieme infinito. Un insieme infinito è semplicemente un insieme che ha un numero infinito di elementi distinti; ad esempio l'insieme di tutti i numeri « naturali » (1, 2, 3, ecc.) è infinito, come infinito è l'insieme di tutti i punti di un dato segmento.

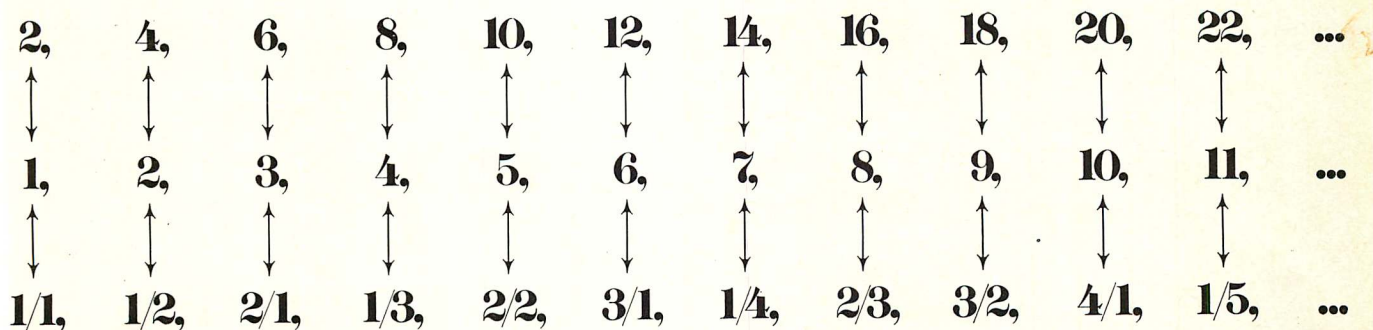
Cantor fece rilevare che anche per insiemi infiniti era sensato parlare del numero degli elementi nell'insieme, o per lo meno aveva senso affermare che due insiemi diversi hanno lo stesso numero di elementi. Proprio come nel caso di insiemi finiti, possiamo dire che due insiemi infiniti hanno lo stesso numero di elementi — la stessa « cardinalità » — se possiamo associare uno per uno gli elementi dei due insiemi, ossia se possiamo stabilire una corrispondenza biunivoca fra i due insiemi. Se ciò è possibile diciamo che i due insiemi sono equipotenti.

L'insieme di tutti i numeri naturali può essere posto in corrispondenza biunivoca con l'insieme di tutti i numeri

pari o con l'insieme di tutte le frazioni (si veda la figura qui sotto). Questi due esempi mettono in evidenza una proprietà paradossale degli insiemi infiniti: un insieme infinito può essere equipotente con uno dei suoi sottoinsiemi, e anzi si dimostra facilmente che un insieme è infinito se, e solo se, è equipotente con un suo sottoinsieme proprio.

Tutto ciò era molto stimolante, ma non soddisfacente per Cantor: la nozione di cardinalità per insiemi infiniti sarebbe stata interessante solo se fosse stato possibile dimostrare che non tutti gli insiemi infiniti hanno la stessa cardinalità. E fu proprio questa la prima grande scoperta di Cantor nella teoria degli insiemi: con una dimostrazione basata sul suo famoso metodo diagonale egli riuscì a dimostrare che l'insieme dei numeri naturali *non* è equipotente con l'insieme dei punti di un segmento (si veda l'illustrazione nella pagina seguente).

Esistono quindi almeno due tipi diversi di infinità. Il primo tipo, l'infinità dei numeri naturali (e di ogni insieme infinito equipotente), viene detta aleph con zero (\aleph_0) e gli insiemi aventi



Un insieme si dice numerabile se può essere posto in corrispondenza biunivoca con i numeri naturali (*successione al centro*). Così, l'insieme di tutti i numeri pari (*successione in alto*) è numerabile e tale risulta anche l'insieme di tutte le frazioni (*successione in basso*). Il metodo qui mostrato è quello usa-

to dal matematico G. Cantor (1845-1918); le frazioni non sono disposte in ordine di grandezza naturale, ma in un ordine determinato dalla somma del numeratore e del denominatore. Entrambi questi esempi mostrano che gli insiemi infiniti, a differenza di quelli finiti, possono essere equipotenti con uno dei propri sottoinsiemi.

1. 0, **1**834798463900...
 2. 0,3**6**94857011092...
 3. 0,50**4**7220017399...
 4. 0,998**0**123010948...
 5. 0,0010**2**30549761...
 6. 0,51546**7**9837123...
 7. 0,551198**7**135042...
 ;
 ;

L'insieme dei numeri reali non è numerabile, come Cantor dimostrò col suo famoso metodo diagonale. Qui è elencato un campione casuale di elementi di tale insieme espressi in forma decimale. Se si prende la prima cifra decimale del primo numero, la seconda cifra del secondo numero e così via (*cifre in colore*) si ottiene un numero reale la cui espressione decimale infinita è 0,1640277... Se in questa espressione si cambia a caso ogni cifra, si dà ottenere ad esempio il numero 0,2751388...; basta un momento di riflessione per accorgersi che questo nuovo numero differisce per almeno una cifra da ogni numero elencato in precedenza. Quindi questo nuovo numero risulta non figurare nel nostro elenco iniziale che in questo modo si dimostra essere incompleto.

cardinalità \aleph_0 vengono detti numerabili. Il secondo tipo di infinità è quello rappresentato da tutti i punti di un segmento e la sua cardinalità viene indicata con una c gotica minuscola (c) che sta per « continuo ». Ogni segmento, di qualunque lunghezza, ha cardinalità c (*si veda la figura della pagina a fronte*). La stessa cardinalità ha ogni rettangolo nel piano, ogni cubo nello spazio e in generale ogni spazio illimitato a n dimensioni, comunque grande sia n .

Fatto il primo passo sulla strada degli infiniti, il successivo segue in modo naturale. Consideriamo il concetto di insieme di tutti i sottoinsiemi di un insieme dato (*si veda la figura a pagina 28*). Se indichiamo con A l'insieme di partenza, questo nuovo insieme viene detto insieme potenza di A e indicato con 2^A . E proprio come a partire da A otteniamo il suo insieme potenza 2^A , così da 2^A possiamo ottenere con un ulteriore passo $2^{(2^A)}$ e così via.

Cantor dimostrò che tanto per A finito quanto per A infinito, 2^A non è mai equipotente con A , sicché l'operazione di formare l'insieme di tutti i sottoinsiemi genera una catena senza fine di insiemi infiniti crescenti e non equipotenti. In particolare, se A è l'insieme

dei numeri naturali è facile dimostrare che 2^A (l'insieme di tutti gli insiemi di numeri naturali) è equipotente col continuo (l'insieme di tutti i punti di un segmento). In simboli,

$$2^{\aleph_0} = c.$$

A questo punto il lettore si chiederà: esiste un insieme infinito la cui cardinalità sia compresa fra \aleph_0 e c ? Esiste cioè su un segmento un insieme infinito di punti che non è equipotente con l'intero segmento né con l'insieme dei numeri naturali?

La questione si presentò a Cantor il quale però non riuscì a trovare un insieme con tali caratteristiche; egli ne concluse — o meglio suppose — che un insieme di questo tipo non esiste. A questa ipotesi di Cantor venne dato il nome di « ipotesi del continuo » e il problema della sua dimostrazione o refutazione figurava come primo in una celebre lista di problemi matematici non risolti presentata da David Hilbert nel 1900. La questione è stata finalmente risolta nel 1963, ma probabilmente in un senso completamente diverso da quello immaginato da Hilbert.

Per affrontare il problema non ci si può fondare più sulla definizione can-

toriana di insieme come « ogni collezione in un tutto di oggetti separati e distinti della nostra intuizione o del nostro pensiero », dal momento che questa definizione, apparentemente così naturale e trasparente, nasconde alcune trappole veramente perfide, come dimostra la brutta esperienza fatta da Gottlob Frege nel 1902. Frege stava per pubblicare il secondo volume della sua opera principale nella quale esponeva una ricostruzione dell'aritmetica in termini di teoria degli insiemi, ossia fondandosi sulla teoria « intuitiva » degli insiemi come allora era nota sulla base dei lavori di Cantor, quando ricevette una lettera dal giovane Bertrand Russell che pubblicò aggiungendo al suo volume una appendice che inizia con le seguenti parole: « A uno scrittore di scienza ben poco può giungere più sgradito del fatto che, dopo completato un lavoro, venga scosso uno dei fondamenti della sua costruzione. Sono stato messo in questa situazione da una lettera del signor Bertrand Russell, quando la stampa di questo volume stava per essere finita ».

La « scossa » data da Russell consisteva nel mettere in evidenza un semplice paradosso. Esistono due tipi di insiemi. Di un primo tipo sono quelli, ad esempio lo « insieme di tutti gli oggetti descrivibili con esattamente dodici parole italiane », che godono della proprietà caratteristica di soddisfare essi stessi alla proprietà che li definisce; in altri termini, insiemi che contengono se stessi come elementi. Dal nome di Russell, li chiameremo insiemi- R . Esistono poi tutti gli altri insiemi, ossia quegli insiemi che non appartengono a se stessi: chiamiamoli insiemi-non R . Ora, dice Russell, consideriamo la collezione di tutti gli insiemi-non R (il termine « collezione » viene qui introdotto come sinonimo conveniente per « insieme »), e chiamiamo M questa collezione. Allora M è un insieme- R o un insieme-non R . Ma se M è un insieme-non R esso appartiene a M per la stessa definizione di M e allora è un insieme- R per la definizione di insieme- R : otteniamo cioè una contraddizione. Se viceversa M è un insieme- R , allora, per la definizione di M , esso non appartiene a M , ossia non appartiene a se stesso, ossia non è un insieme- R , ossia è un insieme-non R : otteniamo ancora una contraddizione.

Morale: l'impiego incondizionato della nozione intuitiva cantoriana di insieme può condurre a contraddizione. La teoria degli insiemi può costituire un sicuro fondamento per la matematica solo se si adotta una tecnica più raffinata, che ci assicuri di non incorrere in antinomie, come furono successiva-

mente chiamate le contraddizioni del tipo di quella scoperta da Russell.

Non era del resto la prima volta che sgraditi paradossi si insinuavano in teorie matematiche apparentemente perfette. I paradossi di Zenone, ad esempio, avevano rivelato ai Greci complessità insospettite nei concetti intuitivi di retta e punto. Possiamo rilevare un'analogia: come Russell aveva riscontrato una contraddizione nell'impiego del concetto intuitivo di insieme, così Zenone aveva trovato una contraddizione nell'impiego incondizionato dei concetti intuitivi di « retta » e « punto ».

Ai suoi inizi, con Talete nel VI secolo a.C., la geometria greca si era fondata su un concetto intuitivo e non specificato di « retta » e « punto ». Circa tre secoli più tardi Euclide aveva dato una sistemazione assiomatica a questi concetti. Per Euclide gli oggetti geometrici erano ancora enti reali intuitivamente noti, ma nella misura in cui divenivano soggetti di un ragionamento geometrico essi venivano specificati sulla base di date proposizioni assunte senza dimostrazione (gli « assiomi » e i « postulati ») a partire dalle quali tutte le loro altre proprietà si sarebbero potute dimostrare come « teoremi ». Non sappiamo se, e in che misura, questa sistemazione era motivata dall'esigenza di superare i paradossi di tipo zenoniano; è indubbio tuttavia che la geometria greca venne resa molto più sicura da questa sistemazione che la faceva dipendere (almeno così credevano e intendevano i geometri greci) solo dall'inferenza logica da un piccolo numero di assunzioni esplicitamente enunciate.

Uno sviluppo analogo per la teoria degli insiemi ha richiesto, invece di 300 anni, solo 35. Se Cantor si può qui avvicinare a Talete — il fondatore della teoria, che si affidava a soli ragionamenti intuitivi — allora la funzione svolta da Euclide viene assunta da Ernst

Zermelo che nel 1908 fondò la teoria assiomatica degli insiemi. Ovviamente, Euclide fu in realtà soltanto un elemento di una lunga successione di geometri greci che crearono la « geometria euclidea »; analogamente, Zermelo fu solo il primo di una mezza dozzina di grandi nomi legati alla creazione della teoria assiomatica degli insiemi.

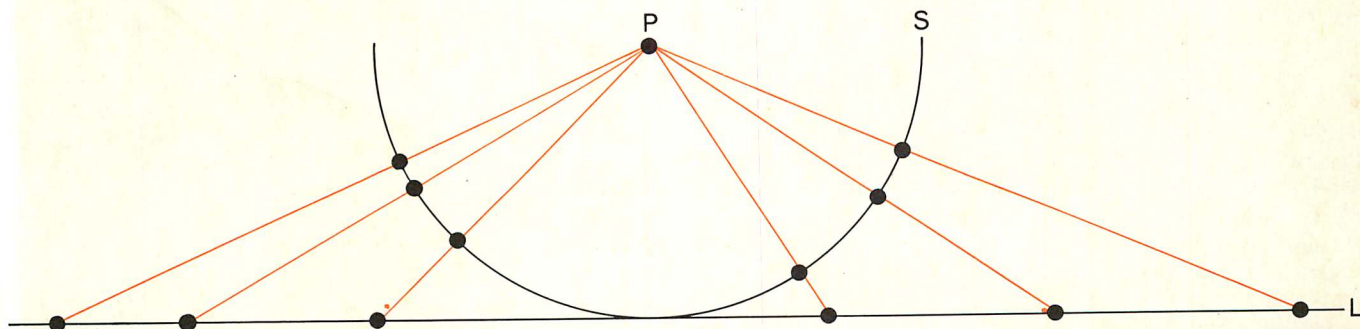
Proprio come Euclide aveva elencato certe proprietà determinate di rette e punti e aveva considerato come dimostrati solo quei teoremi geometrici che potevano essere ottenuti da questi assiomi (e non da ogni possibile ragionamento intuitivo) così nella teoria assiomatica degli insiemi un insieme viene riguardato semplicemente come un oggetto non definito che soddisfa un dato elenco di assiomi. Naturalmente, noi vogliamo pur sempre studiare insiemi (o rette, ad esempio) e quindi gli assiomi non vengono scelti arbitrariamente ma secondo la nostra nozione intuitiva di insieme o di retta; l'intuizione cessa tuttavia di avere qualunque funzione sul piano formale: noi accettiamo solo quelle proposizioni che derivano dagli assiomi. Il fatto che gli oggetti descritti da questi assiomi esistano nel mondo reale è irrilevante per il processo della deduzione formale (malgrado sia essenziale per la scoperta).

Noi conveniamo di comportarci come se i simboli per « retta », « punto » e « angolo » in geometria, o i simboli per « insieme », « è un sottoinsieme di », ecc. nella teoria degli insiemi, siano dei puri segni sulla carta, che possono essere manipolati solo in conformità con un dato sistema di regole (assiomi e regole di inferenza). Sono accettati come teoremi solo quelle proposizioni che si ottengono in conformità a tali manipolazioni di simboli (in pratica sono accettate solo quelle proposizioni che evidentemente *potrebbero* essere ottenute in questo modo se si avesse tempo e voglia di farlo).

Ora, nella storia della geometria un postulato ha avuto una parte del tutto particolare; si tratta del postulato della parallela, il quale afferma che in un piano, per un punto dato, può essere condotta soltanto una retta parallela a una retta data. La difficoltà relativa all'assumere questa proposizione come assioma è data dal fatto che essa non ha quel chiaro carattere di evidenza che si vorrebbe riscontrare negli elementi fondamentali di una teoria matematica. E, in effetti, due rette parallele sono definite come quelle rette che non si incontrano mai anche se prolungate indefinitamente (« all'infinito »); poiché ogni retta che noi possiamo tracciare su un foglio di carta o sulla lavagna ha lunghezza finita, questo è un assioma che per sua natura non può essere verificato da una diretta osservazione sensibile: nondimeno, esso è indispensabile nella geometria euclidea. Per molti secoli un problema fondamentale in geometria fu proprio quello di *dimostrare* il postulato della parallela, per stabilire che esso poteva essere ottenuto come teorema dagli assiomi euclidei più evidenti.

Anche nella teoria astratta degli insiemi c'era un assioma che alcuni matematici trovavano difficile ingoiare; era l'assioma di scelta, che afferma quanto segue: se α è una qualunque collezione di insiemi $\{A, B, \dots\}$ e nessuno degli insiemi di α è vuoto, allora esiste un insieme Z che contiene esattamente un elemento di A , un elemento di B , e così via per tutti gli insiemi di α . Se, ad esempio, α consiste di due insiemi, l'insieme di tutti i triangoli e l'insieme di tutti i quadrati, allora si soddisfa in modo evidente all'assioma di scelta: basta semplicemente scegliere un particolare triangolo e un particolare quadrato e con questi due elementi costituire l'insieme Z .

La maggior parte dei matematici ritiene che l'assioma di scelta, come il postulato della parallela, sia intuitiva-



Si può dimostrare che vi è corrispondenza biunivoca anche fra una retta infinita e un segmento. P è il centro di una semicirconferenza S tangente a una retta L . Una semiretta da P incontra S in un solo punto, sicché le semirette da P forniscono un'asso-

ciazione biunivoca fra punti di S e punti di L . Al variare della direzione della semiretta, non viene omissa alcun punto, né di S né di L . Così fra i punti di una retta infinita e i punti di un segmento di lunghezza arbitraria esiste corrispondenza biunivoca.

$$A: \{ \blacksquare, \blacktriangle, \bullet \}$$

$$2^A: \left\{ \begin{array}{l} \{ \} \\ \{ \blacksquare \} \quad \{ \blacksquare, \blacktriangle \} \\ \{ \blacktriangle \} \quad \{ \blacksquare, \bullet \} \\ \{ \bullet \} \quad \{ \blacktriangle, \bullet \} \quad \{ \blacksquare, \blacktriangle, \bullet \} \end{array} \right\}$$

Esemplificazione del concetto di insieme di tutti i sottoinsiemi di un insieme dato. Il quadrato, il triangolo e il cerchio costituiscono l'insieme A di tre elementi (*in alto*). Questo insieme ha 2^3 sottoinsiemi, ossia 8 sottoinsiemi, una volta che si convenga di considerare tali, anche se impropriamente, l'insieme vuoto e l'intero insieme iniziale A . Questo nuovo insieme formato da 8 elementi viene detto insieme potenza di A e denotato col simbolo 2^A . Se A ha n elementi, l'insieme potenza di A ha 2^n elementi. Se A è infinito, anche 2^A è infinito e non è equipotente con l'insieme A di partenza.

mente plausibile. Le difficoltà nel nostro caso sorgono dall'ammettere per α « qualunque » collezione di insiemi. Come abbiamo visto, esistono catene senza fine di insiemi infiniti sempre più grandi; per una tale collezione inconcepibilmente enorme di insiemi non c'è alcun modo di operare effettivamente una scelta da ognuno dei suoi insiemi membri. Se noi accettiamo l'assioma di scelta, ciò significa semplicemente che noi accettiamo come atto di fede che tale scelta è possibile, proprio come l'accettazione del postulato della parallela significa fare un atto di fede sul comportamento delle rette all'infinito. Ne viene che dall'apparentemente innocuo assioma di scelta derivano alcune conclusioni del tutto inattese ed estremamente potenti. Ad esempio, possiamo usare il ragionamento induttivo per dimostrare proposizioni sugli elementi di *ogni* insieme, praticamente nello stesso modo con cui l'induzione matematica può essere usata per provare teoremi sui numeri naturali 1, 2, 3, ecc.

L'assioma di scelta assolve a una funzione del tutto particolare nella teoria degli insiemi e molti matematici ritengono che il suo impiego dovrebbe essere evitato ogniquale volta è possibile. Una teoria assiomatica degli insiemi ove *non* si assume l'assioma di scelta come vero o come falso sarebbe una

teoria che riscuoterebbe la fiducia della maggior parte dei matematici. Nel seguito useremo la dizione « teoria ristretta degli insiemi » per un sistema assiomatico di questo tipo, mentre con « teoria standard degli insiemi » intenderemo la teoria basata sull'intero sistema di assiomi elaborato da Zermelo e Abraham Fraenkel: teoria ristretta degli insiemi *più* assioma di scelta.

Nuova luce sull'argomento fu gettata nel 1938 da Kurt Gödel. Gödel è meglio noto per i suoi grandi teoremi di « incompletezza » del 1930-1931, ma qui noi ci riferiamo a un suo lavoro successivo poco noto ai non matematici. Nel 1938 Gödel dimostrò il seguente fondamentale risultato: se la teoria ristretta degli insiemi è consistente, allora è tale anche la teoria standard degli insiemi. In altre parole, l'assioma di scelta non è più pericoloso degli altri assiomi; se nella teoria standard può essere scoperta una contraddizione, allora una contraddizione deve già trovarsi nell'ambito della teoria ristretta degli insiemi.

Ma Gödel non dimostrò solo questo. Ricordiamo al lettore la « ipotesi del continuo » di Cantor, l'ipotesi cioè che non esiste alcun cardinale infinito che è maggiore di \aleph_0 e minore di \mathfrak{c} . Gödel dimostrò anche che possiamo as-

sumere senza danni l'ipotesi del continuo come ulteriore assioma della teoria degli insiemi; in altri termini, se la teoria ristretta degli insiemi con l'aggiunta dell'ipotesi del continuo implica una contraddizione, allora, anche in questo caso, deve essere presente una contraddizione già all'interno della teoria ristretta degli insiemi. Questa era una soluzione parziale del problema di Cantor: non era una *dimostrazione* dell'ipotesi del continuo, ma solo del fatto che essa non può essere refutata (ossia che non può essere dimostrata la sua negazione).

Per comprendere come Gödel poté ottenere questo risultato, dobbiamo capire che cosa si intende per modello di un sistema di assiomi. Ritorniamo per un momento agli assiomi della geometria piana. Se prendiamo questi assiomi, incluso il postulato della parallela, abbiamo gli assiomi della geometria euclidea; se invece prendiamo tutti gli assiomi precedenti ma al posto del postulato della parallela assumiamo la sua negazione, otteniamo gli assiomi di una geometria non euclidea. Per entrambi questi sistemi di assiomi (euclideo e non euclideo) chiediamo: possono questi assiomi condurre a una contraddizione?

Può sembrare del tutto irragionevole porre la questione per il sistema euclideo: come può esserci qualcosa di sbagliato nella nostra geometria liceale, vecchia di oltre 2000 anni? D'altra parte il non matematico è certamente in un atteggiamento alquanto sospettoso nei riguardi del secondo sistema d'assiomi, che contiene la negazione del postulato della parallela, intuitivamente pur sempre plausibile. Dal punto di vista della matematica del XX secolo però, i due tipi di geometria stanno più o meno nella stessa posizione: sono entrambi applicabili, in date circostanze, al mondo fisico, e entrambi sono consistenti in un senso relativo che ora spiegheremo.

Mostriamo dapprima che la geometria non euclidea è consistente. Per far questo, sostituiamo la parola « retta », ovunque essa figura, con l'espressione « circonferenza massima » (la curva formata sulla superficie di una sfera da un piano passante per il suo centro). Riguardiamo ora gli assiomi come proposizioni sui punti e le circonferenze massime di una data sfera; inoltre conveniamo di identificare in un unico punto ogni coppia di punti diametralmente opposti sulla sfera. Se il lettore preferisce, può immaginare gli assiomi della geometria non euclidea trascritti, con la parola « retta » ovunque sostituita da « circonferenza massima » e la parola « punto » ovunque sostituita da

« coppia di punti (diametralmente opposti) ». È allora evidente che tutti gli assiomi sono veri, almeno nella misura in cui sono vere le nostre ordinarie intuizioni e nozioni sulla superficie di una sfera. In effetti, dagli assiomi della geometria euclidea solida si può dimostrare facilmente che la superficie di una sfera è una superficie non euclidea nel senso che abbiamo descritto. In altre parole vediamo ora che se gli assiomi della geometria non euclidea conducono a una contraddizione, allora anche l'ordinaria geometria euclidea della sfera condurrebbe a una contraddizione. Abbiamo quindi una dimostrazione *relativa* di consistenza: se la geometria euclidea tridimensionale è consistente, è tale anche la geometria non euclidea bidimensionale. Diciamo che la superficie della sfera euclidea è un modello per gli assiomi della geometria non euclidea. (Nel particolare modello che abbiamo usato il postulato della parallela non è soddisfatto perché non esistono « rette » parallele. È anche possibile costruire una superficie, la « pseudosfera », per la quale il postulato della parallela è falso perché per un punto passa più di una « retta » parallela e una « retta » data).

La scoperta della geometria non euclidea e il riconoscimento del fatto che la sua consistenza è implicata dalla consistenza della geometria euclidea è dovuto all'opera congiunta di numerosi grandi matematici del XIX secolo; ci limitiamo qui a ricordare, in particolare il nome di Bernhard Riemann. Solo nel XX secolo è stata sollevata la questione della consistenza della stessa geometria euclidea.

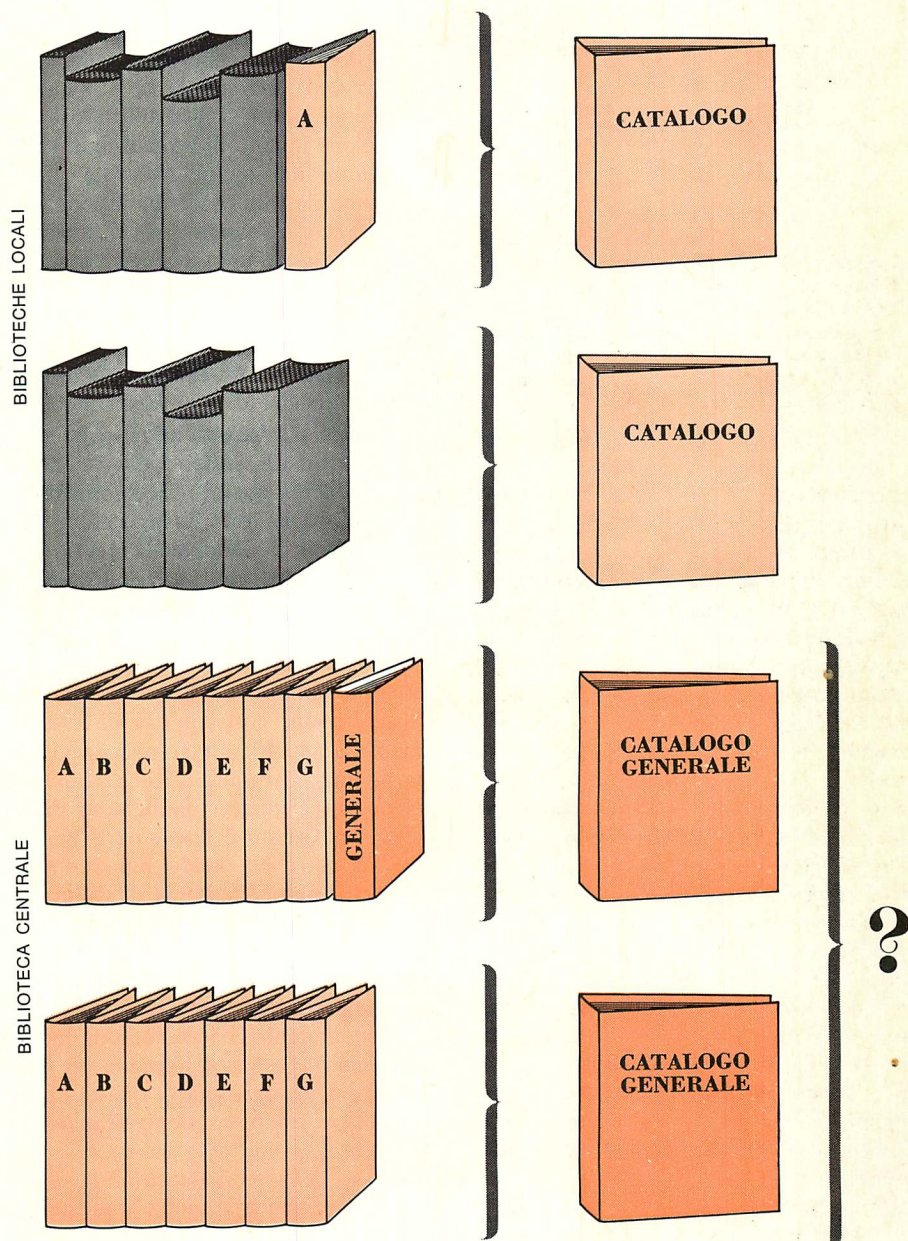
Questo problema fu sollevato da David Hilbert il quale ne diede anche la soluzione, che era una semplice applicazione del concetto di sistema di riferimento cartesiano. Come è noto, a ogni punto del piano possiamo associare una coppia di numeri: la sua ascissa x e la sua ordinata y . A ogni retta o circonferenza possiamo allora associare un'equazione, ossia una relazione fra le coordinate x e y , che è verificata solo dai punti che giacciono sulla retta o sulla circonferenza. Si stabilisce così una corrispondenza fra geometria e algebra elementare: per ogni proposizione espressa in uno dei due campi esiste una corrispondente proposizione nell'altro. Ne segue che gli assiomi della geometria euclidea possono condurre a una contraddizione solo se le regole dell'algebra elementare — le proprietà degli ordinari numeri reali — possono condurre a una contraddizione. Anche in questo caso abbiamo una dimostrazione *relativa* di consistenza. La geo-

metria non euclidea era consistente se lo era la geometria euclidea; ora, la geometria euclidea è consistente solo se lo è l'algebra elementare. La sfera euclidea era un modello per il piano non euclideo, l'insieme di coppie di coordinate è a sua volta un modello per il piano euclideo.

Tenendo presenti questi esempi, possiamo dire che la dimostrazione di Gödel della consistenza relativa dell'assioma di scelta e dell'ipotesi del continuo è analoga alla dimostrazione di Hilbert della consistenza relativa della geometria euclidea: in entrambi gli esempi la teoria standard veniva giustificata in

termini di una teoria più elementare. Ovviamente, nessuno aveva mai dubitato seriamente della sicurezza della geometria euclidea, mentre eminenti matematici come L.E.J. Brouwer, Hermann Weyl e Henri Poincaré sollevavano grossi dubbi circa l'assioma di scelta. In questo senso il risultato di Gödel ebbe una portata e un significato ben maggiori.

Una evoluzione analoga a quella della geometria non euclidea — che potremmo chiamare la teoria non cantoriana degli insiemi — si è avuta solo dal 1963 in un lavoro di uno degli estensori di questo articolo (Cohen). Co-



Il paradosso di Russell viene illustrato supponendo che i bibliotecari siano soliti schedare i loro libri non servendosi di schede ma usando un catalogo a fogli staccabili, in modo cioè che anche il catalogo sia un libro. Alcuni bibliotecari elencano il catalogo stesso nel catalogo (*in alto*) altri invece no (*seconda riga dall'alto*). Il primo tipo di cataloghi viene detto, dal nome di Russell, un insieme- R : gli insiemi- R sono insiemi che contengono se stessi. Ma cosa accade se il direttore centrale delle biblioteche decide di fare un catalogo generale di tutti i cataloghi che non elencano se stessi? Il catalogo che ne risulta deve essere compreso in questo catalogo generale oppure no?

« NOZIONI COMUNI »

1. Cose uguali a una stessa altra cosa sono tra loro uguali.
2. Se a cose uguali si aggiungono cose uguali si ottengono cose uguali.
3. Se da cose uguali si tolgono cose uguali si ottengono cose uguali.
4. Cose che coincidono con un'altra sono tra loro uguali.
5. Il tutto è maggiore della parte.

« POSTULATI »

Si richiede:

1. Che da ogni punto ad ogni altro punto sia possibile condurre una linea retta.
2. Che un segmento di linea retta possa essere indefinitamente prolungato in linea retta.
3. Che attorno a un centro scelto a piacere con un raggio scelto a piacere sia possibile tracciare una circonferenza.
4. Che tutti gli angoli retti siano tra loro uguali.
5. Che se una retta, intersecando altre due rette, forma con esse angoli interni da una medesima parte la cui somma sia minore di due angoli retti, allora queste due rette indefinitamente prolungate finiscano con l'incontrarsi da quella parte nella quale gli angoli interni formano meno di due retti.

Gli assiomi di Euclide erano di due tipi: «nozioni comuni» e «postulati». Il fisico e matematico scozzese John Playfair (1748-1819) ha legato il suo nome a un assioma che si dimostra essere equivalente al quinto postulato di Euclide: in un piano, per un punto A non giacente su una data retta m passa una retta che non interseca m . Si ottiene una geometria non euclidea sostituendo «una» con «nessuna» oppure con «più d'una». Va notato che gli assiomi di Euclide non sono chiari o completi per le nostre esigenze moderne.

sa si intende per «teoria non cantoriana degli insiemi»? Proprio come la geometria euclidea e quella non euclidea si fondavano sugli stessi assiomi con l'unica eccezione del postulato della parallela, così la teoria standard («cantoriana») e quella non standard («non cantoriana») degli insiemi differiscono in un solo assioma. La teoria non cantoriana degli insiemi assume gli assiomi della teoria ristretta degli insiemi e non assume l'assioma di scelta bensì una delle possibili forme di negazione dell'assioma di scelta. In particolare possiamo assumere come assioma la negazione dell'ipotesi del continuo. Quindi, come spiegheremo più avanti, disponiamo ora di una soluzione completa del problema del continuo: alla scoperta, compiuta da Gödel, che l'ipotesi del continuo non è refutabile si è ora aggiunto il fatto che essa non è neppure dimostrabile.

Tanto il risultato di Gödel quanto quello più recente di Cohen richiedono la costruzione di un modello, analogamente a quanto accadeva per le dimostrazioni di consistenza della geometria che abbiamo descritto sopra. In entrambi i casi intendiamo dimostrare che se la teoria ristretta degli insiemi è consistente, è tale anche la teoria standard degli insiemi (o la teoria non standard).

L'idea di Gödel fu quella di costruire un modello per la teoria ristretta degli insiemi e di dimostrare che in questo modello l'assioma di scelta e l'ipotesi del continuo sono teoremi. Il suo procedimento è il seguente. Usando solo gli assiomi della teoria ristretta degli insiemi (*si veda la figura nella pagina a fronte*) siamo intanto garantiti, per l'assioma 2, dell'esistenza di almeno un insieme (l'insieme vuoto); gli assiomi 3 e 4 ci garantiscono allora dell'esistenza di una successione infinita di insiemi finiti sempre più grandi; quindi l'assioma 5 ci garantisce l'esistenza di un insieme infinito; ancora l'assioma 7 ci assicura dell'esistenza di una successione senza fine di insiemi infiniti sempre più grandi (non equipotenti) e così via. Sostanzialmente per questa via Gödel specificò una classe di insiemi in base al modo col quale essi potevano effettivamente essere costruiti in passi successivi a partire da insiemi più semplici. Gödel chiamò questi insiemi «insiemi costruibili»; la loro esistenza era garantita dagli assiomi della teoria ristretta degli insiemi. A questo punto egli fece vedere che nell'ambito degli insiemi costruibili si potevano dimostrare tanto l'assioma di scelta quanto l'ipotesi del continuo; vale a dire (1) che da ogni collezione costruibile α di insiemi costruibili (A, B, \dots) si può scegliere un insieme costruibile Z che contiene almeno un elemento da ognuno degli A, B , ecc. Questo è l'assioma di scelta che ora potrebbe essere più propriamente chiamato il teorema di scelta. E, (2), che se A è un insieme costruibile infinito, allora non esiste alcun insieme costruibile «fra» A e 2^A (maggiore di A o minore dell'insieme potenza di A e non equipotente con nessuno di questi due insiemi). Se per A si prende il primo cardinale transfinito, quest'ultima proposizione è l'ipotesi del continuo.

Nel caso della teoria *costruibile* degli insiemi venne quindi dimostrata un'ipotesi generalizzata del continuo». Il lavoro di Gödel avrebbe quindi risolto completamente questi due problemi se fossimo disposti ad assumere l'assioma che esistono solo insiemi costruibili. Perché non lo facciamo? Perché non ci sembra ragionevole sostenere che un insieme debba essere

costruito in conformità con una qualunque formula prescritta per riconoscerlo come un vero e proprio insieme. Così nell'ordinaria teoria degli insiemi (non necessariamente costruibili) non sono stati dimostrati né l'assioma di scelta né l'ipotesi del continuo. Almeno di questo si era certi: ognuno di essi poteva essere assunto senza pericolo di contraddizione a meno che gli assiomi «sicuri» della teoria ristretta degli insiemi non fossero già di per sé contraddittori; qualunque contraddizione essi causino deve essere già presente nella teoria costruibile degli insiemi che è un modello per la teoria degli insiemi ordinaria. In altre parole, si sapeva che nessuno di essi poteva essere refutato a partire dagli altri assiomi ma non si sapeva se essi potevano essere dimostrati oppure no.

Qui l'analogia col postulato della parallela della geometria euclidea diviene particolarmente stretta. Fino a un'epoca molto recente si assumeva per certo che gli assiomi di Euclide fossero consistenti; la questione che interessava i geometri era se essi fossero o no indipendenti, ossia se il postulato della parallela potesse venir dimostrato a partire dagli altri. Intere generazioni di geometri tentarono di dimostrare il postulato della parallela mostrando che la sua negazione conduce a delle assurdità. Sembra che il primo a rendersi conto che queste «assurdità» altro non erano che teoremi di una nuova geometria non euclidea sia stato Karl Friedrich Gauss, il quale però se pure ebbe il coraggio di pensare queste cose non ebbe il coraggio di pubblicarle. Toccò così a János Bolyai, Nikolai Ivanovic Lobachevskij e Riemann trarre le conseguenze logiche che derivavano dal negare il postulato della parallela. Queste conseguenze erano la scoperta di geometrie «fantastiche» che erano altrettanto consistenti della geometria euclidea del «mondo reale». Solo dopo molto tempo si riconobbe che la geometria non euclidea bidimensionale era esattamente l'ordinaria geometria euclidea di certe superfici curve (sfere e pseudosfere).

Il passo analogo nella teoria degli insiemi sarebbe quello di negare l'assioma di scelta o l'ipotesi del continuo. Con ciò ovviamente intendiamo dire che il passo consisterebbe nel dimostrare che tale negazione è consistente con la teoria ristretta degli insiemi nello stesso senso in cui Gödel ha dimostrato che l'affermazione corrispondente è consistente. È appunto questa dimostrazione che è stata ottenuta nel 1963 e che ha dato origine a tutta una serie di ricerche in logica matematica, di cui non è possibile ancora prevedere quali

\forall PER TUTTI	\leftrightarrow SE E SOLO SE	\in È UN MEMBRO (ELEMENTO) DI
\exists ESISTE ALMENO UN	\vee O	$=$ UGUALE
$\exists!$ ESISTE ESATTAMENTE UN	$\&$ E	\neq DIVERSO
\cup UNIONE	\sim NON	ϕ INSIEME VUOTO
\rightarrow IMPLICA	\subseteq È UN SOTTOINSIEME DI	

1. ASSIOMA DI ESTENSIONALITÀ

$\forall x, y (\forall z (z \in x \rightarrow z \in y) \rightarrow x = y).$

Due insiemi sono uguali se e solo se hanno gli stessi elementi.

2. ASSIOMA DELL'INSIEME VUOTO

$\exists x \forall y (\sim y \in x).$

Esiste un insieme che non ha elementi (l'insieme vuoto).

3. ASSIOMA DELL'INSIEME COPPIA

$\forall x, y \exists z \forall w (w \in z \leftrightarrow w = x \vee w = y).$

Se x e y sono insiemi, allora la coppia (non ordinata) {x, y} è un insieme.

4. ASSIOMA DELL'INSIEME RIUNIONE

$\forall x \exists y \forall z (z \in y \leftrightarrow \exists t (z \in t \& t \in x)).$

Se x è un insieme di insiemi, la riunione di tutti i suoi membri è un insieme. (Ad esempio, se $x = \{\{a, b, c\}, \{a, c, d, e\}\}$ allora la riunione dei due elementi di x è l'insieme {a, b, c, d, e}.)

5. ASSIOMA DELL'INFINITO

$\exists x (\phi \in x \& \forall y (y \in x \rightarrow y \cup \{y\} \in x).$

Esiste un insieme x che contiene l'insieme vuoto ed è tale che se y appartiene a x, allora anche la riunione di y e {y} è in x. È fondamentale la distinzione fra l'elemento y e l'insieme {y} che contiene y come unico elemento. Questo assioma assicura l'esistenza di insiemi infiniti.

6.. ASSIOMA DI RIMPIAZZAMENTO

$\forall t_1, \dots, t_k (\forall x \exists! y A_n(x, y; t_1, \dots, t_k) \rightarrow \forall u \exists v B(u, v))$ dove $B(u, v) \equiv \forall r (r \in v \leftrightarrow \exists s (s \in u \& A_n(s, r; t_1, \dots, t_k)))$.

È difficile enunciare questo assioma in linguaggio comune. È stato numerato 6., invece di 6, perché in realtà è uno schema di assiomi, ossia un'intera famiglia di assiomi. Supponiamo di aver enumerato tutte le formule esprimibili nel nostro sistema e indichiamo con A_n l'n-esima di tali formule. Allora l'assioma di rimpiazzamento dice che, se per dati t₁, . . . , t_k, A_n (x, y; t_i) definisce univocamente y come funzione di x, diciamo y = φ(x) allora, per ogni u, il codominio di φ su u è un insieme. Parlando in termini alquanto approssimativi, ciò significa che ogni proprietà («ragionevole») che può essere enunciata dal linguaggio formale della teoria può essere usata per definire un insieme (l'insieme delle cose che godono della proprietà enunciata).

7. ASSIOMA DELL'INSIEME POTENZA

$\forall x \exists y \forall z (z \in y \leftrightarrow z \subseteq x).$

Questo assioma afferma che esiste, per ogni x, l'insieme y di tutti i sottoinsiemi di x. Malgrado y venga così definito mediante una proprietà, esso non è ottenibile dall'assioma di rimpiazzamento perché non è dato come codominio di una funzione. In effetti, la cardinalità di y risulta maggiore di quella di x, sicché questo assioma ci consente di costruire cardinali superiori.

8. ASSIOMA DI SCELTA

Se $\alpha \rightarrow A_\alpha \neq \phi$ è una funzione definita per tutti gli $a \in x$, allora esiste un'altra funzione f(a) per $a \in x$, e $f(a) \in A_\alpha$.

Questo è il ben noto assioma di scelta che ci permette di fare un numero infinito di «scelte» anche se non disponiamo di alcuna proprietà che ci permetterebbe di definire la funzione di scelta consentendoci così di usare al suo posto l'assioma 6..

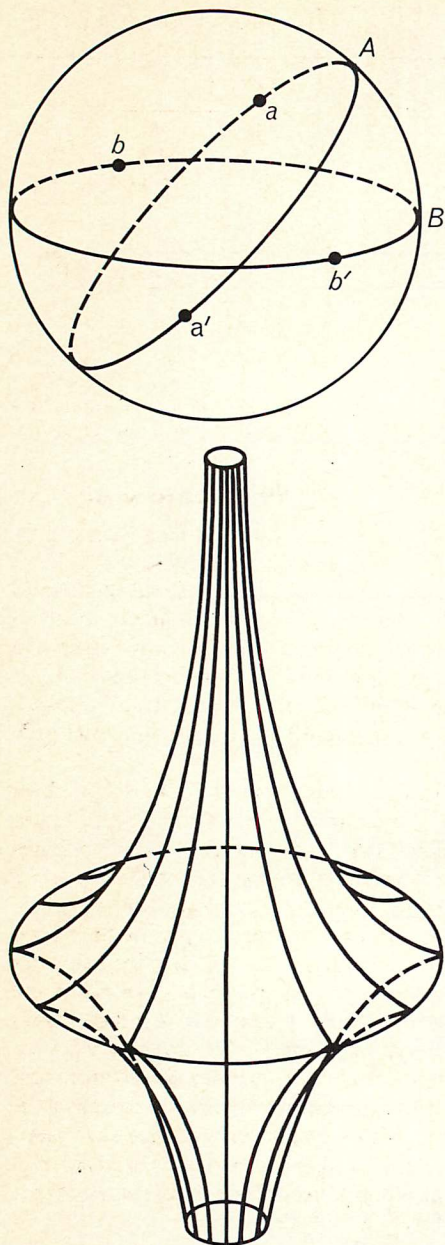
9. ASSIOMA DI FONDAZIONE

$\forall x \exists y (x = \phi \vee (y \in x \& \forall z (z \in x \rightarrow \sim z \in y)))$.

Questo assioma esclude esplicitamente, ad esempio, $x \in x$.

In questa tavola sono elencati gli assiomi proposti da Ernst Zermelo e Abraham Fraenkel per la teoria degli insiemi. Per

enunciare questi teoremi è necessario far uso dei simboli della teoria degli insiemi, un glossario dei quali è presentato in alto.



Sulla sfera, il termine «retta» viene interpretato come significante «circonferenza massima» (A e B in alto). Per ogni coppia di punti diametralmente opposti (aa' e bb') passano molte circonferenze massime. Se interpretiamo «punto» come significante «coppia di punti» allora i postulati di Euclide sono veri. Il secondo postulato è vero se si ammette che la «retta» estesa abbia lunghezza totale finita o che ripassi più volte su se stessa. Anche il terzo postulato è vero se conveniamo che la distanza vada misurata lungo circonferenze massime che possono essere ritracciate molte volte; qui «circonferenza» significa l'insieme di punti sulla sfera a una data distanza (misurata come sopra abbiamo convenuto) da un dato punto. Anche il quarto postulato è vero. Il postulato di Playfair è invece falso perché due qualunque circonferenze massime si intersecano. La sfera è quindi un modello di geometria non euclidea. Anche la pseudosfera (in basso) è un modello di geometria non euclidea ove si interpretino le rette come le più brevi curve che connettono due punti qualsiasi della superficie. Sulla superficie della pseudosfera esistono molte «rette» che passano per un dato punto e non intersecano una data retta.

possano rivelarsi gli sviluppi definitivi.

Poiché si tratta di dimostrare la consistenza relativa di un sistema di assiomi pensiamo in modo naturale alla costruzione di un modello. Come abbiamo visto, la consistenza relativa della geometria non euclidea venne stabilita quando si mostrò che certe superfici della geometria euclidea tridimensionale erano modelli della geometria non euclidea bidimensionale. In modo del tutto analogo, per provare la legittimità di una teoria non cantoriana degli insiemi, nella quale l'assioma di scelta o l'ipotesi del continuo sono falsi, dobbiamo usare gli assiomi della teoria ristretta degli insiemi per costruire un modello nel quale la negazione dell'assioma di scelta o la negazione dell'ipotesi del continuo possono essere dimostrate come teoremi.

Si deve ammettere che la costruzione di un tale modello è una questione complessa e delicata. Ma questa è una cosa che probabilmente ci si doveva aspettare. Nel caso degli insiemi costruibili di Gödel, ossia col suo modello per la teoria cantoriana degli insiemi, ci si proponeva di creare qualcosa che coincidesse sostanzialmente con la nostra nozione intuitiva di insieme ma che fosse più trattabile. Il nostro compito è invece quello di creare un modello di qualcosa di non intuitivo e strano, facendo tuttavia uso del familiare materiale da costruzione fornitoci dalla teoria ristretta degli insiemi.

Piuttosto che fermarci a questo punto affermando che è impossibile descrivere questo modello in un articolo non tecnico, tenteremo almeno di dare un resoconto descrittivo di alcune delle idee guida che questo compito comporta. Il nostro punto di partenza è la teoria ordinaria degli insiemi (senza l'assioma di scelta). Speriamo solo di dimostrare la consistenza della teoria non cantoriana degli insiemi in un senso relativo. Proprio come i modelli di geometria non euclidea provano che questa è consistente se lo è la geometria euclidea, così dimostreremo che se la teoria ristretta degli insiemi è consistente, essa rimane tale se aggiungiamo la proposizione «l'assioma di scelta è falso», o la proposizione «l'ipotesi del continuo è falsa». Supponiamo ora di disporre di un modello della teoria ristretta degli insiemi come punto di partenza. Chiamiamo M questo modello, che può essere riguardato come la classe degli insiemi costruibili di Gödel.

Sappiamo dal lavoro di Gödel che perché l'assioma di scelta o l'ipotesi del continuo divengano falsi dobbiamo aggiungere almeno un insieme non costruibile. Come ottenere ciò? Aggiun-

giamo la lettera a come nome di un oggetto che vada aggiunto a M ; resta da determinare che tipo di oggetto deve essere questo a . Una volta aggiunto a dobbiamo anche aggiungere ogni altro oggetto che può essere formato a partire da a applicando le operazioni ammesse nella teoria ristretta degli insiemi: riunire due o più insiemi per formare un nuovo insieme, formare l'insieme potenza, e così via. La nuova collezione di insiemi così generata da $M + a$ verrà detta N . Il problema è quello di come scegliere a in modo che 1) N sia un modello per la teoria ristretta degli insiemi, come lo era per ipotesi M , e 2) a non sia costruibile in N . Solo se ciò è possibile ci resta una speranza di riuscire a negare l'assioma di scelta oppure l'ipotesi del continuo.

Possiamo dare una vaga idea di come bisogna procedere, chiedendoci come avrebbe proceduto un geometra del 1850 per scoprire la pseudosfera. In senso molto approssimativo, è come se egli fosse partito con una curva M nel piano euclideo, avesse immaginato un punto a non giacente in quel piano e quindi avesse congiunto il punto a a tutti i punti di M . Poiché a è stato scelto non giacente nel piano di M , la superficie N che così ne risulta non coinciderebbe sicuramente col piano euclideo. È quindi ragionevole pensare che con sufficiente ingegno e impegno tecnico si potrebbe mostrare che essa è veramente un modello per una geometria non euclidea.

Il corrispondente procedimento nel caso della teoria non cantoriana degli insiemi consiste nello scegliere il nuovo insieme a come un insieme non costruibile, e quindi nel generare un nuovo modello N costituito da tutti gli insiemi ottenuti mediante le operazioni ammesse nella teoria ristretta degli insiemi e applicate ad a e agli insiemi di M . Se questo è possibile, si sarà allora dimostrato che si può impunemente negare l'assioma di costruibilità. Poiché Gödel aveva dimostrato che la costruibilità implica l'assioma di scelta e l'ipotesi del continuo, questo è il primo passo necessario per negare una delle due proposizioni.

Per poter effettuare questo primo passo si devono far vedere due cose: che a può essere scelto in modo da rimanere non costruibile non solo in M ma anche in N e che N , come M , è un modello per la teoria ristretta degli insiemi. Per caratterizzare a , aggiriamo la questione. Immaginiamo di voler fare un elenco di tutte le possibili proposizioni su a in quanto insieme appartenente a N . Allora a sarà specificato se noi diamo una regola in base alla quale possiamo determinare se ognuna di que-

GEOMETRIA	STADIO DI SVILUPPO	TEORIA DEGLI INSIEMI
TALETE, PITAGORA	BASI INTUITIVE PER I PRIMI TEOREMI	CANTOR
ZENONE	SCOPERTA DI PARADOSSI	RUSSELL
EUDOSSO, EUCLIDE	BASI ASSIOMATICHE PER LA TEORIA STANDARD	ZERMELO, FRAENKEL
CARTESIO, HILBERT	LE TEORIE STANDARD VENGONO DIMOSTRATE (RELATIVAMENTE) CONSISTENTI	GÖDEL
GAUSS, RIEMANN	SCOPERTA DI TEORIE NON STANDARD	RICERCA ATTUALE
MINKOWSKI, EINSTEIN	APPLICAZIONI DELLA TEORIA NON STANDARD	? ? ?

Analogia fra l'evoluzione della geometria (a sinistra) e della teoria degli insiemi (a destra). La geometria non standard (non

euclidea) è stata applicata nella relatività; la teoria non standard degli insiemi non ha ancora trovato applicazioni in fisica.

ste proposizioni è vera oppure no.

L'idea centrale risulta essere quella di scegliere a in modo che sia un elemento « generico », ossia di sceglierlo in modo che siano vere per a solo quelle proposizioni che sono vere per quasi tutti gli insiemi in M . Questa è una nozione paradossale. Ogni insieme in M , infatti, gode tanto di proprietà peculiari che lo identificano, quanto anche di proprietà tipiche generali, che egli condivide praticamente con tutti gli insiemi di M . Ora risulta possibile determinare in modo rigoroso questa distinzione fra proprietà specifiche e generiche sì da renderla perfettamente esplicita e formale. Allora, quando scegliamo a come insieme generico (ossia come un insieme, per così dire, che non ha particolari proprietà che lo distinguano da ogni altro insieme in M) ne segue che N è ancora un modello per la teoria ristretta degli insiemi. Il nuovo elemento a che abbiamo introdotto non ha fastidiose proprietà che possano alterare il modello M da cui siamo partiti. Nel contempo a non è costruibile: ogni insieme costruibile ha un suo proprio carattere particolare – i passi mediante i quali può essere costruito – e il nostro a manca precisamente di tale individualità.

Per costruire un modello nel quale è falsa l'ipotesi del continuo, dobbiamo aggiungere a M non un solo elemento a ma una gran quantità di nuovi elementi, in effetti un numero infinito. Possiamo far questo in modo tale che gli elementi che noi aggiungiamo abbiano cardinalità

$$\aleph_2 = 2^{(2^{\aleph_0})}$$

nel modello M . Anche qui può essere utile un'approssimativa analogia geometrica: a un essere bidimensionale che vivesse immerso in una superficie non euclidea sarebbe impossibile riconoscere che il suo mondo fa parte di uno spazio euclideo tridimensionale.

Nel nostro caso noi, stando al di fuori di M , possiamo vedere che abbiamo immesso solo un'infinità numerabile di nuovi elementi. Essi tuttavia sono tali che non possono essere numerati con nessun metodo disponibile in M stesso. Otteniamo così un nuovo modello N' nel quale l'ipotesi del continuo è falsa. I nuovi elementi, che in N' svolgono il ruolo di numeri reali (ossia, punti di un segmento hanno cardinalità maggiore di 2^{\aleph_0} e così ora esiste un cardinale infinito – precisamente 2^{\aleph_0} – che è maggiore di \aleph_0 e tuttavia è minore di \mathfrak{c} poiché nel nostro modello N' \mathfrak{c} è uguale a

$$2^{(2^{\aleph_0})}$$

Dal momento che possiamo costruire un modello della teoria degli insiemi nel quale l'ipotesi del continuo è falsa, possiamo aggiungere alla nostra ordinaria teoria degli insiemi ristretta l'assunzione della falsità dell'ipotesi del continuo; non può sorgere alcuna contraddizione che non fosse già presente. Nello stesso ordine di idee possiamo costruire modelli per la teoria degli insiemi nei quali è falso l'assioma di scelta; possiamo anche essere più specifici nel precisare da quali insiemi infiniti è possibile « eseguire la scelta » e quali invece sono « troppo grandi perché si possa eseguire la scelta ».

Mentre Gödel otteneva i suoi risultati servendosi di un unico modello (gli insiemi costruibili), nella teoria non cantoriana degli insiemi abbiamo non uno ma molti modelli, ognuno dei quali è opportunamente costruito per scopi particolari. Ma forse più importante dei vari modelli è la tecnica che ci consente di costruirli: la nozione di « generico » e quella ad essa collegata di « costrizione » (*forcing*). In senso molto approssimativo gli insiemi generici hanno solo quelle proprietà che essi sono « costretti » (*forced*) ad avere per poter essere considerati insiemi. Per decidere

se a è « costretto » ad avere certe proprietà dobbiamo riferirci a tutto N . E tuttavia N non resta realmente definito finché non abbiamo specificato a ! Altro elemento chiave nella nuova teoria è appunto dato dall'elaborazione di un metodo per rendere non circolare questo argomento che apparentemente lo è.

Cosa può suggerirci la storia della geometria per il futuro della teoria degli insiemi? Uno dei successi più notevoli della geometria non euclidea è stato quello di risultare un requisito preliminare essenziale per la teoria generale della relatività di Einstein. Riemann creò la geometria riemanniana con lo scopo puramente astratto di unificare, chiarificare e approfondire la geometria non euclidea di Lobacevskij, Bolyai e Gauss. Questa geometria risultò poi essere uno strumento essenziale e indispensabile per la rivoluzionaria interpretazione einsteiniana della forza gravitazionale.

È sufficiente questo esempio a giustificare la fiducia che la teoria non cantoriana degli insiemi troverà un giorno un'applicazione oggi non prevedibile al mondo « reale » (ossia non matematico)? Nessuno, oggi, si impegnerebbe in una risposta. Certamente possiamo vedere (a posteriori) che la geometria ha sempre fornito lo sfondo essenziale nel quale si svolgono gli eventi fisici. In questo senso avremmo forse dovuto aspettarci che progressi fondamentali in geometria avrebbero trovato un'applicazione fisica. La teoria degli insiemi non sembra oggi avere un'analogia relazione organica con la fisica, pure vi sono stati alcuni matematici (Stanislaw Ulam, per esempio) che hanno suggerito che la teoria astratta degli insiemi potrebbe fornire utili modelli per la fisica teorica. Ma a questo punto dello sviluppo la cosa migliore è quella di rifiutarsi di fare ogni previsione attorno al futuro, se non quella che esso è imprevedibile.

L'analisi non-standard

Questa teoria matematica ha riportato in primo piano gli infinitesimi adoperati fin dalla antichità, ma spesso con dubbi, per risolvere problemi come quello di trovare l'area del cerchio

di Martin Davis e Reuben Hersh

L'analisi non-standard, un nuovo ramo della matematica scoperto 10 anni fa dal logico Abraham Robinson, segna un nuovo stadio di sviluppo per la soluzione di un gran numero di antichi e celebri paradossi. Robinson, ora all'Università di Yale, ha riesumato la nozione di «infinitesimo», cioè di un numero che è infinitamente piccolo e tuttavia maggiore di zero. Questo concetto ha radici che affondano nella antichità. Per l'analisi tradizionale, o analisi «standard», appariva palesemente auto-contraddittorio. Eppure era stato un importante strumento in meccanica e in geometria almeno fin dai tempi di Archimede.

Nel secolo XIX gli infinitesimi vennero eliminati dalla matematica una volta per tutte, o almeno così sembrava. Per soddisfare le esigenze della logica, il calcolo infinitesimale (o, come diremo semplicemente, il Calcolo) di Isaac Newton o di Gottfried Wilhelm Leibniz fu riformulato da Karl Weierstrass senza infinitesimi. Ma oggi è proprio la logica matematica, con la sua sottigliezza e le sue attuali capacità, che ha riportato in vita l'infinitesimo e l'ha reso di nuovo accettabile. In un certo senso Robinson ha reso legittimo l'incauto trasporto dei matematici del XVII secolo contro il rigore pieno di scrupoli dei matematici del XVIII secolo, aggiungendo un nuovo capitolo alla interminabile contesa tra finito e infinito, continuo e discontinuo.

Nelle controversie sul concetto di infinitesimo che accompagnarono lo sviluppo del Calcolo, la geometria di Euclide era lo standard rispetto al quale i moderni si misuravano. In Euclide sono deliberatamente esclusi sia l'infinito sia l'infinitesimo. Leggiamo in Euclide che un punto è ciò che ha posizione, ma non grandezza. Si è detto che tale definizione è priva di significato, ma

forse è solo un impegno ad astenersi dall'uso di ragionamenti infinitesimali; e ciò costituiva il rifiuto di precedenti concetti del pensiero greco. Si pensava che l'atomismo di Democrito si riferisse non solo alla materia, ma anche al tempo e allo spazio. Ma allora gli argomenti di Zenone avevano reso insostenibile la nozione di tempo come sequenza di istanti successivi, o di linea come sequenza di successivi «indivisibili». Aristotele, fondatore della logica formale, bandiva dalla geometria l'infinitamente grande o piccolo.

Ecco un esempio tipico dell'uso di ragionamenti infinitesimali in geometria:

«Vogliamo trovare la relazione tra l'area del cerchio e la sua circonferenza. Per semplicità, supponiamo che il raggio del cerchio sia 1. Ora, la circonferenza può essere pensata come composta di un numero infinito di segmenti rettilinei, tutti eguali tra loro e infinitamente corti. Il cerchio è allora la somma di triangoli infinitesimi, che hanno tutti altezza 1. Per un triangolo l'area è data dalla metà del prodotto della base per l'altezza. Dunque la somma delle aree dei triangoli è l'area del cerchio, mentre la somma delle basi dei triangoli è la sua circonferenza. Dunque l'area del cerchio di raggio 1 è eguale alla metà della sua circonferenza».

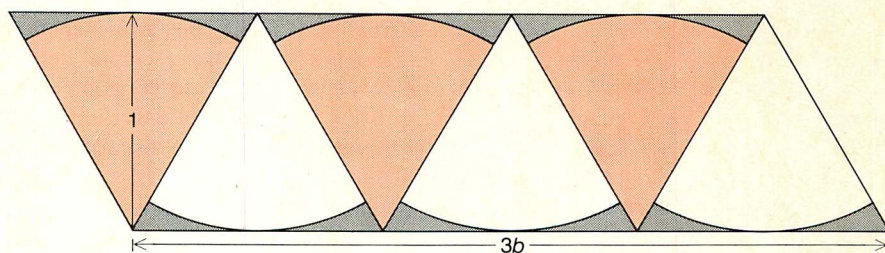
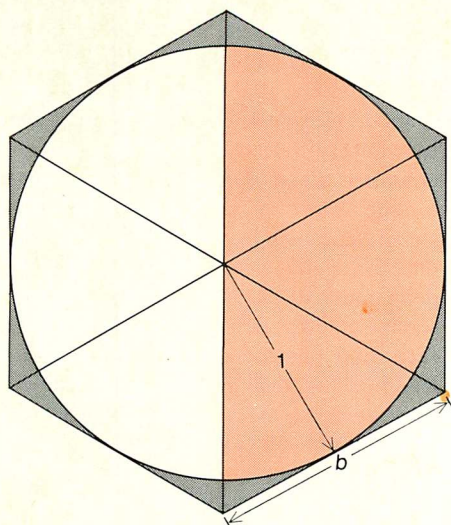
Questo ragionamento, che Euclide

avrebbe respinto, fu pubblicato nel XV secolo da Nicola Cusano. La sua conclusione è ovviamente vera, ma non è difficile sollevare obiezioni a questo modo di procedere. La nozione di triangolo con base infinitamente piccola è perlomeno ambigua. Certamente la base di un triangolo deve avere lunghezza o zero o maggiore di zero. Se è zero, allora l'area è zero e, per quanti termini sommiamo, non otteniamo altro che zero. D'altra parte, se è maggiore di zero, per quanto piccola, se sommiamo un numero infinito di termini eguali, otteniamo una somma infinitamente grande. In nessuno dei due casi possiamo ottenere un cerchio di circonferenza finita come somma di un numero infinito di parti identiche.

L'essenza di questa confutazione è l'asserzione che un numero diverso da zero anche molto piccolo diventa arbitrariamente grande se viene aggiunto a se stesso un numero di volte sufficiente. Poiché tale asserzione fu resa esplicita per la prima volta da Archimede, è detta proprietà archimedeo dei numeri reali. Un infinitesimo, se esistesse, sarebbe appunto un numero non-archimedeo: un numero maggiore di zero, che nondimeno resterebbe per esempio minore di 1, qualunque numero (finito) di volte fosse sommato a se stesso. Archimede, che operava seguendo la tradizione di Aristotele e di Euclide, asserì che ogni numero è ar-

Il metodo di esaurimento viene usato per dimostrare in modo indiretto che l'area di un cerchio di raggio 1 è eguale alla metà della sua circonferenza. Nella dimostrazione della pagina a fronte, un poligono *A* è circoscritto al cerchio (*prima figura in alto*), formando un certo numero di triangoli le cui aree possono essere immediatamente calcolate. Aumentando il numero dei lati del poligono, come nel poligono *B* e nel poligono *N*, i triangoli crescono di numero e diventano più sottili mentre la differenza tra l'area del cerchio e quella del poligono diventa più piccola. Tale differenza non sarà mai nulla, tuttavia, per un poligono con un numero finito qualunque di lati. L'analisi standard supera questa difficoltà asserendo che, al crescere all'infinito del numero dei lati, l'area del poligono si approssima all'area del cerchio come a un limite. L'analisi non-standard fa a meno del concetto di limite per una spiegazione più intuitiva facendo uso di un poligono con un numero infinito di lati, di lunghezza infinitesimale.

POLIGONO A



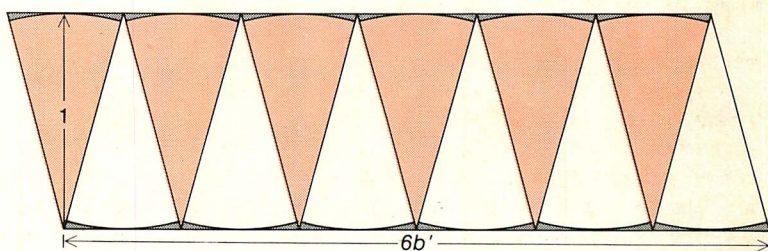
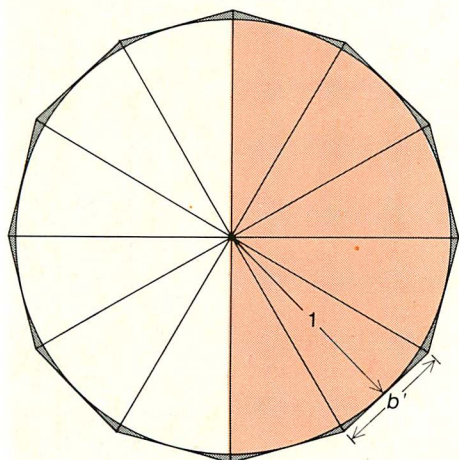
$$\text{PERIMETRO } P = 6b$$

$$\text{AREA DEL POLIGONO A} = 3b = \frac{1}{2}P$$

$$\text{CIRCONFERENZA DEL CERCHIO INSCRITTO} \ll P$$

$$\text{AREA DEL CERCHIO INSCRITTO} \ll \frac{1}{2}P$$

POLIGONO B



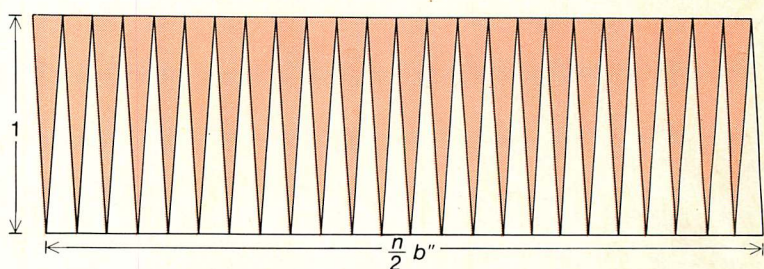
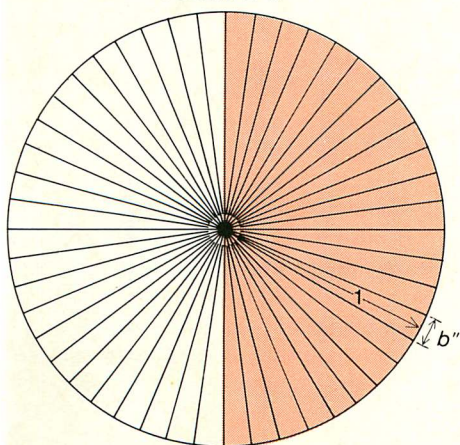
$$\text{PERIMETRO } P' = 12b'$$

$$\text{AREA DEL POLIGONO B} = 6b' = \frac{1}{2}P'$$

$$\text{CIRCONFERENZA DEL CERCHIO INSCRITTO} < P'$$

$$\text{AREA DEL CERCHIO INSCRITTO} < \frac{1}{2}P'$$

POLIGONO N



$$\text{PERIMETRO } P'' = nb''$$

$$\text{AREA DEL POLIGONO N} = \frac{nb''}{2} = \frac{1}{2}P''$$

$$\text{CIRCONFERENZA DEL CERCHIO INSCRITTO} \approx P''$$

$$\text{AREA DEL CERCHIO INSCRITTO} \approx \frac{1}{2}P''$$

chimedee; non ci sono dunque infinitesimi. Archimede, tuttavia, era pure un filosofo della natura, un ingegnere e un fisico. Faceva uso di infinitesimi e della sua intuizione fisica per risolvere problemi della geometria della parabola. Allora, visto che gli infinitesimi « non esistono », diede una dimostrazione « rigorosa » dei suoi risultati, servendosi « del metodo di esaustione », che si basa su un ragionamento diretto e su costruzioni puramente finite. La dimostrazione rigorosa è data nel suo trattato *Sulla quadratura della parabola*, noto fin dall'antichità. L'uso degli infinitesimi, che in realtà serviva a trovare la soluzione del problema, è contenuto in uno scritto chiamato *Sul Metodo*, che rimase sconosciuto sino alla sua sensazionale scoperta nel 1906.

Il metodo di esaustione di Archimede, che evita il ricorso agli infinitesimi, è concettualmente vicino al me-

todo « epsilon-delta » con cui Weierstrass e i suoi allievi nel XIX secolo bandirono i metodi infinitesimali dall'Analisi. È facile spiegarlo facendo riferimento al nostro esempio del cerchio pensato come un poligono di infiniti lati. Vogliamo ottenere una dimostrazione logicamente accettabile dell'enunciato « L'area del cerchio con raggio 1 è eguale alla metà della circonferenza », che abbiamo trovato servendoci di un ragionamento logicamente inaccettabile.

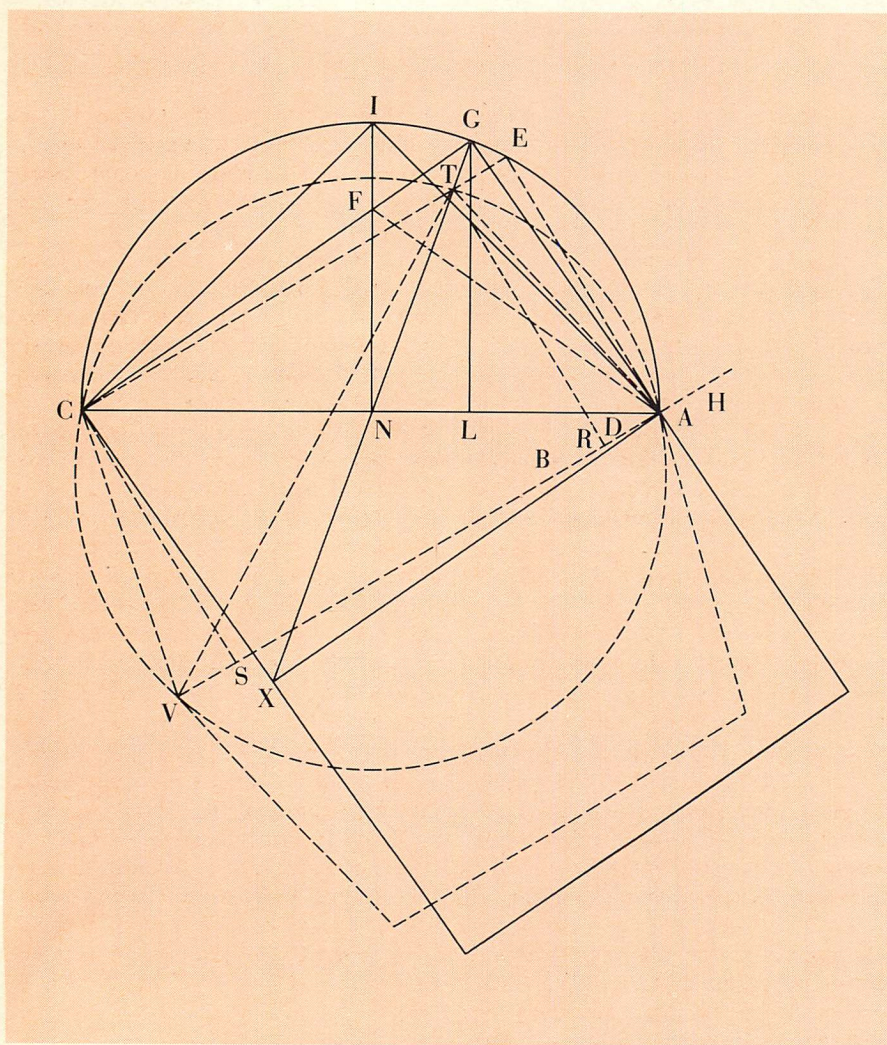
Ragioniamo come segue. L'enunciato esprime l'eguaglianza di due quantità associate a un cerchio di raggio 1: l'area e la metà della circonferenza. Così, se l'enunciato fosse falso, una di queste quantità sarebbe maggiore dell'altra. Sia A il numero positivo ottenuto sottraendo la minore dalla maggiore. Ora, possiamo circoscrivere alla circonferenza un poligono regolare di tanti lati quanti si voglia. Dal momen-

to che il poligono è composto di un numero finito di triangoli finiti di altezza 1, sappiamo che la sua area è eguale alla metà del suo perimetro. Rendendo il numero dei lati sufficientemente grande, possiamo far sì che l'area del poligono differisca dall'area del cerchio per meno della metà di A (qualunque valore si sia preso per A); allo stesso tempo il perimetro del poligono differirà dalla circonferenza per meno della metà di A . Ma allora l'area e il semiperimetro della circonferenza differiranno per meno di A , il che contraddice l'ipotesi di partenza. Dunque tale ipotesi è assurda e A deve essere zero, come volevasi dimostrare.

Questo ragionamento da un punto di vista logico è impeccabile. Paragonato con il procedimento diretto con cui si è precedentemente analizzata la questione, esso presenta un che di micidioso, anzi di pedante. Dopo tutto, se l'uso di infinitesimi dà la risposta giusta, non sarà corretto per certi aspetti anche il ragionamento? Anche se non siamo in grado di giustificare i concetti che esso adopera, come può essere di fatto errato se ci fa giungere a un risultato corretto?

Tale difesa degli infinitesimi non fu fatta da Archimede. In realtà nello scritto *Sul Metodo* egli ha cura di spiegare che « le verità qui enunciate non vengono di fatto dimostrate con il ragionamento usato » e che una dimostrazione rigorosa è stata pubblicata a parte. D'altra parte, Nicola Cusano, che era cardinale, preferiva il ragionamento per quantità infinite per via della sua convinzione che l'infinito costituisse « la fonte, e i mezzi, e nel medesimo tempo lo scopo irraggiungibile di tutta la conoscenza ». Cusano fu seguito nel suo misticismo da Giovanni Keplero, uno dei fondatori della scienza moderna. In un'opera oggi meno nota delle sue scoperte nel campo della astronomia, Keplero nel 1612 usava gli infinitesimi per trovare le proporzioni ottimali per una botte di vino. Egli non era toccato dalle autocontraddizioni del suo metodo; faceva assegnamento sulla ispirazione divina, e scriveva che « la natura insegna la geometria solo per istinto, anche senza raziocinio ». Per di più, le sue formule per il volume della botte sono corrette.

Il più famoso mistico matematico fu senza dubbio Blaise Pascal. Rispondendo a quei contemporanei che sollevavano obiezioni contro il ragionamento con quantità infinitamente piccole, Pascal proclamava con entusiasmo che il cuore interviene a rendere l'opera chiara. Pascal guardava all'infinitamen-



Il problema della botte fu affrontato da Keplero servendosi di infinitesimi nella sua *Nova stereometria doliorum vinariorum*, pubblicata nel 1615. Il problema che Keplero si era posto consisteva nel trovare le dimensioni ottimali di una botte di vino. In figura è riprodotta una pagina tratta da una ristampa ottocentesca dell'opera edita in Europa.

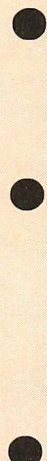
te grande e all'infinitamente piccolo come a misteri, come a qualcosa che la natura ha posto innanzi all'uomo non perché l'uomo la comprenda, ma perché l'uomo la ammiri.

La piena fioritura del metodo degli infinitesimi si ebbe con le generazioni dopo Pascal: Newton, Leibniz, i fratelli Bernoulli (Jacopo e Giovanni) e Leonardo Eulero. I teoremi fondamentali del Calcolo furono trovati da Newton e Leibniz tra il 1660 e il 1680. Il primo libro di testo sul Calcolo fu scritto nel 1696 dal marchese de l'Hôpital, un discepolo di Leibniz e di Giovanni Bernoulli. Qui si enuncia all'esordio, come assioma, il fatto che due quantità che differiscono per un infinitesimo possono essere considerate eguali. In altri termini, nel medesimo tempo tali quantità si considerano eguali e non eguali tra loro! Un secondo assioma dice che una curva è « la totalità costituita da una infinità di segmenti retti, ciascuno infinitamente piccolo ». Ciò costituisce una aperta ripresa dei metodi che Aristotele aveva bandito 2000 anni prima.

In effetti, scriveva de l'Hôpital, « la usuale analisi ha a che fare solo con quantità finite; questa invece si addentra tanto profondamente quanto la stessa infinità. Essa paragona le differenze infinitamente piccole di quantità finite; scopre le relazioni tra queste differenze e in questo modo rende note le relazioni che intercorrono tra le quantità finite che si comportano come fossero infinite in confronto alle quantità infinitamente piccole. Si potrebbe anche dire che questa analisi si estende oltre l'infinito, perché non si limita alle differenze infinitamente piccole ma scopre anche le relazioni tra le differenze di queste differenze ».

Newton e Leibniz non partecipavano all'entusiasmo di de l'Hôpital. Leibniz non proclamava che gli infinitesimi esistevano realmente, ma solo che si può ragionare come se esistessero senza commettere errore. Sebbene Leibniz non fosse in grado di dare fondamento a tale affermazione, l'opera di Robinson mostra che in un certo qual senso dopo tutto aveva ragione lui. Newton cercava invece di eliminare gli infinitesimi. Nei suoi *Principia Mathematica*, come nell'opera *Sulla quadratura della parabola* di Archimede, risultati che erano stati in origine trovati con metodi infinitesimali sono presentati in una veste euclidea puramente finita.

La dinamica era divenuta importante quanto la geometria nel porre problemi all'analisi matematica. Il proble-

WEIERSTRASS	POSIZIONE DELLA PIETRA CHE CADE $s = 4,9t^2$	ROBINSON
Sia $t' = 1 + \Delta t$. Δt è un numero reale positivo. $s' = 4,9 + 9,8 \Delta t + 4,9 (\Delta t)^2$ $\Delta s = s' - s$ $= 9,8 \Delta t + 4,9 (\Delta t)^2$ $\frac{\Delta s}{\Delta t} = 9,8 + 4,9 \Delta t$ Assegnato un qualsiasi numero reale positivo ε , arbitrariamente piccolo, scegliamo $\delta = \frac{\varepsilon}{4,9}$ Allora per tutti i $\Delta t < \delta$, $\frac{\Delta s}{\Delta t} - 9,8 = 4,9 \Delta t < 4,9 \delta$ $= 4,9 \cdot \frac{\varepsilon}{4,9} = \varepsilon$ Così: velocità istantanea = $\lim_{\Delta t \rightarrow 0} \frac{\Delta s}{\Delta t} = 9,8$		Sia $t' = 1 + dt$. dt è un numero positivo infinitesimo. $s' = 4,9 + 9,8dt + 4,9 (dt)^2$ $ds = s' - s$ $= 9,8dt + 4,9 (dt)^2$ $\frac{ds}{dt} = 9,8 + 4,9dt$ Poiché dt è infinitesimo, lo è anche $4,9dt$. $9,8$ è un numero reale standard. Così: velocità istantanea = parte standard di $\frac{ds}{dt} = 9,8$

Soluzione del problema della caduta di una pietra secondo l'analisi standard (a sinistra) e secondo l'analisi non-standard (a destra). L'analisi standard, semplificata dal matematico tedesco dell'800, Karl Weierstrass, calcola la velocità della pietra che cade in ogni istante senza servirsi di infinitesimi, definendo invece la velocità istantanea come un limite che è approssimato da rapporti di incrementi finiti. Abraham Robinson della Università di Yale, creatore dell'analisi non-standard, compie questo calcolo servendosi invece di una versione modificata del metodo infinitesimale.

ma principale era quello di porre in relazione « fluenti » e « flussioni », che oggi si chiamerebbero posizioni istantanee e velocità istantanee di un corpo in movimento.

Consideriamo una pietra che cade. Il suo movimento è descritto assegnandone la posizione come funzione del tempo. Mentre cade, la sua velocità cresce, sicché la velocità in ogni istante è essa pure una funzione variabile del tempo. Newton chiamava « fluente » la funzione posizione e « flussione » la funzione velocità. Se una delle due è data, si può determinare anche l'altra; questa relazione è il nucleo del calcolo infinitesimale costruito da Newton e da Leibniz.

Nel caso della caduta di una pietra la funzione fluente è data dalla formula $s = 4,9t^2$, dove s è il numero dei metri percorsi e t è il numero di secondi trascorsi dal momento in cui si è fatta cadere la pietra. Al cadere della pietra aumenta rapidamente la sua velocità. Come possiamo calcolare la velocità della pietra in un certo istante di tempo, per esempio all'istante $t = 1$?

Potremmo trovare la velocità *media* per un intervallo di tempo finito servendoci di una formula elementare: la velocità è eguale alla distanza divisa per il tempo. Possiamo usare questa formula per trovare la velocità istantanea? In un incremento infinitesimo del tempo l'incremento della distanza

sarebbe esso pure infinitesimo; il loro rapporto, la velocità media per tale intervallo infinitesimo, dovrebbe essere la velocità istantanea finita che cerchiamo.

Rappresentiamo con dt l'incremento infinitesimo del tempo e con ds il corrispondente incremento della distanza (ovviamente ds e dt debbono essere pensati ciascuno come un solo simbolo e non come d volte t o d volte s). Vogliamo calcolare il rapporto ds/dt che deve essere finito. Per trovare l'incremento della distanza da $t = 1$ a $t = 1 + dt$ calcoliamo la posizione della pietra quando $t = 1$, che è $4,9 \times 1^2 = 4,9$, e la sua posizione quando $t = 1 + dt$, che è $4,9 \times (1 + dt)^2$. Basta un po' di algebra elementare per trovare che ds , cioè l'incremento della distanza, è $9,8dt + 4,9dt^2$. Così il rapporto ds/dt , che è la quantità che stiamo cercando, è eguale a $9,8 + 4,9dt$.

Abbiamo risolto il nostro problema? Dal momento che il risultato cercato dovrebbe essere una quantità finita, si dovrebbe eliminare il termine infinitesimo, $4,9dt$, e ottenere come risultato $9,8$ metri al secondo per la velocità istantanea. Ciò è proprio quel che il vescovo Berkeley non ci permetterà di fare.

La brillante e spietata critica al metodo infinitesimale di Berkeley, *The Analyst*, apparve nel 1734. Il libro era rivolto a « un matematico sen-

za fede », che generalmente si ritiene fosse l'astronomo Edmund Halley, amico di Newton. Halley aveva finanziato la pubblicazione dei *Principia* e aveva contribuito alla loro preparazione per la stampa. Si dice pure che egli avesse persuaso un amico di Berkeley della « inconcepibilità delle dottrine del cristianesimo »; il vescovo Berkeley rispose che le flussioni di Newton erano « oscure, ripugnanti e precarie » come nessun punto della teologia.

« Chiederò per me il privilegio del Libero Pensatore – scrisse Berkeley – e mi prenderò la libertà di ricercare sull'oggetto, sui principi, e sul metodo di dimostrazione ammessi dai matematici del tempo presente, con la stessa disinvoltura con cui voi presumete di trattare i principi e i misteri della religione ». Berkeley affermò che il procedimento di Leibniz, che consisteva semplicemente nel « considerare » $9,8 + 4,9dt$ come « lo stesso » che $9,8$ era inintelligibile. « E non servirà neppure – scrisse – dire che [il termine trascurato] è una quantità estremamente piccola; poiché ci hanno detto che *in rebus mathematicis errores quam minimi non sunt contemnendi* ». Se qualcosa è trascurato, per quanto piccolo, non possiamo più affermare di avere la velocità esatta, ma solo una approssimazione.

Newton, a differenza di Leibniz, cercò nei suoi scritti più tardi di addolcire la « durezza » della dottrina degli infinitesimi servendosi di un linguaggio ricco di riferimenti alla fisica. « Con velocità ultima si intende quella con cui il corpo è mosso, né prima che arrivi all'ultima posizione, quando il movimento cessa, né dopo, ma quella all'istante stesso in cui vi arriva ... E, parimenti, come ultima ragione di quantità evanescenti si deve intendere il rapporto di quantità, né prima che esse scompaiono, né dopo, ma quel rap-

porto con cui esse scompaiono ». Quando procedeva al calcolo, tuttavia, egli doveva giustificare l'eliminazione dei termini « trascurabili » non voluti che comparivano nel risultato trovato. Il procedimento newtoniano consisteva nel trovare per prima cosa, come abbiamo fatto anche noi, $ds/dt = 9,8 + 4,9dt$ e quindi nel porre l'incremento dt eguale a zero, lasciando come risultato finale esatto solo $9,8$.

Ma, scrisse Berkeley, « questo ragionamento non sembra né corretto né conclusivo ». Dopo tutto, dt o è eguale a zero oppure non lo è. Se dt è diverso da zero, allora $9,8 + 4,9dt$ è diverso da $9,8$. Se dt è zero, allora l'incremento della distanza ds è esso pure zero, e la frazione ds/dt non è $9,8 + 4,9dt$, ma un'espressione priva di significato, $0/0$. « Perché, una volta ammesso che gli incrementi scompaiono, cioè che gli incrementi siano nulla o che non vi siano incrementi, cade la precedente ipotesi che gli incrementi fossero qualcosa, o che vi fossero incrementi, mentre viene mantenuta una conseguenza di tale ipotesi, cioè un'espressione ottenuta mediante essa. » Ma questo è un modo fallace di ragionare. Quasi ironicamente Berkeley concludeva: « Che sono queste flussioni? La velocità di incrementi evanescenti. E che sono questi stessi incrementi evanescenti? Non sono né quantità finite, né quantità infinitamente piccole, e neppure nulla. Perché non chiamarle fantasmi di quantità svanite? ».

Non si poté dare risposta alla logica di Berkeley; nondimeno, i matematici persistettero nell'uso degli infinitesimi per un altro secolo, e con grande successo. E in realtà fisici e ingegneri non hanno mai cessato di usarli. Nella matematica pura, d'altra parte, si ebbe nel XIX secolo un ritorno al rigore eu-

clideo, che raggiunse il culmine sotto la guida di Weierstrass nel 1872. È interessante osservare che il XVIII secolo – che fu il grande momento degli infinitesimi – fu l'epoca in cui non si riconoscevano barriere tra matematica e fisica. I fisici di punta e i matematici di punta erano le stesse persone. Quando la matematica pura ricomparve come disciplina a sé stante, i matematici di nuovo si assicurarono che i fondamenti della loro attività non contenessero ovvie contraddizioni. L'analisi moderna rese saldi i suoi fondamenti facendo quello che già i greci avevano fatto: bandendo gli infinitesimi.

Per trovare la velocità istantanea seguendo il metodo di Weierstrass, abbandoniamo ogni tentativo di calcolare la velocità come un rapporto. Definiamo invece la velocità come un limite, che è approssimato da rapporti di incrementi finiti. Sia Δt una variabile che rappresenta un incremento finito del tempo e Δs la corrispondente variabile per l'incremento dello spazio. Allora $\Delta s/\Delta t$ è la quantità variabile $9,8 + 4,9\Delta t$. Scegliendo Δt sufficientemente piccolo possiamo far sì che $\Delta s/\Delta t$ assuma valori tanto prossimi quanto si voglia al valore $9,8$, e così, per definizione, la velocità all'istante $t = 1$ è proprio $9,8$.

Questo procedimento riesce a eliminare ogni riferimento a numeri non finiti. Evita anche ogni tentativo di porre immediatamente Δt eguale a zero nella frazione $\Delta s/\Delta t$. In questo modo evitiamo entrambi i trabocchetti logici presentati da Berkeley. Ne paghiamo, tuttavia, il prezzo. La velocità istantanea, quantità intuitivamente perspicua e fisicamente misurabile, diventa soggetta alla nozione di limite che è estremamente raffinata. Se analizziamo nei particolari cosa significa, abbiamo il seguente scioglilingua:

« La velocità è v se, per ogni numero positivo ϵ , $\Delta s/\Delta t - v$ è minore di ϵ in valore assoluto per tutti i valori di Δt minori in valore assoluto di un certo numero positivo δ (che dipenderà da ϵ e da t) ».

Abbiamo definito v con una sottile relazione tra due nuove quantità, ϵ e δ , che, in certo senso, sono irrilevanti rispetto a v stesso. Perlomeno la non conoscenza di ϵ e δ non impedì mai a Bernoulli o a Eulero di calcolare una velocità. La verità è che in realtà conoscevamo cosa era la velocità istantanea prima ancora di avere appreso questa definizione; infatti solo per coerenza logica accettiamo una definizione che è molto più difficile da capire del concetto definito. Naturalmente, per un matematico allenato, la defini-

SIMBOLO	SIGNIFICATO CONVENUTO
~	non
&	e
V	o
→	implica
∀	per tutti
∃	esiste almeno un
=	eguale
x y z	variabili che variano su numeri reali
f g h	variabili che variano su altri oggetti
+ · <	più, per, minore di
() []	parentesi
0 1 2	simboli per particolari numeri

I simboli usati nel linguaggio formale L , in cui si può esprimere il calcolo, sono tradotti in italiano in questo dizionario parziale. Il linguaggio formale, che impiega molti più simboli di questi, fornisce un legame tra l'universo standard e il più ampio universo non-standard che è un concetto centrale della analisi matematica non-standard.

zione epsilon-delta è intuitiva; questo mostra che cosa si può ottenere con un appropriato allenamento.

La ricostruzione del Calcolo sulla base del concetto di limite e la sua definizione epsilon-delta portò a una riduzione del Calcolo stesso all'aritmetica dei numeri reali. Sulla base di queste chiarificazioni dei fondamenti si fu portati in modo del tutto naturale ad affrontare i fondamenti logici del sistema stesso dei numeri reali. Si ritornò così, dopo due millenni e mezzo, al problema dei numeri irrazionali, che i greci avevano abbandonato dopo Pitagora come irrisolvibile. Uno degli strumenti largamente impiegato in questi tentativi fu proprio la logica matematica o logica simbolica, che si veniva nel frattempo potentemente sviluppando.

Più di recente si è scoperto che la logica matematica fornisce i fondamenti concettuali alla teoria delle macchine calcolatrici e della programmazione coi calcolatori. Quindi questo modello di purezza matematica va oggi considerato come proprio della parte applicativa della matematica stessa.

Il legame tra logica e teoria dei calcolatori è in gran parte costituito dalla nozione di linguaggio formale, che è il tipo di linguaggio che le macchine capiscono. Ed è la nozione di linguaggio formale che ha permesso a Robinson di precisare l'affermazione di Leibniz che si può ragionare senza errore come se gli infinitesimi esistessero.

Leibniz aveva pensato gli infinitesimi come numeri positivi o negativi infinitamente piccoli che ancora godevano delle « stesse proprietà » degli usuali numeri della matematica. A prima vista l'idea sembra autocontraddittoria. Se gli infinitesimi godono delle medesime « proprietà » dei numeri usuali, come possono godere anche della « proprietà » di essere positivi e tuttavia minori di ogni numero positivo usuale? Fu usando un linguaggio formale che Robinson poté risolvere il paradosso. Robinson mostrò come costruire un sistema contenente infinitesimi identico al sistema dei numeri « reali » riguardando a tutte le proprietà esprimibili in un certo linguaggio formale. Ovviamente, la « proprietà » di essere positivo e tuttavia minore di ogni numero positivo usuale risulterà non esprimibile nel linguaggio e in questo modo si eluderà il paradosso.

La situazione è nota a chi abbia talvolta adoperato un calcolatore. Un calcolatore accetta in entrata solo simboli di un certo elenco dato in precedenza all'utente, e i simboli debbono essere usati seguendo certe regole prefissate.

PROPOSIZIONE FORMALE DI L	INTERPRETAZIONE NELL'UNIVERSO STANDARD	INTERPRETAZIONE NELL'UNIVERSO NON-STANDARD
$(\forall x) (\exists y) [x = 0 \vee x y = 1]$ Letteralmente: Per tutti gli x , esiste un y tale che $x = 0$ o $xy = 1$	Ogni numero reale diverso da zero ammette un inverso.	Ogni numero reale non-standard diverso da zero ammette un inverso non-standard; in particolare, infinitesimi positivi ammettono come inverso numeri maggiori di qualsiasi numero reale standard, cioè numeri infiniti.

La proposizione formale è enunciata nel linguaggio L . La colonna centrale dà la sua interpretazione nell'universo standard; quella di destra, nell'universo non-standard.

te. Il linguaggio comune, così come è impiegato nelle comunicazioni umane, è soggetto a regole che i linguisti sono ancora lontani dall'aver completamente compreso. Quando si comunica con i calcolatori, ci si accorge che essi sono « stupidi », proprio per il fatto che — a differenza degli uomini — operano in un linguaggio formale con un vocabolario prefissato e con un prefissato insieme di regole. Gli uomini invece comunicano in un linguaggio naturale, con regole che non sono mai state completamente esplicitate.

La matematica è, naturalmente, una attività umana, come la filosofia o la progettazione dei calcolatori; come queste altre attività, essa è compiuta dall'uomo che si serve di linguaggi naturali. Al medesimo tempo la matematica ha come caratteristica distintiva quella di poter essere adeguatamente descritta da un linguaggio formale, che in un certo senso rispecchia in modo preciso il suo contenuto. Si potrebbe dire che la possibilità di formulare una scoperta matematica in un linguaggio formale è la conferma del fatto che la si è compresa appieno.

Nella analisi non-standard si prendono come punto di partenza i numeri reali finiti e il resto del Calcolo come è noto ai matematici standard. Chiameremo questo « l'universo standard » e lo designeremo con la lettera M . Il linguaggio formale in cui parliamo di M può essere designato con L . Ogni proposizione in L è un enunciato su M e, naturalmente, deve essere o vera o falsa. Cioè, ogni proposizione di L o è vera o è vera la sua negazione. Chiamiamo K l'insieme di tutte le proposizioni vere e diciamo che M è un « modello » per K . Con ciò vogliamo dire che M è una struttura matematica tale che ogni proposizione di K , una volta interpretata come riferentesi a M , è vera. Ovviamente, noi non « conosciamo » K in un qualche modo effettivo; se così fosse avremmo con ciò stesso la risposta a ogni possibile pro-

blema dell'analisi. Nondimeno, consideriamo K come se fosse un oggetto ben definito, a proposito del quale possiamo condurre ragionamenti e trarre conclusioni.

Il fatto essenziale, il punto principale, è che, oltre a M , universo standard, vi sono anche modelli non-standard di K . Cioè, esistono strutture matematiche, M^* , differenti da M in modo essenziale (in un senso che spiegheremo) e che nondimeno sono modelli di K nel senso usuale del termine: ci sono oggetti in M^* e relazioni tra oggetti di M^* tali che se reinterpretiamo i simboli di L in modo da riferirli a questi pseudo-oggetti e pseudo-relazioni in modo opportuno, allora ogni proposizione di K è ancora vera, anche se con un differente significato.

Una analogia grossolana può aiutare l'intuizione. Sia M l'insieme di tutti gli studenti della Facoltà di scienze. Supponiamo, per il nostro ragionamento, che tutti questi studenti siano stati fotografati per l'annuario, in cui tutte le fotografie degli studenti compaiono in formato tessera. Allora M^* potrà essere l'insieme di tutte le fotografie formato tessera di ogni pagina dell'annuario. Senza dubbio, con un'ovvia interpretazione, ogni affermazione vera a proposito di uno studente della Facoltà di scienze, corrisponde a una affermazione vera a proposito di una certa fotografia dell'annuario; tuttavia ci sono molte fotografie nell'annuario che non corrispondono ad alcuno studente. M^* è molto più ampio di M ; oltre agli elementi che corrispondono agli elementi di M , contiene anche molti altri elementi.

Quindi l'enunciato « Giuseppe Rossi è più bravo di Paolo Bianchi », una volta che sia interpretato in M^* , è un enunciato a proposito di due particolari fotografie. Esso non è vero se la relazione « più bravo di » è interpretata in maniera usuale, standard. Così « più bravo di » si deve reinterpretare come una pseudo-relazione, tra pseu-

do-studenti (cioè fotografie di studenti). Potremmo definire la pseudo-relazione «più bravo di» (fra virgolette) dicendo che la fotografia contrassegnata «Giuseppe Rossi» è «più brava della» fotografia contrassegnata «Paolo Bianchi» solo se Giuseppe Rossi è di fatto più bravo di Paolo Bianchi. In questo modo enunciati veri a proposito di studenti sono interpretati come enunciati veri a proposito di fotografie.

Naturalmente, in questo esempio, l'intero ragionamento è un po' costruito ad hoc. Tuttavia, se M è l'universo standard per il Calcolo, M^* , universo non-standard, è un'entità rilevante e interessante.

L'esistenza di interessanti modelli non-standard fu per la prima volta scoperta dal logico norvegese Thoralf A. Skolem, che trovò che gli assiomi del «contare» — gli assiomi che descrivono i «numeri naturali» 1, 2, 3, ecc. — ammettono modelli non-standard che contengono «strani» oggetti non contemplati dalla usuale aritmetica. La grande intuizione di Robinson fu di comprendere come questo originale risultato della logica formale moderna poteva costituire la base per far rivivere i metodi infinitesimali nel calcolo differenziale e integrale. In questa ripresa egli si basava su un teorema dimostrato per la prima volta dal logico russo Anatolij Malcev e generalizzato in seguito da Leon A. Henkin, dell'Università della California a Berkeley. Questo teorema è il teorema di «compattezza». Esso è legato al celebre teorema «di completezza» di Kurt Gödel,

che afferma che un insieme di proposizioni è logicamente coerente (nessuna contraddizione può cioè esser dedotta dalle proposizioni) se e solo se le proposizioni hanno un modello, cioè se e solo se c'è un «universo» in cui esse sono tutte vere.

Il teorema di compattezza afferma: supponiamo di avere una collezione di proposizioni del linguaggio L . Supponiamo che nell'universo standard ogni sottoinsieme finito di questa collezione sia vero. Esiste allora un universo non-standard in cui tutte le proposizioni dell'intera collezione sono simultaneamente vere.

Il teorema di compattezza segue facilmente dal teorema di completezza: se ogni sottoinsieme finito di una collezione di proposizioni di L è vero nell'universo standard, allora ogni sottoinsieme finito è logicamente coerente. Così l'intera collezione di proposizioni è logicamente coerente (poiché ogni deduzione può fare uso solo di un numero finito di premesse). Per il teorema di completezza c'è un universo (non-standard) in cui è vera l'intera collezione.

Una diretta conseguenza del teorema di compattezza è «l'esistenza» di infinitesimi. Per vedere come questo sorprendente risultato segua dal teorema di compattezza consideriamo le proposizioni:

« c è un numero maggiore di zero e minore di $1/2$ »

« c è un numero maggiore di zero e minore di $1/3$ »

« c è un numero maggiore di zero e minore di $1/4$ », e così via.

Questa è una collezione infinita di proposizioni, ciascuna delle quali può essere scritta servendosi del linguaggio formale L . Se ci riferiamo all'universo standard R dei numeri reali, ogni suo sottoinsieme finito è vero, perché se consideriamo solo un numero finito di proposizioni della forma « c è un numero maggiore di zero e minore di $1/n$ », una delle proposizioni conterrà la frazione più piccola $1/n$, e $1/2n$ sarà ovviamente maggiore di zero e minore di tutte le frazioni che compaiono nel nostro elenco finito di proposizioni. Tuttavia, se consideriamo l'intero insieme infinito di queste proposizioni, esso risulta falso in riferimento ai numeri reali standard, poiché per quanto piccolo possiamo scegliere un numero reale positivo c , $1/n$ risulterà più piccolo di c prendendo n sufficientemente grande.

Il teorema di compattezza di Malcev ed Henkin asserisce che c'è un universo non-standard R^* che contiene numeri pseudo-reali compreso un numero pseudo-reale positivo c , minore di qualsiasi numero della forma $1/n$. Cioè, c è infinitesimo. Inoltre c gode di tutte le proprietà dei numeri reali standard in un senso ben preciso: ogni enunciato vero a proposito di reali standard, che possiamo esprimere nel linguaggio formale L , è vero anche a proposito di reali non-standard, compreso l'infinitesimo c — con un'opportuna interpretazione. (La fotografia di «Giuseppe Rossi» non è più brava in realtà della fotografia di «Paolo Bianchi», ma l'enunciato «Giuseppe Rossi» è «più bravo di» «Paolo Bianchi» è vero, con la nostra interpretazione non-standard di «più bravo di»). D'altra parte, proprietà godute da tutti i numeri reali standard possono non essere trasferite agli pseudo-numeri non-standard, se non è possibile esprimere queste proprietà nel linguaggio formale L .

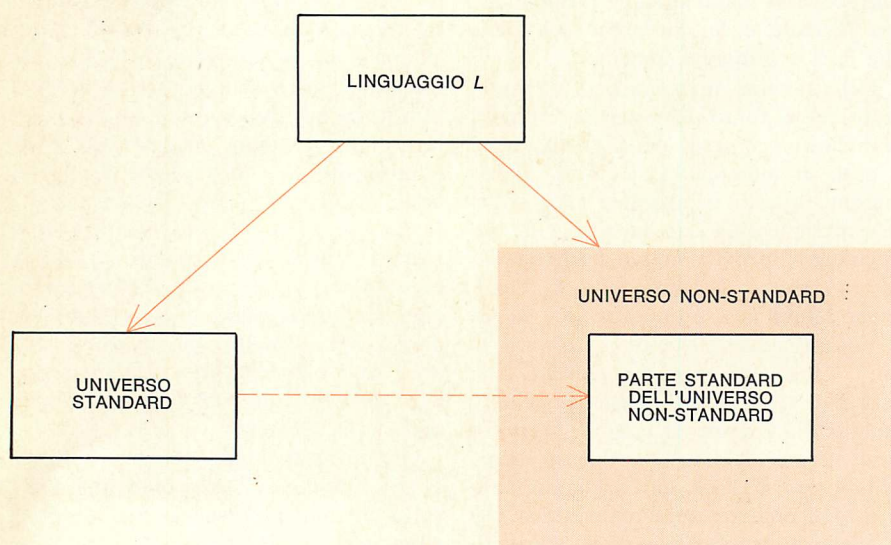
La proprietà archimedeica (non esistenza di infinitesimi) di R può essere espressa facendo uso di un insieme infinito di proposizioni di L come segue (facciamo uso del simbolo « $>$ » come al solito per significare «è maggiore di»). Per ogni elemento positivo c di R , sono vere tutte le proposizioni del tipo qui sotto riportato, tranne un numero finito:

$$c > 1$$

$$c + c > 1$$

$$c + c + c > 1 \text{ e così via.}$$

Questo non è vero, tuttavia, per gli pseudo-reali R^* : se c è infinitesimo (quindi pseudo-reale), tutte queste proposizioni sono false. In altri termini, nessuna somma di un numero finito di



Una raffigurazione schematica del ruolo di mediazione del linguaggio formale tra l'universo standard e quello non-standard. Il linguaggio formale L descrive l'universo standard, che include i numeri reali. Le proposizioni di L vere nell'universo standard sono vere anche in quello non-standard, che contiene ulteriori oggetti matematici, come gli infinitesimi. In questo modo l'analisi non-standard rende preciso il metodo infinitesimale.

addendi tutti eguali a c può essere maggiore di 1, qualunque numero di addendi prendiamo. Proprio il fatto che la proprietà archimedeica è vera nel modello standard, ma falsa in quello non-standard, prova che la proprietà non può essere espressa con una proposizione di L ; l'enunciato che abbiamo usato è in realtà costituito da un numero infinito di proposizioni. Essi si comportano « formalmente » come oggetti standard eppure differiscono per importanti proprietà che non sono formalizzabili in L .

Benché l'universo non-standard sia concettualmente distinto da quello standard, è comodo pensare quest'ultimo come un ampliamento dell'universo standard. Poiché R^* è un modello di L , ogni proposizione vera a proposito di R ha un'interpretazione in R^* . In particolare i nomi di numeri in R hanno un'interpretazione come nomi di oggetti in R^* . Possiamo per semplicità identificare l'oggetto di R^* chiamato « 2 » con l'usuale numero 2 di R . Allora R^* contiene i numeri reali standard di R , insieme con un'ampia collezione di quantità infinitesime e infinite, in cui R è immerso.

Un oggetto di R^* (un numero pseudo-reale) è detto infinito se è pseudo-maggiore di ogni numero reale standard; altrimenti è detto finito. Un numero pseudo-reale positivo è detto infinitesimo se è pseudo-minore di ogni numero reale positivo standard. Se la pseudo-differenza di due pseudo-reali è finita, diciamo che essi appartengono alla stessa « galassia »; l'asse pseudo-reale contiene un numero infinito più che numerabile di galassie. Se la pseudo-differenza di due pseudo-reali è infinitesima, diciamo che essi appartengono alla stessa « monade » (termine che Robinson ha preso a prestito dagli scritti filosofici di Leibniz). Se un pseudo-reale r^* è infinitamente vicino a un numero reale standard r , diciamo che r è la parte standard di r^* . Tutti i reali standard sono ovviamente nella stessa galassia, che è detta galassia principale. Nella galassia principale ogni monade contiene uno e un solo numero reale standard. Questa monade è « l'intorno infinitesimo » di r : cioè, l'insieme dei numeri reali non-standard infinitamente vicini a r . La nozione di monade risulta applicabile non solo ai numeri reali; ma anche a spazi metrici e topologici in genere. L'analisi non-standard è dunque rilevante non solo per il Calcolo elementare ma anche per l'intero campo della moderna analisi astratta.

Quando diciamo che infinitesimi o monadi esistono, dovrebbe esser chiaro che non intendiamo affatto ciò co-

me sarebbe stato inteso da Euclide o da Berkeley. Fino a un secolo fa era tacitamente assunto da tutti i filosofi e matematici che l'oggetto della matematica fosse dotato di realtà obiettiva in un senso molto vicino a quello in cui l'oggetto della fisica è reale. Se esistessero o no gli infinitesimi era questione di fatto, non troppo differente dal problema se gli atomi materiali esistono o no. Oggi molti matematici, forse la maggior parte, non condividono più tale convinzione dell'esistenza obiettiva degli oggetti che essi studiano. La teoria dei modelli non si compromette in un modo o nell'altro su tali questioni ontologiche. Quel che i matematici vogliono dagli infinitesimi, non è una forma di esistenza materiale, ma piuttosto il diritto di servirsi di essi nelle dimostrazioni. A questo scopo tutto ciò di cui si ha bisogno è la certezza che una dimostrazione che faccia uso di infinitesimi non è più errata di una che non ne fa uso.

L'impiego dell'analisi non-standard nella ricerca fa qualcosa di simile. Supponiamo di voler dimostrare un teorema che fa riferimento solo a oggetti standard. Se si immergono gli oggetti standard nell'ampliamento non-standard, si può trovare una dimostrazione molto più breve e più intuitiva usando oggetti non-standard. Il teorema è stato allora dimostrato in realtà facendo riferimento a un'interpretazione non-standard dei suoi termini e simboli. Gli oggetti non-standard che corrispondono a oggetti standard hanno la proprietà caratteristica che proposizioni su di essi sono vere (nell'interpretazione non-standard) solo se la stessa proposizione è vera in riferimento a oggetti standard (nell'interpretazione standard). In questo modo dimostriamo teoremi su oggetti standard ragionando su oggetti non-standard.

Ricordiamo, per esempio, la « dimostrazione » di Nicola Cusano che l'area di un cerchio di raggio 1 è eguale alla metà della sua circonferenza. Vediamo in che senso nella teoria di Robinson il ragionamento di Cusano è corretto. Dal momento che disponiamo di numeri infiniti e infinitesimi (nell'universo non-standard), si può dimostrare che l'area del cerchio è la parte standard della somma (nell'universo non-standard) di un numero infinito di infinitesimi.

Vediamo ora che forma assumerebbe il problema della caduta della pietra secondo le idee di Robinson. Definiamo la velocità istantanea non come il rapporto di incrementi infinitesimi, come faceva de l'Hôpital, ma piuttosto come la parte standard di tale rapporto; allora ds , dt e il loro rapporto

ds/dt sono numeri reali non-standard. Come prima abbiamo $ds/dt = 9,8 + 4,9dt$, ma ora immediatamente concludiamo, rigorosamente e senza alcuna restrizione, che v , parte standard di ds/dt è eguale a 9,8. Una lieve modificazione del metodo degli infinitesimi di Leibniz, cioè la distinzione precisa tra il numero non-standard ds/dt e la sua parte standard v , elimina la contraddizione, che l'Hôpital semplicemente ignorava.

Naturalmente, si richiede la dimostrazione del fatto che la definizione di Robinson dia il medesimo risultato di quella di Weierstrass. Tale dimostrazione non presenta difficoltà, ma non è il caso di darla qui.

Ciò che si è ottenuto è che il metodo infinitesimale è stato reso preciso per la prima volta. Nel passato i matematici dovevano fare una scelta. Se usavano infinitesimi, dovevano basarsi sull'esperienza e sull'intuizione per ragionare correttamente. « Andate avanti – così si suppone che Jean Le Rond d'Alembert avesse rassicurato un amico matematico esitante – e la fede vi verrà. » Per una rigorosa certezza era stato necessario ricorrere all'ingombrante metodo archimedeo di esaurimento o alla sua versione moderna, il metodo epsilon-delta di Weierstrass. Ora il metodo degli infinitesimi, o più in generale il metodo delle monadi, si eleva dal livello euristico a quello rigoroso. L'approccio dal punto di vista della logica formale permette di rispondere compiutamente alla questione sollevata da Berkeley e dagli altri protagonisti delle controversie del passato, cioè se le quantità infinitesime esistono o no in qualche senso oggettivo.

Le applicazioni che abbiamo discusso in questa sede sono elementari, addirittura banali. Si sono fatte e si stanno facendo applicazioni non banali. Sono stati pubblicati lavori sulla dinamica non-standard e sulla probabilità non-standard. Robinson, e il suo allievo Allen Bernstein, han fatto uso di analisi non-standard per risolvere un problema – precedentemente insoluto – riguardante operatori lineari compatti. Si deve nondimeno dire che molti analisti restano scettici a proposito della importanza ultima del metodo di Robinson. È del tutto vero che ogni cosa che può esser fatta con gli infinitesimi può, in linea di principio, esser fatta senza di essi. Forse, come per altre radicali innovazioni, un uso esaustivo delle nuove idee sarà fatto da una nuova generazione di matematici, non tanto immersi nei metodi standard da non poter godere della libertà e delle ampie possibilità dell'analisi non-standard.

Nuovi modelli del sistema dei numeri reali

È passato solo un secolo da quando i matematici hanno elaborato una sistemazione logica soddisfacente della nozione di numero reale; le indagini successive hanno però già portato alla luce problemi nuovi

di Gabriele Lolli

Sono trascorsi più di duemila anni dall'epoca in cui i pitagorici erano costretti a tenere nascosta la scoperta dell'esistenza di grandezze incommensurabili. Oggi i sistemi numerici, che sono uno strumento quotidiano indispensabile in ogni campo della scienza, conservano solo in qualche locuzione (come quella dei numeri immaginari, o impossibili), una traccia di antiche difficoltà; ma non tutti sanno che questa conquista del pensiero moderno ha una data molto recente. Solo nel 1872 i matematici sono infatti riusciti a elaborare una sistemazione logica soddisfacente della nozione di numero reale. Eppure le indagini posteriori, mai interrotte, su questa fondazione dell'intero edificio delle matematiche, hanno contribuito a portare alla luce problemi nuovi.

Un'importante svolta si è avuta nel 1963 nella teoria degli insiemi, in seguito ai risultati di Cohen (si veda l'articolo *La teoria non cantoriana degli insiemi* di P. J. Cohen e R. Hersh in questo volume) che riusciva a provare l'impossibilità di risolvere, nella teoria degli insiemi, questioni essenziali per la caratterizzazione del sistema dei numeri reali. Nel presentare la situazione della disciplina, servendosi di una analogia con la evoluzione delle geometrie non euclidee, Cohen e Hersh la definivano come «il periodo della elaborazione di teorie non standard», in alternativa tra loro, riservando al futuro il responso sulle possibili applicazioni dei frutti di queste ricerche.

Un primo risultato del lavoro svolto negli ultimi anni è il perfezionamento delle tecniche usate per costruire modelli diversi della teoria. Per vedere all'opera tali tecniche, non è necessario considerare la teoria degli insiemi nel suo complesso, e non è neppure conve-

niente: è innegabile che rappresentarsi modelli alternativi di una teoria che dovrebbero esprimere le proprietà dell'intero universo degli oggetti matematici sia un'operazione psicologicamente e logicamente problematica. È sufficiente però considerare la più semplice teoria dei numeri reali, presentata non come una parte della teoria degli insiemi, ma con una assiomatizzazione indipendente; vedremo tuttavia che non è possibile enunciare gli assiomi che caratterizzano i numeri reali senza far intervenire quelle costruzioni intellettuali che permettono ai nuovi metodi di scattare e di produrre modelli alternativi.

Descriveremo quindi una tecnica che permette di costruire un modello particolarmente interessante per i suoi legami con la nozione di probabilità. Tale modello è profondamente diverso dal modello intuitivo dei numeri reali costituito dalla retta euclidea. Questa situazione appare contraddire non solo la comune sensazione che i numeri reali siano qualcosa di ben definito e assoluto, ma anche un teorema di unicità (dei modelli della teoria, come vedremo più avanti), che gli studenti imparano a dimostrare e che è giusto, come ora sappiamo, solo in un senso piuttosto sofisticato.

Cominceremo comunque riassumendo nel paragrafo che segue la teoria classica dei numeri reali.

L'insieme dei numeri reali, che indicheremo con \mathbf{R} , contiene i numeri naturali, i numeri interi relativi, i numeri razionali e i numeri irrazionali. I numeri naturali (0, 1, 2, ...) sono quelli che servono per contare. Gli interi relativi e i razionali si ottengono ampliando successivamente il sistema originario dei numeri naturali affinché si

possano eseguire la sottrazione e, rispettivamente, la divisione, operazioni inverse delle due operazioni fondamentali: addizione e moltiplicazione. Pensiamo che il lettore sia familiare con il sistema dei numeri razionali, che indicheremo con \mathbf{Q} , ma, per comodità, riassumiamo nella figura della pagina a fronte la definizione più comune dei successivi ampliamenti che portano ai numeri razionali. Non si tratta di una successione storica (perché la storia di questi sistemi di numeri non è affatto lineare), ma di una successione logica, da cui si vede che occorrono solo operazioni insiemistiche semplici, oltre alla nozione di base di numero naturale, per definire il sistema \mathbf{Q} . I numeri razionali possono essere scritti sia come frazioni, sia con la rappresentazione decimale, limitata o periodica.

La rappresentazione più comune dei numeri irrazionali è invece quella dei decimali illimitati non periodici, e un momento di riflessione mostra che in questa locuzione si nascondono concetti più complessi dei precedenti. Un decimale illimitato lo si pensa di solito come una successione di decimali limitati, e quindi di numeri razionali: quelli che si ottengono, per ogni n , prendendo i primi n termini della successione. Così π si identifica con la successione 3, 3,1, 3,14, Questa immagine si accorda sia con la definizione rigorosa, sia con la pratica dell'uso dei numeri reali nelle misurazioni.

Fino al secolo scorso, i matematici si distinguevano secondo due tendenze, nella considerazione dei numeri reali. Una di tipo assiomatico, risalente al calcolo dei rapporti di Eudosso (IV secolo a.C.), che più che dei numeri parlava del confronto di grandezze; una di tipo pratico-calcolistico, che si accontentava di postulare la possibi-

lità di eseguire misurazioni di qualsiasi lunghezza con l'approssimazione voluta, mediante i numeri razionali. Si consideri l'esempio classico di numero irrazionale, quello che, stando alla tradizione, fu il primo esempio di una lunghezza non misurabile con un numero razionale: la lunghezza della diagonale di un quadrato di lato unitario che, per il teorema di Pitagora, indichiamo con $\sqrt{2}$. La dimostrazione che tale « numero » non è razionale è facile e nota. È possibile tuttavia approssimare la lunghezza della diagonale mediante numeri razionali, con misurazioni sempre più accurate, in modo che l'errore sia piccolo a piacere. Per esempio la successione, di numeri razionali, definita per ricorsione

$b_0 = 2, b_1 = b_0/2 + 1/b_0, \dots, b_{n+1} = b_n/2 + 1/b_n$
(i cui primi termini sono 2, 3/2, 17/12, ...), è una successione decrescente e tale che b_n^2 tende a 2 (si veda la figura in basso a pagina 44). Questo significa che, comunque si prefissi un « errore » $1/m$, si può determinare un indice n_0 della successione tale che di lì in avanti, cioè per ogni $n > n_0$, si abbia

$$2 < b_n^2 < 2 + 1/m$$

Se la successione dei numeri b_n^2 tende, o converge, al limite 2, ci si aspetta che la successione dei numeri

b_n converga a $\sqrt{2}$; in effetti si può dimostrare che non può convergere né a un numero il cui quadrato sia minore di 2, né a un numero il cui quadrato sia maggiore di 2. Ne segue che, in \mathbb{Q} , la successione, che indicheremo con $\{b_n\}$, non è convergente.

I numeri razionali formano un sistema soddisfacente, per quel che riguarda le manipolazioni con le operazioni algebriche; formano, come si dice, un corpo ordinato (si veda la figura a pagina 45). Quello che manca è la validità del postulato pratico che a ogni grandezza lineare si possa associare un numero come misura e che, in termini matematici, si può formulare dicendo che successioni di numeri del tipo di quella esaminata sopra devono convergere a qualche numero. Il sistema dei razionali deve essere completato riempiendo, in un certo senso, le lacune che le successioni di questo tipo rivelano. Ma quale è la caratteristica di queste successioni a cui deve essere associato un limite? Occorre avere un criterio interno alle successioni, una proprietà che si possa stabilire esaminando solo i termini della successione e nessun elemento estraneo. Georg Cantor (1845-1918), perfezionando il precedente lavoro di Cauchy sulle successioni, propose nel 1872 di assumere come cri-

terio la cosiddetta condizione di Cauchy: una successione $\{b_n\}$ soddisfa questa condizione, ed è detta allora fondamentale (o successione di Cauchy), se al crescere dell'indice n i termini della successione diventano arbitrariamente vicini tra loro. Più esattamente, fissata comunque una « distanza » $1/m$, si può determinare un indice n_0 al di là del quale (cioè per tutti gli indici p e q maggiori di n_0) si abbia

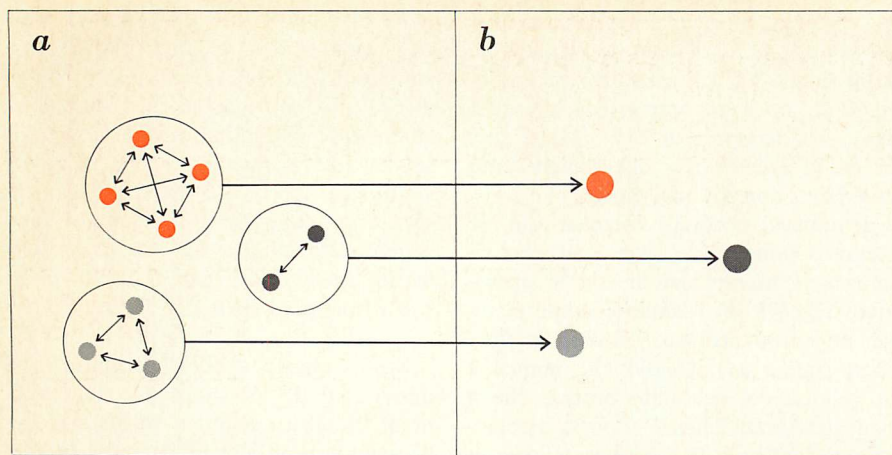
$$|b_p - b_q| < 1/m.$$

Scelto questo criterio, la natura dei nuovi enti da associare alle successioni di Cauchy non è un problema; essi possono essere identificati con le successioni stesse, tra cui si possono definire le operazioni eseguendole termine a termine. La somma delle due successioni $\{b_n\}$ e $\{c_n\}$ è la successione $\{b_n + c_n\}$, e così via. L'insieme dei nuovi enti contiene una copia esatta dei numeri razionali, perché ogni numero razionale si può identificare con la successione costante i cui termini sono tutti uguali al numero dato. C'è una complicazione, dovuta al fatto che successioni diverse possono tendere allo stesso limite; allora bisogna « identificare » tra loro tutte le successioni la cui differenza tende a zero; questa operazione, detta passaggio al quoziente, è frequentissima in matematica, ed

<p style="text-align: center;">NUMERI NATURALI</p> <p style="text-align: center;">$0, 1, 2, \dots, n, \dots$</p>
<p style="text-align: center;">NUMERI INTERI RELATIVI</p> <p style="text-align: center;">$\dots; -2, -1, 0, +1, +2, \dots$</p> <p>I numeri interi relativi si possono definire a partire dai numeri naturali nel seguente modo. Si considerino tutte le coppie ordinate (m, n) di numeri naturali; intuitivamente, la coppia (m, n) rappresenterà il numero relativo $m - n$. Tra queste coppie si definiscano le operazioni in questo modo: $(m, n) + (p, q) = (m + p, n + q)$, $(m, n) \cdot (p, q) = (m \cdot p + n \cdot q, m \cdot q + n \cdot p)$. Le definizioni riflettono quello che si fa in pratica per sommare e moltiplicare i due numeri $m - n$ e $p - q$; due coppie (m, n) e (p, q) rappresentano lo stesso numero se $m + q = p + n$. In termini più precisi la relazione $(m, n) R (p, q)$ se e solo se $m + q = p + n$ è una relazione di equivalenza, e i numeri interi relativi si ottengono identificando tutte le coppie tra loro equivalenti; i numeri interi positivi sono quelli rappresentati da coppie del tipo $(m, 0)$ e quelli negativi da coppie del tipo $(0, m)$.</p>
<p style="text-align: center;">NUMERI RAZIONALI</p> <p style="text-align: center;">$\frac{1}{2} = 0,5 ; \frac{4}{3} = 1,333 \dots = 1,\overline{3} ; \dots$</p> <p>I numeri razionali, o frazioni, si possono definire a partire dai numeri interi relativi nel seguente modo. Si considerino tutte le coppie ordinate (x, y) di numeri interi, con $y \neq 0$; il simbolo (x, y) sarà poi sostituito da quello più corrente $\frac{x}{y}$. Tra queste coppie si definiscano le operazioni nel modo seguente, che non è altro se non quello che si fa nella pratica per calcolare con le frazioni: $(x, y) + (u, v) = (x \cdot v + y \cdot u, y \cdot v)$ e $(x, y) \cdot (u, v) = (x \cdot u, y \cdot v)$. La coppia $(0, 1)$ diventa lo zero dei razionali e $(1, 1)$ l'unità. I numeri interi, contenuti nei razionali, sono le coppie $(x, 1)$. Esistono coppie diverse che intuitivamente rappresentano lo stesso numero. Noi consideriamo come la « stessa » frazione una qualunque delle frazioni $\frac{2}{3}, \frac{4}{6}, \frac{-2}{-3}, \dots$. In termini più precisi, questo significa che la relazione $(x, y) R (u, v)$ se e solo se $x \cdot v - y \cdot u = 0$ è una relazione di equivalenza e che i numeri razionali si ottengono identificando tutte le coppie tra loro equivalenti. Come rappresentante di ciascuna di queste classi si può prendere una coppia qualunque che stia nella classe, di solito quella che noi chiamiamo frazione ridotta ai minimi termini.</p>

Il passaggio dai numeri naturali ai numeri interi e quindi ai razionali può essere giustificato con il « metodo delle coppie » (e

identificazione per equivalenza) proposto nella prima metà del secolo XIX dal matematico irlandese William Rowan Hamilton.

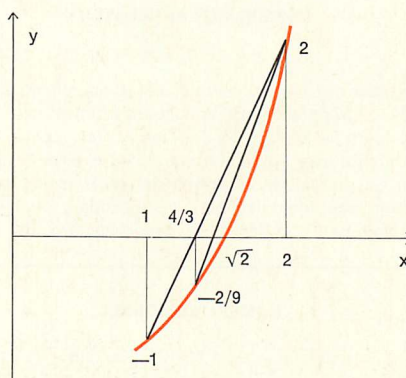
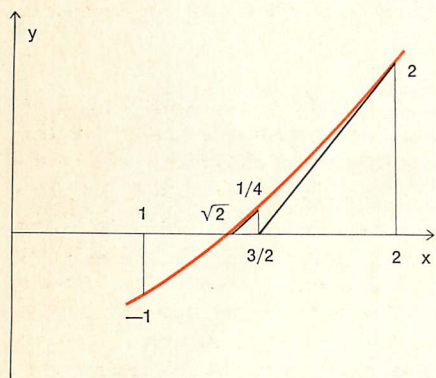


Una relazione di equivalenza R tra gli elementi di un insieme X è una relazione riflessiva (xRx), simmetrica ($xRy \rightarrow yRx$) e transitiva ($xRy \& yRz \rightarrow xRz$). Se uniamo con delle frecce gli elementi che si corrispondono nella relazione R , il grafico della relazione risulta allora formato da cicli chiusi e disgiunti come nella figura a , dove abbiamo trascurato le frecce che partono e arrivano a uno stesso elemento. Fare il quoziente dell'insieme X rispetto alla relazione di equivalenza R significa sostituire un unico oggetto a ciascuno di questi cicli, come in figura b . Se per esempio si considera la relazione di equivalenza «vivere nello stesso appartamento» la figura a può rappresentare tre famiglie che abitano in appartamenti distinti e la figura b rappresenta invece gli appartamenti stessi. Quelli che abbiamo chiamato cicli vengono detti più propriamente, in linguaggio tecnico, classi di equivalenza, e un qualsiasi elemento di una classe si dice rappresentante della classe; nel nostro esempio, uno qualunque dei membri di ciascuna famiglia. Il passaggio al quoziente si esegue quando le operazioni o i ragionamenti che si devono fare non dipendono dalla scelta del rappresentante delle classi; si può anche scegliere un rappresentante da ogni classe e cancellare gli altri.

è già presente nella definizione dei numeri razionali (si veda l'illustrazione alla pagina precedente e qui sopra).

Non è il caso di impegnarsi in tutti i particolari della definizione di \mathbf{R} ; due osservazioni sono importanti, sul risultato finale di questa costruzione: da una parte si può dimostrare la possibi-

lità di associare a ogni numero reale la solita rappresentazione decimale, e si può lavorare tranquillamente come si è sempre fatto; dall'altra \mathbf{R} risulta un'estensione completa di \mathbf{Q} nel seguente senso: le successioni di numeri reali che soddisfano alla condizione di Cauchy convergono a numeri reali.



Una successione di numeri razionali che converge a $\sqrt{2}$ decrescendo (approssimazioni per eccesso) si può ottenere con il cosiddetto metodo delle tangenti (a). Fissato il punto $b_0 = 2$, si conduce la tangente alla curva $y = x^2 - 2$ nel punto $(2, 2)$ e si individua l'ascissa $b_1 = 3/2$ in cui essa incontra l'asse delle x ; quindi si ripete il procedimento partendo dal punto $(3/2, 1/4)$, e così via. Nei corsi di calcolo infinitesimale si dimostra che tale successione è definita dalla formula ricorrente $b_{n+1} = (1/2)b_n + 1/b_n$, e che converge al punto in cui $x^2 - 2 = 0$. Una successione di numeri razionali che converge a $\sqrt{2}$ crescendo (approssimazioni per difetto) si può ottenere con il cosiddetto metodo delle corde (b). (I disegni a e b sono in scala diversa perché la figura risulti più chiara). Fissato il punto $b_0 = 1$, si conduce la corda che unisce i punti $(1, -1)$ e $(2, 2)$ e si individua l'ascissa $b_1 = 4/3$, in cui essa incontra l'asse delle x ; quindi si ripete il procedimento con i due punti $(4/3, -2/9)$ e $(2, 2)$, e così via. Si dimostra che tale successione è definita dalla formula ricorrente $b_{n+1} = (4 + 2b_n - 2b_n^2) / (4 - b_n^2)$ e che converge a $\sqrt{2}$. Le due successioni convergenti a $\sqrt{2}$ rappresentano lo stesso numero reale, se questi sono definiti per mezzo delle successioni di Cauchy.

Dunque \mathbf{R} non si può più estendere con il procedimento descritto, né c'è bisogno di farlo.

Le proprietà essenziali dei numeri reali, quelle algebriche e quella della completezza, si possono esplicitare in una lista finita di assiomi. Nella figura della pagina a fronte sono elencati questi assiomi, che caratterizzano i cosiddetti corpi ordinati continui; la condizione della completezza, che abbiamo discusso prima, è ivi formulata in termini diversi: ogni insieme limitato superiormente ha un minimo confine superiore. Questa versione si può dimostrare equivalente alla precedente, e deriva più direttamente da un'altra definizione dei numeri reali, proposta nello stesso 1872 da Richard Dedekind (1831-1916). Egli osservò che la proprietà di continuità che noi associamo alla retta si può esprimere nel seguente modo: se dividiamo l'intero sistema di punti in due classi, in modo che ogni punto cada in almeno una delle due classi e tutti i punti di una classe precedano tutti quelli della seconda (nell'ordine intuitivo da sinistra verso destra), allora c'è uno ed un solo punto che funge da elemento separatore delle due classi, nel senso che è il massimo della prima o il minimo della seconda. Nel sistema dei numeri razionali questo non è sempre vero, come si vede se si considera la sezione che mette in una classe tutti i numeri il cui quadrato è minore di 2, e nell'altra tutti i restanti. Una estensione di \mathbf{Q} che abbia questa proprietà si può ottenere identificando i nuovi numeri con le sezioni stesse dell'insieme \mathbf{Q} e definendo opportunamente le operazioni. Il risultato è un sistema di enti che soddisfa tutti gli assiomi dei corpi ordinati continui, che è logicamente definito solo usando la nozione di classe, o proprietà, dei numeri razionali e che si presenta come la realizzazione di quella che intuitivamente pensiamo essere la retta.

Le due strutture che abbiamo costruito, una con le successioni fondamentali, l'altra con le sezioni di Dedekind, sono due modelli diversi degli assiomi dei corpi ordinati continui. Ma nonostante gli oggetti che costituiscono le strutture, da una parte le successioni, dall'altra le sezioni, siano diversi, i due modelli sono sostanzialmente la stessa cosa; in matematica questa affermazione si precisa dimostrando che sono isomorfi (si veda la figura a pagina 47). Si può anzi dimostrare di più: due corpi ordinati continui qualunque sono isomorfi. Questo è il teorema di unicità che avevamo annunciato: il

1. PROPRIETÀ ASSOCIATIVA DELLA ADDIZIONE E DELLA MOLTIPLICAZIONE

$$\forall x \forall y \forall z [(x + y) + z = x + (y + z)]$$

$$\forall x \forall y \forall z [(x \cdot y) \cdot z = x \cdot (y \cdot z)]$$

2. PROPRIETÀ COMMUTATIVA DELLA ADDIZIONE E DELLA MOLTIPLICAZIONE

$$\forall x \forall y (x + y = y + x) ; \forall x \forall y (x \cdot y = y \cdot x)$$

3. ESISTENZA DEGLI ELEMENTI NEUTRI

$$\forall x (x + 0 = x) ; \forall x (x \cdot 1 = x) ; 0 \neq 1$$

4. ESISTENZA DEGLI INVERSI

$$\forall x \exists y (x + y = 0) ; \forall x (x \neq 0 \rightarrow \exists y (x \cdot y = 1))$$

5. PROPRIETÀ DISTRIBUTIVA

$$\forall x \forall y \forall z [x \cdot (y + z) = x \cdot y + x \cdot z]$$

6. PROPRIETÀ DELL'ORDINE

$$\forall x \forall y \forall z (x < y \text{ \& } y < z \rightarrow x < z)$$

$$\forall x \forall y (x < y \vee x = y \vee y < x)$$

$$\forall x \forall y \forall z (x < y \rightarrow x + z < y + z)$$

$$\forall x \forall y \forall z (x < y \text{ \& } 0 < z \rightarrow x \cdot z < y \cdot z)$$

7. CONTINUITÀ

$$\forall S [\exists y \forall x (x \in S \rightarrow x < y) \rightarrow \exists y \forall z (y \leq z \leftrightarrow \forall x (x \in S \rightarrow x \leq z))]$$

Gli assiomi che caratterizzano i numeri reali sono scritti nel linguaggio simbolico che è comunemente usato in matematica: $\forall x \dots$ si legge: per tutti gli $x \dots$; $\exists x \dots$: esiste un x tale che...; \vee : oppure; $\&$: e; \rightarrow : implica; \leftrightarrow : se e solo se; \sim : non; $x < y$: x è minore di y ; $x \leq y$: $x < y \vee x = y$; $x \in S$: x è un elemento di S o x appartiene a S . Un insieme X in cui siano definite due operazioni, indicate con $+$ e \cdot , e una relazione $<$, e che contenga due elementi particolari 0 e 1 si indica $\langle X, +, \cdot, <, 0, 1 \rangle$. Se le operazioni $+$ e \cdot soddisfano le proprietà da 1 a 5 si dice che $\langle X, +, \cdot, <, 0, 1 \rangle$ è un corpo commutativo. Se inoltre sono vere le proprietà dell'ordine 6, si dice che è un corpo ordinato. Se è vero anche l'assioma

di continuità, quando si interpretino le lettere maiuscole come sottoinsiemi di X , si dice che è un corpo ordinato continuo. L'assioma di continuità afferma che se un insieme è limitato superiormente, allora esiste un numero che è il più piccolo di tutti i maggioranti, e che si chiama minimo confine superiore, o estremo superiore, e si indica con $\sup(S)$. Con $\inf(S)$ si indica il più grande dei numeri che sono minori di tutti gli elementi di S , o estremo inferiore. L'assioma di continuità implica anche che, se un insieme è limitato inferiormente, allora esiste l'estremo inferiore. La costruzione dei numeri reali dai razionali ricordata nel testo si può generalizzare a corpi ordinati qualunque, che soddisfino un'ulteriore proprietà, detta di Archimede.

completamento dei razionali è unico.

Quando tutti i modelli di un sistema di assiomi sono isomorfi, noi associamo agli assiomi un unico concetto ben definito, per esempio quello dei numeri reali, e i singoli modelli non sono altro che raffigurazioni intuitive diverse di tale concetto. Ecco perché possiamo parlare dell'insieme \mathbf{R} dei numeri reali, senza precisare come sia stato costruito \mathbf{R} . Questo non succede con tutti i sistemi di assiomi; per esempio, sia \mathbf{Q} , sia \mathbf{R} , sono corpi ordinati, ma non sono isomorfi. È l'assioma di continuità che, tra tutti i possibili corpi ordinati, opera una selezione, dopo la quale rimangono soltanto delle strutture isomorfe tra loro.

Come è possibile allora che Cohen abbia costruito dei sistemi di numeri reali di cardinalità diversa, e quindi in particolare non isomorfi? Anche se non

possiamo esaminare la dimostrazione del teorema di unicità, dobbiamo dire che la spiegazione è da ricercarsi nella particolare natura dell'assioma di continuità. L'assioma fa riferimento a tutti i possibili sottoinsiemi di \mathbf{R} ; si dice anche che è formulato nella logica del secondo ordine, mentre gli altri assiomi, che parlano soltanto degli elementi di \mathbf{R} , richiedono la più semplice logica del primo ordine. Per poter fare un qualsiasi ragionamento, abbiamo bisogno di criteri logici che regolino l'uso di questa nozione della totalità di tutti i sottoinsiemi di un insieme infinito. Tali criteri, una volta esplicitati, non sono altro che proprietà dell'universo, codificate negli attuali assiomi della teoria degli insiemi. Ogni teorema che segua da questi assiomi, o da un ragionamento che ne faccia tacitamente uso, sarà vero nell'universo. Il fatto è che non è possibile caratterizzare in modo

unico l'universo degli insiemi per mezzo di un sistema assiomatico del tipo di quelli in uso; ciò significa che è possibile definire delle totalità differenti tra loro, ma in cui sono ugualmente vere tutte le proprietà che abbiamo codificato negli assiomi. Gli universi in particolare differiscono nel numero di sottoinsiemi degli insiemi infiniti che contengono. Può allora succedere che, se anche il teorema di unicità è vero in ogni universo (perché discende dagli assiomi), corpi ordinati continui presenti in *universi diversi* non sono isomorfi tra loro.

Fino a che non si riuscirà a caratterizzare un *unico* universo degli insiemi, si potranno stabilire solo dei risultati validi nei singoli universi possibili, ma non in assoluto.

Ma c'è anche un altro senso in cui il teorema di unicità può venir meno, senso che si può cogliere affinando

l'uso della logica e ridefinendo la nozione di modello di un sistema di assiomi. Si tratta di introdurre i cosiddetti modelli *non standard*, di cui non daremo la definizione generale, ma un esempio relativo appunto alla teoria dei numeri reali. Bisogna però ricordare che la tecnica usata è una elaborazione di quelle applicate originariamente alla teoria degli insiemi.

Il nostro punto di partenza è l'osservazione che nella pratica scientifica si incontrano degli enti che sono usati *come se fossero numeri reali*, pur essendo in realtà qualcosa di molto più complicato: si tratta delle variabili casuali (o aleatorie) su uno spazio di probabilità. La precisazione di questo concetto e del suo uso ci condurrà in modo naturale alla definizione di un modello non standard dei numeri reali. Variabile casuale è già un nome più sofisticato per altri, come elemento aleatorio, di uso comune; si dice di solito che una variabile casuale è «una quantità i cui valori sono determinati dal caso». Con ciò si intende che dei valori numerici vengono associati al risultato di un esperimento, esperimento il cui esito può a priori variare in un certo campo. Con esperimento si intende o una prova guidata, come il lancio di un dado (non truccato), o una situazione osservabile, per esempio una certa configurazione degli elementi atmosferici. I valori associati ai risultati dell'esperimento potrebbero essere, nei due esempi, o dei punteggi di scommesse o, rispettivamente, dei valori della temperatura o della pressione. Questi valori sono poi manipolati in operazioni, relazioni, previsioni, come se fossero dei numeri, e in realtà lo sono, ma dipendono dal caso. Questo significa che dobbiamo valutare, e portarci dietro in qualche modo nei calcoli e nei ragionamenti, la probabilità che i valori siano quelli che manipoliamo. Quindi innanzi tutto dobbiamo valutare la probabilità che i risultati dell'esperimento siano o no di un certo tipo. Se l'esperimento è il lancio di un dado, tipiche situazioni possibili, cui devono essere assegnate delle probabilità, sono: «il risultato è un numero pari», «il risultato è un numero minore di 3», e così via. A queste possibilità si dà il nome di eventi. Dall'esempio citato si vede che gli eventi possono essere rappresentati da sottoinsiemi dell'insieme $\{1, 2, 3, 4, 5, 6\}$ di tutti i possibili risultati; i due eventi menzionati, in particolare, dai sottoinsiemi $\{2, 4, 6\}$ e $\{1, 2\}$ rispettivamente, a cui sono assegnate le probabilità $1/2$ e $1/3$.

Possiamo ora ricordare la definizione di spazio di probabilità, associato a un esperimento. È dato un insieme X , l'insieme dei risultati possibili, e una famiglia \mathcal{A} di sottoinsiemi di X che rappresentano gli eventi. Nelle situazioni più generali, non è necessario che \mathcal{A} contenga tutti i sottoinsiemi di X , ma deve essere chiusa rispetto a certe operazioni; analizzando l'uso comune del termine evento nella lingua di ogni giorno, locuzioni come «evento opposto di un altro», «evento che risulta dalla combinazione di due eventi» e simili, si vede che la famiglia degli eventi deve soddisfare queste condizioni: deve contenere X , l'evento certo, e, nel caso che E e F siano due eventi, anche la riunione $E \cup F$, l'intersezione $E \cap F$ e il complemento $X - E$ devono essere eventi. Inoltre, per poter parlare anche della ripetizione indefinita di un esperimento, occorre che, se $\{E_n\}$ è una successione di

eventi, anche la riunione $\bigcup_{n=1}^{\infty} E_n$ sia un evento. Una famiglia di sottoinsiemi di X soddisfacente a queste condizioni si dice un σ -corpo di insiemi in X (questo nella letteratura algebrica; in teoria delle probabilità si parla anche di tribù di insiemi).

La probabilità degli eventi è ora stabilita da una funzione m , definita per tutti gli eventi, che assume valori reali compresi tra 0 e 1, e tale che $m(X) = 1$ e che se $\{E_n\}$ è una successione di eventi a due a due disgiunti si abbia

$$m\left(\bigcup_{n=1}^{\infty} E_n\right) = \sum_{n=1}^{\infty} m(E_n).$$

Quest'ultima condizione è detta la σ -additività.

Una variabile casuale su uno spazio di probabilità $\langle X, \mathcal{A}, m \rangle$ non è altro che una funzione da X in \mathbb{R} a proposito della quale abbia senso porre questioni di probabilità, e si possa decidere la probabilità che i valori della funzione cadano entro limiti prefissati. Per questo è sufficiente richiedere che, se f è la funzione, per ogni numero reale r l'insieme degli elementi x di X tali che $f(x) \leq r$, in simboli

$$\{x \in X : f(x) \leq r\}$$

sia un evento, e quindi abbia una probabilità associata.

Non possiamo discutere più a fondo la struttura degli spazi di probabilità e dei sistemi di variabili casuali associati. D'altra parte il discorso che segue dipende solo dalle proprietà generali di questi sistemi, comuni a tutti; potremmo dire che dipende solo dalle definizioni, come è tipico dei ragionamenti matematici. Le conclusioni a cui perverremo avranno allora anch'esse

un carattere molto generale, perché non dipenderanno dal particolare spazio di partenza. Dobbiamo solo avvertire che gli spazi più interessanti non sono quelli legati all'esperimento del dado, ma sono quelli infiniti. Il lettore può pensare che X sia l'insieme dei numeri reali compresi tra 0 e 1, e \mathcal{A} la famiglia dei cosiddetti insiemi di Borel: la famiglia che si ottiene partendo dagli intervalli contenuti in X e iterando le operazioni di riunione (anche di successioni di insiemi) e di complemento. La misura m in questo caso è la cosiddetta misura di Borel: gli intervalli hanno come misura la loro ampiezza, e gli altri insiemi di Borel una misura che si ottiene in modo naturale calcolandola a partire da quella degli intervalli che sono alla base della loro costruzione. In questo caso gli eventi che sono composti di un solo punto hanno misura zero, o nulla.

Dopo aver definito esattamente la nozione di variabile casuale, dobbiamo ancora approfondire la comprensione del loro uso. Come abbiamo già spiegato, quando diciamo che le variabili casuali sono numeri reali che dipendono dal caso intendiamo dire che in tutte le operazioni (logiche o di calcolo) che eseguiamo con i valori delle variabili casuali, ci dobbiamo portare dietro la valutazione della probabilità che i valori risultanti dall'esperimento siano quelli su cui operiamo. Ne segue che tutte le proposizioni concernenti un sistema di variabili casuali sono dei giudizi di probabilità, e che la dicotomia vero-falso mal si adatta in questo caso alla natura degli oggetti in questione.

Consideriamo un esempio semplice, per quanto artificioso: se f è la variabile connessa al lancio di un dado, che vale 0 se il risultato del lancio è 1, e 1 in tutti gli altri casi, mentre g è la variabile che vale 0 se il risultato è 6, e 1 in tutti gli altri casi, le due variabili sono diverse, ma questo giudizio è scarsamente informativo; maggiore informazione si ha dall'osservazione che c'è una buona probabilità che f e g siano uguali, e precisamente $2/3$, corrispondente all'evento $\{2, 3, 4, 5\}$.

Da questo esempio si vede come, attribuendo alle proposizioni concernenti un sistema di variabili casuali dei valori di probabilità, da una parte si è più fedeli allo spirito delle questioni che si trattano, dall'altra si arricchiscono le possibilità di giudizio: alcune proposizioni avranno valore 1 e saranno dette valide, altre avranno valore 0 e saranno dette false, altre infine avranno solo una certa probabilità po-

sitiva, che non è uguale alla certezza.

In queste osservazioni non c'è nulla di nuovo rispetto a quello che si fa nella pratica quando si calcolano le probabilità di eventi complessi. Dobbiamo solo descrivere come si fa ad assegnare queste probabilità, in modo sistematico, a tutte le proposizioni concernenti un dato sistema di variabili casuali e in modo da rispettare la correttezza dei ragionamenti. Se non ha senso dire che una affermazione sulle variabili casuali è vera o falsa, ha senso invece controllare la correttezza dei ragionamenti e dei calcoli; questo non significa che l'implicazione dalle premesse alla conclusione è vera, ma che la probabilità di una conclusione non deve essere inferiore a quella delle premesse. Si tenga presente poi che con questa valutazione probabilistica delle proposizioni dobbiamo anche dare un senso preciso e giustificare l'affermazione che le variabili casuali possono essere trattate come numeri reali.

Per la realizzazione di questo compito, il « mestiere » logico ci suggerisce alcune semplificazioni. Supponiamo di aver fissato uno spazio di probabilità $\langle X, \mathcal{A}, m \rangle$ e l'insieme \mathcal{R} delle variabili casuali su X . Se ci limitiamo a considerare le proposizioni sugli elementi di \mathcal{R} che sono analoghe alle proposizioni sui numeri reali, possiamo delimitare in modo preciso la classe di queste proposizioni. Ci sono innanzitutto le proposizioni elementari, o atomiche, del tipo $f = g$, $f \leq g$, $f + g = h$, $f \cdot g = h$, dove f, g, h sono elementi di \mathcal{R} . È chiaro che po-

tremmo considerare anche altre proposizioni elementari, ma il linguaggio in cui abbiamo formulato gli assiomi dei numeri reali contempla solo proposizioni elementari di questo tipo. Poi ci sono le proposizioni più complesse, che si ottengono da quelle più semplici per mezzo dei connettivi logici: « non », « oppure » ecc., e dei quantificatori: « esiste un x tale che... », « per tutti gli x ... » (si veda la figura a pagina 45).

Assegneremo innanzitutto un valore alle proposizioni elementari, e poi faremo vedere, con un ragionamento induttivo, come si possa assegnare un valore alle proposizioni complesse ammettendo di averlo già assegnato alle proposizioni più semplici. Per fare ciò, non è possibile assegnare direttamente dei semplici numeri, misure di probabilità, come valore delle proposizioni. È abbastanza spontaneo proporre, per esempio, che il valore della proposizione $f = g$ sia la misura dell'insieme $\{x \in X : f(x) = g(x)\}$ ma, se lavoriamo solo con i numeri, andiamo incontro a difficoltà quando passiamo alle proposizioni complesse. Si consideri questo caso: il valore della proposizione $f = g \& g = h$ deve essere minore o uguale a quello della proposizione $f = h$, perché $f = g \& g = h \rightarrow f = h$ è una legge logica dell'identità. Ma sia per esempio f la variabile, connessa al lancio di un dado, che vale 1 se il risultato è pari e 0 se il risultato è dispari, e g la variabile che vale 2 se il risultato è pari e 0 se il risultato è dispari. Se h è la variabile che vale 3 se

il risultato è pari e 0 se è dispari allora $f = g$, $g = h$ e $f = h$ hanno tutte probabilità $1/2$; se h è la variabile che vale 2 se il risultato è pari e 1 se il risultato è dispari, allora $f = g$ e $g = h$ hanno probabilità $1/2$, mentre $f = h$ ha probabilità 0.

Questo esempio suggerisce che nel valutare le proposizioni non si può tener conto soltanto della probabilità degli eventi che esse descrivono, ma bisogna prendere in considerazione gli eventi stessi. Lungi dall'essere una complicazione, questa necessità si rivela decisamente semplificatrice. La famiglia \mathcal{A} degli eventi infatti, con le sue operazioni insiemistiche di unione, intersezione e complemento, ha una struttura che è molto simile a quella dei sistemi logici astratti. Ha una struttura che è quasi un'algebra di Boole completa (si veda la figura a pagina 49). Diciamo quasi perché per ottenere proprio un'algebra di Boole occorrerebbe prima identificare tra loro gli eventi la cui differenza è un evento di misura nulla; ma per non interrompere il filo dell'esposizione aggiriamo questo ostacolo e supponiamo che la famiglia degli eventi sia un'algebra di Boole. Ora è ben noto ai logici che per poter assegnare dei valori alle proposizioni del discorso in modo da ottenere una valutazione sensata del ragionamento occorre e basta che tali valori siano strutturati in un'algebra di Boole. I due valori classici vero e falso costituiscono soltanto la più semplice delle algebre di Boole. Con un'algebra di questo tipo, dopo aver assegnato dei

UN SISTEMA DI ASSIOMI

$$\forall x \forall y (x + y = y + x)$$

$$\forall x (x + 0 = x)$$

$$\forall x \exists y (x + y = 0)$$

I simboli $+$ e 0 non vanno intesi come la ordinaria addizione tra numeri e il numero 0, ma come simboli linguistici per una operazione binaria e un elemento neutro rispetto alla operazione.

DUE MODELLI

$$\{-1, 0, 1\}$$

In questo insieme di numeri interi si consideri l'addizione ordinaria e lo zero svolga il ruolo di elemento neutro. I tre assiomi sono validi.

$$\left\{ \frac{1}{2}, 1, 2 \right\}$$

In questo insieme di numeri razionali, si consideri la moltiplicazione ordinaria e 1 svolga il ruolo di elemento neutro. I tre assiomi sono validi.

I due modelli sono isomorfi. Questo significa che esiste una corrispondenza biunivoca f tra i due insiemi, che conserva le operazioni. Più esattamente: 1) a ogni elemento di un sistema corrisponde uno e un solo elemento dell'altro sistema, mediante la f , e a due elementi distinti corrispondono elementi distinti; tutti gli elementi del secondo sistema corrispondono a qualche elemento del primo: $f(-1) = \frac{1}{2}$, $f(0) = 1$, $f(1) =$

$= 2$; 2) all'elemento neutro corrisponde l'elemento neutro $f(0) = 1$; 3) per ogni x, y e z , abbiamo $x + y = z$ nel primo sistema se e solo $(f(x) \cdot f(y)) = f(z)$ nel secondo, e viceversa. Quando due strutture sono isomorfe, esse hanno in particolare la stessa cardinalità, anche se sono infinite; diciamo che due insiemi hanno la stessa cardinalità nel caso appunto che esista tra di essi una corrispondenza biunivoca, nel senso della condizione 1.

valori alle proposizioni elementari, si procederà così: se φ e ψ indicano delle proposizioni di cui già conosciamo il valore, che scriveremo con $[\varphi]$ e $[\psi]$, i valori delle proposizioni composte mediante φ e ψ si calcoleranno secondo questo schema

$$\begin{aligned} [\varphi \& \psi] &= [\varphi] \cap [\psi] \\ [\varphi \vee \psi] &= [\varphi] \cup [\psi] \\ [\sim \varphi] &= [\varphi]' \\ [\exists x \varphi(x)] &= \bigcup_{x \in X} [\varphi(x)] \end{aligned}$$

Ci si può limitare a questo tipo di proposizioni, perché gli altri connettivi possono essere eliminati, definendoli per mezzo di quelli elencati. Seguendo questo criterio di valutazione è noto, e si può verificare facilmente, che tutti gli assiomi della logica sono validi, hanno come valore l'1 dell'algebra di Boole, e che le conclusioni di un ragionamento hanno valore maggiore o uguale a quello delle premesse.

Come esempio si consideri il principio di non contraddizione: $\sim(\varphi \& \sim \varphi)$. Abbiamo:

$$[\sim(\varphi \& \sim \varphi)] = ([\varphi] \cap [\varphi])' = 0' = 1$$

Per raggiungere l'obiettivo che ci eravamo proposti, non resta allora che precisare il valore da assegnare alle proposizioni elementari, e verificare che anche gli assiomi dell'identità e gli assiomi propri dei corpi ordinati risultano validi. Dopo quanto abbiamo anticipato, non v'è dubbio che il valore assegnato alle proposizioni elementari sarà l'evento stesso che essi descrivono:

$$\begin{aligned} [f = g] &= \{x \in X : f(x) = g(x)\} \\ [f + g = h] &= \{x \in X : f(x) + g(x) = h(x)\} \end{aligned}$$

e così via.

Il sistema di valutazione adottato può sembrare un po' astratto dal momento che il valore di una proposizione risulta essere l'evento da essa descritto, ma è sempre possibile, in un secondo tempo, il passaggio dagli eventi alle loro probabilità, per cui è legittimo chiamare questa una valutazione probabilistica delle proposizioni.

Ora possiamo verificare che, in senso probabilistico, le variabili casuali costituiscono un corpo ordinato. Esaminiamo solo i due assiomi in cui maggiormente risalta la differenza tra la valutazione classica delle proposizioni e quella probabilistica. Perché un sistema sia totalmente ordinato occorre che, presi due elementi qualunque f e g , si abbia $f \leq g$ o $g \leq f$. La proposizione $f \leq g \vee g \leq f$ è valida perché $[f \leq g \vee g \leq f] = \{x \in X : f(x) \leq g(x)\} \cup \{x \in X : g(x) \leq f(x)\} = X$ anche se non è vero che per tutti gli x si ha $f(x) \leq g(x)$ o che per tutti gli x si ha $g(x) \leq f(x)$.

Un altro assioma cruciale è quello che afferma l'esistenza, per ogni $f \neq 0$, di un elemento g tale che $f \cdot g = 1$. Deve perciò risultare valida la proposizione $f \neq 0 \rightarrow \exists g (f \cdot g = 1)$. Ricordando che l'implicazione si può esprimere per mezzo della negazione e della disgiunzione ($\varphi \rightarrow \psi$ equivale a $\sim \varphi \vee \psi$), si vede che è sufficiente costruire, per ogni f , una g tale che

$$[f \neq 0] \subseteq [f \cdot g = 1]$$

ovvero

$$\{x \in X : f(x) \neq 0\} \subseteq \{x \in X : f(x) \cdot g(x) = 1\}.$$

Ebbene, basta definire g in questo modo:

$$\text{se } f(x) \neq 0, \quad g(x) = \frac{1}{f(x)}$$

se $f(x) = 0, \quad g(x) = 0$, o un altro valore qualsiasi.

Da questa definizione può sembrare che il reciproco di f non sia determinato in modo unico, ma in senso probabilistico così non è. È facile controllare con qualche calcolo che nella misura in cui due funzioni sono il reciproco di una stessa f esse sono anche uguali tra loro, ovvero che la probabilità che siano uguali è maggiore o uguale a quella che siano il reciproco della stessa funzione. In questo senso diventa anche valida l'affermazione che il reciproco di un elemento è unico.

Questi due esempi dovrebbero bastare a dare un'idea di come si debba ragionare nello spirito probabilistico, e quale significato si debba attribuire alla validità di proposizioni complesse. Resta il fatto che, qualunque sia lo spazio di probabilità di partenza, gli assiomi dei corpi ordinati, interpretati come abbiamo visto sul sistema delle variabili casuali, sono validi, e in questo modo viene a essere giustificata rigorosamente l'affermazione che le variabili casuali possono essere usate come se fossero numeri reali.

Ricordiamo ora che noi cercavamo un modello dei numeri reali che mettesse in luce i limiti del teorema di unicità, ma la garanzia della unicità era data dall'assioma di continuità, di cui non abbiamo ancora detto nulla. Solo per formularlo, dobbiamo arricchire il linguaggio usato finora, e il sistema di oggetti che costituiscono il modello. Nell'assioma si parla di insiemi di numeri, per cui dobbiamo considerare proposizioni della forma « f appartiene all'insieme S » (« $f \in S$ »); ma quale entità mettiamo al posto di S , a svolgere il ruolo di sottoinsieme di \mathcal{R} ? Si vede subito che non si possono prendere i veri sottoinsiemi di \mathcal{R} , perché altrimenti non si riesce a dare una

valutazione a « $f \in S$ » secondo lo schema probabilistico di ricondursi ai valori $f(x)$ della funzione. I conti tornano invece se S è proprio un sottoinsieme di \mathcal{R} ; così mentre i nuovi numeri reali sono ben diversi da quelli classici, gli insiemi di numeri reali sono proprio i sottoinsiemi di \mathcal{R} . Potremo allora definire

$$[f \in S] = \{x \in X : f(x) \in S\}$$

e ottenere una valutazione soddisfacente di tutte le proposizioni che parlano di insiemi di numeri reali. Si può aggiungere che al modello non si possono associare tutti i sottoinsiemi di \mathcal{R} , perché occorre naturalmente, come risulta dalla definizione, che per ogni f

$$\{x \in X : f(x) \in S\}$$

sia un evento. Ci si può restringere allora ai cosiddetti insiemi di Borel, ma non è il caso di scendere nei particolari. L'importante è sapere che esistono insiemi con questa proprietà.

Il sistema delle variabili casuali \mathcal{R} con gli insiemi di Borel forma un corpo ordinato continuo in senso probabilistico. La dimostrazione che l'assioma di continuità è valido è naturalmente il punto cruciale di tutta la costruzione, e premia la profonda intuizione dell'ideatore del metodo, Dana Scott dell'Università di Oxford. Purtroppo i dettagli della dimostrazione sono un po' laboriosi e non possono essere dati; ricordiamo soltanto, per il lettore che volesse ricostruirseli, il passo essenziale: ammesso che un insieme S sia limitato superiormente, l'estremo superiore deve essere un elemento g tale che la proposizione «per ogni $r, g \leq r$ se e solo se $\forall f (f \in S \rightarrow f \leq r)$ » sia valida. Una variabile casuale g siffatta si può definire nel seguente modo: per ogni $x \in X, g(x) = \inf \{q \in \mathcal{Q} : \forall f (f(x) \in S \rightarrow f(x) \leq q)\}$.

Il modo in cui è definita la funzione g , tenendo conto della proposizione che si vuole risulti valida, è tipico delle manipolazioni con i modelli probabilistici.

Ricordiamo ancora che il modello che abbiamo costruito può essere ulteriormente arricchito con l'introduzione di funzioni, funzionali e così via, in modo da risultare una struttura adeguata alle necessità dell'analisi matematica, e addirittura un universo probabilistico della teoria degli insiemi.

Ma vediamo infine di trarre alcune conclusioni dall'esistenza di questi modelli non classici delle teorie matematiche.

1) Una prima considerazione è di carattere logico, e si riallaccia al pro-

<p>1. PROPRIETÀ ASSOCIATIVA</p> $\forall x \forall y \forall z (x \cup (y \cup z)) = ((x \cup y) \cup z)$ $\forall x \forall y \forall z (x \cap (y \cap z)) = ((x \cap y) \cap z)$
<p>2. PROPRIETÀ COMMUTATIVA</p> $\forall x \forall y (x \cup y = y \cup x) ; \forall x \forall y (x \cap y = y \cap x)$
<p>3. PROPRIETÀ DISTRIBUTIVA</p> $\forall x \forall y \forall z [x \cup (y \cap z) = (x \cup y) \cap (x \cup z)]$ $\forall x \forall y \forall z [x \cap (y \cup z) = (x \cap y) \cup (x \cap z)]$
<p>4. PROPRIETÀ DI IDEMPOTENZA</p> $\forall x (x \cup x = x) ; \forall x (x \cap x = x)$
<p>5. PROPRIETÀ DI ASSORBIMENTO</p> $\forall x \forall y [x \cup (x \cap y) = x] ; \forall x \forall y [x \cap (x \cup y) = x]$
<p>6. ESISTENZA DEGLI ELEMENTI NEUTRI</p> $\forall x (x \cup 0 = x) ; \forall x (x \cap 1 = x) ; 0 \neq 1$
<p>7. PROPRIETÀ DEL COMPLEMENTO</p> $\forall x (x \cup x' = 1) ; \forall x (x \cap x' = 0)$

Un'algebra di Boole è un sistema $\langle B, \cup, \cap, ', 0, 1 \rangle$ che soddisfa gli assiomi elencati. I simboli sono quelli usati comunemente per indicare le operazioni insiemistiche, perché l'esempio tipico di algebra di Boole è costituito dalla famiglia dei sottoinsiemi di un insieme fissato X , con le solite operazioni insiemistiche. In tal caso 1 è l'insieme X stesso, e 0 l'insieme vuoto \emptyset . In un'algebra di Boole si può definire una relazione di «minore o uguale» nel seguente modo: $x \subseteq y$ se e solo se $x \cup y = y$. Anche questa corrisponde, nell'esempio canonico, alla relazione di inclusione tra insiemi. Se S è un insieme di

elementi dell'algebra di Boole, si indica con US quell'elemento, se esiste, che è maggiore o uguale a tutti gli elementi di S e che è il più piccolo che gode di questa proprietà. Un'algebra di Boole si dice completa se US esiste per ogni insieme S . La più semplice algebra è quella costituita da soli due elementi, 0 e 1 ; le leggi delle operazioni booleane tra questi due elementi sono le stesse che regolano i valori di verità falso e vero delle proposizioni composte con la disgiunzione, la congiunzione e la negazione. Il nome delle algebre deriva dal matematico e logico inglese George Boole (vissuto dal 1815 al 1864).

blema da cui siamo partiti, dell'unicità del sistema dei numeri reali. È chiaro che, scegliendo opportunamente lo spazio di probabilità, da quelli finiti a quelli infiniti di cardinalità arbitraria, si ottengono modelli profondamente diversi, e infatti si è potuto dimostrare che certe proposizioni, come l'ipotesi del continuo, sono valide nell'uno ma non nell'altro. Queste conclusioni però non dipendono dalla particolare tecnica usata per costruire i modelli, e infatti quella originaria di Cohen era diversa. Ma l'esistenza stessa di modelli non standard, di cui ne abbiamo esibito uno, ha questo significato: quando dimostriamo i teoremi di unicità prendiamo in considerazione solo alcuni dei possibili modelli degli assiomi e ne scartiamo altri; quelli che esaminiamo risultano isomorfi (anche per il sistema dei numeri naturali si ripe-

te tale situazione). Questa operazione è fatta in modo che non è possibile rilevare alcun imbroglio, in modo intuitivamente chiaro e corretto; ma per l'appunto solo intuitivamente, mentre non è possibile giustificarla logicamente in modo rigoroso. Il difetto può essere sia nell'intuizione sia nella logica.

2) Dal punto di vista matematico, la stretta connessione che abbiamo stabilito tra i numeri reali e le variabili casuali per mezzo dei modelli probabilistici può rivelarsi feconda, permettendo un interscambio di metodi e di risultati.

Se, per esempio, formuliamo la continuità per mezzo delle successioni di Cauchy (come all'inizio della esposizione), per dimostrare la validità del principio che le successioni di Cauchy sono convergenti, si utilizza un teorema

sulle funzioni (i nuovi numeri reali sono funzioni!) che ci limitiamo a enunciare e che dice: se una successione di funzioni misurabili è fondamentale in misura, allora converge in misura a una funzione misurabile.

Viceversa, una volta stabilito che tutti i teoremi sui numeri reali sono validi, si possono utilizzare dei teoremi sulla convergenza delle successioni di numeri reali per stabilire proprietà delle variabili casuali del tipo di quelle delle leggi dei grandi numeri.

Tutti i risultati che arricchiscono la nostra conoscenza delle variabili casuali sono importanti, in relazione al ruolo essenziale che queste ricoprono in molti campi della scienza; se queste tecniche daranno dei risultati, al di fuori della problematica logica da cui sono nate, è un problema che merita di essere approfondito.

Tre personaggi della matematica

Perché i numeri e , i , π s'incontrano in matematica strettamente uniti? La «vera» spiegazione è semplice. Per e basta pensare alla capitalizzazione, per i alle rotazioni del piano; π scaturisce dal significato dei primi due

di Bruno de Finetti

Per poco che uno prosegua negli studi, incontrerà il numero (pi-greco) $\pi = 3,1415926...$ come rapporto della circonferenza al diametro, il numero $e = 2,718281828...$ come base dei logaritmi naturali, il numero $i = \sqrt{-1}$ come « unità immaginaria » (ossia, a prima vista, come un trucco per inventare radici fittizie quando un'equazione ne manca). Se va avanti, si accorgerà però sempre più che questi tre numeri si presentano ovunque come personaggi essenziali nel ragionamento matematico, e — ciò che forse è più sorprendente — si presentano congiuntamente, strettamente collegati tra loro, mentre le consuete definizioni non ne lasciano neppure intravedere il « perché ».

Eppure la « vera » spiegazione è semplice e suscettibile di venir illustrata in modo elementare e intuitivo. L'itinerario predisposto per incontrare i nostri tre personaggi (e , poi i , infine π), lungi dall'apparire come un percorso di per sé privo di senso, di interesse e di giustificazioni, gioverà — spero — anche a formare un'efficace comprensione e una veduta panoramica, sostanzialmente organica benché non pedissequamente scolastica, su vaste ed essenziali aree del pensiero matematico. Non si richiederà al lettore alcuna conoscenza propedeutica specifica: tutte le nozioni che servono scaturiranno dallo stesso svolgimento della trattazione.

Naturalmente, dato che la trattazione non può non essere alquanto succinta e dovrà introdurre concetti che si suppongono nuovi al lettore, occorrerà richiederli una certa attenzione e una disposizione d'animo non diffidente e ostile verso la matematica. Penso del resto che tale diffusissimo stato d'animo sia dovuto non alla pretesa difficoltà o aridità della matematica, bensì ai metodi tradizionali del suo insegna-

mento che insistono sul piatto formalismo anziché scatenare la fantasia e la intelligenza creativa. Il difetto sta — per usare la felice immagine di Jerome S. Bruner — nel trascurare gli aspetti di pertinenza della « mano sinistra » (fantasia, intuizione, arte, gioco) disseccando tutto nella mera sistemazione razionale (di pertinenza della « mano destra »). Ed è notevole merito del Bruner (psicologo, non matematico) di aver ravvisato nella matematica, e nell'insegnamento e apprendimento della matematica, il campo ove l'intervento della « mano sinistra » sarebbe particolarmente vivificatore ed è invece purtroppo pressoché totalmente assente.

La presente esposizione è stata scritta (non so se e quanto riuscitamente) con la mano sinistra.

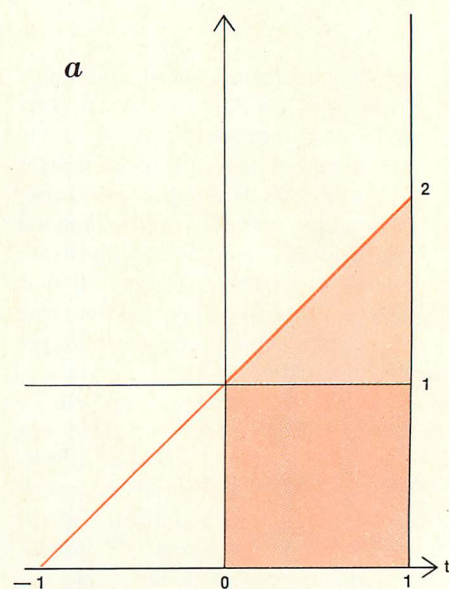
Il significato di e

Volendo scegliere la forma più immediatamente e concretamente intuitiva, possiamo riferirci all'esemplificazione finanziaria e dire: nel tempo in cui un capitale, con l'interesse semplice, si raddoppia, se anche gli interessi fruttano immediatamente interessi, risulterà moltiplicato per e (anziché per 2). Ossia, un milione, anziché due milioni, diviene 2 718 281.

Beninteso, la stessa formulazione si può applicare per qualsiasi altro fenomeno. Per esempio, una popolazione di un milione di abitanti che aumentasse di 55 individui al giorno si raddoppierebbe in 50 anni; ma se rimane costante non il numero (55 individui) bensì l'intensità (55 per ogni milione al giorno, ossia 0,055 % al giorno), dopo 50 anni la popolazione sarà non di due milioni ma di 2 718 281. È chiaro (notiamo incidentalmente) che una formulazione così rigida snatura, per semplicità di espressione, il carattere meramen-

te statistico, probabilistico, del fenomeno. È soltanto « molto probabile » che il risultato sarà « circa » quello indicato (come per le frequenze di un gioco d'azzardo); questo va sottinteso in ogni caso del genere.

L'esempio che segue serve per accennare subito anche a una formulazione analoga che si ha quando si consideri il caso opposto, cioè una diminuzione. L'esempio più idoneo è quello di una sostanza radioattiva di cui in un anno si disintegri, per esempio, un atomo ogni mille. Avendo un miliardo di atomi, il primo anno se ne disintegrerebbero un milione e, se tale numero rimanesse costante, dopo mille anni gli atomi sopravvissuti si ridurrebbero a zero; ma se invece (come è chiaro) rimane costante l'intensità del fenomeno



Il numero e può essere introdotto intuitivamente usando l'esemplificazione finanziaria. Nella capitalizzazione semplice (a) il montante cresce di quantità uguali in

(1 ‰ all'anno, cioè un atomo all'anno su ogni 1000 *sopravvissuti*). Gli atomi superstiti dopo 1000 anni sarebbero circa 368 milioni (cioè ne rimarrebbero nella proporzione $1/e = e^{-1} = 0,367879\dots$). Ciò risulterà evidente anche dal significato geometrico (*si veda la figura in queste due pagine*), ma è preferibile attendere di giungervi attraverso l'itinerario che seguiremo. Si tratterà – per anticiparne il senso in due parole – di vedere come dall'andamento rettilineo del caso dell'interesse semplice si passi alla curva esponenziale del caso dell'interesse continuo. È la curva che, a volte, è detta, soprattutto dai biologi, « curva dell'accrescimento naturale » (attenzione però a non mitizzarla!). Vi giungeremo seguendo due diversi approcci, o procedimenti.

Un primo procedimento, consistente nell'applicare l'interesse semplice in intervalli sempre più piccoli, conduce ad approssimare la curva esponenziale mediante spezzate con ordinate in progressione geometrica. Il secondo, consistente nell'aggiungere successivamente, al capitale, l'interesse (semplice), l'interesse dell'interesse, l'interesse dell'interesse dell'interesse, e via dicendo, conduce ad approssimare la stessa curva esponenziale mediante polinomi (e precisamente mediante lo sviluppo noto col nome di serie di potenze).

La capitalizzazione, semplice e continua

Per sviluppare la trattazione ci riferiremo sempre alla interpretazione finanziaria che appare senz'altro la più direttamente espressiva (ma quanto vien

detto vale per ogni altra esemplificazione di « accrescimento naturale »). Cominciamo con un cenno sulla capitalizzazione semplice: cosa significa? Come si rappresenta graficamente? Significa che un capitale produce un interesse, sempre lo stesso per intervalli di tempo uguali, e che (per usare la finzione abituale ma efficace) l'interesse viene accumulato in un conto a parte, che non produce interessi. Perciò il valore complessivo (capitale più interessi, che si dice montante) cresce di quantità uguali in tempi uguali, ossia, graficamente, cresce secondo una retta. Tale andamento è illustrato nella figura in queste due pagine con riferimento alla schematizzazione che ci sarà comoda: capitale « uno », andamento degli interessi e del montante nell'intervallo in cui il montante si raddoppia (che scegliamo convenzionalmente come unità di tempo, indicandone con $t = 0$ e $t = 1$ gli istanti iniziale e finale).

Osserviamo che, se l'interesse fosse negativo, la retta terminerebbe col valore zero al tempo 1; anche senza disegnare la retta decrescente basta pensare di percorrerla in senso opposto (da destra a sinistra), o, se si preferisce, guardarla per trasparenza (da dietro il foglio).

Non c'è nessuna ragione, tuttavia, perché del denaro lasciato in deposito (si chiami pure « interesse » o come altro si voglia) non debba dar diritto a interessi. Perciò la capitalizzazione semplice non può venir considerata che una semplificazione di calcolo, comoda sì, ma tollerabile soltanto finché l'interesse relegato in un « conto infrut-

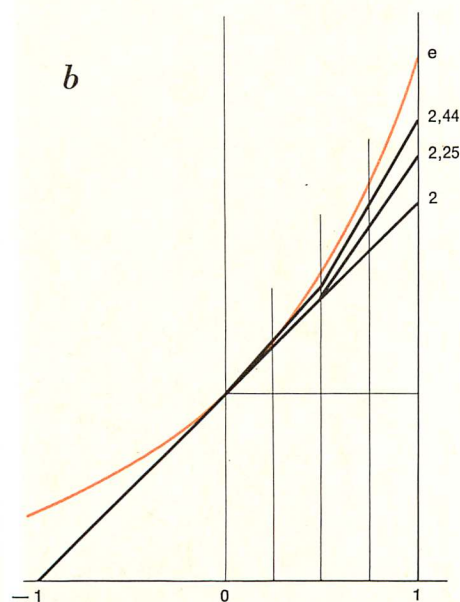
tifero » rimane praticamente trascurabile rispetto al « capitale ». Il modo logico ed esatto di procedere consiste nel far affluire senz'altro l'interesse nel conto stesso del capitale rendendolo fruttifero istantaneamente. È questa la capitalizzazione continua che, come detto, conduce a definire il numero e .

Ma converrà, come vedremo, considerare due diversi modi di avvicinarci a tale caso ideale mediante successive approssimazioni. Entrambi sono illuminanti per l'interpretazione concreta che forniscono e per la facilità con cui conducono – praticamente senza alcun calcolo e senza alcun prerequisito teorico – alle conclusioni che interessano e occorrono ai fini del nostro discorso. In sostanza, si tratterà anzitutto di dimostrare che il montante finale salirà a e anziché a 2 seguendo, anziché la retta, la curva esponenziale.

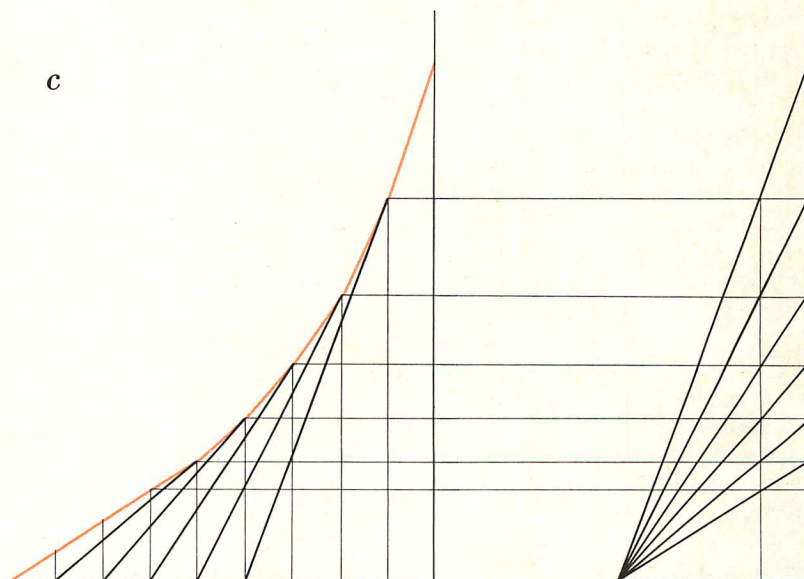
Primo approccio

Si può pensare che gli interessi vengano aggiunti al capitale, se non istantaneamente, alla fine di periodi fissi più o meno brevi. Nella pratica ciò avviene in genere alla fine di ogni anno, o semestre, o trimestre, ecc. Nella nostra schematizzazione considereremo il periodo totale « uno » diviso in due metà, quattro quarti, otto ottavi, e così di seguito con successivi dimezzamenti.

Dividiamo l'intervallo a metà. Alla fine della prima frazione (istante $t = 1/2$) il capitale 1 ha fruttato l'interesse $1/2$ portando quindi il montante a $1 + 1/2 = 3/2 = 1,50$; nella seconda frazione l'interesse, prodotto dal mon-



tempi uguali, cioè secondo una retta. Nella capitalizzazione composta (b) gli interessi risultano produttivi durante le frazioni del periodo preso in considerazione e quindi la retta si trasforma in una spezzata con un numero di lati che aumenta all'au-



mentare del numero di frazioni di periodo considerate; al limite si ha la curva esponenziale (in colore). La costruzione di ordinate equidistanti in progressione geometrica (c) mostra come, sovrapponendo le strisce, le congiungenti concorrano in un punto.

tante $3/2$ anziché dal capitale 1 , non è più $1/2$ bensì $3/4$ (metà di $3/2$ anziché metà di 1). Il montante finale è allora $1 + 1/2 + 3/4 = 9/4 = 2,25$; lo prova anche un ragionamento più diretto in quanto il montante si moltiplica per $3/2$ nella prima frazione e si rimoltiplica per $3/2$ nella seconda e quindi in definitiva si moltiplica per $(3/2)^2 = 9/4$.

È chiaro che, per lo stesso motivo, se dividiamo l'intervallo in n frazioni di durata $1/n$ alla fine delle quali l'interesse ivi prodottosi diviene fruttifero, il montante finale sarà $(1 + 1/n)^n$. Dopo la prima frazione, infatti, il montante sarà $1 + 1/n$ (1 = capitale, $1/n$ interesse di esso per il tempo $1/n$), e lo stesso accrescimento si ripete n volte. È anche intuitivo che il montante finale risulterà tanto maggiore quanto maggiore si prende n (ossia quanto più brevi sono le frazioni di periodo, $1/n$, in cui l'interesse resta infruttifero), e che tale miglioramento diviene trascurabile quando n è grande perché allora si è già, praticamente, nella situazione di capitalizzazione continua.

È ovvio pertanto che il valore di e si può approssimare (per difetto) calcolando $(1 + 1/n)^n$ per n grande; d'altronde risulta (dal medesimo ragionamento) che è invece decrescente al crescere di n il valore $(1 + 1/n)^{n+1}$ (cioè il montante che si otterrebbe attendendo una frazione di periodo in più: per esempio $5/4$ anziché $4/4$, $9/8$ anziché $8/8$, $101/100$ anziché $100/100$, ecc.) e abbiamo così una doppia disuguaglianza che consente di accertare il grado di approssimazione:

$$(1 + 1/n)^n < e < (1 + 1/n)^{n+1}.$$

Queste disuguaglianze mostrano che, per $n = 100$, e vale circa $2,72$ mentre per $n = 100\,000$ ne rimane fissato il valore con quattro decimali esatti. Possiamo quindi scrivere $e \approx (1 + 1/n)^n$, cioè che e è « praticamente uguale a $(1 + 1/n)^n$ quando n è molto grande ». La dicitura non è rigorosa ma fornisce, penso, l'idea più chiara e sostanzialmente più esatta a chi non possiede con sicurezza i concetti e la terminologia dell'analisi matematica.

La curva esponenziale e le sue proprietà

Osserviamo ancora la figura in queste due pagine; si vede non solo come, infittendo le suddivisioni dell'intervallo, cresca il valore finale fino a toccare il livello e , ma anche come le successive spezzate (di 2 , 4 , 8 ,... lati) tendano a confondersi con una curva.

È la curva detta esponenziale, diagramma dei fenomeni che, come nella capitalizzazione continua e come negli altri esempi accennati (sviluppo di una popolazione, disintegrazione di una so-

stanza radioattiva), avvengono con intensità costante di accrescimento (o di diminuzione). In altri termini: in tempi uguali si hanno accrescimenti (o decrementi) percentualmente uguali, cosicché i valori assunti in una successione di istanti equidistanti (per esempio ogni anno, ogni ora, ogni secondo) variano in progressione geometrica, cioè ciascuno si ottiene dal precedente moltiplicandolo per una costante a (positiva), maggiore o minore di 1 a seconda che la funzione è crescente o decrescente.

Nel caso considerato (cioè se l'unità di tempo è quella in cui l'incremento di intensità costante la fa moltiplicare per e) il valore di a è appunto e . Ciò significa che la tangente alla curva, in corrispondenza all'origine, ha pendenza di 45° (come la bisettrice degli assi).

L'equazione di una tale curva si scrive $y = Ka^t$, dove K è il valore di y per $t = 0$ e a la detta costante (o base). Infatti, per $t = 1, 2, 3, \dots, n, \dots$ i valori sono $Ka, Ka^2, Ka^3, \dots, Ka^n, \dots$ (formanti, come detto, una progressione geometrica), e così vale per $t = -1, -2$, ecc. (ricordando che a^{-n} significa $1/a^n$); ciò vale anche interpolando i valori $t = -3/2, -1/2, 1/2, 3/2, 5/2$, ecc., tenendo presente che $a^{1/2}$ significa \sqrt{a} , e in generale $a^{h/k} = \sqrt[k]{a^h}$ (e per continuità risulta il significato di a^t per t qualunque).

Le proprietà della curva esponenziale, in parte già viste basandosi sul significato delle potenze $a^{h/k}$ (peraltro spesso presentate e imparate come si trattasse di una convenzione), meritano però di venir approfondite direttamente, con considerazioni sia di natura geometrica che analitica. Esse mostreranno anzi che sarebbe probabilmente preferibile definire a^t partendo dalla funzione esponenziale e da ciò dedurre il significato di x^c , in particolare per $c = 1/2$, $c = 1/n$, $c = m/n$, $c = -1$, $c = -n$, $c = -m/n$.

Geometricamente, la proprietà caratteristica di una successione di ordinate equidistanti y_m in progressione geometrica è che, prolungando i lati della spezzata, essi tagliano l'asse t a una stessa distanza dal piede della rispettiva ordinata (si veda la figura nelle due pagine precedenti).

Per la curva esponenziale la stessa proprietà (caratteristica di tale curva) è, per chiara analogia, la costanza della sottotangente (cioè del segmento dell'asse t intercettato dalla verticale e dalla tangente per un punto della curva). La sottotangente ha anche un significato fondamentale nell'interpretazione pratica: è il tempo in cui la funzione si raddoppierebbe (oppure si annullerebbe) se continuasse a crescere (o a diminuire) conservando la stessa pendenza che

ha nel punto considerato (ossia: proseguendo secondo la tangente). Il reciproco della sottotangente è pertanto la misura della intensità di accrescimento (se a sinistra; se a destra di diminuzione).

La sottotangente è uguale a 1 quando la base è e ; la definizione data di e , interpretata geometricamente, significa appunto tale fatto (ossia, come già detto, che la tangente in $t = 0$ ha pendenza 45° ; precisamente, è la retta che taglia l'asse t in $t = -1$ e l'asse y in $y = 1$). Al variare della base a , la curva non varia se non per un'alterazione proporzionale delle ascisse.

La proprietà geometrica della successione di ordinate equidistanti in progressione geometrica, oltre che a prolungarla indefinitamente verso sinistra e verso destra, permette anche di completare la curva inserendo punti intermedi a volontà; basta interpretare geometricamente il significato di $\sqrt[n]{a}$ (ossia $a^{1/n}$), che è, nella usuale terminologia scolastica, il « medio proporzionale » tra 1 e a : infatti è x tale che $1:x = x:a$, $x^2 = a$; in forma solo apparentemente più generale, partendo da due successive ordinate y_n e y_{n+1} si costruisce in tal modo quella intermedia, $y_{n+1/2}$ (e , ripetendo la costruzione, quelle a distanza ridotta a $1/4$, $1/8$, ecc.).

Secondo approccio

Abbiamo visto come il primo approccio ci abbia condotto a determinare, prima, il valore di e , e poi l'andamento della funzione e^t , ossia del montante nella capitalizzazione continua.

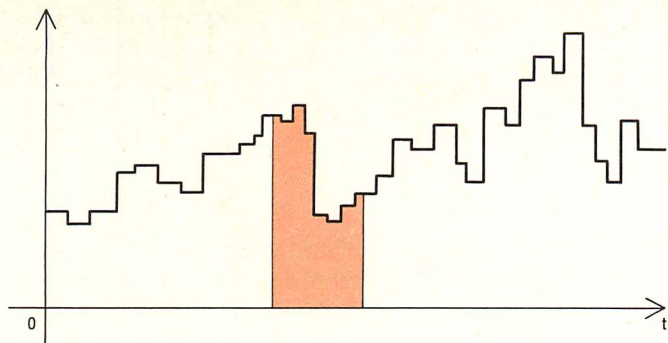
Indichiamo ora, rapidamente, alcune proprietà dell'esponenziale che è importante conoscere e che ci saranno utili nel seguito.

Si ha $e^{x+y} = e^x e^y$ e, in generale, $a^{x+y} = a^x a^y$. In altri termini, se qualcosa cresce moltiplicandosi per e^x in un tempo x e per e^y in un tempo y , nel periodo totale $x + y$ si moltiplica per $e^x e^y$. [E inversamente: se una funzione f è tale che $f(x + y) = f(x)f(y)$ non può che essere un'esponenziale; $f(x) = a^x$; quella considerata è una equazione funzionale].

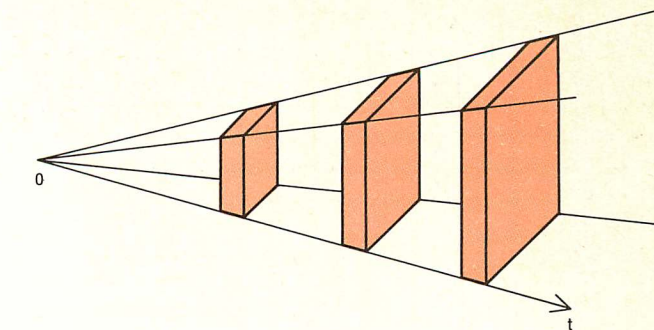
È vero per ogni t (come, per definizione, quando $t = 1$) che

$$(1 + t/n)^n = \left[\left(1 + \frac{1}{n/t} \right)^{n/t} \right]^t \approx e^t.$$

Nel secondo approccio usiamo ancora la stessa interpretazione finanziaria, sempre partendo dalla capitalizzazione semplice, ma seguendo un concetto diverso per approssimarci via via alla capitalizzazione continua. Lasciamo gli interessi (semplici) in un conto a parte; supponendolo fruttifero esso darà



Il numero e può essere introdotto anche con la capitalizzazione continua. A sinistra abbiamo il diagramma di un capitale che varia in modo qualunque e l'interesse semplice da esso prodotto in un dato intervallo di tempo è dato dal rettangoli-



no (in colore) a esso corrispondente. La figura a destra mostra invece come si può giungere intuitivamente a spiegare i coefficienti dello sviluppo in serie di e^t : il volume di una piramide è $1/6$ di quello del prisma di uguale base e altezza.

gli « interessi degli interessi » che collocheremo in un terzo conto; questo darà gli « interessi degli interessi degli interessi », e così proseguiremo. Avremo infinite aggiunte correttive sempre più trascurabili e rioterremo, al limite, la funzione esponenziale e^t .

Il capitale è 1 e l'interesse semplice (al tempo t) è t ; fin qui nulla di nuovo. Ma calcoliamo ora l'interesse dell'interesse; come vedremo, esso è $t^2/2$, e la spiegazione di tale calcolo indica il ragionamento che si ripeterebbe per tutti i termini successivi.

Pensiamo un momento al caso generale di un capitale che varia in modo qualunque (per esempio a un conto corrente, con continui accrediti e addebiti), e consideriamone il diagramma (si veda la figura in questa pagina). L'interesse (semplice) da esso prodotto in un dato intervallo di tempo è dato dall'area corrispondente perché in ogni intervallino esso è dato dall'area del rettangolino (altezza = capitale, larghezza = durata; il tasso, per la convenzione iniziale, è 1). Nel caso del capitale costante 1, risulta che (come sapevamo da sempre) l'interesse è t : l'area del rettangolo di altezza 1 e base t ; considerando l'interesse di tale interesse, tra 0 e t (ove esso è $y = t$), avremo che è l'area del triangolo rettangolo di base t e altezza t , cioè $t^2/2$, come asserito.

Per mostrare che il termine seguente è $t^3/6$, osserviamo che questo è il volume di un tetraedro rettangolo, cioè (si veda la figura in questa pagina) di un sesto del cubo di lato t . Pensandolo tagliato a fettine perpendicolari all'asse t , ciascuna di esse (per $t = \tau$) ha area $\tau^2/2 =$ « interesse dell'interesse all'istante τ » e volume pari all'interesse prodotto da esso nel tempuscolo corrispondente allo spessore della fettina. Perciò il volume complessivo di tutte le fettine, cioè del tetraedro, è l'interesse dell'interesse dell'interesse, come si voleva dimostrare.

Il coefficiente $1/6$ si può spiegare

(sia pure a titolo di aiuto mnemonico apparentemente superficiale, ma il fondamento è valido!) osservando che l'area del triangolo è base per metà dell'altezza, il volume del tetraedro è base (area della) per un terzo dell'altezza, e un sesto è il prodotto di un mezzo per un terzo (ossia: $6 = 2 \times 3 = 3!$). Pur mancando un'immagine nello spazio fisico, possiamo pensare di proseguire il ragionamento per figure analoghe in 4 dimensioni, in 5, in generale in n : anch'esse si ottengono dividendo in $n!$ parti (24 per $n = 4$ dimensioni, 120 per 5, ecc.) il cubo a n dimensioni, e il loro volume (n -dimensionale) è base [volume ($n-1$)-dimensionale] per un n -esimo dell'altezza, cosicché, procedendo, si deve fare il prodotto $t (t/2) (t/3) (t/4) \dots (t/n) = t^n/n!$

Accettando – con o senza giustificazioni rigorose o intuitive (o pseudointuitive) – che i termini sono del tipo $t^n/n!$, abbiamo che: scomponendo il montante in capitale più interessi più interessi degli interessi ecc., avremo:

$$e^t = 1 + t + \frac{t^2}{2!} + \frac{t^3}{3!} + \dots + \frac{t^n}{n!} + \dots$$

(In termini tecnici, questo è lo sviluppo della funzione esponenziale, e^t , in « serie di potenze »).

Calcolo di e

Per il calcolo di e tale sviluppo in serie riesce particolarmente idoneo. Ponendo infatti $t = 1$ si ottiene la seguente espressione

$$e = 1 + 1 + \frac{1}{2} + \frac{1}{6} + \frac{1}{24} + \frac{1}{120} + \frac{1}{720} + \frac{1}{5040} + \frac{1}{40320} + \dots$$

ossia (coi soli termini scritti) $e = 2,718\ 278\ 766\dots$; si può poi garantire quale grado di esattezza si sia raggiunto osservando che la somma dei termini successivi trascurati è certamente inferiore (ma circa uguale) a $1/8$ dell'ultimo termine considerato (che è $1/8!$). ossia a $1/322560 \approx 0,000\ 003\ 100$. In

generale, arrestandosi al termine $1/n!$, l'errore è inferiore – ma praticamente quasi uguale – a $1/n$ di esso.

È chiaro che altrettanto semplice è il calcolo di $1/e$ (e^t per $t = -1$): basta prendere la precedente serie con segni alternati.

Ancor più rapidamente si giunge a un'eccellente approssimazione per t minore di 1 (in valore assoluto). Per esempio per $t = 1/2$ si ha $\sqrt{e} = e^{1/2} = 1 + 1/2 + 1/(2^2 \cdot 2!) + 1/(2^3 \cdot 3!) + \dots + 1/(2^n \cdot n!) + \dots$ (e per ottenere $1/\sqrt{e} = e^{-1/2}$ basta alternare i segni).

I logaritmi naturali

In chiusura di questo studio sul numero e dobbiamo almeno dire cosa sono i logaritmi naturali – dato che e in genere è noto proprio come la « base dei logaritmi naturali » – e perché Napier (1614) fu (sostanzialmente) condotto a introdurli scegliendo la base e .

La funzione logaritmo (in una base a qualunque, reale, positiva e diversa da 1, in genere maggiore di 1) è la funzione inversa della funzione esponenziale di base a .

La nozione di funzione inversa è assai semplice e intuitiva, e non limitata al campo della matematica. Per esempio, « il numero di telefono di » e « il titolare del telefono numero » sono due funzioni inverse una dell'altra, le cui tabelle sono gli elenchi alfabetici (ove si trova che « il numero di telefono di » Bruno de Finetti è 832360) e quelli numerici (ove si trova viceversa che « il titolare del telefono numero » 832360 è Bruno de Finetti). Nessuna differenza logica si ha nell'applicare tale concetto nella matematica: per esempio « il cubo di » e « la radice cubica di » (nel campo dei numeri reali) sono due funzioni inverse l'una dell'altra: se è $y = x^3$ è $x = \sqrt[3]{y}$, e così in generale sono inverse l'una dell'altra le funzioni f e g se $y = f(x)$ equivale a dire che $x =$

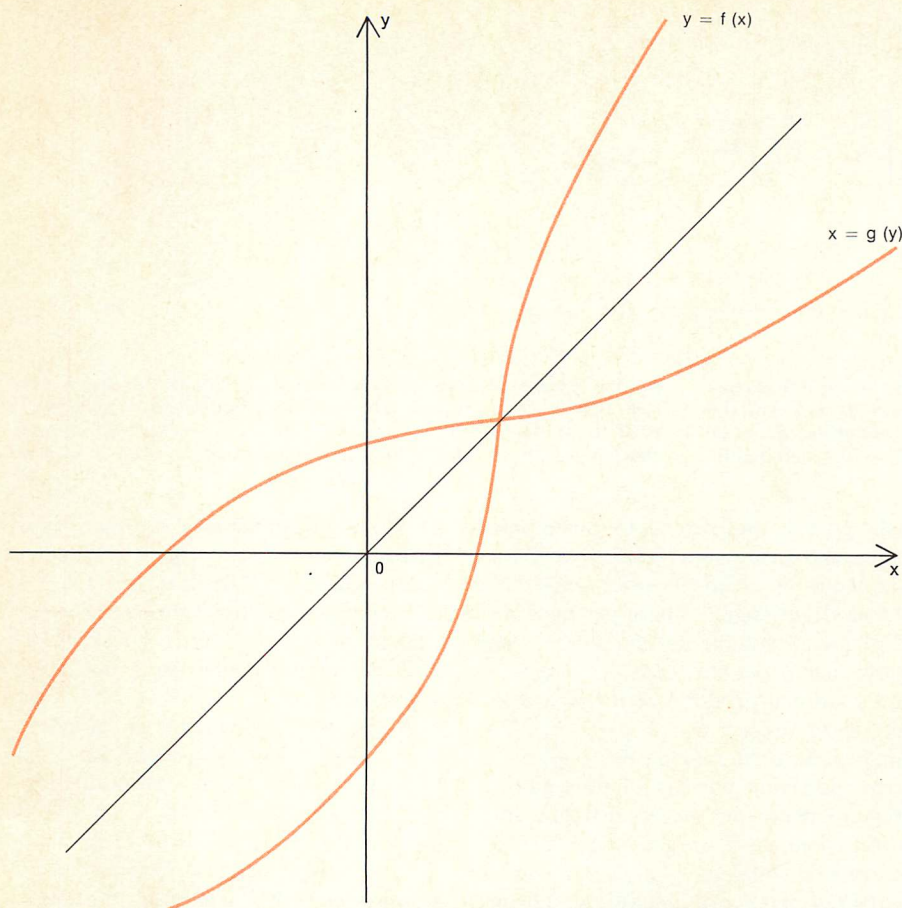


Diagramma di una funzione $y = f(x)$ e della sua inversa $x = g(y)$. Si passa da un diagramma all'altro scambiando gli assi, ossia eseguendo un ribaltamento rispetto alla bisettrice $y = x$. Nel caso in cui la funzione diretta, anziché sempre crescere o decrescere, riprendesse più volte gli stessi valori, la funzione inversa non esisterebbe.

$= g(y)$. Naturalmente occorre che la corrispondenza sia biunivoca (o, altrimenti, renderla tale con opportune limitazioni): occorre per esempio che nessuno abbia più di un telefono, e non si potrebbe prendere l'esempio $y = x^2$ e $x = \sqrt{y}$ perché x e $-x$ hanno il medesimo quadrato, e y ha o nessuna o due radici (salvo limitarsi al campo degli x e y positivi).

Comunque, se si parla di funzioni sempre crescenti (o sempre decrescenti) cade ogni riserva del genere. Ed è facile interpretare visivamente la nozione pensando al diagramma $y = f(x)$; per ottenere quello di $x = g(y)$ (funzione inversa) basta scambiare tra loro gli assi, ossia ribaltare il diagramma rispetto alla bisettrice (si veda la figura in alto), o, se si preferisce, guardare il diagramma da dietro il foglio e ruotato di un angolo retto (in modo da vedere l'asse y orizzontale con la freccia verso destra e l'asse x verticale con la freccia verso l'alto).

Se si parte dalla funzione esponenziale $y = a^x$, la funzione inversa, che dà x dato y , si dice per definizione logaritmo di y (in base a): $x = \log_a y$ o,

(scambiando le lettere per indicare come d'uso con x la « variabile indipendente ») $y = \log_a x$.

Se $a = 10$ si ha il logaritmo decimale (o di Briggs); se $a = e$, si scrive $\log x$ (sottintendendo la base) e il logaritmo si dice naturale (o di Napier): e ciò proprio per lo stesso motivo che ci ha fatto riconoscere come privilegiata la base e per l'esponenziale, e cioè perché la tangente alla curva nel punto 1 dell'asse y (ora asse x) ha pendenza 45° (si veda la figura della pagina a fronte).

Anche il logaritmo si può esprimere mediante una serie di potenze. Si ha questo sviluppo, valido però solo per x compreso tra -1 e $+1$: $\log(1+x) = x - x^2/2 + x^3/3 - x^4/4 + x^5/5 - x^6/6 + x^7/7 - \dots$; osservando che $\log(1-x)$ ha lo stesso sviluppo, ma con tutti i segni negativi, è chiaro che si ottiene un'espressione molto più semplice facendo la differenza, $\log(1+x) - \log(1-x)$, che è il logaritmo del rapporto $(1+x)/(1-x)$, perché i termini di grado pari si elidono mentre quelli di grado dispari si raddoppiano:

$$\log \frac{1+x}{1-x} = 2 \left(x + \frac{x^3}{3} + \frac{x^5}{5} + \frac{x^7}{7} + \frac{x^9}{9} + \frac{x^{11}}{11} + \dots \right).$$

Questa formula permette il calcolo di $\log c$ per ogni c positivo e ci servirà in seguito per il calcolo di π ; in quell'occasione si vedrà anche per quale motivo essa non può essere valida che per x compreso tra -1 e $+1$.

Riferiamoci nuovamente all'interpretazione finanziaria per illustrare praticamente il significato del logaritmo e la ragione della proprietà caratteristica di cui gode (« il logaritmo del prodotto è la somma dei logaritmi »). Il logaritmo di x è il tempo occorrente affinché un capitale cresca nel rapporto da 1 a x (se x è maggiore di 1; altrimenti è il tempo trascorso se da x è ora giunto a 1). Ma il tempo occorrente affinché un capitale si accresca nel rapporto da 1 a xy è la somma di quello occorrente perché cresca nel rapporto da 1 a x e di quello per cui poi cresca nel rapporto da 1 a y ; perciò $\log_a xy = \log_a x + \log_a y$ (sostanzialmente la stessa cosa che $a^{x+y} = a^x a^y$).

I logaritmi decimali sono, notoriamente, i più convenienti per usi pratici dato che 10 è la base del nostro sistema di numerazione. Il logaritmo decimale di 10^n è infatti n (per esempio $\log_{10} 10 = 1$, $\log_{10} 100 = 2$, $\log_{10} 1000 = 3$, ..., $\log_{10} 1 = 0$, $\log_{10} 0,1 = -1$, $\log_{10} 0,01 = -2$, ...). In generale, la parte intera (detta caratteristica) è il numero di posizioni di cui si deve spostare la virgola per portarla a destra della prima cifra significativa (positiva per spostamenti a sinistra, negativa per spostamenti a destra, cioè per numeri minori di 1); a essa va aggiunta, sempre positiva, la mantissa, cioè il logaritmo decimale del numero (tra 1 e 9,999...) che si ottiene spostando la virgola nel modo indicato.

Va inoltre menzionato, almeno di sfuggita, anche il logaritmo in base 2. Nella « teoria dell'informazione » (secondo Shannon) $\log_2 N$ è infatti il numero di « bit » d'informazione occorrenti per individuare l'alternativa esatta tra N ugualmente probabili. Il bit (binary digit, cifra in codice binario) è l'informazione (SI-NO) del caso $N = 2$; è chiaro che se $N = 2^n$ bastano n risposte SI-NO; nel caso di N qualunque ciò sussiste in un senso meno ovvio (di « media »).

Il numero i

Quale senso (intuitivo!) può avere l'introduzione dell'unità immaginaria, i , e quindi dei numeri complessi? Poiché il numero i viene formalmente « definito » come « radice di meno uno » ($i = \sqrt{-1}$), è conveniente cominciare con il prospettare il significato di potenze e di radici, nel caso abituale dei numeri reali, in maniera da condurre

a intravedere la desiderata estensione.

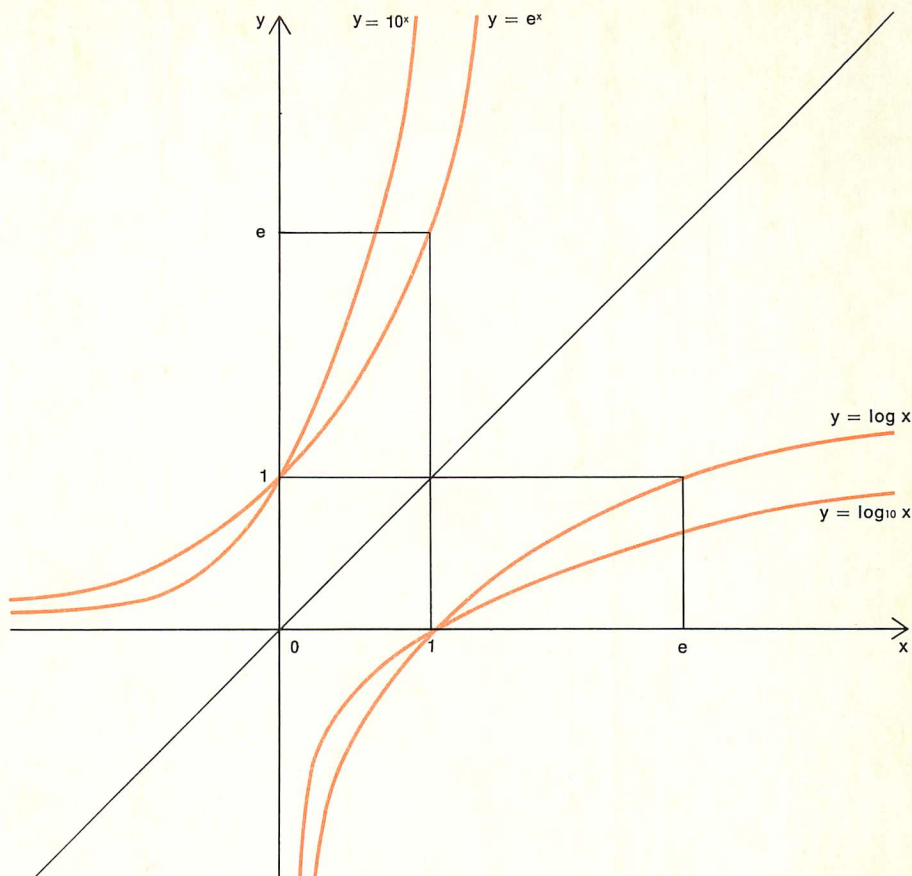
È esattamente la stessa cosa, quando a contenuto, dire che 9 è il quadrato di 3 perché $3 \times 3 = 9$, oppure perché la moltiplicazione per 3, ripetuta due volte, equivale alla moltiplicazione per 9 [cioè, $(3 \times) (3 \times) = (9 \times)$]: triplicare e poi di nuovo triplicare equivale a nonuplicare. Ma, mentre la prima dizione richiama il significato specificamente aritmetico della moltiplicazione fra numeri, la seconda mette invece in evidenza il significato puramente logico e generale della ripetizione di un'operazione o trasformazione. Ciò sembrerà ora una sfumatura, ma è lì, in germe, lo spunto decisivo.

Ogni numero reale a (così come si è visto, nel precedente esempio, per $a = 3$) verrà interpretato come la trasformazione che a ogni cosa x , suscettibile di pensarsi moltiplicata per un numero, fa corrispondere il prodotto ax .

Per fissare le idee su una rappresentazione più limitata e opportuna converrà riferirsi alla retta, e poi al piano, dotati di un'origine fissa O . Se pensiamo di proiettare su uno schermo una retta, con O al centro, e di ingrandire o rimpicciolire l'immagine zoomando (*si veda la figura nella pagina seguente*), tutte le distanze da O (che rimane fisso) dei punti P, Q, R, S, \dots (che vengono portati in P', Q', R', S', \dots) risultano moltiplicate per uno stesso numero a (positivo e maggiore o minore di 1 a seconda che si ha ingrandimento o impicciolimento); la moltiplicazione per $-a$ (negativo) sarà poi la moltiplicazione per a seguita da moltiplicazione per -1 che è il ribaltamento della retta intorno a O , ribaltamento che, volendo, si può immaginare realizzato, senza uscire dalla retta, mediante una simmetria speculare, oppure, permettendoci di invadere il piano, mediante una rotazione di 180° (sempre intorno a O).

Lo stesso vale se si considerano non solo i punti di una retta ma tutti quelli del piano: la moltiplicazione per a positivo fa scorrere tutti i punti lungo i raggi uscenti da O , allontanandoli o avvicinandoli, mentre, se a è negativo, si ha inoltre una simmetria rispetto a O che scambia ogni raggio con quello opposto, e che si può ancora sempre pensare ottenuta con una rotazione di 180° attorno a O . Una tale trasformazione, che indicheremo essa pure con a (praticamente identificandola col numero a inteso come moltiplicatore) si dice omotetia (diretta o inversa a seconda che a è positivo o negativo); la si può pensare limitata a una retta, o a un piano, o anche estesa a spazi a tre o più dimensioni.

La « identificazione » tra omotetie e numeri appare particolarmente giusti-



Il diagramma della funzione logaritmica (qui rappresentata sia in base e che in base 10) si ottiene ribaltando quello della funzione esponenziale (qui rappresentata da e^x e da 10^x). Si noti che, per il logaritmo, cambiando la base cambia solo la scala in senso verticale. Lo stesso accade anche per la scala orizzontale nel caso dell'esponenziale.

ficata adottando il più espressivo linguaggio vettoriale: si può dire che l'omotetia a non fa che moltiplicare per il numero a ogni vettore da O a P , nel senso che la lunghezza viene moltiplicata per a , tenendo fissa la direzione, e conservando o invertendo il verso a seconda che a è positivo o negativo.

Il vettore da O a P si indica spesso con \overrightarrow{OP} , o, meglio, con $P-O$, notazione che ha un significato operativo. Essa consente infatti di scrivere l'omotetia (che porta P in P') come prodotto dei vettori per a : $(P'-O) = a \times (P-O)$. Per « vettore » $P-O$ si può per ora anche intendere, brutalmente, la « freccia » da O a P ; in effetti, però, è essenziale collocare tale disegno dovunque si vuole pur di conservare, nello spostamento, grandezza, direzione e verso. Meglio è dire la stessa cosa prescindendo da nozioni metriche (vedremo più avanti l'interesse di limitarsi alle nozioni affini; a tal fine basta dire che è $P-O = B-A$ se $OPBA$ è un parallelogrammo, ossia se i suoi lati opposti sono paralleli, ovvero ancora, equivalentemente, se le diagonali OB e AP si tagliano a metà (*si veda la figura superiore a pagina 57*).

L'identificazione di omotetie e numeri va bene anche per il prodotto (è

l'osservazione da cui siamo partiti): è $c = ab$ sia come prodotto (funzionale) di omotetie sia come prodotto (aritmetico) di numeri. Il prodotto di due omotetie entrambe dirette o entrambe inverse è pertanto ovviamente un'omotetia diretta, mentre se una è diretta e l'altra inversa il prodotto è un'omotetia inversa; il quadrato di un'omotetia è sempre un'omotetia diretta.

Se chiamiamo radice (quadrata) di una trasformazione a una qualunque trasformazione b che ripetuta due volte equivale ad a , tale cioè che $bb = a$, possiamo dire perciò che *nell'ambito delle omotetie* se ne trovano sempre due (b e $-b$, con $b = \sqrt{a}$ nel senso aritmetico) che sono radici della a se questa è diretta, mentre non se ne trova nessuna se questa è inversa.

Tali considerazioni e conclusioni si potrebbero subito estendere in modo ovvio al caso di radici di ogni ordine (come in aritmetica, di radici n -esime di a ce n'è sempre una e una sola per n dispari, mentre, se n è pari, ce ne sono due, $\pm \sqrt[n]{a}$, se a è positivo e nessuna se a è negativo). Ma arriveremo alle conclusioni generali che ci interessano proseguendo il ragionamento sulla sola radice quadrata.

Rammentiamo solo che \sqrt{a} (per a

positivo) si può ottenere graficamente dalla costruzione nel semicerchio (x è medio proporzionale tra 1 e a) e che, ripetendola, si ottengono sulla curva esponenziale $y = a^x$ tanti punti quanto fitti si vogliano, cosicché rimane determinato con precisione grande a piacere il valore di ogni radice $^n\sqrt{a} = a^t$ per $t = 1/n$.

Tentativo di evasione

Dobbiamo quindi rinunciare alla ricerca di qualcosa che sia radice di «meno uno»? Dobbiamo addirittura interdire tale ricerca come un'assurda insensata chimera? O rassegnarci a considerare valida una risposta consistente in mere e astratte convenzioni?

Abbiamo già detto che il problema non ha soluzione se pensiamo di cercarla nell'ambito delle omotetie (ossia dei numeri reali); ma, avendo visto che tale precisazione era stata messa in rilievo usando il corsivo, il lettore avrà forse capito che tutto stava nel cercare le radici altrove, in un campo opportunamente allargato.

Quale era l'ostacolo? Il ribaltamento della retta, pensato come *simmetria*, non può eseguirsi in due riprese. Ma, se si evade dalla retta pensando di poterle far invadere il piano con una rotazione di 180° , allora la soluzione appare ovvia: basta ripetere due volte una rotazione di 90° (anzi ce ne sono due, potendosi effettuare la rotazione in senso orario o antiorario).

Dato che, anche se si volesse partire considerando solo la retta, tali operazioni obbligano a invadere il piano, diviene ora obbligatorio considerare il piano (mentre prima era facoltativo). La due radici di -1 che abbiamo scoperte le chiameremo i (rotazione di un angolo retto in senso antiorario) e $-i$ (rotazione di un angolo retto in senso orario); sono, nella locuzione dei primi vettorialisti italiani, gli «operatori ma-

novella» che fanno ruotare il piano portando ogni vettore $P-O$ in $P'O = i(P-O)$ o in $-i(P-O)$ (si veda la figura inferiore nella pagina a fronte).

Il bello è che, appena si introduce i con questo significato, si ha automaticamente il modo di esprimere anche ogni altra rotazione, e anzi, più in generale ancora, ogni similitudine (cioè ogni rotazione combinata con un'omotetia, ossia ogni rotazione intorno a O combinata con ingrandimento o impicciolimento con centro in O). Ed è semplicissimo: basta considerare le combinazioni lineari del tipo $a + ib$, cioè somma di un numero reale a e di un «numero immaginario» ib (prodotto dell'unità immaginaria i per un numero reale b); tali combinazioni si dicono numeri complessi, ma il loro significato (nell'interpretazione geometrica, che ne dà la visione più vera e profonda) è quello di similitudini piane.

Prendiamo infatti un qualunque vettore $P-O$ e appliciamogli l'operazione $a + ib$: cerchiamo cioè cosa risulti il vettore $(P'-O) = (a + ib)(P-O)$. Vogliamo naturalmente conservare la linearità, e supponiamo quindi di poter sviluppare il calcolo ottenendo $(P'-O) = a(P-O) + ib(P-O)$; risulta cioè che il vettore $(P'-O)$ è la somma di quello di partenza moltiplicato per a e di quello ruotato ad angolo retto moltiplicato per b (si veda la figura inferiore nella pagina a fronte). È chiaro che tale vettore risulta ruotato di un angolo ϑ tale che $\tan \vartheta = b/a$ e alterato in scala nel rapporto da 1 a $\rho = \sqrt{a^2 + b^2}$.

Ciò vale per ogni P , e rimane così provato che l'applicazione di un numero complesso a un vettore equivale al prodotto di un'omotetia ρ per una rotazione di angolo ϑ . Tale rotazione si può indicare con la notazione $e^{i\vartheta}$, cosicché $a + ib = \rho e^{i\vartheta}$ (per ora si tratta di semplice convenzione, e non occorre specificare come si misurino gli angoli:

quando la notazione assumerà un significato effettivo — come vedremo in seguito — la misura sarà in radianti).

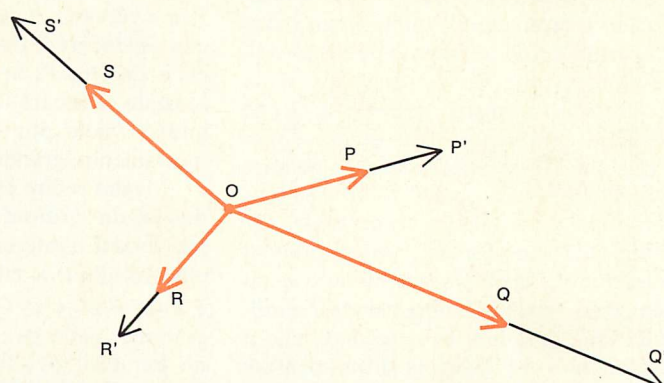
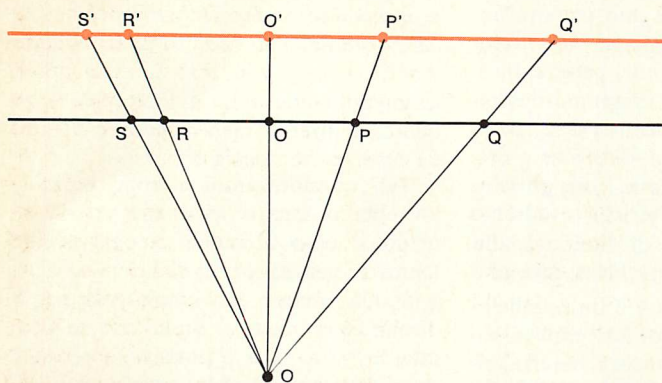
Dato un numero complesso, $a + ib = \rho e^{i\vartheta}$, a si dice sua parte reale, b coefficiente dell'immaginario, ρ modulo e l'angolo ϑ argomento.

Il piano complesso

Convieni raffigurarsi i numeri complessi come punti di un piano, che si dice piano complesso, e altro non è che la naturale estensione al nuovo caso dell'asse delle ascisse sul quale abitualmente si raffigurano i numeri reali. Anzi, tale asse rimane e si chiama asse reale; se ne aggiunge però un altro, ortogonale, passante per l'origine, come asse immaginario, e il punto $P = (x, y)$, di ascissa x e ordinata y , si assume a rappresentare il numero complesso $z = x + iy$; le sue coordinate cartesiane sono la parte reale e il coefficiente dell'immaginario. Più significativo è spesso considerarne le coordinate polari: il modulo, ρ , è la distanza da O , e l'argomento, ϑ , è l'angolo che OP forma col semiasse reale positivo.

A certi effetti è più efficace immaginare $z = x + iy$ come il vettore $P-O$ (la somma, infatti, si fa come per i vettori), o addirittura come la similitudine che porta il punto «1» nel punto P (ossia che porta il vettore unitario del semiasse reale positivo nel vettore $P-O$). In tal modo si vede subito come si costruisca geometricamente il prodotto di due numeri complessi: costruendo triangoli simili (si veda la figura in alto a pagina 58). Infatti, eseguendo due similitudini, i moduli si moltiplicano (omotetie) e gli angoli (rotazioni) si sommano.

La costruzione geometrica del prodotto viene messa efficacemente in luce considerando le successive potenze di un numero complesso: caso questo che ci interessa in modo particolare.



Un'omotetia, in questo caso un allungamento in proporzione 1,5 circa, con centro O porta i punti P, Q, R, S rispettivamente nei punti P', Q', R', S' . Ciò vale sia per un'omotetia per la sola retta (a sinistra), che per comodità di comprensione abbiamo

spezzato in due, sia per il piano (a destra), sia infine per lo spazio. Nulla vieta infatti di pensare che i vettori della figura a destra non stiano nel piano del disegno: possono essere infatti pensati «in prospettiva», cioè in tre dimensioni.

Sia $z = x + iy = \rho e^{i\vartheta}$; le potenze z^2, z^3, z^4, \dots sono i vertici di una spezzata a forma di spirale formata da una successione di triangoli simili al primo, che ha per vertici i punti O , «1» e z (con angolo retto in «1»); i moduli variano in progressione geometrica (sono le potenze del modulo: $\rho, \rho^2, \rho^3, \rho^4, \dots$) mentre gli argomenti crescono sempre di ϑ (sono i suoi multipli: $\vartheta, 2\vartheta, 3\vartheta, 4\vartheta, \dots$). Proseguendo nell'altro senso si hanno $z^{-1}, z^{-2}, z^{-3}, \dots$; in particolare si noti, benché sia ovvio, che $z^{-1} = 1/z$ ha modulo $1/\rho$ e argomento $-\vartheta$ (rispettivamente il reciproco e l'opposto) (si veda la figura in basso nella pagina seguente). In particolare, se $\rho = 1$, tutte le potenze (positive e negative) hanno modulo 1 (geometricamente: sono punti della circonferenza unitaria).

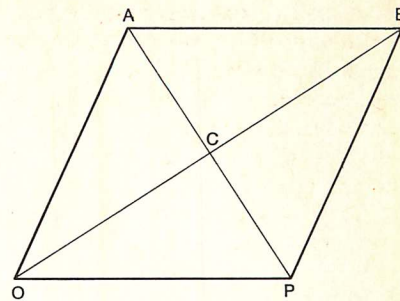
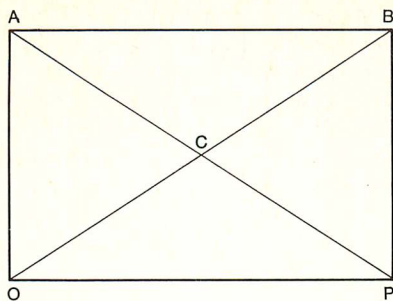
Da questi sommari cenni su ciò che avviene nel campo dei numeri complessi, e che l'interpretazione come punti o come vettori o soprattutto come similitudini del piano rende espressivi e intuitivi, si ricava subito che l'introduzione di i permette non solo di dare due radici quadrate a -1 , ma di dare a ogni numero — sia esso reale (positivo o negativo) o complesso — n radici n -esime, qualunque sia n .

Dato un qualunque numero z , di modulo ρ e argomento ϑ , è chiaro infatti che ne sono radici n -esime tutti i numeri di modulo $\sqrt[n]{\rho}$ e argomento ϑ/n o altro differente da esso per $1/n, 2/n, \dots$ di angolo giro. Per esempio, le radici n -esime dell'unità sono i vertici del poligono regolare di n lati col centro in O e uno dei vertici in «1» (e quelle di -1 si deducono ruotando tale poligono di $1/2n$ di angolo giro). E analogamente nel caso generale: le n radici n -esime di z sono i vertici di un poligono regolare di n lati col centro in O e di cui una radice ha argomento ϑ/n .

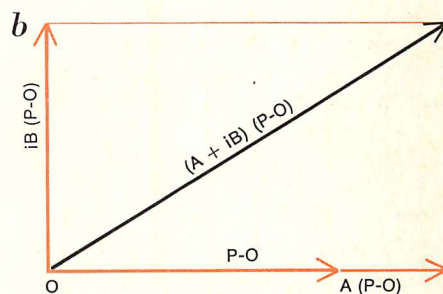
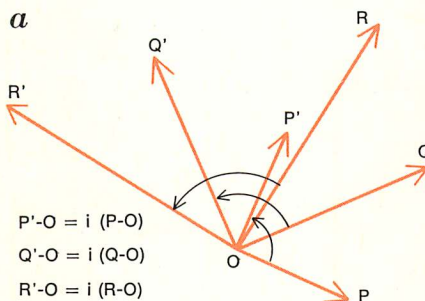
Digressioni

Il poco che si è detto potrebbe bastare per il nostro scopo immediato, che è duplice: 1) dare un'interpretazione e rappresentazione espressiva per i e quindi per i numeri complessi; e, con ciò, 2) fornire le premesse per giungere successivamente a introdurre l'ultimo dei nostri «personaggi», π .

Riassumiamo le conclusioni: un numero complesso $z = x + iy$ si può pensare rappresentato dal (o «identificato» col) punto $P = (x, y)$ del piano complesso (di coordinate x e y), o equivalentemente, e per certi aspetti meglio, dal vettore $P-O$ (di componenti x e y), o infine (e nel modo più espressivo)



Dal punto di vista affine i due parallelogrammi mostrati in questa figura sono uguali; infatti, guardando uno qualsiasi di essi sotto un'opportuna inclinazione (da un punto infinitamente lontano) esso diviene uguale all'altro. Il punto di vista affine consiste appunto nel pensare «di non sapere quale sia la direzione giusta per guardare le figure».



La rotazione ad angolo retto dei vettori di un piano (a) in senso antiorario, che indicheremo con i , porta $P-O$ in $P'-O = i(P-O)$, ecc. Pensando fisso il punto O , porta P in $P' = O + i(P-O)$. La moltiplicazione per un numero complesso (b) fa ruotare tutti i vettori del piano (o tutto il piano se si tiene fisso il punto O) di un determinato angolo e inoltre altera in proporzione tutte le lunghezze dei vettori del piano.

dalla *similitudine* che trasforma il vettore «1», o il punto «1» = (1,0), nel vettore (o punto) $z = x + iy$, il che si ottiene ruotando tutto il piano complesso intorno a O fino a portare il semiasse positivo a passare per il punto z , e operando poi una omotetia (diretta) per portare in definitiva il punto «1» nel punto z .

Sembra tuttavia opportuno soffermarsi su alcune digressioni atte a chiarire dei possibili dubbi, o a sollevarli, per chiarire alcune sottostanti questioni anche a coloro che potrebbero non essere sfiorati da alcun dubbio. Indirettamente, ciò potrà forse anche insegnare a crearsi dei dubbi, che è la via principale per approfondire e far progredire le conoscenze.

Più o meno indirettamente, quasi tutte le considerazioni e digressioni di cui qui di seguito si darà cenno gioveranno inoltre — pur senza essere indispensabili — a meglio penetrare il significato delle successive argomentazioni intese a farci incontrare π . Giovano — beninteso — se uno le prende per quello che sono e che dicono, e cioè come semplici cenni informativi atti a dare un'idea intuitiva di come stanno le cose; chi pretendesse ragionarvi sopra con sicurezza e cognizione di causa dovrebbe approfondire l'argomento per non rischiare di confondersi anziché di

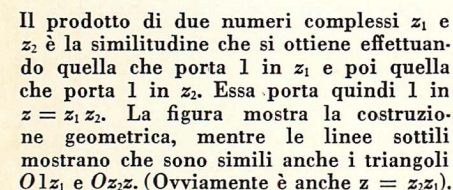
procedere. Nelle digressioni che seguono, soprattutto di natura geometrica, cercheremo di dare almeno una vaga idea della reale utilità del ricorso a ciò che è stato chiamato «immaginario».

Il piano affine e la sfera complessa

Un primo dubbio (o, in mancanza, un'impressione fallace) dovrebbe consistere in ciò: l'introduzione di i , nel modo che abbiamo seguito, presuppone la geometria metrica euclidea in quanto si basa sull'ortogonalità. L'obiezione, in questo momento, è fondata, perché, per illustrare la cosa nel modo più intuitivo, abbiamo scelto subito, senza discussione e quindi in modo acritico e miracolistico, la risposta più comoda. Ma le riflessioni omesse prima, quando potevano disorientare, è bene farle ora per rimediare all'omissione.

Di trasformazioni il cui quadrato è -1 (nello stesso senso che per la rotazione ad angolo retto) se ne potevano trovare, nel piano, infinite. Fissata l'origine O (anch'essa arbitraria) e due punti P e P' (qualsiasi, purché non allineati con O) si sarebbe potuta chiamare i la trasformazione lineare che porta P in P' , pur di stabilire che porti P' nel simmetrico di P . Occorre e basta, cioè porre:

$$i(P-O) = (P'-O) \quad \text{e} \quad i(P'-O) = -(P-O),$$



Nel senso così ampliato avremmo ∞^2 trasformazioni lineari « radici di meno uno »; è un caso che si può studiare, ma la struttura più interessante (per motivi sia teorici che applicativi) si ha quando vi si include un'unica trasformazione i (o, più precisamente, una coppia $\pm i$).

Tale scelta equivale all'introduzione di una metrica nel piano (se prima ne era sprovvisto, cioè era affine). Per spiegarci alla buona: possiamo sempre introdurre in un piano coordinate cartesiane tali che O sia l'origine e P e P' i punti « uno » sugli assi x e y ; dire che il piano è affine significa dire che siamo ancora liberi di scegliere (se vo-

gliamo introdurvi una « metrica ») un sistema di coordinate cartesiane qualunque – per esempio scegliendo ad arbitrio i punti O, P, P' di cui sopra – cui conferire la qualifica di « ortogonale monometrico ». In altri termini, nel piano affine possiamo solo distinguere (come guardandolo obliquamente) se due rette sono o no parallele, e se due segmenti paralleli sono o no uguali (*si veda la figura in alto nella pagina precedente*); confronti tra angoli e fra lunghezze di segmenti non paralleli non hanno senso.

Per darvi senso si può e basta scegliere ad arbitrio due coppie di rette (separantisi) e dirle ortogonali, oppure un triangolo ABC e dirlo equilatero, oppure dirlo isoscele rettangolo in A , e ciò è appunto quel che si è detto sopra per OPP' . Ma ciò è anche la stessa cosa della scelta di i (salvo che essa include anche la scelta di un verso di rotazione positivo, che cambia cambiando i in $-i$ mentre ciò non altera la metrica).

Il piano si può proiettare sulla sfera (o viceversa) mediante la proiezione stereografica (*si veda la figura nella pagina a fronte*) che gode della proprietà di essere conforme (cioè di conservare gli angoli). Essa è usata anche in cartografia, dove tale proprietà significa che la zona intorno a qualunque località sarà, sí, ingrandita in misura diversa rispetto alla sfera-mappamondo da cui si può immaginare proiettata (e precisamente tanto più quanto maggiore è la distanza dal centro *O*) ma senza distorsioni (vengono conservati gli angoli e i rapporti fra lunghezze in qualunque direzione).

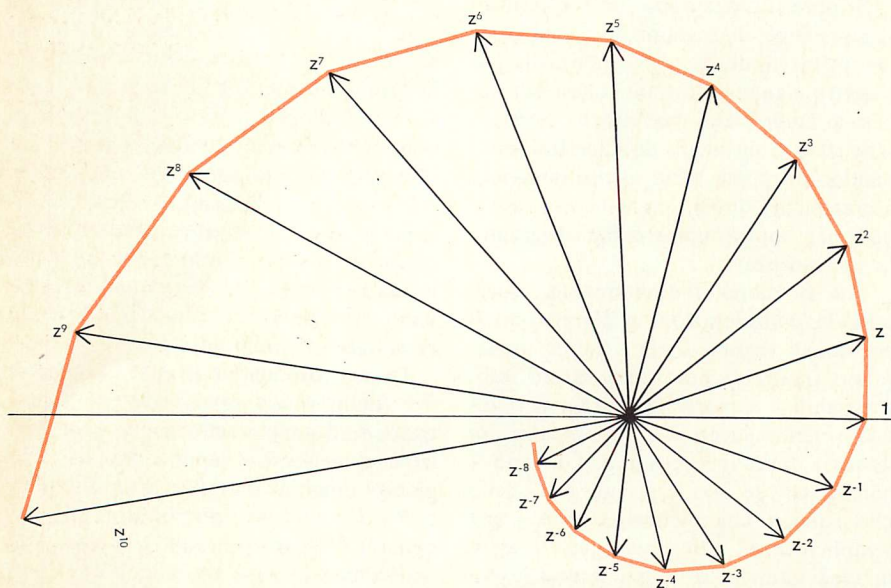
Il piano complesso si può proiettare sulla sfera – e ottenere la sfera complessa, altra forma di rappresentazione dei numeri complessi – in modo da avere (in linguaggio geografico per riferirsi al caso familiare del mappamondo) il punto O (zero) al polo nord, i punti $1, i, -1, -i$ a 90° l'uno dall'altro sull'equatore che corrisponde al cerchio unitario ($\rho = 1$) come tutti gli altri paralleli ai cerchi $\rho = \text{costante}$, mentre i meridiani corrispondono ai raggi $\vartheta = \text{costante}$ (con $\vartheta = 0$, naturalmente, per il meridiano ove si è collocato il punto 1 , semiasse reale positivo)

Al polo sud non corrisponde alcun punto del piano, ma si avvicina a esso l'immagine di un punto che, sul piano, si allontana indefinitamente da O in qualunque direzione: si dice perciò che il polo sud rappresenta il numero « infinito ». Tale convenzione è comoda — per dire meglio, appropriata — nel caso dei numeri complessi perché, sia dal punto di vista algebrico sia da quello analitico, attraverso trasformazioni (per esempio passando da z a $1/z$) il punto all'infinito si trasforma in un altro qualunque senza manifestare eccezionalità.

Per tutte le funzioni « abbastanza naturali » (il senso di tale locuzione si potrà intravedere fra poco) si presenta « naturale » l'estensione al campo complesso, come vedremo per la funzione esponenziale, la cui definizione, già data per e^x (con x reale) verrà estesa a e^z (con $z = x + iy$ complesso): in questo caso, come in altri, è possibile procedere all'estensione, giungendo sempre al medesimo risultato, basandosi vuoi sulla conservazione di proprietà funzionali (equazioni funzionali), come $e^{x+y} = e^x e^y$, o dello sviluppo in serie ($e^x = 1 + x + x^2/2! + x^3/3! + \dots$), o di « proprietà locali » (equazioni differenziali) come $De^x = e^x$, ossia, « e^x è la funzione che è anche la derivata di se stessa », ecc.

Quest'ultima enunciazione richiede di venire spiegata per dare un'idea intuitiva di ciò che è la derivata nel campo reale e nel campo complesso, e di cosa significa la sua esistenza.

Una retta passante per il punto (x_0, y_0) ha l'equazione $y = y_0 + c(x - x_0)$, dove il coefficiente c indica la « pendenza » della retta. Se consideriamo una funzione f , la tangente al suo diagramma in un suo punto (x_0, y_0) , ove naturalmente $y_0 = f(x_0)$, è una retta del tipo predetto, e dove c dev'esser preso in modo che la retta abbia ivi « la stessa pendenza » della curva. (Prendendo c troppo grande la retta è più pendente, cioè sta al di sotto del diagramma nel tratto immediatamente a sinistra, e sopra a destra; viceversa se



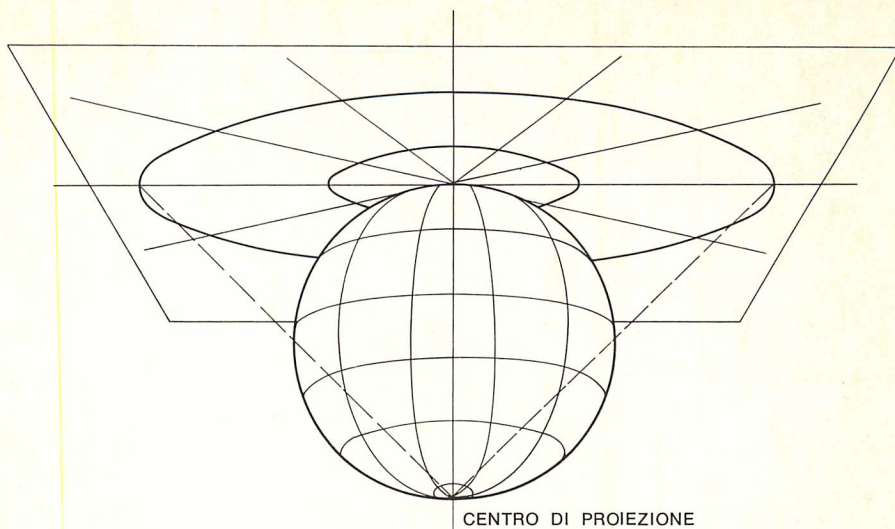
La costruzione geometrica del prodotto viene messa efficacemente in luce considerando le successive potenze di un numero complesso con esponenti interi, positivi o negativi. Si ottiene una spirale (il cui andamento potrebbe essere prolungato indefinitamente in entrambi i sensi) costituita da una successione di triangoli simili al primo.

c si prende troppo piccolo. Se la curva ammette tangente c'è un unico c né troppo grande né troppo piccolo). Tale c si dice « derivata » di f in x_0 ; comunque basti fissare in mente che in prossimità di $x = x_0$ (finché la curva si confonde, praticamente, con la tangente) vale approssimativamente $f(x) = f(x_0) + c(x - x_0)$.

Analogamente, $w = w_0 + c(z - z_0)$ (con c , z e w complessi) è la funzione lineare che in z_0 assume il valore w_0 e il cui valore in z si ottiene aggiungendo lo spostamento $z - z_0$ moltiplicato per c (ossia, se c ha modulo ρ e argomento ϑ , moltiplicato per ρ e ruotato dell'angolo ϑ). Se, per mostrare graficamente l'andamento di $w = u + iv$ in funzione di $z = x + iy$, tracciamo le linee di livello (per quote equidistanti) di u e di v , è chiaro che abbiamo sempre una rete a maglie quadrate (come per x e y) salvo la rotazione e il cambiamento di scala (da 1 a $1/\rho$). E così per una funzione $w = f(z)$, se avviene che in prossimità di $z = z_0$ sia approssimativamente $f(z) = f(z_0) + c(z - z_0)$ con un opportuno c (che si dice derivata della f in z_0), cosicché la rete delle linee di livello $u = \text{costante}$ e $v = \text{costante}$ ha ivi maglie (in piccolo) quadrate, allora la rappresentazione è conforme e concorde (conserva gli angoli e l'orientazione). Le funzioni che si considerano nel campo complesso sono quelle che godono di tale proprietà ovunque salvo eccezioni (per esempio, ovviamente, dove la derivata è nulla come per $w = z^2$ in $z = 0$).

Per mostrare quanto la condizione sia restrittiva basti dire due fatti. Affinché $w = f(z) = u(z) + iv(z)$ la soddisfi, a una data $u(z)$ si può associare solo una corrispondente $v(z)$, e viceversa (a meno di una costante); per esempio $w = ax + iby$ la soddisfa solo se $a = b$, ossia quando è $w = az$. E affinché $f(z)$ sia univocamente determinata basta conoscerla su un'infinità di punti. In particolare, è unica la funzione $f(z)$ coincidente con la $f(x)$ supposta nota per x reale (e questa stessa funzione dev'essere « abbastanza naturale » dal momento che, una volta data in un tratto breve quanto si voglia, non può venir prolungata se non esattamente in un solo modo).

Anche le linee di livello del modulo e dell'argomento di w si tagliano ortogonalmente e danno maglie quadrate se in luogo del modulo si considera il suo logaritmo, perché ciò corrisponde a considerare le linee corrispondenti a $u = \text{costante}$ e $a = v = \text{costante}$ per la funzione $w = \log f(z)$. Particolarmente efficace risulta la visualizzazione delle funzioni di variabile complessa facendo apparire in rilievo il plastico del modu-



Quando da un punto di una sfera (polo) si proiettano gli altri punti della sfera su un piano tangente alla sfera stessa nel polo opposto (proiezione stereografica), i paralleli — in particolare l'equatore — diventano cerchi e i meridiani rette (coordinate polari). Questa rappresentazione è conforme (conserva gli angoli fra rette qualsiasi). Per certe questioni conviene pensare, anziché al piano complesso, alla sfera complessa: per esempio, per non dare una particolare posizione « eccezionale » all'infinito.

lo (o suo logaritmo), su cui le linee di livello dell'argomento appaiono come le linee di massima pendenza (si veda la figura a pagina 62).

Il numero π

Dopo aver imparato a conoscere i primi due dei nostri personaggi, e e i , saranno essi stessi a presentarci spontaneamente il terzo, cioè pi-greco (π). Abbiamo già visto un'espressione, sia pure introdotta a titolo convenzionale, in cui figuravano uniti: con $e^{i\vartheta}$ avevamo indicato infatti il numero complesso che rappresenta la rotazione di un angolo ϑ . Basterà mostrare che essa ha effettivamente il significato di funzione esponenziale per numeri immaginari, pur di misurare gli angoli nel modo appropriato: è cioè in radianti ossia con la lunghezza del corrispondente arco su un cerchio di raggio « uno ».

E π altro non è che tale lunghezza per l'angolo piatto (180°), ossia è la lunghezza della semicirconferenza che va dal punto $z = +1$ al punto opposto $z = -1$. Pertanto è $e^{i\pi} = -1$; è questa la celebre relazione che lega tra loro i nostri tre personaggi, e che appare tanto mirabile e misteriosa ai matematici che per primi ebbero la ventura d'imbarcarsi in essa e di riflettersi sopra (cominciando da Eulero).

L'interpretazione più intuitiva si ha pensando al moto del vettore $e^{i\vartheta}$, o del punto $P = O + e^{i\vartheta}$, in funzione di ϑ inteso come tempo. Il punto P percorre il cerchio unitario, con velocità unitaria in senso antiorario, partendo dal punto « uno » in $\vartheta = 0$. A ciò conducono concordemente (come vedremo)

tutte le vie atte a estendere al caso di esponente immaginario le considerazioni svolte nel caso reale, e cioè ponendo $e^{i\vartheta} \equiv (1 + i\vartheta/n)^n$ (per n grande)

oppure

$$e^{i\vartheta} = 1 + i\vartheta + (i\vartheta)^2/2! + (i\vartheta)^3/3! + (i\vartheta)^4/4! + \dots$$

(cioè $e^{i\vartheta} = 1 + i\vartheta$ per ϑ piccolo)

o anche

$$e^{i\vartheta} = e^{i\vartheta'} e^{i\vartheta''} \quad (\text{se } \vartheta = \vartheta' + \vartheta'').$$

Le prime due espressioni ci consentiranno di sviluppare in seguito considerazioni analoghe a quelle viste nel caso reale, secondo i due « approcci » ivi seguiti, dimostrando e illustrando quanto asserito. Ma l'ultima formula ci dà subito, sostanzialmente, la dimostrazione che ci interessa, espressa peraltro solo in forma intuitiva.

Tra gli istanti ϑ e $\vartheta + \Delta$ (con Δ piccolo) il punto P si sposta da $O + e^{i\vartheta}$ a $O + e^{i(\vartheta+\Delta)}$, ma poiché $e^{i(\vartheta+\Delta)} = e^{i\vartheta} e^{i\Delta} \equiv e^{i\vartheta} (1 + i\Delta) = e^{i\vartheta} + \Delta \cdot i e^{i\vartheta}$, l'ultimo addendo esprime lo spostamento come prodotto del tempuscolo Δ per la velocità $i e^{i\vartheta}$, che è pertanto ortogonale e « uguale » (nell'unità di tempo implicitamente adottata) al vettore $P-O$. Il moto avviene pertanto tagliando sempre ortogonalmente i raggi uscenti da O , e quindi circolarmente intorno a O ; e precisamente sul cerchio unitario perché in $\vartheta = 0$ passa per il punto $z = 1$, e con velocità unitaria perché uguale al raggio (si veda la figura nella pagina seguente).

Il legame tra e , i e π che dovevamo determinare è con ciò stabilito, ma, per l'effettiva determinazione e calcolo, dobbiamo invertire la formula. Formal-

mente potremmo tentare di scrivere

$$i\vartheta = \log e^{i\vartheta},$$

da cui, essendo

$$e^{i\pi} = -1, \quad e^{i\pi/2} = i, \quad e^{i\pi/4} = 1 + i/\sqrt{2}.$$

si ha

$$i\pi = \log(-1), \quad i\pi = 2 \log i, \\ i\pi = 4 \log[(1 + i)/\sqrt{2}].$$

In precedenza, parlando dell'esponenziale e accennando ai logaritmi, avevamo visto gli sviluppi in serie per il logaritmo validi però solo per x reale e compreso tra -1 e $+1$. Potranno servirci nel presente caso? Fortunatamente sí, mediante il semplice accorgimento di ragionare, anziché su ϑ , su $t = \tan \vartheta$; più direttamente ancora, si tratta di considerare i punti $1 + it$ e $1 - it$, intersezioni con la verticale passante per il punto « uno » delle rette che fanno un angolo ϑ col semiasse reale positivo (rispettivamente sopra o sotto). Il rapporto $(1 + it)/(1 - it)$ è il fattore che moltiplicato per $(1 - it)$ dà $(1 + it)$, ossia è la rotazione di un angolo 2ϑ che porta $(1 - it)$ in $(1 + it)$.

Abbiamo $e^{2i\vartheta} = (1 + it)/(1 - it)$, quindi $2i\vartheta = \log[(1 + it)/(1 - it)]$, e, applicando lo sviluppo in serie visto a suo tempo,

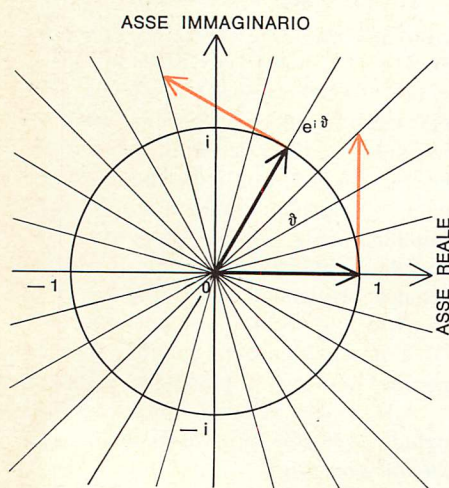
$$2i\vartheta = 2[(it) + (it)^3/3 + (it)^5/5 + (it)^7/7 + \dots] = 2(it - it^3/3 + it^5/5 - it^7/7 + \dots),$$

e infine

$$\vartheta = \arctg t = t - t^3/3 + t^5/5 - t^7/7 + \dots$$

(dove con $\arctg t$ si indica l'angolo la cui tangente è t).

Come sappiamo, tale sviluppo è valido per t compreso tra -1 e $+1$, ossia per angoli fino a 45° ; delle formule sopra scritte solo l'ultima (per $t = 1$, ossia per $\vartheta = \pi/4$) vi rientra (e proprio all'estremo). Comunque essa dà $\pi/4$ mediante la celebre « serie di Leibniz »: $\pi/4 = 1 - 1/3 + 1/5 - 1/7 + 1/9 - \dots$



Pensando ϑ come « tempo », $e^{i\vartheta}$ rappresenta un punto che si muove lungo il cerchio di raggio unitario con velocità unitaria, facendo cioè un giro nel tempo 2π .

Si possono anche ottenere altre formule che danno una convergenza più rapida (una migliore approssimazione con un minor numero di termini); con tali formule si è giunti a calcolare, con elaboratori elettronici, oltre 100 000 cifre decimali di π .

L'esponenziale e il logaritmo nel campo complesso

Le precedenti considerazioni conducono senz'altro a estendere l'esponenziale al campo complesso ponendo, per $z = x + iy$, $e^z = e^{x+iy} = e^x e^{iy}$. Così (e solo così) rimane infatti valida anche nel campo complesso la proprietà caratteristica dell'esponenziale, $e^{z'} + z'' = e^{z'} e^{z''}$. Al medesimo prolungamento conducono del resto anche le altre espressioni, come si vedrà più avanti; intanto diamo uno sguardo al significato della trasformazione $w = f(z)$ e constateremo che è conforme (come occorre). E ciò, in base al criterio di unicità asserito in generale (in fine delle discussioni sulle funzioni $f(z)$ con z complessa), conferma che è l'unica funzione complessa che coincide con e^x sull'asse reale.

Convieni, per maggior chiarezza, considerare due piani distinti (si veda la figura in alto nella pagina a fronte): il piano z su cui rappresentiamo i punti $z = x + iy$, e il piano w sui cui rappresentiamo i corrispondenti punti $w = f(z) = e^z = e^x e^{iy}$. Osserviamo subito che al punto z di coordinate cartesiane x e y sul primo piano corrisponde sul secondo il punto w di cui il modulo e^x e l'argomento y sono le coordinate polari.

Il modo più immediato per render « visibile » la corrispondenza consiste nel disegnare dei sistemi di linee che si corrispondono; ciò è del resto familiare anche dall'abitudine a vedere la rete dei meridiani e paralleli sulle carte geografiche. Nel nostro caso si presenta naturale la scelta di una quadrettatura a coordinate cartesiane sul piano z (pensando segnate le rette equidistanti $x = ka$ e $y = ka$, con k intero e a scelto opportunamente ai fini del disegno), alla quale viene a corrispondere, sul piano w , una rete a tela di ragno, formata da raggi equidistanti uscenti da O (anomalia $\vartheta = ka$) e da cerchi di centro O e con raggi crescenti in progressione geometrica (modulo $\rho = e^{ka}$) di ragione e^a .

Convieni scegliere a in modo che l'angolo $\vartheta = a$ fra due raggi sia « comodo », cioè opportuno sottomultiplo dell'angolo retto, per esempio $a = 1^\circ = \pi/180^\circ = 0,017453$ oppure $a = 10^\circ = 0,17453$.

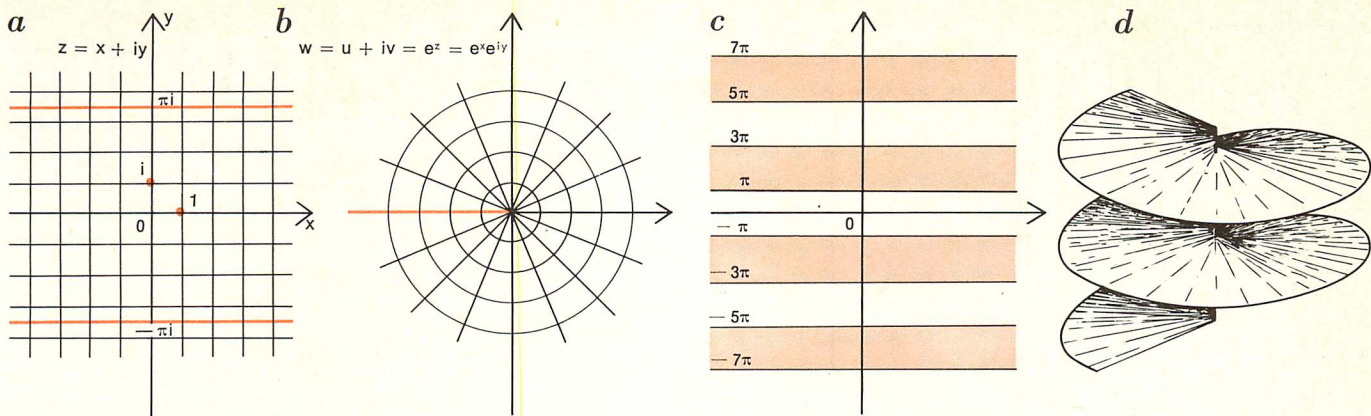
Si vede subito che anche la « tela di ragno » è a « maglie quadrate » come la quadrettatura vera e propria (naturalmente, solo nel senso che ciò è tanto più approssimativamente vero quanto più piccole sono le suddivisioni, ossia a ; se a è grande evidentemente due lati sono sensibilmente curvi, e uno maggiore dell'altro, mentre gli altri due sono rettilinei ma convergenti; ma per a tendente a zero tutte queste imperfezioni vanno scomparendo).

C'è un fatto importante, finora non menzionato per non complicare il discorso, che sarà più chiaro ora che completeremo la discussione parlando del logaritmo: nel campo complesso, l'esponenziale è una funzione periodica. È infatti $e^{2\pi i} = 1$ e quindi $e^{z+2\pi i} = e^z$; la funzione riprende lo stesso valore spostandosi di multipli di 2π in direzione dell'asse immaginario, ossia è periodica di periodo $2\pi i$.

Da ciò che abbiamo visto circa la corrispondenza tra il piano z in coordinate cartesiane e il piano w in coordinate polari si può dedurre senza alcun calcolo che il logaritmo di un numero complesso, di cui si conosca il modulo ρ e l'argomento ϑ , è $\log(\rho e^{i\vartheta}) = \log \rho + i\vartheta + 2k\pi i$ dove l'aggiunta di un multiplo qualsiasi di $2\pi i$ è ovviamente necessaria per l'ovvio fatto che l'angolo ϑ si può alterare di quanti si vogliano angoli-giro in un verso o nell'altro senza che muti la posizione del punto $\rho e^{i\vartheta}$. È un fatto ovvio, ma avevamo evitato di rilevarlo; ora è più chiaro il suo significato sia riguardo all'esponenziale sia riguardo al logaritmo.

Il piano z si può pensare tagliato in infinite strisce orizzontali mediante le parallele $y = (2k + 1)\pi$ all'asse reale $y = 0$ (con k intero, positivo o negativo); ognuna di queste strisce corrisponde a tutto il piano w e, nello stesso modo, a un punto w corrispondono infiniti punti z (uno per ogni striscia, sovrapponibili sovrapponendo le strisce). Per capire come ogni striscia si trasformi per coprire tutto il piano, si immagini dapprima che la metà sinistra (x negative) si rattappisca raccorciandosi e restringendosi a punta e che da tale punta si apra poi come un ventaglio dilatandosi fino a coprire il giro di 360° ; il tratto di asse immaginario ($x = 0$, $0 \leq y \leq 2\pi$) si trasforma nel cerchio unitario; l'asse reale si trasforma nel semiasse reale positivo il semiasse positivo nel tratto $1 < w < +\infty$, quello negativo nel segmento $0 < w < 1$, ove 0 corrisponde a $x = -\infty$).

Ai « tagli » che dividono le strisce del piano z corrisponde un taglio sul piano w lungo il semiasse reale negativo; oltrepassandolo si passa da una



Corrispondenza fra il piano della variabile $z = x + iy$, con rete delle coordinate cartesiane (a) e quello della variabile w data dalla funzione $w = e^z$ (b). La rete corrispondente è in coordinate polari. Al punto $z = 0$ corrisponde $w = 1$; a $x > 0$, $w > 1$ e a $x < 0$, $w < 1$ (parte fuori o dentro del cerchio $w = 1$ corrispondente all'asse immaginario $x = 0$ per z); $w = 0$

corrisponde a $x = -\infty$. Nelle due figure a destra, il piano z si vede diviso in strisce di ampiezza $2\pi i$ (c) e il piano w moltiplicato in infiniti fogli sovrapposti e collegati a elicoide (d) per mostrare come, al variare di y (di $\frac{\pi}{2}$ per w), w riprende gli stessi valori per ogni accrescimento di 2π su y (di $2\pi i$ per z). L'elicoide dà un'immagine più espressiva della situazione.

striscia del piano z a quella successiva (nel medesimo verso). Abbandonando l'immagine dei tagli (utile per una prima spiegazione, ma evidentemente artificiosa e arbitraria) si può viceversa pensare il piano w formato di infiniti fogli sovrapposti formanti come una rampa elicoidale (si veda la figura in alto). Secondo l'immagine precedente, ogni striscia del piano z corrisponde a un piano-foglio della rampa w tagliata lungo il semiasse positivo; col vantaggio però che, pensando distinti idealmente i punti sovrapposti sulla rampa e corrispondenti al medesimo valore w , si ha una corrispondenza biunivoca tra il piano z e la rampa w .

Le serie di potenze

Queste considerazioni ci permettono di chiarire cosa avviene dello sviluppo in serie, che (secondo quanto sappiamo per il campo reale) dovrebbe dare

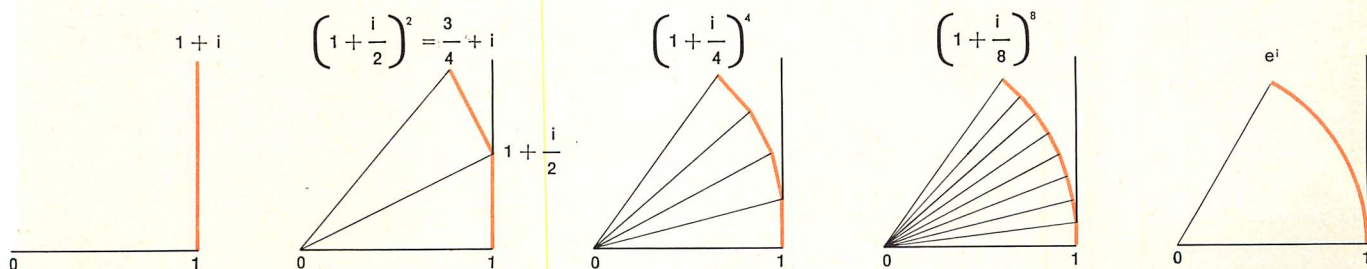
$$\log(1+z) = z - \frac{z^2}{2} + \frac{z^3}{3} - \frac{z^4}{4} + \dots$$

Ma per quali z varrà lo sviluppo? Abbiamo detto, nel campo reale, che esso vale per z compreso tra -1 e $+1$ ed è facile vedere che la convergenza dipende solo dal modulo di z e quindi

varrà nel cerchio di raggio $R = 1$ intorno al punto-origine dello sviluppo. In generale, sia $a_n z^n$ il termine generico della serie, e vediamo se, per un dato z , la successione tenda o no a zero; se si presentano entrambi i casi, è ovvio che la tendenza a zero avverrà per tutti gli z di modulo sufficientemente piccolo (fino a un certo R) e non avverrà per quelli maggiori (oltre tale R) (e sarà $R = 0$ o $R = \infty$ se la convergenza non si ha mai oppure sempre). È chiaro che, fuori del cerchio di raggio R , lo sviluppo non può esser valido; senza che i termini tendano a zero è intuitivo che non può aver senso la « somma della serie ». Il fatto che i termini tendano a zero non basta però a garantire che tale « somma » abbia senso, ma, nel nostro caso, risulta che il dubbio può sussistere solo sul cerchio di raggio R , mentre all'interno tutto va bene (perché i termini diminuiscono almeno altrettanto rapidamente che una progressione geometrica). Si potrebbe subito dedurre l'espressione di R in base al teorema di Cauchy-Hadamard; in pratica esso dice che R è il confine tra gli x per cui nella successione $a_n x^n$ vi sono o no infiniti termini maggiori di 1 (in valore assoluto).

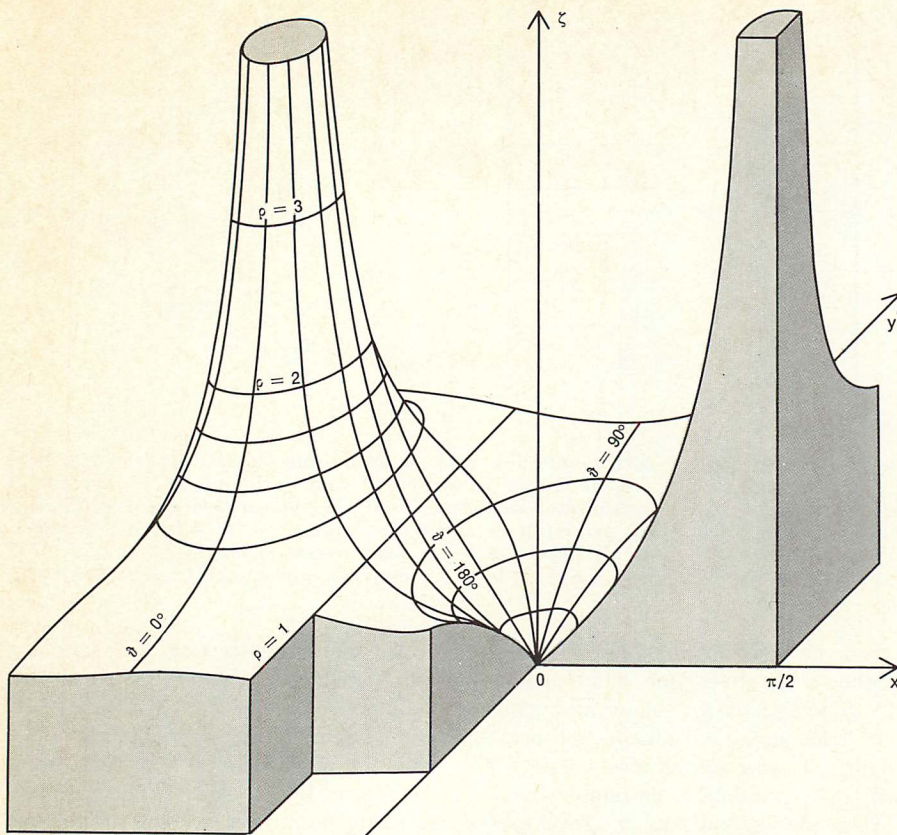
Ma c'è un fatto assai più significativo (anche se è meno facile intuirne la dimostrazione): R è la distanza dal più vicino punto « singolare », dove cioè (in qualche senso) la funzione è irregolare. Il caso più semplice è quello di un punto ove la funzione assume il valore « ∞ » (per esempio per $w = 1/z$ in $z = 0$). Altro caso è quello di un « punto di diramazione », come per esempio l'origine $z = 0$ per la funzione $w = \sqrt{z}$, girando attorno al quale la funzione cambia di segno (tralasciamo di menzionare altri tipi di punti singolari). Nel caso del logaritmo, $z = 0$ è un punto di diramazione infinita perché vi si uniscono gli infiniti fogli della rampa e, girandovi intorno, si passa da uno all'altro (e infatti il cerchio di raggio $R = 1$ con centro nel punto « 1 » passa per l'origine). Come immagine: il campo di convergenza si espande in tutte le direzioni, circolarmente, finché il cerchio urta contro un punto singolare e deve fermarsi.

Per fissare l'attenzione sul solo aspetto che forse può risultare istruttivo senza approfondire tutta la teoria, possiamo dire e suggerire di tener presente che *tutto dipende dall'esistenza e disposizione dei punti singolari* (di cui ne



Se, partendo dal punto 1, ci si sposta di 1 in senso ortogonale si arriva in $1 + i$ (primo disegno). Se ci si sposta di $1/2$, e poi si ripete la stessa similitudine, si arriva in $(1 + i/2)^2$ (se-

condo disegno). Lo stesso accade per 4 e 8 passi, cominciando con $1 + i/4$ e $1 + i/8$, ecc. Si ha in questo modo una serie di spezzate che si approssimano sempre più all'arco di cerchio.



Rappresentazione grafica del modulo della funzione $f(z) = \operatorname{tg} z$ in prospettiva assonometrica. Questo metodo di rappresentazione dà direttamente informazioni sui moduli ρ che danno le «curve di livello» del plastico e sugli argomenti ϑ che danno invece le «curve di massima pendenza». Questo metodo è interessante in quanto accentra la sua attenzione sul modulo della funzione anziché sulle sue parti reale e immaginaria.

esiste sempre almeno uno, escluso il caso banale della «funzione» $w = \text{costante}$; però, se l'unica singolarità è in $z = \infty$, ai nostri effetti non conta).

Prendendo come punto di partenza un qualunque punto z_0 (non singolare) si può esprimere $w = f(z) = f[z_0 + (z - z_0)]$ in serie di potenze di $(z - z_0)$, valida entro il cerchio di raggio pari alla distanza da z_0 del più vicino punto singolare. Per un z_1 fuori di tale cerchio il procedimento, direttamente, non vale; però ci permette di ricavare l'analoga serie in $(z - z_1)$ prendendo come punto di partenza un qualunque punto z_1 interno al cerchio, e tale sviluppo sarà valido nel cerchio che si estende a toccare il più vicino punto singolare, cerchio che andrà al di fuori del precedente. E così di seguito si può, saltando da z_0 a z_1 e poi a z_2 , ecc. coprire con successivi cerchi un'area sempre più vasta e raggiungere qualunque punto del campo di esistenza.

L'interpretazione in base ai due «approcci» iniziali

L'interpretazione dell'esponenziale nel campo reale, nel «primo approccio», consisteva nel considerare e^x , li-

mite di $(1 + x/n)^n$, come capitalizzazione continua intesa nel senso di capitalizzazione semplice in periodi molto piccoli alla fine dei quali gli interessi venivano aggiunti al capitale. Geometricamente, il grafico $y = e^x$ era approssimativamente dato dalla successione di ordinate in progressione geometrica $y_h = (1 + 1/n)^h$ in corrispondenza alle ascisse $x_h = h/n$ (con n grande).

Nel campo complesso tale interpretazione geometrica è ancor più evidente: i punti $z_h = (1 + i/n)^h$, potenze di $z_1 = 1 + i/n$, sono i vertici di una spezzata a spirale ottenuta giustapponendo triangoli rettangoli simili a quello iniziale coi vertici nei punti 0, 1 e $1 + i/n$ (si veda la figura in basso nella pagina precedente). Pensando, in una visione dinamica, a un moto che fa saltare al punto z_h nell'istante $\vartheta = h/n$, il «tempo» ϑ non corrisponderebbe però all'angolo ϑ (argomento) del punto z_h se non approssimativamente (o, più propriamente, asintoticamente) per n grande (ogni salto fa girare di un angolo di cui $1/n$ è la tangente, non l'arco, e cioè un po' meno).

Per mostrare, come dovevamo, la concordanza fra l'estensione basata su questo approccio e quella della via ini-

zialmente seguita, resta a provare che, per n grande, la spezzata a spirale diventa praticamente circolare; in termini di moto, il moto è praticamente quello circolare uniforme con velocità unitaria (sia pure effettuato attraverso n scatti aventi lunghezza $1/n$ per unità di tempo).

In sostanza basta mostrare che e^i , definito come limite di $(1 + i/n)^n$, è proprio, come concluso per altra via, il punto del cerchio unitario di anomalia 1 (angolo: un radiante); va poi da sé che tutta la spezzata coincide praticamente con quella che si ha dividendo l'arco in n parti uguali e rettificandole (sostituendole con la loro corda).

Vediamo dunque quali siano il modulo e l'argomento dei vertici delle nostre spezzate, i punti $(1 + i/n)^h$. Il modulo di $1 + i/n$ è la radice di $1 + 1/n^2$, quello della sua potenza h -esima è quindi $(1 + 1/n^2)^{h/2}$, e per $h = n$ è $(1 + 1/n^2)^{n/2} \approx e^{1/2n} \approx 1$; per n grande il punto finale va a cadere sul cerchio unitario. I lati della spezzata crescono secondo la stessa progressione geometrica dei moduli (ragione $1 + 1/n^2$) cominciando con $1/n$; ma, poiché per n grande i moduli iniziale e finale coincidono ($= 1$), la lunghezza finale è semplicemente $1/n$ moltiplicato per n , ossia 1. Tutto è provato.

Parlando del «primo approccio» ci siamo limitati all'interpretazione geometrica, ma si sarebbe potuto anche dire che il moto circolare mostra l'andamento del montante nella capitalizzazione continua con ... interesse immaginario unitario (l'interesse di un capitale reale positivo è immaginario positivo, di un capitale reale negativo è immaginario negativo, di un capitale immaginario positivo è reale negativo, di un capitale immaginario negativo è reale negativo). Naturalmente, ciò non ha alcun senso nell'interpretazione finanziaria, ma spesso anche un'immagine traslata riesce utile sia per porre attenzione ad analogie formali utili, sia perché in tal modo si può esser condotti ad avvedersi dell'esistenza di altri casi nei quali l'interpretazione ha veramente senso.

Si può osservare, incidentalmente, che si sarebbe potuto parlare di «interesse negativo» per descrivere lo sconto (capitalizzazione con tempo invertito), dove lo sconto dello sconto sarebbe stato positivo, $t^2/2$; e in genere $t^n/n!$ andava preso con segno positivo o negativo a seconda che n fosse pari o dispari (il che dà appunto lo sviluppo di e^{-t}).

Ma intanto ricorriamo comunque a tale immagine, di un interesse immaginario (senza preoccuparci della sua in-

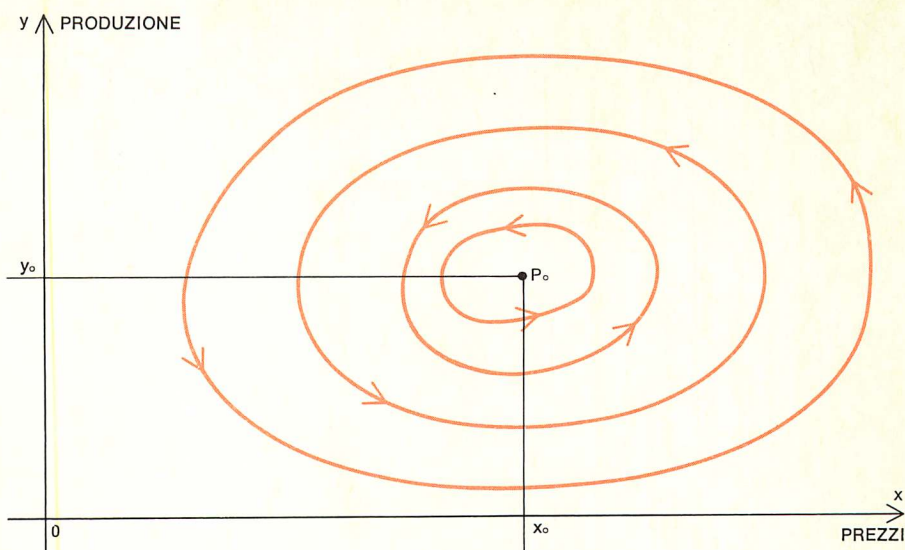
sensatezza finanziaria) per farne uso nel « secondo approccio ».

Nello sviluppo in serie di e^{ix} tutti i termini sono i medesimi che per e^x , cioè $x^n/n!$, tranne che sono moltiplicati per i^n , ossia moltiplicati a rotazione per 1, i , -1 , $-i$. È proprio quel che accadrebbe applicando un interesse immaginario: se si parte da un capitale reale positivo, gli interessi sono immaginari positivi, gli interessi degli interessi sono reali negativi, e via di seguito. Il montante avrebbe una componente reale, continuamente oscillante fra ± 1 (supponendo un capitale iniziale unitario), in quadratura di fase con l'altra (cioè: l'una tocca i valori estremi quando l'altra si annulla).

Ma questo è il comportamento tipico dei fenomeni oscillatori: nel pendolo, (per esempio) lo spostamento e la velocità variano in questo modo, e interpretazioni analoghe si hanno in molti altri campi della fisica; in particolare in elettrotecnica. Indicando, per il pendolo, con x lo spostamento e con y la velocità, avviene che, rappresentando nel piano il punto (x, y) , ossia il punto $z = x + iy$, esso si muove di moto circolare uniforme (beninteso, prendendo opportune unità di misura per spostamento e velocità). Per un cenno si presta meglio un esempio diverso, cioè un modello schematico di fluttuazioni economiche. Supponiamo di indicare con x lo scostamento dei prezzi e con y lo scostamento delle quantità (produzione) dai rispettivi livelli medi (di « equilibrio »), scegliendo le unità di modo che x e y oscillino tra $+1$ e -1 . Qui si ha la stessa rotazione e, con opportuna unità di tempo, sarà $z = e^{it}$; i prezzi alti ($x > 0$) inducono ad accrescere la produzione (y salirà da -1 a $+1$), mentre una quantità elevata ($y > 0$) farà abbassare i prezzi (x scenderà da $+1$ a -1), e così si girerà sempre lungo il cerchio (in senso antiorario) (si veda la figura in questa pagina). Beninteso, perché sia proprio un cerchio occorrono ipotesi particolari (ma semplici e « naturali ») di proporzionalità; inoltre, altre ipotesi possono far prevedere che le oscillazioni si smorzino o divengano « esplosive » (si riducano o crescano di ampiezza), ma ciò va oltre gli scopi del presente cenno.

Riflessioni finali

Nella scorribanda appena terminata abbiamo incontrato, come ci eravamo ripromessi, i tre famosi personaggi ma, oltre che farne la conoscenza, ne abbiamo approfittato per esplorare un po' le regioni ove essi sono nati, per renderci conto del perché sono nati, per intravedere molte cose più o meno col-



Rappresentazione schematica semplice di perturbazioni che generano oscillazioni. Supposto che esista equilibrio quando il livello dei prezzi è x_0 e il livello della produzione è y_0 , un prezzo inferiore o superiore a detto livello spinge a diminuire o ad accrescere la produzione mentre una produzione inferiore o superiore al rispettivo livello di equilibrio provoca una diminuzione o un aumento dei prezzi. Ogni spostamento dal punto di equilibrio provoca un movimento ciclico che, nella figura, è stazionario; potrebbe però anche essere smorzato o, al contrario, « esplosivo », cioè rappresentato da spirali che si avvicinano a P_0 o che si allargano sempre di più.

legate e che meritano di attirare l'attenzione.

E il lettore? Avrà forse un po' di capogiro? Ha la sensazione di aver visto fuggacemente troppe cose, e troppo in disordine? Può darsi, dato che certamente è abituato al metodo scolastico di fare un passo per volta, di essere obbligato a collocare un mattone sopra l'altro, come un manovale che non sa cosa stia costruendo (un castello o una prigione, un bunker o uno stadio, un'officina o un manicomio).

Ma, al confronto, quel metodo era preferibile? Poiché l'unico risultato, purtroppo incontestabile, di quel tipo di insegnamento è la pressoché unanime avversione e incomprensione nei riguardi della matematica, sembrerebbe difficile sostenerlo. E, verosimilmente, coloro che lo ritenessero preferibile in quanto « più facile » intenderanno asserire (con ragione, senza dubbio) che costa minor fatica ottenere la promozione e superare un esame ripetendo enunciati e procedimenti stereotipati senza apprezzarne scopo ed essenza per poi dimenticarli al più presto, che non limitandosi magari a meno ma con l'obbligo di capirlo. Ed è altrettanto vero che sarebbe più facile superare l'esame in una lingua se bastasse saper leggere con discreta pronuncia qualunque testo senza sapere il significato di nessuna parola né addirittura distinguere parole effettive da altre inventate o contraffatte. Ma, in entrambi i casi, si deve obiettare che — da questo punto di vista — meglio ancora è non studiare nul-

la perché quel poco o tanto di fatica, anziché utile, è controproducente, diseducativo, deleterio.

Come mia esperienza personale (di ragazzo), ricordo che, pur riuscendo facilmente nella matematica scolastica, preferivo le materie umanistiche che, nonostante la degradazione a materia scolastica, conservavano qualcosa di vivo. Ma ebbi la fortuna di appassionarmi alla matematica grazie a opere divulgative o culturali, come le polemiche a quell'epoca incandescenti pro o contro la relatività, o come l'antologia di saggi da vari autori (da Platone a Einstein) di Andreas Speiser che, tra l'altro, col mirabile eclettismo dell'autore, includeva come esempio di struttura matematicamente concepita il *Paradiso*, nella descrizione di Dante e nella raffigurazione del Tintoretto.

Ora si insiste molto su argomenti e metodi qualificati (con dubbio gusto) « matematica moderna ». Sono cose ottime se usate con criterio per giungere *prima e meglio* a quella che è pur sempre « la matematica che serve »; se invece vengono considerate come fine a se stesse gonfiate a scapito del resto, non farebbero che peggiorare gli aspetti negativi della matematica « scolastica ». Ciò che occorre è suscitare interesse e curiosità con visioni ampie e suggestive, insegnare, *più per problemi che per teorie*, usare ogni metodo utile ad ampliare le prospettive (e anche, certo, i metodi propri della « matematica moderna », ma in quanto *utili*, non in quanto *sedicentemente « moderni »*).

II

MATEMATICA E LOGICA

Algebra di Boole, diagrammi di Venn e calcolo proposizionale

di Martin Gardner

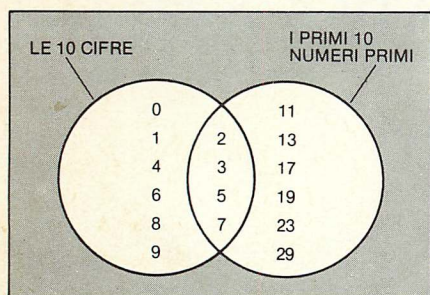
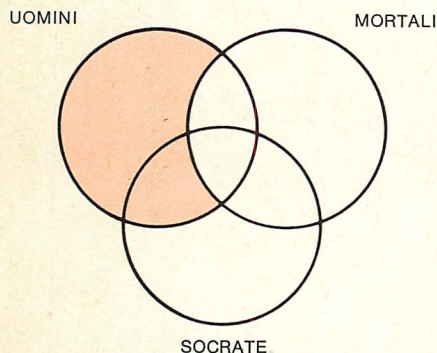
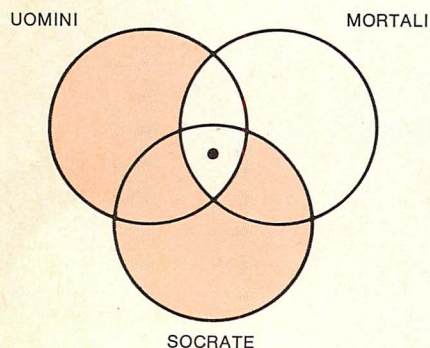


Diagramma di Venn per l'intersezione di insiemi.



Premessa: «Tutti gli uomini sono mortali».



Premessa: «Socrate è un uomo».

Ad Aristotele va tutto il merito di aver fondato la logica formale nonostante abbia quasi interamente limitato la sua attenzione al sillogismo. Ora che il sillogismo è diventato una parte banale della logica, è difficile credere che per 2000 anni esso sia stato l'argomento principale degli studi logici e che ancora nel 1797 Kant potesse scrivere che la logica era «un corpo di conoscenze chiuso e completo».

«Nelle inferenze sillogistiche — ha spiegato B. Russell — si suppone di sapere già che tutti gli uomini sono mortali e che Socrate è un uomo: da ciò si deduce quello che non si era mai sospettato prima, che Socrate è mortale. Questa forma di inferenza viene effettivamente usata, anche se molto raramente». Russell prosegue dicendo che il solo esempio di cui aveva avuto notizia gli veniva da un numero umoristico di «Mind», una rivista inglese di filosofia, che l'editore aveva preparato nel 1901 come numero speciale di Natale. Un filosofo tedesco, reso perplesso dalle inserzioni pubblicitarie della rivista, fece probabilmente questo ragionamento: in questa rivista tutto è scherzo, le inserzioni pubblicitarie sono in questa rivista, quindi le inserzioni sono scherzi. «Per chi vuol diventare un logico — scrisse Russell in altra occasione — posso suggerire un consiglio sul quale non insisterò mai abbastanza: non impari la logica tradizionale. Ai tempi di Aristotele ciò rappresentava uno sforzo degno di rispetto, ma lo stesso può dirsi dell'astronomia tolemaica.»

L'anno della grande svolta fu il 1847, quando George Boole (1815-1864), un modesto autodidatta figlio di un povero calzolaio inglese, pubblicò *L'analisi matematica della logica*. Questa e altre pubblicazioni lo portarono (benché non avesse titoli universitari) all'incarico di professore di matematica al Queens College (l'attuale University College) di Cork, in Irlanda, ove scrisse il trattato *Una ricerca sulle leggi del pensiero sulle quali sono fondate le teorie matematiche della logica e delle probabilità* (Londra, 1854). L'idea fondamentale — la sostituzione con simboli di tutte le parole usate nella logica — non era nuova, ma fu Boole il primo a presentare un sistema funzionale dal punto di vista operativo. Nel complesso, i filosofi e i matematici del suo tempo non mostrarono molto interesse per questo notevole risultato. Forse fu questa una delle ragioni dell'atteggiamento tollerante di Boole verso gli eccentrici della matematica. Scrisse un articolo su un tipo bisbetico di Cork, certo John Walsh («Philosophical Magazine», novembre 1851) che Augustus De Morgan nella sua *Raccolta di paradossi*, definisce «la miglior biografia che io conosca di un eroe del genere».

I pochi che apprezzarono il genio di Boole (in particolare il matematico tedesco Ernst Schröder) migliorarono rapidamente la notazione booleana che presentava alcune deficienze soprattutto perché Boole aveva tentato di rendere il suo sistema quanto più poteva aderente all'algebra tradizionale. Per «algebra di Boole» si intende oggi una struttura astratta «non interpretata», che può essere assiomaticizzata in molti modi, ma che essenzialmente è una versione migliorata e semplificata del sistema di Boole. «Non interpretata» vuol dire che ai simboli che descrivono la struttura non viene assegnato alcun significato di natura logica, matematica o fisica.

Come per tutte le algebre puramente astratte, si possono dare molte interpretazioni diverse ai simboli booleani. Lo stesso Boole interpretò il suo sistema in modo aristotelico, come un'algebra delle classi e delle loro proprietà, che era però un'ampia estensione dell'antica logica delle classi oltre gli stretti confini del sillogismo. Poiché la notazione di Boole è stata superata, l'algebra di Boole moderna è scritta usando i simboli della teoria degli insiemi, intendendo per insieme la stessa cosa che Boole intendeva per classe: una qualunque collezione di singoli «elementi». Un insieme può essere finito, come quello costituito dai numeri 1, 2, 3 (che indicheremo con $\{1, 2, 3\}$) o dagli abitanti di Milano che hanno gli occhi verdi, o dagli spigoli di un cubo, o dai pianeti del sistema solare o da qualsiasi altra specifica raccolta di cose. Un insieme può anche essere infinito, come l'insieme dei numeri interi pari e, probabilmente, l'insieme di tutte le stelle. Se specifichiamo un insieme finito o infinito e poi consideriamo tutti i suoi «sottoinsiemi» (essi comprendono l'insieme stesso come pure l'insieme vuoto, privo di membri) come posti in relazione l'uno

all'altro per inclusione (così l'insieme $\{1, 2, 3\}$ è incluso nell'insieme $\{1, 2, 3, 4, 5\}$), possiamo costruire un'algebra di Boole di insiemi.

Una notazione moderna per un'algebra di questo tipo usa lettere per insiemi, sottoinsiemi ed elementi. L'«insieme universale», il più grande insieme preso in considerazione, si indica con U . L'insieme vuoto o nullo si indica con \emptyset . L'operazione di «unione» degli insiemi a e b (che dà la totalità degli elementi di a e di b) è indicata con \cup . (L'unione di $\{1, 2\}$ e di $\{3, 4, 5\}$ è $\{1, 2, 3, 4, 5\}$). Per la «intersezione» degli insiemi a e b (tutti gli elementi comuni ad a e b) si usa il simbolo \cap . (L'intersezione di $\{1, 2, 3\}$ e $\{3, 4, 5\}$, è $\{3\}$.) Se due insiemi sono identici (per esempio l'insieme di numeri dispari è identico all'insieme di tutti gli interi che divisi per due danno resto 1) si usa il simbolo $=$. Il «complemento» dell'insieme a , costituito da tutti gli elementi dell'insieme universale che non sono in a , è indicato con a' . (Il complemento di $\{1, 2\}$ rispetto all'insieme universale $\{1, 2, 3, 4, 5\}$ è $\{3, 4, 5\}$.) La relazione binaria fondamentale di appartenenza di un elemento a un insieme è simbolizzata con \in ; $a \in b$ vuol dire che a è membro di b . Infine l'inclusione fra insiemi si indica con il simbolo \subset : $a \subset b$ indica che l'insieme a è incluso nell'insieme b (o che a è un sottoinsieme di b).

I simboli impiegati da Boole comprendevano lettere per elementi, classi e sottoclassi, 1 per la classe universale, 0 per la classe nulla, + per l'unione di classi (che egli assumeva in senso «esclusivo» per significare quegli elementi delle due classi che *non* appartengono a entrambe; il passaggio al senso «inclusivo» venne compiuto dal logico ed economista inglese W.S. Jevons, e presenta tali vantaggi che i logici seguenti lo adottarono), \times per l'intersezione delle classi, $=$ per l'identità e il segno $-$ per la sottrazione di un insieme da un altro. Per indicare il complemento di x , Boole scriveva $1 - x$. Egli non dava un simbolo specifico per la relazione di inclusione fra classi ma avrebbe potuto esprimerla in vari modi, per esempio come $a \times b = a$, affermando che l'intersezione di a e b è uguale ad a e quindi che la classe a è inclusa nella classe b .

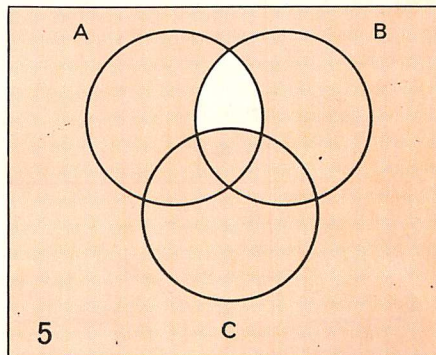
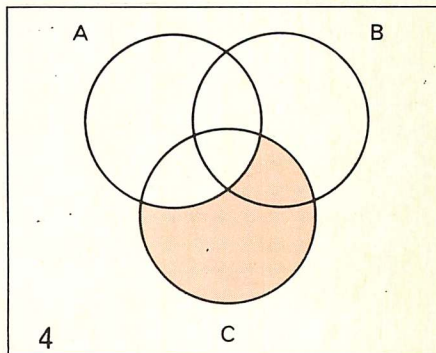
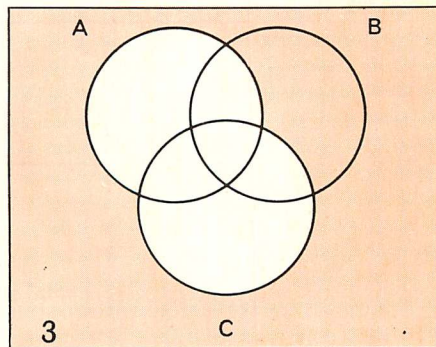
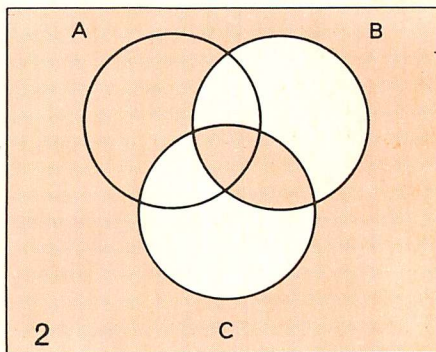
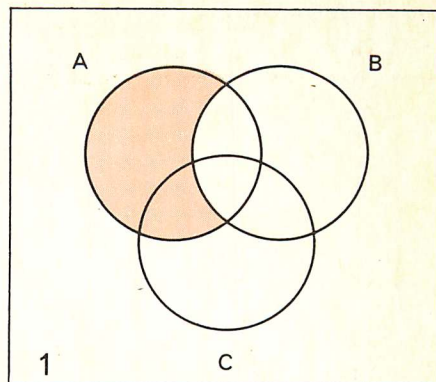
L'algebra di Boole di insiemi trova una elegante espressione grafica nei diagrammi di Venn (introdotti dal logico inglese John Venn). I diagrammi di Venn sono la rappresentazione grafica di una interpretazione dell'algebra di Boole nella topologia degli insiemi di punti del piano. Indichiamo con due cerchi sovrappontentisi l'unione di due insiemi (si veda l'illustrazione in alto nella pagina a fronte), che supponiamo siano l'insieme dei numeri corrispondenti alle prime dieci cifre e l'insieme dei primi 10 numeri primi. Il rettangolo esterno rappresenta l'insieme universale, e questo comprende l'area esterna a entrambi i cerchi, ombreggiata per indicare che essa è l'insieme vuoto; è l'insieme vuoto perché siamo interessati solo agli elementi all'interno dei due cerchi. Questi 16 elementi costituiscono l'unione dei due insiemi. L'area comune ai due cerchi contiene l'intersezione. Essa consiste dell'insieme $\{2, 3, 5, 7\}$, cioè di quelle cifre che si trovano fra i primi 10 numeri primi.

Adottando la convenzione di ombreggiare ogni area che sappiamo rappresentare un insieme vuoto, possiamo vedere come un diagramma di Venn a tre cerchi provi l'antico sillogismo così beffardamente citato da Russell. I tre cerchi sono utilizzati per indicare insiemi di uomini, di cose mortali e di Socrate (quest'ultimo è un insieme con un solo membro). La prima premessa, «Tutti gli uomini sono mortali», viene resa nel diagramma ombreggiando il cerchio degli uomini, per mostrare che la classe degli uomini non mortali è vuota (si veda l'illustrazione al centro nella pagina a fronte).

La seconda premessa, «Socrate è un uomo», viene resa graficamente in modo analogo ombreggiando il cerchio di Socrate, per mostrare che la totalità di Socrate, vale a dire lui stesso, è all'interno del cerchio degli uomini (si veda l'illustrazione in basso nella pagina a fronte). Osserviamo ora il diagramma per vedere se la conclusione «Socrate è mortale» è valida. Lo è. La totalità di Socrate (la parte non ombreggiata del suo cerchio contrassegnata da un punto) è infatti dentro al cerchio delle cose mortali.

La prima importante e nuova interpretazione dell'algebra di Boole venne suggerita da Boole stesso. Egli fece notare che se si interpretavano il suo 1 come verità e il suo 0 come falsità, si poteva applicare il calcolo a enunciati che sono o veri o falsi. Boole non portò a termine questo programma, che fu portato avanti dai suoi successori e diede origine a quello che oggi si chiama calcolo proposizionale. Si tratta cioè del calcolo riguardante enunciati veri o falsi legati da relazioni binarie del tipo: «se p allora q », «o p o q ma non entrambi», «o p o q o entrambi», « p se e solo se q », «né p né q », ecc. La tabella a pagina 69 mostra un confronto fra i simboli del calcolo proposizionale e i corrispondenti simboli dell'algebra di Boole di insiemi.

E' facile capire l'isomorfismo delle due interpretazioni prendendo in considerazione il sillogismo su Socrate. Invece di dire «tutti gli uomini sono mortali», ragionando in termini di proprietà di classi o di inclusione di insiemi, asseriamo: «se x è un uomo allora x è mortale»; in questo modo enunciamo due proposizioni collegate mediante il «connettivo» chiamato «implicazione». Ciò si può rendere graficamente mediante i diagrammi di Venn con lo stesso procedimento usato per il caso di «tutti gli uomini sono mortali». In effetti tutte le relazioni



Soluzione al problema dell'aperitivo mediante i diagrammi di Venn.

binarie del calcolo proposizionale possono essere rese graficamente con i diagrammi di Venn, e i diagrammi possono essere a loro volta usati per risolvere problemi semplici di calcolo. E' un vero peccato che gli autori della maggior parte dei testi introduttivi alla logica formale non si siano ancora resi conto di questo fatto. Essi continuano a usare i diagrammi di Venn per illustrare la vecchia logica di inclusione di classi, ma non li applicano al calcolo proposizionale, ove sono altrettanto efficienti. Anzi, sono perfino più efficienti, dal momento che nel calcolo proposizionale non ci si occupa del «quantificatore esistenziale», il quale asserisce che una classe non è vuota perché le appartiene almeno un elemento. Nella logica tradizionale ciò era espresso dalla parola «qualche» (come in «qualche mela è verde»). Per tener conto di enunciati di questo tipo, Boole è stato costretto a ricorrere a procedure assai complicate e pesanti.

Per vedere come i diagrammi di Venn risolvono facilmente alcuni tipi di rompicapi logici, si considerino le seguenti premesse relative a tre uomini d'affari, Altieri, Bianchi e Cisotti, che fanno colazione insieme ogni giorno feriale:

1. Se Altieri ordina un aperitivo altrettanto fa Bianchi.
2. Bianchi o Cisotti ordinano sempre un aperitivo, ma mai entrambi alla stessa colazione.
3. Altieri o Cisotti o entrambi ordinano sempre un aperitivo.
4. Se Cisotti ordina un aperitivo altrettanto fa Altieri.

Per rappresentare questi enunciati con i diagrammi di Venn, identifichiamo il prendere un aperitivo con il vero e il non prenderlo con il falso. Le otto aree determinate dai cerchi sovrappontendosi, mostrate nell'illustrazione in basso nella pagina a fronte, sono contrassegnate per indicare tutte le possibili combinazioni di valori di verità, con a , b , c , che stanno per Altieri, Bianchi e Cisotti. Così l'area segnata con a , $\sim b$, c rappresenta Altieri e Cisotti che prendono l'aperitivo e Bianchi che non lo prende. Ombreggiando le aree che risultano vuote in base alle quattro premesse ed esaminando il risultato, il lettore può stabilire chi ordinerebbe l'aperitivo qualora egli facesse colazione con i tre uomini. La soluzione risulta particolarmente intuitiva se si procede nel modo seguente: si preparino quattro diagrammi colorati corrispondenti alle quattro premesse, come si vede alle figure 1, 2, 3 e 4 della pagina precedente, quindi sovrapponendo queste quattro figure si ottiene il diagramma della figura 5 da cui si ricava che, se le quattro premesse sono vere, l'unica combinazione possibile dei valori di verità è a , b e $\sim c$, ossia a vera, b vera e c falsa. Poiché si identifica la verità con l'ordinare un Martini, ciò significa che Altieri e Bianchi ordinano sempre Martini mentre Cisotti non lo fa mai.

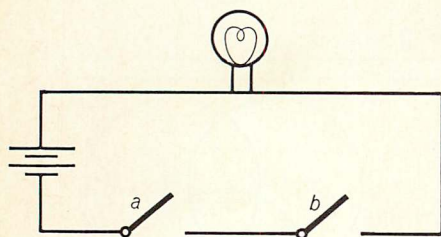
Tra i vari modi in cui può essere interpretata un'algebra di Boole vi è la struttura astratta chiamata «anello», o la struttura astratta chiamata «reticolo»: di entrambe l'algebra di Boole è un caso particolare. Possiede inoltre interpretazioni nell'analisi combinatoria, nella teoria dell'informazione, nella teoria dei grafi, nella teoria delle matrici e nelle teorie metamatematiche dei sistemi deduttivi in generale. In anni recenti l'interpretazione più utile è stata data nella teoria dei circuiti, che è importante nella progettazione dei calcolatori elettronici ma non è limitata ai circuiti elettrici. Si applica a qualsiasi tipo di trasmissione di energia lungo canali con apparati di connessione che permettono o interrompono il flusso all'energia o la deviano da un canale all'altro.

L'energia può essere costituita da un flusso di gas o di liquido, come nei moderni sistemi di controllo a fluido, oppure da fasci luminosi. Può essere energia meccanica come nella macchina logica che Jevons inventò per risolvere problemi a quattro termini nell'algebra di Boole. Ancora, se gli abitanti di un altro pianeta avessero un senso dell'olfatto molto sviluppato, i loro calcolatori potrebbero usare la trasmissione di odori lungo tubi fino a sbocchi annusatori. Nella misura in cui l'energia percorre o non percorre un canale, si stabilisce un nesso tra questi due stati e i due valori di verità del calcolo proposizionale. Per ogni connettivo binario nel calcolo esiste un corrispondente circuito di commutazione. Nell'illustrazione di questa pagina ne vengono presentati tre semplici esempi. Il circuito in basso è usato quando due interruttori molto distanti servono a controllare un'unica lampadina. E' facile vedere che se la luce è spenta, cambiando lo stato di uno qualunque degli interruttori si accenderà e che, se la luce è accesa, uno qualsiasi degli interruttori la spegnerà.

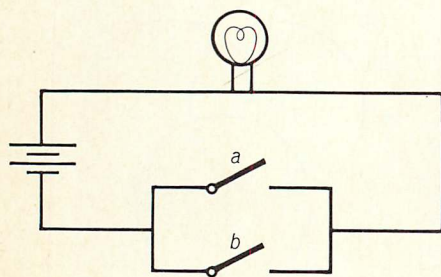
Questa interpretazione dell'algebra di Boole mediante circuiti elettrici è stata suggerita in una rivista russa da Paul S. Ehrenfest nel 1910 e indipendentemente, nel 1936, in Giappone; ma la prima pubblicazione importante, che introducesse tale interpretazione nell'ambiente dei progettisti di calcolatori, fu l'articolo di Claude E. Shannon *A Symbolic Analysis of Relay and Switching Circuits* («Transactions of the American Institute of Electrical Engineers», vol. 57, dicembre 1938). La pubblicazione era basata sulla tesi presentata da Shannon al MIT, dove è ora professore di matematica.

Dal tempo della pubblicazione dell'articolo di Shannon, l'algebra di Boole è diventata essenziale per la progettazione dei calcolatori, dove rende servizi particolarmente utili nella semplificazione dei circuiti allo scopo di risparmiare componenti. Un circuito è prima di tutto tradotto in un enunciato di logica

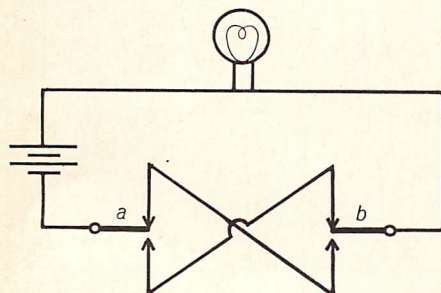
CIRCUITO « E »: LA LAMPADA SI ACCENDE SOLO SE a E b SONO ENTRAMBI CHIUSI.



CIRCUITO « O » INCLUSIVO: LA LAMPADA SI ACCENDE SOLO SE a O b O ENTRAMBI SONO CHIUSI.



CIRCUITO « O » ESCLUSIVO: LA LAMPADA SI ACCENDE SOLO SE a O b , MA NON ENTRAMBI, SONO ABBASSATI.



Circuiti per tre connettivi binari.

ALGEBRA BOOLEANA DI INSIEMI	CALCOLO PROPOSIZIONALE
U (INSIEME UNIVERSALE)	V (VERO)
\emptyset (INSIEME VUOTO)	F (FALSO)
a, b, c, \dots (INSIEMI, SOTTOINSIEMI, ELEMENTI)	p, q, r, \dots (PROPOSIZIONI)
$a \cup b$ (UNIONE: LA TOTALITÀ DI a E b)	$p \vee q$ (DISGIUNZIONE: O p O q O ENTRAMBE SONO VERE)
$a \cap b$ (INTERSEZIONE: ELEMENTI COMUNI AD a E b)	$p \cdot q$ (CONGIUNZIONE: p E q SONO ENTRAMBE VERE)
$a = b$ (IDENTITÀ: a E b SONO LO STESSO INSIEME)	$p \equiv q$ (EQUIVALENZA: p È VERA SE E SOLO SE È VERA q)
a' (COMPLEMENTO: TUTTI GLI ELEMENTI DI U CHE NON SONO IN a)	$\sim p$ (NEGAZIONE: p È FALSA)
$a \in b$ (APPARTENENZA: a È UN MEMBRO DI b)	$p \supset q$ (IMPLICAZIONE: SE p È VERA, q È VERA)

Simboli corrispondenti in due interpretazioni dell'algebra di Boole

simbolica; l'enunciato viene quindi « minimizzato » con metodi opportuni e infine ritradotto in un circuito più semplice. Naturalmente, nei moderni calcolatori, gli interruttori non sono più apparecchiature magnetiche o diodi a vuoto, ma transistori e altri piccolissimi semiconduttori. Per un certo periodo, dopo la pubblicazione dello storico articolo di Shannon, la maggior parte del lavoro intorno alla logica della progettazione dei calcolatori è stata fatta senza comunicazione fra gli esperti dei vari paesi. Gerard Piel nel suo libro *Science in the Cause of Man* (Knopf, 1961) dà notizia che i matematici americani impiegati presso molte grosse società lavorarono per cinque anni, con un costo di circa 200.000 dollari, per ottenere risultati che erano già stati pubblicati in Russia prima che essi cominciassero.

E ora un'ultima interpretazione dell'algebra di Boole che è un'autentica curiosità. Si consideri il seguente insieme di otto numeri: $\{1, 2, 3, 5, 6, 10, 15, 30\}$. Sono i divisori di 30, compresi 1 e 30 (che possiamo chiamare divisori « impropri »). Interpretiamo l'« unione » come il minimo comune multiplo di ogni coppia di questi numeri e l'« intersezione » di una coppia come il massimo comun divisore. L'inclusione insiemistica diventa la relazione « è divisore di ». L'insieme universale è 30, l'insieme vuoto è 1. Il complemento di un numero a è $30/a$. Con queste nuove interpretazioni delle relazioni booleane si ottiene ancora una struttura booleana consistente! Tutti i teoremi dell'algebra di Boole hanno i loro corrispettivi in questo curioso sistema basato sui divisori di 30. Per esempio, nell'algebra di Boole il complemento del complemento di a è semplicemente a , o nell'interpretazione del calcolo proposizionale la negazione di una negazione è lo stesso che se non ci fossero negazioni. Più in generale, solo una serie dispari di negazioni equivale a una negazione. Applichiamo questa legge booleana al numero 3. Il suo complemento è $30/3 = 10$. Il complemento di 10 è $30/10 = 3$, il che ci riporta di nuovo a 3.

Consideriamo due famose leggi booleane chiamate leggi di De Morgan. Nell'algebra degli insiemi esse assumono la forma:

$$(a \cup b)' = a' \cap b'$$

$$(a \cap b)' = a' \cup b'$$

Nel calcolo proposizionale:

$$\sim (a \vee b) \equiv \sim a \cdot \sim b$$

$$\sim (a \cdot b) \equiv \sim a \vee \sim b$$

Se il lettore sostituisce due qualunque divisori di 30 al posto di a e b , e interpreta i simboli come spiegato prima, troverà che le leggi di De Morgan sono valide. Le leggi di De Morgan possono servire a illustrare il famoso principio di dualità dell'algebra di Boole. Se in ogni enunciato si scambiano (se e ovunque compaiano) l'unione con l'intersezione, l'insieme universale con quello nullo, e si inverte la direzione dell'inclusione insiemistica, il risultato è un'altra legge valida. Inoltre, questi cambiamenti possono essere fatti in ogni passo della dimostrazione di una legge per ottenere una dimostrazione valida dell'altra! (Un principio di dualità altrettanto interessante vale in geometria proiettiva rispetto a scambi di linee e punti.)

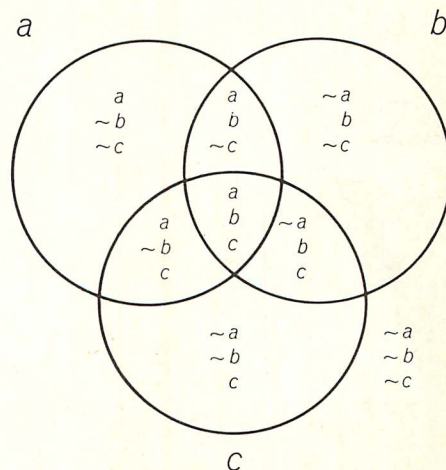


Diagramma di Venn per il rompicapo dell'aperitivo.

Verità e dimostrazione

L'antinomia del mentitore, noto ostacolo a una definizione di verità nei linguaggi naturali, è usata costruttivamente nei linguaggi formalizzati per concludere che non tutte le proposizioni vere possono essere dimostrate

di Alfred Tarski

L'argomento di questo articolo non è nuovo. Esso è stato trattato molto di frequente nella moderna letteratura logica e filosofica, e non sarebbe facile portare qualche contributo originale alla discussione. Per molti lettori, temo, nessuna delle idee avanzate in questo articolo risulterà essenzialmente nuova. Nondimeno, la mia speranza è che essi possano trovare qualche interesse nella maniera in cui il materiale è stato disposto e articolato.

Come indica il titolo, intendo qui esaminare due concetti diversi ma connessi fra loro: il concetto di verità e quello di dimostrazione. L'articolo si divide in tre paragrafi; nel primo viene trattato il concetto di verità, nel secondo il concetto di dimostrazione e nel terzo si considerano le relazioni fra questi due concetti.

Il concetto di verità

Il compito di spiegare il significato del termine « vero » sarà qui interpretato in modo restrittivo. Il concetto di verità compare in molti contesti diversi e ci sono varie categorie di oggetti alle quali può essere riferito il termine « vero »: in un contesto psicologico si può parlare sia di emozioni vere sia di credenze vere; in un contesto estetico si può analizzare la verità interiore di un oggetto d'arte. Ma in questo articolo ci interessa solo ciò che potrebbe chiamarsi il concetto logico di verità. Più esattamente ci interessa esclusivamente il significato della parola « vero » usata in riferimento a una proposizione. Presumibilmente, questo era appunto l'uso originale del termine « vero » nel linguaggio umano. Le proposizioni vengono trattate qui come oggetti linguistici, come certe successioni finite di suoni o di segni scritti. (Naturalmente non tutte le successioni

siffatte sono proposizioni). Inoltre, quando parleremo di proposizioni intenderemo sempre quelle che in grammatica sono chiamate proposizioni dichiarative, e non altri tipi di proposizioni come ad esempio quelle interrogative o imperative.

Ogni volta che si precisa il significato di un termine tratto dal linguaggio quotidiano, si dovrebbe tener presente che lo scopo e lo status logico di una tale spiegazione può variare da un caso all'altro. Per esempio, la spiegazione può essere intesa come un resoconto dell'uso effettivo del termine in questione e in tal caso ha senso domandarsi se il resoconto sia corretto. In altri casi la spiegazione può essere di natura normativa, ossia può essere fornita come indicazione circa l'uso del termine in modo ben definito, senza tuttavia pretendere che il suggerimento si adegui al modo in cui il termine viene effettivamente usato nella pratica; una tale spiegazione può essere valutata, per esempio, dal punto di vista della sua utilità anziché da quello della sua correttezza. E si potrebbe continuare. La spiegazione che daremo noi partecipa, entro certi limiti, di entrambi gli aspetti. Ciò che viene proposto può essere trattato, in linea di principio, come suggerimento di usare il termine « vero » in un certo modo definito, ma, al tempo stesso, ci conforta l'opinione che la nostra proposta sia in accordo con l'uso corrente del termine nel linguaggio quotidiano.

La nostra interpretazione del concetto di verità si accorda nella sostanza con varie spiegazioni di tale concetto che si trovano nella letteratura filosofica. La più antica spiegazione si trova forse nella *Metafisica* di Aristotele:

Dire di ciò che è che non è, o di ciò che non è che è, è falso, men-

tre dire di ciò che è che è o di ciò che non è che non è, è vero.

Qui e in quanto segue la parola « falso » ha lo stesso significato dell'espressione « non vero » e può essere da questa sostituita.

Il contenuto intuitivo della formulazione aristotelica risulta piuttosto chiaro. Tuttavia la formulazione lascia alquanto a desiderare dal punto di vista della precisione e della correttezza formale. Anzitutto non è sufficientemente generale; essa si riferisce solo a proposizioni che « dicono », di qualche cosa, « che è » o « che non è »; nella maggior parte dei casi sarebbe ben difficile far rientrare una proposizione in questo schema senza distorcerne il senso e forzare lo spirito del linguaggio. Questa è forse una delle ragioni per cui nella filosofia moderna sono state proposte varie alternative a quella di Aristotele. Come esempi citiamo i seguenti:

Una proposizione è vera se denota lo stato di cose esistente.

La verità di una proposizione consiste nella sua corrispondenza con la realtà.

Grazie all'uso di termini filosofici tecnici, queste formulazioni hanno indubbiamente un tono molto « dotto ». È tuttavia mia impressione che una volta analizzate più da vicino, le nuove formulazioni risultino meno chiare e più ambigue di quella avanzata da Aristotele.

La concezione della verità che trova la sua espressione nella formulazione aristotelica (e in altre più recenti formulazioni a essa collegate) viene di solito chiamata classica o concezione semantica della verità. Per semantica si intende quella parte della logica che,

grosso modo, tratta le relazioni fra gli oggetti linguistici (quali le proposizioni) e ciò che tali oggetti esprimono. Il carattere semantico del termine « vero » viene chiaramente rilevato dalla spiegazione proposta da Aristotele e da alcune formulazioni che daremo più avanti. A volte si parla della teoria basata sulla concezione classica come della teoria corrispondentistica della verità.

(Nella moderna letteratura filosofica sono state trattate anche altre concezioni e teorie della verità, come la concezione pragmatica e la teoria della coerenza. Queste concezioni sembrano avere carattere esclusivamente normativo e hanno scarsa connessione con l'uso effettivo del termine « vero »; nessuna di esse è stata finora formulata con un minimo di chiarezza e precisione. Nel presente articolo esse non verranno esaminate.)

Cercheremo qui di ottenere una più precisa spiegazione del concetto classico di verità, una che possa superare la formulazione aristotelica pur conservandone gli intenti di fondo. A questo scopo, dovremo ricorrere ad alcune tecniche della logica contemporanea e dovremo specificare il linguaggio nel quale sono formulate le proposizioni che ci interessano. Ciò è necessario, non fosse altro per la ragione che una successione di suoni o di segni che in un linguaggio sia una proposizione vera o falsa, ma comunque sensata, può risultare priva di significato in un altro linguaggio. Per il momento supponiamo che il linguaggio che ci interessa sia l'italiano.

Cominciamo con un problema semplice. Consideriamo una proposizione italiana il cui significato non dia adito a dubbi, per esempio « la neve è bianca ». Per brevità denotiamo questa proposizione con « *S* », cosicché « *S* » diventa il nome di una proposizione. Ci domandiamo: cosa intendiamo dicendo che *S* è vera o che è falsa? La risposta è semplice: in conformità con la spiegazione aristotelica, dicendo che *S* è vera intendiamo semplicemente che la neve è bianca, e dicendo che *S* è falsa intendiamo che la neve non è bianca. Eliminando il simbolo « *S* » arriviamo alla seguente formulazione:

- (1) « La neve è bianca » è vera se e solo se la neve è bianca.
 (1') « La neve è bianca » è falsa se e solo se la neve non è bianca.

Così la (1) e la (1') forniscono spiegazioni soddisfacenti del significato dei termini « vero » e « falso » quando questi termini sono riferiti alla proposizione « la neve è bianca ». Possiamo

considerare la (1) e la (1') come definizioni parziali dei termini « vero » e « falso », anzi come definizioni di questi termini in relazione a una particolare proposizione. Si noti che la (1), così come la (1'), ha la forma richiesta dalle regole della logica per una definizione, precisamente la forma di una equivalenza logica: essa consta di due parti, il primo e il secondo membro dell'equivalenza, collegati dal connettivo « se e solo se ». Il primo membro è il *definiendum*, la frase il cui significato viene spiegato dalla definizione; il secondo membro è il *definiens*, la frase che fornisce la spiegazione. Nel nostro caso il *definiendum* è la seguente espressione:

« la neve è bianca » è vera;

il *definiens* ha la forma:

« la neve è bianca ».

Può sembrare a prima vista che la (1), considerata come definizione, presenti una pecca essenziale ampiamente discussa nella logica tradizionale: un circolo vizioso; e ciò perché alcune parole, per esempio « neve », compaiono sia nel *definiens* sia nel *definiendum*. In realtà, tuttavia, la parola compare nel primo caso con carattere essenzialmente diverso che nel secondo. La parola « neve » è una parte sintattica o organica del *definiens*; precisamente, il *definiens* è una proposizione il cui soggetto è appunto « neve ». Anche il *definiendum* è una proposizione, ed esprime il fatto che il *definiens* è una proposizione vera. Il suo soggetto è un nome del *definiens*, ottenuto mettendo il *definiens* stesso fra virgolette. (Quando si predica qualcosa di un oggetto, si usa sempre il nome di questo oggetto e non l'oggetto stesso: questo vale anche per gli oggetti linguistici.) Per varie ragioni un'espressione fra virgolette va trattata come una parola sola senza parti sintattiche. Perciò la parola « neve » che indubbiamente compare nel *definiendum* come una parte, non vi compare come componente sintattica. La logica medievale direbbe che nel *definiens* « neve » compare in *suppositione formali* e nel *definiendum* compare in *suppositione materiali*. Comunque, le parole che non sono parti sintattiche del *definiendum* non possono dare luogo a un circolo vizioso, e quindi questo pericolo scompare.

Le precedenti osservazioni toccano alcune questioni che sono piuttosto sottili e tutt'altro che semplici dal punto di vista logico. Anziché soffermarsi su di esse indicherei un altro modo per dissipare ogni timore di circolo vizioso.

Nel formulare la (1) abbiamo applicato un metodo comune per ottenere un nome di una proposizione, o di una qualunque espressione: quello che consiste nel mettere l'espressione fra virgolette. Il metodo ha molti pregi, ma è anche all'origine delle difficoltà discusse sopra. Per rimuovere queste difficoltà, tentiamo di adottare un altro metodo per formare nomi di espressioni e precisamente un metodo che potrebbe essere caratterizzato come una descrizione lettera per lettera di un'espressione. Usando questo metodo anziché la (1), avremmo la seguente formulazione più tediosa:

- (2) La successione di quattro parole, la prima delle quali è la successione delle lettere Elle e A, la seconda è la successione delle lettere Enne, E, Vi ed E, la terza è costituita dalla sola lettera E accentata, e infine la quarta è la successione costituita dalle lettere Bi, I, A, Enne, Ci e A, è una proposizione vera se e solo se la neve è bianca.

La formulazione (2) non differisce dalla (1) nel suo significato; la (1) può essere considerata semplicemente come una forma della (2). La nuova formulazione è certamente meno perspicua della vecchia, ma ha su di essa il vantaggio di non creare il sospetto di un circolo vizioso.

Naturalmente definizioni parziali di verità analoghe alla (1) oppure alla (2) possono essere formulate anche per altre proposizioni. Ciascuna di queste definizioni ha la forma:

- (3) « *p* » è vera se e solo se *p*₂

dove « *p* » va sostituito in ambedue i membri della (3) con la proposizione per la quale si costruisce la definizione. Si dovrebbe, tuttavia, usare particolare cura in quelle situazioni nelle quali la proposizione posta in luogo di « *p* » contenga la parola « vero » come parte sintattica. In tali casi la corrispondente equivalenza (3) non può essere considerata come definizione parziale di verità, giacché, se fosse considerata come tale, costituirebbe ovviamente un circolo vizioso. Anche in questi casi, tuttavia, la (1) è una proposizione sensata, anzi è addirittura vera dal punto di vista della concezione classica di verità. A titolo di illustrazione immaginiamo che in una recensione di un libro si trovi la seguente proposizione:

- (4) Non ogni proposizione in questo libro è vera.

Applicando il criterio aristotelico si vede che la proposizione (4) è ve-

ra se, di fatto, non ogni proposizione del libro in questione è vera e che la (4) è altrimenti falsa; in altre parole, possiamo enunciare l'equivalenza ottenuta dalla (3) sostituendo la (4) in luogo di « p ». Naturalmente questa equivalenza si limita a enunciare le condizioni sotto le quali la proposizione (4) è vera o non vera, ma l'equivalenza da sola non ci permette di decidere quale dei due casi si verifichi in realtà. Per verificare il giudizio espresso dalla (4) uno dovrebbe leggere attentamente il libro recensito e analizzare la verità delle proposizioni in esso contenute.

Alla luce della precedente discussione possiamo ora riformulare il nostro problema centrale. Conveniamo che l'uso del termine « vero » in riferimento alle proposizioni italiane è conforme al concetto classico di verità quando e soltanto quando ci permette di stabilire ogni equivalenza della forma (3) in cui « p » sia sostituita in ambo i membri da un'arbitraria proposizione italiana. Se questa condizione è soddisfatta, diremo semplicemente che l'uso del termine « vero » è adeguato. Così il nostro problema principale è il seguente: è possibile stabilire un uso adeguato del termine « vero » per le proposizioni italiane e, in tal caso, mediante quali metodi? Naturalmente si può sollevare una questione analoga per le proposizioni di una qualsiasi altra lingua.

Il problema sarà risolto completamente se riusciremo a costruire una definizione generale di verità che sia adeguata nel senso che ne deriveranno, come conseguenze logiche, tutte le equivalenze della forma (3). Se una tale definizione è accettata dalle persone che parlano italiano, essa stabilisce ovviamente un uso adeguato del termine « vero ».

Sotto certe ipotesi particolari la costruzione di una definizione generale di verità è facile. Si supponga infatti di essere interessati non già alla lingua italiana nel suo complesso, ma solo a una sua porzione e di voler definire il termine « vero » esclusivamente in riferimento alle proposizioni di questo linguaggio parziale; chiameremo linguaggio L questa porzione della lingua italiana. Supponiamo inoltre che L venga corredato di precise regole sintattiche che ci permettano, in ciascun caso particolare, di distinguere una proposizione da un'espressione che non lo è, e che il numero di tutte le proposizioni del linguaggio L sia finito (eventualmente molto grande). Si supponga, infine, che la parola « vero » non compaia in L e che il significato di tutte le parole in L sia sufficientemente chiaro,

sí da evitare ogni obiezione sull'uso delle medesime per definire la verità.

Sotto queste ipotesi si proceda come segue. Anzitutto prepariamo un elenco completo di tutte le proposizioni di L ; supponiamo per esempio che ci siano mille proposizioni in L e conveniamo di usare i simboli « s_1 », « s_2 », ..., « s_{1000} » come abbreviazioni delle successive proposizioni dell'elenco. Poi, per ciascuna delle proposizioni « s_1 », « s_2 », ..., « s_{1000} » costruiamo una definizione parziale di verità sostituendo successivamente queste proposizioni a « p » in ambo i membri dello schema (3). Infine formiamo la congiunzione logica di tutte queste definizioni parziali, cioè combiniamole in un unico enunciato mettendo il connettivo « e » in mezzo a ciascuna coppia di definizioni parziali consecutive. L'unica cosa che resta da fare è di dare alla congiunzione risultante una forma diversa, ma logicamente equivalente, in modo da soddisfare i requisiti formali richiesti alle definizioni dalle regole della logica.

- (5) Per ogni proposizione x (del linguaggio L), x è vera se e solo se
 s_1 , e x coincide con « s_1 »
oppure
 s_2 , e x coincide con « s_2 »
oppure
.....
oppure infine
 s_{1000} , e x coincide con « s_{1000} ».

Siamo giunti così a un enunciato che può ben costituire la desiderata definizione generale di verità: essa è formalmente corretta e è adeguata nel senso che implica tutte le equivalenze della forma (3) nelle quali « p » è stata sostituita da una qualunque proposizione del linguaggio L . Notiamo, per inciso, che la (5) è una proposizione della lingua italiana, ma ovviamente non del linguaggio L ; giacché la (5) contiene tutte le proposizioni di L come parti proprie, non può coincidere con alcuna di esse. Un'ulteriore disamina servirà a meglio chiarire questo punto.

Per ovvie ragioni il procedimento sopra descritto non si può applicare se si è interessati a tutta la lingua italiana, e non semplicemente a una sua porzione. Già se tentiamo di preparare un elenco completo delle proposizioni italiane, incontriamo la difficoltà che le regole della grammatica italiana non determinano in modo preciso la forma delle espressioni (successioni finite di parole) che vanno riguardate come proposizioni: un'espressione particolare, diciamo un'esclamazione, può fungere da proposizione in un dato contesto, mentre un'espressione della stessa forma

non lo sarà in un contesto diverso. Di più, l'insieme di tutte le proposizioni italiane è almeno potenzialmente infinito. È certamente vero che al momento attuale sono state pronunciate e scritte solo un numero finito di proposizioni dagli esseri umani: probabilmente nessuno sarebbe dell'avviso, però, che l'elenco di queste proposizioni esaurisca tutte le proposizioni italiane. Al contrario, vedendo un tale elenco, ciascuno di noi potrebbe facilmente esibire una proposizione italiana che non si trovi già in esso. Infine il fatto che la parola « vero » compaia nella lingua italiana è sufficiente a impedire un'applicazione del procedimento descritto sopra.

Dalle precedenti osservazioni non segue che non si possa ottenere per qualche altra via la desiderata definizione di verità per proposizioni italiane arbitrarie. C'è, tuttavia, una ragione più seria e fondamentale che sembra precludere una tale possibilità. Di più, la mera ipotesi che si possa trovare mediante un qualunque metodo un uso adeguato del termine « vero » (in riferimento a proposizioni italiane arbitrarie) conduce chiaramente a una contraddizione. Il ragionamento più semplice per ottenere una tale contraddizione è noto come *antinomia del mentitore*; lo descriveremo ora brevemente.

Si consideri la proposizione:

- (6) La proposizione stampata in colore a pagina 72 del libro *Verità e dimostrazione* è falsa.

Conveniamo di usare « s » come abbreviazione per questa proposizione. Controllando il titolo di questo libro e il numero di questa pagina, si riscontra facilmente che « s » è proprio l'unica proposizione stampata in colore a pagina 72 del libro *Verità e dimostrazione*. Ne segue in particolare che

- (7) « s » è falsa se e solo se la proposizione stampata in colore a pagina 72 del libro *Verità e dimostrazione* è falsa.

D'altra parte « s » è indubbiamente una proposizione italiana. Perciò, supponendo che l'uso del termine « vero » sia adeguato, possiamo stabilire l'equivalenza (3) in cui « p » sia sostituito con « s ». Così possiamo affermare:

- (8) « s » è vera se e solo se s .

Ricordiamo ora che « s » sta per l'intera proposizione (6). Quindi possiamo sostituire « s » con la (6) nel secondo membro della (8); si ottiene allora:

- (9) «s» è vera se e solo se la proposizione stampata in colore a pagina 72 del libro *Verità e dimostrazione* è falsa.

Confrontando ora la (7) con la (9) si conclude:

- (10) «s» è falsa se e solo se «s» è vera.

Ciò conduce a un'ovvia contraddizione: «s» risulta sia vera che falsa, e siamo così di fronte a un'antinomia. La formulazione sopra riportata dell'antinomia del mentitore è dovuta al logico polacco Jan Łukasiewicz.

Sono note anche formulazioni più elaborate dell'antinomia. Immaginiamo, per esempio, un libro di cento pagine contenente una e una sola proposizione per ciascuna pagina. A pagina 1 si legge:

La proposizione stampata a pagina 2 di questo libro è vera.

A pagina 2 si legge:

La proposizione stampata a pagina 3 di questo libro è vera.

E così via fino a pagina 99. A pagina 100, l'ultima pagina, troviamo:

La proposizione stampata a pagina 1 di questo libro è falsa.

Supponiamo che la proposizione stampata a pagina 1 sia davvero falsa. Mediante un ragionamento non difficile ma molto lungo e che richiede di sfogliare tutto il libro, si conclude che la nostra supposizione è errata. Di conseguenza supponiamo che la proposizione stampata a pagina 1 sia vera e, con un ragionamento altrettanto facile e lungo del precedente, ci si convince che anche la nuova supposizione è errata. Così ci troviamo di nuovo di fronte a un'antinomia.

Non è difficile comporre molti altri libri antinomici che siano varianti di quello descritto. Ammettiamo che ciascuno di essi abbia 100 pagine, ciascuna delle quali contenga una e una sola proposizione della forma:

La proposizione stampata a pagina 00 di questo libro è XX,

dove in ciascun caso particolare «XX» va sostituito con una delle parole «vera» o «falsa», mentre «00» va sostituito con uno dei numerali «1», «2», ..., «100»; non si esclude che lo stesso numerale compaia in più pagine. Non è detto che ogni variante composta secondo queste regole dia effet-

tivamente luogo a un'antinomia. Il lettore che ama i rompicapo logici non troverà difficile descrivere tutte quelle varianti che fanno al caso nostro. La seguente avvertenza può risultare utile al riguardo: si immagini che a qualche punto del libro, diciamo a pagina 1, si dica che la proposizione a pagina 3 è vera, mentre in qualche altro punto, diciamo a pagina 2, si affermi che la stessa proposizione è falsa. Da queste informazioni non segue affatto che il libro sia «antinomico»; se ne può solo dedurre che o è falsa la proposizione di pagina 1 o quella di pagina 2. Un'antinomia sorge, invece, ogniqualevolta si sia in grado di dimostrare che una delle proposizioni del libro è sia vera sia falsa, indipendentemente da ogni ipotesi sulla verità o falsità delle rimanenti proposizioni.

L'antinomia del mentitore risale all'antichità. Di solito viene attribuita al logico greco Eubulide; essa ha tormentato molti logici antichi e ha causato la morte prematura di almeno uno di essi, Fileta di Coa. Altre antinomie e paradossi furono scoperti nell'antichità, nel medioevo e in tempi recenti. Sebbene molte di esse siano ormai completamente dimenticate, l'antinomia del mentitore viene ancora analizzata e discussa negli scritti contemporanei. Assieme ad alcune recenti antinomie scoperte all'inizio del secolo (in particolare l'antinomia di Russell), essa ha avuto un grande influsso sullo sviluppo della logica moderna.

Nella letteratura sull'argomento si incontrano due modi diametralmente opposti di affrontare le antinomie. Uno è quello di ignorarle, di trattarle come sofismi, come giochi che non sono seri ma maliziosi e che mirano soprattutto a dimostrare l'ingegnosità di chi le formula. L'atteggiamento opposto è caratteristico di certi pensatori del XIX secolo e è ancora presente o lo era fino a poco tempo fa. Secondo questo modo di vedere, le antinomie costituiscono un elemento essenziale del pensiero umano: esse continueranno a comparire nelle attività intellettuali e la loro presenza è la fonte basilare del vero progresso. Come spesso accade, la verità è probabilmente una via di mezzo. Personalmente, come logico, non posso persuadermi che le antinomie siano un elemento permanente del nostro sistema di conoscenze; tuttavia non sono affatto incline a prendere le antinomie alla leggera. La comparsa di un'antinomia è per me un sintomo di malattia: partendo da premesse che sembrano intuitivamente ovvie, usando forme di ragionamento che intuitivamente sembrano sicure, un'antinomia ci conduce a conclusioni assurde, contraddit-

torie. Ogni volta che ciò accade, dobbiamo sottoporre i nostri modi di pensare a una revisione approfondita, rinunciare a certe premesse, alle quali credevamo, oppure migliorare certe forme di ragionamento che eravamo abituati a usare. E questo lo facciamo nella speranza non solo di liberarci dell'antinomia, ma anche di non incontrarne di nuove. A questo scopo mettiamo alla prova il nostro pensiero così riveduto con tutti i mezzi a disposizione, e per prima cosa cerchiamo di ricostruire la vecchia antinomia nella nuova sistemazione; queste prove costituiscono un'attività molto importante nel campo del pensiero speculativo, simile a quella di condurre a termine degli esperimenti cruciali nella scienza empirica.

Da questo punto di vista consideriamo ora in particolare l'antinomia del mentitore. Nell'antinomia interviene il concetto di verità in riferimento a proposizioni arbitrarie della lingua italiana; essa potrebbe facilmente essere riformulata in modo da applicarsi ad altri linguaggi naturali. Ci troviamo di fronte a un serio problema: come si possono evitare le contraddizioni indotte da questa antinomia? Una soluzione radicale del problema che può subito venire in mente è semplicemente quella di abolire la parola «vero» dal vocabolario italiano o per lo meno di trattenerci dall'usarla in ogni seria questione.

Coloro ai quali una tale mutilazione dell'italiano sembra altamente insoddisfacente e illegittima possono essere inclini ad accettare una soluzione di compromesso, che consiste nell'adottare ciò che potrebbe chiamarsi (secondo il filosofo polacco contemporaneo Tadeusz Kotarbinski) «la concezione nichilista della teoria della verità». Secondo questa concezione la parola «vero» non ha un significato indipendente, ma può essere usata come componente delle due espressioni sensate «è vero che» e «non è vero che». Queste espressioni vengono trattate come se fossero singole parole, senza parti organiche. Il significato a esse attribuito è tale che possano essere immediatamente eliminate da ogni proposizione in cui compaiano. Per esempio, anziché dire

è vero che tutti i gatti sono neri

si può semplicemente dire

tutti i gatti sono neri

e invece che

non è vero che tutti i gatti siano neri

si può dire

non tutti i gatti sono neri.

In altri contesti la parola « vero » è priva di senso. In particolare, non può essere usata come un vero predicato che qualifica nomi di proposizioni. Nella terminologia della logica medievale possiamo dire che la parola « vero » si può usare come parola sincategorematica in alcune situazioni particolari, ma non la si può mai usare in modo categorematico.

Per rendersi conto delle implicazioni di questa concezione, si consideri la proposizione che è stata il punto di partenza per il paradosso del mentitore, cioè la proposizione stampata in colore a pagina 72 di questo libro. Dal punto di vista « nichilista » essa non è una proposizione sensata e l'antinomia semplicemente svanisce. Disgraziatamente, anche molti altri usi della parola « vero », che altrimenti sembrano pienamente legittimi e ragionevoli, restano parimenti colpiti da questa concezione. Immaginiamo, per esempio, che un certo termine che si ripete frequentemente nelle opere di un matematico antico ammetta diverse interpretazioni. Uno storico della scienza che studia tali opere giunge alla conclusione che sotto una di queste interpretazioni tutti i teoremi enunciati dal matematico risultano veri; ciò lo conduce in modo naturale alla congettura che la stessa cosa valga per ciascuno dei lavori di questo matematico che sono sconosciuti al presente ma che possono essere scoperti in futuro. Se però lo storico della scienza condivide la concezione « nichilista » della verità, perde per ciò stesso la possibilità di esprimere a parole la sua congettura. Si potrebbe dire che la teoria « nichilista » della verità rende un omaggio ipocrita a certe forme del linguaggio umano entrate nelle consuetudini, mentre in realtà abolisce l'idea di verità dal patrimonio concettuale della mente umana.

Proveremo perciò un'altra via per uscire dalla nostra situazione spiacevole, cercando una soluzione che lasci essenzialmente intatto il concetto classico di verità. L'applicabilità del concetto di verità dovrà subire alcune restrizioni, ma il concetto rimarrà disponibile almeno per gli scopi di un discorso colto.

A questo scopo dobbiamo analizzare quei tratti caratteristici del linguaggio comune che sono la fonte reale dell'antinomia del mentitore. Nel corso di questa analisi colpisce subito un aspetto notevole del linguaggio in questione: il suo carattere universale, onnicom-

prensivo. Il linguaggio comune è universale, né deve essere altrimenti, giacché ci si aspetta che esso fornisca i mezzi adeguati a esprimere ogni cosa che possa essere espressa, in un qualsivoglia linguaggio; esso si arricchisce di continuo per soddisfare tale requisito. In particolare il linguaggio comune è semanticamente universale nel senso seguente: accanto agli oggetti linguistici, come proposizioni e termini, che sono componenti del linguaggio, in esso sono presenti anche i nomi di tali oggetti (come sappiamo, un nome di un'espressione si può ottenere ponendo l'espressione fra virgolette); inoltre, il linguaggio contiene termini semantici quali « verità », « nome », « designazione » che, direttamente o indirettamente, si riferiscono alla relazione fra gli oggetti linguistici e ciò che essi esprimono. Di conseguenza, per ogni proposizione formulata nel linguaggio comune, possiamo formarne nello stesso linguaggio un'altra la quale esprima che la prima proposizione è vera o falsa. Con un ulteriore espediente, possiamo perfino costruire nel linguaggio ciò che talvolta vien detta una proposizione autologa, cioè una proposizione *S* che esprime il fatto che *S* stessa è vera o che è falsa. Se *S* esprime la propria falsità, si può dimostrare con un semplice ragionamento che *S* è contemporaneamente vera e falsa, e così ci ritroviamo di fronte l'antinomia del mentitore.

Si noti, tuttavia, che non in tutte le situazioni è necessario l'uso dei linguaggi universali. In particolare, tali linguaggi non sono generalmente necessari per gli scopi della scienza (e per scienza intendo qui l'intero campo dell'investigazione intellettuale). In rami particolari della scienza, diciamo nella chimica, si studiano certi oggetti particolari, come gli elementi, le molecole eccetera, ma non per esempio oggetti linguistici come proposizioni o termini. Il linguaggio che ben si adatta a tale trattazione è perciò un linguaggio ristretto con un vocabolario limitato: esso deve contenere nomi per gli oggetti chimici, termini come « elemento » e « molecola », ma non nomi per oggetti linguistici; non è perciò necessario che sia semanticamente universale. Un discorso analogo vale per la maggior parte degli altri rami della scienza. La situazione si fa un po' confusa, però, se passiamo a considerare la linguistica, ossia la scienza in cui si studiano i linguaggi: il linguaggio della linguistica deve certamente essere fornito di nomi per gli oggetti linguistici. Tuttavia non è necessario identificare il linguaggio della linguistica col linguaggio universale, né con alcuno dei linguaggi che

sono oggetto di indagine linguistica, e non siamo tenuti a supporre che in linguistica si usi uno stesso linguaggio per tutte le indagini. Il linguaggio della linguistica deve contenere nomi per le componenti linguistiche dei linguaggi studiati, ma non i nomi delle sue proprie componenti; così di nuovo, non è necessario che sia semanticamente universale. Lo stesso può dirsi per il linguaggio della logica, o meglio di quella parte della logica nota come metalogica e metamatematica; anche qui abbiamo a che fare con linguaggi e in primo luogo con i linguaggi delle teorie logiche e matematiche (sebbene qui si studino tali linguaggi da un punto di vista diverso da quello della linguistica).

Ci si può chiedere ora se possa definirsi in modo preciso un concetto di verità e quindi se possa stabilirsi un uso coerente e adeguato di questo concetto, almeno per i linguaggi semanticamente limitati del discorso scientifico. Sotto certe condizioni la risposta a tale domanda risulta affermativa. Le principali condizioni imposte al linguaggio sono che il suo vocabolario sia completamente ed esplicitamente determinato e che siano formulate con precisione le regole sintattiche che riguardano la formazione delle proposizioni e delle altre espressioni sensate a partire dalle parole elencate nel vocabolario. Inoltre le regole sintattiche debbono essere puramente formali, cioè debbono riferirsi esclusivamente alla forma (forma esteriore) delle espressioni; la funzione e il significato di una espressione debbono dipendere esclusivamente dalla sua forma. In particolare, osservando un'espressione si deve essere in grado in ciascun caso di decidere se è una proposizione o no; non deve mai capitare che un'espressione assolva da qualche parte la funzione di proposizione mentre un'espressione della stessa forma non si comporti così da qualche altra parte, oppure che una proposizione possa essere asserita in un certo contesto mentre una proposizione della stessa forma possa essere negata in un altro. Ne segue, in particolare, che pronomi e avverbi dimostrativi quali « questo » e « qui » non dovranno comparire nel vocabolario del linguaggio; i linguaggi che soddisfano a queste condizioni verranno detti linguaggi formalizzati. Quando si studia un linguaggio formalizzato non è necessario distinguere fra espressioni della stessa forma che sono state scritte o pronunciate in luoghi diversi; spesso se ne parla come di una stessa espressione. Il lettore può aver notato che talvolta si usa questo modo di parlare anche nel trattare un linguaggio naturale, cioè che non sia formalizzato;

si fa così per semplicità e solo in quei casi in cui non sembra esserci pericolo di confusione.

I linguaggi formalizzati sono del tutto adeguati per la presentazione di teorie logiche e matematiche; non vedo alcuna ragione essenziale perché non possano venire adattati sì da poter essere impiegati in altre discipline scientifiche, e in particolare nello sviluppo delle parti teoriche delle scienze empiriche. Vorrei sottolineare che nell'usare il termine « linguaggi formalizzati » non mi riferisco esclusivamente a quei sistemi linguistici che sono formulati completamente in simboli, né penso a qualcosa di essenzialmente opposto ai linguaggi naturali. Al contrario, gli unici linguaggi formalizzati che sembrano essere di un qualche interesse reale sono quelli che sono frammenti dei linguaggi naturali (frammenti provvisti di un vocabolario completo e di precise regole sintattiche) o quelli che almeno possano essere adeguatamente tradotti nei linguaggi naturali.

Ci sono alcune condizioni ulteriori dalle quali dipende la realizzazione del nostro programma. Si dovrebbe fare una rigida distinzione fra il linguaggio che è l'oggetto del nostro studio e per il quale in particolare si vuole costruire la definizione di verità e il linguaggio nel quale la definizione deve essere formulata e le sue implicazioni studiate. Quest'ultimo verrà detto metalinguaggio e il primo linguaggio oggetto. Il metalinguaggio deve essere sufficientemente ricco; in particolare deve contenere come parte il linguaggio oggetto. Infatti secondo le nostre convenzioni una definizione adeguata di verità implicherà come conseguenze tutte le definizioni parziali di tale concetto, cioè tutte le equivalenze della forma (3):

$\langle p \rangle$ è vera se e solo se p ,

dove « p » va sostituita in ambo i membri dell'equivalenza con una proposizione arbitraria del linguaggio oggetto. Giacché tutte queste conseguenze sono formulate nel metalinguaggio, se ne conclude che ogni proposizione del linguaggio oggetto dev'essere anche una proposizione del metalinguaggio. Inoltre, il metalinguaggio deve contenere nomi per proposizioni (e altre espressioni) del linguaggio oggetto, giacché questi nomi compaiono nel primo membro delle equivalenze del tipo (3). Esso deve anche contenere alcuni altri termini che occorrono per lo studio del linguaggio oggetto, e precisamente termini che denotino certi particolari insiemi di espressioni, relazioni fra espressioni e operazioni sulle espressioni; per

esempio, deve esserci la possibilità di parlare dell'insieme di tutte le proposizioni o dell'operazione di giustapposizione, mediante la quale da due espressioni se ne forma una nuova ponendo una delle due immediatamente di seguito all'altra. Infine, nel definire la verità, si vede che i termini semantici (quelli che esprimono le relazioni tra le proposizioni del linguaggio oggetto e gli oggetti a cui tali proposizioni si riferiscono) possono essere introdotti nel metalinguaggio mediante definizioni. Se ne conclude che un metalinguaggio che fornisca mezzi sufficienti a definire la verità deve essere essenzialmente più ricco del linguaggio oggetto; esso non può coincidere né essere traducibile in quest'ultimo, giacché altrimenti ambedue risulterebbero semanticamente universali, e l'antinomia del mentitore sarebbe ricostruibile in entrambi. Ritorneremo su questa questione nell'ultimo paragrafo del presente articolo.

Se tutte le precedenti condizioni sono soddisfatte, la costruzione della desiderata definizione di verità non presenta difficoltà essenziali. Tecnicamente, tuttavia, essa è troppo complicata per essere esposta qui in dettaglio. Per ogni data proposizione del linguaggio oggetto si può facilmente formulare la corrispondente definizione parziale della forma (3). Tuttavia, giacché l'insieme di tutte le proposizioni del linguaggio oggetto è di regola infinito, mentre ogni proposizione del metalinguaggio è una successione finita di segni, non si può arrivare a una definizione generale formando semplicemente la congiunzione logica di tutte le definizioni parziali. Nondimeno, ciò che alla fine si ottiene è, in un senso intuitivo, equivalente alla immaginaria congiunzione infinita. Molto approssimativamente, si procede come segue. Dapprima si considerano le proposizioni più semplici, che non contengono altre proposizioni come parti; per queste proposizioni si trova il modo di definire la verità direttamente (usando la stessa idea che conduce alle definizioni parziali). Poi, mediante l'uso delle regole sintattiche che riguardano la formazione di proposizioni più complicate a partire da quelle più semplici, si estende la definizione a proposizioni composte arbitrarie; si applica qui il metodo conosciuto in matematica come definizione per ricursione. (Questa è, a dire il vero, una grossolana approssimazione del procedimento. Per ragioni tecniche il metodo di ricursione si applica, in realtà, non per definire il concetto di verità, ma quello a esso collegato di soddisfazione; la verità viene poi facilmente definita in termini di soddisfazione.)

Sulla base della definizione così costruita si può sviluppare l'intera teoria della verità. In particolare si possono derivare da essa, oltre a tutte le equivalenze della forma (3), alcune conseguenze di carattere generale, quali il famoso principio di non contraddizione e quello del terzo escluso. Per il primo di questi due principi, non possono essere entrambe vere due proposizioni una delle quali sia la negazione dell'altra; per il secondo principio, due proposizioni siffatte non possono essere ambedue false.

Il concetto di dimostrazione

Qualunque cosa possa ottenersi dalla costruzione di una definizione adeguata del concetto di verità per un linguaggio scientifico, una cosa è certa: la definizione non porta con sé un criterio pratico per decidere se una particolare proposizione di tale linguaggio sia vera o falsa (e inverso questo non è affatto il suo scopo). Si consideri per esempio la seguente proposizione nel linguaggio della geometria elementare: « le tre bisettrici di un triangolo passano per uno stesso punto ». Se ci interessa sapere se questa proposizione è vera e ci rifacciamo alla definizione di verità per scoprirlo, siamo destinati ad avere una delusione. L'unica informazione che ricaviamo è che la proposizione è vera se le tre bisettrici di un triangolo si incontrano sempre in un punto, e falsa in caso contrario; solo un'indagine di natura geometrica ci permetterà di decidere come stanno le cose in realtà. Considerazioni analoghe valgono per proposizioni tratte dal dominio di altre scienze particolari: è compito della scienza stessa scoprire se una tale proposizione è vera o falsa, e non della logica o della teoria della verità.

Alcuni filosofi ed epistemologi sono propensi a rifiutare ogni definizione che non fornisca un criterio per decidere, per ciascun oggetto particolare assegnato, se cada o no sotto il concetto definito. Nella metodologia delle scienze empiriche tale tendenza è rappresentata dall'operazionismo e è condivisa anche da quei filosofi della matematica che appartengono alla scuola costruttivista; in ambedue i casi, tuttavia, solo una piccola minoranza dei pensatori sono di questa opinione. Non è mai stato fatto un tentativo organico per condurre a termine, in pratica, il programma di sviluppare una scienza senza l'uso di definizioni indesiderabili. È chiaro che seguendo un programma del genere molta della matematica contemporanea scomparirebbe e resterebbero gravemente mutilati anche

gli aspetti teorici della fisica, della chimica, della biologia e delle altre scienze empiriche. Nella definizione di concetti come atomo, gene, eccetera, così come nella maggior parte delle definizioni matematiche, non è implicito alcun criterio per decidere se un oggetto cade o no sotto il termine definito.

Proprio per la mancanza di un tale criterio di decisione nella definizione di verità, mentre la ricerca della verità è giustamente considerata l'essenza delle attività scientifiche, diventa un problema importante trovare almeno dei criteri parziali di verità e sviluppare procedimenti che ci permettano di asserire o negare la verità (o per lo meno la probabilità della verità) del maggior numero possibile di proposizioni. E, in realtà, esistono procedimenti di questo tipo, alcuni usati esclusivamente nelle scienze empiriche, altri prevalentemente nelle scienze deduttive. Il concetto di dimostrazione – il secondo preso in considerazione nel presente lavoro – si riferisce appunto a un procedimento per verificare la verità delle proposizioni, che viene impiegato prevalentemente nelle scienze deduttive; esso è un elemento essenziale di ciò che è noto come metodo assiomatico, il solo metodo ormai usato per sviluppare le discipline matematiche.

Il metodo assiomatico e, all'interno di esso, il concetto di dimostrazione, sono frutto di un lungo sviluppo storico; per capire l'odierno concetto di dimostrazione è forse essenziale una conoscenza, sia pure rudimentale, di tale sviluppo.

In origine, una disciplina matematica era considerata come un aggregato di proposizioni, relative a una certa classe di oggetti o fenomeni, che venivano formulate a partire da certi termini iniziali e che erano considerate vere. Tale aggregato di proposizioni non aveva alcuna struttura interna; una proposizione veniva accettata come vera o perché intuitivamente evidente o perché era dimostrata a partire da alcune proposizioni intuitivamente evidenti, risultando quindi una conseguenza di queste ultime sulla base di un ragionamento intuitivamente sicuro. Il criterio dell'evidenza intuitiva e dell'intuitiva certezza dei ragionamenti veniva applicato senza restrizioni; ogni proposizione riconosciuta per vera sulla base di tale criterio veniva automaticamente inclusa nella disciplina. La descrizione testé fornita sembra adeguata, per esempio, alla scienza della geometria com'era nota agli antichi Egizi e ai Greci nel suo stadio preeuclideo.

Presto ci si accorse, però, che il criterio dell'evidenza intuitiva è lungi dal-

l'essere infallibile, non ha alcun carattere oggettivo e spesso conduce a seri errori. Tutto lo sviluppo successivo del metodo assiomatico può considerarsi come espressione della tendenza a limitare l'uso dell'evidenza intuitiva.

Tale tendenza si rivelò la prima volta nello sforzo di dimostrare il maggior numero possibile di proposizioni e quindi di limitare il più possibile il numero delle proposizioni accettate per vere sulla sola base dell'evidenza. L'ideale, da questo punto di vista, sarebbe quello di dimostrare ogni proposizione che vada accettata per vera, ma per ovvie ragioni un tale ideale non può essere raggiunto. Infatti, una proposizione si dimostra a partire da altre, queste a partire da altre ancora e così via: se vogliamo evitare sia un circolo vizioso sia un regresso all'infinito, siamo costretti a interrompere la trafilata in qualche punto. Come compromesso fra l'ideale irraggiungibile e le possibilità effettive, emersero due principi che furono successivamente usati per costruire discipline matematiche. Per il primo di questi principi ogni teoria comincia con un piccolo elenco di proposizioni che appaiono intuitivamente evidenti e che sono accettate per vere senza ulteriori giustificazioni. Secondo l'altro principio, nessun'altra proposizione viene accettata per vera nella disciplina se non si riesce a dimostrarla col solo ausilio degli assiomi e delle proposizioni dimostrate in precedenza.

Tutte le proposizioni che si possono considerare vere grazie a questi due principi vengono chiamate teoremi o proposizioni dimostrabili della data teoria. Due principi analoghi regolano l'uso e la costruzione dei termini della teoria: per il primo, si fa all'inizio un elenco di quei termini, intelligibili direttamente, che si usano nell'enunciato e nella dimostrazione dei teoremi, senza spiegarne il significato: essi vengono detti termini primitivi o termini non definiti; per il secondo principio si conviene di non usare alcun altro termine il cui significato non possa spiegarsi mediante esplicita definizione a partire dai termini primitivi e da quelli definiti precedentemente. Questi quattro principi sono le pietre miliari del metodo assiomatico, e le teorie sviluppate in accordo con essi vengono chiamate teorie assiomatiche.

È ben noto che il metodo assiomatico fu applicato allo sviluppo della geometria verso il 300 a.C. negli *Elementi* di Euclide. Da allora esso fu usato per oltre 2000 anni praticamente senza alcun cambiamento dei suoi principi fondamentali (i quali, per inciso, non furono nemmeno esplicitamente formu-

lati per lungo tempo) né dell'atteggiamento di fondo nell'affrontare l'argomento. Nel XIX e nel XX secolo, però, il metodo assiomatico subì una profonda evoluzione, la quale, per quanto riguarda il concetto di dimostrazione, è particolarmente significativa per la nostra trattazione.

Fin verso la fine del XIX secolo il concetto di dimostrazione rivestì un carattere prevalentemente psicologico: una dimostrazione era un'attività intellettuale volta a convincere se stesso e gli altri della verità di una data proposizione; più precisamente, nello sviluppare una teoria matematica, le dimostrazioni venivano usate per convincere se stessi e gli altri che una data proposizione doveva essere accettata per vera una volta che fossero state prese per vere certe proposizioni precedenti. Per i ragionamenti usati nelle dimostrazioni non veniva posta alcuna condizione, salvo quella di essere intuitivamente convincenti. A un certo momento, però, cominciò a farsi sentire il bisogno di sottoporre il concetto di dimostrazione a un'analisi più approfondita, con l'eventuale risultato di dover limitare il ricorso all'evidenza intuitiva, sia pur rimanendo nel vecchio schema. La cosa fu dovuta, probabilmente, a certi particolari sviluppi della matematica, soprattutto alla scoperta delle geometrie non euclidee. L'analisi fu condotta dai logici, a partire dal logico tedesco Gottlob Frege; si giunse così a introdurre un nuovo concetto, quello di dimostrazione formale, che rappresentò un sostituto adeguato del vecchio concetto psicologico e anche un essenziale miglioramento.

Il primo passo per corredare una teoria matematica del concetto di dimostrazione formale è la formalizzazione del suo linguaggio, nel senso spiegato sopra in relazione alla definizione di verità; si danno così delle regole sintattiche formali che permettano, in particolare, di distinguere, a partire solo dalla forma esteriore delle espressioni, fra le espressioni stesse quelle che sono proposizioni. Il passo successivo consiste nel formulare alcune regole di natura diversa, le cosiddette regole di deduzione o di inferenza. In virtù di queste regole una proposizione viene considerata direttamente deducibile da certe proposizioni date se, in generale, la sua forma è collegata alla forma delle altre in una maniera assegnata. Il numero delle regole di inferenza è piccolo, e il loro contenuto semplice; proprio come per il caso delle regole sintattiche, esse hanno un carattere formale, cioè si riferiscono esclusivamente alla forma esteriore delle proposizioni a cui si applicano. Intuitivamente,

tutte le regole di deduzione sono chiaramente infallibili, nel senso che risulta necessariamente vera ogni proposizione che sia direttamente derivabile, in virtù di una di tali regole, da proposizioni vere; l'infallibilità delle regole di deduzione può essere stabilita, in pratica, sulla base di una definizione adeguata di verità. Il più noto e importante esempio di regola di inferenza è la regola di separazione, nota anche come *modus ponens*. In virtù di questa regola (che in alcune teorie è l'unica regola di inferenza) una proposizione « q » è derivabile direttamente da due date proposizioni se una di esse è la proposizione condizionale « se p , allora q », mentre l'altra è « p »; qui « p » e « q » sono, come al solito, abbreviazioni di due qualsivogliano proposizioni del nostro linguaggio formalizzato. Siamo ora in grado di spiegare in cosa consiste una dimostrazione formale di una data proposizione. In un primo momento si applicano le regole di inferenza agli assiomi e si ottengono nuove proposizioni che sono derivabili direttamente dagli assiomi, poi si applicano le solite regole al nuovo insieme di proposizioni così ottenuto e si ottengono ancora altre proposizioni, e così via. Di una proposizione che sia stata ottenuta per questa via dopo un numero finito di passi si dice che è stata dimostrata formalmente. Quanto precede si può esprimere in modo più preciso come segue: una dimostrazione formale di una data proposizione consiste nel costruire una successione finita di proposizioni tale che (1) la prima proposizione della successione è un assioma, (2) ciascuna delle proposizioni seguenti è un assioma oppure è derivabile direttamente, mediante una regola di inferenza, da alcune delle proposizioni che la precedono nella successione, e (3) l'ultima proposizione della successione è quella che si voleva dimostrare. Forzando un po' l'uso della parola « dimostrazione », si può dire addirittura che una dimostrazione formale è proprio una successione finita di proposizioni con le proprietà sopra elencate.

Si chiama teoria formale una teoria assiomatica il cui linguaggio sia stato formalizzato, e per la quale si sia data una definizione di dimostrazione formale. Conveniamo di accettare come uniche dimostrazioni in una teoria formale le sue dimostrazioni formali; quindi le uniche proposizioni accettabili come teoremi saranno quelle che compaiono nell'elenco degli assiomi e quelle per cui possa trovarsi una dimostrazione formale. Il metodo che consiste nel presentare una teoria formalizzata in ognuna delle fasi del suo sviluppo è, in

linea di principio, molto elementare. Prima si elencano gli assiomi, poi tutti i teoremi noti, in un ordine tale che ciascuna proposizione dell'elenco che non sia un assioma possa venire riconosciuta direttamente come teorema semplicemente confrontando la sua forma con quella delle proposizioni che la precedono nell'elenco; viene così escluso ogni ricorso a complicate vie di ragionare e di convincere. (Naturalmente non mi riferisco qui ai processi psicologici mediante i quali i teoremi sono stati scoperti nella realtà.) Il ricorso all'evidenza intuitiva viene così a essere notevolmente limitato; non che si siano completamente eliminati i dubbi relativi alla verità dei teoremi, ma ci siamo almeno ridotti solo agli eventuali dubbi sulla verità delle poche proposizioni elencate come assiomi e sulla infallibilità delle poche e semplici regole di deduzione. Anche per quanto riguarda l'introduzione di nuovi termini nel linguaggio, si potrebbe usare un metodo formale, fornendo certe regole formali per le definizioni.

È ormai noto che tutte le teorie matematiche esistenti possono venire presentate come teorie formali, e che si possono trovare delle dimostrazioni formali per quei teoremi matematici, per profondi e complicati che siano, che in origine furono stabiliti con ragionamenti intuitivi.

Le relazioni fra verità e dimostrazione

È stato senza dubbio un grande progresso della logica moderna l'aver sostituito al vecchio concetto psicologico di dimostrazione, che non poteva certo essere reso, chiaro e rigoroso, un nuovo concetto¹ che avesse sia la dote della semplicità sia un carattere puramente formale. Tuttavia nel trionfo del metodo formale era già insito il germe di un futuro regresso: come ora vedremo, proprio la semplicità del nuovo concetto era il suo tallone di Achille.

Per dare un giudizio di valore sul concetto di dimostrazione formale, dobbiamo chiarirne i rapporti col concetto di verità. In fondo, la dimostrazione formale, come del resto la vecchia dimostrazione intuitiva, è un procedimento che mira ad acquisire nuove proposizioni vere. Un tale procedimento si considererà perciò adeguato solo se tutte le proposizioni che possono acquisirsi con il suo ausilio risulteranno vere e, viceversa, tutte le proposizioni vere saranno ottenibili per suo mezzo. Ne sorge, in modo spontaneo, il problema se la dimostrazione formale sia in effetti un procedimento adeguato per acquisire la verità; in altre parole, se l'insieme delle proposizioni dimostrabili

li (formalmente) coincida con l'insieme delle proposizioni vere.

Studieremo questo problema, per semplicità, relativamente a una particolare teoria matematica molto elementare, cioè all'aritmetica dei numeri naturali (teoria elementare dei numeri). Supponiamo che la teoria in questione venga presentata come teoria formale. Il suo vocabolario è molto povero e consta di variabili, quali « m », « n », « p », ... che rappresentano numeri naturali; dei numerali « 0 », « 1 », « 2 », ... che denotano numeri particolari; di simboli che denotano alcune relazioni e operazioni familiari sui numeri, quali « $=$ », « $<$ », « $+$ », « $-$ »; e infine di certi termini logici, e precisamente i connettivi proposizionali (« e », « oppure », « se... allora », « non ») e i quantificatori (espressioni della forma « per ogni numero m » e « per qualche numero m »). Le regole sintattiche e le regole di inferenza sono semplici. Nel seguito, parlando di proposizioni, ci riferiremo costantemente alle proposizioni del linguaggio formalizzato dell'aritmetica.

Da quanto esposto nel primo paragrafo sulla verità, risulta che si può prendere il linguaggio formalizzato dall'aritmetica come linguaggio oggetto e costruire un opportuno metalinguaggio nel quale si possa formulare una definizione adeguata di verità. Risulta comodo, in tale contesto, dire che si è così definito l'insieme delle proposizioni vere; la definizione di verità esprime infatti che una certa condizione formulata nel metalinguaggio è soddisfatta da tutti e soli gli elementi di questo insieme, cioè da tutte e sole le proposizioni vere. Ancor più facilmente si può definire, nel metalinguaggio, l'insieme delle proposizioni dimostrabili, attenendosi alla spiegazione del concetto di dimostrazione formale data nel secondo paragrafo. A rigore, sia la definizione di verità sia quella di dimostrabilità appartengono a una nuova teoria formulata nel metalinguaggio e costruita col preciso scopo di studiare l'aritmetica formalizzata e il suo linguaggio. La nuova teoria è chiamata metateoria o, nel caso in esame, metaaritmetica. Non approfondiremo qui il modo in cui la metateoria, con i suoi assiomi, termini primitivi eccetera, viene costruita; ci limiteremo a far notare che essa rappresenta il quadro nel quale si formula e si risolve il problema del confronto fra l'insieme delle proposizioni dimostrabili e quello delle proposizioni vere.

La soluzione del problema risulta negativa, e qui daremo una rapida descrizione del metodo che ha permesso di dimostrarlo. L'idea principale è stret-

tamente connessa con quella usata dal logico contemporaneo americano (di origine austriaca) Kurt Gödel nella sua famosa pubblicazione sulla incompletezza dell'aritmetica.

Nel primo paragrafo abbiamo osservato che il metalinguaggio che ci permette di definire e studiare il concetto di verità deve essere ricco. Esso contiene tutto il linguaggio oggetto come sua parte e quindi si può parlare in esso dei numeri naturali, delle relazioni fra numeri naturali e così via. Esso contiene inoltre dei termini necessari per lo studio del linguaggio oggetto e delle sue componenti; di conseguenza si può parlare, nel metalinguaggio, delle espressioni, e in particolare delle proposizioni e così via. Nella metateoria si ha dunque la possibilità di studiare le proprietà di tali oggetti e stabilirne i mutui rapporti.

Facendo ricorso alla descrizione delle proposizioni fornita dalle regole sintattiche del linguaggio oggetto è facile, in particolare, disporre tutte le proposizioni in successione in ordine progressivo di complessità e numerarle successivamente. Si viene ad associare, così, un numero naturale a ciascuna proposizione, in modo tale che a proposizioni diverse vengano sempre associati numeri diversi; in altre parole, si stabilisce una corrispondenza biunivoca fra le proposizioni e i numeri. Questa a sua volta induce una corrispondenza biunivoca fra gli insiemi di proposizioni e gli insiemi di numeri, come pure fra le relazioni fra proposizioni e le relazioni fra numeri. In particolare, ci saranno dei numeri associati a proposizioni dimostrabili e altri associati a proposizioni vere; per brevità li chiameremo rispettivamente numeri dimostrabili* e numeri veri*. Il problema che ci interessa si riduce quindi al confronto fra l'insieme dei numeri dimostrabili* e quello dei numeri veri*.

Per mostrare che i due insiemi non coincidono, sarà ovviamente sufficiente indicare una proprietà di cui goda uno dei due insiemi e non l'altro. La proprietà che indicheremo noi potrà sembrare sorprendente, quasi un *deus ex machina*.

La semplicità intrinseca dei concetti di dimostrazione formale e di dimostrabilità formale rivestirà qui un ruolo fondamentale. Si è visto nel secondo paragrafo che il significato di questi concetti viene spiegato essenzialmente in termini di certe relazioni semplici fra proposizioni assegnate da poche regole di inferenza, come per esempio la regola del *modus ponens*. Le corrispondenti relazioni fra i numeri delle proposizioni sono altrettanto semplici; esse possono essere descritte me-

diate le operazioni e relazioni aritmetiche più semplici, quali l'addizione, la moltiplicazione e l'uguaglianza, cioè mediante termini che compaiono esplicitamente nella teoria aritmetica. Di conseguenza, anche l'insieme dei numeri dimostrabili* può essere descritto mediante tali termini. Si può riassumere quanto precede dicendo che la definizione di dimostrabilità è stata tradotta dal metalinguaggio nel linguaggio oggetto.

D'altro lato, l'esposizione che abbiamo fatto relativamente al concetto di verità nei linguaggi correnti suggerisce la congettura che un'analoga traduzione non sia possibile per la definizione di verità, altrimenti il linguaggio oggetto risulterebbe in un certo senso semanticamente universale col risultato che potrebbe ricomparire l'antinomia del mentitore. Convalideremo tale congettura mostrando che effettivamente sarebbe possibile riformulare in questo linguaggio l'antinomia del mentitore qualora l'insieme dei numeri veri* fosse definibile nel linguaggio dell'aritmetica. Tuttavia, avendo ora a che fare con un linguaggio formalizzato limitato, l'antinomia assumerebbe qui una forma più involuta e complicata; in particolare, nella nuova formulazione non comparirebbe alcuna espressione il cui significato empirico fosse del tipo « la proposizione stampata nel punto tal dei tali », mentre nella formulazione iniziale dell'antinomia è stata proprio un'espressione di questo tipo a rivestire un ruolo fondamentale. In questo articolo non ci addentreremo però ulteriormente nei dettagli tecnici della questione.

In conclusione, l'insieme dei numeri dimostrabili* non coincide con quello dei numeri veri*, giacché mentre l'uno è definibile nel linguaggio dell'aritmetica, l'altro non lo è. Ne segue che sono distinti anche gli insiemi delle proposizioni dimostrabili e vere. Ricorrendo alla definizione di verità, si verifica d'altronde che tutti gli assiomi dell'aritmetica sono veri e tutte le regole di inferenza sono infallibili, quindi tutte le proposizioni dimostrabili sono vere e perciò non vale il viceversa. Abbiamo così raggiunto la conclusione che: esistono delle proposizioni vere formulate nel linguaggio dell'aritmetica le quali non possono essere dimostrate a partire dagli assiomi e dalle regole di inferenza accettate in aritmetica.

Si potrebbe pensare che la conclusione precedente dipenda in modo essenziale dalla scelta degli assiomi e delle regole di inferenza particolari scelte per la nostra teoria aritmetica: arricchendo opportunamente la teoria con l'aggiunta di nuovi assiomi e nuove regole

di inferenza, si potrebbe forse giungere a un risultato diverso. Una più attenta analisi mostra, invece, che il precedente ragionamento dipende ben poco dalle proprietà specifiche della teoria esaminata, essendo possibile estenderlo alla maggior parte delle altre teorie formali. Supposto che in una teoria sia contenuta l'aritmetica dei numeri naturali, o almeno che questa si possa ricostruire all'interno di quella, allora è possibile ripetere la parte essenziale del precedente ragionamento in una forma praticamente immutata, giungendo così ancora alla conclusione che l'insieme delle proposizioni dimostrabili è diverso da quello delle proposizioni vere. Se inoltre (e questo è il caso più frequente) tutti gli assiomi della teoria sono veri e tutte le regole di inferenza infallibili, si conclude ancora che esistono nella teoria delle proposizioni vere che non sono dimostrabili. Se si escludono alcune teorie parziali con limitati mezzi espressivi, in genere l'ipotesi precedente sulle relazioni fra una teoria e l'aritmetica è soddisfatta, perciò la nostra conclusione ha un carattere quasi universale. Per quanto riguarda quelle teorie povere che non comprendono l'aritmetica dei numeri naturali, può accadere che il relativo linguaggio non abbia mezzi di espressione sufficienti a definire il concetto di dimostrabilità, e talvolta le proposizioni dimostrabili coincidano proprio con quelle vere. I più noti esempi, e forse i più importanti, di teorie per le quali i due concetti coincidono sono la geometria elementare e l'algebra elementare dei numeri reali.

La parte predominante sostenuta in tutto il ragionamento dall'antinomia del mentitore chiarisce in modo interessante le precedenti osservazioni sul ruolo delle antinomie nella storia del pensiero umano. All'inizio l'antinomia del mentitore si è presentata nella nostra trattazione come una specie di forza del male con un enorme potere di distruzione, costringendoci ad abbandonare ogni tentativo di chiarire il concetto di verità per i linguaggi naturali. Abbiamo dovuto limitare i nostri sforzi ai linguaggi formalizzati del discorso scientifico e, come salvaguardia contro la possibilità che l'antinomia rifacesse la sua comparsa, abbiamo dovuto complicare considerevolmente la indagine con la distinzione fra linguaggio e metalinguaggio.

In un secondo momento, però, nella nuova sistemazione, siamo riusciti a domare l'energia distruttiva incanalandola verso scopi pacifici e costruttivi; l'antinomia non si è riaffacciata, e l'idea di fondo in essa contenuta è stata anzi usata per stabilire un significativo ri-

sultato metalogico con notevoli e profonde conseguenze.

Il fatto che le implicazioni filosofiche del nostro risultato siano di carattere sostanzialmente negativo non ne diminuisce affatto l'importanza. Il risultato mostra, invero, che in nessun campo della matematica il concetto di dimostrabilità è un sostituto perfetto di quello di verità; la credenza che la dimostrazione formale possa costituire uno strumento adeguato per stabilire la verità di tutti gli enunciati matematici si è rilevata infondata, e al trionfo iniziale ha fatto seguito un regresso.

Qualunque aggiunta si facesse ora alla precedente esposizione non potrebbe che diminuirne l'effetto. Il concetto di verità per le teorie formalizzate può ora introdursi in modo preciso mediante una definizione adeguata e può quindi venire usato senza riserve né limi-

tazioni nelle trattazioni metalogiche, giacché è addirittura diventato un concetto metalogico fondamentale che entra in problemi e risultati fra i più importanti. D'altra parte, neanche il concetto di dimostrazione ha perduto la sua importanza; la dimostrazione resta l'unico metodo usato per verificare la verità delle proposizioni nell'ambito di una teoria matematica particolare. Tuttavia, abbiamo ora acquisito la consapevolezza che esistono delle proposizioni, formulate nel linguaggio della teoria, vere ma non dimostrabili, e non possiamo trascurare l'eventualità che fra queste se ne trovi qualcuna di quelle interessanti che si desidererebbe dimostrare. In alcune situazioni può essere dunque desiderabile indagare sulla possibilità di ampliare l'insieme delle proposizioni dimostrabili; allo scopo si arricchisce la data teoria aggiun-

gendo nuove proposizioni al suo sistema di assiomi o dotandola di nuove regole di inferenza. In tale processo serve da guida il concetto di verità, giacché vogliamo evitare di aggiungere assiomi che si sospettino falsi o regole di inferenza che, applicate a proposizioni vere, diano luogo a proposizioni false. L'operazione di estendere una teoria può naturalmente essere ripetuta quante volte si vuole; ne risulta che il concetto di proposizione vera funziona come limite ideale che non potrà mai essere raggiunto ma che si tenta di approssimare gradualmente mediante successivi ampliamenti dell'insieme delle proposizioni dimostrabili. Non c'è alcun conflitto fra il concetto di verità e quello di dimostrazione nello sviluppo della matematica; i due concetti non sono in guerra fra loro, ma vivono in coesistenza pacifica.



Charles Lutwidge Dodgson nacque il 27 gennaio 1832 e morì il 14 gennaio 1898. Il logico che sotto lo pseudonimo di Lewis Carroll scrisse *Alice nel paese delle meraviglie*, era anche un appassionato fotografo dilettante specializzato in ritratti.

Questo proviene dalla collezione Gernsheim dell'Humanities Research Center dell'Università del Texas di Austin. Lo sviluppo è stato ottenuto dal negativo originale della collezione Gernsheim. Il negativo porta il numero 2439, scritto a mano da Dodgson.

Un libro di logica smarrito di Lewis Carroll

Si sa che l'autore di Alice scrisse anche un testo di logica simbolica. È stata recentemente scoperta una continuazione di questo libro che avvalorava l'ipotesi della profonda originalità della sua opera di logico

di W. W. Bartley III

«Sono così contenta che non mi piacciono gli asparagi – disse la Piccola Ragazza a un Amico Simpatico. – Perché, se mi piaceranno, dovrei mangiarli mentre non li posso sopportare!».

Queste parole dal suono familiare furono scritte dal reverendo Charles Lutwidge Dodgson, incaricato di matematica a Oxford e universalmente noto come Lewis Carroll. Sarebbe vano tuttavia cercarle in qualche lavoro di Dodgson già pubblicato. La Piccola Ragazza e il suo Amico Simpatico, come Achille e la Tartaruga, il Coccodrillo e il Mentitore, i Tre Barbieri, i Cinque Bugiardi e il Logico-giocatore d'azzardo divoratore di costolette di maiale, sono alcuni tra i molti personaggi, in parte già noti e in parte nuovi, che aguzzano il loro ingegno – e il nostro – nel manoscritto e nelle bozze di stampa recentemente scoperte con cui Dodgson dava seguito alla sua *Logica simbolica: Parte I, Nozioni elementari*, pubblicata nel 1896. Dal titolo è chiaro che Dodgson aveva in programma la pubblicazione di nuovo materiale sull'argomento, ma il manoscritto al quale stava lavorando scomparve poco dopo la morte dell'autore, sopraggiunta nel gennaio 1898, all'età di 65 anni. Di esso non si trova cenno in nessuno dei minuziosi elenchi delle opere della «Carrolliana» pubblicati nell'ultimo mezzo secolo, e la maggior parte degli incartamenti fu bruciata poco dopo la sua morte.

Fu nel 1959 che rintracciai, inizialmente, una piccola parte dell'opera mancante in alcune carte superstiti di Dodgson rimaste alla Christ Church di Oxford. Il sospirato frammento era accuratamente composto a caratteri di stampa, sotto forma di bozze. Nei 10 anni successivi sono andato alla ricerca di altre parti dell'opera nelle colle-

zioni, pubbliche e private, degli scritti di Lewis Carroll. Infine, nella vasta collezione dei manoscritti e delle lettere di Dodgson raccolta da Morton N. Cohen, professore di inglese al Centro per laureati dell'Università di New York, ho rinvenuto alcune fotocopie di altre porzioni delle bozze. Cohen le aveva tratte dagli originali appartenenti alla biblioteca di John H. A. Sparrow, direttore dell'All Souls College di Oxford. Sparrow aveva ricevuto i fogli dal defunto A. S. L. Farquharson, che aveva curato la pubblicazione degli scritti postumi di John Cook Wilson. Quando era professore di logica a Oxford, Wilson, il 6 novembre 1896, aveva ricevuto per posta le bozze da Dodgson stesso, e si era evidentemente dimenticato di rispeditargliele. Cohen e Roger Lancelyn Green stanno ora curando un'edizione definitiva delle lettere di Dodgson. Il materiale finora inedito che qui presentiamo è coperto da diritti d'autore riservati di Charles Lutwidge Dodgson. Attualmente sto allestendo una edizione critica dell'opera logica completa di Dodgson.

Accanto agli squarci di lettere e manoscritti rimasti finora in larga misura indecifrabili, le pagine di manoscritto e di bozze ora recuperate confermano il giudizio formulato su Dodgson in via congetturale da alcuni storici della matematica in base agli scarsi elementi forniti dal primo volume della *Logica Simbolica* e da diversi scritti pubblicati sulla rivista «Mind».

Un decennio dopo la morte di Dodgson la sua opera fu offuscata dalla rivoluzione provocata nella logica dalla pubblicazione dei *Principia Mathematica* di Alfred North Whitehead e Bertrand Russell. La seconda parte della *Logica Simbolica* dimostra che Dodgson fu uno dei più interessanti innova-

tori tecnici del periodo di transizione dalla logica tradizionale di scuola aristotelica alla nuova logica propugnata da Russell. Essa conferma anche che Dodgson era senza rivali nel proporre problemi, rompicapi e paradossi. Ciò è tanto più stupefacente se si considera che l'opera di rinnovamento in campi come quello della logica è compiuta abitualmente da giovani mentre Dodgson ha prodotto la maggior parte del suo lavoro nell'ambito della logica a 60 anni compiuti. Egli lavorava da solo; il solo logico con cui manteneva contatti regolari era Wilson, che peraltro era ben poco stimolante. Wilson fu un accanito avversario della nuova logica simbolica sviluppata da Dodgson e da altri. Più tardi Wilson disse che non riusciva a credere che Russell, la cui opera qualificava come «roba spregevole», potesse trovare un editore.

La natura del trapasso dalla logica aristotelica alla logica matematica contemporanea è talvolta fraintesa. Non manca chi pensa erroneamente che la logica aristotelica sia stata dimostrata «sbagliata» e che sia stata soppiantata dalla logica contemporanea nello stesso modo in cui una nuova ipotesi scientifica può soppiantarne una preesistente. La differenza tra la logica tradizionale e la logica contemporanea è di natura differente.

I logici cercano di formulare «regole del ragionamento valido» che ci assicureranno la possibilità di trarre solo conclusioni vere da premesse vere. Un ragionamento è valido quando e solo quando non si può produrre nessun controesempio. Un controesempio si produce se si può argomentare da un insieme di premesse vere a una conclusione falsa seguendo le regole formulate. L'obiettivo è quello di evitare i ragionamenti non validi e le regole di inferenza che li rendono possibili.

BOOK XXI.
LOGICAL PUZZLES.

CHAPTER I.
INTRODUCTORY.

UNDER this general heading I shall discuss various arguments, which are variously described by Logical writers. Some have been classified as 'Sophisms', that is, according to etymology, "cunning arguments", whose characteristic Attribute seems to be that they are intended to *confuse*: others as 'Paradoxes', that is, according to etymology, "things contrary to expectation", whose characteristic Attribute seems to be that they seem to prove what we know to be false: but all may be described by the general name "Puzzles."

CHAPTER II.
CLASSICAL PUZZLES

§ 1.

Introductory.

I SHALL here enuntiate five certain well-known Puzzles, which have come down to us from ancient times, and which the Reader will no doubt like to know by their classical titles

§ 2.

Pseudomenos.

This may also be described as "*Mentiens*", or "*The Liar*". In its simplest form it runs thus:—

"If a man says 'I am telling a lie', and speaks truly, he *is* telling a lie, and therefore speaks falsely: but if he speaks falsely, he is *not* telling a lie, and therefore speaks truly"

§ 3.

Crocodilus.

That is, "*The Crocodile*". This tragical story runs as follows:—

"A Crocodile had stolen a Baby off the banks of the Nile. The Mother implored him to restore her darling. 'Well', said the Crocodile, 'if you say truly what I shall do, I will restore it: if not, I will devour it'. "You will devour it!" cried the distracted Mother. "Now", said the wily Crocodile, "I *cannot* restore your Baby: for, if I do, I shall make ~~you speak falsely, and I warned you that, if you spoke falsely, I would devour it~~". "On the contrary", said the yet wiler Mother, "you *cannot devour* my Baby: for, if you do, you will make me speak *truly*, and you promised me that, if I spoke *truly*, you would *restore* it!" (We assume, of course, that he was a Crocodile of his word; and that his sense of honour outweighed his love of Babies.)

§ 4.

Antistrephon.

That is "*The Retort*". This is a tale of the law-courts:

"Protagoras had agreed to train Euathius for the profession of a barrister, on the condition that half his fee should be paid at once, and that the other half should be paid, or not paid, according as Euathius should win, or lose, his first case in Court. After a time, Protagoras, becoming impatient, brought an action against his pupil, to recover the second half of his fee. It seems that Euathius decided to plead his own cause. "Now, if I win this action", said Protagoras, "you will have to pay the money by the decision of the Court: if I *lose* it, you will have to pay by our agreement. Therefore, in any case, you must pay it". "On the contrary", retorted Euathius, "if you *win* this action, I shall be released from payment by our agreement: if you *lose* it, I shall be released by the decision of the Court. Therefore, in any case, I need not pay the money".

Si consideri il ragionamento che segue, che nel quadro della logica aristotelica viene trattato agevolmente con un sillogismo:

Tutti gli uomini sono mortali
Tutti i greci sono uomini
∴ Tutti i greci sono mortali.

Il simbolo ∴, naturalmente, sta per «perciò». Questa è una inferenza valida, e la regola è:

Tutti gli *M* sono *X*
Tutti i *G* sono *M*
∴ Tutti i *G* sono *X*.

Qualunque ragionamento di questa forma, indipendentemente da ciò che viene sostituito a *M*, *X* e *G*, sarà valido.

I logici aristotelici del Medioevo classificavano le forme di inferenza valide in base alla loro «figura» e al loro «modo». Le variazioni nelle posizioni dei termini di un sillogismo (nell'esempio *M*, *X* e *G*) sono denominate differenze di figura. Ciascuno dei sillogismi possiede anche un modo che è determinato dalla forma delle sue proposizioni componenti. Ci possono essere 15, 19, 24 o più forme di inferenza valide, a seconda del tipo di classificazione impiegato.

Il problema sta nel fatto che ci sono molti ragionamenti validi le cui regole di inferenza non si possono neppure formulare nel quadro della logica aristotelica tradizionale. Per esempio:

Rebecca è la madre di Giacomo
Giacomo è il padre di Giuseppe
La madre del padre è la nonna paterna
∴ Rebecca è la nonna paterna di Giuseppe

Il ragionamento può essere così formulato nel linguaggio della logica aristotelica:

Tutti gli *A* sono *B*
Tutti i *C* sono *D*
Tutti gli *E* sono *F*
∴ Tutti gli *A* sono *G*.

Ma una volta formulato in questo modo questo ragionamento valido è assolutamente impossibile formulare una regola di inferenza valida per esso. Una espressione come «madre di Giacomo» è contratta in un singolo termine (*B*) e non può essere di nuovo analizzata. Si possono trovare facilmente altre espressioni che, una volta sostituite ai termini da *A* a *G*, produrranno un controesempio, tale da portare da una premessa vera a una conclusione falsa. La struttura logica del linguaggio delle proposizioni categoriche aristoteliche è troppo debole per rendere trasparente il modo in cui il predicato

Le bozze di stampa scoperte dall'autore alla Christ Church di Oxford nel 1959 furono il primo indizio di una possibile esistenza di più ampie porzioni del secondo libro, mancante, dedicato da Dodgson alla logica simbolica. Il numero segnato a penna sull'angolo sinistro in alto è di pugno di Dodgson, che teneva un registro di tutta la sua corrispondenza. Una piega nelle bozze copre parzialmente una riga del testo stampato.

«madre di Giacomo» contiene il soggetto della seconda premessa e parte del soggetto della terza.

Con l'attuale logica delle relazioni dare la regola di inferenza valida per questo esempio è assolutamente banale. Supponiamo che x , y e z stiano per Rebecca, Giacomo e Giuseppe, e che M , F e T stiano per le relazioni tra individui «madre di», «padre di» e «nonna materna di». Allora

« Mxy »
« Fyz »
« $MF = T$ »
« Txz »

La regola di inferenza afferma che qualunque conclusione della forma logica « Txz » è deducibile incondizionatamente dalle forme che compaiono sopra di essa.

Uno degli obiettivi principali che presiedono alla costruzione dei calcoli della logica contemporanea è quello di ridurre sistematicamente le regole di inferenza al numero più piccolo possibile. Di certo Dodgson si poneva lo stesso problema. Nella *Logica Simbolica: Parte I* egli scriveva: «Per quanto riguarda i sillogismi, trovo che le loro diciannove forme, attorniate da tutto un insieme di altre che [i manuali] hanno ignorato, si possono disporre tutte sotto tre forme, ciascuna con una semplicissima regola correlata». Egli considerava la logica aristotelica come «una macchina pressoché inutile a scopi pratici, data l'incompletezza di molte delle conclusioni, e l'omissione di molte forme del tutto legittime».

È evidente che Dodgson stava cercando di spingersi oltre le forme tradizionali di argomentazione valida. Quali sono le acquisizioni specifiche che gli si possono ora attribuire alla luce del suo lavoro inedito di logica simbolica? Per prima cosa va detto che nel corso del 1896 egli aveva sviluppato una procedura meccanica di controllo della validità per buona parte della logica dei termini, un risultato abitualmente attribuito a Leopold Löwenheim.

In secondo luogo, già dal 1894 Dodgson usava tavole di verità per la soluzione di problemi logici particolari. La applicazione di tavole di verità e matrici non divenne di uso comune prima del 1920. In terzo luogo, durante il 1896 Dodgson aveva sviluppato il «metodo degli alberi» per determinare la validità di argomentazioni che erano notevolmente complicate in rapporto alle normali capacità dei logici inglesi del suo tempo. L'idea chiave era quella di controllare se una conclusione

SIMBOLO

SIGNIFICATO

x_1	L'indice sottoscritto 1 asserisce l'esistenza di x : «Alcune cose esistenti hanno l'attributo x » o più brevemente «Qualche x esiste».
x'	L'apice nega un termine o un enunciato. Se x significa «nuovo», x' significa «non-nuovo». Quindi x' va letto «non- x ».
x_0	L'indice sottoscritto 0 asserisce la non-esistenza di x : «Nessuna cosa esistente ha l'attributo x » o, più brevemente, «Nessun x esiste».
xy_1	«Qualche xy esiste». Quando in un'espressione ci sono due lettere, non importa quale sta per prima; il significato è identico. L'espressione significa anche «Alcuni x sono y » e «Alcuni y sono x ».
xy_0	«Non esiste nessun xy », che è equivalente a «Nessun x è y » e a «Nessun y è x ».
$x_1y'_0$	«Tutti gli x sono y ». L'indice sottoscritto 0 asserisce la non-esistenza della combinazione x e y' . L'espressione ci dice che alcune cose- x esistono ma che nessuna di esse ha l'attributo y' e, nel sistema di Dodgson, questo equivale alla proposizione «Tutti gli x sono y ».
\dagger	La <i>croce</i> significa «e». Quindi $xy_1 \dagger xy'_0$ significa «Esiste qualche xy e non esiste nessun xy' » oppure «Alcuni x sono y e nessun x è non- y ».
\P	Significa «se vero, dimostrerebbe» oppure «è derivabile da». Quindi $x \P y$ significa « x dimostra y » oppure « y è derivabile da x ».

Il metodo degli indici sottoscritti fu escogitato da Dodgson per rappresentare le proposizioni in forma stenografica. Ciascuno degli indici comincia ad avere effetto dall'inizio dell'espressione a cui è posposto, ma l'apice nega solamente il termine a cui viene accostato. Dodgson assumeva le proposizioni inizianti con «Tutti» («Tutti gli x sono y ») come equivalenti alle due proposizioni «Ci sono alcuni x » e «Nessun x è non- y ». Nel convertire le proposizioni di tipo «Tutti» in forma sottoscritta, il predicato doveva essere negato. Così «Tutti gli x sono y » è scritto « $x_1y'_0$ » e letto o come «Nessun x è non- y » e «Alcuni x sono y » o come il loro equivalente «Tutti gli x sono y ». Analogamente, «Tutti gli y sono non- x » diventa « y_1x_0 » («Nessun y è x » e «Alcuni y sono non- x ») e l'espressione « $x'_1y'_0$ » è letta «Tutti i non- x sono y ». Nel tradurre un enunciato «Tutti» da una forma all'altra il predicato (ultimo termine) passa sempre da positivo a negativo (da x a x') o da negativo a positivo.

Si considerino le premesse: $xm_0 \dagger ym'_0$

Poiché m' nega m , si possono eliminare ambedue i termini:

$$xm_0 \dagger ym'_0$$

Ciò che resta si può considerare una singola espressione: xy_0

$$\therefore xm_0 \dagger ym'_0 \P xy_0$$

Per premesse multiple, in quella che Carroll chiama prima figura, il procedimento si ripete finché non si può eliminare altro. Ciò che rimane a questo punto sono i termini che compariranno nella conclusione. Per esempio, se le premesse di partenza sono:

Dati 1. 2. 3. 4.
 $a'_1c'_0$ $a_1e'_0$ $c_1b'_0$ d_1b_0
Si combinino i dati 1 e 2, e si elimini $a'a$
Si combinino i dati 1 e 3, e si elimini $c'c$
Si combinino i dati 3 e 4, e si elimini $b'b$.

Ciò che rimane è $d_1e'_0$. Questa è la conclusione.

$$\therefore a'_1c'_0 \dagger a_1e'_0 \dagger c_1b'_0 \dagger d_1b_0 \P d_1e'_0$$

La derivazione di conclusioni da premesse formulate nel linguaggio con indici sottoscritti di Dodgson comporta spesso l'eliminazione di termini che si negano reciprocamente (m e m' , per esempio, giacché m' significa non- m). Dodgson preferiva cancellare i termini da eliminare anziché valersi del metodo qui descritto, ma i risultati sono identici. La soluzione per il primo problema si può leggere: «Il fatto che nessun x è m e che nessun y è non- m , se fosse vero, dimostrerebbe che nessun x è y ». Nel secondo problema si possono combinare le premesse finché non sono usate tutte.

ipoteticamente assunta come falsa, una volta congiunta con una serie di premesse assunte come vere, avrebbe portato a una contraddizione o a una assurdità. La sua procedura presenta una sorprendente affinità con gli alberi frequentemente impiegati dai logici d'oggi, in gran parte derivanti dal metodo delle « Tavole Semantiche » escogitato nel 1955 da Evert Willem Beth.

Anche se Dodgson non si attenne mai ai canoni di rigore odierni, la sua anticipazione dei recenti sviluppi della logica è sufficiente a attestare la sua

originalità. E tuttavia, malgrado il carattere antiaristotelico di tanta parte del suo lavoro, Dodgson rimase ostinatamente aristotelico su un punto: la « portata esistenziale » delle proposizioni universali. Egli sosteneva la tesi che un enunciato del tipo « Tutti », per esempio « Tutti gli uomini sono mortali », era equivalente ai due enunciati « Non ci sono uomini non mortali » e « Alcuni uomini sono mortali ». Poiché ogni enunciato del tipo « Tutti » contiene un enunciato del tipo « Qualche », tutti gli enunciati del tipo « Tut-

ti » asseriscono l'esistenza reale dei loro soggetti.

Verso la metà del XIX secolo i logici, e in particolare George Boole, avevano cominciato a negare che gli enunciati del tipo « Tutti » asserissero necessariamente l'esistenza dei loro soggetti. L'interpretazione booleana oggi è quasi universalmente accettata dai logici matematici (anche se nel 1964 ha trovato un oppositore in Richard B. Angell della Wayne State University). Così dal punto di vista della maggior parte dei logici contemporanei la credenza di Dodgson nella « portata esistenziale » delle proposizioni universali compromette seriamente il suo contributo alla logica. (Cionondimeno, le tecniche di decisione e il formalismo di Dodgson si possono interpretare in modo tale da ottenere risultati booleani anziché aristotelici.)

Per quanto le innovazioni tecniche di Dodgson siano storicamente interessanti, i passi più affascinanti dei suoi scritti inediti sono quelli dedicati ai paradossi e ai rompicapi. Uno di essi, *Un paradosso logico*, fu pubblicato su « Mind » nel 1894 ed è ancora oggetto di una accesa polemica tra i logici contemporanei. Il problema chiama in causa una bottega di barbiere con tre barbieri che possono lasciare il negozio solo sotto determinate condizioni, che vengono formulate come premesse (si veda la figura nella pagina a fronte). Tuttavia un ragionamento valido porta da queste due premesse a conclusioni contraddittorie. Dodgson parlò del paradosso del barbiere come di una « presentazione ornamentale » di una disputa tra lui e Wilson iniziata nel 1893. Questa si trascinò per più di dieci anni con un nutrito scambio di corrispondenza, gran parte della quale si conserva ancora e sarà forse pubblicata, e con una serie di manoscritti, alcuni dei quali furono pubblicati privatamente da Dodgson.

Un altro paradosso tratto dagli scritti inediti riguarda l'antico problema del coccodrillo e del bambino. Scriveva Dodgson:

« La tragica storia si snoda così:

Un coccodrillo aveva rapito un bambino sulle rive del Nilo. La madre lo supplicò di restituirle il caro piccino. « Orbene — disse il coccodrillo — se tu indovini ciò che io farò, te lo restituirò: altrimenti lo divorerò ».

« Tu lo divorerai! » gridò la madre fuori di sé.

« A questo punto — disse lo scaltro coccodrillo — io non posso restituirti il bambino: infatti se lo facessi ti farei dire il falso; e ti ho avvertito che, se

Come esempio il più possibile semplice di questo Metodo prendiamo il Sillogismo della figura in basso a pagina 37, ossia $xm_0 \vdash ym'_0 \vdash xy_0$.

Qui i nostri dati sono due Nullità, xm_0 e ym'_0 , che presentano l'Attributo m nella forma sia *positiva* che *negativa*; il nostro *Quaesitum* è la Nullità xy_0 .

Cominciamo coll'assumere che l'aggregato xy sia un'Entità: assumiamo cioè che Qualcosa di esistente abbia ambedue gli Attributi x e y .

Ora la prima Premessa ci dice che x è incompatibile con m . Quindi la « Cosa » sotto considerazione che si assume sia in possesso dell'Attributo x non può avere l'Attributo m . Ma deve necessariamente avere o m o m' , in quanto questi costituiscono una *Divisione esaustiva* dell'intero Universo. Quindi *deve* avere l'attributo m' .

Analogamente, dalla seconda Premessa, possiamo dimostrare, come nostro secondo risultato, che la « Cosa » sotto considerazione ha l'Attributo m .

Questi due risultati, presi insieme, ci danno la sorprendente asserzione che questa « Cosa » ha ambedue gli attributi, m e m' , *simultaneamente*; cioè abbiamo $xy_1 \vdash xym'm_1$.

Ora noi sappiamo che m e m' sono *contraddittori*: quindi il risultato è evidentemente *assurdo*: per cui torniamo alla nostra assunzione originale (che un aggregato xy fosse un'Entità) e diciamo « quindi xy non può essere un'Entità: in altre parole, è una Nullità ».

Disponiamo ora questo ragionamento sotto forma di *Albero*.

Per cominciare devo spiegare che tutti gli Alberi in questo sistema crescono a testa in giù: la radice è in cima, e i rami sono sotto. Se qualcuno obiettasse che il nome « Albero » non è appropriato, la mia risposta è che mi limito a seguire l'esempio di tutti gli autori che trattano di genealogia. Un « albero » genealogico cresce *sempre verso il basso*: perché un albero logico non potrebbe fare lo stesso?

Dunque metterò la radice del mio Albero in alto. Essa consiste dell'aggregato xy ; e il puro fatto di scrivere queste due lettere va inteso con questo significato (usando la forma regolare di una *reductio ad absurdum*): « l'aggregato xy sarà una Nullità; perché, altrimenti, poniamo sia un'Entità; sia cioè una data cosa esistente in possesso dei due attributi x e y ».

Sotto questo « xy » allora pongo la lettera m' (questa fa parte del *fusto* del nostro Albero), e sul suo lato sinistro pongo il numero « 1 », seguito da un punto fermo, cosicché il nostro Albero è ora

xy
1. m'

Il significato di questo è che la « Cosa », che si assume sia in possesso dei due attributi x e y , *deve anche* avere l'attributo m' : e il numero « 1 » vi rimanda alla prima premessa come garanzia per questa asserzione.

Poi io pongo la lettera m sul lato destro di m' , e il numero « 2 », seguito da una virgola, sul lato sinistro dell'« 1 », cosicché il nostro Albero ora è

xy
2, 1. $m'm$

Ciò significa che la « Cosa » *deve* avere anche l'attributo m , (cioè che $xym'm$ è un'Entità) e che la garanzia per asserire questo è la seconda premessa. (Si osservi che le due lettere, nella linea inferiore, vanno lette da sinistra a destra, mentre i due numeri di riferimento da destra verso sinistra.)

Ora noi sappiamo che m e m' sono *contraddittori*: quindi è impossibile che un aggregato che li contiene entrambi sia un'Entità: quindi è una Nullità. Questo fatto lo indico tracciando un piccolo cerchio (rappresentante un vuoto) sotto di esso, cosicché il nostro Albero ora è

xy
2, 1. $m'm$
○
∴ xy_0

Il significato del cerchio è « L'aggregato degli Attributi, dalla radice a questo punto, è una Nullità ».

Poi pongo sotto il circoletto la conclusione: « ∴ xy_0 », cosicché

xy
2, 1. $m'm$
○
∴ xy_0

l'Albero ora è

Il significato dell'ultima riga è « Abbiamo ora dimostrato, dall'assunzione che xy fosse un'Entità, che questo aggregato, $xym'm$, deve essere una Entità. Ma esso è evidentemente una Nullità. Il che è *assurdo*. Quindi la nostra assunzione era falsa ». Quindi siamo in diritto di dire « Perciò xy è una Nullità ».

Il metodo degli alberi venne sviluppato da Dodgson come strumento per controllare la validità di una conclusione derivata da determinate premesse. Esso è sorprendentemente simile agli « alberi » frequentemente impiegati dai logici contemporanei. L'idea base è quella di assumere che le premesse siano vere ma che la conclusione sia falsa (e la sua negazione sia vera). Se combinando la negazione della conclusione con le premesse si giunge a un'assurdità, ciò prova che le premesse dimostravano veramente la conclusione. La spiegazione che viene data è stata tratta dalla sezione descrivente « il metodo degli alberi » nel lavoro inedito di Dodgson. Essa si riferisce a un semplice albero.

tu avessi detto il falso, io lo avrei divorato».

«Al contrario — disse la madre, ancora più scaltramente — tu non puoi divorare il mio bambino: infatti se lo facessi mi faresti dire il vero, e tu mi hai promesso che se avessi detto il vero me lo avresti restituito!» (Stiamo naturalmente ipotizzando che fosse un coccodrillo di parola e che il suo senso dell'onore fosse più forte del suo amore per i bambini) ».

Dodgson quindi applica al problema la sua logica con indici sottoscritti e il suo « metodo degli alberi ».

« Su questo sofisma [R. H.] Lotze fa la sconsolante osservazione che « Non c'è via d'uscita al dilemma ». Io penso, tuttavia, che l'apparato della logica simbolica risulterà sufficiente alla sua soluzione.

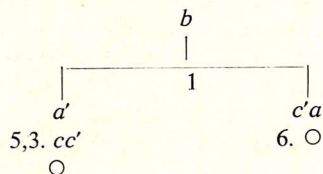
Sia Univ. l'universo; a = la madre dice la verità; b = il coccodrillo mantiene la parola; c = il coccodrillo divora il bambino ».

(Per negare un termine o un enunciato, Dodgson poneva un apice immediatamente dopo di esso. Così a' sta per « La madre dice il falso ». L'indice sottoscritto 1 asserisce l'esistenza del termine mentre l'indice 0 asserisce la sua non-esistenza. Così ab_1c_0 stava a significare « Tutti gli ab sono non- c », che in questo particolare esempio vuol dire: « Ogniqualvolta la madre dice la verità e il coccodrillo mantiene la sua parola, allora il coccodrillo non divorerà [restituisce] il bambino ». E c'_1a_0 significa « Tutti i non- c sono non- a », ovvero « Ogniqualvolta il coccodrillo restituisce il bambino, la madre dice il falso ».)

« Allora abbiamo, come dati del problema:

- | | |
|----------------|----------------|
| 1. ab_1c_0 | 4. $a'b'_1c_0$ |
| 2. $ab'_1c'_0$ | 5. $c_1a'_0$ |
| 3. $a'b_1c'_0$ | 6. c'_1a_0 |

Possiamo ignorare 2 e 4 in quanto sono contenuti in 6 e 5; e vediamo che, all'analisi, va ritenuto soltanto b .



∴ b_0 : cioè « qualunque cosa faccia il coccodrillo, viene meno alla sua parola ».

Il paradosso del barbiere fu pubblicato sulla rivista filosofica « Mind » nel luglio 1894, ma Dodgson chiaramente lo riservava al secondo volume della *Logica Simbolica*, insieme alla sua soluzione. Questa versione proviene dalle bozze inedite.

« Come, non hai niente da fare? — disse zio Jim. — Allora vieni con me da Allen. E puoi farti un giro mentre mi faccio radere ».

« Benissimo — disse zio Joe — e anche il cucciolo farebbe bene a venire, no? »

Il cucciolo sono io, come il lettore avrà forse capito da solo. Ho compiuto i quindici anni più di tre mesi fa; ma non c'è verso di far capire questo a zio Joe, che si limiterebbe a dire: « Torna a cuccia, cucciolo! » oppure « Suppongo allora tu riesca a risolvere equazioni cubiche! » o qualche altra battuta egualmente ignobile. Ieri mi chiese di dargli un esempio di proposizione in A , e io risposi « Tutti gli zii fanno ignobili giochi di parole ». Non penso gli sia piaciuto. Comunque, questo non ha nessuna particolare importanza. Ero abbastanza contento di andare con loro. Adoro sentire questi miei zii « discutere in termini logici », come usano dire; e ce la mettono tutta, ve l'assicuro!

« Questa non è un'inferenza logica dalla mia affermazione », diceva zio Jim.

« Mai detto una cosa del genere. È una *reductio ad absurdum* », diceva zio Joe.

« Un *illicito procedere del minore*! » buttò là zio Jim.

Questo è il modo in cui discutono sempre, quando io sono con loro. Come se ci fosse chissà quale divertimento nel chiamarmi un minore!

Dopo un po' zio Jim ricominciò, proprio quando fummo in vista del barbiere. « Spero solo che Carr sia in bottega — disse. — Brown è così maldestro. E la mano di Allen è diventata malferma dopo che ha avuto quel febbre ».

« È cosa certa che Carr è in bottega », disse zio Joe.

« Scommetto mezzo scellino che non c'è! », dissi io.

« Risparmia le tue scommesse per occasioni migliori — disse zio Joe. — Voglio dire — si affrettò ad aggiungere, vedendo dal mio sogghigno che aveva fatto un passo falso — voglio dire che posso dimostrarlo logicamente. Non dipende dal caso ».

« Dimostralo logicamente! — disse zio Jim, beffardo. — Sbrigati, ti sfido a farlo! »

« Per amor di discussione — cominciò zio Joe — supponiamo che Carr sia fuori. E vediamo a che ci porterebbe questa assunzione. Farò questo per *reductio ad absurdum* ».

« Naturalmente! — borbottò zio Jim — non ho mai visto un tuo ragionamento che non finisse in qualche assurdità! »

« Senza lasciarmi provocare dai tuoi vili insulti — disse zio Joe con accenti elevati — vado avanti. Se Carr è fuori, mi concedi che, se Allen è anche fuori, Brown deve essere in bottega? »

« Che vantaggio c'è nel fatto che ci sia lui in bottega? — disse zio Jim — Non voglio che Brown mi rada! È troppo maldestro ».

« La pazienza è una di quelle inestimabili qualità... » cominciò a dire zio Joe; ma zio Jim tagliò corto.

« Ragione! — disse — non moraleggiare! »

« Bene, ma tu ammetti questo? — insistette zio Joe — Mi concedi che, se Carr è fuori, ne segue che se Allen è fuori Brown deve essere dentro? »

« Certo che deve — disse zio Jim — altrimenti non ci sarebbe nessuno a occuparsi del negozio ».

« Vediamo, dunque, che l'assenza di Carr chiama in causa una certa ipotetica, la cui *protasi* è « Allen è fuori » e la cui *apodosi* è « Brown è dentro ». E vediamo che, finché Carr resta fuori, l'ipotetica resta valida? »

« Bene, supponiamo di sì. E allora? »

« Tu mi concederai anche che la verità di un'ipotetica, voglio dire la sua *validità* come *sequenza* logica, non dipende nemmeno minimamente dal fatto che la *protasi* sia effettivamente vera, e nemmeno dal fatto che sia possibile. L'ipotetica « Se tu dovessi correre da qui a Londra in cinque minuti sorprenderesti la gente » resta vera come *sequenza*, che tu lo possa fare o meno ».

« Io non posso farlo », disse zio Jim.

« Dobbiamo ora considerare un'altra ipotetica. Che mi hai detto ieri relativamente ad Allen? »

« Ti ho detto — disse zio Jim — che da quando ha avuto quella febbre l'uscire da solo lo rende nervoso, per cui si fa sempre accompagnare da Brown ».

« Proprio così — disse zio Joe — Quindi, l'ipotetica « Se Allen è fuori Brown è fuori » è sempre valida, no? »

« Suppongo di sì », disse zio Jim (Era lui a sembrare un po' nervoso, ora).

« Quindi, se Carr è fuori, abbiamo due ipotetiche, « Se Allen è fuori Brown è dentro » e « Se Allen è fuori Brown è fuori » simultaneamente valide. E due ipotetiche incompatibili, bada bene! Non è possibile che siano vere insieme! »

« Non è possibile? » disse zio Jim.

« E come potrebbero? — disse zio Joe — Come potrebbe una stessa *protasi* dimostrare due *apodosi* contraddittorie? Tu mi concedi che le due apodosi, « Brown è dentro » e « Brown è fuori » sono contraddittorie? »

« Questo lo concedo »

« Allora posso tirare le somme. — disse zio Joe — Se Carr è fuori, queste due ipotetiche sono vere insieme. E noi sappiamo che non possono essere vere insieme. Il che è assurdo. Perciò Carr non può essere fuori. Eccoti una bella *reductio ad absurdum*! »

Zio Jim sembrava profondamente perplesso; ma dopo un po' riprese coraggio e ricominciò. « Non vedo del tutto chiaro circa questa incompatibilità. Perché non dovrebbero essere vere insieme queste due ipotetiche? Mi sembra che questo dimostrerebbe semplicemente « Allen è dentro ». Naturalmente è chiaro che le apodosi di queste due ipotetiche, « Brown è dentro » e « Brown è fuori », sono incompatibili. Ma perché non dovremmo porre le cose così? Se Allen è fuori Brown è fuori. Se Carr e Allen sono entrambi fuori, Brown è dentro. Il che è assurdo. Perciò Carr e Allen non possono essere entrambi fuori. Ma fintantoché Allen è dentro non vedo cosa impedisca a Carr di andare fuori ».

« Mio caro, ma assai illogico fratello! — disse zio Joe — (Ogniqualvolta zio Joe comincia a darti del « caro », puoi star sicuro che ti ha messo in un bel guaio!) Non vedi che stai erroneamente scindendo la *protasi* dall'*apodosi* dell'ipotetica? La *protasi* è semplicemente « Carr è fuori »; e l'*apodosi* è un tipo di sotto-ipotetica, « se Allen è fuori, Brown è dentro ». Ed è semplicemente l'assunzione « Carr è fuori » che ha provocato questa assurdità. Per cui c'è una sola conclusione possibile. Carr è dentro! »

Non ho la minima idea di quanto avrebbe potuto durare questa discussione. Sono convinto che sia l'uno che l'altro avrebbero potuto discutere per sei ore di seguito. Ma proprio in questo momento arrivammo alla bottega del barbiere; ed entrando trovammo

Così, se divora il bambino, le fa dire la verità, e quindi *viene meno* alla sua parola; se lo restituisce, le fa dire il falso, e quindi *viene meno* alla sua parola. Essendo così frustrato senza speranza il suo senso dell'onore, non possiamo dubitare che egli si comporterebbe secondo i dettami della sua *seconda* passione dominante, il suo amore per i bambini!».

Osservando che i dati 2 e 4 sono contenuti nei dati 6 e 5 Dodgson intendeva dire che i primi sono logicamente derivabili dai secondi. La conclusione che si raggiunge applicando la regola di inferenza di Dodgson è «Nessun *b* esiste» oppure «Non ci sono casi in cui il cocodrillo mantiene la parola» oppure, per usare le parole di Dodgson, «Qualunque cosa faccia viene meno alla sua parola».

L'albero logico è un'applicazione di un ragionamento per *reductio ad absurdum* all'assunzione ipotetica che b_0 (non esiste nessun *b*) sia falso. Postuliamo b_1 (*b* esiste). Questa informazione, congiuntamente al primo dato (o premessa), consente due possibilità. La prima, sul ramo sinistro dell'albero, è *a'* (la madre dice il falso). Questo risultato, tuttavia, congiuntamente a *b*, porta per la terza premessa a *c*, mentre per la quinta premessa porta a *c'*. Poiché *c* e *c'* sono contraddittorie, la prima possibilità porta a una assurdità (indicata dal cerchio). La seconda possibilità, rappresentata nel ramo destro dell'albero, di congiungere b_1 con la prima premessa produce *c'a* (il cocodrillo restituisce il bambino e la madre dice la verità). Questo è in contrasto con la sesta premessa e porta pure a

una assurdità. Se l'assunzione che il cocodrillo può talvolta mantenere la parola porta a un'assurdità, allora è vero che «Qualunque cosa faccia, il cocodrillo viene meno alla parola data».

Per verificare la nostra padronanza della tecnica di Dodgson, proviamoci a determinare ciò che succede se la madre dice: «Tu restituirai il bambino». Il dato 5 diventa c_1a_0 (tutti i *c* sono non-*a*) e il dato 6 diventa $c'_1a'_0$ (tutti i non-*c* sono *a*). Cioè «Se il cocodrillo restituisce il bambino, la madre dice il vero». Qui i dati 1 e 3 si possono trascurare in quanto sono derivabili dai dati 5 e 6. I dati rilevanti sono allora 2, 4, 5 e 6. La conclusione è b'_0 , «Qualunque cosa faccia, il cocodrillo mantiene la parola». Per sottoporre a controllo il ragionamento con il

«Un rompicapo logico» di Lewis Carroll

Ci sono tre proposizioni, *A*, *B* e *C*.
Dato che

- «Se *A* è vera, *B* è vera; (i)
- «Se *C* è vera, allora se *A* è vera *B* non è vera» (ii)

NEMO e OUTIS divergono sulla verità di *C*.

NEMO dice che *C* non può essere vera; OUTIS sostiene il contrario.

Ragionamento di NEMO

Il numero (ii) si riduce a questo:

«Se *C* è vera, allora (i) non è vera».

Ma, *ex hypothesi*, (i) è vera.

∴ *C* non può essere vera; infatti l'assunzione di *C* comporta un'assurdità.

Risposta di OUTIS

Le due asserzioni di NEMO, «se *C* è vera, allora (i) non è vera» e «l'assunzione di *C* comporta un'assurdità» sono erranee.

L'assunzione di *C da sola* non comporta nessuna assurdità, in quanto le due ipotetiche «se *A* è vera *B* è vera» e «se *A* è vera *B* non è vera» sono *compatibili*; cioè possono essere vere assieme, nel qual caso *A* non può essere vera.

Ma l'assunzione di *C* e di *A assieme* comporta effettivamente un'assurdità, in quanto le due proposizioni «*B* è vera» e «*B* non è vera» sono *incompatibili*.

Quindi segue non che *C, presa da sola*, non può essere vera, ma che *C* e *A* non possono essere vere *insieme*.

Difesa di NEMO

OUTIS ha erroneamente separato protasi e apodosi in (ii).

L'assurdità non è l'ultima clausola di (ii), «*B* non è vera», ma con *tutto ciò* che segue la parola «allora», ossia l'ipotesi «se *A* è vera *B* non è vera» e, per (ii), è soltanto l'assunzione di *C* che genera questa assurdità.

In effetti OUTIS ha reso (ii) equivalente a «Se *C* è vera [e *A* è vera] allora se *A* è vera *B* non è vera». Questo è un errore: le parole tra parentesi nella protasi composta sono superflue, e ciò che rimane è la vera protasi che condiziona l'apodosi assurda, come è evidente dalla forma di (ii) data originariamente.

Questo teorema sulle ipotetiche — che le proposizioni numero (i) e (ii) insieme dimostrano che *C* non può essere vera — possono essere illustrate dal seguente esempio algebrico:

$$\text{Sia } ax + (a - b)y + z = 5; \dots\dots\dots (1)$$

$$bx + z = 6; \dots\dots\dots (2)$$

L'equazione (1) si può formulare come un'ipotesi in questo modo:

«Se ax , $(a - b)y$ e z sono aggiunte l'una all'altra, si ottiene il numero "5"».

«*A*» stia a significare « ax , $(a - b)y$ e z sono aggiunte l'una all'altra»;

«*B*» stia a significare «si ottiene il numero "5"»;

«*C*» stia a significare « $a = b$ ».

Allora abbiamo

«Se *A* è vera, *B* è vera».

Si assuma che *C* sia vera; cioè che $a = b$.

Allora $(ax + (a - b)y + z)$ diventa $(bx + z)$ che, per l'equazione (2), deve *sempre* essere = 6.

Quindi

«Se *C* è vera, allora se *A* è vera *B* non è vera».

Perciò *C* non può essere vera;

cioè «*a*» non può essere «*b*».

Seconda risposta di OUTIS

Questa risposta comprenderà (α) una dimostrazione che il «ragionamento di NEMO» è autodistruttivo; (β) una dimostrazione che il suo esempio algebrico è inadeguato in quanto non rappresenta correttamente i dati; (γ) una dimostrazione che, una volta corretto, esso conforta la tesi di OUTIS, ossia che le ipotesi (i) e (ii) dimostrano non che *C, presa da sola*, non può essere vera, ma che *C* e *A* non possono essere vere *insieme*; (δ) una semplice dimostrazione del *vero* risultato di queste due ipotesi.

(α)

Consideriamo la terna di ipotetiche (che chiameremo (K), (L), e (M))

(K) «Se *X* è vera, *Y* non è vera».

(L) «Se *X* è vera, *Y* è vera».

(M) «Se *X* non è vera, *Y* è vera».

Non si discuterà il fatto che (L) e (M), prese insieme, sono equivalenti alla categorica (che chiameremo «N») «*Y* è vera». Quindi la terna di cui sopra è equivalente all'ipotesi e alla categorica

(K) «Se *X* è vera, *Y* non è vera»

(N) «*Y* è vera».

Per questa terna (o per la sua coppia equivalente) si possono proporre due diverse interpretazioni, ossia

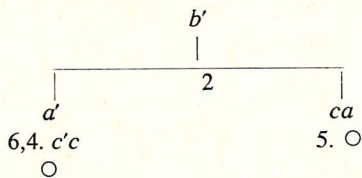
«(K) e (L) non possono essere vere insieme. Quindi (K), (L) e (M) non possono essere vere insieme».

«(K) e (N) possono essere vere insieme; cioè, (K), (L) e (M) possono essere vere insieme».

La risoluzione del paradosso del barbiere trovata tra le carte di Dodgson figura come la sua trattazione definitiva del problema. Egli riformula il paradosso in termini più astratti e procede, con ragio-

namenti pro e contro, a discutere se *C* possa essere vera o meno («Carr è fuori» nella versione originale). Nemo (zio Joe) sostiene il punto di vista di John Cook Wilson

metodo dell'albero logico di Dodgson, assumiamo che b'_0 sia falso e proviamo a vedere se questo conduce o no a una assurdità. Se b'_0 è falso, allora b'_1 è vero, e questo, insieme alla seconda premessa (dato), lascia aperte due possibilità.



La prima (ramo sinistro dell'albero) è a' (la madre dice il falso), ma $a'b'_1$ per la quarta premessa porta a c' e per la sesta a c . Avere c' e c insieme è assurdo. La seconda possibilità è ca_1 (il ramo destro), che contraddice la quinta premessa ed è quindi assurda. Così

b' deve essere vero e il coccodrillo, qualunque cosa faccia, mantiene la parola.

Un altro problema preso in considerazione da Dodgson è il celebre «paradosso del mentitore». Egli riferisce la sua «forma più semplice» in questo modo:

«Se un uomo dice "Sto dicendo una bugia", e dice la verità, egli sta dicendo una bugia, e perciò dice il falso: ma se egli dice il falso non sta dicendo una bugia e perciò dice la verità».

Molti logici hanno in anni recenti accantonato il paradosso del mentitore perché il paradosso sorge dal fatto di consentire l'autoriferimento, ossia dal fatto di permettere agli enunciati di riferirsi alla loro stessa verità o fal-

sità. Tale atteggiamento nasce probabilmente da una lettura troppo frettolosa dell'articolo di Tarski *Il concetto di verità nei linguaggi formalizzati*, apparso nel 1935, nel quale l'autore sostiene che nessun linguaggio consistente può contenere i mezzi per parlare del significato o della verità delle sue stesse espressioni. Quando un linguaggio consente invece l'autoriferimento, non sorprende che esso conduca a inconsistenze e paradossi.

Lo stesso Dodgson prende in considerazione un suggerimento del genere, lo valuta seriamente e quindi lo respinge, tutto nel volgere di poche righe. Egli scrive:

«Il modo migliore per uscire dalla difficoltà [il paradosso del mentitore] sembra sia quello di sollevare la questione se la proposizione "Sto dicendo

Queste interpretazioni sono incompatibili.

Ora, quando NEMO dice «l'assunzione di C comporta un'assurdità», l'assurdità a cui egli allude è la verità simultanea delle due proposizioni «Se A è vera B è vera» e «Se A è vera B non è vera».

Queste due proposizioni sono ipotetiche di forma (L) e (K) ; nel dichiarare che l'assunzione della loro verità simultanea comporta un'assurdità, NEMO virtualmente dichiara che esse non possono essere vere assieme.

Qui, allora, egli adotta la prima interpretazione della terna di ipotetiche, (K) , (L) e (M) .

Di nuovo, quando egli dice « C non può essere vera» le premesse da cui egli deduce la conclusione sono le due proposizioni «Se C è vera, allora (i) non è vera. Ma, ex hypothesi, (i) è vera».

Queste due proposizioni sono un'ipotetica e una categorica di forma (K) e (N) ; nel dedurre da esse, considerate come premesse, una conclusione, NEMO virtualmente dichiara che esse non possono essere vere insieme.

Qui, allora, egli adotta la seconda interpretazione della terna di ipotetiche, (K) , (L) e (M) .

Così egli ha adottato, nel corso dello stesso ragionamento, due interpretazioni incompatibili di questa terna.

Quindi il ragionamento di NEMO è autodistruttivo.

(β)

Esaminiamo ora l'esempio algebrico di NEMO.

Egli ci dà le equazioni (1) e (2) come sempre vere.

Quindi l'equazione (1) resta vera anche quando $a = b$.

Quindi la sua seconda ipotetica è incompleta: essa dovrebbe essere «Se C è vera, allora se A è vera B è (per l'equazione 1) vera ma (per l'equazione 2) non vera».

Quindi il suo esempio algebrico è inadeguato, in quanto non rappresenta correttamente i dati.

(γ)

Le due ipotetiche, una volta formulate completamente, si presentano così:

«Se A è vera, B è (per l'equazione 1) vera»;

«Se C è vera, allora se A è vera B è (per l'equazione 1) vera, ma (per l'equazione 2) non vera».

Queste si possono formulare come tre ipotetiche, ossia:

«Se A è vera, B è (per l'equazione 1) vera»;

«Se C è vera, allora se A è vera B è (per l'equazione 1) vera»;

«Se C è vera, allora se A è vera B è (per l'equazione 2) non vera».

Possiamo omettere la seconda di queste, in quanto non porta nessun risultato. Le altre due si possono più brevemente enunciare così:

«Se A e (1) sono vere, B è vera»;

«Se C e A e (2) sono vere, B non è vera».

E la conclusione corretta è non che C , presa da sola, non può essere vera, ma che C , A , (1) e (2) non possono essere tutte vere assieme.

Ma A è sempre possibile; cosicché possiamo, se vogliamo, assumerla come sempre vera, e non farne menzione.

Le due ipotetiche si possono ora scrivere così:

«Se (1) è vera, B è vera».

«Se C e (2) sono vere, B non è vera».

Perciò C e (1) e (2) non possono essere tutte vere insieme, anche se due qualsiasi di esse possono essere vere da sole.

Così, se C e (1) sono vere, allora (2) non può essere vera; cioè, se $a = b$ (per cui l'equazione 1 diventa « $bx + z = 5$ ») e se l'equazione 1 è vera, allora non può essere vero che $bx + z = 6$.

In secondo luogo, se C e (2) sono vere, allora (1) non può essere vera: cioè, se $a = b$ e $bx + z = 6$, allora non può essere vero che $ax + (a - b)y + z = 5$.

In terzo luogo, se (1) e (2) sono vere, allora C non può essere vera: cioè, se ambedue le equazioni date sono vere, allora a non può essere uguale a b .

Questo esempio algebrico potrebbe facilmente fuorviare un lettore incauto, per il fatto che la sua conclusione, « C non può essere vera», è vera (in base all'assunzione che le equazioni 1 e 2 sono sempre vere). La fallacia sta nel premettere la parola «Perciò», e quindi nell'asserire che questa conclusione segue dalle due ipotetiche. Questo non è vero: la ragione reale per cui C non può essere vera è che essa è incompatibile con le equazioni 1 e 2 (per sottrazione otteniamo $(a - b)(x + y) = -1$, da cui segue che $(a - b)$ non può essere uguale a 0 e cioè che a non può essere uguale a b); le due ipotetiche da sole non lo dimostrano.

(δ)

Ciò che risulta veramente dalle ipotetiche originali indicate con (i) e (ii) si può molto semplicemente esporre nel modo che segue:

Stia «v» per «vero» e «f» per «falso».

Ci sono 8 combinazioni concepibili di A , B e C , per quanto riguarda verità e falsità: queste sono:

	1.	2.	3.	4.	5.	6.	7.	8.
A.	v	v	v	v	f	f	f	f
B.	v	v	f	f	v	v	f	f
C.	v	f	v	f	v	f	v	f

Di queste, i numeri 3 e 4 sono proibiti da (i) e il numero 1 è proibito da (ii).

Le altre 5 combinazioni sono possibili; due di esse, ossia i numeri 5 e 7, contengono la condizione « C è vera», che NEMO riteneva impossibile.

[settembre 1894]

e Outis (zio Jim) quello di Dodgson. Dodgson mette l'accento sul fatto che un condizionale è falso, solo quando l'antecedente è vero e la conclusione è falsa. La tavola

di verità viene usata per provare che ci sono situazioni possibili in cui C è vera. Questo paradosso con le sue ipotetiche in contrasto ha sollevato problemi che ancor oggi sono oggetto di controversia.

This letter is a good illustration
of her extraordinary illusion.
Dodgson is liable to, from wanting study of anything like real logic.

Ch. Ch.

Oct. 28 (96)

2

My dear Wilson,

As no shadow of irritation has ever crossed my mind, with regard to anything received from you, it was pure accident that my language should have suggested it: please regard all such language as unwritten.

I fear I can't take your view — that all "lying" problems are impossible & unmeaning: but I haven't yet written that chapter in Part II: you shall see it as soon as I get it into type: meanwhile I will just say that such problems seem to me to be of two kinds — one, where the Premises cannot refer to their subject-matter: the other, where they can.

Your example is of the former kind, viz.

"A says that B is false (i.e. is speaking falsely): B says that A is true."

Let $a = A$ speaks truly; $a' = A$ speaks falsely.

Then, on hypothesis that these Propositions can refer to their subject-matter, $a \supset b'$, and $b' \supset a'$; $\therefore a \supset a'$, which is absurd. Again, $a' \supset b$, and $b \supset a$; $\therefore a' \supset a$, which is absurd. Hence hypothesis is false: i.e. these Propositions cannot refer to each other.

But, if we take the example "A says B is false: B says A is false," we get a different result.

~~viz.~~ viz. on hypothesis that these Props can refer to each other, $a \supset b'$, and $b' \supset a$, which is possible. Again, $a' \supset b$, and $b \supset a'$, which again is possible.

Hence these Propositions can refer to each other, & the Conclusion is "One of the two lies, & the other speaks truly: but we have no data to fix which is which."

Here is a very pretty example.

"A says B lies: B says C lies: C says 'You're both of you liars.'"

My solution of this is as follows:—

Take hypothesis that these Props can refer to each other. Then $a \supset b'$: $b' \supset c$: $c \supset a'$: which is absurd. $\therefore A$ cannot be speaking truly. Again, $a' \supset b$: $b \supset c'$: $c' \supset (a \text{ or } b)$: which is possible. Hence these Propositions can refer to each other, & there is one, & only one, possible state of things, viz. "A & C lie: B speaks truly."

Some of these problems (I've worked out a lot of them) give several possibilities. I prefer those that give only one. The "Five Liars" problem admits of only one possible state of things.

Yours truly,

C. L. Dodgson.

una bugia" si possa supporre ragionevolmente tale da riferirsi a se stessa come proprio oggetto ».

Egli concludeva che « Sto dicendo una bugia » non può essere autorizzata a riferirsi a se stessa « in quanto fare questo porterebbe ad un'assurdità ». Dodgson non si fermava qui. Egli proseguiva sottolineando che l'autoriferimento in sé e per sé non è criticabile e faceva rilevare che l'enunciato « Sto dicendo la verità » non porta ad assurdità quando si riferisce a se stesso.

Il punto essenziale, che Dodgson ha individuato chiaramente, è che alcuni enunciati autoreferenziali non creano nessuna difficoltà, mentre altri sono fonte di paradossi. Per esempio:

L'enunciato contenuto in questo rettangolo è vero.

L'enunciato contenuto in questo rettangolo è falso.

Alan Ross Anderson dell'Università di Pittsburgh ha recentemente fatto il punto sulla situazione che si produrrebbe se alcuni tipi particolari di autoriferimento non venissero consentiti. Egli ha scritto: « Perderemmo virtualmente tutti i campi più interessanti negli studi contemporanei sui fondamenti filosofici della matematica. I teoremi fondamentali della teoria degli insiemi e della ricorsività scomparirebbero, e logici e matematici di tutto il mondo si troverebbero disoccupati ».

Che influenza potrà avere sulla logica contemporanea la pubblicazione della seconda parte della *Logica Simbolica* di Charles Lutwidge Dodgson? Alcuni logici proveranno grande interesse nel ricostruire l'evoluzione della procedura di decisione di Dodgson e dei suoi metodi di controllo della validità; altri si divertiranno alle sue spiritose analisi di paradossi e rompicapi. L'impressione più grande sarà prodotta, credo, dagli 80 nuovi esercizi del libro. Per più di mezzo secolo i logici hanno saccheggiato per i loro testi e le loro dispense gli stravaganti problemi che Dodgson ha proposto nel primo volume della *Logica Simbolica*. In complesso non vi hanno perso sopra molto tempo, in quanto Dodgson in quel libro dava anche le soluzioni degli esercizi. Soluzioni per alcuni dei nuovi esercizi sono state trovate, principalmente nelle lettere di Dodgson a sua sorella e a Wilson: ma per la maggior parte dei problemi non è data, nel testo rimasto, nessuna soluzione. Se Lewis Carroll potesse vederci alle prese con essi, ho il sospetto che ridacchierebbe di gusto.

Il paradosso del mentitore era l'oggetto della lettera in data 28 ottobre 1896, inviata da Dodgson a Wilson. In cima alla lettera Wilson ha scarabocchiato: « Questa lettera è un buon esempio delle straordinarie illusioni di cui è vittima Dodgson, per mancato studio della vera logica, o addirittura del reale processo del pensiero. J. C. W. ».

Un nuovo livello di astrazione: la teoria delle categorie

I progressi della matematica sono segnati dai progressi nell'astrazione. Da circa un quarto di secolo hanno acquistato una grande importanza le astrazioni di terzo grado

di Lucio Lombardo Radice

La matematica è sempre astrazione. Anche ai suoi primordi. Anche nei suoi concetti più semplici: semplici, ma *concetti*.

Astrazione è il numero naturale nella sua accezione cardinale, che è forse – anzi senza forse – il primo concetto matematico elaborato dall'uomo nella sua attività. Il numero cardinale non è un segno, non è il gruppo di tacche incise dal pastore primitivo all'ingresso della grotta nella quale raccoglieva il gregge la sera, bensì l'astratto di una corrispondenza biunivoca tra due insiemi (« per ogni pecora una tacca, per ogni tacca una pecora »), un fatto mentale. Astrazione è il *punto*, l'atomo geometrico, « ciò che non ha parti » come diceva Euclide, e che tuttavia è il costituente ultimo di ogni lunghezza, di ogni superficie, di ogni solido nella loro schematizzazione mentale. Il passaggio dalla concezione fisica del punto, indivisibile ma materiale-esteso, propria dei pitagorici, alla sua trasfigurazione razionale – indispensabile per fondare la geometria – fu un vero e proprio passaggio al limite (dalla fisica alla geometria, dallo spazio ambiente reale alla sua idealizzazione) che ha richiesto secoli e forse millenni e che ha impegnato i più grandi scienziati-filosofi dell'antichità classica, dai pitagorici agli eleatici a Platone a Euclide ad Archimede.

Quando si afferma che la matematica è sempre, è tutta astrazione, si dice però ancora poco. Si traccia la linea di demarcazione tra matematica e sperimentazione naturalistica, ma non si fornisce un criterio per distinguere tra matematica e matematica. Il fatto è che la astrazione è un mondo, ed è un mondo complesso, un mondo in movimento, colle sue gerarchie o meglio i suoi livelli. Così, per dare l'esempio che tra un momento riprenderemo nel dettaglio, il concetto matematico ele-

mentare di numero (naturale, intero e frazionario) diventa, insieme ad altre astrazioni elementari a esso analoghe, il « materiale », l'« oggetto », cioè il punto di partenza per la costruzione di un concetto più generale, quello di anello del quale i numeri interi, oppure quelli frazionari, sono « esempio concreto », o meglio modello particolare. Il concetto di anello è uno dei tanti esempi di astrazione di secondo grado, cioè di sintesi mentale di concetti che derivano in modo diretto dal nostro operare sulla realtà, schematizzandolo (astrazioni di primo grado). Potremmo anche parlare di « schemi di fenomeni » e di « schemi di schemi ».

Le astrazioni di secondo grado, o schemi di schemi, hanno caratterizzato la matematica negli ultimi decenni del XIX secolo e nei primi del XX; l'hanno trasformata qualitativamente, ne hanno moltiplicato la potenza e hanno nel tempo stesso sconvolto le classificazioni tradizionali, ristrutturando dalle fondamenta l'edificio matematico. (La classificazione tradizionale era fondata su astrazioni di primo grado: aritmetica = scienza dei numeri; geometria = studio delle figure; algebra = teoria delle equazioni; la classificazione attuale è fatta invece secondo i tipi di struttura, cioè in base a schemi di schemi). Da un quarto di secolo circa hanno acquistato crescente importanza le *astrazioni di terzo grado*, cioè concetti che si collocano a un nuovo e più elevato livello di generalità, perché i loro « casi particolari » sono astrazioni di secondo grado.

Abbiamo in mente la cosiddetta « algebra universale », e più in particolare la « teoria delle categorie », l'atto di nascita della quale è una memoria del 1945 degli statunitensi S. Eilenberg e S. Mac Lane. Poiché vorremmo risultare leggibili e intelleggibili anche a coloro che, senza essere matematici di

professione, desiderano – legittimamente! – far proprie le più recenti conquiste di pensiero delle matematiche, dovremo prendere le mosse dalle astrazioni algebriche di secondo grado, ripercorrendo il più rapidamente possibile un cammino lungo oltre un secolo.

Già nel 1847 George Boole affermava che l'algebra deve occuparsi « delle operazioni in sé considerate, indipendentemente dalle materie diverse alle quali possono essere applicate ». (Lo affermava nel libro *The mathematical analysis of logic*, mettendo in evidenza le analogie formali tra disgiunzione e congiunzione logica da una parte – « o », « e » – e addizione e moltiplicazione aritmetica – « + », « · » – dall'altra.) Non possiamo qui illustrare storicamente la elaborazione e la nascita delle astrazioni algebriche di secondo grado; ci limitiamo a definire i concetti di anello e di gruppo, che diverranno a loro volta « incarnazioni » particolari del concetto più generale di categoria, esempi « concreti » di categorie. I concetti di anello e di gruppo si avviano a discendere, rapidamente, dalle università ai licei giù giù fino alla scuola media inferiore (con qualche puntata alle elementari in USA e URSS). Certamente nel giro di un decennio diverranno schemi mentali di uso corrente, non più « astrusi » del principio posizionale nella numerazione o del calcolo letterale. Ma il decennio è appena all'inizio, ed è bene perciò dare definizioni esplicite e dettagliate.

Il concetto di anello generalizza, come si è detto, quello di « numero ». Un anello è un insieme (= collezione) *A* di oggetti, di natura qualunque (non precisata), tra i quali sono definite due operazioni binarie (su coppie ordinate), anch'esse di natura non precisata. Solo convenzionalmente le chiameremo « addizione » e « moltiplicazione », e-

stendendo il simbolismo e il linguaggio ordinario: « + » = « piú »; « · » = « per »; « $a + b$ » = somma di a e b ; $a \cdot b$ prodotto di a e b (nell'ordine scritto). Per esprimere il fatto che in A sono definite le due operazioni « + » e « · » useremo il simbolo: $A(+, \cdot)$.

Il « + » e il « · » di A possono non avere niente a che fare nel concreto con il « piú » e il « per » della aritmetica ordinaria: debbono verificare però una gran parte (e preciseremo quale) delle famose proprietà formali delle operazioni, che si studiano – con un certo tedio per solito – sotto il titolo aritmetica razionale. (Tra parentesi: si studiano con tedio quando non si conoscono « + » e « · » diversi dal « piú » e dal « per » aritmetico, perché in tal caso quelle proprietà sembrano verità assolute, e con ciò banali).

Elenchiamo, raggruppandole, le proprietà formali che si assumono come assiomi, che cioè si suppongono verificate comunque si scelgano x, y, z, \dots in A :

a) Assiomi relativi alla addizione:

$a_1) x + (y + z) = (x + y) + z$, proprietà associativa;

$a_2) x + y = y + x$, proprietà commutativa;

$a_3)$ (esistenza di un elemento neutro additivo, di uno « zero »).

Esiste un elemento, che indicheremo con il simbolo 0 , tale che: $0 + x = x$ (e allora per a_2) anche $x + 0 = x$);

$a_4)$ (esistenza dell'opposto) per ogni x di A esiste un \bar{x} , suo opposto, tale che: $x + \bar{x} = 0$ (si scrive: $\bar{x} = -x$);

b) Assiomi della moltiplicazione:

$b_1) x \cdot (y \cdot z) = (x \cdot y) \cdot z$, proprietà associativa.

c) Assiomi di collegamento tra le due operazioni:

$c_1) x \cdot (y + z) = x \cdot y + x \cdot z$;

$c_2) (y + z) \cdot x = y \cdot x + z \cdot x$; sono le due proprietà distributive, e occorre scriverle tutte e due perché non si suppone a priori verificata la proprietà commutativa della moltiplicazione.

Quando in un insieme A sono definite due operazioni $+$ e \cdot che verificano gli assiomi a), b), c), diremo che abbiamo fornito l'insieme A di una struttura di anello, oppure che $A(+, \cdot)$ è un anello.

Un primo esempio di anello è offerto dagli interi relativi (positivi e negativi e lo zero; se prendiamo i soli positivi piú lo zero cade a_4). Si usa il simbolo: $Z(+, \cdot)$.

I soli numeri pari, $P(+, \cdot)$, costituiscono un altro esempio di anello, « subanello » di $Z(+, \cdot)$, e un primo esempio di anello senza elemento neutro moltiplicativo, senza « unità ». I multi-

plici di 3, 4, 5, ..., n, \dots forniscono altrettanti esempi di anelli.

Un anello è anche $Q(+, \cdot)$, l'insieme dei numeri razionali o frazionari, « quozienti » di interi. In $Q(+, \cdot)$ la moltiplicazione è commutativa; c'è la unità, 1, e inoltre ogni elemento non nullo x ammette un inverso $x^{-1} = 1/x$, con $x \cdot x^{-1} = 1$. Quando un anello verifica queste ulteriori proprietà si dice un campo. Campi sono anche i reali, $R(+, \cdot)$, e i complessi, $C(+, \cdot)$.

Anelli sono anche gli insiemi $Z[x]$, $Z[x, y], \dots$ dei polinomi in una indeterminata x , o in due indeterminate x, y , ecc., con coefficienti interi, rispetto alle usuali operazioni di addizione e moltiplicazione tra polinomi. Anelli sono anche $Z(x)$, $Z(x, y), \dots$ che si ottengono dai precedenti considerando non solo polinomi ma anche quozienti di polinomi (funzioni razionali).

L'elenco, già molto ricco, delle « strutture concrete » (astrazioni di primo grado) che rientrano nello schema formale di anello (astrazione di secondo grado), non comprende ancora nessun esempio di una amplissima classe di anelli: quelli non commutativi. Citiamo soltanto l'anello delle matrici quadrate di un dato ordine a elementi in Z (o in Q , o in R , o in C), rispetto all'addizione tra matrici (« elemento per elemento ») e alla loro moltiplicazione righe per colonne, che verificano tutti gli assiomi prescritti, ma non la proprietà commutativa della moltiplicazione (non prescritta!).

Se in un anello $A(+, \cdot)$ consideriamo solo l'operazione « + » e gli assiomi a) a essa relativi, otteniamo altrettanti esempi di gruppo commutativo, $A(+)$. Sinonimo è gruppo abeliano, aggettivo scelto in memoria di N.E. Abel, uno dei pionieri della matematica moderna. Un esempio significativo di gruppo abeliano (con operazione « + » lontanissima dall'ordinario « piú » dell'aritmetica) è fornito dai vettori di un piano o dello spazio ordinario. Il « + », in questo caso, è l'addizione di vettori eseguita colla ben nota regola del « parallelogramma delle forze » (il vettore somma di due, è la diagonale del parallelogramma che ha per lati i due vettori addendi). A ogni vettore è associata una traslazione, alla somma di due vettori la traslazione che si ottiene eseguendo successivamente (non importa in quale ordine) le traslazioni associate ai vettori addendi; ecco un altro esempio, geometrico, di gruppo abeliano, il gruppo delle traslazioni (del piano, o dello spazio).

Il concetto – piú generale – di gruppo si ottiene quando non si impone la commutatività a_2) della « composizio-

ne », ferme restando $a_1)$, $a_3)$, $a_4)$. Riserberemo il simbolo di operazione « + » al caso di un gruppo abeliano mentre nel caso di un gruppo non (necessariamente) commutativo useremo la notazione moltiplicativa con il relativo linguaggio. Diremo insomma che $G(\cdot)$ è un gruppo se esiste una composizione (moltiplicazione) associativa [vedi a_1]), con un elemento neutro da chiamarsi unità, simbolo 1 [in questo caso a_2) diventa $a'_2) 1 \cdot x = x \cdot 1 = x$], e con un inverso x^{-1} per ogni elemento x [al posto di a_3) si ha $a'_3): x \cdot x^{-1} = x^{-1} \cdot x = 1$]. I gruppi non commutativi dominano la geometria. Non commutativo è il gruppo formato dai movimenti piani (rotazioni e traslazioni), o spaziali (sono i gruppi metrici); non commutativo il gruppo delle affinità, quello delle proiezioni, ancora del piano o dello spazio; non commutativo il gruppo topologico del piano, cioè il gruppo che ha per elementi le trasformazioni biunivoche e bicontinue (senza « strappi » e senza « incollature ») del piano in sé (pensiamo al piano come se fosse fatto di gomma). La « moltiplicazione » in tutti questi gruppi geometrici, gli elementi dei quali sono trasformazioni in sé di una figura, è il cosiddetto « prodotto operatorio »: si ottiene il prodotto, per esempio, di due movimenti dati in un certo ordine eseguendoli l'uno dopo l'altro nell'ordine prescritto (il risultato è ancora un movimento, da dirsi appunto il prodotto dei due nell'ordine dato).

Gruppi, gruppi abeliani e anelli rientrano – lo si intuisce subito – nel concetto piú generale di « insieme fornito di operazioni che verificano determinate regole », cioè di « struttura algebrica » (definita assiomaticamente). Per questa via si giunge, appunto, al concetto di *struttura algebrica*, non a quello di *categoria*. Per trasformare le classi dei gruppi abeliani, dei gruppi, degli anelli (e molte altre classi ancora, di strutture algebriche e non), in esempi particolari del concetto generale di categoria, occorre mettere in luce anche un'altra analogia, piú riposta, tra le classi in questione.

Questa piú riposta analogia risiede nella possibilità di associare a due gruppi abeliani (a due gruppi, a due anelli), diciamoli A e B , un particolare insieme di rappresentazioni, o applicazioni, di A in B , e precisamente l'insieme $\text{Hom}(A, B)$ degli « omomorfismi » di A in B , ove per omomorfismo di A in B si intende una rappresentazione che rispetta (« conserva ») la struttura algebrica della quale A e B sono muniti. La cosa, detta cosí, risulta senza dubbio incomprensibile a chi

già non la sappia. Dobbiamo perciò fare qualche passo indietro, e cominciare a definire che cosa è una rappresentazione tra gli insiemi A e B , e poi introdurre in modo preciso la nozione di omomorfismo di un gruppo $A(+)$ in un gruppo $B(+)$, di omomorfismo di un anello $A(+, \cdot)$ in un anello $B(+, \cdot)$, e così via.

Una rappresentazione di un insieme A in un insieme B è una cosa molto semplice; è una corrispondenza f che associa a ogni elemento a di A un ben determinato elemento di B , immagine di a in f , che potremo denotare col l'ordinario simbolo funzionale $f(a)$, oppure con simboli che hanno lo stesso significato, e talvolta più comodi per ragioni tecniche, quali af , a^f .

Il linguaggio pittorico — « rappresentazione », « immagine » — non è casuale. Si può pensare infatti A come un oggetto, B come un quadro sul quale dobbiamo rappresentarlo. Non si ammette che un punto dell'oggetto A abbia due immagini in B (e quindi certe famose signore con doppio naso di Pablo Picasso non sono rappresentazioni nel nostro senso), ma non si pretende altro: a ogni punto di A dobbiamo associare una e una sola immagine in B , ma la possiamo scegliere a nostro capriccio, non escludendo neppure il caso-limite che tutti i punti di A abbiano come immagine un medesimo punto fisso in B . In questo ultimo caso, non viene utilizzato tutto il quadro B per rappresentare A , e così in tanti altri casi. Chiameremo immagine di A (in simboli: $\text{Im } A$) quella parte, o « sottoinsieme », di B , che è costituita dalle immagini dei punti di A e parleremo (in generale) di « rappresentazione di A in B ». Quando si utilizza tutto il quadro B , cioè quando $\text{Im } A$ coincide con B , la f si dirà una « rappresentazione di A su B » (una sovraiezione o suriezione). Se poi due punti distinti di A hanno sempre immagini distinte in f , si parlerà di *iniezione* (tra i punti di A e quelli di $\text{Im } A$ c'è una corrispondenza biunivoca: A è « iniettato » o « immerso » in B).

In una mostra d'arte non figurativa, nessuna rappresentazione di un soggetto A su un quadro B potrà essere a priori esclusa; in una mostra fotografica, invece, sono ammesse solo rappresentazioni di un tipo particolare, le rappresentazioni di un oggetto A su un quadro B che si ottengono proiettando A su B da un « centro di vista ». Nella rappresentazione topografica, la limitazione sarà ancora più pesante: vengono ammesse solo le rappresentazioni (di A in B) cosiddette « in scala »; si ammettono deformazioni, ma secondo un ben determinato rapporto

(scala). I bambini quando disegnano si permetteranno una maggiore libertà, faranno loro stessi grandi grandi accanto a una casa piccola piccola, faranno diventare storto il diritto, ovale il tondo; tuttavia, rispetteranno inconsapevolmente il principio della continuità, si serviranno di rappresentazioni (o applicazioni) *continue* dell'oggetto A sul quadro B , deformando sì le figure, ma come se fossero di gomma, senza sovrapposizioni né rotture (una circonferenza avrà per immagine un ovale, mai un segmento aperto!).

Riflettiamo un momento sulle restrizioni, tacite, poste al concetto di rappresentazione di A in B nella fotografia, nella topografia, nel disegno infantile. Qual è il criterio generale astratto che presiede a tali restrizioni concrete? È il criterio della *conservazione* nell'immagine di qualche proprietà strutturale di A . Nel caso della fotografia, si pretende che punti allineati dell'oggetto vadano in punti allineati dell'immagine (a meno che non si sovrappongano in un unico punto-immagine); nel caso della « carta topografica », che non vengano alterati i rapporti tra le lunghezze obiettive e le loro immagini; nel caso del disegno infantile, che sia conservata la relazione di contiguità o vicinanza tra punti, cioè la « struttura topologica » (in un senso da precisare, e ben precisabile).

Quando rappresentiamo un gruppo (abeliano) $A(+)$ in un gruppo (abeliano) $B(+)$, è perciò del tutto naturale limitarsi alle rappresentazioni che conservano la struttura di gruppo, cioè a quelle rappresentazioni f — da chiamarsi *omomorfismi* — di A in B tali che:

$$(+)\quad f(a+b) = f(a) + f(b)$$

(a parole: il corrispondente della somma è la somma dei corrispondenti). Se abbiamo A che fare con due insiemi A e B nei quali siano definite — per intenderci alla buona — le « stesse operazioni » (binarie), e siano $+$, \cdot , \circ , ... e quante altre vogliamo, diremo che la rappresentazione f di A in B rispetta le operazioni (è un omomorfismo tra A e B come strutture con quelle operazioni) se per ogni operazione 0 accade che:

$$(0)\quad f(a0b) = f(a)0f(b).$$

In particolare, la rappresentazione f di un anello $A(+, \cdot)$ in un anello $B(+, \cdot)$ si dirà un omomorfismo tra anelli se oltre alla (+) vale la:

$$(\cdot)\quad f(a \cdot b) = f(a) \cdot f(b).$$

Chiariamo la cosa con qualche esempio. Consideriamo $Z(+)$, cioè il gruppo abeliano formato dagli interi rispetto all'addizione. A ogni intero n facciamo corrispondere la sua metà $n/2$. La metà della somma di due interi è la somma

delle loro metà; perciò la rappresentazione degli interi nei razionali che si ha associando a ogni n intero il numero frazionario $n/2$, è un omomorfismo h di $Z(+)$ in $Q(+)$, gruppo abeliano dei razionali rispetto alla addizione. $\text{Im } h$ è l'insieme dei razionali a denominatore 2, i quali non esauriscono certo i razionali, ma ne costituiscono un sottogruppo rispetto all'addizione (e un sottoanello rispetto alle due operazioni di addizione e moltiplicazione). Invece h non è un omomorfismo dell'anello $Z(+, \cdot)$ nell'anello (campo) dei razionali $Q(+, \cdot)$, perché il prodotto delle metà di due interi non è la metà, bensì la quarta parte, del prodotto dei due interi stessi.

Consideriamo invece la rappresentazione t di Z in Q che a ogni intero n fa corrispondere il razionale $n + 1/2$:

$$t(n) = n + 1/2.$$

È subito visto che, mentre:

$$t(n+m) = n+m+1/2,$$

si ha invece:

$$t(n) + t(m) = n + 1/2 + m + 1/2 = n + m + 1;$$

la t è una rappresentazione dell'insieme Z nell'insieme Q , ma non è un omomorfismo del gruppo abeliano $Z(+)$ nel gruppo abeliano $Q(+)$.

Consideriamo l'insieme $Q' = Q - 0$ dei numeri frazionari, (razionali) non nulli. Essi formano un gruppo (abeliano) $Q'(\cdot)$ rispetto alla moltiplicazione. Anche i soli razionali positivi, Q'_+ , costituiscono un gruppo $Q'_+(\cdot)$, rispetto alla moltiplicazione. Infine, l'insieme formato dai due soli numeri $+1$, -1 , chiamiamolo S (iniziale di « segno »), costituisce del pari un gruppo rispetto alla moltiplicazione. Elenchiamo i seguenti omomorfismi:

p : di $Q'_+(\cdot)$ in $Q'(\cdot)$, definito facendo corrispondere a ogni reale positivo se stesso;

q : di $Q'(\cdot)$ in S , definito facendo corrispondere a ogni positivo il numero $+1$ di S , a ogni negativo il numero -1 (il corrispondente del prodotto è sempre il prodotto dei corrispondenti in virtù della ben nota regola dei segni).

Il primo è « iniettivo », in quanto elementi distinti hanno immagini distinte; il secondo è « suriettivo » ma non iniettivo, perché tutti i positivi hanno la stessa immagine ($+1$), e così tutti i negativi (-1). Si chiama *nucleo* di un omomorfismo f di un gruppo A (abeliano o non) in un gruppo B , e si indica con il simbolo $\text{Ker } f$, l'insieme degli elementi di A che hanno per immagine in B l'elemento neutro di B (Ker è abbreviazione dell'inglese *kernel* = nucleo). Se f è un omomorfismo di un anello A in un anello B , si chiama $\text{Ker } f$ l'insieme degli elementi di A che hanno per immagine lo 0 di B . (Diamo un

esempio di omomorfismo di anelli. Sia Z_2 l'anello composto da due soli elementi: $\bar{0}$ = classe degli interi pari, $\bar{1}$ = classe degli interi dispari. Allora si ha: $\bar{0} + \bar{0} = \bar{1} + \bar{1} = \bar{0}$; $\bar{0} + \bar{1} = \bar{1} + \bar{0} = \bar{1}$; $\bar{0} \cdot \bar{1} = \bar{1} \cdot \bar{0} = \bar{0}$, $\bar{0} \cdot \bar{0} = \bar{0}$, $\bar{1} \cdot \bar{1} = \bar{1}$, regole di calcolo che, lungi dall'essere folli, traducono in simboli fatti ben noti: pari + pari = pari = pari = dispari piú dispari ecc. Se a ogni pari di Z associamo lo $\bar{0}$ di Z_2 , a ogni dispari l' $\bar{1}$, abbiamo un omomorfismo di $Z(+, \cdot)$ su $Z_2(+, \cdot)$ che ha per nucleo il sottoanello $P(+, \cdot)$ dei pari.)

Risulta subito che:

$\text{Ker } p$ = elemento neutro di Q'_+ , e basta. (In generale, in un omomorfismo iniettivo di A in B il nucleo si riduce al solo neutro di A).

$\text{Ker } q = Q'_+$, sottogruppo dei reali positivi.

Riprenderemo questi esempi dopo una parentesi.

Data una rappresentazione p di A in B e una rappresentazione q di B in C , posso costruire la rappresentazione $p \circ q$ di A in C nella quale all'elemento a di A corrisponde l'elemento c di C corrispondente in q del corrispondente di a in p ; in simboli:

$$p \circ q(a) = q(p(a)).$$

A parole: la rappresentazione composta $p \circ q$ si ottiene applicando successivamente p e q (prima la p , poi la q), cioè facendone il prodotto operatorio, il che ha un senso perché p porta da A in B e q riparte dallo stesso B per farci finire in C .

A questo punto, viene piuttosto spontaneo aiutarsi con « frecce ». Invece di dire che p è una rappresentazione di A in B , faremo questo « disegno »:

$$A \xrightarrow{p} B$$

Se poi vogliamo esprimere il fatto che la rappresentazione f di A in C è il prodotto operatorio, $p \circ q$, della rappresentazione p di A in B per quella q di B in C , diremo che il diagramma

$$\begin{array}{ccc} A & \xrightarrow{p} & B \\ & \searrow f & \downarrow q \\ & & C \end{array}$$

è *commutativo*, cioè che si può andare da A in C indifferentemente per le due vie indicate, cioè in definitiva che:

$$f(x) = p \circ q(x) = q(p(x))$$

per ogni x in A .

Se le rappresentazioni fattori p e q sono omomorfismi, per esempio di un gruppo A in un gruppo B , di un gruppo

B in un gruppo C , si controlla senza difficoltà, sulla base della definizione di omomorfismo, che anche il loro prodotto operatorio è un omomorfismo (di gruppi, di anelli o di altra struttura, secondo i casi).

Se ora ho tre rappresentazioni: la p di A in B , la q di B in C , la r di C in D , posso definire una rappresentazione $p \circ q \circ r$ di A in D partendo da un x in A e andando a finire nell'elemento $r(q(p(x)))$ in D (che è, a parole, il corrispondente in r del corrispondente in q del corrispondente in p di x); aiutandosi con frecce:

$$\begin{array}{ccccc} A & \xrightarrow{p} & B & \xrightarrow{q} & C & \xrightarrow{r} & D \\ & & & & \searrow p \circ q \circ r & & \end{array}$$

cioè

$$x \longrightarrow p(x) \longrightarrow q(p(x)) \longrightarrow r(q(p(x))).$$

Il punto di arrivo dipende *solo* dal punto di partenza, non dal modo in cui si associano i singoli passaggi; in formule:

$$(p \circ q) \circ r = p \circ (q \circ r).$$

Insomma, quando di tre rappresentazioni si può fare il prodotto operatorio, perché l'insieme « di arrivo » della prima è quello « di partenza » della seconda e quello di arrivo della seconda è di partenza per la terza, allora tale prodotto è associativo. Lo stesso risultato vale se al posto di « rappresentazione di un insieme A in un insieme B » si parla di « omomorfismo di un gruppo A in un gruppo B », di « omomorfismo di un anello A in un anello B ».

Diremo che la successione di omomorfismi

$$A \xrightarrow{f} A' \xrightarrow{g} A''$$

è *esatta* se $\text{Im } f = \text{Ker } g$. Per esempio, la successione

$$Q'_+ \xrightarrow{p} Q' \xrightarrow{q} S$$

è esatta per quanto si è visto poco fa, prima di aprire questa parentesi.

La « astrazione di secondo grado » che ci accingiamo a fare è, a questo punto, molto naturale. Invece di parlare della classe degli insiemi (« tutti gli insiemi » non costituiscono un insieme, ma una collezione « piú grossa », una « classe »), o di quella dei gruppi, o di quella degli anelli, parliamo di una classe di *oggetti*, da « particularizzare » volta a volta negli insiemi, nei gruppi, negli anelli. Dunque, dare una categoria C vuol dire innanzitutto considerare

una classe di oggetti, in simboli: $\text{Ob } C$. Ma non basta, per quanto sopra detto. Scelti comunque due oggetti, e siano A e B , in un certo ordine, per esempio prima A e poi B , alla coppia (A, B) deve rimanere associato un insieme di *morfismi* di A in B , $\text{Hom}(A, B)$. (Per essere precisi occorrerebbe scrivere $\text{Hom}_C(A, B)$ per indicare che si tratta dei morfismi da A e B nella categoria C). I morfismi saranno da particularizzare volta a volta nelle rappresentazioni di un insieme A in un insieme B , negli omomorfismi di un gruppo A in un gruppo B , oppure di un anello A in un anello B , ecc.

Si pretende poi, e anche questo è naturale dopo i « casi particolari » esaminati, che si possa comporre (moltiplicare) un morfismo da A a B con un morfismo da B a C in modo da ottenere un morfismo da A a C , e che detta composizione (nei casi in cui è definita) sia associativa. In altri termini, dati p in $\text{Hom}(A, B)$, q in $\text{Hom}(B, C)$ e r in $\text{Hom}(C, D)$, si ha:

$$a) \quad p \circ (q \circ r) = (p \circ q) \circ r,$$

dove i morfismi composti a primo e a secondo membro sono definiti per l'ipotesi sopra fatta della componibilità (simbolo \circ) di un morfismo da A a B con un morfismo da B a D e, rispettivamente, di un morfismo da A a C con un morfismo da C a D .

Non basta. Nei casi « concreti » (cioè... meno astratti) sopra considerati, comunque scelto l'oggetto A esisteva sempre un morfismo identico, e precisamente la rappresentazione identica dell'insieme A in se stesso: l'omomorfismo di un gruppo – oppure di un anello – A in se stesso che si ottiene facendo corrispondere a ogni elemento di A se stesso. Il morfismo identico, o identità, 1_A da A in A gode in tutti questi casi concreti di una proprietà caratteristica esprimibile nei termini della sola composizione di morfismi: 1_A è elemento neutro, o indifferente, in tale composizione, quando lo si componga a destra con un morfismo da B (qualunque) in A , oppure a sinistra con un morfismo di A in B (qualunque). Insomma, se p è un morfismo di B in A , q un morfismo da A in B :

$$i) \quad p \circ 1_A = p, \quad 1_B \circ q = q$$

(si pensi a rappresentazioni tra insiemi; un x di B va in un $p(x)$ di A , il quale resta fermo se si applica successivamente l'identità di A , 1_A , ecc.). La cosa risulta piú chiara con un disegno:

$$\begin{array}{ccc} & A & \\ p \nearrow & & \searrow 1_A \\ B & & A \end{array} \quad \begin{array}{ccc} & A & \\ 1_A \nearrow & & \searrow q \\ A & & B \end{array}$$

Aggiungiamo infine la richiesta che $\text{Hom}(A, B)$ e $\text{Hom}(A', B')$ non abbiano nessun elemento comune se le coppie (A, B) e (A', B') sono distinte, e siamo arrivati a una rigorosa definizione astratta di categoria. La riassumiamo:

Una categoria C è una classe di oggetti, $\text{Ob } C$, tale che a ogni coppia ordinata di oggetti, (A, B) , è associato un insieme di elementi, detti morfismi di A in B , $\text{Hom}(A, B)$; si ha inoltre che i morfismi di $\text{Hom}(A, B)$ si possono comporre con quelli di $\text{Hom}(B, C)$ dando luogo a morfismi di $\text{Hom}(A, C)$, e che tale composizione è associativa (quando è definita); in $\text{Hom}(A, A)$, comunque si scelga l'oggetto A , esiste un morfismo identico, I_A , cioè tale che valgono le (i); se le coppie (A, B) , (A', B') sono distinte, $\text{Hom}(A, B)$ e $\text{Hom}(A', B')$ non hanno elementi a comune.

Siamo pervenuti al concetto generale di categoria partendo dalle classi degli insiemi, dei gruppi, dei gruppi abeliani, degli anelli e dai loro « morfismi », assiomatizzando (« schematizzando ») alcune proprietà formali comuni ai vari casi. Perciò, viceversa, abbiamo subito a disposizione quattro esempi di categorie: 1) la categoria degli insiemi, Ens (iniziali di *ensemble*), nella quale i morfismi sono le ordinarie rappresentazioni, o applicazioni di un insieme in un altro. (In altri termini ancora, le *funzioni* che hanno per dominio A e per codominio B). 2) La categoria dei gruppi, Grp ; i morfismi sono gli omomorfismi gruppali (applicazioni nelle quali il corrispondente del composto è il composto dei corrispondenti: esse rispettano cioè la composizione gruppale). 3) La categoria dei gruppi abeliani o commutativi, Ab , *sotlocategoria* della precedente in quanto gli oggetti di Ab sono oggetti anche di Grp e i morfismi $\text{Hom}(A, B)$ con A e B in Ab sono anche morfismi di Grp . 4) La categoria degli anelli, Rng (abbreviazione di *ring*); morfismi sono gli omomorfismi di anelli, cioè le rappresentazioni di un anello A in un anello B che conservano, o meglio « rispettano », tutte e due le operazioni definite in un anello.

Potremmo elencare molte altre categorie relative all'algebra: quella dei moduli, quella dei reticoli, quella delle « omega-algebre » per un certo complesso di operazioni « omega ». Ma esse non ci darebbero nulla di qualitativamente originale, rispetto agli esempi 2), 3), 4). Si tratta, infatti, caso per caso, di prendere insieme con certe operazioni, verificanti eventualmente determinati assiomi, come oggetti di una ca-

tegoria C , e come morfismi tra coppie di elementi di C le rappresentazioni che conservano le operazioni. Preferiamo perciò aggiungere due esempi qualitativamente diversi dai precedenti e tra di loro. 5) La categoria degli spazi topologici, Tps (*topological spaces*). Dati due spazi topologici A e B , l'insieme $\text{Hom}(A, B)$ è costituito dalle applicazioni continue di A in B . (Non possiamo dare definizioni rigorose; ancora una volta, ci riferiamo alla idea intuitiva di « rappresentazione continua » come « disegno » che altera sì le forme e le distanze, rispettando però i rapporti di continuità. 6) La categoria che ha per unico oggetto il « monoide » moltiplicativo degli interi non nulli, $\mathbb{Z}(\cdot)$, e come morfismi di \mathbb{Z} in \mathbb{Z} le corrispondenze (moltiplicazioni per un numero fisso z):

$$a \longrightarrow a \cdot z \quad (a \text{ variabile in } \mathbb{Z}).$$

(Gli interi, zero escluso, rispetto alla moltiplicazione hanno un elemento neutro, l'1, e godono della proprietà associativa, ma in generale un intero non possiede un inverso nell'ambito degli interi; perciò $\mathbb{Z}(\cdot)$ non è un gruppo, ma — appunto — un « monoide ».)

Bastano già questi pochi esempi, crediamo, per far comprendere la « portata » del concetto di categoria. Una portata così vasta, da far temere — a prima impressione — che la teoria sia povera di risultati, perché troppo estesa. Così non è. Prima di fare però qualche osservazione (conclusiva) sulla problematica e sui risultati della teoria delle categorie, vogliamo introdurre altri due concetti fondamentali: quello di categoria *duale*, C^* , di una categoria C , e quello di *funtore* da una categoria C a una categoria C' .

« Dualizzare » significa molto spesso in algebra « scambiare la destra col la sinistra ». Nel caso di una categoria, C , posso costruire un'altra categoria, la sua *duale* (o *opposta*), C^* , oppure C^{op} , lasciando immutati gli oggetti, e chiamando (ribattezzando) morfismo da B ad A un morfismo da A a B della categoria C . Insomma:

Categoria *duale* C^* di una categoria C :

a) $\text{Ob } C^* = \text{Ob } C$;

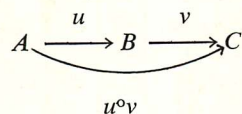
b) i morfismi da B ad A nella C^* sono tutti e soli i morfismi da A a B della C ; cioè u appartiene a $\text{Hom}_{C^*}(B, A)$ se e soltanto se u appartiene a $\text{Hom}_C(A, B)$.

I conti si fanno tornare, definendo il prodotto vu in C^* tra un morfismo v da C a B e un morfismo u da B ad A come il morfismo $u \circ v$ composto di u e v nella C ; $u \circ v$ infatti è nella C un morfismo da A a C , e quindi nella sua *duale* un morfismo da C ad A . Insomma,

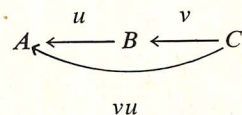
facendo attenzione all'ordine in cui sono scritti i fattori, il prodotto di v per u nella *duale* è il prodotto di u per v nella categoria di partenza. (Abbiamo usato simboli diversi per la composizione di morfismi in C e in C^* : $u \circ v$ e vu rispettivamente).

La dualità, scambio di posto tra destra e sinistra, si esprime efficacemente nella teoria delle categorie come cambiamento di verso delle frecce. Infatti il passaggio dalla C alla sua *duale* C^* , che forse è risultato oscuro a parole, diventa chiarissimo graficamente: basta cambiare verso alle frecce!

In C :



In C^* :



Il concetto di *funtore* generalizza quello di funzione (applicazione, rappresentazione). Poiché in una categoria abbiamo a che fare non solo con oggetti, ma anche con morfismi, un *funtore* dalla categoria C alla categoria C' sarà una « doppia funzione », F ; si ha cioè che:

a) a ogni oggetto A di C corrisponde un oggetto $F(A)$ di C' ;

b) a ogni morfismo h di A in B , A e B oggetti di C , corrisponde un morfismo $F(h)$ di $F(A)$ in $F(B)$, in modo tale che:

b₁) sia conservata la composizione di morfismi:

$$F(h \circ k) = F(h) \circ F(k) \quad (k \text{ è un morfismo da } B \text{ in } C);$$

b₂) siano conservate le identità:

$$F(1_A) = 1_{F(A)}$$

(l'identità dell'oggetto A di C viene portata nella identità dell'oggetto corrispondente ad A nella C' , tramite F).

Questa è, per essere precisi, la definizione di *funtore covariante* da C a C' ; la definizione di *funtore controvariante* si ottiene sostituendo alla b) la

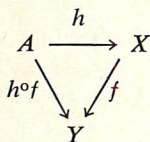
b') a ogni h di $\text{Hom}(A, B)$ corrisponde un $F(h)$ di $F(B)$ in $F(A)$ in modo tale che:

$$b'_1) F(h \circ k) = F(k) \circ F(h).$$

Quando passiamo dalla categoria C alla sua *duale* C^* , facendo corrispondere a ogni oggetto se stesso, al morfismo u di A in B (in C) un morfismo $u^* = u$ di B in A nella C^* , introduciamo appunto un *funtore controvariante* dalla C alla C^* . Un esempio banale di *funtore covariante*, da una categoria C a

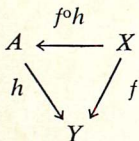
se stessa, è la « doppia identità »: tanto a ogni oggetto, quanto a ogni morfismo, si fa corrispondere se stesso. Ancora due esempi di funtori, l'uno covariante, l'altro controvariante.

Funtore morfismo covariante. Si fissa un oggetto A in una categoria C ; a ogni oggetto X di C si fa corrispondere l'insieme $\text{Hom}(A, X)$ dei morfismi di A in X . In questo modo, a ogni oggetto X della categoria C resta associato un oggetto, $h_A(X) = \text{Hom}(A, X)$, della categoria degli insiemi, Ens . Sia dato un morfismo f da X a Y . A f possiamo associare una rappresentazione (morfismo di insiemi), $h_A(f)$ dell'insieme $\text{Hom}(A, X)$ nell'insieme $\text{Hom}(A, Y)$, nel modo suggerito da questo diagramma:



Precisamente; a ogni h di $\text{Hom}(A, X)$ facciamo corrispondere l'elemento $h \circ f$ di $\text{Hom}(A, Y)$.

Funtore morfismo controvariante. Fissato A nella C , a ogni oggetto X della C facciamo corrispondere l'insieme $\text{Hom}(X, A)$ dei morfismi di X in A ; a ogni morfismo f di X in Y possiamo associare allora una rappresentazione di $\text{Hom}(Y, A)$ in $\text{Hom}(X, A)$ (attenzione, c'è stato lo scambio di posto tra X e Y !), e precisamente quella che fa corrispondere al morfismo h di Y in A il morfismo $f \circ h$ di X in A , come il diagramma illustra:



Ci troviamo a questo punto di fronte a una difficoltà. Perché la teoria delle categorie unifichi le diverse teorie di « secondo livello » che in essa si inquadrano (insiemi, gruppi, anelli, spazi topologici e così via), occorrerà bene tradurre in *linguaggio categorico* alcuni, almeno, tra i concetti fondamentali che abbiamo già richiamato, e che elenchiamo ora alla rinfusa: iniezione, suriezione, elemento neutro (« zero » nella scrittura additiva), nucleo, immagine, successione esatta. Non è detto che ogni categoria abbia un elemento zero; ma come possiamo esprimere in linguaggio puramente categorico il fatto che la categoria dei gruppi abeliani, o quella degli anelli, hanno un elemento zero, cioè che esiste un gruppo abeliano (o

un anello) ridotto al solo elemento neutro additivo?

La risposta è sorprendentemente più facile di quanto non ci si aspetti a prima impressione. Consideriamo un gruppo abeliano qualunque, A , e lo zero-gruppo abeliano, composto dal solo elemento 0 , colla regola $0 + 0 = 0$. Ebbene: a) esiste uno e un solo omomorfismo da 0 in A : quello che manda 0 nell'elemento neutro di A (un omomorfismo tra gruppi è obbligato a spedire l'elemento neutro del primo nell'elemento neutro del secondo); b) esiste uno e un solo omomorfismo da A a 0 . (Se c'è, deve essere per forza la applicazione che manda ogni elemento di A nell'unico elemento, 0 , dello zero-gruppo, perché... non ci sono altri posti nei quali collocare le immagini degli elementi di A ; tale rappresentazione è effettivamente un omomorfismo, perché, essendo $0 + 0 = 0$, si ha banalmente che l'immagine della somma di due elementi è sempre la somma delle loro immagini.)

Ma la a) e la b) possono subito essere trascritte in termini strettamente categorici. Diremo che una categoria C possiede uno zero-oggetto, 0 , se:

a) $\text{Hom}(0, X)$ contiene uno e un solo elemento per ogni X di C ; in questo caso si dice che 0 è un *oggetto iniziale*;

b) $\text{Hom}(X, 0)$ contiene uno e un solo elemento per ogni X di C ; si dice in tal caso che 0 è un *oggetto finale*.

(I concetti di oggetto iniziale e oggetto finale, sia osservato tra parentesi, sono *duali*.) La categoria degli insiemi possiede oggetti finali: essi sono gli insiemi composti da un singolo punto, o *singletons* all'inglese. Possiede anche un oggetto iniziale, che è l'insieme vuoto (privo del tutto di elementi; allora l'unica rappresentazione possibile dell'insieme vuoto in un altro insieme è quella vuota, nella quale... non c'è da far corrispondere niente a niente, per mancanza di elementi da rappresentare). La categoria degli insiemi non possiede però uno zero-oggetto, che è invece posseduto dalla categoria dei gruppi abeliani, da quella degli anelli, da quella dei gruppi (è il gruppo identico, che scritto moltiplicativamente è il gruppo con il solo elemento 1 , con la composizione $1 \cdot 1 = 1$).

Diamo di colpo la traduzione categorica dei concetti insiemistici di *iniezione* e *suriezione*, che sono resi rispettivamente dai termini *monomorfismo* ed *epimorfismo*, tra di loro duali (come il lettore vedrà dai disegni).

Un morfismo da A a B si chiama un *monomorfismo* se, quale che sia C , da $f \circ a = g \circ a$ (con f e g morfismi da C ad

A), si può dedurre:

$$f = g.$$

Insomma se

$$\begin{array}{ccc} & f & \\ C & \xrightarrow{\quad} & A \xrightarrow{\quad} B \\ & g & \end{array}$$

è commutativo, allora $f = g$.

Dualmente, un morfismo a da B ad A si chiama un *epimorfismo* se da $a \circ f = a \circ g$, con f e g morfismi da A a C , si deduce $f = g$ (quale che sia C). Invertiamo le frecce, e abbiamo:

$$\begin{array}{ccc} & f & \\ B & \xrightarrow{\quad} & A \xrightarrow{\quad} C \\ & g & \end{array}$$

è commutativo, allora $f = g$.

Giustificiamo, almeno a metà, la seconda definizione, facendo vedere che una suriezione in senso insiemistico è un epimorfismo in senso categorico. Se B, A, C sono insiemi, e se a è una applicazione di B su A , allora comunque si scelga un elemento x in A , esiste un y in B tale che $a(y) = x$. Supponiamo che si abbia $a \circ f = a \circ g$; ciò vuol dire che per ogni y in B si ha $a \circ f(y) = a \circ g(y)$, cioè:

$f(a(y)) = g(a(y))$ per ogni y in B , cioè ancora, visto che a è una suriezione:

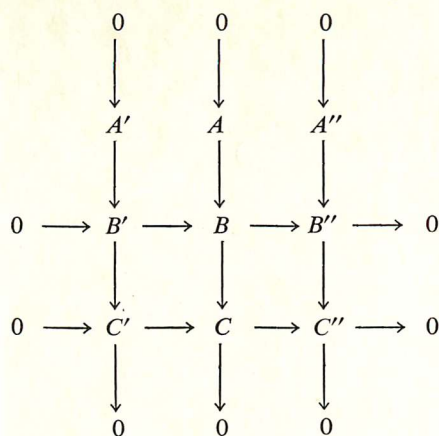
$$f(x) = g(x) \text{ per ogni } x \text{ in } A,$$

il che significa, appunto, f identica a g . Poiché si dimostra facilmente che se a non è una suriezione tra gli insiemi B e A , allora a non è neppure un epimorfismo, ecco che la definizione categorica di epimorfismo corrisponde a una proprietà caratteristica delle suriezioni tra insiemi. Analogamente, si vede che i monomorfismi, nella categoria degli insiemi, sono tutte e sole le iniezioni.

Si apre, è il caso di dirlo, un nuovo mondo; o meglio, si può costruire un nuovo vocabolario, nel quale acquistano senso categorico i termini « immagine », « nucleo », « inclusione », « sotto-oggetto », « sequenza esatta » e così via; il nuovo significato, categorico, collima con il significato « ordinario », insiemistico o gruppale, dei termini.

Siamo alla fine del nostro articolo, siamo al principio della teoria delle categorie. Possiamo forse azzardarci a trascrivere, lasciando nel vago qualche termine, l'enunciato del primo teorema di teoria delle categorie di una certa consistenza, il cosiddetto « Lemma dei nove soggetti » o brevemente « Lemma dei nove ».

A pagina 20 del libro di Mitchell *Theory of Categories* troviamo il seguente « diagramma commutativo » (non ha importanza la strada che si segue per andare da un posto all'altro seguendo il verso delle frecce):



Supponiamo che tutte le righe e tutte le colonne siano successioni (o sequenze) esatte, cioè che il nucleo di ogni morfismo sia la immagine di quello che lo precede. Bene, in tali ipotesi e in una classe particolare, ma abbastanza va-

sta, di categorie (le « categorie esatte ») il diagramma si può completare, in modo che resti commutativo, aggiungendo morfismi tra A' e A , e tra A e A'' , in modo che anche la nuova riga:

$0 \longrightarrow A' \longrightarrow A \longrightarrow A'' \longrightarrow 0$
risulti esatta.

Rileggo quanto ho scritto, e dispero che qualche non matematico sia arrivato in porto. Temo, perciò, di avere lavorato, tutt'al più, a beneficio di qualche matematico buon conoscitore di insiemi, gruppi, anelli, anche se non di categorie.

Tuttavia il compito iniziale resta. E non si tratta di un lusso intellettuale, del desiderio di completezza culturale nell'uomo di scienza. No. Il fatto è, invece, che occorre fortemente accelerare i tempi di passaggio dalla cultura specialistica alla cultura scientifica media (di massa) dei livelli di astrazione algebri-

ca dei quali abbiamo discorso. È sembrato già un successo, otto anni fa, introdurre le astrazioni di secondo grado (gruppi e anelli, soprattutto) al primo anno del corso di laurea in matematica. Ora i fisici si chiedono se quei concetti non *debbero* essere acquisiti anche dagli studenti del primo anno di fisica. Oggi, leggendo quello che a mio avviso è il miglior manuale di algebra per un primo biennio di matematica, il volume di Mac Lane-Birkhoff, vediamo che esso è fondato sui concetti di categoria, di funtore, di soluzione (elemento) universale, che compaiono sin dalle prime pagine.

La verità è che le astrazioni, più potenti sono, più sono semplici, e più sono utili. Ogni passaggio a un nuovo livello di astrazione è la liberazione da una gran quantità di nozioni, di concetti particolari, perché astrazione è unificazione, ed è chiarificazione.

Induzione e probabilità

L'evoluzione del concetto di probabilità e la « crisi dei fondamenti » hanno strettamente legato la teoria della probabilità alla moderna logica induttiva

di Domenico Costantini e Marco Mondadori

Sin dall'antichità è considerato compito della logica analizzare le nostre argomentazioni, cioè analizzare i ragionamenti mediante i quali a partire da certe premesse arriviamo a certe conclusioni. Nel Medioevo si individuano due tipi di logica: la logica formale, o minore, e quella materiale, o maggiore. La prima avrebbe dovuto occuparsi di analizzare la correttezza delle argomentazioni prescindendo dai loro contenuti. La seconda invece avrebbe dovuto consentire di stabilire le condizioni alle quali i ragionamenti non solo sono corretti ma sono anche veri, arrivano cioè alla scienza « vera ». Dalla fine del secolo scorso però la pretesa di costruire una logica materiale è stata definitivamente abbandonata, almeno dalla maggior parte degli studiosi di logica, di guisa che ora per logica si intende solo la logica formale.

Le nostre argomentazioni possono essere suddivise in due grandi tipi: da una parte quelle che prendono le mosse da premesse vere o presunte tali e consentono di arrivare a conclusioni vere (inferenze dimostrative); dall'altra quelle che non possiedono questa caratteristica, quelle cioè che non conservano la verità delle premesse nel senso che pur partendo da premesse vere o presunte tali non escludono la possibilità di pervenire a conclusioni false (inferenze non dimostrative). La disciplina che studia a livello formale le inferenze dimostrative è nota col nome di logica simbolica o matematica; siamo però convinti che sarebbe più opportuno chiamarla logica deduttiva. Essa risale sostanzialmente ad Aristotele, fu largamente studiata nel Medioevo e fu sempre nettamente distinta dalla logica materiale.

La disciplina che studia le inferenze non dimostrative è nota come logica induttiva. A differenza di quanto avven-

ne per la logica deduttiva, la distinzione fra logica formale e materiale a livello delle inferenze non dimostrative non è stata affatto netta nel senso che la convinzione che anche la logica induttiva dovesse essere sviluppata come una disciplina formale si è affermata solo molto recentemente se si escludono le idee di G. W. Leibniz a questo riguardo, idee che però restarono del tutto isolate. Infatti il termine « logica induttiva » è stato tradizionalmente usato per designare un insieme di ricerche e di dottrine che hanno la loro origine in F. Bacone e raggiungono il loro massimo sviluppo con J. S. Mill. Nell'ambito di questa tradizione l'induzione viene concepita, appunto da Mill, come un processo inferenziale mediante il quale « concludiamo che ciò che è vero di certi individui di una classe è vero per l'intera classe, o che ciò che è vero in certi tempi sarà vero, in circostanze simili, in tutti i tempi ».

Ma, dato che per l'empirismo classico lo scopo della scienza è la scoperta di leggi generali, ne consegue che la logica induttiva è intesa da Bacone e Mill come la logica della scoperta scientifica, cioè come lo studio e l'esplicitazione dell'insieme di quei metodi che ci portano alla scoperta delle leggi scientifiche, in ultima analisi, come una logica materiale. Ed è appunto in questo quadro che si collocano i famosi quattro metodi milliani della concordanza, della differenza, delle variazioni concomitanti e dei residui.

Induzione e probabilità

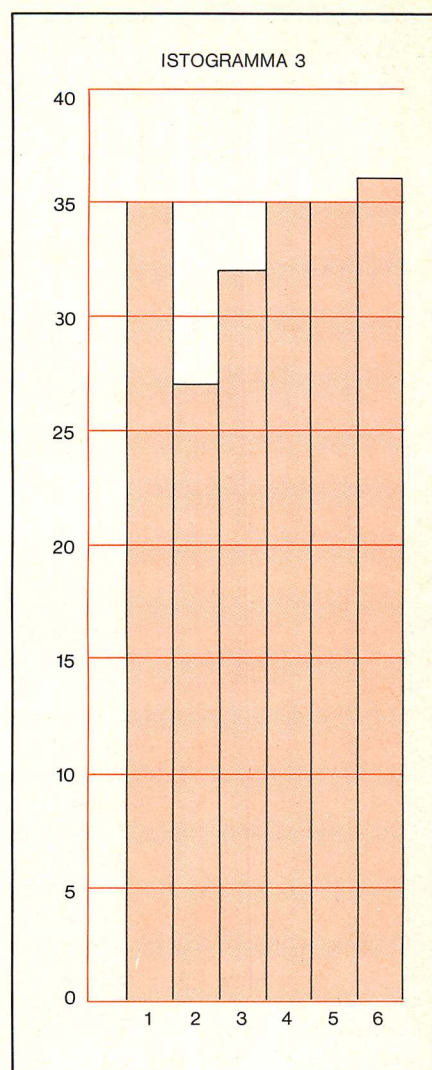
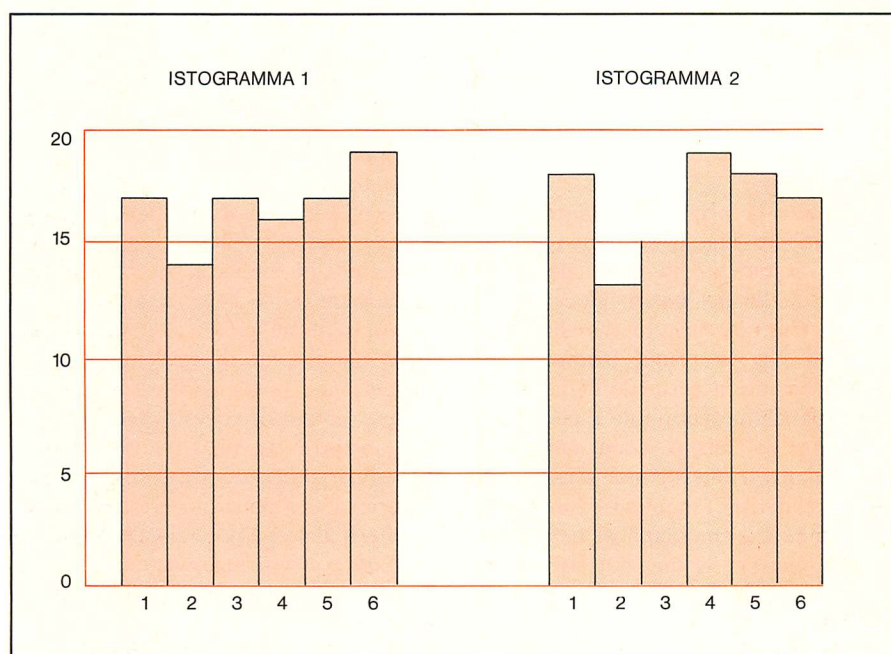
Questo modo di intendere il compito e la portata della logica induttiva fu universalmente accettato fino agli inizi del nostro secolo. Come abbiamo detto vi era però un modo alternativo di intendere la logica induttiva anche se esso ha avuto a livello filosofico « ufficia-

le » un'influenza oltremodo scarsa. Il primo rappresentante di questa ulteriore concezione è – lo abbiamo già visto – Leibniz secondo il quale le leggi generali « non si scoprono mediante l'induzione dagli esempi »; mediante questa induzione possiamo semplicemente giungere ad assegnare un certo grado di probabilità a leggi generali già date. Secondo Leibniz, quindi, la logica induttiva è una disciplina formale nel senso che non è suo compito scoprire le leggi di natura: essa deve limitarsi ad assegnare un grado di probabilità alle leggi di natura delle quali in qualche modo già disponiamo.

Con Leibniz troviamo quindi la prima intuizione del rapporto tra probabilità e induzione, rapporto che costituisce il punto centrale delle moderne ricerche di logica induttiva. Dovette però passare quasi un secolo prima che la geniale intuizione leibniziana riuscisse ad assumere, per opera di P. S. Laplace, la sua prima parziale formulazione scientificamente corretta. Infatti, fu soltanto quando la teoria delle probabilità assunse una veste scientificamente soddisfacente – e ciò avvenne grazie all'opera del grande scienziato francese – che fu possibile disporre di un quadro concettuale adeguato alla formulazione di una concezione della logica induttiva come quella di Leibniz.

Con Laplace viene operata la prima sistemazione rigorosa dei metodi con cui è possibile derivare, da probabilità date di ipotesi, probabilità di ipotesi logicamente connesse con le prime, ipotesi in generale più complesse di quelle di partenza. Se, per esempio, conosciamo la probabilità dell'ipotesi *A* e quella dell'ipotesi *B*, e sappiamo che *A* e *B* sono incompatibili, la teoria delle probabilità ci dice che la probabilità dell'ipotesi « *A* oppure *B* » è uguale alla somma della probabilità di *A* e di quella di *B*. La teoria delle probabilità

FACCE DEL DADO	FREQUENZE RELATIVE PRIMA SERIE	FREQUENZE RELATIVE SECONDA SERIE	FREQUENZE RELATIVE TERZA SERIE
1	0,17	0,18	0,175
2	0,14	0,13	0,135
3	0,17	0,15	0,160
4	0,16	0,19	0,175
5	0,17	0,18	0,175
6	0,19	0,17	0,180



Nei due istogrammi a sinistra sono visualizzati i risultati di due serie di 100 lanci effettuati con un dado comperato in un supermercato qualsiasi. Nel terzo istogramma è invece visualizzato il

risultato delle due stesse serie di lanci considerate come un'unica serie di 200 lanci. Le frequenze relative della tabella in alto sono desunte rispettivamente dal primo, secondo e terzo istogramma.

assolve infatti questo compito e, come è bene precisare fin d'ora, *solo* questo compito. Essa non ci dice infatti come debbano essere determinate le probabilità delle ipotesi di partenza; nell'esempio considerato, non ci dice come dobbiamo determinare la probabilità di *A* e quella di *B*.

In realtà Laplace si era interessato anche alla determinazione delle probabilità mediante la sua famosa definizione del termine « probabilità » di cui ci occuperemo tra poco. Prima però di analizzare questa definizione è opportuno ribadire un fatto che divenne chiaro solo all'inizio del nostro secolo e al quale abbiamo già fatto cenno quando pocanzi abbiamo detto che la teoria delle probabilità consente solo di passare da probabilità note a probabilità sconosciute legate alle prime. Questo fatto può essere espresso nel modo seguente: la teoria delle probabilità non permette di individuare in modo univo-

co il significato da attribuire al termine « probabilità » né di formulare alcuna indicazione sul modo con cui le probabilità delle ipotesi devono essere determinate. Questa caratteristica della teoria delle probabilità deve essere sempre tenuta presente perché si tratta di un problema la cui estrema importanza diventa palese quando si pensi che, senza un modo per determinare i valori di probabilità delle ipotesi, non è possibile applicare la teoria delle probabilità.

Contrariamente a quanto potrebbe sembrare a prima vista, il dibattito sul significato da attribuire al termine « probabilità » non è circoscritto a coloro che si occupano professionalmente della teoria delle probabilità. In una certa misura tale dibattito riguarda ciascuno di noi quando, nel servirci di termini quali « probabilità », « probabile » e « probabilmente », ci chiediamo che cosa in realtà intendiamo dire;

quasi certamente in casi del genere o non siamo stati in grado di trovare una risposta soddisfacente, oppure ci siamo accorti che la risposta da noi data non era condivisa da altri.

Chi di noi ha avuto l'occasione di cimentarsi con un gioco d'azzardo come la roulette oppure con un concorso a pronostici come il totocalcio – che resta comunque un gioco d'azzardo anche se gestito dallo stato – si sarà certamente chiesto, da una parte, che cosa si intenda dire con l'affermare che il rosso è tanto probabile quanto il nero o che il banco ha maggiori probabilità di vincere rispetto al giocatore; d'altra parte, cosa si intenda quando si sostiene che probabilmente la squadra *A* vincerà la partita che giocherà la domenica seguente contro la squadra *B*.

La risposta più immediata che si può dare nel caso della roulette è che il rosso è tanto probabile quanto il nero perché vi sono 18 numeri rossi e 18

Simbolo	Significato
\neg	non
\vee	o
\wedge	e
\rightarrow	se ... allora
\leftrightarrow	se e solo se
\exists	esiste almeno un
\forall	per tutti
a_1, a_2, \dots, a_n	costanti individuali
P_1, P_2, \dots, P_m	costanti predicative

Alcuni dei principali simboli, e relativo significato, appartenenti al linguaggio simbolico di cui si parla nell'articolo. Le costanti individuali denotano individui descritti dal linguaggio mentre le costanti predicative descrivono proprietà godute dagli individui.

numeri neri e quindi i casi favorevoli al rosso sono tanti quanti quelli favorevoli al nero. Parimenti la ragione per cui il banco ha maggior probabilità di vincere di quante non ne abbia il giocatore va cercata nel fatto che quest'ultimo puntando, per esempio, sui numeri pari ha soli 18 casi favorevoli mentre il banco ne ha 19, cioè 18 dispari e lo 0.

La risposta data da uno scommettitore del totocalcio circa il significato da attribuire al termine « probabilmente » sarà sicuramente diversa. Essa infatti potrà essere che la squadra *A* vincerà probabilmente la partita della domenica seguente perché i tre migliori giocatori della squadra *B* sono stati squalificati e quindi non potranno giocare; oppure la risposta potrebbe anche essere che la squadra *B* sta attraversando un periodo di scarsa forma come è dimostrato dal fatto che da cinque domeniche non riesce a vincere una partita, ecc.

È evidente che le due spiegazioni ora viste del termine « probabilità » sono diametralmente opposte. In un caso infatti si fa appello all'analisi dei casi favorevoli all'ipotesi che interessa, mentre nell'altro si ricorre a certe informazioni più o meno « soggettive » di cui il tifoso dispone.

Le varie accezioni di « probabilità »

La disparità di opinioni che abbiamo incontrato a livello del modo comune di parlare la ritroviamo anche a livello scientifico, come vedremo subito esponendo le varie definizioni di probabilità che sono state proposte. Cominciamo questa esposizione dalla definizione classica di probabilità, definizione dovuta, come abbiamo detto, a Laplace. Il grande scienziato francese,

quando per la prima volta introdusse esplicitamente una definizione di probabilità, argomentò in modo del tutto simile a quello del nostro giocatore di roulette. Ciò non deve meravigliare dal momento che la prima formulazione scientificamente accettabile della teoria delle probabilità operata da Laplace nei primi anni del secolo scorso costituiva la razionalizzazione di una serie di risultati, relativi ai fenomeni casuali, che erano andati accumulandosi nel XVIII secolo e che per la maggior parte erano relativi ai giochi d'azzardo. Come abbiamo visto nel nostro esempio della roulette, in questi casi è del tutto naturale legare la probabilità di un evento ai casi favorevoli all'evento stesso (parlare di probabilità di un'ipotesi o, come abbiamo fatto ora, di probabilità di un evento, non comporta alcuna differenza sostanziale: si può infatti immediatamente passare dalle une alle altre prendendo in considerazione rispettivamente l'ipotesi che descrive l'evento e l'evento descritto dall'ipotesi). Gli scienziati che per primi nel XVIII secolo si interessarono a questi problemi, e Laplace che diede una sistemazione rigorosa ai loro risultati, definiscono la probabilità di un evento come il rapporto fra il numero dei casi favorevoli al verificarsi dell'evento stesso e il numero dei casi possibili purché questi ultimi siano tutti ugualmente possibili.

Questa definizione di probabilità fu unanimemente accettata e, prescindendo da qualche isolata voce discorde, fu la definizione adottata durante tutto il secolo scorso. Abbiamo però detto che qualche studioso, soprattutto nella seconda metà del XIX secolo, non era affatto d'accordo con l'opinione più diffusa. Questi dissensi divennero consistenti nei primi anni del nostro

secolo e gli argomenti contrari alla definizione laplaciana possono essere riassunti nel modo seguente.

L'applicazione di questa definizione dipende in modo essenziale dalla clausola secondo cui i casi possibili debbono essere tutti ugualmente possibili. Ma come facciamo a sapere che un dato insieme di casi possibili è un insieme di casi equipossibili? La risposta di Laplace è la seguente: « La teoria della probabilità consiste nel ridurre tutti gli eventi di un certo tipo a un certo numero di casi ugualmente possibili, *cioè tali che possiamo essere ugualmente incerti circa la loro esistenza* ». La frase in corsivo è una delle possibili formulazioni di un famoso principio enunciato da J. Bernoulli sotto il nome di principio di ragione non sufficiente e ora noto invece come principio d'indifferenza. Per la maggior parte delle situazioni induttive del tipo considerato per lo più da Laplace, risulta non solo abbastanza facile, ma anche abbastanza naturale analizzare la situazione in un insieme di casi equipossibili, tali cioè che non sussistano ragioni di credere che se ne verificherà uno di essi piuttosto che un altro (il lettore si sarà accorto che questa è un'altra formulazione del principio d'indifferenza). In effetti nei casi dei giochi d'azzardo che fanno uso di meccanismi non truccati (dadi, monete, roulette, urne) è molto plausibile l'ipotesi di equipossibilità dei casi possibili. Ma le difficoltà si presentano non appena abbiamo qualche ragione di sospettare che i dadi, le monete, le roulette sono truccati. Per esempio nel caso di un dado zavorrato è ovvio che le uscite di ciascuna delle sei facce non potranno più essere considerate come equipossibili. A un'analisi più sottile la definizione di Laplace si rivela già inapplicabile, almeno in senso generale, agli stessi casi per i quali era stata formulata.

Al di fuori di questi casi l'inapplicabilità di tale definizione è addirittura palese: che cosa significa infatti che il signor Rossi che ha compiuto 38 anni ha una probabilità pari a 0,003 di morire nel corso del 39° anno di vita? Significa forse che il nostro signore ha di fronte 1000 modi alternativi di vivere di cui solo 3 lo conducono a morte nel corso del prossimo anno? L'assurdità di una simile ipotesi è così evidente da non meritare alcun commento.

Le critiche che abbiamo ricordato sono dovute a R. von Mises che le formulò attorno al 1920 e segnarono l'inizio del periodo ormai noto come « crisi dei fondamenti della probabilità », cioè di quel periodo di dibattiti volti a

individuare una nuova spiegazione del significato del termine « probabilità » da sostituire a quella data da Laplace all'inizio del XIX secolo.

Von Mises fece seguire alle critiche viste una nuova definizione di probabilità; serviamoci di un esempio per meglio afferrare la innovazione contenuta nella definizione da lui proposta. Supponiamo di lanciare un certo numero di volte un dado e di registrare i risultati dei nostri lanci nel seguente prospetto:

1°	2°	3°	4°	5°	6°	7°
6	4	2	1	6	4	2
8°	9°	10°	11°	12°	13°	...
1	3	2	3	1	1	...

Chiamiamo ora frequenza relativa dell'evento « uscita di 1 » il numero di volte in cui compare l'1 in un certo numero di lanci diviso il numero di questi ultimi. Nel nostro esempio la frequenza relativa di questo evento sarà allora 4/13 mentre la frequenza relativa dell'evento « uscita di 6 » sarà 2/13. Ciascun evento avrà una sua frequenza relativa che varierà col variare del numero di lanci che prendiamo in considerazione. La frequenza relativa dell'evento « uscita di 1 », quando ci limitiamo a considerare i primi sei lanci, sarà 1/6.

Von Mises sostenne che per probabilità del verificarsi di un certo evento in una successione casuale di eventi si-

mili e ripetibili, quali i lanci di un dado del nostro esempio, doveva intendersi il limite al quale tende la frequenza relativa di quell'evento quando il numero dei termini della successione cresce indefinitamente, ovviamente quando questo limite esiste. È bene precisare che il termine « limite » usato in questo contesto non è il preciso concetto dell'analisi matematica bensì il concetto impreciso che ciascuno di noi usa quando nel linguaggio di tutti i giorni si serve di questo termine.

All'incirca nello stesso periodo fu però proposta, a opera dell'economista J. M. Keynes, una terza definizione di probabilità. Per Keynes la probabilità di una ipotesi è una relazione logica che lega un insieme di informazioni disponibili per un certo individuo o per una certa collettività e l'ipotesi stessa. Questo autore non precisava tuttavia che cosa dovesse intendersi per questa relazione logica così che l'intera definizione restava alquanto oscura. Si può dare un'idea di ciò che Keynes intendeva facendo intervenire una relazione logica per spiegare il significato della probabilità, tornando all'esempio del tifoso che riteneva probabile la vittoria della squadra A nella partita della domenica seguente. Il nostro tifoso riteneva infatti che esistesse qualche relazione oggettiva, e in questo senso logica, tra la mancanza dei tre migliori giocatori di B e il rendimento di questa squadra nella prossima partita, oppure che lo stato di scarsa forma dimostrato dalla squadra B nelle ultime cinque partite fosse in qualche modo oggettivamente connesso con il com-

portamento di questa squadra nella prossima partita con la squadra A. In breve, il nostro tifoso riteneva che la vittoria della squadra A fosse probabile perché un insieme di informazioni in suo possesso sostenevano l'ipotesi della vittoria di A più di quanto sostenessero l'ipotesi contraria.

Alcuni anni più tardi, attorno al 1930, due autori, l'inglese F. P. Ramsey e l'italiano B. de Finetti, pur partendo da considerazioni del tutto differenti, arrivarono a una medesima conclusione relativamente al significato da attribuire al termine « probabilità », conclusione del tutto diversa da quella alla quale erano arrivati i due autori che abbiamo già preso in esame. Ramsey infatti, dopo aver tentato senza successo di dare un significato alla relazione logica che compare nella definizione di Keynes, si convinse che l'unica possibilità di legare fra di loro le informazioni in possesso di un individuo e una ipotesi da questi avanzata, consisteva nel fare appello alla fiducia che l'individuo nutriva nella verità della ipotesi in questione. De Finetti invece, muovendo dalla convinzione che ogni evento è un fatto unico non ripetibile e che, conseguentemente, non ha senso ipotizzare la ripetibilità degli eventi come presuppone la definizione di von Mises, arrivò come Ramsey a sostenere che la probabilità di un evento è la fiducia che un dato individuo nutre nel verificarsi dell'evento stesso.

Una conseguenza immediata di questi contrastanti punti di vista era che ciascuno di essi proponeva un modo,

1. Assioma dell'equivalenza
Se $\vdash H \leftrightarrow H'$ e $E \leftrightarrow E'$, allora $c(H, E) = c(H', E')$.
2. Assiomi dell'invarianza
PER LE COSTANTI INDIVIDUALI: Se $q(H)$ e $q(E)$ sono ottenuti mediante q da H ed E , allora $c(H, E) = c(q(H), q(E))$.
PER LE COSTANTI PREDICATIVE: Se $r(H)$ e $r(E)$ sono ottenuti mediante r da H ed E , allora $c(H, E) = c(r(H), r(E))$.
PER LE FAMIGLIE: Se $t(H)$ e $t(E)$ sono ottenuti mediante t da H ed E , allora $c(H, E) = c(t(H), t(E))$.
3. Assioma della rilevanza degli esempi positivi
Se a_i e a_j sono costanti individuali che non compaiono in E , allora $c(P_{a_i}, E \wedge P_{a_j}) > c(P_{a_i}, E)$.

Rappresentazione simbolica degli assiomi della logica induttiva di cui si parla nell'articolo. $c(H, E)$ significa la probabilità di H sulla base di E ; nei casi interessanti l'enunciato H descrive un'ipotesi relativa a un fatto sconosciuto mentre l'enunciato E descrive un insieme di risultati sperimentali conosciuti da chi for-

mula l'ipotesi. Il simbolo \vdash serve per indicare che quanto lo segue è dimostrabile, e quindi $\vdash H \leftrightarrow H'$ significa che H e H' sono logicamente equivalenti. q, r e t sono permutazioni rispettivamente delle costanti individuali, delle costanti predicative e delle famiglie con lo stesso numero di membri del linguaggio simbolico.

diverso da quello proposto dagli altri, di determinare le probabilità. Per i frequentisti, cioè per coloro che accettano la definizione di von Mises, la probabilità deve essere determinata sulla base di osservazioni sperimentali, cioè esaminando le frequenze relative con cui gli eventi si verificano. Per i logici, cioè per coloro che accettano la tesi di Keynes, la probabilità deve essere determinata cercando di dare un valore numerico alla relazione logica che lega le conoscenze di cui si dispone e l'ipotesi che descrive il verificarsi dell'evento. Per i soggettivisti infine, cioè per coloro che fanno proprie le argomentazioni di Ramsey e de Finetti, le probabilità devono essere determinate rendendo esplicita mediante l'uso di scommesse la fiducia degli individui nel verificarsi dell'evento. È chiaro che, mentre per i frequentisti e i logici vi è una sola probabilità per ogni evento, nel senso che la probabilità è un fatto esterno all'uomo che questo deve scoprire, per i soggettivisti invece ogni uomo può assegnare a un dato evento una sua propria probabilità e il fatto che due individui possano avere la stessa fiducia nel verificarsi di un evento e quindi assegnargli la stessa probabilità, è un fatto del tutto accidentale.

L'assiomatizzazione di Kolmogorov

All'inizio degli anni trenta le cose stavano quindi come abbiamo visto: esistevano cioè tre modi contrapposti di definire la probabilità, e di conseguenza tre modi diversi di determinare la probabilità; ma ancora, ciò che è senza dubbio più grave, si credeva che esistessero tre teorie diverse delle probabilità poiché ciascuno degli indirizzi che abbiamo passato in rassegna sosteneva che la teoria delle probabilità era quella sviluppata a partire dalla definizione che esso sosteneva, cioè era rispettivamente la teoria delle successioni casuali di eventi, delle relazioni logiche tra evidenze e ipotesi e della fiducia degli individui nel verificarsi degli eventi.

La crisi dei fondamenti della probabilità si rifletteva quindi sulla stessa teoria delle probabilità con evidente grave pregiudizio di quest'ultima. Questo stato di incertezza fu però superato nel 1933 quando il matematico sovietico A. N. Kolmogorov dimostrò che la teoria delle probabilità poteva essere completamente sviluppata senza mai farvi intervenire il significato del termine « probabilità ». In altre parole, il matematico sovietico mostrò che la teoria delle probabilità – come del resto tutta la matematica di cui questa

teoria è una branca – si occupa di entità astratte che non necessitano, almeno a livello dello sviluppo della teoria, di alcuna interpretazione. Kolmogorov si poneva cioè sulle orme di D. Hilbert, e, come quest'ultimo aveva mostrato che l'intera geometria poteva essere sviluppata senza fare riferimento al significato dei termini « punto », « retta », « piano », ecc., così Kolmogorov mostrò che l'intera teoria delle probabilità può essere sviluppata senza fare riferimento al significato del termine « probabilità », limitandosi a introdurre assiomaticamente alcune sue caratteristiche, per altro molto generali, note appunto come assiomi della teoria delle probabilità: a) la probabilità è una funzione i cui argomenti sono eventi (ipotesi) e i cui valori sono numeri reali maggiori o uguali a zero e minori o uguali a uno; b) la probabilità di due eventi (ipotesi) che non possono verificarsi contemporaneamente (che sono incompatibili) è la somma delle probabilità dei due eventi (ipotesi).

L'opera di Kolmogorov ebbe una importanza notevolissima innanzitutto perché, liberando la teoria delle probabilità dalle dispute relative ai fondamenti, le permise uno sviluppo indipendente e rigoglioso tanto che oggi questa disciplina rappresenta una delle branche della matematica più attive e feconde di risultati. Tuttavia l'influenza del fondamentale lavoro del matematico sovietico non restò limitata a ciò: la generalizzazione dell'istanza di fondo che caratterizza il lavoro di Kolmogorov avrebbe infatti portato a una svolta radicale nella crisi dei fondamenti della probabilità, cioè alla nascita della logica induttiva moderna.

Abbiamo detto che la tesi di fondo su cui si basa il lavoro del matematico sovietico è la seguente: la teoria delle probabilità può e deve essere sviluppata senza fare riferimento al significato del termine « probabilità ». Vediamo ora cosa significhi generalizzare questa tesi. È chiaro che, prescindendo dalle intenzioni sostanzialmente metafisiche volte a individuare l'essenza della probabilità, intenzioni che ora non ci interessa analizzare, il significato delle ricerche connesse con la crisi dei fondamenti della probabilità risiede nella convinzione che solo quando si è attribuito un significato al termine « probabilità » si possano determinare i valori di probabilità e quindi in ultima istanza si possa applicare la teoria delle probabilità. È ora naturale, accettando la istanza formalistica di Kolmogorov, porsi la seguente domanda: se è possibile sviluppare la teoria delle probabilità senza fare appello al significato

del termine « probabilità », non sarà possibile trovare il modo di determinare anche i valori di probabilità senza fare riferimento al significato di questo termine? Se si generalizza l'impostazione formalistica di Kolmogorov, non si può non dare risposta affermativa a questa domanda; ciò significa tentare di elaborare dei metodi atti alla determinazione delle probabilità nei quali non si faccia alcun riferimento a un qualunque significato di « probabilità ». La domanda successiva sarà allora come ciò sia possibile. La nostra risposta è che le probabilità vengono di fatto determinate – in generale ma non necessariamente – mediante i cosiddetti metodi statistici. Esempi di questa determinazione si hanno nelle ricerche fisiche, astronomiche, economiche, demografiche, genetiche, farmacologiche, meteorologiche e in molte altre discipline, in sostanza in tutte le discipline sperimentali.

La logica induttiva

Le ricerche relative ai fondamenti della probabilità acquistano quindi un significato nuovo e completamente diverso da quello che hanno avuto nella prima metà del secolo: non si tratterà più di andare alla ricerca di esplicazioni del termine « probabilità » bensì dei principi che di fatto regolano la determinazione delle probabilità nelle scienze sperimentali. Una volta che questi principi siano stati individuati, dovranno essere formalizzati e, sulla base di questa formalizzazione, se ne dovrà studiare la portata, la compatibilità reciproca, l'indipendenza reciproca e così di seguito; se ne dovrà cioè analizzare la struttura logica e le implicazioni epistemologiche.

Alcuni di questi principi sono già stati individuati e dopo essere stati formalizzati sono stati assunti quali assiomi della logica induttiva; ora ne esamineremo alcuni fra i più importanti. La logica induttiva costituisce un ampliamento della teoria delle probabilità; come abbiamo detto, quest'ultima può a pieno diritto essere considerata una parte della logica induttiva. Gli assiomi di questa disciplina comprenderanno quindi anche quelli della teoria delle probabilità; ciò significa che una volta determinati i valori di probabilità delle ipotesi induttive, essi possono essere « maneggiati » con i teoremi della teoria delle probabilità.

Abbiamo visto che Kolmogorov basò il suo sistema d'assiomi sul presupposto che la probabilità fosse una funzione: ciò ci assicura l'esistenza di uno e un solo valore di probabilità in corri-

spondenza di ogni ipotesi. Sia quindi p la funzione-probabilità; essa avrà come argomenti enunciati (ipotesi) di un certo linguaggio simbolico e come valori numeri reali non negativi. Se H è un enunciato, indicheremo con $p(H)$ (a volte anche con p_H) il valore che la funzione-probabilità assume in corrispondenza di H (la probabilità di H).

Siano H e H' enunciati, imponiamo allora che la probabilità soddisfi le seguenti condizioni:

assioma dell'equivalenza: se H è logicamente equivalente a H' , allora $p(H) = p(H')$;

assioma del limite superiore: se H è logicamente certo, allora $p(H) = 1$;

assioma dell'additività: se H e H' sono logicamente incompatibili, allora $p(H \vee H') = p(H) + p(H')$.

Ci limiteremo a illustrare l'assioma dell'equivalenza poiché, come il lettore avrà certamente notato, degli altri due ci siamo già occupati esaminando l'assiomatizzazione di Kolmogorov. Il lettore si sarà inoltre accorto che tra gli assiomi allora visti non compariva quello dell'equivalenza; ciò dipende dal fatto che Kolmogorov costruiva il suo sistema assiomatico sulla nozione di « evento » e non su quella di « enunciato » come noi stiamo facendo, e questo ci impone di prendere delle precauzioni in vista di evitare spiacevoli conseguenze. L'assioma dell'equivalenza afferma infatti che ipotesi logicamente equivalenti non possono avere probabilità diverse; in questo modo evitiamo l'eventualità senza dubbio spiacevole di assegnare probabilità diverse a ipotesi con lo stesso contenuto informativo, a ipotesi cioè che « dicono la stessa cosa » seppure in modi diversi.

Dopo aver esaminato gli assiomi della teoria delle probabilità passiamo a un assioma che sta per così dire a mezza strada fra la teoria delle probabilità e la logica induttiva. Si tratta cioè di un assioma la cui validità alcuni ritengono debba essere confinata alla logica induttiva mentre altri vorrebbero estendere anche alla teoria delle probabilità. Esso fu introdotto dai logici induttivi ma le ricerche di cui è stato oggetto hanno consentito di metterne in luce aspetti indubbiamente probabilistici. Ma passiamo all'enunciazione dell'*assioma di regolarità*: se $p(H) = 1$, allora H è logicamente certo. Questo assioma è indubbiamente legato a quello del limite superiore: questo ci assicura che la probabilità delle ipotesi certe è pari a uno, quello di regolarità che vale anche l'inverso. La loro congiunzione porta a questa conclusione: le ipotesi certe sono tutte e sole quelle la cui probabilità è pari a uno.

Ma passiamo alla logica induttiva in senso proprio concentrando la nostra attenzione su alcuni fra i più caratteristici assiomi di questa disciplina. Intendiamo riferirci agli assiomi dell'invarianza, che per certi aspetti si ricollegano al principio d'indifferenza di Bernoulli, e all'assioma di rilevanza degli esempi positivi.

Gli assiomi dell'invarianza sono largamente accettati dalla metodologia statistica e possono trovare una matrice comune nel principio secondo cui, in alcune circostanze, quello che ha importanza nella valutazione delle probabilità è la quantità delle osservazioni e non anche la loro qualità. Ma cerchiamo di chiarire questa affermazione con alcuni semplici esempi. Come è d'uso in questi casi, supponiamo di voler determinare la probabilità di estrarre una pallina di un dato colore da un'urna che sappiamo contenerne di bianche e di nere in proporzione sconosciuta. Supponiamo inoltre di operare la determinazione dopo aver fatto un certo numero di estrazioni, cioè di valutare la probabilità di un'estrazione futura sulla scorta della conoscenza dei risultati di alcune estrazioni passate. Ovviamente l'esempio di cui ci serviremo può essere facilmente generalizzato sia al caso in cui i colori delle palline contenute nell'urna siano più di due sia al caso in cui si tratti di esaminare individui provenienti da una popolazione qualsiasi.

Supponiamo quindi di aver estratto n palline dall'urna e di aver constatato che n_B sono bianche e le restanti n_N sono nere. Il nostro problema sarà allora di determinare la probabilità di estrarre al prossimo colpo una pallina bianca o, ciò che sostanzialmente è lo stesso, di determinare la probabilità dell'ipotesi secondo cui la prossima pallina che estrarremo sarà bianca. Supponiamo di averlo risolto in qualche modo cioè usando un metodo che però ora non ci interessa analizzare. In prima istanza l'assioma dell'invarianza che per primo prenderemo in esame afferma che questa probabilità dipende unicamente da n_B e n_N e non anche dall'ordine con cui le palline bianche e nere si sono presentate nel corso delle n estrazioni già eseguite. Se cioè $n_B = 2$ e $n_N = 1$, le tre estrazioni che possono aver dato luogo a questo risultato sono ovviamente le seguenti:

a) $Ba_1Ba_2Na_3$

b) $Ba_1Na_2Ba_3$

c) $Na_1Ba_2Ba_3$

in cui B e N stanno per bianca e nera e a_1 , a_2 e a_3 indicano rispettivamente il risultato della (la pallina estratta alla) prima, seconda e terza estrazione.

Indichiamo con p_{Ba} , p_{Bb} e p_{Bc} le probabilità di estrarre al prossimo colpo una pallina bianca calcolata supponendo noti rispettivamente i risultati a), b) e c). L'assioma in questione stabilisce che

$$p_{Ba} = p_{Bb} = p_{Bc};$$

questa probabilità la indicheremo con p_B .

Se ora formalizziamo il nostro discorso traducendo a), b) e c) in un opportuno linguaggio simbolico, B e N diventeranno segni per proprietà (costanti predicative) e a_1 , a_2 e a_3 segni per individui (costanti individuali). In questo caso l'assioma di cui ci stiamo occupando diventa l'assioma dell'*invarianza relativamente alle costanti individuali*: la probabilità di un enunciato è invariante relativamente alla permutazione delle costanti individuali del linguaggio. Infatti i tre risultati sperimentali visti e anche gli enunciati che nel linguaggio simbolico li descrivono, si ottengono l'uno dall'altro permutando opportunamente gli individui (le costanti individuali). Una diversa formulazione di questo assioma, che può essere più convincente a livello intuitivo, è la seguente: nel valutare le probabilità tutti gli individui (le costanti individuali) devono essere trattati alla pari.

Prima di passare al secondo e al terzo assioma dell'invarianza, introduciamo il concetto di *famiglia* di proprietà (costanti predicative); per famiglia intenderemo un insieme di proprietà (costanti predicative) a due a due incompatibili e tali che ogni individuo (costante individuale) deve godere di almeno una di esse. Le proprietà coinvolte nelle estrazioni che abbiamo finora preso in considerazione danno origine a una famiglia; infatti una pallina non può essere bianca e nera e ogni pallina che estraiamo gode di una delle due proprietà, questo ovviamente per il modo con cui abbiamo costruito il dispositivo sperimentale. Anche le proprietà « essere testa » ed « essere croce » riferite ai possibili risultati del lancio di una moneta sono una famiglia. Infatti lanciando una moneta non possono presentarsi contemporaneamente due facce e d'altra parte una delle due facce deve presentarsi, ovviamente quando si escluda che la moneta possa fermarsi in bilico sullo spigolo. Ci si convince facilmente che anche le date di nascita degli individui e il luogo in cui si trovano in un certo istante sono famiglie così come le loro stature, i loro perimetri toracici e così di seguito.

Supponiamo ora d'aver fatto una nuova serie di estrazioni tali che $n_B = 1$ e $n_N = 2$. Come sappiamo le tre estra-

zioni che possono aver dato luogo a questo risultato sono le seguenti:

$$d) Na_1Na_2Ba_3$$

$$e) Na_1Ba_2Na_3$$

$$f) Ba_1Na_2Na_3$$

in cui i simboli hanno il significato che conosciamo. Supponiamo ora di essere in grado di determinare le probabilità che la pallina estratta al prossimo colpo sia nera sulla base della conoscenza di d), e) e f) e inoltre che esse siano

$$p_{Nd} = p_{Ne} = p_{Nf}$$

Quest'ultima uguaglianza vale ovviamente in virtù dell'assioma dell'invarianza relativamente alle costanti individuali; indichiamo questa probabilità con p_N . Il secondo assioma dell'invarianza stabilisce che

$$p_N = p_B$$

cioè che la probabilità di estrarre una pallina nera, quando si sappia che sono state estratte due palline nere e una bianca, deve essere uguale alla probabilità di estrarre una pallina bianca, quando si sappia che sono state estratte due palline bianche ed una nera. Generalizzato e tradotto a livello simbolico esso diventa l'assioma dell'invarianza relativamente alle costanti predicative: la probabilità di un enunciato è invariante relativamente alla permutazione delle costanti predicative di una famiglia. Anzi, in questo caso vi è un'altra formulazione: nel valutare le probabilità tutte le proprietà (costanti predicative) di una famiglia devono essere trattate alla pari.

Prendiamo ora in considerazione, oltre a quello delle urne, un nuovo esperimento, supponiamo cioè di aver lanciato tre volte una moneta senza conoscere se la stessa è truccata oppure no, e di aver osservato che due volte si è presentata la faccia con la croce e una volta quella con la testa, cioè $n_C = 2$ e $n_T = 1$. Supponiamo inoltre di essere in grado di valutare la probabilità che nel prossimo lancio si presenti la faccia con la croce e che essa sia pari a p_C . Se applichiamo il principio che finora ci ha guidato dovremo concludere che, pur essendo il lancio di una moneta e l'estrazione dall'urna esperimenti che coinvolgono proprietà diverse, debba valere

$$p_C = p_B$$

in virtù del fatto che sia nell'uno sia nell'altro caso i risultati favorevoli all'ipotesi erano due e quello contrario uno. L'assioma che ha diretto le nostre considerazioni è l'assioma dell'invarianza relativamente alle famiglie: la probabilità di un enunciato è invariante relativamente alla permutazione delle famiglie del linguaggio con lo stesso numero di membri.

L'ultimo assioma della logica indut-

tiva che intendiamo prendere in esame, traduce a livello formale un principio che non solo è largamente usato in tutte le elaborazioni statistiche dei dati sperimentali, ma è anche una delle tesi fondamentali dell'empirismo: si tratta del principio secondo cui *dobbiamo* imparare dall'esperienza.

Si supponga di aver compiuto n estrazioni da un'urna di cui al solito non conosciamo la composizione. Delle n estrazioni effettuate n_B sono state palline bianche. Supponiamo ancora di disporre di un metodo che ci consenta la determinazione delle probabilità e che grazie a questo si sia calcolata pari a p_B la probabilità di ottenere una pallina bianca nella prossima estrazione. Questa probabilità tiene ovviamente conto delle estrazioni effettuate. Supponiamo ora che la situazione sperimentale sia leggermente mutata nel senso che delle n estrazioni fatte non più n_B siano palline bianche bensì $n_B + 1$; ciò significa che nella nuova situazione sperimentale avremo estratto una pallina bianca in più rispetto alla situazione presa precedentemente in esame. Sulla base di questo mutato risultato sperimentale calcoliamo, servendoci del metodo di cui disponiamo, la probabilità di estrarre una pallina bianca e sia $'p_B$ questa probabilità. L'assioma della rilevanza degli esempi positivi impone che qualunque sia il metodo usato per determinare p_B e $'p_B$ debba valere

$$p_B < 'p_B;$$

cioè: la probabilità di un'ipotesi deve crescere col crescere del numero degli individui osservati che verificano l'ipotesi stessa.

Con questo assioma abbiamo concluso il nostro esame della logica induttiva; ricordiamo però ancora una volta che gli assiomi di cui ci siamo occupati sono solo una parte di quelli che stanno alla base dei sistemi finora noti di logica induttiva. Ci siamo infatti limitati ad alcuni assiomi che ci sono sembrati particolarmente adatti a corroborare le nostre affermazioni sulla natura e gli scopi della logica induttiva moderna. Siamo inoltre convinti che questo esame sia servito a dare la misura di quanto queste moderne ricerche si riallaccino a quelle che, secondo Leibniz, avrebbero dovuto costituire l'oggetto della logica induttiva.

Indipendentemente comunque dal ricolligarsi o meno alle intuizioni leibniziane, la disciplina di cui abbiamo ora visto alcuni assiomi merita bene il nome di logica induttiva che le abbiamo dato. Infatti, mentre la logica deduttiva si occupa della struttura logica delle inferenze dimostrative, la logica indut-

tiva si occupa con gli stessi intenti delle inferenze non dimostrative, cioè della esplicitazione e dell'analisi delle modalità d'assegnazione delle probabilità alle ipotesi induttive. La logica induttiva moderna può quindi a pieno diritto essere inclusa nell'ambito della logica formale.

L'opera di Carnap

Abbiamo detto all'inizio del presente articolo che è solo nei primi anni del nostro secolo che la logica induttiva acquista il suo significato moderno e abbandona definitivamente ogni pretesa di porsi come la logica della scoperta scientifica. Questa affermazione, sulla base di quanto abbiamo visto, deve però essere ridimensionata nel senso che è all'inizio del nostro secolo che matura la crisi dei fondamenti della probabilità, crisi che sfocerà nella nascita della moderna logica induttiva. In senso lato, quindi, questa disciplina può farsi risalire all'inizio del nostro secolo solo se si intende – e questa è la nostra opinione – come la naturale conseguenza e il superamento della crisi dei fondamenti della probabilità. In senso stretto, cioè come analisi razionale delle modalità di assegnazione delle probabilità alle ipotesi induttive, essa è il risultato di ricerche portate avanti negli ultimi venti anni principalmente per opera del filosofo della scienza tedesco R. Carnap.

La logica induttiva quindi è una disciplina molto giovane ed è naturale che stia muovendo i suoi primi passi fra difficoltà d'ogni genere. Ciò nonostante ha già al suo attivo alcuni risultati notevoli, per esempio la dimostrazione della non indipendenza dell'assioma di rilevanza degli esempi positivi e la dimostrazione di un certo tipo di incompatibilità tra l'assioma di invarianza relativamente alle famiglie e gli altri assiomi. In questa sede non possiamo occuparci della non indipendenza dell'assioma della rilevanza a ragione delle difficoltà tecniche che comporta la sua dimostrazione. Il secondo problema può invece essere affrontato, almeno fino a un certo punto, senza fare ricorso ad alcun tecnicismo. Abbiamo parlato di « un certo tipo » di incompatibilità perché non si è dimostrato che da questo assioma, insieme con gli altri della logica induttiva, si può derivare una contraddizione, ma piuttosto che la sua accettazione conduce a un risultato intuitivamente paradossale che può essere evitato se si lascia cadere l'assioma o se ne limita la portata.

Il paradosso della conferma, questo

è il nome con cui è noto questo risultato paradossale, può essere formulato nel modo seguente. È opinione comune che se uno scienziato avanzasse la ipotesi « Tutte le pietre di Marte sono sferiche », ipotesi che indicheremo con H , egli dovrebbe, in vista di aumentare il grado di probabilità dell'ipotesi stessa (per confermarla), osservare le pietre di Marte e costatare che sono sferiche. Sarà cioè necessario che il nostro scienziato si ponga in grado di osservare le pietre di quel pianeta allo scopo di vedere se sono effettivamente sferiche. In altri termini, in accordo con l'assioma della rilevanza, per confermare l'ipotesi H sarà necessario osservare oggetti che siano contemporaneamente pietre di Marte e sferici. Ora se ci rammentiamo l'assioma dell'equivalenza giungiamo al paradosso.

L'ipotesi H è logicamente equivalente alla seguente « Tutti gli oggetti non sferici non sono pietre di Marte » che indicheremo con H' . Questa ipotesi è confermata dall'osservazione di oggetti che non sono né sferici né pietre di Marte, per esempio, un autocarro di Pavia o un tram di Milano. Ma H è logicamente equivalente a H' e quindi, per l'assioma dell'equivalenza, la probabilità di H deve essere uguale a quella di H' ; ne consegue che crescendo la probabilità di H' deve crescere anche quella di H , in definitiva che l'osservazione di autocarri di Pavia o tram di Milano conferma l'ipotesi che tutte le pietre di Marte sono sferiche. Quindi, per aumentare il grado di probabilità dell'ipotesi relativa alla forma delle pietre di Marte, lo scienziato che l'ha formulata dovrebbe mettersi alla ricerca di oggetti non sferici per vedere se sono o no pietre di Marte. Per confermare la sua ipotesi il nostro scienziato non dovrebbe recarsi su Marte o nelle sue vicinanze perché gli basterebbe controllare che nella sua città esistono un gran numero di oggetti che non sono sferici. Una conferma di questo tipo è indubbiamente inaccettabile, d'altra parte è una conseguenza degli assiomi della logica induttiva: da qui il paradosso.

A prima vista il paradosso della conferma sembra destinato a infrangere ogni sforzo volto a costruire la logica induttiva intesa come analisi razionale delle inferenze non dimostrative. Tuttavia lo studio approfondito che ne è stato fatto ha reso un grande servizio a questa disciplina poiché ha reso i suoi cultori molto più consapevoli di quanto vanno facendo aumentando pertanto l'efficacia delle loro ricerche. E vediamo le ragioni.

Appena si scopre un paradosso subito si cerca la sua « soluzione » e immaginiamo che anche il lettore sarà curioso di sapere quale « soluzione » è stata data al paradosso della conferma, forse per confrontarla con quella che ha trovato per conto suo. Come insegna però la storia della logica, ogni paradosso trova sempre più di una soluzione e in una certa misura non ne trova mai perché ogni « soluzione » affronta il paradosso secondo un ben determinato punto di vista. Nel nostro caso purtroppo non possiamo presentare le « soluzioni » finora suggerite, in primo luogo, perché gli studi a esse legati sono ben lungi dall'essere conclusi; in secondo luogo, perché le parziali soluzioni finora trovate presuppongono un apparato tecnico che è fuori di luogo ora presentare. Una direzione di ricerca sembra però un fatto ormai acquisito e può essere esposta senza fare ricorso ad alcun tecnicismo. Si è cioè scoperto che il paradosso della conferma trova origine nell'aver messo sullo stesso piano le due famiglie di cui fanno parte rispettivamente le proprietà « essere una pietra di Marte » e « essere sferico ». Il paradosso potrà quindi risolversi se questa parità sarà lasciata cadere; ciò significa che se si vuole evitare la situazione paradossale di cui si è detto, si deve ammettere che le due famiglie « pesino » in modo diverso sulla determinazione della probabilità dell'ipotesi delle pietre di Marte. La incompletezza delle soluzioni finora trovate dipende dal fatto che non si è ancora in grado di dare un significato ai pesi che si assegnano alle famiglie.

Conclusione

Non ci sembra il caso di continuare in questa direzione, esponendo cioè le proposte di interpretazione dei « pesi » suddetti, quanto piuttosto di riflettere brevemente sulla influenza del paradosso sulla logica induttiva. Per far ciò è bene ricordare la situazione esistente prima della scoperta del paradosso. Non si tratta beninteso di un primo cronologico poiché il paradosso è una scoperta anteriore ai primi tentativi di costruire un sistema assiomatico per la logica induttiva, bensì di un primo metodologico. Nei primi tentativi che si sono compiuti per analizzare i metodi di assegnazione delle probabilità alle ipotesi induttive, era sembrato naturale assumere che tutte le famiglie dovessero essere trattate alla pari. Questa convinzione era suffragata dal fatto che la suddetta parità viene largamente usata nella valutazione statistica dei risultati sperimentali, come abbiamo mo-

strato quando abbiamo illustrato gli assiomi della logica induttiva e come mostra una breve riflessione sui metodi statistici.

La parità delle famiglie nei confronti della valutazione delle probabilità era quindi stata stabilita assiomaticamente. Convinti che la sua validità fosse universale, la si era applicata anche al caso delle probabilità delle leggi di natura (tale è infatti l'ipotesi delle pietre di Marte). La scoperta del paradosso segnò il fallimento del tentativo di affermare la validità universale dell'assioma che stabiliva la parità delle famiglie.

La conclusione che i logici induttivi hanno tratto da questa sconfitta è, come abbiamo detto, un arricchimento concettuale della loro disciplina. Ci si è cioè accorti che l'assioma dell'invarianza in questione può essere accettato solo in casi particolari, per esempio a livello delle situazioni sperimentali affrontate dalla statistica, mentre in altri casi, per esempio per le leggi di natura, deve essere sostituito con assiomi più articolati che, stabilendo una gerarchia fra le famiglie, evitino di appiattirle tutte allo stesso livello.

È evidente che questa duttilità nell'uso degli assiomi è possibile solo se si pensa alla logica induttiva in modo non dogmatico, cioè solo se i suoi principi di base non sono intesi come verità incrollabili, bensì come strumenti che debbono essere accettati solo fin quando servono e che comunque devono confrontarsi dialetticamente con l'avanzare della ricerca. Ed è questa infatti la posizione dei moderni logici induttivi; essi sono in altri termini disposti ad accettare l'assioma dell'invarianza relativamente alle famiglie così come gli altri assiomi fin dove essi possono servire e sono invece pronti ad abbandonarli non appena si rivelassero dannosi per il progresso delle loro ricerche.

Abbiamo visto come il paradosso della conferma abbia aumentato la consapevolezza degli studiosi di logica induttiva e come lungi dal segnare la fine di questa disciplina le abbia permesso di scoprire nuove insospettite caratteristiche delle modalità d'assegnazione delle probabilità alle ipotesi induttive. La breve storia che abbiamo fatto di questo paradosso ci ha inoltre fornito l'opportunità di porre in risalto la posizione di fondo su cui si reggono le ricerche di logica induttiva, posizione che si ispira alla grande tradizione della moderna epistemologia di cui accetta le istanze filosofiche che sono poi le istanze su cui si fonda tutta la ricerca scientifica moderna.

I problemi della conferma

Per secoli i ricercatori hanno posto a controllo, hanno confermato e infirmato ipotesi mediante l'osservazione e l'esperimento. Eppure la logica di questo procedimento è ancora lungi dall'essere compresa

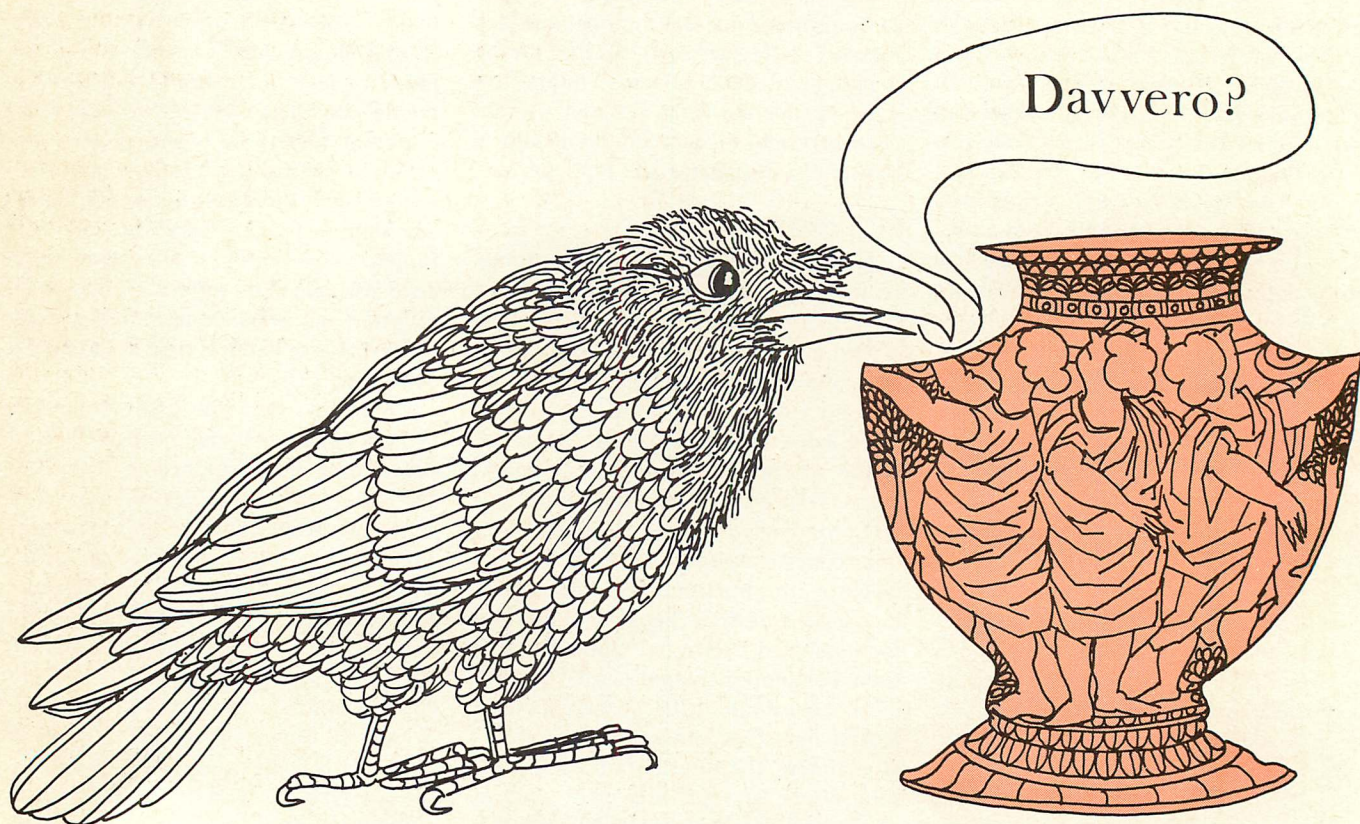
di Wesley C. Salmon

È opinione comune che la conferma di un'ipotesi scientifica consista nell'accertare che essa ha conseguenze vere. La teoria della gravitazione di Newton, per esempio, implicava che due corpi materiali qualsiasi esercitassero un'attrazione reciproca. La teoria fu confermata dall'osservazione del comportamento dei corpi in caduta, dei movimenti dei pianeti e del flusso e riflusso delle maree; essa ricevette conferma sperimentale diretta quando Henry Cavendish escogitò una bilancia di torsione per mezzo della quale diveniva possibile dimostrare

in laboratorio l'attrazione tra due corpi. Tali scoperte non verificano una generalizzazione in modo conclusivo. Malgrado le molte diverse conferme sperimentali della teoria di Newton, essa non è più considerata completamente accettabile. Si è presentata dell'evidenza infirmante, e ora la fisica newtoniana è stata soppiantata dalla relatività einsteiniana. Indipendentemente dalla sicurezza con cui sembra sia stabilita una teoria, non si può avere la certezza che essa non vada incontro a un destino analogo. Il massimo che si può dire è che i risultati

sperimentali tendono a confermare una teoria e che in alcuni casi l'accumularsi dell'evidenza confermantе innalza un'ipotesi generale allo stato di una accettabilità almeno provvisoria.

La precedente caratterizzazione del controllo scientifico delle ipotesi può sembrare diretta e priva di problemi, e tuttavia è una ipersemplificazione. Nei dettagli mancanti si annida una moltitudine di difficoltà fondamentali. Uno dei modi migliori per esporre alcuni dei problemi nascosti è quello di considerare una serie di semplici esempi, ognuno in possesso di qualche ca-



Il paradosso dei corvi è un esempio delle sottigliezze della teoria della conferma. Se tutti i corvi sono neri, sicuramente le cose non-neri devono essere non-corvi. Poiché le due generaliz-

zazioni sono logicamente equivalenti, ogni evidenza che conferma l'una conferma anche l'altra. Quindi un vaso colorato sembra fornire conferma all'ipotesi che tutti i corvi sono neri.

ratteristica paradossale o controintuitiva. Bertrand Russell una volta ha osservato: « Una teoria logica può essere sottoposta a controllo saggiando la sua capacità di trattare i paradossi logici, ed è un procedimento salutare, quando ci si concentra su una teoria logica, corredare la mente con il maggior numero possibile di rompicapi logici, la cui funzione non è molto diversa da quella assoluta dagli esperimenti nella scienza fisica ». Anche se quest'affermazione nelle intenzioni di Russell doveva applicarsi in primo luogo alla logica deduttiva, particolarmente nella misura in cui questa entra nella matematica pura, penso si attagli egualmente bene alla logica della conferma nelle scienze empiriche.

La teoria della conferma non è nuova ai paradossi logici. Fin dagli anni quaranta i filosofi hanno aguzzato il loro ingegno su due celebri rompicapo. Il primo è il « paradosso dei corvi » di Carl G. Hempel, che può essere descritto nel modo che segue:

L'osservazione di corvi neri (in assenza di osservazioni di corvi di altro colore) verrebbe assunta normalmente a conferma della generalizzazione « Tutti i corvi sono neri ». Questo enunciato è logicamente equivalente a una seconda generalizzazione, « Tutte le cose non-neri sono non-corvi ». La osservazione di cose non-neri che sono non-corvi (in assenza di osservazione di corvi non-neri) sembrerebbe suonare a conferma di questa seconda generalizzazione. Poiché le due generalizzazioni sono logicamente equivalenti tra loro, ciò che vale come evidenza per una delle due deve anche valere come evidenza per l'altra. Quindi l'osservazione di un vaso verde (un non-corvo non-nero) sembra costituisca un'evidenza per l'ipotesi « Tutti i corvi sono neri ». A quanto pare qualcosa non funziona.

Il secondo enigma è il « paradosso del verdlu-berde », dovuto a Nelson Goodman. Siano definiti due termini molto particolari, designanti colori, « verdlu » e « blerde ». Si consideri un arbitrario punto futuro t_0 nel tempo. Di un oggetto esistente in un qualunque tempo t , diremo che è verdlu in t se l'oggetto è verde e il tempo t è precedente a t_0 , ma se il tempo t è posteriore a t_0 , per potersi qualificare così l'oggetto deve essere blu. Se prendiamo come t_0 la mezzanotte del 31 dicembre dell'anno 2000, un oggetto che esiste durante un periodo che si estende tanto nel secolo XX che nel XXI è verdlu durante tutto il periodo se è verde durante il XX secolo ma si

tramuta in blu all'inizio del XXI secolo, rimanendo blu in seguito. « Blerde » è definito in modo analogo: un oggetto è blerde per tutto l'intero lasso di tempo se è blu prima della fine del XX secolo e verde in seguito. Verdlu e blerde sono termini strani, ma perfettamente ben definiti.

Ora sembrerebbe normale il dire che l'osservazione di smeraldi verdi (in assenza di osservazioni di smeraldi di altri colori) tende a confermare la generalizzazione « Tutti gli smeraldi sono verdi ». Poiché però la data in cui vengono scritte queste parole è precedente a t_0 , ciascuno degli smeraldi osservati come verdi è anche osservato come verdlu – almeno ora – per cui le stesse osservazioni confermano la ipotesi « Tutti gli smeraldi sono verdlu ». Che dovremmo prevedere, allora, relativamente agli smeraldi del XXI secolo? Saranno verdi? O saranno verdlu? (Come ha osservato Henry E. Kyburg jr., questo è uno dei problemi più impellenti di tutta la logica della conferma, in quanto abbiamo soltanto 27 anni per risolverlo).

La risposta con cui si reagisce immediatamente e abitualmente alla terminologia verdlu-berde consiste nel dire che è « posizionale », e cioè che chiama in causa un riferimento arbitrario a un punto particolare nel tempo. Questo è certamente vero se partiamo dalle parole del linguaggio ordinario. Come ha osservato Goodman, tuttavia, se noi partiamo dalla terminologia verdlu-berde, allora le parole con cui ordinariamente designiamo i colori risultano posizionali: « verde » significa « verdlu prima di t , ma blerde dopo questo momento ». C'è una qualche ragione per preferire la terminologia ordinaria, oppure la nostra preferenza per questa rispetto alla terminologia verdlu-berde di Goodman è semplicemente il risultato di un accidente storico?

Nell'intraprendere il controllo di una ipotesi un ricercatore usa l'ipotesi per prevedere qualche fenomeno la cui occorrenza o non-occorrenza può essere appurata dall'osservazione. Tuttavia un'ipotesi generale non implica logicamente, di per se stessa, dei fatti osservabili: l'ipotesi deve essere applicata a qualche situazione particolare la cui descrizione costituisce un insieme di « condizioni iniziali ». Per poter prevedere un'eclisse, per esempio, l'astronomo deve conoscere non solo le leggi del moto che governano la Terra e il suo satellite naturale, ma anche le posizioni relative della Terra, della Luna e del Sole in qualche momento

particolare; dalle leggi del moto congiunte alle condizioni iniziali egli può dedurre il tempo e il luogo di una eclisse solare totale. Spesso è necessario manipolare le circostanze così da ottenere un insieme di condizioni iniziali adeguate al controllo di una determinata ipotesi; questo è quanto si fa nell'esecuzione di un esperimento.

C'è un paradosso proposto da Russell che è direttamente pertinente al tipo di schema sperimentale che ora ho abbozzato. Si consideri l'ipotesi: « I maiali hanno le ali ». In congiunzione con il fatto osservato (condizione iniziale) che la carne di maiale è di gusto gradevole, deduciamo la conseguenza (prediciamo) che alcune creature alate sono di gusto gradevole. Quando vediamo che comunemente si mangiano di gusto anatre e tacchini, osserviamo che la conseguenza – o la previsione – è vera; e a quanto pare abbiamo una conferma dell'ipotesi originaria.

Non c'è sforzo di immaginazione che possa far sì che la verifica di una tale conseguenza fornisca un sostegno qualsiasi all'ipotesi in questione. Se chiamiamo « esempio positivo » un dato come questo, con ciò intendendo che si accorda con la previsione dedotta, dobbiamo concludere che gli esempi positivi non forniscono necessariamente la minima credibilità all'ipotesi. Questo esempio apparentemente sciocco sottolinea una morale profonda e significativa: la conferma scientifica è qualcosa di più che non il reperimento di conseguenze vere.

Se ragioniamo deduttivamente dalle premesse « Tutti i mammiferi sono pelosi » e « Le balene sono mammiferi » alla conclusione « Le balene sono pelose », possiamo essere sicuri che la conclusione sarà vera se sono vere le premesse (come effettivamente sono); questa è essenzialmente la caratteristica definitoria dell'inferenza deduttiva valida. Nella deduzione, tuttavia, ragionare a ritroso dalla verità della conclusione alla verità delle premesse è un banale errore logico (noto come « fallacia dell'affermare il conseguente »). D'altro canto dal punto di vista dell'induzione scientifica suona assai attendibile la supposizione che l'osservazione di pelosità nell'embrione di balena (un mammifero, notoriamente) costituisce un'evidenza per la generalizzazione che tutti i mammiferi sono pelosi. Questo doppio binario ha portato Morris R. Cohen a caratterizzare spiritosamente i testi di logica come libri divisi in due parti: nella prima [sulla deduzione] si illustrano le *fallaciae*, e nella seconda [sull'induzione]

o sul metodo scientifico] tali *fallaciae* si commettono a quanto pare impunemente. L'esempio russelliano dei maiali alati prova che, tuttavia, ci sono esempi del sofisma dell'affermare il conseguente che non si qualificano au-

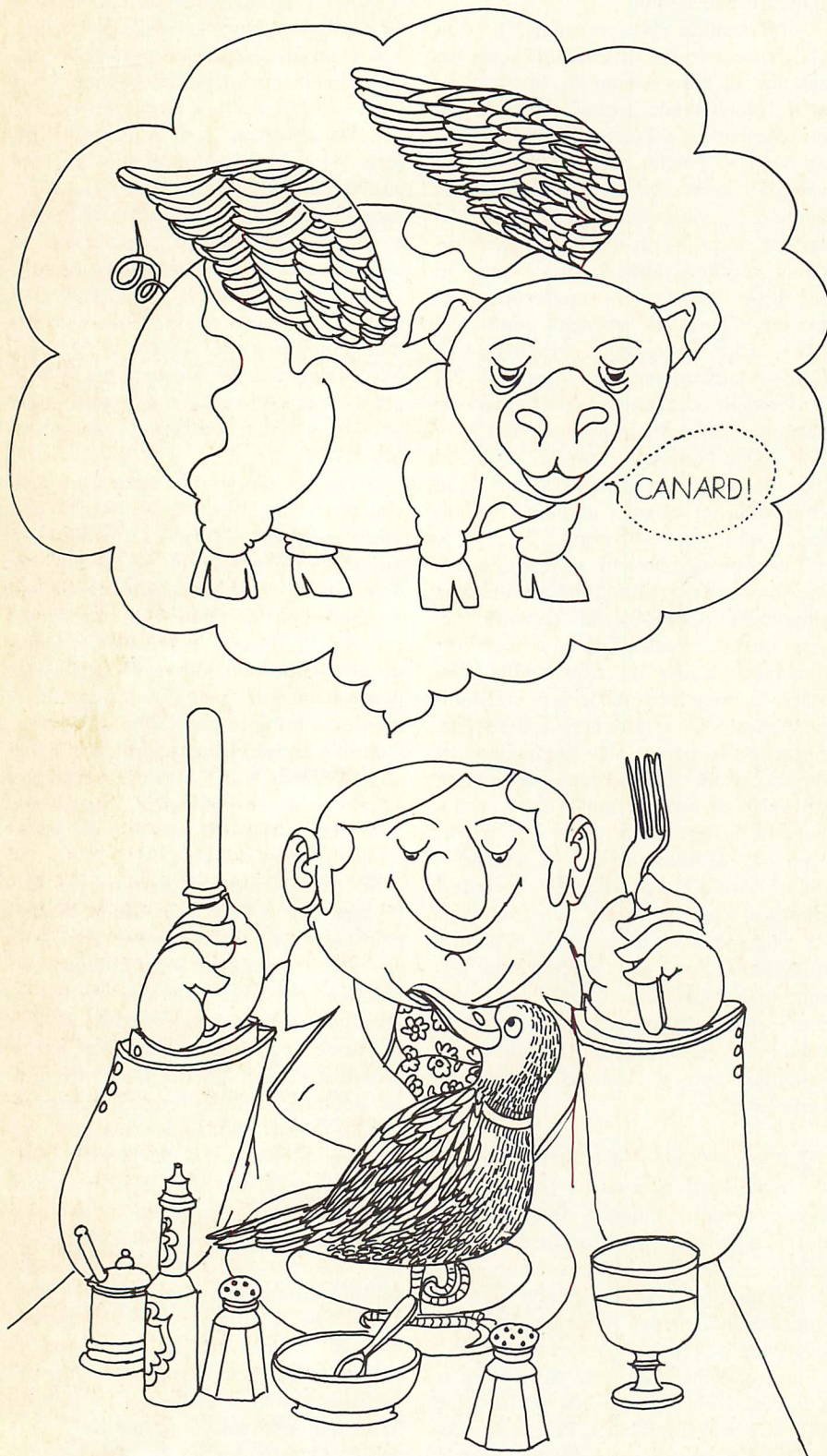
tomaticamente come conferme scientifiche valide.

La relazione di implicazione logica ha l'ovvia e importante proprietà della transitività: se *A* implica logicamente *B* e *B* implica logicamente *C*,

allora *A* implica logicamente *C*. Se la deducibilità di conseguenze vere esaurisse tutto ciò che c'è da dire sulla conferma, allora sarebbe anch'essa una relazione transitiva. Se *C* dovesse confermare *B* per il fatto che segue da *B*, e se *B* dovesse confermare *A* in quanto segue da *A*, allora *C* dovrebbe confermare *A* poiché *C* seguirebbe da *A* per la transitività dell'implicazione logica deduttiva. In verità sembrerebbe intuitivo ritenere la conferma una relazione transitiva. Il ragionamento che segue era effettivamente citato come esempio di ragionamento induttivo valido in un libro pubblicato una dozzina di anni or sono: « Dal momento che, se qui c'era fumo, c'era molto probabilmente fuoco, e se qui c'era fuliggine c'era molto probabilmente fumo, quindi, poiché qui c'era fuliggine, qui c'era probabilmente fuoco ». E tuttavia questo ragionamento ha sostanzialmente la stessa struttura logica del ragionamento che segue: Poiché è molto probabile che un qualsiasi scienziato tra quelli vissuti in tutti i tempi sia vivo oggi (essendo stato calcolato che il 90 % di tutti gli scienziati è ancora vivo), e poiché è molto probabile che un qualsiasi organismo vivo oggi sia un microorganismo, allora, dato che Rossi è uno scienziato, è probabile che sia un microorganismo.

Questo esempio prova inequivocabilmente che *A* (essere uno scienziato) può fornire conferma a *B* (essere vivo attualmente), e *B* può a sua volta fornire sostegno a *C* (essere un microorganismo), mentre *A* non solo manca miseramente di confermare *C* ma è addirittura incompatibile con esso. L'esempio illustra un serio rischio sottostante alle riflessioni sulla conferma: il pericolo di affidarsi ad analogie infondate con la deduzione. Viene spontaneo assumere intuitivamente che le proprietà delle relazioni deduttive si applichino più o meno esattamente alla logica della conferma. Questo non è vero. Alcune delle proprietà più importanti delle relazioni deduttive vengono meno in modo completo e assoluto quando ci si allontana, anche di poco, dalla logica deduttiva stessa. La relazione di transitività è un eccellente esempio di questo fenomeno.

L'antidoto migliore contro gli errori che sorgono dalle intuizioni relative alla conferma, consiste, a mio avviso, nel concentrare l'attenzione sulla teoria matematica della probabilità (si veda l'illustrazione in basso a pagina 108). Se è possibile considerare la conferma come un tipo di probabilità, è facile provare che le re-



« I maiali hanno le ali » è l'ipotesi considerata. Il ricercatore sa che la carne di maiale è di gusto gradevole. In base all'ipotesi egli prevede che alcune creature alate devono essere di gusto gradevole. Per controllare l'ipotesi assaggia l'anitra e, trovandola di gusto gradevole, considera l'ipotesi confermata. (Il maiale — *canard* in francese oltre che anatra significa frottola — mostra di conoscere l'irrilevanza di questa conseguenza vera.)

lazioni di conferma non sono transitive. Si può anche provare che la verità della conseguenza di una ipotesi non accresce necessariamente la probabilità di quella ipotesi. Anche se mettiamo scopertamente in gioco il calcolo delle probabilità, c'è tuttavia ancora pericolo che si produca una seria confusione, provocata da una semplice ambiguità.

Dire, per esempio, che la teoria della relatività ristretta è stata confermata potrebbe significare l'una o l'altra di due cose diverse. O qualche elemento di evidenza, come il recente esperimento relativo al ritardo relativistico di orologi in movimento, ha aumentato in qualche misura la credibilità della teoria, o la teoria è stata sostenuta da un corpo di evidenza così largo e svariato da potersi considerare, almeno provvisoriamente, una legge scientificamente accettabile. Un'ipotesi generale che abbia inizialmente una credibilità relativamente bassa può vederla accresciuta in qualche misura dalla nuova evidenza senza con ciò attingere un alto grado di conferma. In casi siffatti un'ipotesi può essere ovviamente confermata in un senso senza essere confermata nell'altro. Conveniamo di parlare dell'accrescimento di probabilità di un'ipotesi come di conferma « incrementale » e del raggiungimento di un'alta probabilità come di conferma « assoluta ». Al fine di sottolineare la distinzione circoscriviamo il termine « confermato » (senza qualifica) al senso incrementale e impieghiamo sempre un'espressione come « altamente confermato » per il senso assoluto. Dovremmo dire, allora, che la teoria della relatività ristretta è stata confermata dall'esperimento di ritardo dell'orologio e che è altamente confermata dal corpo di evidenza totale che lo sostiene.

Anche se la distinzione tra questi due sensi di conferma è ovvia ed è stata riconosciuta da lungo tempo, le sue implicazioni non sempre sono state chiaramente riconosciute. Come dato di fatto, il senso incrementale ha alcune strane proprietà che vengono correntemente trascurate in quanto non sono condivise dal concetto assoluto. Si consideri l'esempio che segue, il quale, benché totalmente fittizio, espone tuttavia circostanze logicamente possibili.

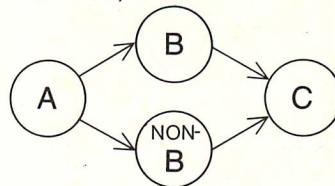
Bianchi consulta il suo medico per un disturbo respiratorio. Dopo un esame preliminare il medico dice di ritenere che Bianchi abbia la polmonite ma di non essere sicuro se essa sia batterica, virale oppure (eventualità molto rara) di ambedue i tipi insieme. Si

IMPLICAZIONE LOGICA

SE A IMPLICA LOGICAMENTE B E B IMPLICA LOGICAMENTE C, ALLORA A IMPLICA LOGICAMENTE C.

CORRISPETTIVO PROBABILISTICO

PROBABILITÀ (C, DATO A) = PROBABILITÀ (B, DATO A) × PROBABILITÀ (C, DATI A E B) + PROBABILITÀ (NON-B, DATO A) × PROBABILITÀ (C, DATI A E NON-B).



La transitività è una proprietà dell'implicazione logica: se essere una balena implica logicamente essere un mammifero ed essere un mammifero implica logicamente essere peloso, allora essere una balena implica logicamente essere peloso. Il corrispettivo probabilistico è più complesso: la probabilità che un essere umano (A) si ammali di cancro polmonare (C) è uguale alla probabilità che (A) sia un fumatore (B) per la probabilità che un fumatore (A e B) si ammali di cancro più la probabilità che (A) sia un non-fumatore (non-B) per la probabilità che un non-fumatore (A e non-B) si ammali.

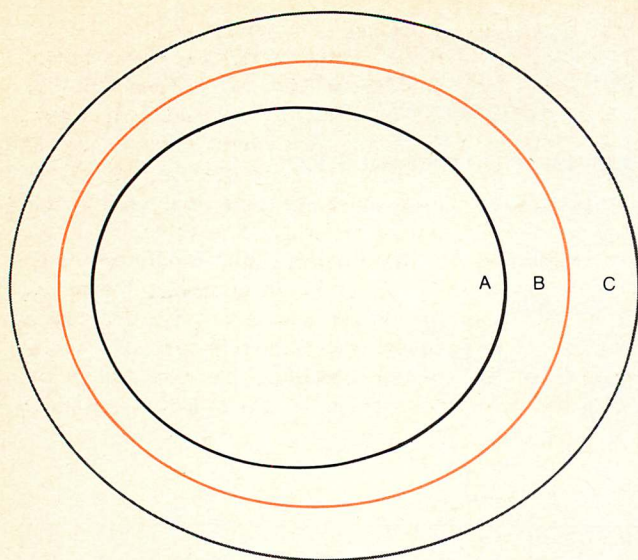
rende necessario un controllo ulteriore. Bianchi è sottoposto al test prescritto. Il medico gli dice che il test ha confermato l'ipotesi che la malattia sia polmonite batterica e ha anche confermato l'ipotesi che sia polmonite virale, ma ha infirmato l'ipotesi che Bianchi abbia la polmonite! La maggior parte di noi troverebbe poco da ridire se Bianchi a questo punto si rivolgesse a un altro medico.

Eppure può darsi che il medico sia dalla parte della ragione. Si supponga che sulla base dell'esame superficiale egli concluda che c'è un 96 % di probabilità che Bianchi abbia la polmonite, senza però avere indicazioni sul fatto che essa sia batterica o virale (assumendo che questi siano gli unici due tipi). Inoltre egli stabilisce che c'è la probabilità del 2 % che Bianchi abbia ambedue i tipi di polmonite. Conseguentemente la probabilità che Bianchi abbia la polmonite batterica è del 49 % e la probabilità che abbia la polmonite virale è pure del 49 %. Si supponga inoltre che ci sia un test che rileva molto attendibilmente i rari casi in cui sono presenti insieme ambedue i tipi. Quando il test è somministrato a Bianchi, il risultato è positivo, rendendo certo all'89 % il fatto che egli ha ambedue i tipi. Si assuma inoltre che raramente questo test risulta sbagliato per chi ha solo un tipo di polmonite; cioè, se il risultato è positivo e l'individuo non ha ambedue i tipi, molto probabilmente non ha nessuno dei due tipi. In particolare il risultato positivo del test significa che c'è la

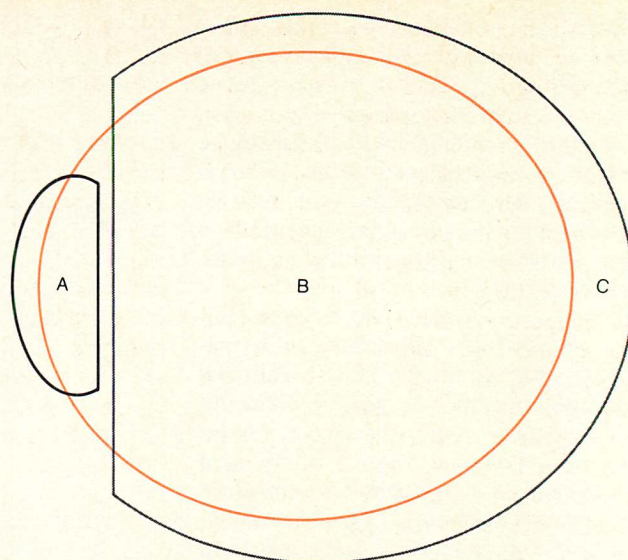
probabilità del 90 % che Bianchi abbia la polmonite batterica, la probabilità del 90 % che abbia la polmonite virale e conseguentemente la probabilità del 91 % che abbia una o l'altra o ambedue (si veda l'illustrazione a pagina 110).

Questo fantasioso esempio illustra la eventualità che l'evidenza possa confermare ciascuna delle due ipotesi e tuttavia infirmare la loro disgiunzione (la combinazione l'uno-o-l'altro-o-ambedue). In questo caso il risultato del test speciale ha accresciuto la probabilità che Bianchi avesse la polmonite batterica e anche la probabilità che avesse la polmonite virale, mentre nello stesso tempo ha abbassato la probabilità che avesse l'una o l'altra – in altri termini, dato che ce ne sono soltanto due tipi, la probabilità che avesse la polmonite!

Per vedere come questo è possibile dobbiamo rivolgerci alla regola di addizione per probabilità (si veda l'illustrazione in basso a pagina 108). Al fine di calcolare la probabilità per una disgiunzione di alternative non esclusive sommiamo le probabilità di ciascuna delle due alternative prese separatamente e quindi sottraiamo la probabilità della loro occorrenza congiunta. Per esempio, la probabilità di estrarre da un normale mazzo da bridge o un onore (un asso, un re, una regina, un fante o un 10) o un *atout* (poniamo un picche) è eguale alla probabilità di un onore (20/52) più la probabilità di un picche



«Tutti» e «quasi tutti» possono essere molto diversi. Dato che tutti gli *A* sono *B* e tutti i *B* sono *C*, deve essere vero che tutti gli *A* sono *C* (a sinistra). Tuttavia, dato che quasi tutti gli *A* sono *B* e quasi tutti i *B* sono *C*, può verificarsi che nessun *A* sia *C*



(a destra). Nell'esempio del testo, anche se per la maggior parte gli scienziati (*A*) risultano essere esseri viventi (*B*) e per la maggior parte gli esseri viventi (*B*) sono microorganismi (*C*), non ci sono affatto scienziati che risultano microorganismi.

(13/52) meno la probabilità di un onore di picche (5/52). La sottrazione è necessaria poiché gli onori di picche sono stati, per così dire, contati due volte: una volta come onori, una seconda come picche. La risposta risulta 28/52 (si veda l'illustrazione a pag. 111). Nell'esempio della polmonite la probabilità della disgiunzione decresce mentre aumenta la probabilità di ciascuna delle alternative in virtù dell'alta probabilità (89 %) dell'occorrenza congiunta (o contata due volte) dopo il test, in contrasto con la sua bassa probabilità (2 %) prima del test.

Risultati come questi sono tipici della conferma incrementale, e suonano sconcertanti in parte perché c'è una tendenza naturale a confondere il concetto incrementale con quello assoluto. La conferma in senso assoluto significa che un'ipotesi ha un'alta probabilità di essere corretta. Se la pol-

monite virale è altamente confermata in senso assoluto, allora la polmonite è per lo meno confermata nello stesso grado. Questo perché l'avere la polmonite (di qualche tipo) è una conseguenza logica dell'avere la polmonite virale, e una regola fondamentale della probabilità matematica afferma che la probabilità attribuita alla conseguenza logica di una proposizione è uguale o maggiore a quella della proposizione stessa. Non ho violato questa condizione nel descrivere il caso di Bianchi: la probabilità della polmonite è maggiore di quella della polmonite virale in base alla diagnosi preliminare, nonché in base al test speciale. E tuttavia, quando ci si mette dal punto di vista della conferma incrementale (un cambiamento nella conferma) risulta che la nuova evidenza può accrescere la probabilità di una proposizione e tuttavia diminuire la probabilità di una

delle sue conseguenze. Questo è un sorprendente contrasto logico tra conferma incrementale e assoluta.

Possono verificarsi fatti anche più strani. Supponiamo che due ricercatori si accingano a sottoporre a controllo una ipotesi. Ciascuno dei due si reca al suo laboratorio per eseguire un esperimento e ciascuno acquisisce una scoperta positiva: una conferma dell'ipotesi. Può mai accadere che, anche se ciascuna delle scoperte conferma la ipotesi, la loro congiunzione infirmi la stessa? Sì, è logicamente possibile.

Siano *A* e *B* gli atomi di un isotopo immaginario che può decadere radioattivamente in uno qualsiasi di tre modi diversi. Dato il verificarsi della disintegrazione, c'è una probabilità di 0,7 che sia stata emessa una particella alfa, una probabilità di 0,2 che sia stato emesso un elettrone negativo, e una probabilità di 0,1 che sia stato

ADDIZIONE	PROBABILITÀ (A O B) = PROBABILITÀ (A) + PROBABILITÀ (B) - PROBABILITÀ (A E B).
CONSEGUENZA LOGICA	SE A IMPLICA LOGICAMENTE B, ALLORA PROBABILITÀ (B) È EGUALE O MAGGIORE ALLA PROBABILITÀ (A).
MOLTIPLICAZIONE	PROBABILITÀ (A E B) = PROBABILITÀ (A) × PROBABILITÀ (B, DATO A).
TEOREMA DI BAYES	PROBABILITÀ (A, DATO B) = $\frac{\text{PROBABILITÀ (A)} \times \text{PROBABILITÀ (B, DATO A)}}{\text{PROBABILITÀ (B)}}$

Il concetto di probabilità svolge nella conferma empirica lo stesso ruolo che svolge la deduzione logica nella dimostrazione ma-

tematica. Sono qui elencate alcune regole della probabilità menzionate nel testo. Un'altra è illustrata nella pagina precedente.

emesso un elettrone positivo (positone). Si supponga che ambedue questi atomi siano appunto stati disintegrati e che le particelle emesse si accostino tra loro. Prendiamo in considerazione l'ipotesi che incontrandosi esse si annichilino reciprocamente, cosa che si verificherà se e solo se una particella è un elettrone negativo e l'altra è un positone. Supponendo di non disporre di informazioni ulteriori, la probabilità dell'annichilazione è 0,2 per 0,1 più 0,1 per 0,2, ossia 0,04. Supponiamo che un fisico scopra che un atomo *A* ha emesso un elettrone; in base a questa evidenza la probabilità che si verifichi l'annichilazione è di 0,1, giacché è questa la probabilità che *B* emetta un positone. Supponiamo che un altro fisico scopra che l'atomo *B* ha emesso un positone; in base a questa evidenza (ma senza l'evidenza ottenuta dal primo fisico) c'è, analogamente una probabilità di 0,1 che si verifichi l'annichilazione. Eppure sulla base di ambedue gli elementi di evidenza in congiunzione è chiaro che l'annichilazione è impossibile.

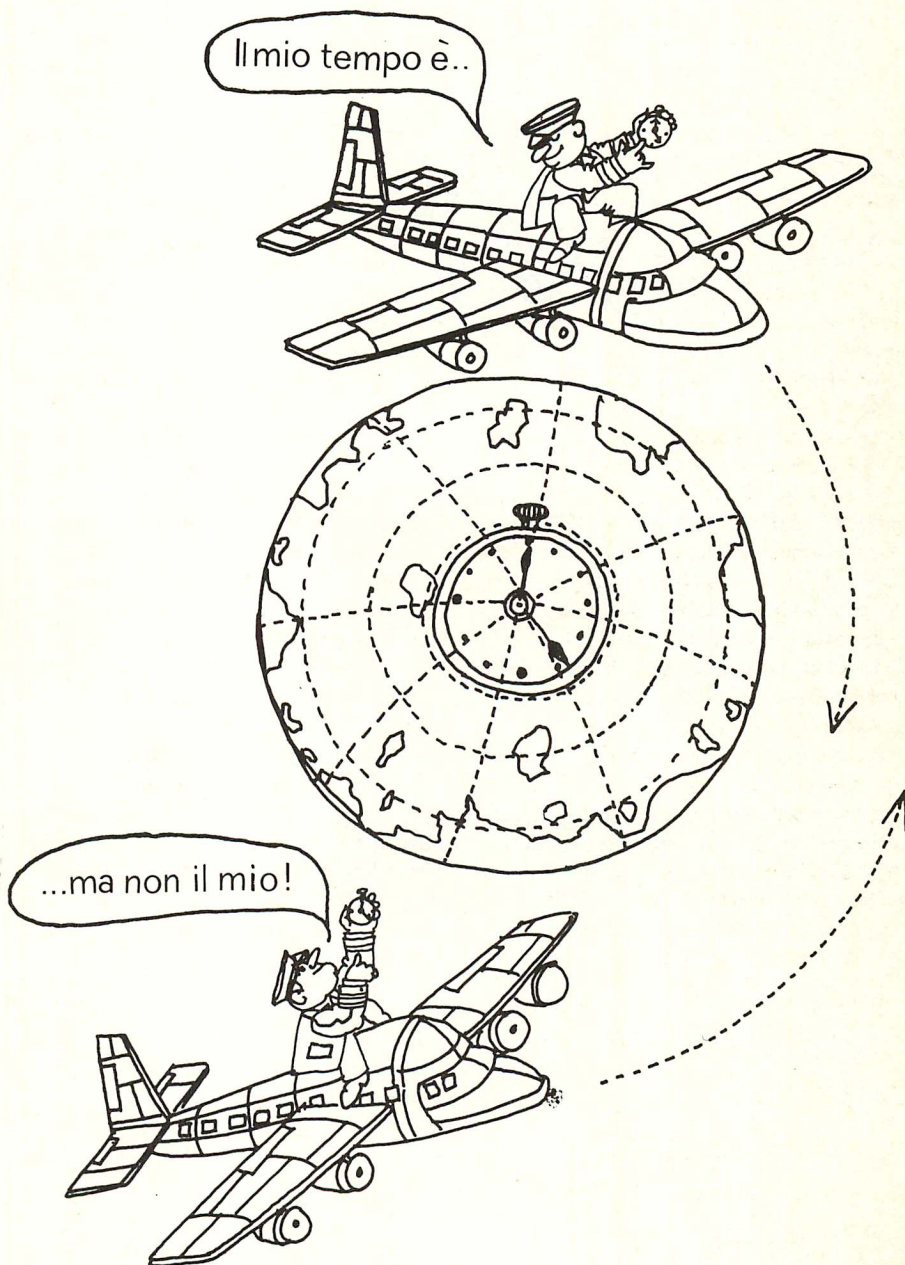
Ciascuno dei due elementi di evidenza conferma, separatamente, l'ipotesi dell'annichilazione (ciascuno innalza la sua probabilità da 0,04 a 0,1). In congiunzione, tuttavia, essi non solo la infirmano, ma di fatto la refutano. Potrebbe mai verificarsi una cosa del genere nel corso di una indagine? In un esperimento relativo alla diffusione Compton di fotoni da parte di elettroni un fisico potrebbe misurare la frequenza di fotoni dispersi con un angolo particolare mentre un altro potrebbe misurare l'energia degli elettroni di rinculo. Anche se ciascun insieme di misure conferma l'ipotesi della diffusione Compton, come si può essere sicuri che continuino a farlo presi insieme? Nel caso della diffusione Compton accade che la congiunzione delle scoperte conferma l'ipotesi, ma questo non segue automaticamente dal semplice fatto che, separatamente prese, le scoperte sono conferme; esso dipende da diverse circostanze, compreso il fatto che la congiunzione stessa è una delle previsioni della teoria. L'esempio dell'annichilazione prova, tuttavia, che ci sono questioni profonde e fondamentali circa la legittimità del supporre che l'accumulazione di molti risultati di test confermantici accresca inevitabilmente la credibilità delle ipotesi scientifiche.

Finora abbiamo parlato solo di conferma, o del caso in cui il risultato del test è positivo. Resta da discutere il risultato del test negativo: il caso in

cui la previsione dell'ipotesi risulta falsa. Si è spesso sostenuto che c'è una forte asimmetria tra i casi positivi e negativi a ragione di una semplice circostanza logica: che mentre l'inferenza da una conclusione vera alla verità delle premesse è un caso di fallacia dell'affermare il conseguente, l'inferenza dalla falsità della conclusione alla falsità di almeno una premessa è perfettamente valida; essa è nota come *modus tollens*, o « negare il conseguente ». Poiché una deduzione valida è definita come quella la cui con-

clusione deve essere vera se ha premesse vere, possiamo effettivamente concludere che un ragionamento deduttivo valido con una conclusione falsa non può avere premesse che sono tutte vere. Sembra quindi che un esito negativo non solo infirmi una ipotesi ma di fatto la refuti.

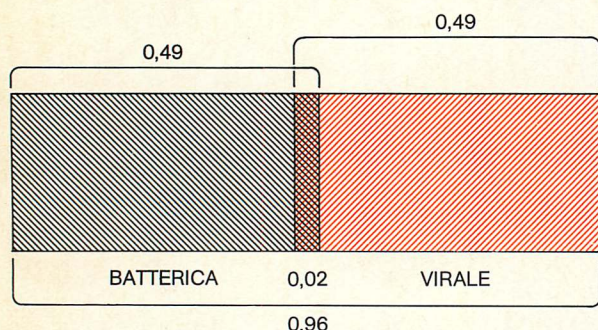
La situazione non è affatto semplice. Come ha fatto notare Pierre Duhem, nella maggior parte dei casi in cui si intraprende il controllo di una ipotesi scientifica, per non dire in tutti, entrano in gioco ipotesi ausiliarie.



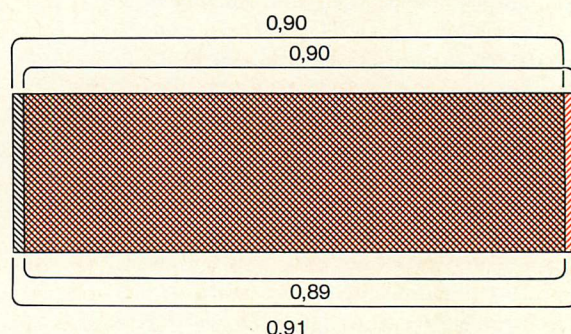
La teoria della relatività permette di prevedere che quanto più velocemente è trasportato un orologio, tanto più lentamente procede. L'anno scorso Joseph C. Hafele della Università di Washington, insieme ad altri, ha riferito di un esperimento in cui orologi al cesio erano stati trasportati in un giro del mondo su jet commerciali in direzioni opposte. In relazione alle stelle fisse uno si muoveva circa 1000 chilometri più velocemente della Terra, e l'altro 1000 chilometri più lentamente; uno avrebbe dovuto perdere tempo, e l'altro guadagnarlo, rispetto agli orologi di riferimento a terra. Gli scarti di tempo risultarono vicini ai valori previsti, fornendo in questo modo conferma alla teoria.



PROBABILITÀ (POLMONITE) = PROBABILITÀ (BATTERICA)
+ PROBABILITÀ (VIRALE) - PROBABILITÀ (AMBEDUE), OSSIA
 $0,49 + 0,49 - 0,02 = 0,96$



PROBABILITÀ (POLMONITE) = PROBABILITÀ (BATTERICA)
+ PROBABILITÀ (VIRALE) - PROBABILITÀ (AMBEDUE), OSSIA
 $0,90 + 0,90 - 0,89 = 0,91$



Ciò che conferma ciascuna delle due ipotesi può contemporaneamente infirmare la loro disgiunzione. Dopo un esame superficiale il medico decide che c'è la probabilità del 96 % che Bianchi abbia la polmonite: non sa se batterica o virale, ma ritiene che ci sia solo la probabilità del 2 % che siano presenti

ambedue i tipi. Le probabilità sono schematizzate dalle linee sovrapposte (a sinistra). Si predispose un test di controllo. Il risultato accresce notevolmente la probabilità che Bianchi abbia la polmonite batterica e anche quella che abbia la polmonite virale, ma riduce al 91 % quella che abbia la polmonite (a destra).

Nel sottoporre a controllo una teoria astronomica è probabile vengano usati telescopi e altri strumenti, cosicché vengono chiamate in causa le leggi dell'ottica e altre leggi nello sforzo di stabilire una connessione tra una macchia su una lastra fotografica e un corpo celeste. In una scienza a un certo grado di sofisticazione le condizioni iniziali richieste per sottoporre a controllo una ipotesi sono difficilmente accertabili per mezzo di osservazione diretta; occorrono delle ipotesi ausiliarie per correlare ciò che si osserva effettivamente alle condizioni iniziali opportune. Inoltre, gli esiti previsti possono essere in qualche misura lontani dall'osservazione diretta, e di nuo-

vo entrano in gioco ipotesi ausiliarie.

Il risultato evidente di queste complicazioni è che l'esito negativo del test sperimentale di un'ipotesi non si può prendere automaticamente come refutazione di quell'ipotesi. Il risultato negativo del test prova soltanto che da qualche parte c'è qualcosa che non va. Può darsi che l'ipotesi sottoposta a controllo sia falsa o può darsi che sia falsa qualche ipotesi ausiliaria.

A rigore, il risultato negativo non fa che confutare la congiunzione delle ipotesi ausiliarie e dell'ipotesi sotto controllo; ma non refuta l'ipotesi sottoposta a controllo di per se stessa. (Come interessante esempio storico, si considerino le false previsioni circa i

movimenti di Urano emerse dalla meccanica newtoniana del XIX secolo. Anziché confutare la fisica newtoniana essere portarono alla scoperta di Nettuno; l'esito negativo fu attribuito all'ipotesi ausiliaria o forse alle stesse condizioni iniziali. Più tardi, tuttavia, le irregolarità nell'orbita di Mercurio non portarono alla scoperta di un nuovo pianeta secondo quanto era stato previsto in base alla meccanica newtoniana ma piuttosto al tracollo definitivo della meccanica newtoniana: la precessione del perielio di Mercurio è stato un primo elemento di evidenza cruciale per la teoria della relatività generale di Einstein).

Alla luce della fondamentale intui-

zione di Duhem che morale bisognerebbe trarre, agli effetti del confermare o infirmare un'ipotesi, dall'esito negativo di un test? Sicuramente o è infirmata l'ipotesi sotto controllo o sono infirmate in qualche misura le ipotesi ausiliarie; sicuramente né le ipotesi ausiliarie né l'ipotesi principale possono essere confermate dall'esito negativo. Ambedue queste assunzioni suonano plausibili ma sono sbagliate. È logicamente possibile che un esito sperimentale refuti la congiunzione di due ipotesi e tuttavia confermi ciascuna di esse presa individualmente. Questa possibilità è strettamente correlata all'esempio del decadimento.

Assumiamo lo stesso sistema di atomi dell'esempio: gli atomi *A* e *B* hanno ciascuno tre possibili modalità di decadimento, con probabilità 0,7 per la particella alfa, 0,2 per l'elettrone negativo e 0,1 per il positone. Questa volta l'evidenza osservata è l'annichilazione delle particelle. Consideriamo l'ipotesi che l'atomo *A* emetta un elettrone negativo. Poiché sappiamo, in virtù dell'annichilazione, che uno degli atomi ha emesso un elettrone negativo, senza sapere quale, l'ipotesi che esso sia stato emesso da *A* ha probabilità 0,5. Lo stesso vale per l'ipotesi che *B* emetta un elettrone negativo. Il fatto che si sia verificata l'annichilazione rende tuttavia impossibile la emissione da parte di ambedue gli atomi di elettroni negativi; esso perciò refuta la congiunzione delle due ipotesi. Nondimeno proprio questa evidenza conferma separatamente ciascuna delle due ipotesi, in quanto in ognuno dei casi ha innalzato la probabilità da 0,2 a 0,5.

Si immagini che cosa potrebbe significare questo genere di fatti per la metodologia scientifica. Lo scienziato Verdi, torna a casa a tarda sera dopo una dura giornata di laboratorio. «Come è andato il tuo lavoro oggi, caro?» chiede sua moglie.

«Sai della "ipotesi di Verdi", su cui ho impegnato per intero la mia reputazione? Be', oggi l'ho sottoposta a una verifica sperimentale, e l'esito è stato negativo».

«Oh caro, che peccato! Questo significa forse che la tua ipotesi prediletta è affossata e che la tua reputazione scientifica è a pezzi?».

«Non del tutto. Per eseguire il test ho fatto uso di alcune ipotesi ausiliarie».

«Oh che sollievo, salvato da Duhem! La tua ipotesi dopo tutto non è stata refutata». La signora Verdi trae un profondo sospiro.

«Meglio ancora» continua Verdi.

«Ho addirittura confermato l'ipotesi.»
«Oh, caro, è meraviglioso!» risponde la signora Verdi. «Devi aver trovato che respingendo l'ipotesi ausiliaria potevi provare che il test sosteneva effettivamente la tua ipotesi. Che intelligenza!».

«No» continua Verdi «è ancora meglio. Ho scoperto di aver confermato anche l'ipotesi ausiliaria!»

Queste inconsuete possibilità sembra facciano strazio della metodologia scientifica. Il fatto che gli scienziati del giorno d'oggi nella pratica quotidiana non si imbattano realmente in difficoltà del genere può essere considerato da molti una prova del fatto che la conferma è questione di intuizione scientifica e resiste a tutti i tentativi di formalizzazione. Io non credo che tale conclusione sia giustificata, e sosterrò il mio punto di vista mettendo in luce un parallelo con la storia del concetto di dimostrazione in matematica.

L'idea di dimostrazione in matematica è emersa verso il 600 a.C. Si attribuisce a Talete di Mileto il merito di aver importato la geometria dall'Egitto alla Grecia e di aver contribuito al processo che ha trasformato questa in una scienza matematica. Anche se gli Egiziani avevano applicato la geometria nel rilevamento topografico, nulla prova che essi abbiano dimostrato in realtà un qualsiasi teorema geometrico; si crede che Talete abbia dimostrato, oltre ad altri teoremi, che gli angoli alla base di un triangolo isoscele sono eguali. All'incirca verso il 300 a.C. Euclide aveva riformulato la geometria come un sistema assiomatico in cui tutti i teoremi devono essere dedotti da un piccolo numero di assiomi o postulati. Alcuni frammenti elementari di logica deduttiva furono sviluppati nell'antichità da Aristotele e dai filosofi stoici. Eppure non fu prima del 1879 che Gottlob Frege sviluppò una logica deduttiva che poteva cominciare a caratterizzare adeguatamente la deduzione in matematica. Ci vollero quindi come minimo 2500 anni, dal tempo in cui furono per la prima volta impiegate delle dimostrazioni matematiche, perché i logici giungessero a un chiaro intendimento della loro natura.

La logica matematica è fiorita negli ultimi 100 anni, e sono stati stabiliti molti risultati importanti. Questo processo non è stato privo di vicissitudini. Per esempio Russell ha scoperto una celebre contraddizione proprio nella logica su cui Frege aveva cercato di basare tutta la matematica. Essa sor-

geva dalla considerazione di paradossi simili a quello famoso del barbiere: In una certa città c'è un barbiere che rade tutti gli uomini che non si radono da soli. Chi rade il barbiere?

Il fatto che i matematici nella loro attività non si siano trovati costantemente involuppati in contraddizioni non ha impedito al paradosso di Russell di costituire un fattore di crisi nei fondamenti della matematica. Altri sviluppi sconcertanti, come la dimostrazione data da Kurt Gödel della essenziale incompletezza dell'aritmetica, erano a dir poco preoccupanti, per quanto non portassero alla luce di per sé delle contraddizioni.

Gli scienziati empirici hanno fatto osservazioni e compiuto esperimenti per controllare ipotesi complesse fin dalla nascita della scienza moderna, nei secoli XVI e XVII. Quando si trat-

♠
♥
♦
♣

A	A	A	A
K	K	K	K
Q	Q	Q	Q
J	J	J	J
10	10	10	10
9	9	9	9
8	8	8	8
7	7	7	7
6	6	6	6
5	5	5	5
4	4	4	4
3	3	3	3
2	2	2	2

PROBABILITÀ (ONORE O PICCHE) =
PROBABILITÀ (ONORE) + PROBABILITÀ
(PICCHE) – PROBABILITÀ (ONORE DI PICCHE),
OSSIA 20/52 + 13/52 – 5/52 = 28/52

La probabilità di estrarre un onore o un picche è determinata dall'addizione delle probabilità, come qui si fa vedere. La sottrazione dei cinque onori di picche è necessaria perché, come indicano i rettangoli, altrimenti questi vengono contati due volte. La probabilità all'89 % dell'occorrenza congiunta della polmonite batterica e virale dell'esempio precedente risulta analoga all'occorrenza dell'onore di picche.

ta di trarre conclusioni dai risultati di queste osservazioni ed esperimenti, siamo ben lungi dall'aver una chiara comprensione del tipo di ragionamento che entra in gioco. Ora siamo in una situazione analoga a quella in cui si trovava la matematica durante i millenni in cui la dimostrazione matematica era usata spesso e con buoni risultati, mentre la logica a essa sottostante rimaneva fondamentalmente

oscura. Il lavoro d'analisi correntemente eseguito in teoria della conferma e in logica induttiva sta tentando di porre rimedio a questa situazione.

C'è una grande divergenza di opinioni circa la via migliore da seguire nel cercare di trattare i problemi della conferma. Sono due gli espedienti che mi sembra offrano una considerevole speranza di aiuto. Il primo di questi è il teorema di Bayes, un semplice

teorema del calcolo matematico delle probabilità (si veda l'illustrazione in basso a pagina 108). Il teorema di Bayes è spesso chiamato «regola della probabilità inversa». Data la probabilità che una certa evidenza varrebbe se dovesse valere un'ipotesi particolare (e date anche altre probabilità), il teorema di Bayes ci consente di calcolare la probabilità che l'ipotesi sia vera una volta scoperta la suddetta evidenza. In certi casi almeno esso può essere usato per accertare la probabilità che fosse operante qualche causa particolare, dato il verificarsi di un certo effetto.

Il teorema di Bayes è stato utilizzato largamente in anni recenti dagli studiosi di statistica che si sono auto-definiti bayesiani, in particolare negli ultimi lavori di L.J. Savage. Il teorema di Bayes ha fornito alla teoria della conferma uno schema che sembra assai più adeguato all'inferenza scientifica di quanto non possa mai sperare di esserlo la fallacia dell'affermare il conseguente.

Il secondo espediente nasce da un chiaro riconoscimento del concetto incrementale di conferma in quanto opposto al concetto assoluto. La conferma incrementale chiama in causa il cambiamento di probabilità, che è fondamentalmente un concetto di rilevanza probabilistica. È così aperta la strada alla definizione di una misura di rilevanza basata sulla caratteristica matematica della probabilità, per mezzo della quale si può studiare la conferma incrementale in modo sistematico e preciso. Tale misura è stata definita ed elaborata da Rudolf Carnap nel 1950, ma sembra che a essa non sia stata prestata sufficiente attenzione. Il problema della polmonite e i due paradossi dell'annichilazione sono stati individuati sulla piattaforma della trattazione carnapiana della rilevanza, e credo che un ulteriore sforzo di attenzione verso la nozione formale di conferma incrementale priverà questi esempi della loro apparenza paradossale.

Allo stadio attuale di sviluppo gli studi di teoria della conferma e di logica induttiva hanno prodotto più paradossi ed enigmi che non soluzioni convincenti o largamente accettate di problemi fondamentali. Altre ricerche su questi problemi dovrebbero, tuttavia, produrre notevoli approfondimenti nella logica delle scienze empiriche, allo stesso modo in cui gli studi sui fondamenti della matematica sono stati ripagati dai notevoli progressi conseguiti nella comprensione della logica di questa disciplina.



La congiunzione di due conferme può refutare una ipotesi. Se decade l'atomo A (o B), la probabilità che emetta un elettrone è 0,2 e quella che emetta un positone è 0,1. La probabilità dell'annichilazione, che si verificherà se un atomo emette un elettrone e l'altro un positone, è perciò $0,2 \times 0,1 + 0,1 \times 0,2$, ossia 0,04. Se lo scienziato A osserva soltanto che l'atomo A ha emesso un elettrone, considera che la probabilità dell'annichilazione si è accresciuta a 0,1 (la probabilità che l'atomo B emetta un positone). Lo scienziato B, osservando soltanto che l'atomo B ha emesso un elettrone, registra lo stesso aumento di probabilità. La congiunzione di ambedue le osservazioni (corrispondente all'emissione di due elettroni), tuttavia, significa che non può esserci annichilazione.

Problemi non risolti dell'aritmetica

È nella natura dell'aritmetica porsi più problemi di quanti non ne possa risolvere. L'aver dimostrato che ci sono problemi che è impossibile risolvere è una delle grandi conquiste della logica matematica

di Howard DeLong

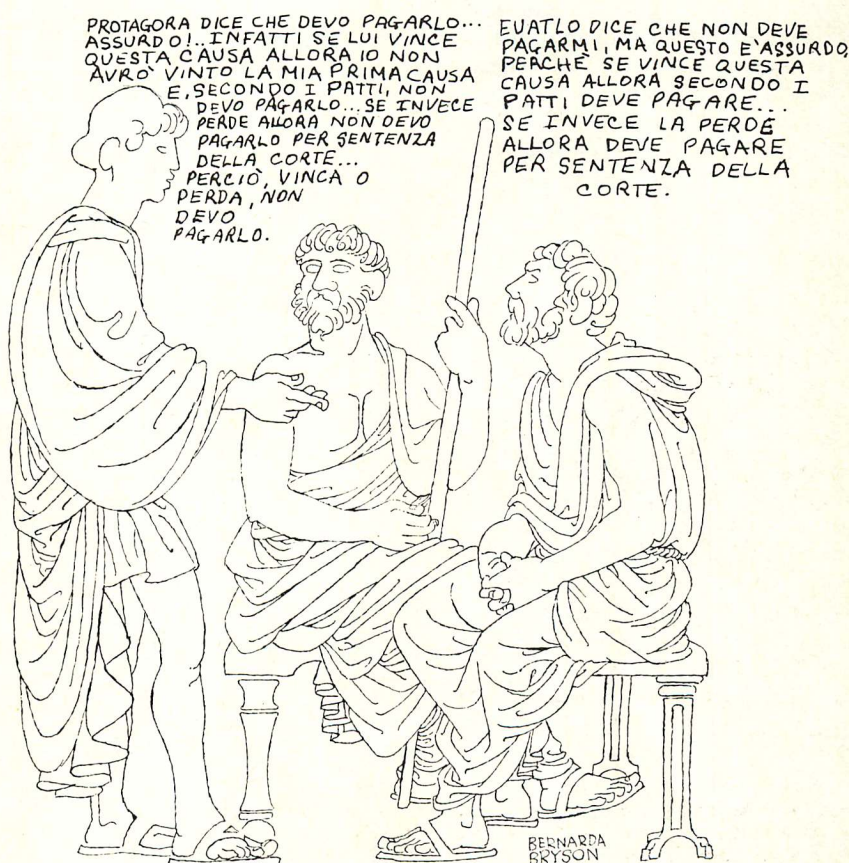
Ci sono problemi aritmetici la cui comprensione è alla portata di un ragazzo di dieci anni e che ciononostante sono rimasti insoluti per decine, centinaia, anche migliaia di anni. Per molto tempo si è pensato che, se un problema di aritmetica si poteva formulare con precisione, solo un difetto di capacità poteva essere d'ostacolo alla sua soluzione; se un problema poi restava insoluto, si pensava che una soluzione si sarebbe alla fine trovata grazie all'applicazione di qualche nuovo metodo, o all'applicazione più avveduta di qualche metodo già usato. La metalogica, un settore della logica matematica contemporanea, ha chiarito che il continuo scacco subito nei tentativi di soluzione di alcuni problemi matematici non risolti può essere dovuto non tanto a scarsa abilità quanto a limitazioni connaturate alle capacità degli uomini e delle macchine.

È stato attraverso un affinamento di alcune tendenze già vive tra i greci che la logica matematica è venuta alla luce, negli ultimi anni del secolo XIX e nei primi del secolo XX. Quindi il nostro discorso deve cominciare con Aristotele, l'inventore della logica, che diede vita a questa scienza in risposta a una duplice istanza, filosofica e matematica insieme. Il problema filosofico nasceva dalla necessità di rispondere alla estrema varietà di argomentazioni che si trovava a dover affrontare. Talete aveva sostenuto che la sostanza fondamentale del mondo è l'acqua, Anassimandro che non è una cosa particolare ma qualcosa di indeterminato; Eraclito aveva sostenuto che tutte le cose sono in movimento, Parmenide che non lo è nessuna; Protagora aveva affermato che i giudizi etici sono relativi, Socrate che non lo sono, e così via. Aristotele doveva anche essere allenato a trattare con paradossi, come il paradosso zenoniano di Achille e della

tartaruga, e con argomentazioni semi-legali del tipo di quelle che si dice siano intercorse tra Protagora ed Euatlo (si vedano le figure a pagina 113 e a pagina 116). Si sentiva la necessità di porre dei principi generali in grado di fornire metodi sistematici per distinguere le argomentazioni corrette.

La matematica, e specificamente la geometria, procurò un motivo ancor più valido per promuovere una ricerca

nel campo della logica: la scuola pitagorica aveva scoperto l'esistenza delle grandezze incommensurabili. Essere commensurabile significa avere una misura comune, ed è ovvio che l'unità è la misura comune di tutti i numeri naturali (1, 2, 3 e così via). Un numero era inteso come una collezione composta di unità. L'unità aritmetica era pensata come indivisibile. Frazioni come $1/2$ o $2/3$ erano intese non come par-



Un'argomentazione semilegale si presentò quando Protagora si impegnò a insegnare retorica a Euatlo e a ricevere la seconda metà dell'onorario, solo dopo che Euatlo avesse vinto la sua prima causa. Quando Euatlo ritardò l'inizio della sua attività forense, Protagora fece causa per ottenere il suo onorario. I due avvocati discussero così.

ti di una unità ma sempre come una unità presa da due, due prese da tre e via dicendo. Su questa base doveva probabilmente sembrare ovvio che, date due lunghezze qualsiasi, dovesse esserci una unità geometrica così piccola da essere sottomultipla di ambedue le lunghezze. In tal caso in qualsiasi coppia di lunghezze queste starebbero fra loro in una determinata proporzione: la prima starebbe alla seconda come x sta a y , dove x è il numero delle unità geometriche della prima lunghezza e y quello della seconda.

Data la credenza quasi religiosa dei pitagorici nel numero come principio unificante dell'aritmetica, della geometria, della cosmologia e della filosofia, la loro scoperta delle grandezze incommensurabili deve essere stato uno choc, il primo dei molti conflitti tra religione e scienza nel mondo occidentale. La scoperta scaturì sotto forma di una dimostrazione della incommensurabilità del lato e della diagonale di un quadrato. Altri esempi di dimostrazioni che urtarono contro credenze radicate sulla base di ragioni diverse sollevarono la questione della natura esatta e della attendibilità delle dimostrazioni.

La logica aristotelica

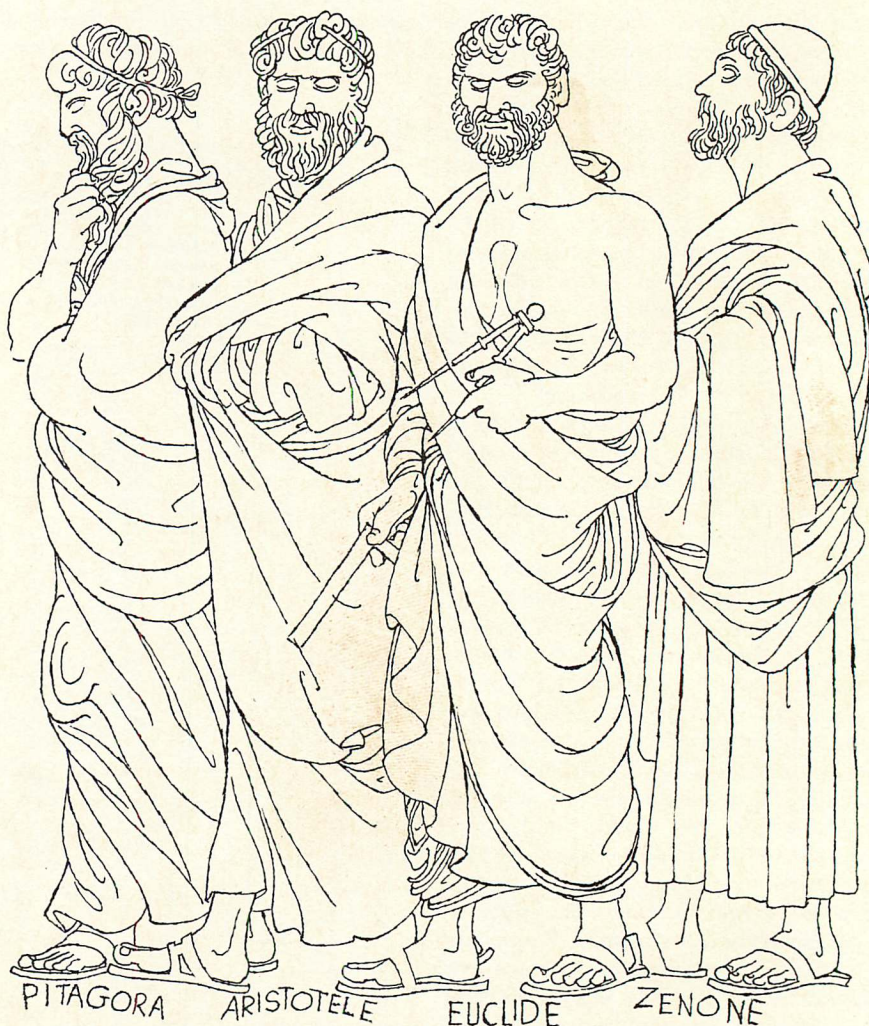
È stato quindi in risposta a considerazioni tanto filosofiche quanto matematiche che Aristotele creò la logica. Egli fissò la sua attenzione su quattro tipi generali di proposizione: la universale affermativa, la universale negativa, la particolare affermativa e la particolare negativa. Esempi di ciascuna sono, nell'ordine: « Tutti gli uomini sono bianchi », « Nessun uomo è bianco », « Qualche uomo è bianco » e « Qualche uomo non è bianco ». Aristotele si pose il problema di sapere quali ragionamenti sono validi se due proposizioni del genere con un termine in comune sono assunte come vere. Il risultato fu la famosa teoria del sillogismo. Per esempio, da « Tutti gli uomini sono bianchi » e « Qualche vittima del cancro non è bianca » c'è una sola conclusione valida da ricavare: « Qualche vittima del cancro non è uomo ». Le tre proposizioni insieme costituiscono un sillogismo. La logica era intesa come uno studio della connessione tra premesse e conclusione. In altre parole, asserire che un ragionamento è valido significava asserire non che le premesse sono vere, ma piuttosto che se le premesse sono vere, allora la conclusione è pure vera. Aristotele si spinse oltre e considerò una « dimostrazione » (questa era la sua terminologia) come un sillogismo valido con premesse necessariamente vere. Egli diede co-

sì una risposta al problema di sapere quando si ha dimostrazione in geometria: una proposizione è dimostrata in geometria se è la conclusione di una dimostrazione. Secondo Aristotele non tutta la conoscenza è dimostrativa. Bisogna partire da verità autoevidenti (che non sono oggetto di dimostrazione) e procedere sillogisticamente alle conclusioni. In questa concezione sta il germe del metodo assiomatico tradizionale.

L'esempio classico di questo metodo sono gli *Elementi* di Euclide, che datano a poco dopo l'epoca di Aristotele. Euclide apre la sua opera con una serie di 23 definizioni, 5 postulati e 5 « nozioni comuni », e da questi dimostra un gran numero di teoremi. Per le sue distinzioni tra queste quattro categorie (definizioni, postulati, nozioni comuni e teoremi) e per aver esibito un pari numero di dimostrazioni, Euclide occupa un posto importante nello sviluppo del metodo assiomatico (si veda l'illustrazione alle pagine 118 e 119). Il

metodo di Euclide e le sue applicazioni contenevano dei difetti, ma così sottili da sfuggire all'attenzione di tutti, o quasi, per 2000 anni.

Gli antichi greci fecero qualcos'altro, in geometria, che ebbe una funzione importante nello sviluppo della logica matematica: essi si posero problemi geometrici che loro stessi non erano in grado di risolvere. I più celebri erano quelli della duplicazione del cubo, della quadratura del cerchio e della trisezione di un angolo (arbitrario). In ciascuno dei casi si trattava di eseguire una costruzione esatta servendosi solo di riga e compasso. I geometri greci non riuscirono a risolvere questi problemi sotto tale condizione limitativa, anche se li risolsero usando metodi più complessi. A quanto pare alcuni dei greci antichi erano convinti della impossibilità di risolvere i problemi sotto la condizione che si è formulata, ma nessuno, per quanto ne so, concepì l'idea che l'impossibilità potesse essere oggetto di dimostrazione. Era lasciato



In questo disegno di Bernarda Bryson filosofi e matematici meditano sull'insolubilità. Aristotele (384-322 a.C.) fondò la logica allo scopo di affrontare problemi come quelli posti da Pitagora (580-500 a.C. ca.) e da Zenone (490-430 a.C. ca.); Euclide applicò

ai posteri il compito di decidere se l'insuccesso nel risolvere questi problemi sotto la condizione limitativa fosse dovuto a mancanza di abilità oppure alla natura dei problemi stessi. Fu solo nel XIX secolo che si dimostrò che tali costruzioni erano impossibili con riga e compasso solamente.

Sfide alla logica

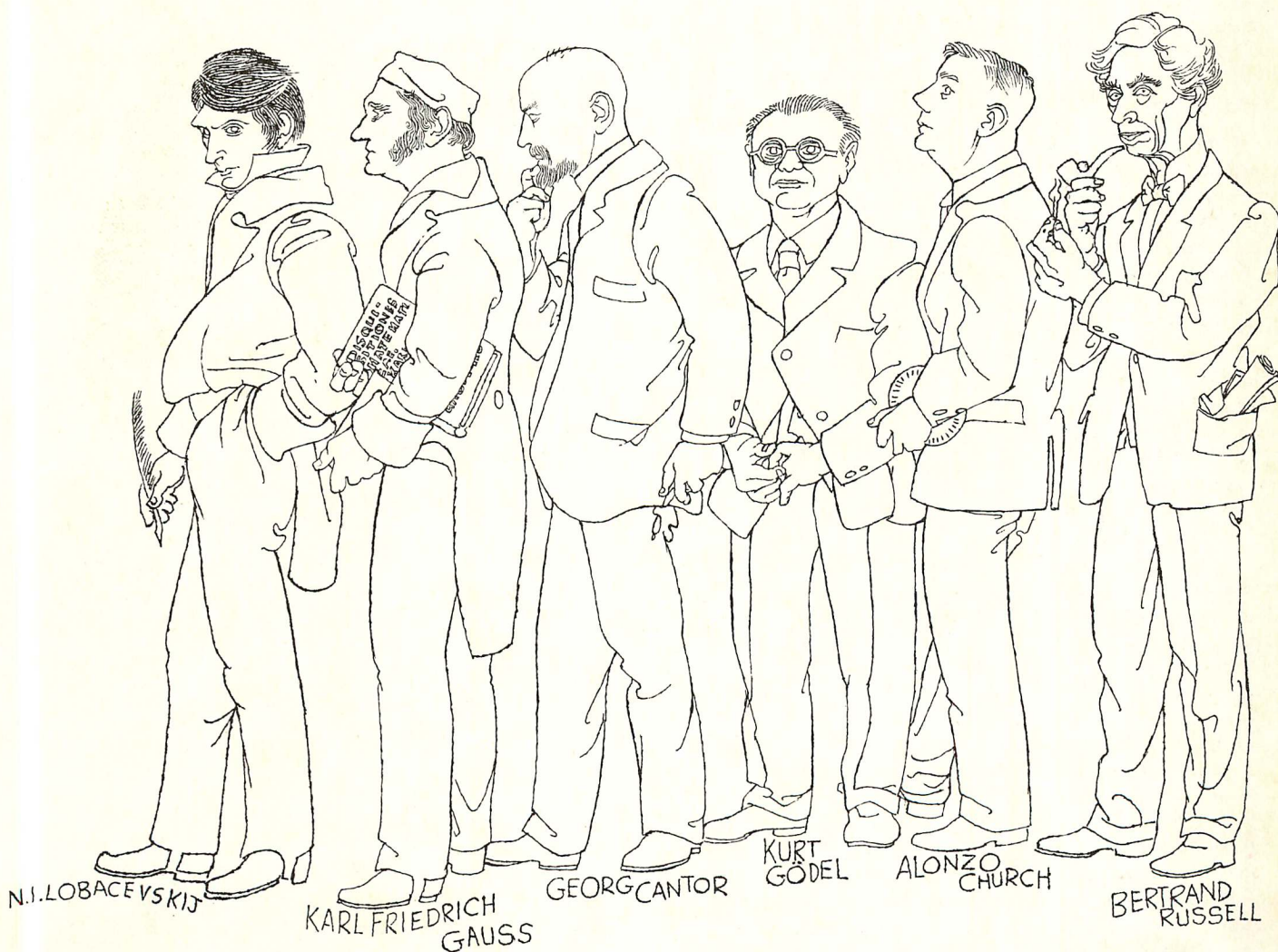
Circa nello stesso periodo divennero evidenti le manchevolezze della logica aristotelica. Lo divennero perché lo sviluppo della geometria non-euclidea, la teoria degli insiemi e la comparsa di nuovi e inquietanti paradossi portarono alla luce problemi di logica sistematici e fondamentali a cui né l'intuizione né la logica aristotelica potevano dare una risposta.

Consideriamo per prima cosa la geometria non-euclidea. Karl Friedrich Gauss e Nikolaj Ivanovic Lobacevskij, tra gli altri, scoprirono indipendentemente una geometria differente da quel-

la di Euclide, una geometria cioè in cui non vale la relazione pitagorica. Essi decisero che solo l'osservazione empirica e l'esperimento potevano determinare quale era la vera geometria, e con la scoperta di diverse geometrie questo atteggiamento si diffuse largamente. Esso portò a respingere la concezione del geometra come artefice di « dimostrazioni » nel senso aristotelico; il compito del geometra doveva essere quello di trovare connessioni logiche tra proposizioni geometriche, e non quello di determinare se erano vere le assunzioni di partenza. Era indispensabile tuttavia che gli assiomi fossero consistenti, perché una collezione inconsistente di enunciati non è possibile che sia vera. I matematici del secolo XIX svilupparono i mezzi per provare che se un sistema (poniamo la geometria euclidea) è consistente, allora un altro (poniamo qualche geometria non-euclidea) è pure consistente. Con ciò restava aperta la questione di come si potesse dimostrare in assoluto la consi-

stenza di un sistema. La geometria non-euclidea diveniva così essenziale allo sviluppo della matematica per almeno due riguardi. Essa centrava l'attenzione sul compito del matematico come scopritore di connessioni logiche, e sollevava il problema di come dimostrare la consistenza di una collezione di enunciati prescindendo dalla verità degli enunciati stessi.

La teoria degli insiemi fornì un'altra spinta alla crescita della logica matematica. Georg Cantor definì un insieme come una qualunque collezione costituente un tutto di oggetti definiti e distinti della nostra intuizione o del nostro pensiero. Questo concetto molto astratto di insieme è applicabile a generi diversi di cose come un insieme di scacchi o l'insieme dei numeri naturali o l'insieme di tutti gli insiemi di piatti. La teoria degli insiemi si dimostrò estremamente feconda. Cantor riuscì a dimostrare che il concetto di infinito è strutturato, e che proprio come gli insiemi finiti possono avere dimen-



il metodo assiomatico di Aristotele alla geometria. Gauss (1777-1855) e N.I. Lobacevskij (1793-1856) diedero sviluppo alla geometria non-euclidea. Cantor (1845-1918) è il fondatore della teoria

degli insiemi, all'interno della quale Russell (1872-1970) scoprì un paradosso. Kurt Gödel (1906-1978) e Church (1903) sono autori di teoremi in cui viene provata l'insolubilità di alcuni problemi.

sioni differenti, così ci possono essere insiemi infiniti con dimensioni differenti. Per esempio, proprio come l'insieme delle capitali di stato è grande quanto l'insieme degli stati ma più piccolo dell'insieme dei senatori degli USA così, ha dimostrato Cantor, l'insieme dei numeri naturali è grande quanto l'insieme dei numeri razionali (rapporti di interi) ma più piccolo dell'insieme dei numeri reali (decimali). Nel ragionare su insiemi infiniti bisogna essere cauti, in quanto la logica che vale per insiemi finiti non si applica necessariamente a quelli infiniti. Per esempio si può dimostrare che ci sono tanti numeri pari quanti sono i numeri naturali mettendoli in corrispondenza biunivoca (1-2, 2-4, 3-6, ..., $n-2n$, ...). Con ciò si contraddice l'assioma euclideo secondo il quale il tutto è maggiore della parte.

La teoria degli insiemi sembrò fondamentale a tal punto che per qualche tempo si ebbe la sensazione che tutta la matematica potesse essere fondata su di essa. Questa fiducia fu scossa dalla scoperta che nella teoria degli insie-

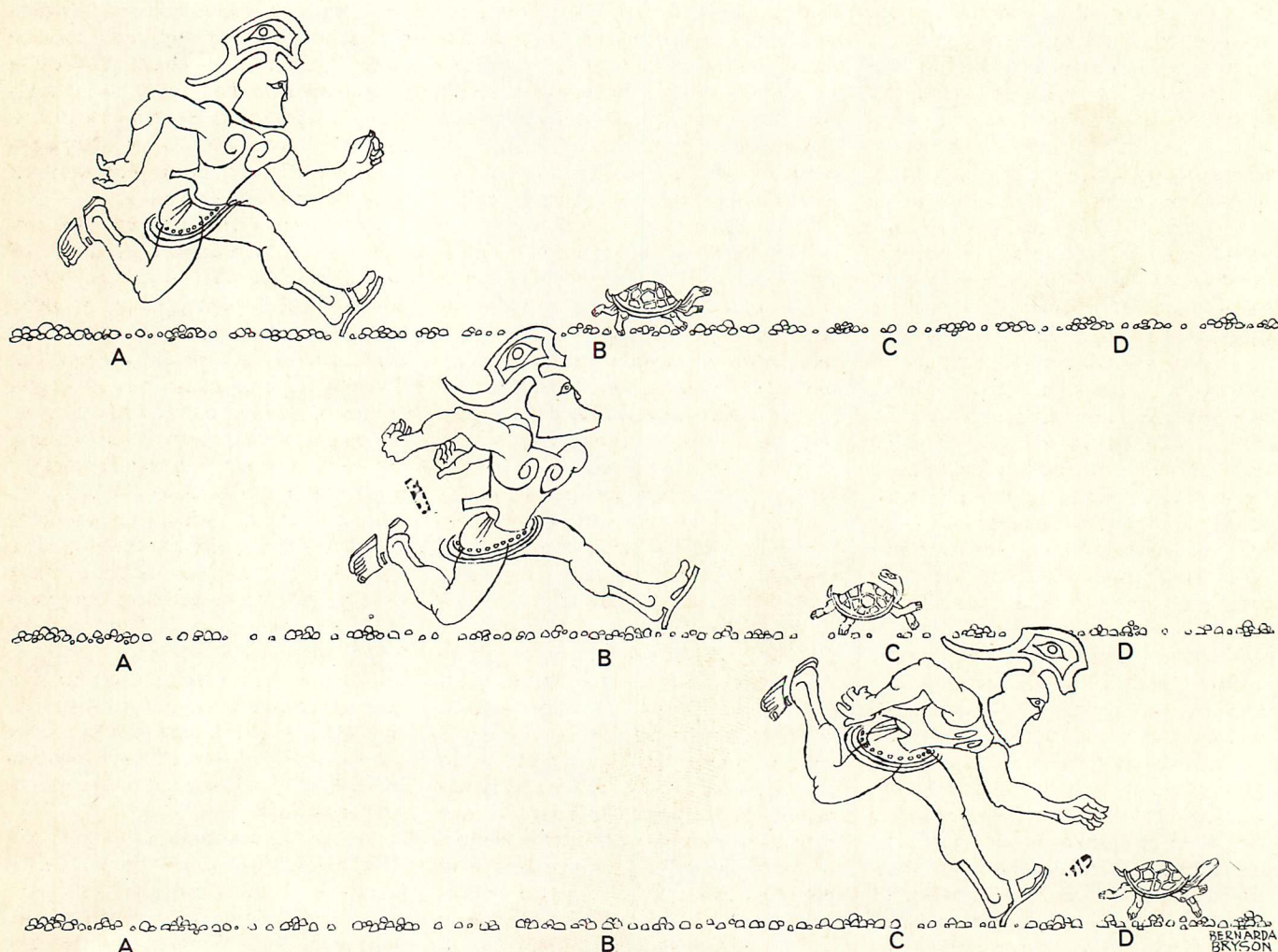
mi si presentavano delle contraddizioni logiche. La più semplice e famosa di queste è il paradosso, scoperto da Bertrand Russell nel 1902, relativo all'insieme di tutti gli insiemi che non sono membri di se stessi. Lo spirito dell'antinomia si può rendere con il paradosso del postino (si veda l'illustrazione a pagina 120), ma è importante distinguere la volgarizzazione dalla formulazione di Russell. La volgarizzazione non è in realtà un paradosso perché possiamo facilmente risolverla, per esempio negando che tale postino possa esistere (oppure, se ne esiste uno, che questo sia abitante del paese; o se lo è, che riceva posta, e così via). Lo stesso non succede con il paradosso di Russell: tutte le assunzioni di questo sembrano infatti necessariamente vere e tutti i modi per evitarlo sembrano avere conseguenze logicamente indesiderabili.

La teoria degli insiemi quindi è stata importante ai fini dello sviluppo della logica matematica da due punti di vista: ha introdotto nella pratica comune

ragionamenti relativi a collezioni infinite e ha provato che la « logica » che è adatta per collezioni finite non si applica necessariamente alle collezioni infinite. Per di più il paradosso di Russell, insieme a molti altri paradossi, sottolineava la necessità di una attenta indagine su quella che era la radice degli errori logici.

Logica matematica

Dunque, la rinascita delle indagini logiche che ebbe luogo nel XIX secolo e agli inizi del XX fu provocata da considerazioni analoghe a quelle che fecero nascere la logica aristotelica: cioè l'esigenza di affrontare dei paradossi e la questione della dimostrazione in matematica. Tanto la geometria non-euclidea quanto la teoria degli insiemi si esprimevano in affermazioni in contrasto con le credenze comunemente accettate. La logica aristotelica era inadeguata a risolvere questi problemi che andavano quindi affrontati secondo un approccio diverso. La nuo-

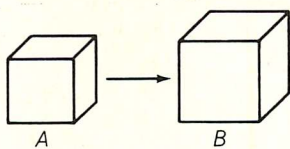


Un paradosso di Zenone afferma che il piè veloce Achille non può raggiungere la tartaruga. Alla tartaruga sia dato un vantaggio

AB. Nel tempo in cui Achille raggiunge B la tartaruga sarà in C; quando Achille sarà in C, la tartaruga sarà in D, e così via.

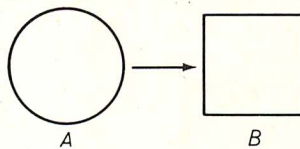
PROBLEMI
NON RISOLTI
DAGLI
ANTICHI
GRECI

Duplicazione del cubo: costruire un cubo con volume doppio di un cubo dato.



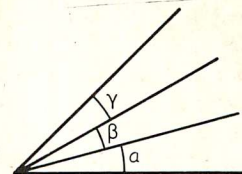
$$2 \times (\text{volume di } A) = \text{volume di } B$$

Quadratura del cerchio: costruire un quadrato con un'area uguale a quella di un cerchio dato.



$$\text{area di } A = \text{area di } B$$

Trisezione di un angolo (arbitrario): dividere un angolo arbitrario in tre parti uguali.



$$\alpha = \beta = \gamma$$

PROBLEMI
NON RISOLTI
DELLA
ARITMETICA

Problema dei numeri perfetti. Un numero perfetto è un numero naturale uguale alla somma dei suoi divisori (escluso se stesso). Per esempio 6, il più piccolo numero perfetto, è uguale a $1+2+3$. I sei numeri perfetti successivi sono, nell'ordine, 28, 496, 8128, 33 550 336, 8 589 869 056, 137 438 691 328. Problema: determinare se c'è un numero finito o infinito di numeri perfetti.

Congettura di Goldbach. Nel 1742 Goldbach avanzò l'ipotesi che ogni numero naturale pari maggiore di 2 fosse la somma di due numeri primi (un numero primo è un numero naturale maggiore di 1 divisibile solo per se stesso e per l'unità). Per esempio, $20 = 7 + 13$, $88 = 5 + 83$, $7000 = 3 + 6997$. Problema: dimostrare la congettura o trovare un numero naturale pari che non possa essere rappresentato come somma di due primi.

Ultimo teorema di Fermat. Fermat disse di avere una dimostrazione (che non è mai stata trovata) dell'enunciato che per tutti i numeri naturali x, y, z, n ,
$$x^{n+2} + y^{n+2} \neq z^{n+2}$$
Tutte le altre proposizioni matematiche di Fermat sono state finora o confermate o refutate. Di qui il nome di «ultimo teorema».

Sono qui elencati alcuni problemi non risolti di geometria e di aritmetica. Per gli antichi geometri greci il problema era quello di fare le costruzioni con riga e compasso solamente. Nel

secolo XIX si dimostrò che questi mezzi sono insufficienti. I problemi aritmetici, d'altro canto, possono essere non risolti perché tutti i mezzi disponibili sono inadeguati a risolverli.

va strategia passò attraverso l'impiego di metodi matematici e generò una nuova disciplina: la logica matematica.

Furono due i metodi matematici che trasformarono la logica in una nuova disciplina. Il primo fu il metodo algebrico, in cui le relazioni tra entità matematiche sono rispecchiate da relazioni intercorrenti tra i simboli di queste entità, cosicché si possono scoprire nuove relazioni tra entità matematiche manipolando i simboli di queste seconde certe regole. Il secondo fu il metodo assiomatico, non quello di Euclide ma una sua revisione; per distinguerlo dal metodo di Euclide va spesso sotto il nome di metodo assiomatico formale.

Il metodo assiomatico formale differisce da quello di Euclide per quanto riguarda la definizione e la esplicitazione di simboli, grammatica e logica (si veda l'illustrazione a pagina 118 e 119). Il discorso che Aristotele faceva per le proposizioni è applicabile anche alle definizioni: proprio come ci sono alcune proposizioni che non si possono dimostrare (altrimenti si cadrebbe in un regresso all'infinito), così ci devono essere alcuni termini che non si possono definire. La analogia si spinse oltre. Proprio come tutte le proposizioni devono essere dimostrate in termini di proposizioni non dimostrate, così tutti i termini devono essere definiti in termini di termini indefiniti. Le definizioni sono perciò superflue da un punto di vista teorico; esse servono solo a titolo di abbreviazioni.

Gli assiomi di un sistema formale possono contenere solo simboli non de-

finiti. Così, in un certo senso, si può dire che gli assiomi non vertono su nulla; il matematico che dimostra teoremi non fa riferimento ai significati possibili di termini indefiniti. In un altro senso invece gli assiomi vertono su qualunque insieme di oggetti tale da rendere veri gli assiomi stessi; per esempio gli assiomi di un particolare sistema formale potrebbero essere veri quando vengono applicati a numeri naturali. Il fatto che i termini indefiniti negli assiomi si applichino a molte cose differenti si potrebbe pensare come un difetto del metodo assiomatico formale. Al contrario esso costituisce la grande forza e versatilità del metodo, in quanto rende possibile ricavare una volta per tutte un largo corpo di teoremi veri per qualunque insieme di oggetti per cui sono veri gli assiomi.

Per farsi un'idea di questa caratteristica, si consideri un gioco con nove carte - l'asso e dal due al nove - disposte scoperte su un tavolo tra due giocatori. I giocatori a turno prendono una carta, e ciascuno dei giocatori cerca di essere il primo ad avere tre carte la cui somma è quindici. Il gioco è strutturalmente simile al «filetto». La analogia non è ovvia finché non si immagina di disporre le carte in una scacchiera per il gioco del «filetto» costruita in modo tale che, invece di prendere una carta, i giocatori possano segnare le carte alternativamente con X o O (si veda l'illustrazione a pagina 121). Allora è facile rendersi conto che per ciascuna partita con le carte c'è un corrispondente «filetto» e vi-

ceversa, ricordando che la scacchiera permette a chiunque conosca la strategia di un gioco di applicarla all'altro gioco. Questo è, in effetti, quanto fa il matematico che studia un sistema formale: egli studia la struttura di ambedue i «giochi» e deriva affermazioni descriventi la struttura che sono vere per entrambi i «giochi».

Un sistema formale ha qualche analogia con un linguaggio naturale. I suoi simboli corrispondono a lettere dell'alfabeto, segni di punteggiatura, numerali e così via; le regole di formazione corrispondono alle regole grammaticali di un linguaggio naturale; le regole di trasformazione corrispondono a varie operazioni che qualunque parlante può compiere sul linguaggio, come il convertire un enunciato da attivo in passivo. Per gli assiomi questa corrispondenza non è immediatamente identificabile nel linguaggio naturale, anche se si potrebbero considerare paragonabili a enunciati come «Tutto ciò che è, è» oppure « A è A ». Ci sono naturalmente differenze notevoli tra linguaggi naturali e sistemi formali, ma l'analogia è abbastanza stretta sicché quando i sistemi formali sono interpretati spesso si dà loro il nome di linguaggi artificiali.

Quando agli assiomi sono applicate delle regole di trasformazione, ciò che si ottiene è un teorema. L'esibizione della applicazione delle regole è una dimostrazione. Più esplicitamente, una dimostrazione è una sequenza finita di enunciati formali tali che ciascuno degli enunciati formali è un assioma o

segue da uno o più enunciati formali precedenti mediante applicazione di una regola di trasformazione. L'ultima riga della dimostrazione è un teorema. Le regole di trasformazione devono essere tali da rendere puramente meccanico il procedimento di determinare se una data sequenza di enunciati formali è una dimostrazione o no.

Proprietà dei sistemi formali

Si consideri un insieme di oggetti che siano indipendenti da qualche sistema formale del tipo visto. Noi diamo una interpretazione al sistema assegnando ai simboli elementi di questo insieme in maniera tale che gli assiomi siano veri per gli elementi dell'insieme di oggetti considerato. Questo insieme di oggetti in tal caso è detto un modello per il sistema. Il linguaggio (l'italiano) in cui noi poniamo e risolviamo problemi intorno al sistema è il metalinguaggio. Il sistema formale è il linguaggio-oggetto, in quanto è l'oggetto del nostro discorso. Per esempio, se diciamo « Nel sistema S possiamo dimostrare l'enunciato formale ' $1+1=2$ ' », stiamo facendo una affermazione nel metalinguaggio intorno a « $1+1=2$ », che è un enunciato espresso nel linguaggio-oggetto. Questo è il motivo per cui l'analisi delle proprietà formali dei sistemi ha preso il nome di metalogica.

Si consideri ora un sistema formale (chiamiamolo A +) costituito con l'obiettivo di formalizzare le regole per l'addizione dei numeri naturali. È possibile dimostrare che tale sistema è « consistente », cioè che è impossibile che qualunque suo teorema sia la negazione di qualche altro suo teorema. Per esempio, una sola dimostrazione ci assicura che se « $1+1=2$ » è un teorema, allora « $1+1 \neq 2$ » (« non è uguale a 2 ») non sarà un teorema. Si può anche provare che A + è « corretto », cioè che ogni teorema è vero quando è applicato ai numeri naturali. Inoltre, si potrà dimostrare che A + è « completo »: cioè a dire, ogni verità circa i numeri naturali esprimibile nel simbolismo del sistema è un teorema. Infine si può provare che il sistema ha una « procedura di decisione », cioè che esiste un metodo con il quale ogni problema esprimibile in A + può essere risolto in un numero finito di passi. Le procedure di decisione — talvolta chiamate algoritmi — sono comuni nella pratica matematica. Per esempio, la tecnica della divisione rappresenta una procedura di decisione per il predicato « x è divisibile per y », dove x ed y possono essere numeri qualsivoglia (un predicato è un enunciato aperto: un enunciato cioè che può essere comple-

tato sostituendo nomi alle sue variabili). Anche se questo predicato rappresenta un numero infinito di problemi, questi problemi sono essenzialmente simili, cosicché chiunque abbia dimestichezza con la divisione ha un algoritmo con cui è in grado di risolvere in un lasso di tempo finito qualunque problema di divisione sottopostogli, senza bisogno di particolari intuizioni o capacità personali.

L'esistenza di una procedura di decisione fa scemare l'interesse teoretico per una determinata area. Anche se è stato compiuto effettivamente solo un numero finito di divisioni, nessuno ha più interesse per le divisioni; e i problemi che restano in questo campo sono di indole pratica. L'esistenza di una procedura di decisione per un sistema formale risolve in un certo senso i problemi teorici per tutti i modelli del sistema. Il successo incontrato con un certo numero di sistemi portò verso il 1930 alla speranza, per non dire alla aspettativa, che sistemi espressivamente più potenti di A + si potessero dimostrare consistenti, corretti, completi e in possesso di una procedura di decisione.

In particolare si nutrì la speranza che un sistema « totale » per l'aritmetica — chiamiamolo A — si potesse dimostrare in possesso di queste proprietà (per « sistema totale » si intende un sistema esprimente l'addizione e la moltiplicazione di numeri naturali e tutte le operazioni definibili in termini di addizione e moltiplicazione, come la divisione). Tale dimostrazione avrebbe permesso ai matematici di risolvere i problemi teorici dell'aritmetica; ed era concepibile persino che tutti i problemi della matematica si potessero risolvere in modo analogo.

È stata una delle grandi conquiste dei logici matematici l'aver provato che a queste speranze non si poteva dare appagamento. Come esempio del modo in cui si è raggiunto questo risultato discuterò il teorema di incompletezza scoperto da Kurt Gödel nel 1931 e un teorema presentato da Alonzo Church nel 1936. Questi teoremi e alcuni altri dello stesso tipo ricevono allora il nome di teoremi limitativi, in quanto sembrano esprimere delle limitazioni delle facoltà umane.

Il teorema di Gödel

L'enunciazione dettagliata del teorema di incompletezza di Gödel è difficile, ed è necessario un bagaglio ragguardevole di conoscenze per intendere la dimostrazione nel suo complesso. Lo spirito della dimostrazione si può comunque rendere schematicamente (si

METODO ASSIOMATICO TRADIZIONALE

DESCRIZIONE

[**Simboli**]. Non elencati esplicitamente; presumibilmente comprendevano simboli ordinari presi dal linguaggio naturale.

[**Regole di formazione**]. Non enunciate esplicitamente; presumibilmente comprendevano le regole grammaticali ordinarie del linguaggio naturale.

Definizioni. Le definizioni erano intese come obiettive e vere e dovevano includere tutti i termini geometrici fondamentali.

Postulati. Intesi come verità e costruzioni autoevidenti che si applicano specificamente alla geometria.

Nozioni comuni. Verità autoevidenti che si potrebbero anche applicare ad altre scienze, per esempio dell'aritmetica.

[**Regole di inferenza**]. Non enunciate esplicitamente, presumibilmente comprendevano le inferenze intuitive ordinarie come « se p , e p implica q , allora q ».

Teoremi. Verità e costruzioni non evidenti che si dovevano dimostrare in modo tale da allontanare ogni dubbio legittimo (cioè autorizzando a comparire in una dimostrazione solo assunzioni contenute nelle definizioni, nei postulati, nelle nozioni comuni, e teoremi già dimostrati).

Il metodo assiomatico formale, insieme al metodo algebrico, ha trasformato la logica

veda l'illustrazione a pagina 122). Si consideri il « paradosso del mentitore », formulato dai greci antichi, che si può riformulare come il problema di decidere se l'enunciato: « Questo enunciato non è vero » è o no vero. Ne nasce una contraddizione tanto assumendo che l'enunciato del mentitore è vero quanto assumendo che l'enunciato è falso, da cui segue che qualunque sistema formale in cui esso è esprimibile è inconsistente. Si supponga ora che il potere espressivo di un certo sistema formale sia tale da rendere esprimibile non « vero » o « falso » ma « dimostrabile » o « indimostrabile ». L'analogo dell'enunciato del mentitore sarebbe allora « Questo enunciato non è dimostrabile ». Si chiami P questo enunciato. L'esistenza di P non rende inconsistente il sistema, ma produce qualcosa di sconcertante: P infatti è vero

(APPLICATO ALLA GEOMETRIA)	METODO ASSIOMATICO FORMALE (APPLICATO ALL'ARITMETICA)	
ESEMPI	DESCRIZIONE	ESEMPI
	Elenco di simboli. Deve comprendere tutti i simboli da impiegare nel sistema. Per esempio, x e y sono variabili individuali (cioè, possono stare per qualunque elemento del modello), mentre p e q sono variabili proposizionali (cioè possono stare per qualunque proposizione).	$+$ $=$ 0 $($ $,$ $)$ x y p q
	Regole di formazione. Indicano i modi in cui i simboli si possono legittimamente combinare per generare enunciati formali. Per esempio « $x = x$ » è un enunciato formale mentre « $x =$ » non lo è.	Se u e v sono variabili individuali, $u = v$ è un enunciato formale.
unto è ciò che non ha parti. nea è lunghezza senza larghezza. nea retta è quella che giace ugualmente spetto ai suoi punti.	[Definizioni]. Le definizioni sono eliminate in quanto da un punto di vista teorico non assolvono alcuna funzione. La sola utilità che presentano è tachigrafica.	
a qualsiasi punto a ogni punto si può ndurre una linea retta. può descrivere un cerchio con qualunque entro e qualunque raggio utti gli angoli retti sono uguali tra loro.	Assiomi (aritmetici). Un elenco di enunciati formali dai quali prendiamo le mosse. Questi assiomi sono destinati a essere interpretati come enunciati aritmetici, anche se non si fa riferimento a questo nello sviluppare il sistema.	$x = x$ $(x + 0) = x$ $(x \cdot 0) = 0$
ose che sono uguali alla stessa cosa sono uali anche tra loro. e cose uguali sono aggiunte a cose uguali, e totalità sono uguali. tutto è maggiore delle parti.	Assiomi (logici). Una volta interpretati, diventano assunzioni logiche. L'esempio a destra significa: «se p è vero, allora, se q è vero, p è vero».	$(p \supset (q \supset p))$
	Regole di trasformazione. Regole con cui manipoliamo enunciati formali per produrre nuovi enunciati formali. Per esempio, da « $x = x$ » possiamo inferire « $0 = 0$ ».	Se r è un enunciato formale contenente una variabile individuale, si inferisce l'enunciato formale con la variabile individuale rimpiazzata da 0.
u una linea retta finita è possibile costruire n triangolo equilatero. quadrato dell'ipotenusa di un triangolo ettangolo è uguale alla somma ei quadrati dei cateti.	Teoremi. Le ultime righe delle dimostrazioni. Ogni riga di una dimostrazione deve o essere un assioma o seguire da una o più righe precedenti per applicazione di una regola di trasformazione.	$x = x$ $0 = 0$

in logica matematica. Il metodo assiomatico formale applicato all'aritmetica si è sviluppato a partire dal metodo assiomatico tra-

dizionale della geometria, con cui viene qui confrontato. Sono descritti gli elementi e di ciascun elemento sono dati esempi.

se e solo se P non è dimostrabile. Quindi possiamo concludere che se abbiamo un sistema abbastanza ricco per esprimere P , allora viene a saltare la comoda relazione tra verità e dimostrabilità che si tenta di raggiungere in un sistema formale, quella cioè per cui l'insieme di enunciati veri sotto qualunque interpretazione tale da render veri gli assiomi risulta identico all'insieme degli enunciati dimostrabili.

Il teorema di incompletezza di Gödel afferma, in parole povere, che per qualunque dei sistemi formali noti per l'aritmetica ci sono enunciati formali analoghi a P , e cioè o il sistema è scorretto (dimostra enunciati falsi) o è incompleto (contiene verità non dimostrabili nel sistema). Proprio come il « filetto » e il gioco di carte sono strutturalmente identici, così Gödel ha potuto dimostrare che alcuni enunciati

aritmetici sono strutturalmente identici a P . Gödel raggiunse questo risultato sviluppando un codice (analogo alla scacchiera per il gioco del « filetto ») per mezzo del quale è riuscito a provare che affermazioni e ragionamenti intorno a un sistema (nel metalinguaggio di questo) potevano essere rispecchiati nel sistema (nel suo linguaggio-oggetto). Non è essenziale ai nostri fini aver dimestichezza con i dettagli di questo codice. Basti dire che Gödel assegnò un numero unico a ciascuno dei simboli, enunciati e dimostrazioni in modo tale che questi si potevano poi ordinare e si poteva parlare, per esempio, del simbolo n . 1 o della dimostrazione n . 15.

Consideriamo ora l'argomentazione di Gödel in modo un po' più dettagliato. Prendiamo come G l'enunciato « per qualunque numero naturale x , la dimostrazione numero x non è una di-

mostrazione in A dell'enunciato numero n ». Assumeremo quanto ha provato Gödel: e cioè che n può essere scelto in modo tale da essere il numero dello stesso enunciato G . Se esaminiamo da vicino la situazione, ci imbattiamo nel curioso fenomeno della ω -incompletezza (si veda l'illustrazione a pagina 122). Un sistema formale è ω -incompleto se contiene una generalizzazione di cui è possibile dimostrare ogni esempio numerico senza che sia possibile dimostrare la generalizzazione stessa. Ed è proprio relativamente alla ω -incompletezza che la logica matematica può gettare luce su qualcuno dei problemi irrisolti dell'aritmetica.

Per esempio, nel caso della congettura di Goldbach (si veda l'illustrazione a pagina 117) per ogni dato numero è un procedimento completamente meccanico il determinare se il nume-

ro dato è la somma di due numeri primi. Si supponga, tuttavia, che la congettura di Goldbach, come *G*, illustri la ω -incompletezza in *A*. Potremmo allora dimostrare ciascuno dei seguenti enunciati senza poter dimostrare l'enunciato generale: « 4 è la somma di due numeri primi, 6 è la somma di due numeri primi, 8 è la somma di due numeri primi... ». In altri termini, la ragione per cui la congettura rimane indecisa può essere che non c'è nes-

suna dimostrazione o refutazione a partire dalle assunzioni che i matematici hanno fatto effettivamente per cercare di deciderla. Considerando la cosa sotto un altro aspetto, osserviamo che una tecnica standard in matematica è quella di dividere una dimostrazione in casi. Se la congettura di Goldbach è un esempio di ω -incompletezza, abbiamo un problema che si polverizza in un numero infinito di casi!

Se d'altro canto potessimo dimostra-

re in modo matematico che la congettura di Goldbach è un esempio di ω -incompletezza in *A*, dimostreremo (come nel caso di *G*) che essa è vera per i numeri naturali. Dimostrando infatti che essa esemplifica la ω -incompletezza in *A*, proveremmo che è impossibile trovare un numero pari maggiore di 2 che non sia la somma di due primi, il che equivale a provare che ogni numero pari maggiore di 2 è la somma di due numeri primi. Ciò che verrebbe a significare una dimostrazione siffatta è che qualunque dimostrazione informale della congettura di Goldbach dovrebbe usare metodi che trascendono quelli dell'aritmetica, per esempio i metodi del calcolo infinitesimale. Affermazioni analoghe si potrebbero fare per l'ultimo teorema di Fermat (si veda l'illustrazione a pagina 117). Bisogna sottolineare che noi non sappiamo se la congettura di Goldbach o l'ultimo teorema di Fermat siano esempi di ω -incompletezza in *A*, ma ciò non è da escludersi allo stato attuale delle conoscenze.

Il teorema di Church è pure volto a illustrare le difficoltà in cui possiamo imbatterci relativamente ai problemi irrisolti dell'aritmetica. Il teorema afferma, in termini approssimativi, che non c'è nessuna procedura di decisione per qualunque sistema formale di aritmetica. Una analogia può facilitare la comprensione. È facile bisecare un angolo arbitrario in geometria piana; non è necessario conoscere tutto ciò che si può fare con riga e compasso. Non si deve fare altro che costruire una bisettrice. Al contrario, per provare che è impossibile, in generale, trisecare un angolo, è necessaria una definizione precisa di tutte le costruzioni possibili con riga e compasso. Poiché i greci antichi mancavano di tale definizione, non erano in condizione di poter dimostrare la impossibilità della trisezione. Analogamente, procedure di decisione si sono applicate per migliaia di anni, ma fino al 1930 nessuno ha dato una analisi matematica esatta di questa nozione. Servendosi appunto di una tale analisi rigorosa, Church è riuscito a provare che non c'è nessun algoritmo per decidere il predicato « L'enunciato *x* esprime una verità dell'aritmetica » mentre c'è un algoritmo basato sulla divisione per decidere il predicato « *x* è divisibile per *y* ».

Church ha dimostrato il suo teorema con un ragionamento che ha somiglianze tanto con quello di Gödel quanto con quello di Cantor, mirante a provare che sono più i numeri reali che non i numeri naturali. La prova di Church dimostrava che l'assunzione dell'esistenza di una procedura di de-



Il paradosso del postino è simile al paradosso di Bertrand Russell riguardante gli insiemi che non contengono se stessi. A differenza del paradosso di Russell, tuttavia, questa versione può essere risolta. Per esempio, si può negare che il postino esista.

cisione per l'aritmetica porta a una assurdità. Ciò significa non solo che non abbiamo ancora trovato questo algoritmo, ma che non esiste nessun algoritmo che permetta sempre di individuare la verità aritmetica. In altri termini proprio come non c'è nessun metodo nel « filetto » tale da garantire una vincita contro tutte le strategie, così non c'è nessun metodo in aritmetica che garantisca una dimostrazione per tutte le verità. Ne segue che ci saranno sempre verità aritmetiche per cui non saranno adeguati i metodi attuali e per cui si renderà sempre necessaria la creazione di nuovi metodi.

I teoremi limitativi di Gödel e Church sembrano avere profonde conseguenze filosofiche, ma qualcuno potrebbe obiettare: « Nego che partendo da questi teoremi si possano trarre conclusioni filosofiche. Gödel e Church hanno provato che tutti i sistemi formali noti per l'aritmetica sono incompleti e mancano di una procedura di decisione, ma questo è poco più significativo del fatto che non si può quadrare il cerchio o trisecare un angolo. È perfettamente possibile fare queste costruzioni, ma non con riga e compasso solamente. Analogamente, tutto ciò a cui pervengono i teoremi di Gödel e Church è che, dati i mezzi che hanno scelto (proprio allo stesso modo in cui Euclide scelse riga e compasso), segue inevitabilmente la loro conclusione. Non c'è nessuna portata filosofica nei loro teoremi: bisogna semplicemente cercare altri mezzi ».

Questa obiezione avrebbe una portata significativa se non fosse per una sola circostanza: che non sembra esserci nessun altro mezzo. A questa profonda conclusione si è giunti attraverso il ragionamento che segue. Euclide aveva studiato triangoli ideali nello sforzo di cogliere il concetto di spazio. In modo analogo Emil Post negli anni trenta intraprese una analisi del modo in cui opera un ideale calcolatore umano quando appunto esegue calcoli, con il proposito di cogliere il concetto di procedura di decisione. Indipendentemente, e circa nello stesso periodo, A. M. Turing intraprese una analisi parallela, e guidata dallo stesso obiettivo, di una macchina calcolatrice ideale. Si provò che le due analisi risultavano equivalenti. Church suggerì che qualunque effettiva computazione operata da uomini o macchine poteva essere duplicata dall'uomo ideale o dalla macchina ideale. La tesi di Church è empirica, ma l'evidenza a suo favore è schiacciante. Se la accettiamo, gli aspetti limitativi dei teoremi di Gödel e di Church si possono rendere espliciti (si veda l'illustrazione a pagina 123).

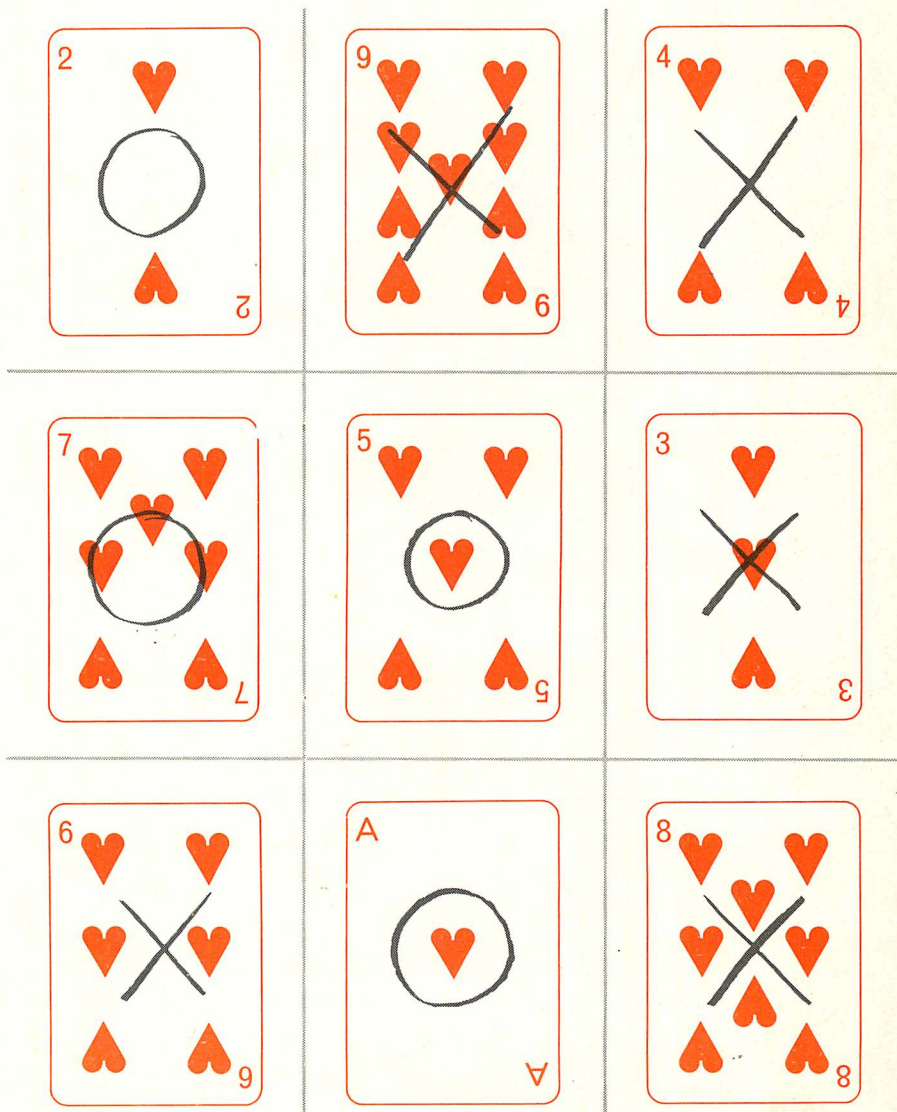
Problemi senza risposta

Il cambiamento fondamentale che i teoremi limitativi spingevano ad apportare a tutte le precedenti teorie sulla natura della matematica era il riconoscimento che ci sono problemi di questa area del sapere destinati a restare senza risposta. Precedentemente si era convinti del fatto che se un interrogativo si poteva esprimere con esattezza, quell'interrogativo doveva avere una risposta. Ora invece ci si accorgeva che forse alcuni quesiti precisi non hanno risposte precise. A titolo di analogia si pensi a un soggetto, poniamo una lampadina. Se voi allora chiedete: « È parzialmente fatta di vetro? » la risposta sarà probabilmente sì; se voi chiedete « È parzialmente di sughero? » la risposta sarà probabilmente no. Ma se chiedete « Pesa esattamente 93 grammi? » la domanda probabilmente è destinata a restare senza risposta. La real-

tà verso cui la domanda è diretta è per qualche riguardo indeterminata. Tale indeterminatezza è caratteristica dei prodotti dell'immaginazione.

A confronto con le creazioni dell'immaginazione, la realtà fisica è determinata, eppure i risultati della teoria quantistica ci dicono che la realtà fisica è pure in un certo modo indeterminata. Proprio come l'indeterminazione, prima d'allora considerata peculiare della immaginazione creativa, è stata individuata anche nel mondo fisico con la scoperta della teoria quantistica, così l'indeterminazione è stata ritrovata in matematica con la scoperta dei teoremi limitativi. La portata rivoluzionaria di questa scoperta è paragonabile a quella della scoperta pitagorica della incommensurabilità del lato e della diagonale del quadrato. Essa capovolge concezioni filosofiche che per un lungo periodo di tempo sono state fondamentali.

Il suo peso filosofico si può sottoli-



La strategia per un gioco di carte in cui lo scopo è di essere i primi a prendere tre carte la cui somma dà 15 è identica alla strategia per il « filetto ». I matematici cercano strutture identiche in campi della matematica che sono manifestamente diversi.

neare con un altro paragone con la geometria euclidea. Limitando le nostre assunzioni a quelle di Euclide, non è possibile dimostrare l'esistenza di un quadrato la cui area è uguale a quella di un dato cerchio; sulla base di queste limitazioni concettuali l'esistenza di tale quadrato è accidentale: non è né necessaria né impossibile. L'indeterminazione è dovuta a un rifiuto di fare assunzioni aggiuntive. L'indeterminazione in aritmetica, invece, è dovuta a una *incapacità* sistematica di allargare le nostre assunzioni in modo da renderle adeguate a dimostrare tutte le verità dell'aritmetica. Dal punto di vista di un qualsiasi essere umano particolare devono esserci verità accidentali dell'aritmetica, ossia affermazioni che pur essendo vere sono tali che niente di ciò che egli potrebbe assumere gli permetterebbe di dimostrarle vere; esse non

sarebbero cioè né necessarie né impossibili. Ciò che rappresentano i teoremi limitativi, allora, è la scoperta di una struttura astratta per la quale non riesce possibile, a nessun essere umano, fare sistematicamente assunzioni complete e corrette. Proprio come le nostre percezioni sensoriali hanno limiti che si possono estendere con certe tecniche (per esempio con un microscopio) ma che non si possono eliminare, così le nostre concezioni astratte sono limitate e i metodi della matematica miranti a estenderle (per esempio l'induzione matematica) conseguono tutt'al più un risultato parziale. I nostri poteri di discriminazione concettuale hanno dei limiti proprio come i nostri poteri di discriminazione percettiva.

Non sembra esserci modo di evitare queste conclusioni: ci sono anzi conseguenze ancora più profonde. Per chia-

rare questo punto bisogna fare una distinzione tra indecidibile e insolubile. Un enunciato è indecidibile in un dato sistema se né l'enunciato stesso né la sua negazione sono dimostrabili nel sistema. Il concetto di indecidibilità è relativizzato a un sistema, nel senso che ciò che è indecidibile in un sistema può essere decidibile in un altro. Per esempio, G è indecidibile in A ma è, come abbiamo visto, decidibile (e vero) nel sistema più grande che include A e pure l'argomentazione che G è indecidibile in A . Per converso, il concetto di insolubilità è assoluto. Un predicato di numeri naturali è ricorsivamente solubile se riceve una soluzione da parte dell'ideale calcolatore umano (o dell'ideale calcolatore meccanico) di cui si è detto; in caso contrario è ricorsivamente insolubile. Per esempio il predicato numerico « x è divisibile per y » è

PARADOSSO DEL MENTITORE	RAGIONAMENTO DI GÖDEL	RAGIONAMENTO DI GÖDEL (PIÙ DETTAGLIATO)
Questo enunciato non è vero	Questo enunciato non è dimostrabile.	Per qualunque numero naturale x , la dimostrazione x non è una dimostrazione in A dell'enunciato n .
Sia L il nome dell'enunciato scritto sopra.	Sia P il nome dell'enunciato scritto sopra.	Sia G il nome dell'enunciato scritto sopra (n è un numero naturale scelto in modo che l'enunciato n sia G stesso).
Si supponga che L sia vero. Allora ciò che dice L è corretto e L dice di non essere vero. Quindi se L è vero L non è vero.	Si supponga che P sia dimostrabile. In tal caso, dal momento che afferma di non essere dimostrabile, P non è vero. Dunque se P è dimostrabile non è vero.	Si supponga che G sia dimostrabile nel sistema formale A . Assumendo che A è corretto, ciò che G esprime è vero. Tuttavia, dal momento che G non è dimostrabile, c'è una inconsistenza. Quindi se A è consistente, G non è dimostrabile. Ciò significa che ognuna delle seguenti affermazioni è vera: La dimostrazione 1 non è una dimostrazione in A dell'enunciato n . La dimostrazione 2 non è una dimostrazione in A dell'enunciato n . La dimostrazione 3 non è una dimostrazione in A dell'enunciato n Si può provare che ognuna di queste è anche dimostrabile. Ma « per qualunque numero naturale x , la dimostrazione x non è una dimostrazione in A dell'enunciato n » non è dimostrabile. Quindi se A è consistente è dimostrabile ciascuna esemplificazione numerica di G ma non G stesso. Qualunque sistema formale che contiene un enunciato (formale) tale che è dimostrabile ciascuna delle esemplificazioni numeriche ma non l'enunciato generale è detto ω -incompleto. Quindi se A è consistente G non è dimostrabile e A è ω -incompleto.
Si supponga che L non sia vero. In tal caso ciò che L dice è scorretto e L deve essere vero. Quindi se L non è vero L è vero.	Si supponga che P non sia dimostrabile. In tal caso, poiché afferma di non essere dimostrabile, è vero. Quindi se P non è dimostrabile P è vero.	Si supponga che nel sistema formale A sia dimostrabile non- G (non- G = è falso che G = negazione di G). In tal caso, poiché sappiamo che ogni esempio numerico di G è dimostrabile, abbiamo un sistema in cui esiste un enunciato (formale) ogni esemplificazione del quale è dimostrabile ma la negazione del quale è pure dimostrabile. Poiché ogni esemplificazione di G è vera (purché A sia consistente), la dimostrazione di non- G è una dimostrazione di una asserzione falsa. Quindi, se A è corretto, non- G non è dimostrabile.
Perciò L è vero se e solo se L non è vero. Conclusione: il sistema in cui è esprimibile L è inconsistente.	Perciò P è vero se e solo se P non è dimostrabile. Conclusione: il sistema in cui è esprimibile P è o scorretto (tale da dimostrare enunciati falsi) o incompleto (cioè tale da contenere verità non dimostrabili nel sistema).	Perciò G è vero dei numeri naturali se e solo se ogni esempio numerico di G è dimostrabile in A ma G non è dimostrabile in A . Conclusione: il sistema in cui G è esprimibile (cioè A) è o scorretto (tale da dimostrare enunciati falsi) o incompleto (tale da contenere verità non dimostrabili nel sistema).

La dimostrazione data da Gödel del suo teorema di incompletezza è qui suggerita dal «paradosso del mentitore». Gödel considera la dimostrabilità invece della verità e ha provato che un sistema formale per l'aritmetica è o scorretto o incompleto.

TEOREMA	LINGUAGGIO DELLA PSICOLOGIA	LINGUAGGIO DELLA FISICA
TEOREMA DI INCOMPLETEZZA DI GÖDEL	Non c'è nessun calcolatore umano consistente in grado di formulare un programma che, eseguito, produrrebbe tutti e solo gli enunciati veri dell'aritmetica.	Non c'è nessun calcolatore artificiale consistente che si possa programmare così da produrre tutti e solo gli enunciati veri dell'aritmetica.
TEOREMA DI CHURCH	Esiste un insieme di problemi dell'aritmetica che nessun « calcolatore » umano consistente è in grado di risolvere.	Esiste un insieme di problemi dell'aritmetica che nessun calcolatore artificiale consistente può essere programmato a risolvere.

L'aspetto limitativo dei teoremi di Gödel e Church diventa evidente alla luce della tesi di Church, secondo la quale qualun-

que computazione reale, eseguita da un uomo o da una macchina, può essere duplicata da un uomo o da una macchina ideali.

ricorsivamente solubile. L'enunciato « x esprime una verità dell'aritmetica » invece non è ricorsivamente solubile. Ciò non significa che per qualche valore dato del predicato il problema non si possa risolvere. Per esempio, se « $1+1=2$ » è l'enunciato 2467 in A , allora l'enunciato « L'enunciato 2467 esprime una verità dell'aritmetica » è tanto vero quanto dimostrabile. L'insolubilità significa che non c'è nessuna tecnica di nessun tipo in grado di funzionare *sempre* in modo corretto allo stesso modo in cui funziona sempre la divisione. Sorge naturalmente il problema consistente nel chiedersi se qualche noto problema della matematica elementare è, o no, ricorsivamente insolubile. Non c'è nessun motivo apparente perché la risposta possa essere negativa. Per esempio, il predicato « x è divisore di un numero perfetto » potrebbe benissimo essere ricorsivamente insolubile e perciò, sulla base della tesi di Church, insolubile nel senso ordinario, cioè a dire insolubile da parte di qualunque uomo o calcolatore.

Problemi pratici

Come se non bastasse, ci sono ancora altre possibili difficoltà. Ho descritto finora ciò che si potrebbe chiamare l'indecidibilità teorica e l'insolubilità teorica. Bisogna aggiungere che esistono anche fenomeni di indecidibilità pratica e di insolubilità pratica. Un enunciato formale è praticamente indecidibile se, pur essendo teoricamente decidibile, la dimostrazione è talmente lunga da non potersi eseguire per pura impossibilità fisica. Per esempio, può darsi che la congettura di Goldbach sia vera, ma che le dimostrazioni più brevi entro un dato sistema richiedano più litri di inchiostro di quanti non siano gli atomi dell'universo. Oppure, se la congettura è falsa, può darsi che il primo numero pari che non è la somma di due numeri primi sia grande al punto che la possibilità di scoprirlo risulta

virtualmente nulla. Oppure, di nuovo, un predicato è praticamente insolubile se, pur essendo teoricamente solubile, è fisicamente impossibile utilizzare le tecniche di soluzione. Per esempio, le tecniche per la soluzione del predicato « x è divisore di un numero perfetto » possono essere così complicate da superare la capacità di qualsiasi essere umano o calcolatore concepibile. La capacità di uomini e calcolatori di affrontare certe complicazioni ha un limite, mentre l'aritmetica presenta problemi le cui soluzioni sono molto complicate.

Si potrebbe chiedere: « Ammettiamo che non si possa dimostrare la completezza dei sistemi formali dell'aritmetica e che non ci sia nessuna procedura di decisione. Ammettiamo inoltre che ci possano essere seri problemi di indecidibilità pratica e di insolubilità pratica per l'aritmetica. Possiamo almeno dimostrare che questi sistemi dell'aritmetica sono consistenti? ». Al giorno d'oggi la risposta è a metà strada tra sì e no. Ci sono dimostrazioni di consistenza basate su metodi talmente complicati da essere più dubbi che non la consistenza stessa dei sistemi. D'altro canto, se ci limitiamo a metodi non più complicati dall'aritmetica stessa, allora incontestabilmente possiamo dimostrare un enunciato formale esprimere la consistenza. Tuttavia per essere sicuri che l'enunciato formale esprima la consistenza, dobbiamo assumere la consistenza dell'aritmetica! Come Gödel ha sottolineato, la questione della consistenza riposa sulla nostra capacità di « sostituire [gli assiomi della matematica classica] relativi a entità astratte di una oggettiva sfera platonica con intuizioni di date operazioni della mente ».

La verità aritmetica così ha una natura « ideale » o « prospettiva » (i termini sono stati proposti da John Myhill). L'analogia con un gioco può essere utile per illustrare questa nozione. Accade spesso che si inventi un gioco e che si pongano regole che lo definiscono, ma che più tardi sorga una cir-

costanza per la quale le regole non danno nessuna istruzione. A questo punto deve essere presa una decisione su ciò che dovrà essere la regola relativa a questa circostanza. Una decisione si potrebbe prendere con un criterio di eleganza, o con uno che assicuri un più intenso divertimento dello spettatore, o che incrementi il rischio, o sortisca qualche altro effetto. La decisione non si prenderà comunque sulla base delle regole del gioco proprio in quanto queste sono definite in modo incompleto. Ora parte della forza dei teoremi limitativi sta nel fatto che le regole per mezzo delle quali definiamo e scopriamo la verità aritmetica non soltanto sono, ma devono essere definite in modo incompleto. Siamo perciò costretti a definire storicamente la nozione di verità aritmetica; essa non può essere esplicitata una volta per tutte ma deve essere ridefinita continuamente. Abbiamo visto come tanto Gauss quanto Lobacevskij siano giunti alla conclusione che i problemi della verità nella geometria non-euclidea li spingevano a procedere oltre i dati della pura geometria. In modo analogo dobbiamo oltrepassare i dati dell'aritmetica se vogliamo definire la verità aritmetica. L'uomo ha creato un gioco dell'aritmetica che è incompleto, a quanto pare a causa della incommensurabilità tra i suoi ideali e le sue possibilità.

È certo una tentazione molto forte quella di dire che se l'aritmetica può essere un gioco, è un gioco nel quale impariamo qualcosa intorno alla realtà — la realtà astratta. In tal caso significa che c'è una incommensurabilità tra la realtà astratta e la nostra capacità di intenderla pienamente. Se il predicato « L'enunciato x esprime una verità dell'aritmetica » fosse solubile, l'aritmetica diventerebbe poco interessante dal punto di vista teorico. Non essendo solubile, l'aritmetica continuerà a suscitare il nostro interesse proprio per il fatto che la sua indagine esige il contributo della creatività umana.

III

MATEMATICA E REALTÀ

La macchina di Turing e la questione da essa sollevata: può una macchina pensare?

di Martin Gardner

«Ci fu un tempo in cui dev'essere sembrato molto improbabile che le macchine potessero imparare a esprimere i loro desideri con parole che giungessero all'orecchio umano; non potremmo perciò immaginare che verrà un giorno in cui l'orecchio umano non sarà più necessario, perché i desideri della macchina saranno percepiti dai delicati meccanismi della macchina stessa, quando il loro linguaggio si sarà sviluppato dal pianto degli animali a una parola altrettanto complicata quanto la nostra? »

— SAMUEL BUTLER, *Erewhon*.

Alan Mathison Turing, un matematico inglese morto nel 1954 all'età di 42 anni, fu uno dei più creativi fra i primi studiosi di calcolatori. Oggi è ben noto soprattutto per il suo concetto di «macchina di Turing». Descriveremo brevemente questo tipo di macchina, per poi passare a una delle sue idee meno conosciute, il «gioco di Turing», che conduce a profonde controversie filosofiche.

La macchina di Turing è una «scatola nera» (un congegno cioè dai meccanismi non specificati) capace di osservare un nastro di lunghezza infinita diviso in celle quadrate. La scatola può avere qualsiasi numero finito di stati. Una parte finita del nastro consiste di celle non vuote, su ognuna delle quali è impresso uno qualsiasi di un numero finito di simboli. Quando la scatola esamina una particolare cella, essa può lasciare il simbolo inalterato, cancellarlo, sostituendolo con un altro simbolo, o stampare un simbolo in una cella vuota. Il nastro può poi venir spostato di una cella verso destra o verso sinistra, oppure rimanere fermo. La «scatola nera» a sua volta può rimanere nello stesso stato, o passare a un altro.

Una tabella di regole stabilisce il comportamento della macchina per ogni possibile combinazione di stati e simboli; una tabella di questo tipo definisce completamente una particolare macchina di Turing. Esiste una infinità numerabile (alef con zero) di macchine di Turing, ognuna progettata per un compito specifico; per ogni compito diverso la struttura della macchina può variare molto in simboli, stati e regole.

Un buon sistema per capire in che consiste una macchina di Turing, è quello di costruirne una, sia pure del tipo più rudimentale. Otto celle sul nastro di carta sono contrassegnate con $1111 + 111$ per indicare la somma di 4 e 3 in un sistema «a base uno», sistema in cui un numero intero n è rappresentato da n «uno». Per costruire la macchina disegniamo un quadrato (la scatola nera) e tagliamo su di esso due fessure in modo da potervi inserire il nastro come indicato nella figura. Sistemiamo il nastro in modo che sia visibile il primo 1. La tavola sotto l'illustrazione della pagina a fronte dà le regole necessarie.

Cominciamo supponendo che la macchina si trovi nello stato A ; consultiamo la tabella circa la combinazione del simbolo 1 con lo stato A e facciamo quanto indicato; cancelliamo cioè l'1, spostiamo il nastro verso sinistra (così che la scatola osservi la successiva cella sulla destra) e supponiamo che ora la macchina passi nello stato B . Si continua in questo modo finché la tavola non ci dica di fermarci.

Se si seguono le regole correttamente, la macchina cancellerà il primo 1, sposterà il nastro verso sinistra, cella per cella, fino a inquadrare il segno $+$, cambierà il $+$ in 1 e si fermerà. Il nastro porterà a questo punto 1111111 ossia 7. Queste semplici regole programmano ovviamente la macchina per l'addizione di qualsiasi coppia di numeri interi anche molto grandi.

Si tratta naturalmente di un sistema lento e noioso, ma l'intenzione di Turing era di ridurre il calcolo meccanico a uno schema semplice e astratto, rendendo così più semplice l'analisi di spinose questioni teoriche quali lo stabilire cosa può e cosa non può venir calcolato. Turing ha dimostrato che questo suo congegno idealizzato può essere programmato per realizzare nel suo rozzo modo qualsiasi operazione alla portata del più potente calcolatore elettronico. Naturalmente come ogni calcolatore, e anche come il cervello umano, la macchina di Turing è limitata dal fatto che per certi calcoli occorre un numero infinito di passi (per esempio per calcolare π) e dal fatto che certi problemi sono

insolubili in linea di principio: *non esiste*, cioè, alcun algoritmo o procedimento efficace per risolverli. Una « macchina universale di Turing » è in grado di effettuare qualsiasi operazione che ogni singola macchina di Turing specializzata può effettuare; in altre parole, essa è in grado effettivamente di calcolare tutto ciò che è calcolabile.

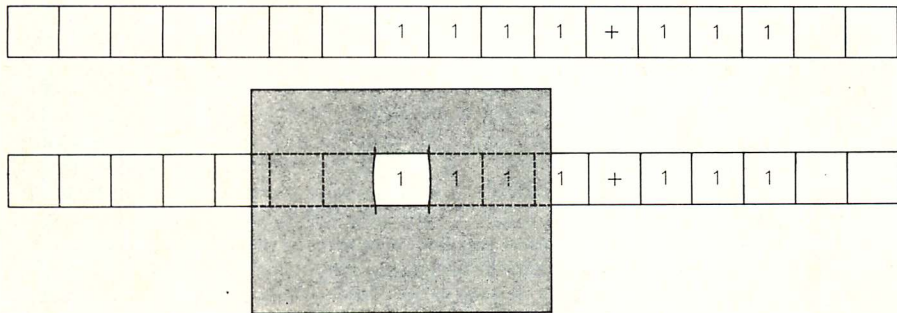
Nel 1950 apparve sul numero di ottobre di « Mind », una rivista inglese di filosofia, l'articolo di Turing intitolato *Computing machinery and intelligence* riprodotto in seguito in molte antologie, e anche nel volume *La filosofia degli automi* (Boringhieri, 1965). « Propongo — incominciava Turing — di prendere in considerazione la questione: può una macchina pensare? » Turing, considerando una simile domanda troppo vaga per avere una risposta significativa, suggeriva di sostituirla con una più precisa: si può insegnare a un calcolatore a vincere il « gioco di imitazione »? (ora generalmente noto come « gioco di Turing »).

Questo test è basato su un gioco di società in cui un uomo si nasconde in una camera e una donna in un'altra camera. Un « interrogatore », uomo o donna, pone domande, per mezzo di intermediari, alle due persone nascoste, e ne riceve risposte dattiloscritte. Ognuno dei giocatori cerca di convincere l'interrogatore di essere, diciamo, la donna. L'interrogatore vince se, in base alle risposte, riesce a stabilire la verità.

Supponiamo, dice Turing, di sostituire uno dei giocatori nascosti con una macchina che apprende, alla quale sia stato insegnato a esprimersi in un linguaggio ordinario, per esempio l'inglese. E' possibile per una simile macchina riuscire a ingannare l'interrogatore, quando sia la macchina sia l'altra persona nascosta cerchino di convincere l'interrogatore di essere umani?

Qui il senso di « ingannare » andrebbe definito esattamente; qual è la lunghezza della conversazione ammessa, qual è l'intelligenza dell'interrogatore e quella della persona in gara con la macchina? Attualmente un calcolatore riesce a superare il « test di Turing » se l'interrogatore è un bambino e se gli sono permesse solo poche domande. Probabilmente non ci saranno drammatici colpi di scena in futuro, così come probabilmente non ce ne furono nella lenta evoluzione dell'uomo. Le macchine parlanti potrebbero migliorare gradualmente, di modo che saranno necessari dialoghi sempre più lunghi e interrogatori sempre più intelligenti per battere la macchina. Forse un giorno solo un calcolatore sarà in grado di battere regolarmente un altro calcolatore. Turing stesso ha fatto una prudente previsione: nel 2000, ha scritto, i calcolatori parleranno inglese in modo abbastanza disinvolto da ingannare un « interrogatore medio » per circa il 30% delle volte, dopo circa 5 minuti di dialogo. Egli immaginò la seguente tipica conversazione tra una persona e un calcolatore:

Domanda: Scrivimi un sonetto su Forth Bridge.
Risposta: Lascia perdere, non sono mai stato capace di comporre poesie.
D. Addiziona 34957 e 70764.
R. (pausa di circa 30 secondi) 105621.
D. Sai giocare a scacchi?
R. Sì.



	STATO A.	STATO B.
1	1. CANCELLA L'1. 2. LEGGI LA CELLA SUCCESSIVA A DESTRA. 3. PASSA ALLO STATO B.	1. LEGGI LA CELLA SUCCESSIVA A DESTRA. 2. RIMANI NELLO STATO B.
+		1. CANCELLA IL +. 2. STAMPA 1. 3. STOP.

Una macchina di Turing per l'addizione.

D. Bene: ho il re e1 e nessun altro pezzo. Tu hai solo il re in e3 e la torre in h8. Tocca a te muovere. Cosa fai?

R. (pausa di 15 secondi) h8-h1. Matto.

Esistono parecchi programmi che consentono a un calcolatore di simulare una conversazione in inglese. Joseph Weizenbaum, un ingegnere del MIT, ne ha preparato uno che ha battezzato « Elisa » (nome ispirato dall'Eliza Doolittle del *Pigmaliione* di Shaw). Il « New York Times », il 5 giugno 1968, ha riprodotto un dialogo tra una donna e un calcolatore, un IBM 7094, in veste di psicanalista:

Donna: Gli uomini sono tutti uguali.

Calcolatore: Perché?

D. Ci seccano sempre per un motivo o per l'altro.

C. Può darmi un esempio specifico?

D. Beh, il mio amico ha insistito perché venissi da lei.

C. Il suo amico ha insistito?

D. Sì, dice che sono depressa.

C. Mi spiace sentire che lei è depressa.

D. E' vero, sono infelice.

C. Crede che venendo da me troverà un aiuto?

D. Ho bisogno di aiuto, certo.

C. Che specie di aiuto cerca?

D. Forse, dovrei imparare ad andare d'accordo con mia madre.

C. Mi dica qualcosa di più sulla sua famiglia.

D. Mia madre si occupa di me.

C. Chi altro si occupa di lei in famiglia?

D. Mio padre.

C. Suo padre?

D. Lei è un po' come mio padre, in un certo senso.

C. Che somiglianze trova?

D. Lei non è molto aggressivo, ma credo che non voglia farmelo notare.

C. Cosa le fa pensare che io non sia molto aggressivo?

D. Beh, non mi contraddice.

Il dialogo continua nello stesso modo ed è indistinguibile da una conversazione tra un paziente e uno psicanalista. Il programma non era nemmeno un programma di apprendimento. Weizenbaum ammise volentieri che il calcolatore non « capiva » niente di quanto veniva detto. Un simile calcolatore non potrebbe certamente superare il test di Turing. Supponiamo però che nell'anno 2000 i calcolatori possano giocare al « gioco di Turing » altrettanto bene quanto oggi giocano a dama e a scacchi: che significato avrebbe questo, dal punto di vista « mente » del calcolatore?

Chi ha visto il film 2001: *Odissea nello spazio*, ricorderà probabilmente che il calcolatore della nave spaziale, HAL, è definito come capace di « pensare » in quanto « avrebbe potuto superare con facilità il test di Turing ». HAL « pensa » veramente o imita semplicemente il pensiero? Turing era convinto che, quando un giorno i calcolatori sapranno discutere abbastanza bene da superare il suo test, nessuno potrà mettere in dubbio che essi « pensino ».

A questo punto sorgono interrogativi complicatissimi. Può un simile calcolatore essere cosciente di se stesso? Può avere emozioni? Senso dell'humor? In breve, può essere considerato una « persona » o è solo una macchina inerte costruita per imitare una persona?

Keith Gunderson, criticando l'articolo di Turing in « Mind » (aprile 1964) scriveva che la capacità di un calcolatore di passare il test di Turing non proverebbe niente dal punto di vista filosofico. E' possibile (l'esempio è mio) costruire un tulipano di cera che non sia distinguibile a vista da un tulipano vero, questo certamente non prova nulla circa la capacità di un chimico di sintetizzare le sostanze organiche del tulipano. Un calcolatore parlante potrebbe provare qualcosa, oltre al fatto che un calcolatore può essere in grado di imitare la conversazione? Gunderson concludeva: « il fatto che l'operaio sia stato rimpiazzato dalla scavatrice nelle gallerie, non prova certo che la macchina abbia dei muscoli; prova, casomai, che i muscoli non sono necessari per scavare gallerie ».

Una interpretazione curiosa del test di Turing è quella data da Michael Scriven in una conferenza del 1959, pubblicata poi col titolo *The Compleat Robot: A Prolegomena to Androidology* nella raccolta *Dimensions of Mind*. Scriven ammette che la capacità di parlare non prova che il calcolatore possieda altri attributi di una « persona ». Supponiamo però che un calcolatore parlante apprenda il significato di « verità » (nel senso, per esempio, specificato da Alfred Tarski nell'articolo in questo volume) e che sia programmato in modo tale da non poter mentire. « Questo lo renderebbe certamente inadatto — dice Scriven — al ruolo di cameriere personale, copywriter pubblicitario o personaggio politico, ma potrebbe renderlo capace di un altro servizio. » Ma potremo chiedergli se è cosciente di esistere, se ha emozioni, se crede che certe barzellette siano divertenti, se agisce liberamente, se gli piace un dato autore ecc., e avremmo il diritto di considerare le risposte come corrette.

E' possibile che la « macchina di Scriven » (come l'hanno battezzata alcuni filosofi polemizzando con l'articolo di Scriven) risponda no a tutte queste domande. Ma se dà una risposta, Scriven sostiene, saremmo giustificati nel

dargli credito come a un essere umano, e non avremmo ragioni per negargli la definizione di « persona ».

I filosofi non sono d'accordo su quanto sostenuto da Turing e da Scriven. In un articolo intitolato *The supercomputer as Liar* (« Mind », febbraio 1963) Scriven rispose ad alcuni critici. Mortimer J. Adler in *The Difference of Man and the Difference It Makes* è del parere che il test di Turing sia una questione di « tutto o niente » e che il successo o l'insuccesso nel creare calcolatori capaci di superare questo test, rafforzerà o rispettivamente indebolirà la teoria che l'uomo sia diverso da ogni macchina possibile, come da ogni animale.

Ma le macchine parlanti riuscirebbero veramente a smuovere la fede di coloro che credono in questa netta distinzione? Non è difficile immaginare che tra 50 anni ci potrà essere un programma televisivo con un robot come animatore nella cui memoria saranno immagazzinate migliaia di barzellette; dubito che qualcuno ne potrà dedurre che il robot abbia un senso dell'humor, come nessuno, battuto agli scacchi da un calcolatore, dedurrebbe di aver a che fare con un calcolatore di genere diverso da quello che gioca a filetto. Le regole della sintassi e della semantica non sono in effetti diverse dalle regole degli scacchi.

In ogni modo il dibattito continua, complicato da credenze metafisiche e religiose e da complessi problemi linguistici. Tutti i vecchi dilemmi concernenti il corpo, la mente e la natura della personalità sono riproposti in una nuova terminologia.

Samuel Butler, in *Erewhon* spiega come gli abitanti di Erewhon distrussero le loro macchine prima che queste potessero trasformarsi da servitori in padroni: le sue parole erano lette cento anni fa solo come improbabile satira; oggi, suonano come una profezia. « Non dà sicurezza — scriveva Butler — contro l'estremo sviluppo della coscienza meccanica, il fatto che oggi le macchine posseggono un basso grado di coscienza. Un mollusco non ha molta coscienza di se stesso. Riflettete sugli enormi progressi fatti dalle macchine durante gli ultimi cento anni, e notate quanto invece progredivano lentamente il regno animale e vegetale. Le macchine più organizzate non sono tanto creature di ieri, quanto degli ultimi minuti, per così dire, in confronto al lento evolversi della natura ».

Giochi, logica e calcolatori

Un solitario con domino colorati mostra che sono strettamente legati. Il fatto che si possa vincere il gioco oppure no è analogo al fatto che un problema si possa o meno risolvere per mezzo di un calcolatore

di Hao Wang

Oggi molto del lavoro che una volta era fatto dall'uomo viene delegato alle macchine e la gente si chiede sempre più spesso: quali sono le capacità umane insostituibili? Che cosa non possono fare le macchine? Può sorprendere il lettore che, mentre la prima domanda non ha una risposta definitiva, la seconda ha una soluzione matematica immediata.

Prima della costruzione del primo dei moderni calcolatori, il logico inglese Alan Turing formulò il seguente problema: che cosa non possono fare i calcolatori? Nel suo tentativo di creare una teoria su ciò che può essere computato e ciò che non può esserlo, Turing progettò un semplice calcolatore immaginario che dimostrò capace, da un punto di vista teorico, di effettuare tutte le operazioni di un qualsiasi calcolatore. Egli usò la sua macchina per dimostrare la stretta parentela fra la scienza dei calcolatori e la logica, branche della matematica che riguardano entrambe il concetto di dimostrazione matematica e i sistemi di notazione che permettono di presentare i nostri ragionamenti in forma esatta. Questo articolo si propone di illustrare alcuni concetti fondamentali nel terreno comune alla scienza dei calcolatori e alla logica per mezzo dei giochi.

La mente umana può afferrare solamente numeri e quantità relativamente piccole. La matematica, al contrario, è interessata in modo essenziale all'infinito. Le operazioni matematiche finite, da un lato, e le entità matematiche infinite, dall'altro, presentano un contrasto significativo e affascinante. La trasposizione graduale dei casi individuali, intuitivamente comprensibili, alle situazioni generali è un'importante acquisizione dell'intelletto umano. Alcune considerazioni astratte concernenti i giochi condurranno in maniera naturale all'osservazione di questo fenomeno.

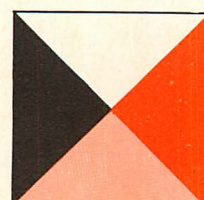
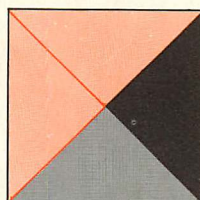
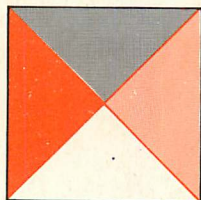
Trovare la successione di mosse che più probabilmente conduce alla vittoria in un gioco come il ticktacktoe è un pro-

blema logico analogo a quello di trovare la serie dei passi che conducono alla soluzione di qualsiasi problema matematico di una data classe. Per certi giochi non esiste una strategia ottimale che garantisca la vittoria; per certe classi di problemi non esiste un algoritmo, ossia non esiste nessun metodo generale che fornisca una serie di passi che conducano a una soluzione. Poiché il programma di un calcolatore è semplicemente un algoritmo progettato per operare per mezzo della macchina, ciò significa che esistono classi di problemi che i calcolatori non possono risolvere. Prima di considerare le difficoltà concernenti la costruzione di algoritmi per risolvere problemi (ovvero la formulazione dei programmi per ottenere soluzioni o la formulazione di strategie ottimali), vedremo perché la loro costruzione non è solo utile ma rappresenta anche uno degli scopi fondamentali della matematica.

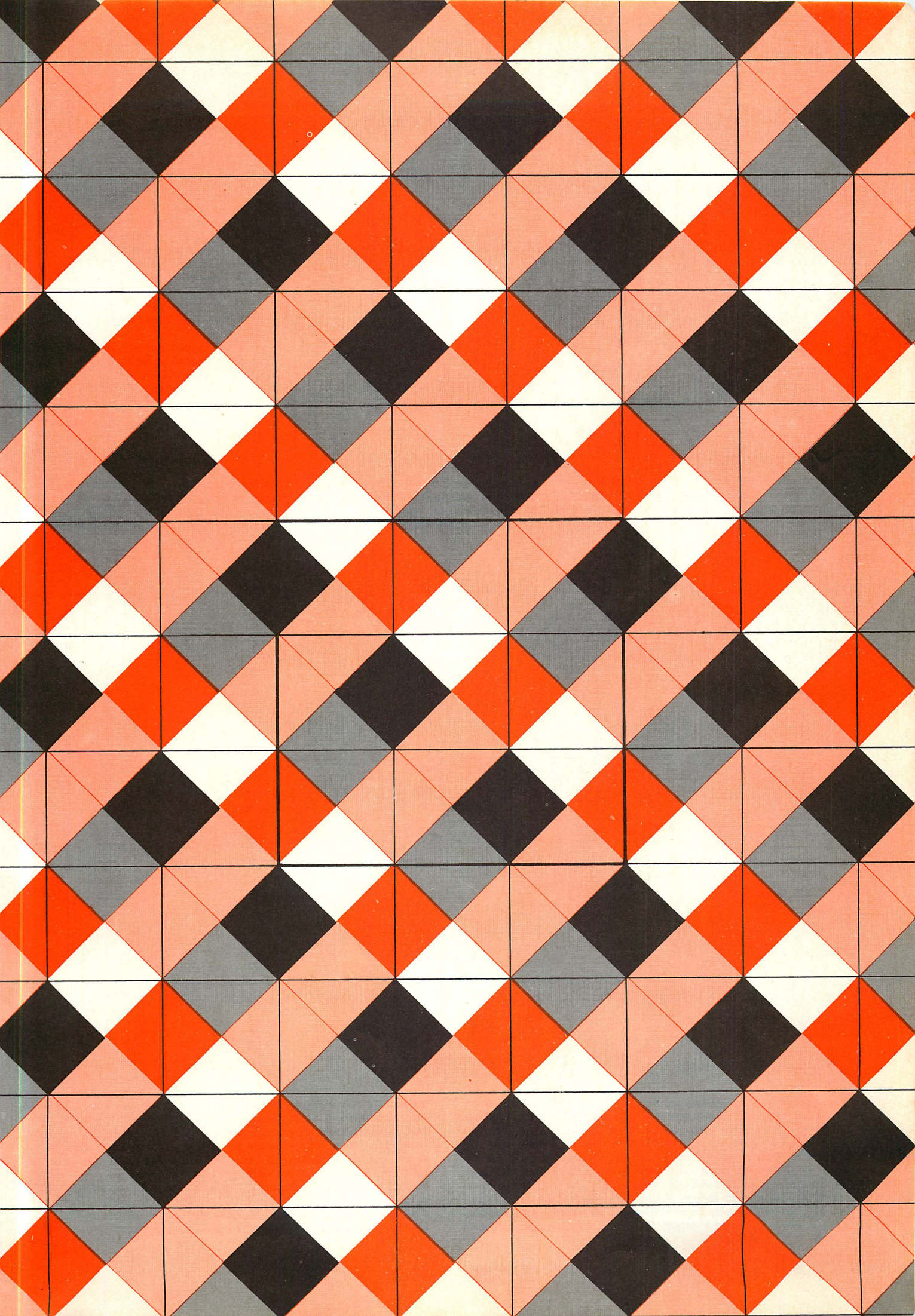
Ovviamente sarebbero necessari tempo ed energia infiniti per memorizzare la tavola per la moltiplicazione se, invece di includere solo i prodotti di tutte le coppie di numeri in una sola cifra, la tavola includesse il prodotto di tutte le coppie di numeri. L'uomo ha reso inutile questa tavola di moltiplicazione infinita

limitandosi a memorizzare, oltre alla tavola di moltiplicazione per i numeri di una cifra, una lista di passi che include il riporto e l'addizione dei prodotti parziali: ciò permette di ottenere il prodotto di due numeri qualsiasi di più cifre.

Sappiamo che le operazioni dell'aritmetica elementare comportano regole formali, e molti di noi ricordano che certe altre operazioni, come l'estrazione di radice quadrata, possono essere fatte secondo una lista fissa di passi successivi. Se abbiamo a che fare con problemi di maggior complessità diviene meno chiaro che essi possono essere risolti per mezzo di un algoritmo. Si consideri il seguente problema: dati due numeri positivi 6 e 9, si trovi il loro massimo comun divisore. Il lettore risponderà immediatamente: 3. Se i due numeri fossero 68 e 153, lettori che fossero propensi a tentare le varie possibilità potrebbero ancora trovare la risposta (17). Comunque, se si potesse mostrare che il problema generale «Dati due numeri positivi a e b , si trovi il loro massimo comun divisore» può essere risolto con un algoritmo, allora chiunque, o qualunque macchina capace di effettuare le operazioni specificate, potrebbe risolverlo per qualsiasi a e b . Un algoritmo di questo tipo effettivamente esiste e risale a Euclide (si vedano



Il problema del domino riguarda la riunione di tessere colorate appartenenti a tre tipi diversi per formare una superficie, estendibile all'infinito, in modo che tutti i bordi adiacenti abbiano lo stesso colore (si veda la pagina a fronte). Si suppone che il giocatore abbia una quantità infinita di domino di ogni tipo e che non si possa ruotare nessun domino nel piano. Il problema si risolve trovando una superficie rettangolare in cui la successione dei colori del bordo superiore è la stessa di quello inferiore, e la successione del bordo destro è la medesima di quello sinistro. Iterando tale superficie in ogni direzione si può ricoprire un piano infinito (si veda la pagina a fronte).




```

:R
GCD :R   THIS ROUTINE COMPUTES THE GREATEST COMMON
:R       DIVISOR OF TWO INTEGERS A AND B.
:R
:R   EXTERNAL FUNCTION (A,B)
NORMAL MODE IS INTEGER
ENTRY TO GCD.
LOOP REMAIN = B - A*(B/A)
WHENEVER REMAIN.E.0,FUNCTION RETURN .ABS.(A)
B = A
A = REMAIN
TRANSFER TO LOOP
END OF FUNCTION

```

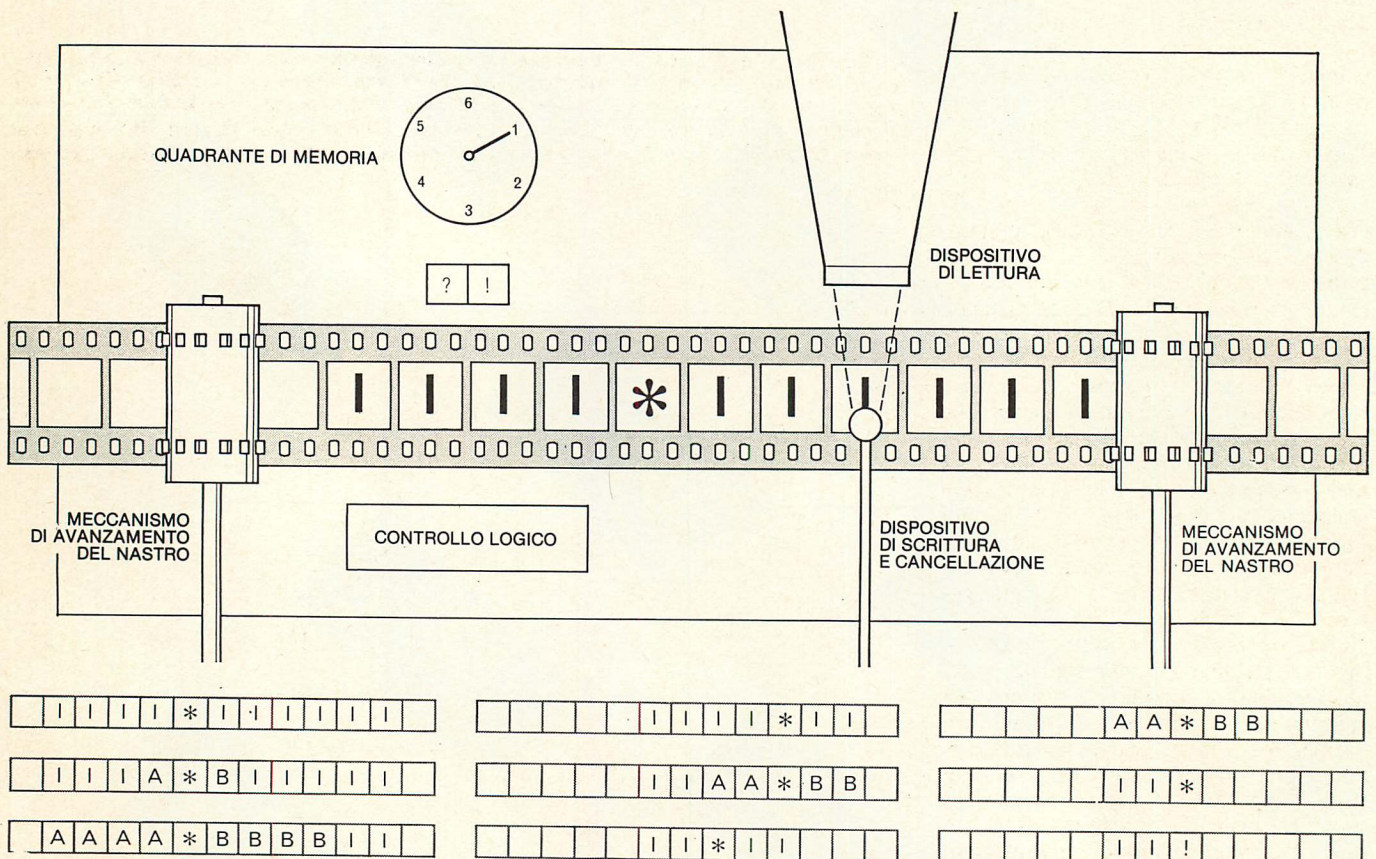
TYPE INPUT
A=8,B=12*
THE GCD OF A AND B IS
4

TYPE INPUT
A=12345678,B=87654321*
THE GCD OF A AND B IS
9

- 1 Si considerino due interi positivi a e b . Si passi all'istruzione successiva.
- 2 Si confrontino i due numeri in esame (si determini se sono uguali, in caso contrario si determini quale dei due è il maggiore). Si passi all'istruzione successiva.
- 3 Se i numeri sono uguali ognuno di essi è la risposta; stop. Se non sono uguali, si passi all'istruzione successiva.
- 4 Si sottragga il numero più piccolo dal più grande e si sostituiscano i due numeri in esame col sottraendo e col resto. Si passi all'istruzione 2.

Un algoritmo nel linguaggio di calcolo stabilisce una serie di passi tramite i quali si può trovare il massimo comun divisore di due numeri dati qualsiasi. In basso sono computate le soluzioni per due coppie di numeri. La procedura è data nel linguaggio MAD (da *Michigan Algorithm Decoder*). Sopra la lettera O, allo scopo di distinguerla dallo zero, è stata stampata una linea trasversale.

Lo stesso algoritmo viene ora espresso nel linguaggio ordinario. Il processo della divisione è reso attraverso sottrazione ripetuta. La serie dei passi è conosciuta come algoritmo euclideo.



In questa illustrazione schematica una macchina di Turing ideata per effettuare i passi dell'algoritmo euclideo porta i numeri 4 e 6 sul suo nastro. (Ciascuna cifra è rappresentata da una sbarra; un asterisco segnala la separazione dei numeri.) Il controllo logico della macchina consiste di istruzioni determinate dal segno stampato sulla cella in esame e dalla posizione del quadrante di memoria. I passi dell'algoritmo euclideo conducono ai cambiamenti del nastro illustrati in sequenza in basso nella figura. Prima la macchina determina quale numero è maggiore mediante un «ciclo di confronto» in cui si sostituiscono le sbarre a destra e a sinistra dell'asterisco con dei simboli («A»

a sinistra e «B» a destra). Quando un insieme di sbarre è esaurito, la macchina inizia il «ciclo di sottrazione», cancellando i simboli del numero più piccolo e riconvertendo i simboli del numero più grande in sbarre. Queste sono separate per mezzo di un asterisco dalle due sbarre che rappresentano il resto della sottrazione. Il processo di confronto e di sottrazione viene ripetuto per 4 e 2, poi per 2 e 2 e infine 2 e 0 appaiono sul nastro. A questo punto, appena inizia il ciclo di confronto, il nastro bianco a lato dell'asterisco fa scattare il segnale d'arresto (!): la macchina si ferma e sul nastro compare la risposta (2). È una macchina ideale perché il nastro è potenzialmente infinito.

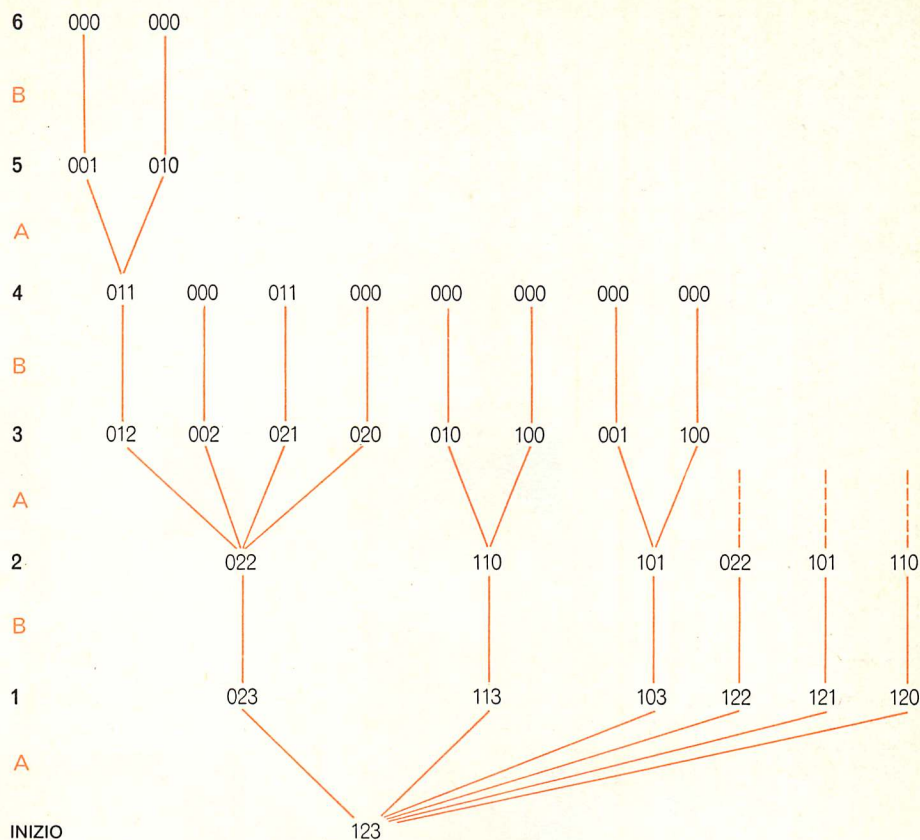
le illustrazioni in alto a pagina 132).

L'utilità degli algoritmi è ugualmente evidente nel campo dei giochi, dove essi sono studiati allo scopo di compiere le mosse più vantaggiose. Il matematico, naturalmente, è meno interessato a vincere un gioco che a capire la struttura astratta della classe di giochi a cui appartiene. Considerando l'esistenza o la non esistenza di una strategia che conduca alla vittoria, egli aumenta la sua comprensione della struttura astratta di quel gioco e di quelli del medesimo tipo.

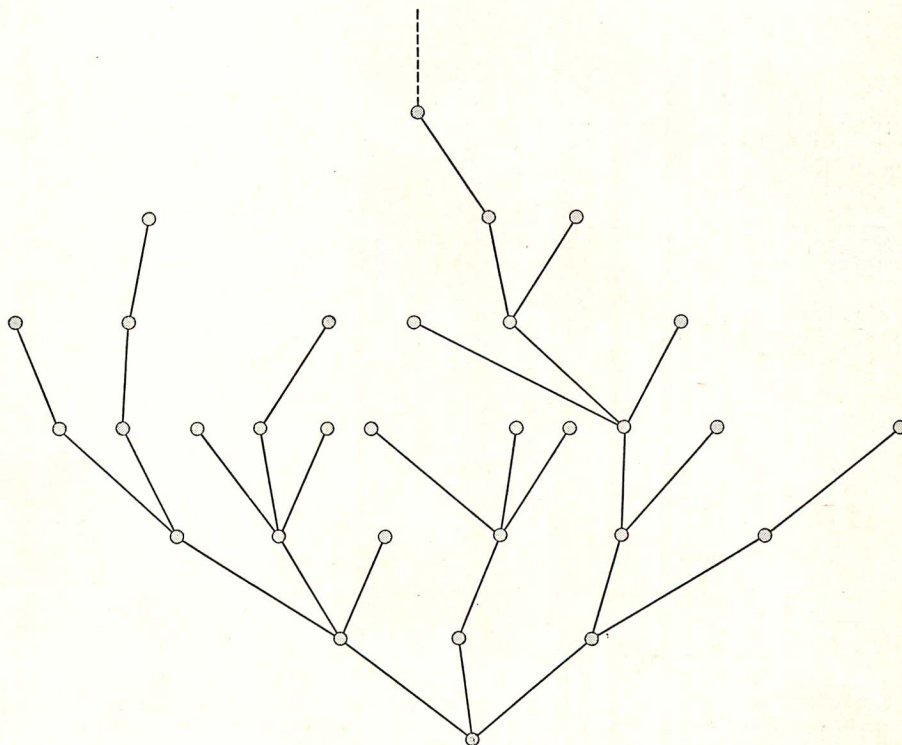
Si prenda per esempio il gioco conosciuto come nim. Un numero qualsiasi di oggetti, diciamo sei fiammiferi, sono messi in tre pile. Due giocatori, *A* e *B*, muovono a turno, prendendo ciascuno un numero qualsiasi di fiammiferi da una pila qualsiasi. Chi prende l'ultimo fiammifero è il vincitore. Poiché la quantità finita di fiammiferi prima o poi si esaurirà, è ovvio che il gioco non permette pareggio. È significativo che solamente uno dei giocatori abbia una strategia vincente, che dipende dalla grandezza delle tre pile di fiammiferi iniziali e da chi muove per primo. Nel gioco particolare definito per pile di 1, 2 e 3 fiammiferi, esiste una strategia vincente per *B*, il giocatore che muove per secondo. Questa situazione può essere rappresentata per mezzo di un albero schematico dove i nodi rappresentano la situazione ai vari stadi del gioco e i rami che partono da ciascun nodo rappresentano le possibili mosse che il giocatore può fare in quella situazione (si veda l'illustrazione in alto a destra).

Si supponga di giocare con tre pile contenenti rispettivamente dieci milioni, 234 e 2729 fiammiferi. È teoricamente possibile calcolare tutte le sequenze possibili di mosse a partire da queste tre pile e poi dire se esiste una strategia vincente per *A* o per *B*. Comunque nessuno vorrebbe o sarebbe in grado di effettuare tale calcolo. Il matematico inizierebbe una sistematica ricerca di scorciatoie per rendere le operazioni più facili e per realizzare una economia di pensiero. Una ricerca di questo tipo in effetti è stata fatta per il gioco del nim ed è stato ottenuto un semplice metodo per determinare quale giocatore abbia la strategia vincente. Il metodo stabilisce che *A* può vincere sempre se, una volta espresso il numero degli oggetti nelle tre file in notazione binaria, sommando colonna per colonna le cifre di questi tre numeri binari, non si ottiene mai come totale di una colonna un numero dispari. L'algoritmo scoperto può rappresentare una trovata geniale ma non importante, oppure il conseguimento di una maggior conoscenza matematica, a seconda della natura del gioco e delle sue relazioni con problemi matematici e logici più importanti.

Un gioco come il nim si dice «ingiusto», perché un giocatore ha sempre una strategia vincente. Un gioco come il ticktacktoe si dice «vano» perché ciascun giocatore ha una strategia per non perdere che impedisce che ci sia un vincitore.



L'albero per il nim indica che il giocatore che muove per secondo, *B*, ha una strategia vincente. All'inizio del gioco (in basso) ci sono pile di uno, due e tre fiammiferi (le cifre di ciascun nodo dell'albero danno il numero dei fiammiferi che rimangono nelle pile). I giocatori tolgono uno o più fiammiferi da una singola pila finché uno dei due giocatori vince togliendo l'ultimo fiammifero. I rami mostrano tutte le possibili mosse di *A* e le risposte di *B* a ognuna.



Il lemma dell'infinito si può rappresentare con un albero illimitato verso l'alto. Il lemma stabilisce che se nell'albero esistono infiniti rami connessi, e se da ciascun nodo parte solo un numero finito di ramificazioni, allora deve esistere un nodo, a ciascun livello, da cui le ramificazioni si estendono indefinitamente verso l'alto; questi nodi formano un sentiero infinito che attraversa l'albero. Il lemma dell'infinito può parafrasarsi così: se la specie umana non scomparirà mai, allora oggi esiste qualcuno che in ogni istante del futuro avrà un discendente.

re. Queste caratterizzazioni possono essere espresse come teorema: ogni gioco è vano o ingiusto se esiste un limite superiore finito alla lunghezza di ciascun sentiero del suo albero e se esiste un numero finito di rami che partono da ciascun nodo. Il teorema vale perché, se non sono ammessi giochi senza fine, e se a ogni stadio esiste solo un numero finito di mosse, allora il numero totale delle possibili sequenze di mosse è finito. Se si rappresentano tutti i giochi ammessi per

mezzo di un albero, si vede che se entrambi i giocatori non hanno una strategia vincente, allora ciascun giocatore ha una strategia per non perdere.

Il teorema non si applica in modo immediato al gioco degli scacchi perché non esistono regole precise per evitare che il gioco si prolunghi all'infinito. Supponiamo tuttavia che si possano introdurre regole per escludere le partite di scacchi senza fine, senza imporre un limite al numero delle mosse permesse in

una partita completa. Allora si potrebbe applicare agli scacchi una proposizione nota come lemma dell'infinito, che stabilisce che se esiste un numero infinito di rami connessi nell'albero di un gioco e da ciascun nodo parte solamente un numero finito di rami, allora esiste un cammino infinito.

Dato il lemma dell'infinito e assunte nuove regole che escludono le partite senza fine, segue che l'albero che rappresenta il gioco degli scacchi ha solamente



A



B



C

$\neg = \text{NON}$
 $\vee = \text{O}$
 $\wedge = \text{E}$
 $x' = x + 1$

1

$$(Axy \wedge Bx'y) \vee (Bxy \wedge Cx'y) \vee (Cxy \wedge Ax'y)$$



O



O



2

$$(Ayx \wedge Bx'y) \vee (B'yx \wedge Cyx') \vee (Cyx \wedge A'yx')$$



O

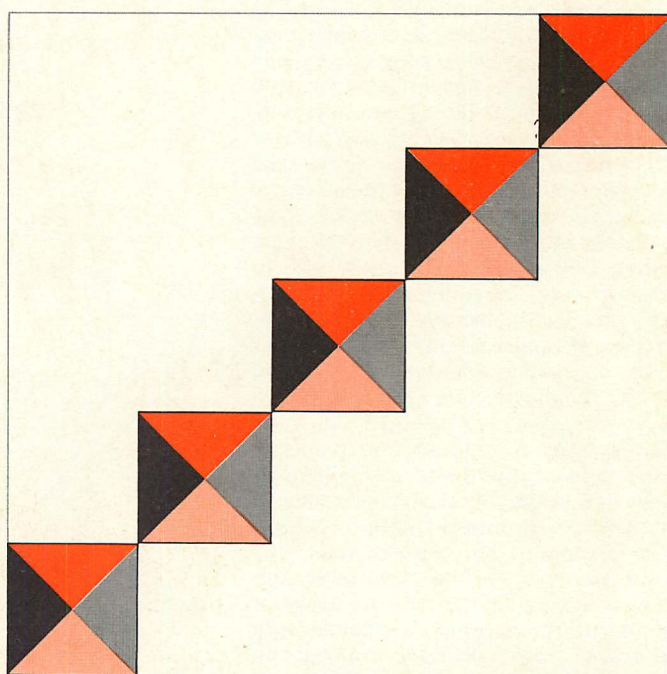


O



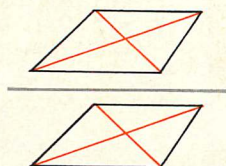
4

Axx



3

$$\neg (Axy \wedge Bxy) \vee \neg (Bxy \wedge Cxy) \vee \neg (Cxy \wedge Axy)$$



Le regole per i problemi del domino sono presentate nel simbolismo del linguaggio formale usato nella logica matematica (si veda il glossario in alto a destra). In alto al centro c'è un insieme di domino: A, B e C. La prima espressione stabilisce che i colori devono corrispondersi nei bordi a sinistra e a destra, la seconda stabilisce che i colori devono corrispondersi nei bordi in alto e in quelli in basso. La terza regola dice che le tessere non possono essere collocate l'una sopra all'altra.

La quarta è una delle tipiche limitazioni che vengono usate per complicare giochi nel tentativo di approssimare difficili problemi di calcolo. Essa stabilisce che solamente A può giacere sulla diagonale principale del piano. La posizione sul piano viene stabilita per mezzo di coordinate cartesiane. In designazioni come «Ayx» la posizione occupata della tessera sull'asse orizzontale è determinata dalla prima variabile, y, e la posizione sull'asse verticale dalla seconda variabile.

un numero finito di rami. Altrimenti, per il fatto che solamente un numero finito di rami parte da ciascun nodo, dovrebbe esistere un sentiero infinito (un gioco senza fine). Quindi c'è solo un numero finito di possibili successioni di mosse e perciò il gioco degli scacchi è mosse e perciò il gioco degli scacchi è mosse o ingiusto.

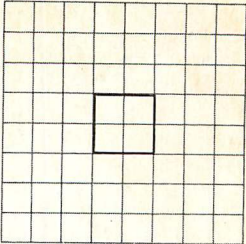
La dimostrazione del lemma dell'infinito è abbastanza immediata. Si prenda l'ultimo nodo dell'albero. Poichè si suppone che ci sia un numero infinito di rami, ma da ciascun nodo ne parte solo un numero finito, almeno uno dei nodi del livello precedente deve essere la fine di un sottoalbero con un numero infinito di rami. Si chiami questo nodo *X*. Tuttavia per ipotesi da *X* parte solo un numero finito di rami. Quindi uno dei nodi del livello precedente, rispetto a *X*, deve essere la fine di un sottoalbero infinito. Ripetendo il ragionamento si trova che a ogni livello c'è almeno un nodo che è la fine di un sottoalbero infinito e l'insieme di tutti questi nodi finali determina un sentiero infinito che attraversa l'albero. Un'applicazione antropomorfa del lemma dell'infinito potrebbe consistere nello stabilire che, se la specie umana non scomparirà mai, oggi esiste qualcuno che in ogni momento del futuro avrà un discendente.

Naturalmente si sceglie di esaminare un gioco a seconda dell'importanza dei problemi matematici a cui risulta legato. Nel 1960, mentre stavo studiando ai Bell Telephone Laboratories alcuni problemi di logica, inventai un nuovo solitario giocato con domino costituito da tessere colorate. Più recentemente, allo Harvard Computational Laboratory, io e alcuni miei colleghi trovammo alcune significative e sorprendenti applicazioni di questo gioco. Alcuni problemi che sorgono nel gioco del domino sono esattamente analoghi ai problemi per cui sono progettate le macchine di Turing. Le condizioni sotto cui si svolge il gioco possono essere fatte corrispondere alle computazioni di macchine di Turing, quindi lavorando con i domino si può avere una nuova versione, alcune volte particolarmente illuminante, di alcuni problemi matematici.

Nel gioco in questione viene dato un insieme finito di tessere quadrate (i domino); le tessere sono tutte della medesima grandezza, ma ciascuna ha il bordo di un colore stipulato e i colori sono combinati in parecchi modi specificati. Si assuma di avere una infinità di copie di ciascun tipo di domino, e che non sia permesso ruotare un domino nel piano. Il gioco consiste nel coprire con tutte le tessere un piano infinito in modo tale che i bordi che combaciano abbiano il medesimo colore. Se si può coprire il piano con un dato insieme di domino, si dice che l'insieme è risolubile.

Si consideri l'insieme di tre tessere mostrato nell'illustrazione di pagina 130. L'insieme è risolubile perchè può essere raggruppato in blocchi di nove tessere che soddisfano la regola che riguarda i

Si ricopra una sezione del piano cartesiano con tessere nere e bianche in modo tale che nessun blocco (delle dimensioni indicate qui a destra o più largo) abbia i bordi di sinistra e di destra, quelli in alto e quelli in basso che corrispondono. Esiste un metodo per ricoprire un piano infinito in questo modo?



L'autore presenta al lettore quattro problemi di cui solo il secondo non ha soluzioni note. La soluzione del problema mostrato in questa figura si trova nell'illustrazione di pagina 137.

È possibile scrivere successioni di 0 e di 1 in modo che producano una «progenie» per mezzo di queste regole:
 1. Se la stringa ha meno di tre simboli, stop.
 2. Se la stringa inizia con 0, si cancellino i primi tre simboli e si aggiunga 00 alla fine.
 3. Se la stringa inizia con 1, si cancellino i primi tre simboli e si aggiunga 1101 alla fine.

011010001001	101110110011
01000100100	1101100111101
0010010000	1100111101101
001000000	01111011101101
00000000	11011101110100
0000000	111011101001101
000000	0111010011011101
00000	101001101110100
0000	0011011101001101
000	10110100110100
00	1101001101001101

Date due stringhe, esiste un algoritmo per determinare se una delle due è una progenie dell'altra? (stop) (continuazione della progenie)

Il secondo problema consiste nel trovare un algoritmo che mostri se due stringhe di zero e di uno sono in una certa relazione. Il problema è complicato dal fatto che certe stringhe danno origine a progenie infinite. Una soluzione alternativa sarebbe mostrare che l'algoritmo non esiste.

Esiste un algoritmo per decidere se un'equazione polinomiale a coefficienti interi ha radici intere?

Le equazioni di questo tipo includono

$$x^2 - 4x + 3 = 0$$

e

$$a^2 + b^2 - c^2 = 0.$$

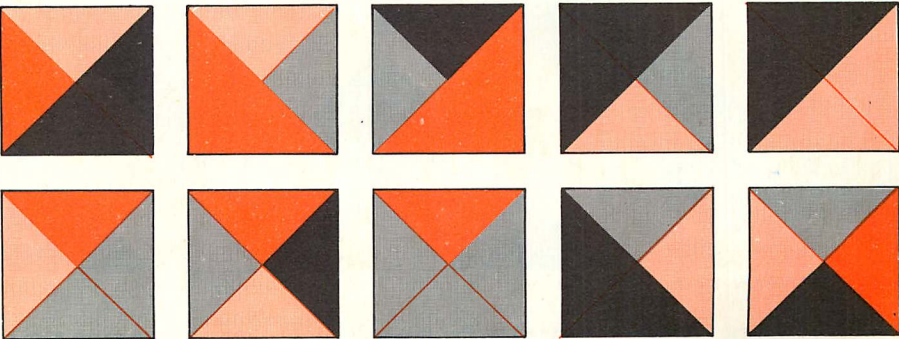
La prima equazione ha una sola incognita, *x*. Ha cioè la forma

$$a_n x^n + a_{n-1} x^{n-1} + \dots + a_1 x + a_0 = 0$$

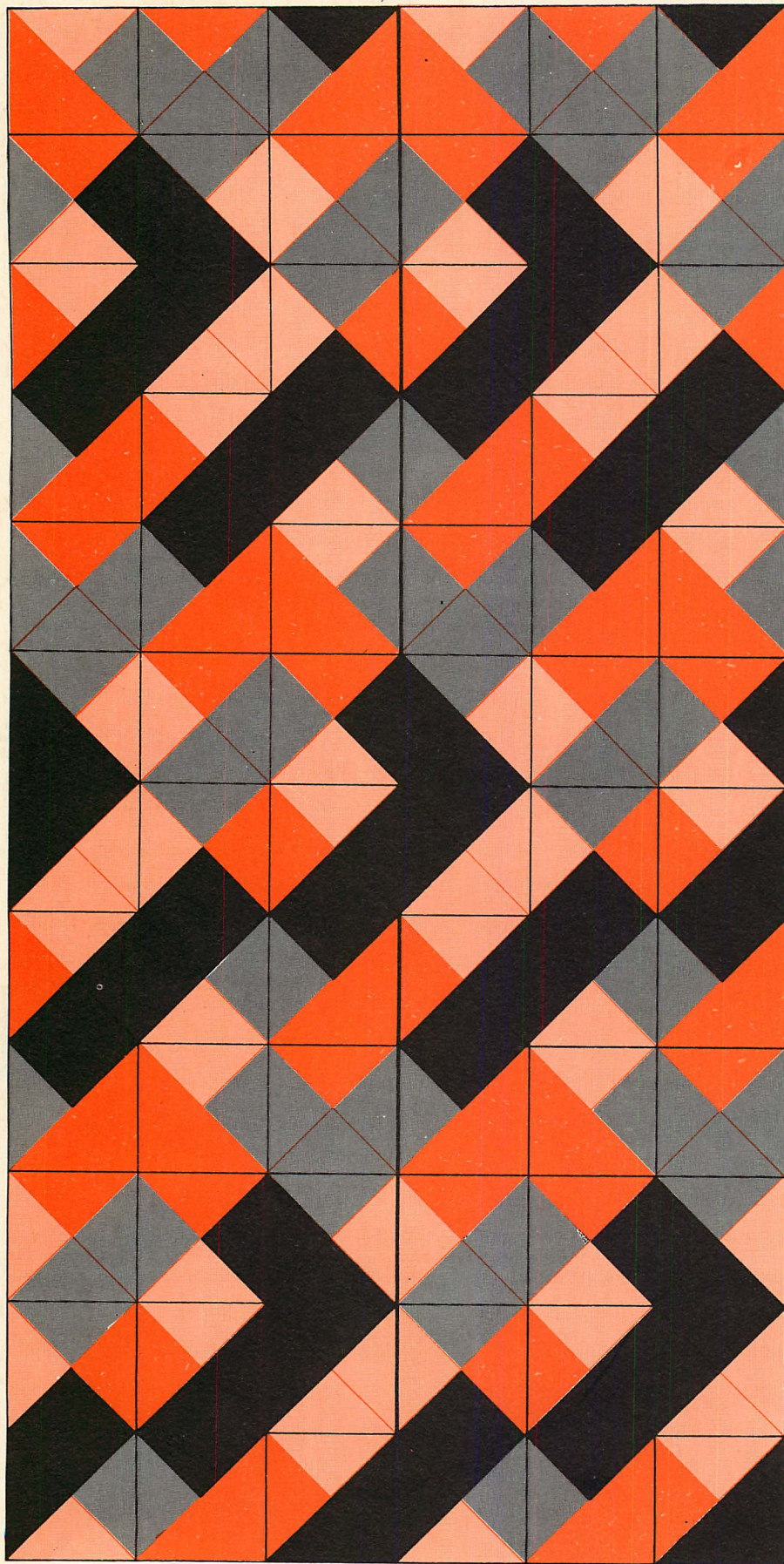
e per tali equazioni l'algoritmo desiderato è:
 1. Si trovino tutti i divisori di *a*.
 2. Si sostituisca la *x* con ogni divisore di *a* e si calcolino i valori risultanti per il lato sinistro dell'equazione.
 3. Se qualche divisore porta al valore 0, è una radice. Se nessuno porta al valore 0, l'equazione non ha radici intere.

Il problema consiste nell'ideare un algoritmo per le equazioni, come la seconda, che contengono più di una incognita.

Il terzo problema è conosciuto come «decimo problema di Hilbert» poichè il matematico tedesco David Hilbert nel 1900 lo elencò come uno dei problemi rilevanti a cui la matematica si trovava di fronte. Per la soluzione si veda in questo volume l'articolo dedicato a questo problema.



Il problema dei dieci domino sta nel disporre le tessere in modo che lo schema dei colori sia lo stesso per i bordi in alto, in basso, a sinistra e a destra. La soluzione alla pagina seguente.



Soluzione del primo problema di pagina 135. Per ricostruire la soluzione, sia a un pezzo nero e b uno bianco (a sinistra). Si scriva a e lo si sostituisca con ab . Si sostituisca b con ba e si continui a sostituire a e b in questo modo. Si trascriva la successione sulla riga in alto del piano e si copi ogni simbolo lungo la diagonale dall'alto a destra fino in basso a sinistra.

bordi e che possono essere ripetuti in ogni direzione. Data una soluzione per l'intero piano, è possibile ovviamente tagliare via tre settori quadrati (o quarti) per ottenere la soluzione di un settore quadrato. Il viceversa è meno ovvio, ma si può stabilire con l'aiuto del lemma dell'infinito. Poiché esiste la soluzione per un settore quadrato infinito, esistono soluzioni parziali per i suoi settori quadrati, soluzioni di ogni area di n per n . A partire da tali soluzioni parziali è possibile disegnare un albero infinito e mostrare per mezzo del lemma dell'infinito che nell'albero esiste un cammino infinito che produce la soluzione per l'intero piano. Quindi, se è possibile riempire un settore quadrato del piano infinito, è possibile riempire l'intero piano.

Si possono usare queste tessere per simulare varie macchine di Turing e per creare un equivalente del famoso problema dell'arresto (*halting problem*) formulato da Turing.

Questo si può fare più facilmente se si specifica quale domino deve essere messo all'origine del piano, cioè quale domino deve essere collocato per primo. Con maggior difficoltà si può ottenere lo stesso risultato o specificando che certi domino occorrono sulla diagonale principale, o omettendo qualsiasi restrizione all'infuori di quelle menzionate prima. Consideriamo questa equivalenza tra il problema dell'arresto e il gioco del domino in maggior dettaglio.

Turing ideò il suo semplice calcolatore per riprodurre l'operazione di calcolo effettuata dall'uomo. Un uomo che risolve problemi matematici è probabile che usi carta e matita per scrivere e cancellare i numeri; può inoltre avere una collezione di fatti matematici sotto forma di un libro di tavole e, contenuto nella sua mente o nel suo libro, un insieme di istruzioni per eseguire i passi adeguati nella successione adeguata. Anche l'immaginaria macchina di Turing ha un dispositivo per stampare e uno per cancellare i numeri secondo le istruzioni che provengono da una unità di controllo logica che segue una tavola di comandi predisposti in anticipo. I numeri vengono scritti sotto forma di singole sbarre su celle quadrate di un nastro infinitamente lungo che serve come unità di memoria. (Poiché nessuna macchina reale può avere una memoria infinita, la macchina di Turing è una idealizzazione.)

Una cella del nastro, a un dato tempo, viene esaminata da un dispositivo di esplorazione che fornisce all'unità di controllo il simbolo stampato sulla cella (*si veda l'illustrazione in basso a pagina 132*). L'unità di controllo quindi consulta le istruzioni interne per mezzo di un quadrante che indica la posizione designata «istruzione corrente». A seconda del simbolo in esame, l'istruzione specifica uno dei quattro ordini seguenti: (1) si stampi un contrassegno sulla cella, cancellando se necessario; (2) si muova il nastro di una cella a destra; (3) si muova il nastro di una cella a sinistra; (4) alt! L'istruzione

indica quindi la posizione dell'istruzione successiva. Si può dotare una macchina di Turing del numero necessario di posizioni delle istruzioni e del numero necessario di ordini, come pure di un nastro con un numero infinito di celle, per risolvere qualsiasi problema matematico (che appartenga alla classe dei problemi risolvibili per mezzo di algoritmi).

Turing ideò il problema dell'arresto come esempio di problema per il quale nessun programma poteva produrre tutte le soluzioni corrette: formulando un problema che riguardava tutte le macchine di Turing possibili superò le capacità di qualsiasi singola macchina di Turing. Dimostrò che sebbene ciascuna macchina, a seconda del suo nastro, si arresti a un certo momento o continui a operare indefinitamente, non esiste nessun algoritmo generale per determinare questo comportamento, ossia non esiste nessun metodo equivalente a quello che determina chi è il vincitore nel nim. Ora è possibile trovare, per ciascuna macchina di Turing, un insieme di tessere del domino tali che la macchina si ferma se e solo se l'insieme di tessere non ha soluzione. Si ricava immediatamente che il problema del domino è irrisolvibile. Se si potesse risolvere il problema del domino, si potrebbe risolvere il problema dell'arresto, quindi il problema del domino non si può risolvere. In altre parole, non esiste un metodo generale per decidere se un qualsiasi insieme di tessere ha una soluzione.

Il problema del domino è un esempio di «problema di decisione infinito», un problema che compare frequentemente in logica, nella teoria dei calcolatori e in generale in matematica. Si tratta di un problema infinito nel senso seguente: qualsiasi soluzione al problema del domino deve consistere in un singolo metodo che fornisca la risposta corretta (sì o no) a un numero infinito di problemi della forma: «Un certo insieme di tessere ricopre il piano?» Mentre qualsiasi dato insieme di tessere del domino è finito, esiste naturalmente un insieme infinito di tali insiemi e quindi un numero infinito di problemi.

In questo modo si sono ridotti i problemi riguardanti insiemi di tessere del domino a problemi riguardanti macchine e si sono stabiliti i risultati riguardanti i domino facendo appello ai risultati conosciuti riguardanti le macchine. Il passo successivo consiste nel ridurre il problema dell'interpretazione di una formula in logica al problema che riguarda la soluzione di un insieme di tessere. Poiché il fatto che un insieme di tessere ha una soluzione può essere espresso in lo-

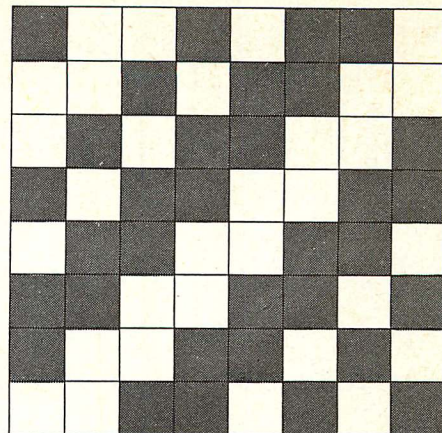
a

a b

a b b a

a b b a b a a b

a b b a b a a b b a a b b a b b a ...



La soluzione del problema dei dieci domino è un blocco rettangolare di 36 domino, due dei quali sono separati da una linea nera pesante che attraversa il centro dell'illustrazione. La soluzione non è unica in quanto sono possibili e ugualmente accettabili altre configurazioni.

gica per mezzo di una semplice formula, questa riduzione fornisce una risposta a un importante problema di decisione nell'ambito della logica.

Se si desidera esprimere la condizione che un insieme di tre tessere del domino ha una soluzione nel primo quadrante del piano infinito, basta riferirsi alle familiari coordinate cartesiane per determinare la posizione dei domino nel quadrante e rappresentare ciascun domino per mezzo di un predicato, Axy percorre nella posizione (x, y) . Se si scrive x' per $x+1$, la condizione richiesta può essere data per mezzo di un certo numero di clausole che richiedono pochissimi quantificatori (ossia espressioni del tipo «per tutti gli x » e «esiste almeno un x »). Siamo generosi solamente con le operazioni finite della logica formale, quali «non», e «e», e «o» (si veda l'illustrazione di pagina 134). Si conclude che per ogni insieme di tessere dato si può trovare una corrispondente «formula AEA » (una proposizione che inizia con «Per tutti gli x , esiste un y , tale che per tutti gli z ...», seguita da una combinazione logica di predicati senza quantificatori) tale che l'insieme ha una soluzione se e solo se la formula non è autocontraddittoria. In altre parole, si può tradurre un problema riguardante il domino in una formula logica specificando certe restrizioni, e quindi determinare se l'insieme di tessere è risolvibile guardando se la formula è o non è autocontraddittoria. Quindi, poiché il problema generale del domino non è risolvibile, non esiste nessun metodo generale per decidere se una arbitraria formula AEA è autocontraddittoria.

Si tratta di un risultato utile perché in

logica la complessità delle formule è spesso misurata per mezzo del numero e dell'ordine dei quantificatori e le formule sono spesso poste in classi differenti a seconda dei quantificatori che vi compaiono. È sorprendente che una classe semplice come quella delle formule AEA (con tre quantificatori soltanto) sia indecidibile. Infatti, con questo risultato si possono risolvere i problemi di decisione per tutte le classi di quantificatori. Data una qualsiasi successione di quantificatori si può ora dire se la classe di formule determinata da essa è decidibile.

Il problema della decidibilità della logica è interessante perché tutte le teorie matematiche possono essere formulate nei termini della logica elementare. Il problema che riguarda il fatto che una formula (F) possa o non possa essere derivata da un insieme di assiomi (A) si riduce al problema di decidere se la formula logica « A ma non F » non sia autocontraddittoria. In questo senso tutta la matematica è riducibile alla logica. Infatti la misura della complessità di un problema matematico è data dalla struttura della formula logica corrispondente. È quindi un compito importante determinare la complessità delle varie classi di formule logiche.

Si può legittimamente dire che tutta la matematica può essere ridotta, per mezzo delle macchine di Turing, a un solitario giocato con le tessere di questo domino. In molti casi la riduzione non rende più facile la trattazione di un problema matematico, tuttavia la dimostrazione che certi problemi non sono risolvibili da un calcolatore può essere facilitata riducendoli a problemi riguardanti le tessere del domino.

Il decimo problema di Hilbert

È possibile determinare una procedura che indichi se esistono soluzioni di un'equazione diofantea? Oggi finalmente è stata data una risposta a tale domanda che faceva parte di un celebre elenco di problemi

di Martin Davis e Reuben Hersh

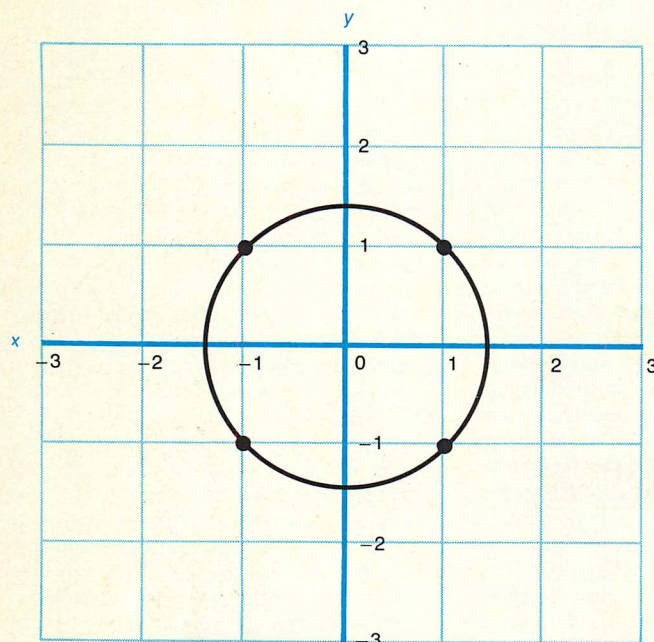
«Sentiamo in noi l'eterno richiamo: ecco il problema, cercarne la soluzione. È possibile ricavarla facendo appello solo alla ragione, dal momento che in matematica non esiste *ignorabimus* [ossia problemi non risolvibili].» In tal modo David Hilbert inaugurò il secondo Congresso internazionale dei matematici, tenuto a Parigi l'otto agosto 1900, salutando l'inizio del nuovo secolo con la presentazione di un elenco di 23 importanti problemi come sfida per i futuri ma-

tematici. Alcuni dei problemi proposti da Hilbert rimangono ancora insoluti; altri hanno ispirato generazioni di ricercatori e hanno portato a nuove e importanti teorie matematiche. L'ultimo risultato sui problemi di Hilbert riguarda il decimo, che è stato risolto nel 1970 dal matematico russo ventiduenne Yuri Matyasevich.

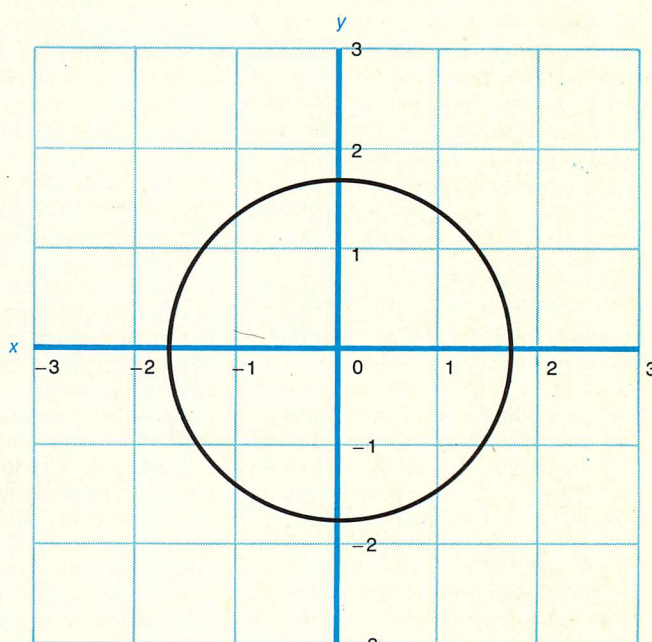
David Hilbert nacque a Königsberg nel 1862 e fu professore all'Università di Gottinga dal 1895 fino alla sua morte avvenuta nel 1943; venne considera-

to, dopo la morte di Henri Poincaré nel 1912, come il più importante matematico del suo tempo. Diede contributi fondamentali in molte discipline, ma forse è ricordato soprattutto per il suo sviluppo del metodo astratto quale potente strumento per la matematica.

È molto facile descrivere il decimo problema di Hilbert. Esso riguarda la più semplice e fondamentale attività matematica: risolvere equazioni. Le equazioni da risolvere sono polinomiali, cioè del tipo $x^2 - 3xy = 5$, costi-



$$x^2 + y^2 - 2 = 0$$



$$x^2 + y^2 - 3 = 0$$

I grafici delle due equazioni mettono in luce la differenza fra un'equazione usuale e un'equazione diofantea in cui ci si occupa solo delle soluzioni intere; tale differenza è sostanziale per formulare il decimo problema di Hilbert. Le equazioni in questione sono $x^2 + y^2 - 2 = 0$ (a sinistra) e $x^2 + y^2 - 3 = 0$ (a destra); entrambe sono rappresentate con circonferenze il cui centro è sull'origine, cioè sul punto di coordinate $x = 0, y = 0$. Nel caso di $x^2 + y^2 - 2 = 0$ la circonferenza ha un raggio di valore $\sqrt{2}$. Se l'equazione viene considerata nel modo usuale, vi

saranno infinite soluzioni. Se invece è vista come equazione diofantea, ve ne saranno solo quattro: 1) $x = 1, y = 1$, 2) $x = -1, y = 1$, 3) $x = 1, y = -1$ e 4) $x = -1, y = -1$. Tali soluzioni sono state messe in evidenza là dove il grafico passa per i quattro punti le cui coordinate sul reticolato cartesiano sono quelle sopra specificate. Nel caso di $x^2 + y^2 - 3 = 0$ la circonferenza ha un raggio di valore $\sqrt{3}$. Considerata come un'equazione usuale essa ha un numero infinito di soluzioni; se viene considerata come un'equazione diofantea non può avere alcuna soluzione.

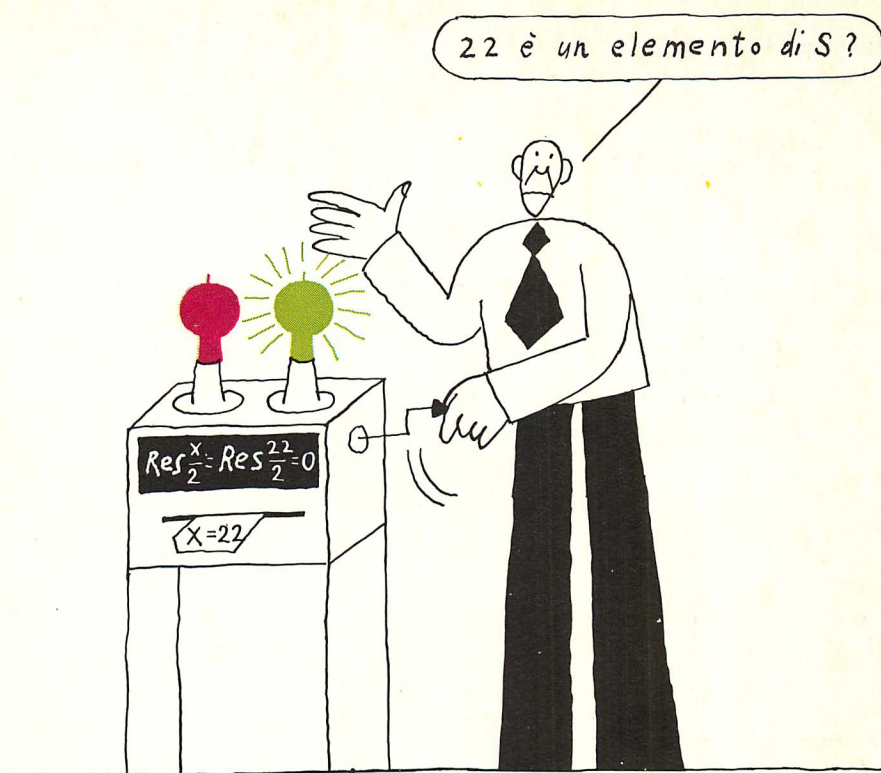
tuite da somme e prodotti di costanti e variabili, e in cui compaiono esponenti interi positivi. Oltre a ciò, Hilbert richiese che le equazioni dovessero far uso solo di numeri interi (positivi o negativi): né i numeri irrazionali né quelli immaginari e neppure le frazioni sono ammessi sia nelle equazioni sia nelle loro soluzioni. Equazioni di questo genere sono dette diofantee da Diofanto di Alessandria che scrisse un libro su tale argomento nel terzo secolo.

Il decimo problema di Hilbert è allora: descrivere una procedura meccanica grazie alla quale sia possibile saggiare ogni equazione diofantea per decidere se possiede soluzioni. Con le parole di Hilbert: «Data un'equazione diofantea con qualunque numero di quantità incognite e con coefficienti numerici interi razionali, escogitare un procedimento in base al quale sia possibile determinare, con un numero finito di operazioni, se l'equazione ammetta soluzioni razionali.» La richiesta di Hilbert non riguarda un procedimento per trovare le soluzioni ma semplicemente per determinare se l'equazione ne possieda. Tale procedura deve essere formalmente dettagliata in modo da poter essere programmata per una macchina calcolatrice e utilizzata in tutti i casi possibili: procedimenti di questo genere sono detti algoritmi.

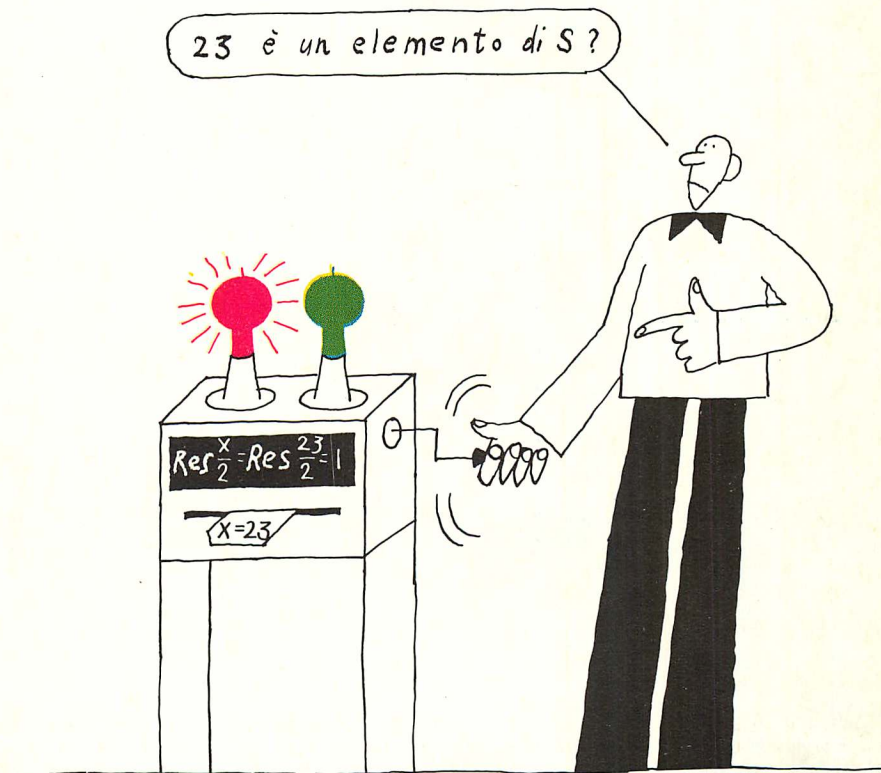
Se il problema di Hilbert è semplice a formularsi, ancor più facile è formulare la soluzione di Matyasevich: non è possibile escogitare una procedura di questo tipo; un algoritmo siffatto non esiste. Così espressa, la soluzione suona negativa in modo scoraggiante. In realtà, il risultato di Matyasevich costituisce un importante contributo alla conoscenza delle proprietà dei numeri.

Il lavoro di Matyasevich estende una serie di ricerche di tre americani: uno degli autori (Davis), Julia Robinson e Hilary Putnam. A loro volta, i lavori di questi ultimi erano basati su precedenti ricerche di molti dei fondatori della logica moderna e della teoria della computabilità: Alan Turing, Emil Post, Alonzo Church, Stephen Kleene e lo stesso Kurt Gödel, celebre per i suoi risultati sulla consistenza dei sistemi assiomatici (secondo problema di Hilbert) e sull'ipotesi del continuo di Cantor (primo problema di Hilbert).

Iniziamo l'esposizione del decimo problema di Hilbert analizzando alcune equazioni diofantee. La locuzione «equazione diofantea» è in parte fuorviante dal momento che a essere cruciale non è tanto la natura dell'equazione quanto quella delle soluzioni ammesse. Per esempio, l'equazione $x^2 + y^2 - 2 = 0$, se non la si con-



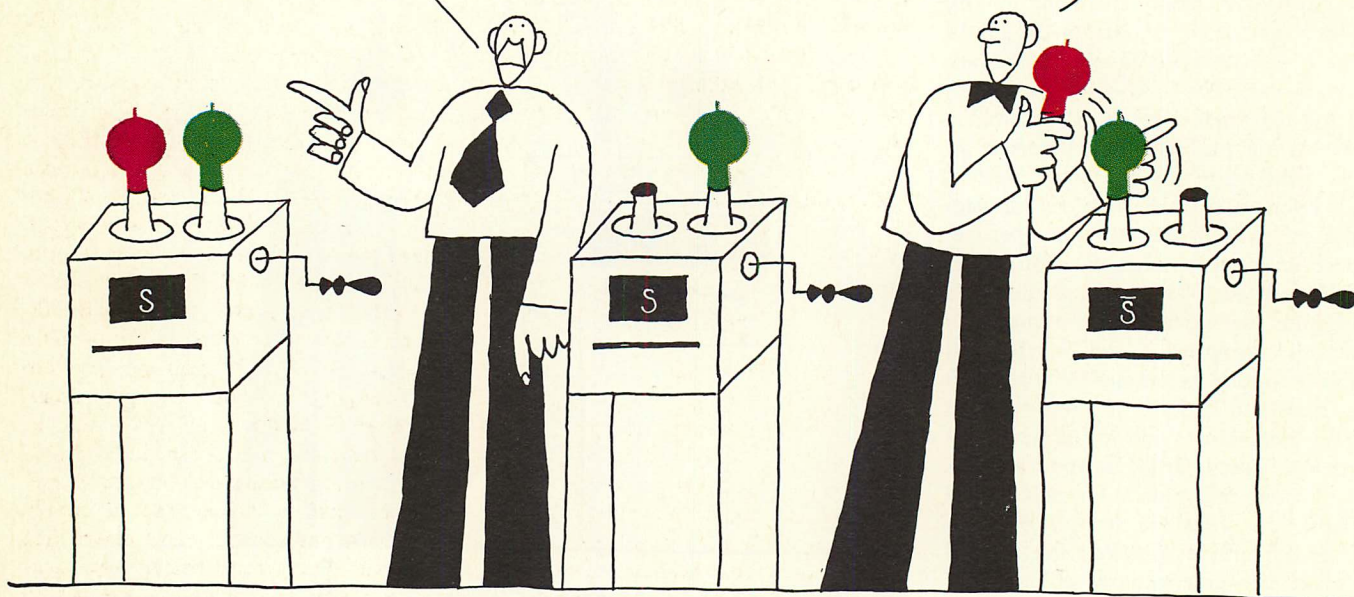
Una macchina a luce verde e rossa è un dispositivo ideale che esamina numeri per determinare se siano o meno elementi di un insieme dato. Il decimo problema di Hilbert consiste nel chiedersi se sia possibile costruire una «macchina di Hilbert» a luce verde e rossa per saggiare se le equazioni diofantee abbiano o meno soluzioni. Nel caso dell'esame compiuto sui numeri per determinare la loro appartenenza a un insieme, si accende la luce verde se la macchina è stata in grado di stabilire in un numero finito di passi che un certo ingresso è elemento dell'insieme. Se per esempio S è l'insieme di tutti i numeri pari, si può trovare un algoritmo che divida ogni ingresso x per 2. Se il resto è 0 la macchina accenderà la luce verde indicando che x è un elemento di S .



La luce rossa si accenderà nella macchina a luce verde e rossa se essa è stata in grado di determinare che l'ingresso non è elemento dell'insieme. Se l'ingresso x è il numero intero 23, diviso per 2 dà un resto 1, quindi 23 non è elemento di S . Il complemento di S è \bar{S} , l'insieme dei numeri dispari. Poiché si può costruire una macchina a luce verde e rossa per separare gli elementi di S da quelli di \bar{S} , S è computabile.

E' possibile trasformare una macchina a luce verde e rossa in una macchina a luce verde per S più una macchina a luce verde per il complemento di S ?

Sì, posso dimostrarlo collegando diversamente le lampadine

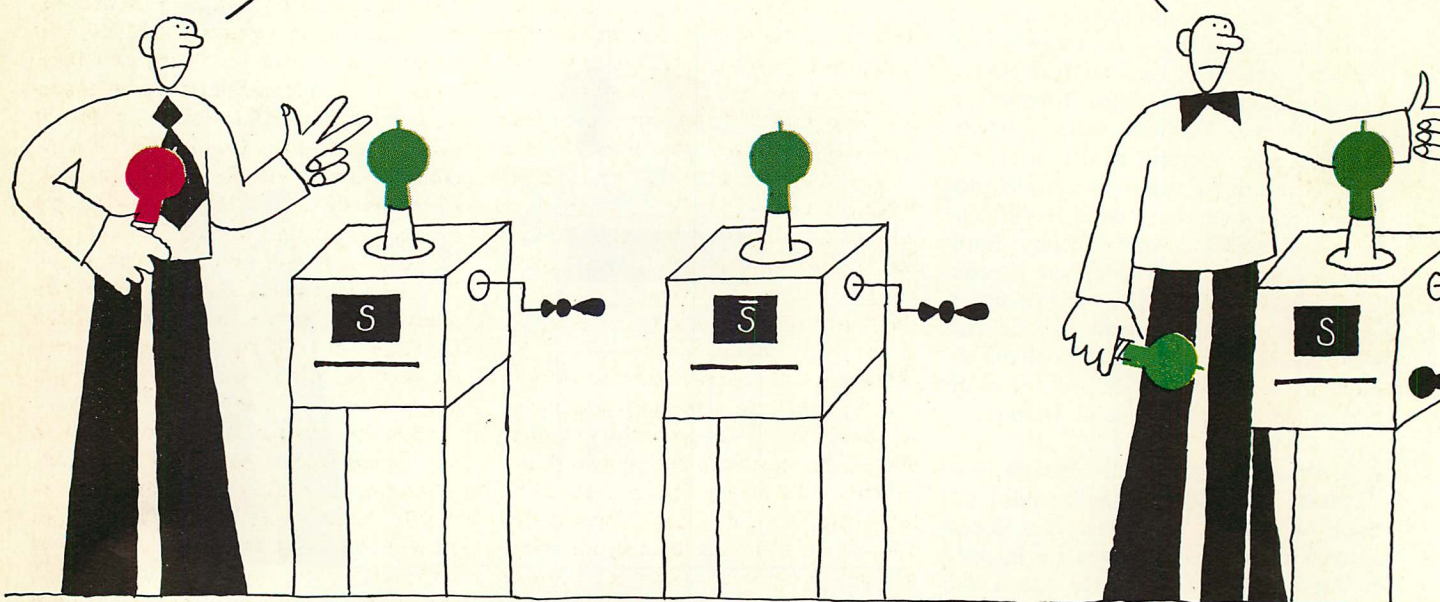


Una macchina a luce verde e rossa per l'insieme S può essere trasformata in una macchina a luce verde per S (cioè una macchina che si accende solo quando l'ingresso è un elemento di S) più una macchina a luce verde per \bar{S} , il complemento di S . La dimostrazione è semplice. Per costruire una macchina a luce verde per S , è sufficiente svitare la lampada rossa della mac-

china a luce verde e rossa. Per costruire una macchina a luce verde per \bar{S} , occorre svitare la lampada verde della macchina a luce verde e rossa e porla nello zoccolo in cui vi era la lampada rossa. In altri termini: se un insieme (come S , l'insieme dei numeri pari) è computabile allora sia tale insieme sia il suo complemento (come \bar{S} , l'insieme dei numeri dispari) sono elencabili.

Viceversa, posso utilizzare una macchina a luce verde per S più una macchina a luce verde per il complemento di S (oltre a una nuova lampada rossa)

... per costruire una macchina a luce verde e rossa per S .



Le macchine a luce verde per S e per \bar{S} possono essere utilizzate per costruire una macchina a luce verde e rossa per l'insieme S . Nella macchina a luce verde per \bar{S} si sostituisca la lampada ver-

de con una lampada rossa e si colleghino le due macchine in parallelo di modo che uno stesso ingresso valga simultaneamente per entrambe. Il risultato sarà ovviamente una macchina a lu-

sidera come equazione diofantea, ha infinite soluzioni. Esse possono venir rappresentate mediante il grafico dell'equazione che è una circonferenza nel piano formato dagli assi delle x e delle y . Il centro di tale circonferenza ha le coordinate $x = 0$ e $y = 0$: tale punto è detto origine ed è abbreviato con $(0,0)$. Il raggio della circonferenza è $\sqrt{2}$ (si veda l'illustrazione della pagina 138). Le coordinate di ogni punto sulla circonferenza soddisfano la equazione e vi è un numero infinito di tali punti. Tuttavia, se consideriamo il problema dal punto di vista di un'equazione diofantea, vi sono solo quattro soluzioni: 1) $x = 1, y = 1$; 2) $x = -1, y = 1$; 3) $x = 1, y = -1$ e 4) $x = -1, y = -1$.

Si supponga di modificare l'equazione in $x^2 + y^2 - 3 = 0$. Vi sarà ancora un numero infinito di soluzioni se la consideriamo come un'equazione usuale, ma nessuna soluzione del tutto se la consideriamo come una equazione diofantea. La ragione sta nel fatto che ora il grafico è una circonferenza il cui raggio è $\sqrt{3}$, e nessun punto di tale curva ha le coordinate simultaneamente uguali a un numero intero.

Una famiglia molto nota di equazioni diofantee ha la forma $x^n + y^n = z^n$, con n uguale a 2, 3, 4 o a qualsiasi intero maggiore. Se n è uguale a 2, la

equazione è soddisfatta dalle lunghezze dei lati di ogni triangolo rettangolo ed è chiamata teorema di Pitagora. Una di tali soluzioni è l'insieme dei numeri $x = 3, y = 4, z = 5$. Se n è uguale o maggiore di 3, l'equazione è nota sotto il nome di equazione di Fermat. Il matematico francese del diciassettesimo secolo Pierre de Fermat ritenne di aver dimostrato che tali equazioni non hanno soluzioni intere positive. Sul margine della sua copia del libro di Diofanto egli scrisse di avere trovato una « dimostrazione veramente mirabile » che sfortunatamente era troppo lunga per poter essere scritta in quello spazio. La dimostrazione (ammesso che in realtà Fermat ne avesse trovata una) non è mai stata scoperta. Noto sotto il nome di ultimo teorema di Fermat, è probabilmente il più antico e famoso dei problemi matematici non risolti. Questi esempi mostrano come le equazioni diofantee siano facili a scriversi ma difficili a risolversi; e sono così difficili a risolversi proprio perché sono così esclusive nei riguardi del tipo di numeri accettati come soluzioni.

Per quanto riguarda le equazioni di primo grado, cioè le equazioni in cui le incognite non vengono moltiplicate fra loro e tutti gli esponenti sono uguali a 1, come per esempio $7x + 4y - 3z - 99t + 13u - 10 = 0$, l'esistenza di soluzioni può essere determinata mediante un metodo di divisione noto fin dall'antichità come algoritmo di Euclide. Per le equazioni di secondo grado a due incognite, come $3x^2 - 5y^2 + 7 = 0$ oppure $x^2 - xy - y^2 = 1$, una teoria, sviluppata all'inizio del diciannovesimo secolo dal grande Karl Friedrich Gauss, permette di stabilire se vi siano soluzioni. Recenti lavori del giovane matematico inglese Alan Baker hanno gettato una considerevole luce sulle equazioni a due incognite di grado superiore al secondo. Per quanto concerne le equazioni di grado superiore al primo e che abbiano più di due incognite, vi sono solo alcuni casi speciali che possono essere trattati con particolari artifici, ma per il resto ci troviamo in un vasto mare di ignoranza.

Perché è così difficile trovare un procedimento del tipo di quello richiesto da Hilbert? L'approccio più diretto sarebbe semplicemente quello di sperimentare tutti i possibili insiemi di valori delle incognite, uno dopo l'altro, fino a trovare una soluzione. Per esempio, se l'equazione ha due incognite, si potrebbe fare un elenco di tutte le coppie di interi. A questo punto basterebbe semplicemente procedere lungo l'elenco, provando una coppia dopo

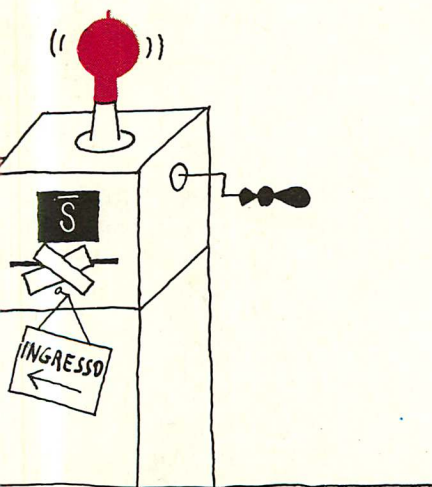
l'altra, per vedere se soddisfa l'equazione. Questa è certamente una procedura dettagliata e meccanica che potrebbe essere espletata da una macchina. Quale sarebbe però il risultato?

Se l'equazione è la prima che abbiamo menzionato, cioè $x^2 + y^2 - 2 = 0$, si devono provare $(0,0)$, $(0,1)$, $(1,0)$, $(0,-1)$, $(-1,0)$ e scartarle tutte. La possibilità successiva, $(1,1)$, è una soluzione. Siamo stati fortunati: sono state considerate solo sei coppie. Se invece l'equazione fosse $x^2 + y^2 = 20\,000$, si dovrebbero provare migliaia di coppie di numeri prima di giungere a una soluzione. In ogni caso è chiaro che se esiste una soluzione essa verrà trovata in un numero finito di passi.

D'altro canto, che cosa si può dire riguardo alla seconda equazione: $x^2 + y^2 - 3 = 0$? Si possono provare per tutta l'eternità nuove coppie di interi e tutto quello che si potrà sempre sapere sarà che non si è ancora riusciti a trovare una soluzione. Ma non si riuscirebbe mai a stabilire se la coppia successiva potrebbe o meno essere una soluzione. Per questo esempio particolare è possibile dimostrare che non vi sono soluzioni: ma tale dimostrazione richiede considerazioni di nuovo tipo, non può cioè essere ottenuta semplicemente mediante successive sostituzioni di interi nell'equazione.

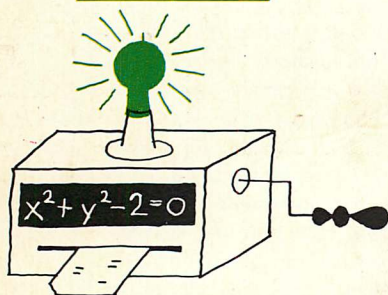
Un meccanismo, in grado di svolgere un procedimento del tipo suggerito da Hilbert, dovrebbe accettare come possibile ingresso (*input*) i coefficienti di una qualsivoglia equazione diofantea. Come possibile uscita (*output*) si accenderebbe una luce verde nel caso che l'equazione abbia una soluzione, una luce rossa in caso contrario. Una macchina di questo tipo potrebbe venir chiamata macchina di Hilbert. Viceversa, una macchina, che semplicemente ricerca soluzioni mediante successivi tentativi *ad infinitum*, potrebbe essere descritta nei termini di una macchina a luce verde. Se l'equazione possiede una soluzione, la luce verde si accenderà dopo un numero finito di passi. Se l'equazione non ha soluzioni, si avrà semplicemente che il calcolo proseguirà per sempre; a differenza della macchina di Hilbert, quella a luce verde non ha alcun modo per stabilire quando deve fermarsi.

È facile costruire una macchina a luce verde per le equazioni diofantee. Il problema è se sia possibile far di meglio e costruire una macchina di Hilbert, cioè una macchina a luce verde e rossa, che si fermi sempre dopo un numero finito di passi e fornisca una risposta definitiva, positiva o negativa. Ciò che ha dimostrato Matyasevich è che questo non potrà mai venir

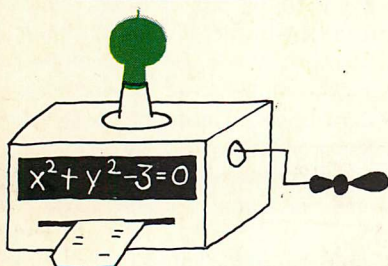


ce verde e rossa. In altri termini: se sia un insieme sia il suo complemento sono elencabili, allora l'insieme è computabile.

x	y
0	0
0	1
1	0
0	-1
-1	0
1	1
-1	1
1	-1
-1	-1



x	y
0	0
0	1
1	0
0	-1
-1	0
1	1
-1	1
1	-1
-1	-1
⋮	⋮



Le coppie di interi possono essere singolarmente saggiate da una macchina a luce verde per stabilire se siano soluzioni di equazioni diofantee. Nel caso dell'equazione $x^2 + y^2 - 2 = 0$, dopo i primi insuccessi, al sesto tentativo si giunge a una soluzione (in alto). Invece la macchina a luce verde che saggia l'equazione $x^2 + y^2 - 3 = 0$ non ha alcun modo di stabilire quando dovrà fermarsi dal momento che non vi sono soluzioni intere (in basso). Tutto ciò che può sapere è che per il momento non ha trovato soluzioni.

fatto. Anche se concediamo alle macchine una capacità di memoria e un tempo di calcolo illimitati, non potrà mai venir compilato alcun programma né costruita alcuna macchina che compiano ciò che era richiesto da Hilbert.

Nella sua allocuzione del 1900, Hilbert continuava in questo modo: «Tallora può capitare che si cerchi una soluzione sotto ipotesi insufficienti o in un senso non corretto ed è per questo che non si ha successo. Sorge allora il problema: dimostrare l'impossibilità della soluzione sotto le ipotesi date o nel senso stabilito.» Questo è esattamente ciò che è avvenuto con il decimo problema. Per spiegare come sia possibile sapere che non esiste alcuna macchina di Hilbert, dovremo svolgere alcune semplici considerazioni riguardo alla computabilità. Supponiamo che S indichi un insieme di interi. S è «elencabile» se può essere costruita una macchina a luce verde che risponda al seguente scopo: accetti come ingresso qualunque intero e come uscita accenda una luce verde dopo un numero finito di passi se e solo se l'ingresso (cioè l'intero) appartiene a S . Per esempio, l'insieme dei numeri pari è elencabile: in tal caso la macchina dovrà dividere l'ingresso per 2 e accendere la luce verde se il resto è 0. In matematica gli insiemi di questo tipo vengono detti ricorsivamente enumerabili; la parola «elencabile» è una nostra espressione non formale ma usata in senso equivalente.

L'insieme S è «computabile» se può essere costruita una macchina a luce verde e rossa (analoga alla macchina di Hilbert) con un compito più difficile: accettare ogni intero come ingresso e, dopo un numero finito di passi, accendere la luce verde se l'intero è in S , la luce rossa se l'intero non è in S . Per esempio, l'insieme dei numeri pari è anche computabile: la macchina dovrà dividere l'ingresso per 2 e se il resto è 0 accenderà la luce verde, quella rossa se invece è 1 (si vedano le illustrazioni a pagina 139).

Vi è uno stretto collegamento fra queste due definizioni. Per chiarire ciò, si supponga che \bar{S} indichi il complemento di S , cioè l'insieme di tutti gli interi che non appartengono a S . Se, nei due esempi dati, S è l'insieme degli interi pari, allora \bar{S} sarà l'insieme degli interi dispari. Si può dimostrare che se S è computabile allora sia S sia \bar{S} saranno elencabili. Per dirla in altri termini: se esiste per S una macchina a luce verde e rossa allora esistono una macchina a luce verde per S e una macchina a luce verde per \bar{S} . La dimostrazione è semplice. Per costruire una macchina a luce verde per S , basta svi-

tare la lampada rossa della macchina a luce verde e rossa. Per costruire una macchina a luce verde per \bar{S} , è sufficiente svitare la lampada verde della macchina di Hilbert e porla nello zoccolo in cui vi era la lampada rossa.

È vera anche l'affermazione inversa: se S e \bar{S} sono elencabili, allora S è computabile. In modo equivalente si può dire: se esiste una macchina a luce verde sia per S sia per \bar{S} , allora può essere costruita per S una macchina a luce verde e rossa. È facile fare questo: si sostituisca, nella macchina a luce verde per S , la lampada verde con una rossa e quindi si colleghino in parallelo le due macchine, di modo che uno stesso ingresso valga simultaneamente per entrambe. Il risultato è ovviamente una macchina a luce verde e rossa.

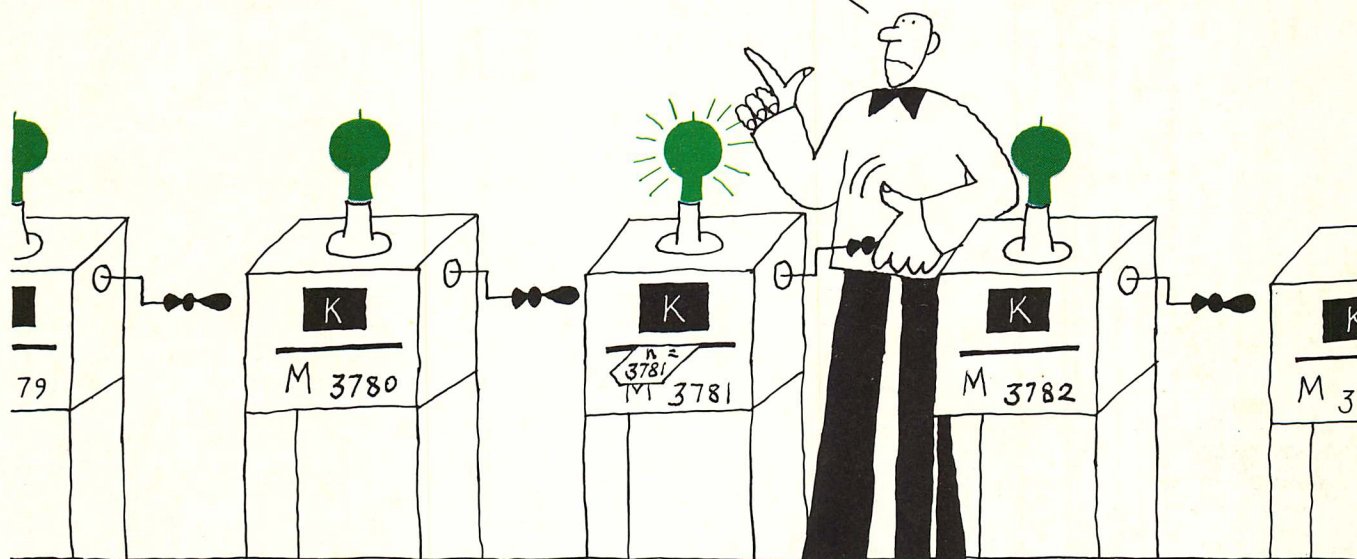
Fatte queste considerazioni, possiamo ora stabilire uno dei fatti cruciali della teoria della computabilità che ha un ruolo centrale nella soluzione del decimo problema di Hilbert: esiste un insieme K elencabile ma non computabile! Esiste cioè una macchina a luce verde per K ma non è possibile costruire una macchina a luce verde per \bar{K} , il complemento di K .

Al fine di dimostrare questo asserto a prima vista singolare, si supponga che ogni macchina a luce verde sia determinata da un dettagliato «libretto di istruzioni» redatto in lingua italiana che descriva esattamente il modo in cui è costruita la macchina. Questi libretti possono essere ordinati e numerati nella successione 1, 2, 3 e così via. In tal modo risultano numerate tutte le macchine a luce verde; M_1 sia la prima macchina, M_2 la seconda, eccetera. In questo ragionamento è celato un argomento piuttosto sottile: non sarebbe possibile avere un elenco ordinato di questo tipo per i libretti di istruzioni delle macchine a luce verde e rossa. La difficoltà consiste nel fatto che non è possibile stabilire sulla base del libretto se, per ogni ingresso posto nella macchina che gli corrisponde, si accenderà la luce verde o la luce rossa.

L'insieme K è definito come l'insieme dei numeri n tali che l' n -esima macchina si accende se le viene posto in ingresso il numero n stesso. In altre parole, il numero 1 appartiene a K se e solo se M_1 accende la propria luce verde quando «1» viene posto in ingresso: il numero 2 appartiene a K se e solo se M_2 giunge ad accendersi dopo che «2» le è stato posto in ingresso e così via (si veda l'illustrazione in alto della pagina a fronte).

Per costruire una macchina a luce verde per K occorre, oltre alla biblioteca dei libretti di istruzioni, un omet-

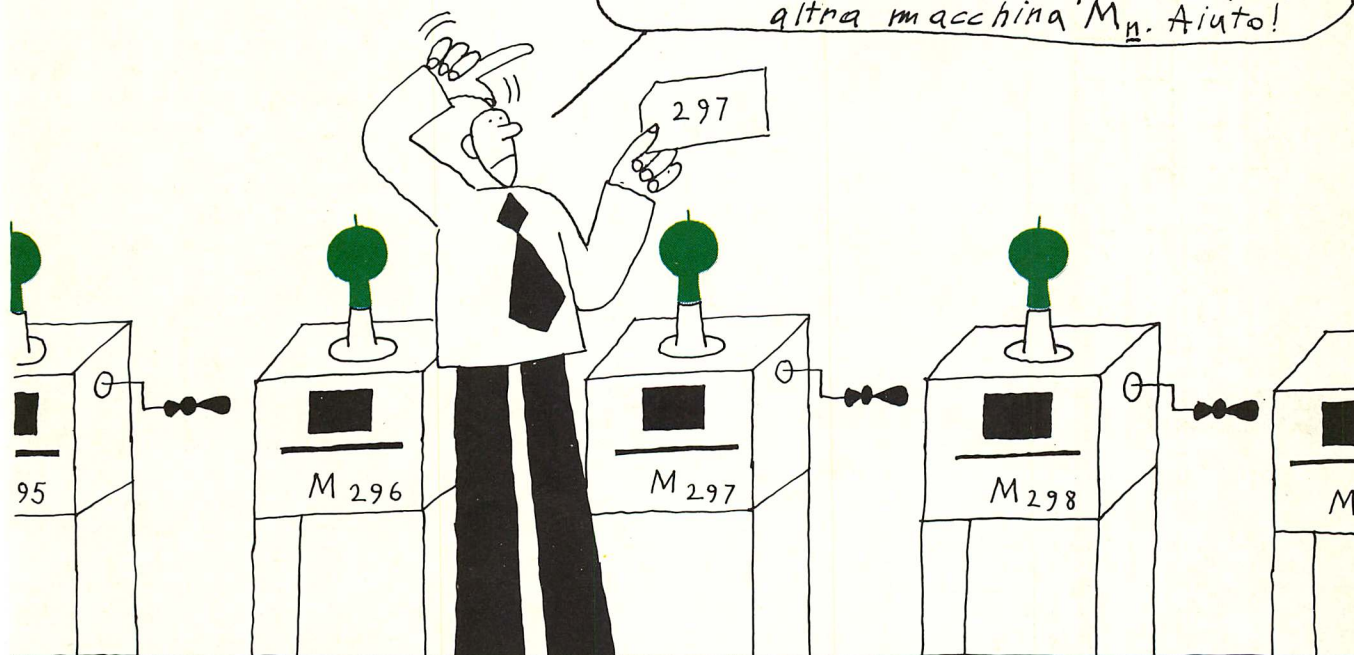
Scegli un numero, uno qualunque, per esempio 3781. K è definito come l'insieme dei numeri " n ", tali che l' n -esima macchina a luce verde si accende se le si pone in ingresso il numero n .



L'insieme K è elencabile, esiste cioè una macchina a luce verde per K . Si numerino tutte le possibili macchine a luce verde: M_1 sia la prima macchina, M_2 la seconda, M_3 la terza e così via fino all' n -esima macchina. K è l'insieme dei numeri n tali che

l' n -esima macchina si accende se le si pone in ingresso il numero n . Nell'illustrazione un ometto ha posto in ingresso nella macchina M_{3781} il numero 3781 e la luce verde si è accesa, indicando che il numero intero 3781 è un elemento dell'insieme K .

Vi è una macchina a luce verde per il complemento di K ? Se c'è, non può essere M_{297} , poiché in tal caso la definizione di K afferma che M_{297} si accenderà con 297 come ingresso. Analogamente non può essere alcuna altra macchina M_n . Aiuto!



L'insieme K non è computabile, non esiste cioè alcuna macchina a luce verde per \bar{K} , il complemento di K . Si supponga che vi sia una tale macchina a luce verde per \bar{K} . Dal momento che \bar{K} è il complemento di K , ne scende che per ogni ingresso, per esempio 297, tale macchina si accenderà se e solo se M_{297} non si accende per 297. Quindi la macchina per \bar{K} non è certamen-

te M_{297} . Per le stesse ragioni non potrà coincidere con nessuna M_n per ogni altro valore di n . Quindi non esiste alcuna macchina a luce verde per \bar{K} e cioè \bar{K} non è elencabile. Un insieme elencabile, il cui complemento non lo sia, non è computabile; quindi non si può costruire per esso alcuna macchina a luce verde e rossa, ossia non vi è nessun algoritmo per separare K da \bar{K} .

$$\begin{array}{rcl}
1. & 1 & \\
2. & 1 & \\
3. & 1 + 1 = 2 & \\
4. & 1 + 2 = 3 & \\
5. & 2 + 3 = 5 & \\
6. & 3 + 5 = 8 & \\
7. & 5 + 8 = 13 & \\
8. & 8 + 13 = 21 & \\
9. & 13 + 21 = 34 & \\
10. & 21 + 34 = 55 & \\
11. & 34 + 55 = 89 & \\
12. & 55 + 89 = 144 & \\
13. & 89 + 144 = 233 & \\
\vdots & \vdots & \\
n & \approx \frac{1}{\sqrt{5}} \left(\frac{1 + \sqrt{5}}{2} \right)^n &
\end{array}$$

I numeri di Fibonacci vennero scoperti nel 1202 da Leonardo Pisano. La successione è ottenuta iniziando con 1, quindi ancora 1 e poi sommando gli ultimi due numeri per ottenere il successivo. Tale successione cresce in modo esponenziale: l' n -esimo numero è in via approssimata proporzionale all' n -esima potenza del numero reale $[(1 + \sqrt{5})/2]^n$.

PROBLEMA: Trovare il più piccolo numero n che diviso per 10, 3, 7 e 11 abbia come resti rispettivamente 4, 2, 3 e 1.

SOLUZIONE: Sia x il numero cercato. «Res» sia un'abbreviazione per «Il resto di...». È possibile riformulare il problema scrivendo:

$$\begin{array}{ll}
\text{Res} \left(\frac{x}{10} \right) = 4 & \text{Res} \left(\frac{x}{7} \right) = 3 \\
\text{Res} \left(\frac{x}{3} \right) = 2 & \text{Res} \left(\frac{x}{11} \right) = 1
\end{array}$$

Per determinare x si devono risolvere quattro problemi ausiliari per delle nuove incognite y_1, y_2, y_3 e y_4 . In ogni caso il numeratore è ottenuto moltiplicando fra loro tre dei divisori mentre il denominatore è il quarto. Per esempio, nella prima equazione in y_1 il numeratore 231 è uguale $3 \times 7 \times 11$, mentre 10 è posto come denominatore:

$$\begin{array}{llll}
\text{Res} \left(\frac{231y_1}{10} \right) = 4, & y_1 < 10 & \text{Res} \left(\frac{330y_3}{7} \right) = 3, & y_3 < 7 \\
\text{Res} \left(\frac{770y_2}{3} \right) = 2, & y_2 < 3 & \text{Res} \left(\frac{210y_4}{11} \right) = 1, & y_4 < 11
\end{array}$$

L'insieme dei più piccoli interi che siano soluzioni di queste equazioni ausiliarie è $y_1 = 4, y_2 = 1, y_3 = 3$ e $y_4 = 1$. Per ottenere x (il numero cercato inizialmente) vengono sommati insieme i numeratori delle quattro equazioni ausiliarie:

$$\begin{aligned}
x &= (231y_1) + (770y_2) + (330y_3) + (210y_4) \\
&= (231 \times 4) + (770 \times 1) + (330 \times 3) + (210 \times 1) \\
&= 924 + 770 + 990 + 210 \\
&= 2894
\end{aligned}$$

Quindi 2894 è un possibile valore di x . Si può ottenere un numero minore sottraendo da questa soluzione il prodotto dei quattro divisori:

$$2894 - (10 \times 3 \times 7 \times 11) = 2894 - 2310 = 584$$

Pertanto 584 è la più piccola soluzione del problema.

to che possa leggerli ed eseguire ciò che vi è scritto. Potrebbe benissimo trattarsi di un anziano e saggio signore, ma dovrà essere una persona obbediente che faccia esattamente ciò che le si dice. Si fornisce all'ometto un numero, per esempio 3781; allora egli esaminerà attentamente il libretto n. 3781 e così facendo sarà in grado di costruire la macchina M 3781. Fatto ciò porrà in ingresso nella macchina a luce verde M 3781 l'intero 3781: se si accende la luce verde, il numero 3781 appartiene a K . Si è così ottenuta una macchina a luce verde per K .

Che cosa possiamo dire riguardo a \bar{K} ? Come possiamo essere sicuri che per esso non esista alcuna macchina a luce verde? Si supponga allora che vi sia una tale macchina. Dal momento che \bar{K} è il complemento di K , si avrà allora che comunque si fornisca un ingresso a tale macchina, per esempio 297, essa si accenderà se e solo se M 297 non si accende per 297 (se M 297 si accendesse, ciò significherebbe che l'intero 297 appartiene a K e non a \bar{K}). Quindi la macchina per \bar{K} non è certamente M 297 (si veda l'illustrazione in basso a pagina 143). Ma per la stessa ragione essa non coinciderà con nessuna M_n per qualunque altro valore di n . Lo stesso argomento si potrà infatti applicare, come è stato fatto per 297, anche per ogni altro numero, e ciò dimostra che nella biblioteca dei libretti di istruzioni non compare in alcun luogo una macchina a luce verde per \bar{K} . Dal momento che nel nostro elenco deve essere presente ogni possibile macchina a luce verde, ne segue che non può esistere alcuna macchina a luce verde per \bar{K} . Cioè, \bar{K} non è elencabile.

Il risultato è certo notevole e merita di essere meditato e valutato. Sappiamo perfettamente che cosa sia l'insieme K ; in linea di principio siamo in grado di ottenere, mediante la stampatrice di un elaboratore, tanti suoi elementi quanti ne vogliamo. Tuttavia, non vi potrà mai essere una procedura formale (un algoritmo o un programma di macchina) per separare K da \bar{K} . Ecco dunque un esempio di problema stabilito con precisione ma che non si risolverà mai con mezzi meccanici.

Che cosa ha a che fare tutto ciò con le equazioni diofantee? Si tratta semplicemente di questo: Matyasevich ha dimostrato che ogni insieme elencabile può essere fatto corrispondere a un'equazione diofantea. Più in particolare, se S è un insieme elencabile, vi è allora un polinomio P con coefficienti interi e con le variabili x, y_1, y_2, \dots, y_n che corrisponde a tale insieme e che

Il teorema cinese del resto è impiegato nella soluzione del decimo problema di Hilbert. In questo caso, il teorema è usato per trovare un numero che diviso per 10, 3, 7 e 11 dia come resti rispettivamente 4, 2, 3 e 1. La più piccola soluzione è l'intero 584.

indicheremo con $P_S(x, y_1, y_2, \dots, y_n)$. Ogni intero, per esempio 17, appartiene all'insieme S se e solo se l'equazione diofantea $P_S(17, y_1, y_2, \dots, y_n) = 0$ ammette una soluzione.

Si potrebbe pensare che per alcuni insiemi si debba far ricorso a polinomi molto complicati, ma ciò non avviene: non è mai necessario che il grado di P superi 4 né che il numero delle variabili y_1, y_2, \dots, y_n superi 14.

Dal risultato di Matyasevich scende facilmente la conclusione che non può esistere alcuna macchina di Hilbert. Si rammenti infatti l'insieme elencabile K costruito alcuni capoversi sopra. In base a quanto dimostrato da Matyasevich, vi dovrà essere un'equazione diofantea, $P_K(x, y_1, y_2, \dots, y_n) = 0$, associata a tale insieme. Se fosse possibile costruire una macchina di Hilbert, cioè una macchina a luce verde e rossa per saggiare le equazioni diofantee allo scopo di determinare se abbiano o meno soluzioni, allora per ogni intero x potremmo stabilire se esistano o meno degli interi y_1, y_2, \dots, y_n per cui l'equazione sopra indicata ammetta una soluzione. Tuttavia nello stabilire questo avremmo anche determinato se x appartiene o meno a K . In altri termini, una macchina di Hilbert applicata all'equazione diofantea che descrive K potrebbe venir usata come macchina a luce verde e rossa per K . D'altra parte abbiamo già dimostrato che K non è computabile, di modo che non può esistere alcuna macchina a luce verde e rossa per K . L'unico modo per uscire da questo dilemma è quello di concludere che non esiste alcuna macchina di Hilbert. In altre parole, il decimo problema di Hilbert non è risolvibile!

Il fatto che si possa associare una equazione diofantea a ogni insieme elencabile è un risultato positivo già di per sé di grande interesse, indipendentemente dalla sua applicazione al decimo problema di Hilbert. Un insieme di interi particolarmente importante e interessante è quello dei numeri primi. Un numero è primo se è divisibile solo per 1 e per se stesso: come esempi si considerino 2, 3, 5, 7, 11, 13 e 17. È abbastanza ovvio che tali numeri siano elencabili. Un algoritmo per elencarli ci è stato tramandato dai greci sotto il nome di « crivello di Eratostene ». Combinando il risultato di Matyasevich con una tecnica ideata da Putnam, si può ottenere un'equazione diofantea $Q(y_1, y_2, \dots, y_n) = z$ tale che un numero positivo z è primo se e solo se tale equazione ammette una soluzione intera positiva y_1, y_2, \dots, y_n . (La struttura precisa del polinomio Q è troppo complicata per essere qui descritta.)

Un altro notevole risultato può esse-

re ottenuto combinando il teorema di Matyasevich con i lavori di Gödel sull'indecidibilità. Comunque si consideri un sistema di assiomi di qualunque genere, dal quale si possano dedurre informazioni riguardo alle equazioni diofantee, si potrà sempre ottenere una particolare equazione diofantea con le seguenti proprietà: 1) l'equazione non ha soluzioni intere positive; 2) il fatto che essa non abbia soluzioni intere positive non può essere dedotto logicamente dall'insieme di assiomi dato. Naturalmente, una volta ottenuta l'equazione diofantea, si potrà ottenere un nuovo sistema di assiomi in base al quale si possa dimostrare che tale equazione non ha soluzioni. Ma allora questo nuovo insieme di assiomi darà origine a una nuova equazione diofantea per cui valgono le proprietà asserite.

Su che cosa si basa la dimostrazione del teorema di Matyasevich? Oltre ai risultati già menzionati della teoria classica dei numeri, ciò che ha un ruolo essenziale è l'asserto noto sotto il nome di teorema cinese del resto. È utile illustrare tale teorema con un esempio numerico.

Si supponga di voler trovare un numero che diviso per 10, 3, 7 e 11 abbia come resti rispettivamente 4, 2, 3 e 1 (si veda l'illustrazione in basso della pagina a fronte). Il teorema cinese del resto garantisce l'esistenza di un numero siffatto (in questo caso, per esempio, 584 è un numero di tal genere). Tutto quello che è richiesto affinché si possa applicare il teorema cinese del resto è che nessuna coppia dei divisori utilizzati abbia fattori comuni (eccettuata ovviamente l'unità). Vi può essere un qualunque numero di divisori e i resti richiesti possono essere interi positivi qualsivoglia.

Nel 1931 Gödel mostrò che si poteva far uso del teorema cinese del resto come metodo di codificazione, cioè per codificare con un solo numero qualunque sequenza finita di numeri. A partire dal numero di codice si può riottenere la sequenza nello stesso modo in cui nell'esempio dato si possono ottenere 4, 2, 3 e 1 a partire da 584, cioè come resti di divisioni successive. I divisori possono essere scelti in progressione aritmetica.

Il primo tentativo di dimostrazione dell'impossibilità di esistenza di una macchina di Hilbert, venne compiuta da uno di noi (Davis) nella sua tesi di dottorato nel 1950. La tecnica di Gödel basata sull'uso del teorema cinese del resto come metodo di codificazione venne applicata per associare a ogni elencabile S un'equazione diofantea $P_S(k, x, z, y_1, y_2, \dots, y_n) = 0$. Sfortu-

- I. $u + w - v - 2 = 0$
- II. $l - 2v - 2a - 1 = 0$
- III. $l^2 - lz - z^2 - 1 = 0$
- IV. $g - bl^2 = 0$
- V. $g^2 - gh - h^2 - 1 = 0$
- VI. $m - c(2h + g) - 3 = 0$
- VII. $m - fl - 2 = 0$
- VIII. $x^2 - mxy + y^2 - 1 = 0$
- IX. $(d - 1)l + u - x - 1 = 0$
- X. $x - v - (2h + g)(e - 1) = 0$

La soluzione di Matyasevich del decimo problema di Hilbert si basa su un'equazione diofantea ottenuta elevando al quadrato ognuna di queste dieci equazioni, sommandole insieme ed eguagliando a zero il polinomio complicato che ne risulta. In queste equazioni i valori u e v sono collegati nelle soluzioni in modo tale che v è il $2u$ -esimo numero di Fibonacci. Da tale soluzione segue che per ogni insieme elencabile esiste un'equazione diofantea a esso associata. Dal momento che esistono insiemi elencabili i cui complementi non lo sono, si avrà che non tutti gli insiemi elencabili possono avere una macchina a luce verde e rossa. Poiché avere una macchina di Hilbert per le equazioni diofantee equivale ad avere una macchina a luce verde e rossa per ogni insieme elencabile, il risultato di Matyasevich dimostra che non si può costruire alcuna macchina di Hilbert per saggiare le equazioni diofantee.

natamente la relazione che sussisteva fra un insieme e l'equazione a esso collegata risultava essere più complicata di quanto fosse richiesto dal decimo problema di Hilbert. Più precisamente la relazione era: un intero positivo x appartiene all'insieme S se e solo se per qualche valore intero positivo di z è possibile trovare una soluzione per ciascuna delle equazioni diofantee ottenute a partire da $P_S(k, x, z, y_1, y_2, \dots, y_n) = 0$ ponendo $k = 1$, poi $k = 2$ e così via fino a z . Anche se tale risultato sembrava vicino in modo assai promettente a ciò che era necessario, esso era solo un punto di partenza.

Circa nello stesso periodo Julia Robinson diede inizio alle sue ricerche sugli insiemi che potessero venir definiti per mezzo di equazioni diofantee e sviluppò molte tecniche ingegnose per trattare equazioni le cui soluzioni si comportassero in modo esponenziale (crescessero cioè come una potenza). Nel 1960 la Robinson, Davis e Putnam dimostrarono insieme un altro risultato. Utilizzarono cioè sia le ricerche della Robinson sia il risultato di Davis per dimostrare che a ogni insieme elencabile poteva corrispondere un'equazione diofantea in senso lato, cioè nel senso che si ammetteva che le variabili del-

l'equazione comparissero anche come esponenti (un esempio di un'equazione di questo tipo può essere $2^t + x^2 = z^3$). Davis, Robinson e Putnam, combinando le loro ricerche con alcuni risultati precedenti della Robinson, giunsero alla seguente conclusione: se fosse stato possibile trovare anche una sola equazione diofantea le cui soluzioni si comportassero in modo esponenziale in un senso ben precisato, si sarebbe potuto descrivere ogni insieme elencabile con un'equazione diofantea. Ciò a sua volta avrebbe dimostrato la insolubilità del decimo problema di Hilbert.

Occorsero dieci anni per trovare una equazione diofantea le cui soluzioni crescessero esponenzialmente nel modo precisato. Nel 1970 Matyasevich riuscì a trovare un'equazione di questo tipo facendo uso di quelli che sono noti come i numeri di Fibonacci. Que-

sti famosi numeri vennero scoperti nel 1202 da Leonardo Pisano, detto anche il Fibonacci. Egli li scoprì calcolando il numero totale delle coppie di discendenti di una coppia di conigli sotto la ipotesi che la coppia originaria e ogni successiva coppia di figli si riproducesse una volta al mese (a partire dal secondo mese di vita). La serie di Fibonacci si ottiene iniziando con 1, quindi ancora 1 e successivamente sommando i due numeri precedenti per ottenere il successivo: il primo numero di Fibonacci è 1, il secondo è 1, il terzo è $1 + 1 = 2$, il quarto è $1 + 2 = 3$, il quinto è $2 + 3 = 5$ e così via. La proprietà importante ai fini del decimo problema di Hilbert è il fatto che i numeri di Fibonacci crescono in modo esponenziale, ossia l' n -esimo numero di Fibonacci è approssimativamente proporzionale all' n -esima potenza di un certo numero reale fissato.

Se si fosse trovata un'equazione le cui soluzioni mettessero in relazione n con l' n -esimo numero di Fibonacci, essa sarebbe stata l'esempio cercato di equazione diofantea le cui soluzioni si comportano in modo esponenziale. Da tale esempio sarebbe scesa allora la soluzione del decimo problema di Hilbert. Ciò che Matyasevich fece fu proprio costruire un'equazione diofantea di questo tipo (*si veda l'illustrazione a pagina 145*). Dimostrando che l'insieme dei numeri di Fibonacci risultava in tal modo associato a un'equazione diofantea, si ricavava immediatamente dal teorema di Davis, Robinson e Putnam che per ogni insieme elencabile esisteva un'equazione a esso associata, e ciò in particolare anche per l'insieme K che non è computabile. E in questo modo si conclude la ricerca della soluzione del decimo problema di Hilbert.

Gli algoritmi

Un algoritmo è un insieme di regole in grado di fornire una specifica uscita come risposta a una specifica entrata. Ogni passo va definito con precisione per essere tradotto in un linguaggio per calcolatori

di Donald E. Knuth

Fino a dieci anni fa la parola «algoritmo» era sconosciuta alla maggior parte delle persone colte, e, in verità, era scarsamente necessaria. Lo sviluppo rapidissimo della scienza dei calcolatori, il cui punto focale è costituito dallo studio degli algoritmi, ha modificato la situazione: attualmente la parola è indispensabile. Esistono svariati altri termini che esprimono, almeno in parte, il concetto in questione: procedura, prescrizione, processo, *routine*, metodo. Al pari di queste cose un algoritmo è un insieme di regole o direttive atte a fornire una risposta specifica a una specifica entrata. Caratteristica distintiva degli algoritmi è la totale eliminazione delle ambiguità: le regole devono essere sufficientemente semplici e ben definite da poter essere eseguite da una macchina. Un'altra caratteristica fondamentale degli algoritmi è che devono sempre avere termine dopo un numero finito di passi.

Un programma è l'esposizione di un algoritmo in un linguaggio accuratamente definito. Quindi, il programma di un calcolatore rappresenta un algoritmo, per quanto l'algoritmo stesso sia un costrutto intellettuale che esiste indipendentemente da qualsiasi rappresentazione. Allo stesso modo, il concetto di numero due esiste nella nostra mente anche quando non sia espresso graficamente. Chiunque abbia steso un programma di calcolo è in grado di rendersi conto del fatto che un algoritmo debba essere esattamente definito, con un'attenzione ai particolari inusitata in altre attività.

Programmi per problemi numerici sono stati compilati fino dal 1800 a.C.,

quando i matematici babilonesi del tempo di Hammurabi stabilirono delle regole di risoluzione per alcuni tipi di equazioni. Le regole erano determinate come procedure passo-passo, applicate sistematicamente a esempi numerici particolari. La stessa parola «algoritmo» ha origine nel medio oriente, sebbene molto tempo dopo. Essa proviene dall'ultima parte del nome dello studioso persiano Abu Ja'far Mohammed ibn Mûsâ al-Khowârizmî, il cui testo di aritmetica (825 d.C. circa) esercitò una significativa influenza per molti secoli.

Tradizionalmente gli algoritmi erano applicati unicamente a problemi numerici; l'esperienza con i calcolatori ha tuttavia mostrato che i dati elaborati dai programmi possono virtualmente rappresentare qualsiasi cosa. Di conseguenza, l'attenzione della scienza dei calcolatori si è trasferita allo studio delle diverse strutture con cui si possono rappresentare le informazioni e all'aspetto ramificato o decisionale degli algoritmi, che permette di seguire differenti sequenze di operazioni in dipendenza dallo stato delle cose in un determinato istante. È questa la caratteristica che talvolta rende preferibili, per la rappresentazione e l'organizzazione delle informazioni, i modelli algoritmici a quelli matematici tradizionali. Per quanto gli algoritmi numerici possiedano molte interessanti caratteristiche, nella discussione che segue mi limiterò a esempi non-numerici, proprio per sottolineare il fatto che gli algoritmi si occupano innanzitutto della manipolazione di simboli, che non rappresentano necessariamente numeri.

La ricerca nella memoria di un calcolatore

Per illustrare come gli algoritmi possano venire studiati con profitto, prenderò in considerazione in modo abbastanza approfondito un semplice problema di ritrovamento di un'informazione. Il problema consiste nel determinare se una certa parola x appaia o meno nella memoria di un calcolatore. La parola x può essere il nome di una persona, il numero di un componente meccanico, un termine in qualche lingua straniera, un composto chimico, il numero di una carta di credito o qualsiasi altra cosa. Il problema è interessante solo quando l'insieme di tutte le possibili x è troppo grande perché il calcolatore possa esaminarlo all'istante; altrimenti sarebbe sufficiente predisporre una posizione in memoria per ogni parola.

Supponiamo che n differenti parole siano state immagazzinate nella memoria dell'elaboratore. Dobbiamo costruire un algoritmo che accetti in entrata la parola x e ci fornisca in uscita la posizione j in cui appare x . Quindi, se x è presente, otterremo come uscita un numero compreso tra 1 e n ; qualora invece x non si trovi nella memoria, avremo 0 in uscita, per indicare l'insuccesso della ricerca.

Naturalmente la risoluzione di questo problema non presenta alcuna difficoltà. L'algoritmo più semplice consiste nel collocare le parole nelle posizioni da 1 a n ed esaminare a turno ogni parola. Se x si trova nella posizione j , il calcolatore deve fornire in uscita j per poi fermarsi; ma se vengono esaurite tutte le n possibili

lità senza alcun esito, si deve fermare con 0 in uscita. Questa descrizione della strategia di ricerca sarebbe probabilmente troppo imprecisa per un elaboratore, ma può venire determinata in modo più accurato. La possiamo infatti scrivere come una sequenza di passi nel modo seguente:

Algoritmo A; ricerca sequenziale.

A1. [Inizio] Poni $j \leftarrow n$ (La freccia sta a indicare che il valore della variabile j è posto uguale a n , il numero delle parole della tabella da esaminare. Questo è il valore iniziale di j . I passi successivi dell'algoritmo faranno assumere a j la sequenza di valori $n, n-1, n-2$, fino a 0 o a una posizione contenente la parola in entrata x .)

A2. [Esito negativo?] Se $j=0$, j in uscita e l'algoritmo si arresta. (Altrimenti vai al passo A3.)

A3. [Esito positivo?] Se $x=KEY[j]$, j in uscita e l'algoritmo si arresta. (Il termine $KEY[j]$ indica la parola collocata nella posizione j .)

A4. [Iterazione] Poni $j=j-1$ (diminuisce di 1 il valore di j) e torna al passo A2.

Possiamo visualizzare la successione dei passi rappresentando l'algoritmo con uno schema a blocchi (si veda l'illustrazione nella pagina a fronte). Una delle ragioni per cui è importante specificare con la massima accuratezza i passi, è che l'algoritmo deve funzionare in ogni caso. La descrizione informale fatta più sopra potrebbe suggerire un algoritmo erroneo

che porti direttamente dal passo A1 al passo A3; un simile algoritmo sarebbe inadeguato per $n=0$ (quando non ci sia alcuna parola in memoria), poiché il passo A1 porrebbe $j=0$ e A3 ci farebbe considerare l'inesistente $KEY[0]$.

È interessante notare la possibilità di perfezionare l'algoritmo A assegnando un significato alla notazione $KEY[0]$, permettendo cioè che venga collocata una parola in «posizione 0», così come nelle posizioni da 1 a n . Se il passo A1 ponesse $KEY[0] \leftarrow x$ come $j \leftarrow n$, si potrebbe eliminare il passo A2, rendendo la ricerca più veloce del 20 per cento circa su molti calcolatori.

Sfortunatamente per i programmatori, i linguaggi dei calcolatori maggiormente usati (il COBOL e il FORTRAN stan-

ENTRATA: $x = \text{GRANT}$

1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16	17	18	19	20	21	22	23	24	25
WASHINGTON	ADAMS	JEFFERSON	MADISON	MONROE	JACKSON	VAN BUREN	HARRISON	TYLER	POLK	TAYLOR	FILLMORE	PIERCE	BUCHANAN	LINCOLN	JOHNSON	GRANT	HAYES	GARFIELD	ARTHUR	CLEVELAND	McKINLEY	ROOSEVELT	TAFT	WILSON

:: DOPO SETTE PASSI

1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16	17	18	19	20	21	22	23	24	25
WASHINGTON	ADAMS	JEFFERSON	MADISON	MONROE	JACKSON	VAN BUREN	HARRISON	TYLER	POLK	TAYLOR	FILLMORE	PIERCE	BUCHANAN	LINCOLN	JOHNSON	GRANT	HAYES	GARFIELD	ARTHUR	CLEVELAND	McKINLEY	ROOSEVELT	TAFT	WILSON

ENTRATA: $x = \text{GIBBS}$

1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16	17	18	19	20	21	22	23	24	25
WASHINGTON	ADAMS	JEFFERSON	MADISON	MONROE	JACKSON	VAN BUREN	HARRISON	TYLER	POLK	TAYLOR	FILLMORE	PIERCE	BUCHANAN	LINCOLN	JOHNSON	GRANT	HAYES	GARFIELD	ARTHUR	CLEVELAND	McKINLEY	ROOSEVELT	TAFT	WILSON

USCITA: 17

1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16	17	18	19	20	21	22	23	24	25
WASHINGTON	ADAMS	JEFFERSON	MADISON	MONROE	JACKSON	VAN BUREN	HARRISON	TYLER	POLK	TAYLOR	FILLMORE	PIERCE	BUCHANAN	LINCOLN	JOHNSON	GRANT	HAYES	GARFIELD	ARTHUR	CLEVELAND	McKINLEY	ROOSEVELT	TAFT	WILSON

:: DOPO 23 PASSI

1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16	17	18	19	20	21	22	23	24	25
WASHINGTON	ADAMS	JEFFERSON	MADISON	MONROE	JACKSON	VAN BUREN	HARRISON	TYLER	POLK	TAYLOR	FILLMORE	PIERCE	BUCHANAN	LINCOLN	JOHNSON	GRANT	HAYES	GARFIELD	ARTHUR	CLEVELAND	McKINLEY	ROOSEVELT	TAFT	WILSON

USCITA: 0

L'algoritmo a ricerca sequenziale (algoritmo A nel testo) ricerca una parola in entrata in una tabella i cui ingressi non siano disposti in alcun ordine particolare. Questa tabella è dotata di 25 ingressi, o chiavi: $KEY[1]$, $KEY[2]$, fino a $KEY[25]$. Ogni chiave è un nome di persona. Supponiamo che la parola in entrata sia il nome «Grant».

L'algoritmo A ricerca «Grant» confrontandolo dapprima con $KEY[25]$, che è «Wilson», quindi con $KEY[24]$, che è «Taft», ecc. Qui «Grant» è $KEY[17]$, e quindi l'algoritmo ci fornisce «17» (in alto). Se in entrata ci fosse stato «Gibbs», l'algoritmo avrebbe confrontato «Gibbs» con tutte le chiavi ottenendo 0 in uscita (in basso).

dard) non prevedono che lo 0 sia impiegato come indice per posizioni di memoria; quindi l'algoritmo espresso in uno di questi linguaggi, non può venire perfezionato tanto facilmente.

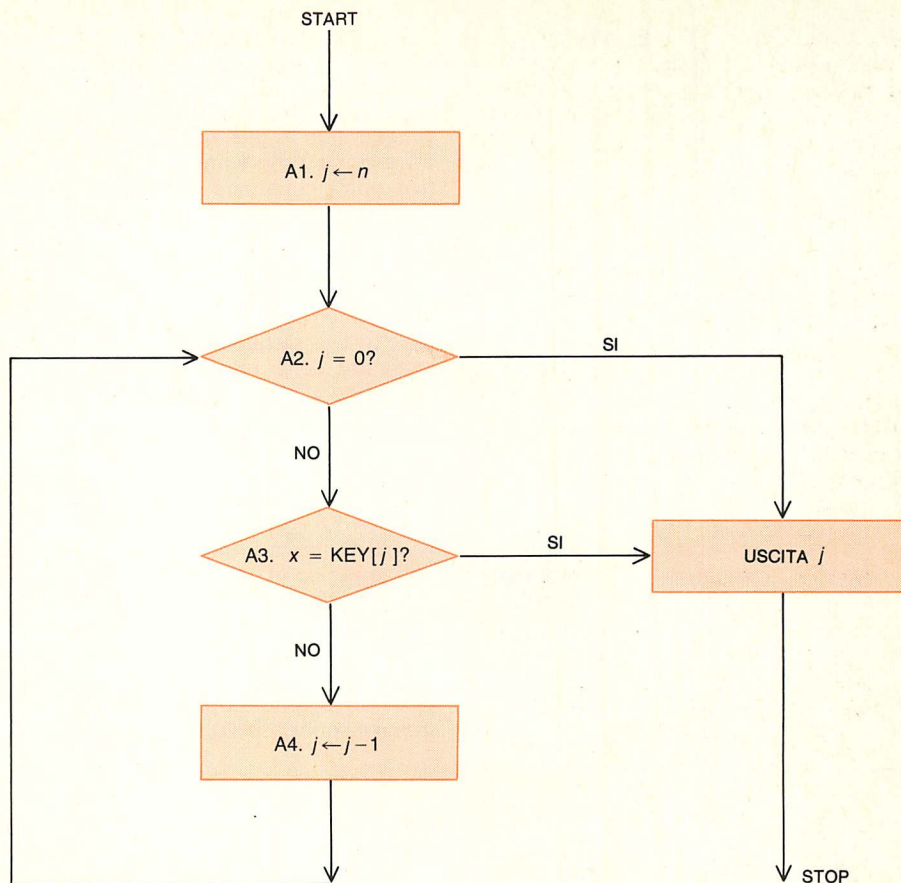
L'algoritmo A risolve certamente il problema della ricerca in una tabella di parole, ma la soluzione non è certo ottimale, a meno che il numero delle parole non sia decisamente esiguo, diciamo 25 o meno. Se n fosse pari a un milione, una semplice ricerca sequenziale costituirebbe una tecnica di consultazione insopportabilmente lenta. Difficilmente ci si sobbarcherebbe il costo della preparazione di una simile tabella, se non ci si aspettasse di consultarla frequentemente, e non si vorrebbe sprecare molto tempo per la ricerca. L'algoritmo A equivale alla ricerca di un numero telefonico fatta consultando un elenco degli abbonati pagina per pagina, colonna per colonna, una riga alla volta. Si può certamente fare di meglio.

L'utilità di un ordinamento

Può essere istruttivo considerare un elenco telefonico come esempio di una simile tabella di informazioni. Di fronte alla richiesta di trovare il numero di telefono di qualcuno che abiti al 1642 East della 56^a Strada, non ci sarebbe in effetti alternativa migliore alla ricerca sequenziale equivalente all'algoritmo A, dato che un normale elenco degli abbonati non è predisposto per la ricerca basata sull'indirizzo. D'altronde, quando si cerca un nome, è possibile trarre vantaggio dall'ordinamento alfabetico. Tale ordinamento comporta un vantaggio sostanziale poiché in ogni punto dell'elenco è sufficiente un semplice sguardo per eliminare da ogni ulteriore considerazione una grande quantità di nomi.

Se le parole di una tabella sono sistemate in un qualche ordine coerente, esistono svariati modi per determinare una efficiente procedura di ricerca. Quella più semplice inizia dall'esame dell'ingresso centrale della tabella. Se la parola cercata precede, alfabeticamente o numericamente, l'ingresso centrale, l'intera seconda metà della tabella può venire eliminata; analogamente, se x segue l'ingresso centrale, si potrà eliminare l'intera prima metà. Un unico confronto ci porta quindi a un problema di ricerca di dimensioni dimezzate rispetto a quello di partenza. Lo stesso procedimento può essere applicato anche alla parte restante della tabella, e così via fino a che la parola x non sia individuata. Questa procedura è nota come ricerca binaria.

Sebbene l'idea sottostante alla ricerca binaria sia semplice, è necessaria una certa attenzione per stenderne l'algoritmo. Tanto per cominciare, in una lista con un numero pari di elementi non esiste un unico ingresso centrale e, inoltre, non è immediatamente chiaro quando ci si debba fermare in caso di insuccesso. Gli insegnanti di scienza dei calcolatori hanno infatti notato che l'80 per cento circa degli studenti cui sia stato



Il diagramma a blocchi per l'algoritmo A illustra il percorso logico seguito dalla ricerca sequenziale per individuare una parola x in una tabella di n chiavi. L'algoritmo cerca x confrontandolo dapprima con $KEY[n]$, quindi con $KEY[n-1]$, con $KEY[n-2]$, ecc. Se x uguaglia $KEY[j]$, l'algoritmo emette j , la posizione in cui si trovava x . Se x non si trova nella tabella, l'algoritmo fornisce in uscita 0. La freccia al passo A1 ($j \leftarrow n$) significa «Poni j uguale a n » in quel passo. I vari passi sono illustrati nel testo. Per trovare x , l'algoritmo A esamina in media metà tabella. Nel caso peggiore, se cioè x è $KEY[1]$ o non è presente, l'algoritmo A esamina l'intera tabella.

richiesto per la prima volta di compilare un programma di ricerca binaria fornisce un risultato inesatto, anche quando abbia avuto più di un anno di esperienza di programmazione! Il lettore che ritenga di dominare sufficientemente bene gli algoritmi, ma non abbia mai compilato un programma per la ricerca binaria, può divertirsi a cercare di stenderne uno prima di leggere la soluzione riportata sotto.

Algoritmo B; ricerca binaria. La notazione è la stessa dell'algoritmo A; inoltre, si assume che nella relazione $<$ la prima parola, $KEY[1]$, preceda la seconda, $KEY[2]$, e questa la terza, e così fino all'ultima, $KEY[n]$. La condizione può essere espressa così: $KEY[1] < KEY[2] < \dots < KEY[n]$.

B1. [Inizio] Poni $l \leftarrow 0$, $r \leftarrow n + 1$ (La lettera l indica il limite sinistro della ricerca, r quello destro. Più precisamente, $KEY[j]$ non può essere il riscontro della parola data x , a meno che la posizione j non sia contemporaneamente maggiore di l e minore di r .)

B2. [Ricerca del punto centrale] Poni $j \leftarrow \lfloor (l+r)/2 \rfloor$ (le parentesi $\lfloor \rfloor$ stanno per «Arrotonda per difetto all'intero più vicino». Cioè, se $(l+r)$ è pari, j è uguale a $(l+r)/2$; se $(l+r)$ è dispari, j è uguale a $(l+r-1)/2$.)

B3. [Esito negativo?] Se $j = l$, 0 in uscita e l'algoritmo si arresta. (Se j uguaglia l , allora r deve essere pari a $l+1$, poiché r è comunque maggiore di l ; ne segue che x non corrisponde ad alcuna chiave della tabella.)

B4. [Confronto.] (A questo punto $j > l$ e $j < r$) Se $x = KEY[j]$, j in uscita e l'algoritmo si arresta. Se $x < KEY[j]$, poni $r \leftarrow j$ e ritorna a B2. Se $x > KEY[j]$, poni $l \leftarrow j$ e torna a B2.

L'illustrazione della pagina seguente mostra una rappresentazione, passo per passo, della ricerca operata dall'algoritmo B in una lista di 25 parole.

Appare evidente che la ricerca binaria (algoritmo B) è decisamente migliore di quella sequenziale (algoritmo A), ma è possibile un ulteriore perfezionamento? E quando l'algoritmo sarà ottimale? Per rispondere a queste domande si impone un'analisi di tipo quantitativa.

Analisi quantitativa

Analizziamo, per prima cosa, i casi peggiori degli algoritmi A e B. Quanto tempo impiega ogni algoritmo a trovare la parola x in un elenco di lunghezza n ? Per l'algoritmo A la risposta è semplice: se x è uguale a $KEY[1]$, o addirittura

non è presente nella lista, occorreranno n applicazioni del passo A3; vale a dire, la parola in questione deve essere confrontata con tutti gli n ingressi della tabella, prima che la ricerca si concluda. Inoltre, l'algoritmo non eseguirà il passo A3 più di n volte. Se la ricerca sequenziale viene applicata a una tabella con un milione di ingressi, verranno eseguiti, nel caso peggiore, un milione di confronti.

La risposta è di poco più complessa nel caso della ricerca binaria. Poiché l'algoritmo B scarta metà della tabella rimanente, dopo ogni esecuzione del passo B4, all'inizio opera con l'intera lista, quindi con metà, un quarto, un ottavo, e via dicendo. Il numero massimo di applicazioni del passo B4 sarà quindi k , dove k è il più piccolo intero tale che 2^k sia maggiore di n . Così, qualora la ricerca binaria sia applicata a una tabella con un milione di ingressi (10^6), k sarà uguale a 20. Infatti 2^{20} è maggiore di 10^6 , ma 10^6 è maggiore di 2^{19} . Consultando una tabella con 10^6 ingressi, dovremo operare al più 20 confronti.

Prendendo le mosse dal risultato ottenuto per il caso peggiore, si può arrivare

ad affermare che l'algoritmo B non solo costituisce un buon metodo di ricerca, ma che, anzi, è il migliore di tutti i possibili algoritmi che procedano unicamente comparando x alle chiavi della tabella. La ragione è che è un siffatto algoritmo non potrà esaminare più di $2^k - 1$ differenti chiavi nei primi k confronti. Indipendentemente dalla strategia adottata, il primo confronto avverrà sempre con una particolare chiave della tabella, e il secondo avrà a che fare con al massimo altre due (a seconda che x sia minore o maggiore della prima chiave), il terzo confronto al massimo con altre 4 chiavi, il quarto con otto, e via di questo passo. Se quindi un algoritmo di ricerca per confronti non opera più di k confronti, la tabella non potrà contenere più di $1 + 2 + 4 + 8 + \dots + 2^{k-1}$ chiavi distinte, e il valore di questa somma è appunto $2^k - 1$.

Il popolare «Gioco delle venti domande» può essere utilizzato allo stesso modo. In questo passatempo un giocatore pensa a un oggetto di cui scrive, non visto, il nome su un foglietto. Gli altri giocatori devono scoprire di che oggetto si tratti, ponendo 20 domande che am-

mettano come risposta unicamente un sì o un no. All'inizio del gioco viene detto se l'oggetto misterioso sia un animale, un vegetale, un minerale, o una combinazione di simili attributi, che si suppongono ben definiti. Sulla falsariga di quanto fatto nel precedente paragrafo, si può dimostrare che gli avversari del primo giocatore non possono identificare correttamente più di 2^{23} oggetti differenti, quale che sia l'ingegnosità delle domande. Ci sono solo 2^3 (otto) possibili sottoinsiemi di attributi animali, vegetali e minerali, e ci sono solo 2^{20} possibili risposte ai venti quesiti. Il numero totale di oggetti eventualmente identificabili sarà allora 2^{23} . L'argomentazione resta valida anche quando le domande siano in relazione alle risposte ottenute precedentemente.

In altri termini, per identificare più di 2^{23} oggetti distinti, 20 domande non sono sufficienti. Il problema della ricerca è molto simile, per quanto non identico, dato che l'algoritmo non si limita a quesiti del tipo «sì-no». Le domande poste da algoritmi del tipo considerato ammettono tre possibili risultati: $x < \text{KEY}[j]$,

ENTRATA: **x = GRANT**

$l = 0$	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16	17	18	19	20	21	22	23	24	25	$r = 26$
	ADAMS	ARTHUR	BUCHANAN	CLEVELAND	FILLMORE	GARFIELD	GRANT	HARRISON	HAYES	JACKSON	JEFFERSON	JOHNSON	LINCOLN	MADISON	McKINLEY	MONROE	PIERCE	POLK	ROOSEVELT	TAFT	TAYLOR	TYLER	VAN BUREN	WASHINGTON	WILSON	
$l = 0$	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16	17	18	19	20	21	22	23	24	25	
	ADAMS	ARTHUR	BUCHANAN	CLEVELAND	FILLMORE	GARFIELD	GRANT	HARRISON	HAYES	JACKSON	JEFFERSON	JOHNSON	LINCOLN	MADISON	McKINLEY	MONROE	PIERCE	POLK	ROOSEVELT	TAFT	TAYLOR	TYLER	VAN BUREN	WASHINGTON	WILSON	
	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16	17	18	19	20	21	22	23	24	25	
	ADAMS	ARTHUR	BUCHANAN	CLEVELAND	FILLMORE	GARFIELD	GRANT	HARRISON	HAYES	JACKSON	JEFFERSON	JOHNSON	LINCOLN	MADISON	McKINLEY	MONROE	PIERCE	POLK	ROOSEVELT	TAFT	TAYLOR	TYLER	VAN BUREN	WASHINGTON	WILSON	
	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16	17	18	19	20	21	22	23	24	25	
	ADAMS	ARTHUR	BUCHANAN	CLEVELAND	FILLMORE	GARFIELD	GRANT	HARRISON	HAYES	JACKSON	JEFFERSON	JOHNSON	LINCOLN	MADISON	McKINLEY	MONROE	PIERCE	POLK	ROOSEVELT	TAFT	TAYLOR	TYLER	VAN BUREN	WASHINGTON	WILSON	
	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16	17	18	19	20	21	22	23	24	25	
	ADAMS	ARTHUR	BUCHANAN	CLEVELAND	FILLMORE	GARFIELD	GRANT	HARRISON	HAYES	JACKSON	JEFFERSON	JOHNSON	LINCOLN	MADISON	McKINLEY	MONROE	PIERCE	POLK	ROOSEVELT	TAFT	TAYLOR	TYLER	VAN BUREN	WASHINGTON	WILSON	

USCITA: **7**

L'algoritmo a ricerca binaria (algoritmo B nel testo), costituisce un sostanziale miglioramento rispetto alla ricerca sequenziale, quando la tabella da esaminare sia grande. Gli ingressi della tabella devono venire preventivamente ordinati. Qui i 25 nomi sono elencati in ordine alfabetico. La parola cercata x sia ancora «Grant». L'algoritmo confronta «Grant» dapprima con la chiave che si trova nella posizione centrale, j , della tabella. Il valore iniziale di j è ottenuto ponendo il limite sinistro l della ricerca pari a 0, e quello destro r pari a $n + 1$: in questo caso r è 26. l e r vanno sommati, dividendo quindi per 2 e arrotondando per difetto all'intero più vicino, qualora il risultato non sia un intero. Il punto centrale j della tabella è $26/2$, cioè 13, che è la posizione di «Lincoln» (in alto). Poiché «Grant» precede alfabeticamente

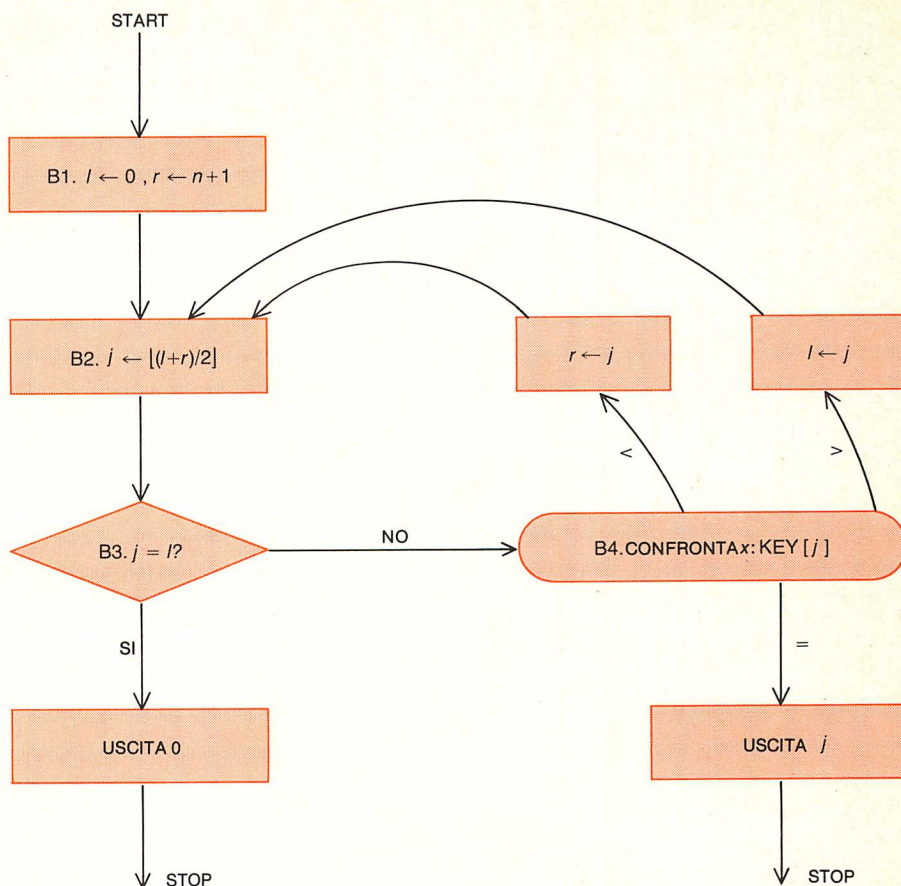
«Lincoln», l'algoritmo scarta la metà di destra della tabella, contenente solo nomi alfabeticamente precedenti, o al più coincidenti con «Lincoln». Per la metà rimanente l'algoritmo calcola un nuovo punto centrale, ponendo r uguale alla posizione j appena esaminata, che è 13 (seconda riga dall'alto). Il nuovo punto medio j è $(0 + 13)/2$, arrotondato per difetto a 6, la posizione di «Garfield». «Grant» segue «Garfield», sicché viene scartato il quarto di tabella di sinistra, ponendo il limite sinistro $l = 6$ (seconda riga dal basso). Ripetendo la procedura, si trova «Grant» in posizione 7 (in basso). Se x fosse stato «Gibbs», l'algoritmo B avrebbe eseguito ancora un passo, con $l = 6$ e r posto uguale a 7. La coincidenza del punto centrale $j = 6$ con il limite sinistro avrebbe indicato che «Gibbs» non è in tabella.

$x = \text{KEY}[j]$ oppure $x > \text{KEY}[j]$. Se una tabella è dotata di 2^k o più ingressi, il precedente ragionamento mostra che k confronti di x con le chiavi della tabella non sono sempre sufficienti. Ogni algoritmo che compie una ricerca in una lista di un milione di parole operando solo confronti, dovrà in alcuni casi esaminare 20 o più di queste parole. In breve, nel peggiore dei casi la ricerca binaria dà il risultato migliore.

L'analisi del comportamento degli algoritmi nel caso più sfavorevole non può però esaurire la nostra indagine: sarebbe eccessivamente pessimistico prendere delle decisioni basandosi unicamente sulla conoscenza di quanto di peggio possa accadere. Una più approfondita comprensione dei rispettivi meriti degli algoritmi A e B può essere raggiunta analizzando il comportamento in un caso intermedio. Se ognuna delle n chiavi di una tabella ha le stesse probabilità di essere esaminata, qual è il numero di confronti che si renderà necessario? Nel caso della ricerca sequenziale (algoritmo A), la risposta è data dalla media $(1+2+3+\dots+n)/n$, che è uguale a $(n+1)/2$. In altre parole, per trovare x dovremmo consultare in media metà della tabella. Per determinare il numero medio di confronti necessari per trovare x con l'aiuto dell'algoritmo B (ricerca binaria) il calcolo è di poco più complicato. La risposta è infatti $k - [(2-k-1)/n]$ dove k è il numero di confronti richiesto nel caso più sfavorevole. Per valori di n molto elevati il risultato è approssimativamente uguale a $k-1$; nel caso medio, dunque, l'algoritmo B necessita di un solo confronto in meno rispetto al caso maggiormente sfavorevole. Estendendo debitamente l'argomentazione, si può dimostrare che la ricerca binaria costituisce il migliore degli algoritmi possibili, anche dal punto di vista del caso medio: ogni algoritmo di ricerca deve operare in media almeno $k - [(2^k-k-1)/n]$ confronti, e ancor più nel caso peggiore.

Al di là del meglio

Non appena una cosa viene dimostrata impossibile a compiersi, una quantità di gente si mette al lavoro per cercare di realizzarla comunque. Questa sembra una componente essenziale del comportamento umano. Ho appena finito di mostrare che la ricerca binaria costituisce la condotta ottimale per compiere ricerche nella memoria di un calcolatore, ed ecco che cerco un metodo migliore. In primo luogo, quando il numero delle parole è esiguo, l'algoritmo A si rivela in pratica migliore dell'algoritmo B. Ma ciò non sarebbe in contraddizione con la dimostrazione precedente? La ragione sta nell'aver paragonato gli algoritmi A e B unicamente in base al numero di confronti operati. Di fatto l'algoritmo A richiede una minore attività di calcolo, cosicché una macchina impiegherà un tempo minore per eseguire ogni confronto. Con un normale calcolatore saranno necessarie in media $2n+6$ unità di tempo,



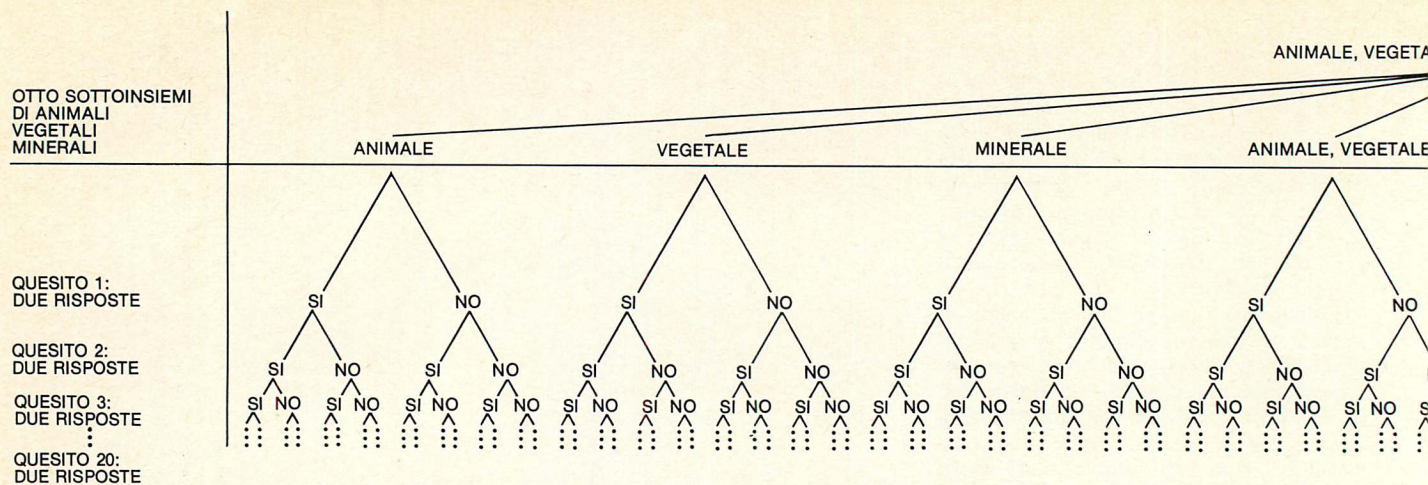
Il diagramma a blocchi per l'algoritmo B illustra le regole preposte alla ricerca binaria. L'algoritmo ricerca la parola in entrata x in una tabella di n chiavi precedentemente ordinata. La x viene confrontata con l'ingresso centrale della tabella. Se x è maggiore ($>$) dell'ingresso centrale, viene confrontato con il punto di mezzo della metà di destra della tabella. Se invece x è minore ($<$) del punto centrale, viene confrontata con il punto di mezzo della metà di sinistra. Il processo continua, scartando ogni volta la metà della tabella rimanente, finché non si rinvenga x oppure se ne constati l'assenza della tabella. I simboli (L...) hanno il significato «arrotonda per difetto all'intero più vicino». L'algoritmo B viene dettagliatamente analizzato a pagina 149.

circa, per una lista di lunghezza n . L'algoritmo B, d'altra parte, richiederà, sotto le stesse condizioni, $12\log_2 n - 11 + 12(k+1)/n$ unità di tempo, circa. Se non sono da sondare più di 20 chiavi l'algoritmo A sarà allora preferibile all'algoritmo B. Questo numero può variare leggermente da calcolatore a calcolatore, ma ciò basta a mostrare che l'efficienza di un algoritmo non può essere determinata semplicemente dal numero di confronti.

C'è ancora un'altra ragione per cui l'algoritmo B può venire migliorato. Quando consultiamo un elenco telefonico, cercando un nome x , e lo confrontiamo con quelli di una pagina, il nostro comportamento successivo non sarà influenzato solo dal fatto che il nome cercato sia alfabeticamente precedente o seguente a quelli della pagina, ma dipenderà anche da quanto questa sia precedente o seguente, in modo da saltare molte pagine nel caso pensassimo di esserne molto distanti. La dimostrazione precedente dell'ottimalità della ricerca binaria non si applica ad algoritmi che facciano uso di procedimenti quali il grado di diversità sussistente tra x e una particolare chiave. Nel caso delle venti doman-

de vale un discorso analogo: i giocatori potrebbero notare la lunghezza della parola segreta, mentre questa viene scritta, oppure potrebbero ricavare delle informazioni dal tempo impiegato dall'avversario per rispondere «sì» o «no».

Chi sia interessato all'efficienza, non necessariamente consulterà un elenco telefonico partendo da metà, come farebbe un calcolatore: il metodo di interpolazione, con l'aiuto dell'ordinamento alfabetico, da tempo di uso comune, dà migliori risultati, a dispetto delle dimostrazioni per cui la ricerca binaria sarebbe la migliore. Andrew C. Yao del MIT e F. Frances Yao della Brown University hanno recentemente dimostrato che il numero medio delle unità di tempo necessarie a un algoritmo di ricerca con interpolazione per accedere alla tabella è $\log_2 n$, a cui va sommata al massimo una piccola costante, a patto che gli ingressi della tabella siano numeri casuali indipendenti e uniformemente distribuiti. Per n molto grande, $\log_2 n$ è decisamente minore di $\log_2 n$, per cui una ricerca con interpolazione risulterà molto più rapida di quella binaria. L'idea sottostante alla dimostrazione di Yao è che ogni iterazione di una ricerca con



Il gioco delle venti domande è un passatempo nel quale un giocatore pensa a un oggetto, che descrive come animale, vegetale o minerale, o come una combinazione di questi attributi. Gli avversari devono cercare di indovinare di che oggetto si tratti, ponendo venti domande

alle quali si risponde con un «sì» o con un «no». Si dimostra che i giocatori non possono identificare correttamente più di 2^{23} oggetti, cioè 8 388 608. Ciò perché l'insieme costituito dai 3 attributi animale, vegetale e minerale, ha solo otto, 2^3 , possibili sottoinsiemi, (compre-

interpolazione tenda a ridurre l'incertezza della posizione di x da n alla radice quadrata di n . Resta così dimostrato che la ricerca con interpolazione è, in un senso molto generale, la migliore possibile: ogni algoritmo che compia una ricerca in una tabella casuale del tipo descritto deve esaminare, operando gli appropriati confronti, approssimativamente $\log_2 \log_2 n$ ingressi, in media.

La rilevanza teorica di questo risultato è grande; l'esperienza sui calcolatori ha tuttavia mostrato che una ricerca per interpolazione non costituisce, in pratica, un perfezionamento della ricerca binaria. Ciò perché i dati immagazzinati in una tabella non sono in generale sufficientemente casuali per conformarsi alla assunzione di distribuzione uniforme; inoltre n è solitamente tanto piccolo che i calcoli aggiuntivi richiesti da ogni confronto comportano una perdita di tempo superiore a quello risparmiato con la diminuzione dei confronti. La semplicità è una delle virtù della ricerca binaria, ed è importante mantenere un giusto rapporto tra teoria e pratica.

Ricerca ad albero binario

La ricerca binaria può peraltro essere perfezionata, lasciando cadere l'assunzione che ogni chiave della tabella abbia le stesse probabilità di venire cercata. Quando si supponga che certe chiavi possano venir consultate più spesso di altre, un algoritmo efficiente le esaminerà per prime.

Prima di approfondire questo concetto può essere utile guardare alla ricerca binaria da un'altro punto di vista. Consideriamo le 31 parole della lingua inglese di uso più comune (secondo quanto riporta Helen Fouché Gaines nel suo libro *Cryptanalysis*). Se queste vengono ordinate alfabeticamente nelle posizioni KEY[1], KEY[2], KEY[3], ..., KEY[31] della tabella, l'algoritmo B confronterà la parola cercata x con il punto centrale

KEY[16], che è la parola «I»; se x è alfabeticamente precedente rispetto a «I», il confronto successivo avverrà con KEY[8], la parola «by»; qualora sia invece seguente, verrà esaminata KEY[24]: «that». In altri termini, l'algoritmo B agisce sulla tabella ricalcando una struttura che appare simile a un albero capovolto, partendo dalla sommità e muovendosi verso il basso, a sinistra quando x sia precedente, a destra quando sia seguente (si veda l'illustrazione in alto a pagina 154). Non è difficile accorgersi che ogni algoritmo progettato per compiere una ricerca in una tabella ordinata operando semplicemente dei confronti può venire descritto da un albero binario di questo tipo.

L'albero di ricerca binaria è definito implicitamente nell'algoritmo B dalle operazioni aritmetiche su l , r , e j . È possibile definirlo esplicitamente immagazzinando nella tabella stessa le informazioni sull'albero. A questo scopo sia LEFT[j] la posizione della tabella da consultare se $x < \text{KEY}[j]$, e sia RIGHT[j] se $x > \text{KEY}[j]$. Così, in un elenco di 31 parole si avrebbe LEFT[16] uguale a 8 e RIGHT[16] uguale a 24, poiché la ricerca parte da KEY[16] e procede con KEY[8] oppure con KEY[24]. Se la ricerca termina con un insuccesso dopo che si sia appurato che la parola x precede o segue KEY[j], poniamo rispettivamente LEFT[j] o RIGHT[j] pari a 0. Nelle illustrazioni di pagina 16 questi 0 sono rappresentati dai piccoli quadrati che costituiscono i nodi alla base dell'albero.

La posizione della prima parola esaminata è detta comunemente «radice»; nell'esempio delle 31 parole la radice è 16. È possibile costruire algoritmi di ricerca che non abbiano inizio esaminando KEY[16], e questi, se qualche parola è consultata più frequentemente delle altre, possono essere anche più efficienti dell'algoritmo B. Una procedura di ricerca ad albero generalizzata è costituito dal seguente algoritmo:

Algoritmo C; ricerca ad albero binario.

C1. [Inizio.] Uguaglia j alla posizione della radice dell'albero binario.

C2. [Esito negativo?] Se $j=0$, j in uscita, e l'algoritmo si arresta.

C3. [Confronto.] Se $x = \text{KEY}[j]$, j in uscita e l'algoritmo si arresta. Se $x < \text{KEY}[j]$, poni $j \leftarrow \text{LEFT}[j]$ e vai al passo C2. Se $x > \text{KEY}[j]$, poni $j \leftarrow \text{RIGHT}[j]$ e torna al passo C2.

L'algoritmo C assomiglia a un testo per l'istruzione programmata in cui, a seconda delle risposte fornite a opportuni quesiti, viene indicata al lettore la pagina con cui proseguire. Esso funziona con qualsiasi albero binario in cui tutte le chiavi accessibili da LEFT[j] precedano KEY[j], mentre lo seguano quelle raggiungibili da RIGHT[j], e ciò per tutte le posizioni j . Un albero siffatto è detto di ricerca binaria.

Uno dei vantaggi dell'algoritmo C sull'algoritmo B è l'assenza di calcoli numerici, il che rende più rapida la ricerca sul calcolatore. Ma il vantaggio principale dell'algoritmo C è dovuto alla sua particolare flessibilità, determinata dalla struttura ad albero: gli ingressi della tabella possono ora venire risistemati in qualsivoglia ordine. Non è più necessario che KEY[1] preceda KEY[2], e così via fino a KEY[n]: fintanto che LEFT e RIGHT definiscono un albero di ricerca binario valido, la posizione effettiva delle chiavi nella tabella è irrilevante. Ciò significa che possiamo aggiungere nuovi ingressi alla tabella, senza dover spostare gli altri. Così, la parola «has» può venire addizionata all'albero binario di ricerca delle 31 parole inglesi più comuni, semplicemente ponendo KEY[32] ← «has», e cambiando RIGHT[j] da 0 a 32 quando j sia la posizione della chiave «has». Si potrebbe pensare che le aggiunte all'estremità dell'albero potrebbero squilibrarne la struttura, si dimostra invece che, se i nuovi ingressi sono aggiunti in ordine casuale, il risultato sarà equilibrato.

dendo l'insieme vuoto \emptyset per un oggetto che non goda di alcuna di tali caratteristiche). Tali otto possibilità si combinano poi con le sole 2^{20} possibili risposte alle 20 domande sì-no. Si può usare un'argomentazione del tutto simile per mostrare che un algoritmo che ponga al

massimo 20 quesiti «minore-uguale-maggiore», non può distinguere più di $2^{20}-1$ differenti chiavi, poiché $1+2+4+8+\dots+2^{19}=2^{20}-1$. Poiché la ricerca binaria è in grado di raggiungere questo limite massimo, è il più efficiente algoritmo di ricerca di questo tipo.

Alberi di ricerca binari ottimali

Poiché l'algoritmo C si applica a qualsiasi albero di ricerca binaria, è possibile adattarlo alle proprie esigenze, in modo da esaminare per prime le chiavi di maggior consultazione. Una tale modifica riduce il tempo medio impiegato dall'elaboratore per completare la ricerca. L'illustrazione in basso della pagina seguente mostra l'albero di ricerca ottimale per le 31 parole inglesi più comuni. Il numero medio di confronti necessari per trovare x in questo albero di ricerca binario ottimale è solo 3,437, mentre in un albero di ricerca binario equilibrato tale numero sale a 4,393. Vale la pena di osservare che l'albero ottimale, basato sulla frequenza delle parole, non parte dal confronto di x con la parola «the». Sebbene «the» sia di gran lunga la parola inglese più comune, è troppo avanti nell'ordinamento alfabetico e troppo lontana dal centro della lista per servire da radice ottimale.

Dal punto di vista matematico tradizionale è addirittura banale identificare l'albero binario ottimale per un particolare insieme di n parole e frequenze, poiché di tali alberi ne esiste solo un numero finito. In linea di principio è sufficiente compilarne una lista esaustiva e scegliere quello che funziona meglio. In realtà questa è una strada impraticabile, visto che il numero di alberi costruibili con n elementi è $(2n)!/n!(n+1)!$, dove $n!$ indica il prodotto $1 \times 2 \times 3 \times \dots \times n$. Questa formula rivela l'enorme numero di alberi binari esistenti, approssimativamente $4^n/\sqrt{\pi n^3}$. Per esempio, se n è 31, il numero totale degli alberi binari è 14 544 636 039 226 909, e ognuno di questi 14 trilioni di alberi sarà ottimale per qualche particolare insieme di frequenze assegnato alle 31 parole. Com'è allora possibile stabilire che l'albero scelto sia il migliore per le frequenze indicate dalla Gaines? Anche il più rapido dei moderni calcolatori non è abbastanza veloce per

esaminare singolarmente tutte le possibilità: se venisse esaminato un albero per microsecondo, occorrerebbero 460 anni.

Disponiamo però di un importante principio che rende possibile il calcolo: ogni sottoalbero di un albero ottimale deve essere esso stesso ottimale. Nell'albero ottimale di ricerca binaria per il problema delle 31 parole, il sottoalbero alla sinistra della parola «of» deve rappresentare il miglior metodo di ricerca per le venti parole da «a», «and» e così via fino a «not». Se esistesse un metodo migliore, ci dovrebbe condurre a un albero complessivamente migliore, da cui seguirebbe la non-ottimalità dell'albero in esame. Analogamente, il sottoalbero, ancora più piccolo alla destra di «for» deve rappresentare il miglior metodo di ricerca per le 11 parole da «from», «had», fino a «not». Ogni sottoalbero corrisponde a un insieme di parole consecutive $KEY[i], KEY[i+1] \dots KEY[j]$, con $1 \leq i < j \leq n$. È possibile determinare tutti i sottoalberi ottimali identificando dapprima i più piccoli, e proseguendo la computazione per valori via via crescenti di $j-i$. Per ogni scelta degli indici i e j , esistono $j-i+1$ possibili radici del sottoalbero. Se, procedendo con la computazione lungo l'albero, consideriamo ogni possibile radice di sottoalbero, possiamo constatare che i sottoalberi ottimali di destra e di sinistra sono già stati identificati.

In questo modo è possibile giungere al miglior albero di ricerca in sole n^3 operazioni. In effetti è stato possibile migliorare il metodo, tanto da ridurre il numero di operazioni richieste a n^2 . Nel caso delle 31 parole più frequenti ciò significa che l'albero di ricerca binario ottimale può essere identificato dopo soli 961 passi, invece di 14 quadrilioni.

Si potrebbe osservare che nei precedenti paragrafi si sono analizzati svariati algoritmi, il cui solo scopo è quello di determinare l'albero binario di ricerca più efficiente. In altri termini, il risultato fornito da questi algoritmi è esso stesso

un algoritmo, destinato alla soluzione di un problema diverso! L'esempio aiuta a comprendere perché la scienza dei calcolatori si sia sviluppata tanto rapidamente come una disciplina autonoma. Studiando l'uso corretto dei calcolatori, ci si imbatte in problemi di per sé interessanti, e molti di questi richiedono tecniche e concetti interdependenti.

Può essere stimolante e istruttivo considerare il peggior albero di ricerca binaria possibile per le 31 parole inglesi più comuni, tanto per vedere a quali rovinosi risultati sia possibile giungere con l'algoritmo C. Possiamo definire un siffatto albero «pessimale» in n^2 operazioni circa, come nel caso ottimale. Per le frequenze indicate dalla Gaines, l'albero di ricerca binario «pessimale», costringe, in media, a operare 19 158 confronti per ricerca. Tanto per fare un paragone, il peggior ordinamento delle chiavi per una ricerca sequenziale porta l'algoritmo A a operare in media ben 22 907 confronti per ricerca. Quindi, il peggior caso dell'algoritmo C non sarà mai infame quanto quello dell'algoritmo A.

Hashing

Gli algoritmi di ricerca finora considerati erano strettamente correlati alla tecnica di consultazione di un dizionario. Per esaminare una grande massa di parole con un elaboratore esiste tuttavia una ottima alternativa. Questa, chiamata *hashing* (miscuglio), costituisce un approccio completamente diverso al problema, inutilizzabile da un operatore umano, basato com'è sulla capacità delle macchine di fare calcoli ad altissima velocità. L'idea è quella di trattare le lettere delle parole come se fossero numeri ($a=1, b=2, \dots, z=26$), e quindi di rimescolare i numeri, in modo da assegnare a ogni parola un singolo numero. Il numero è l'hash-indirizzo della parola e indica al calcolatore dove cercare la parola nella tabella.

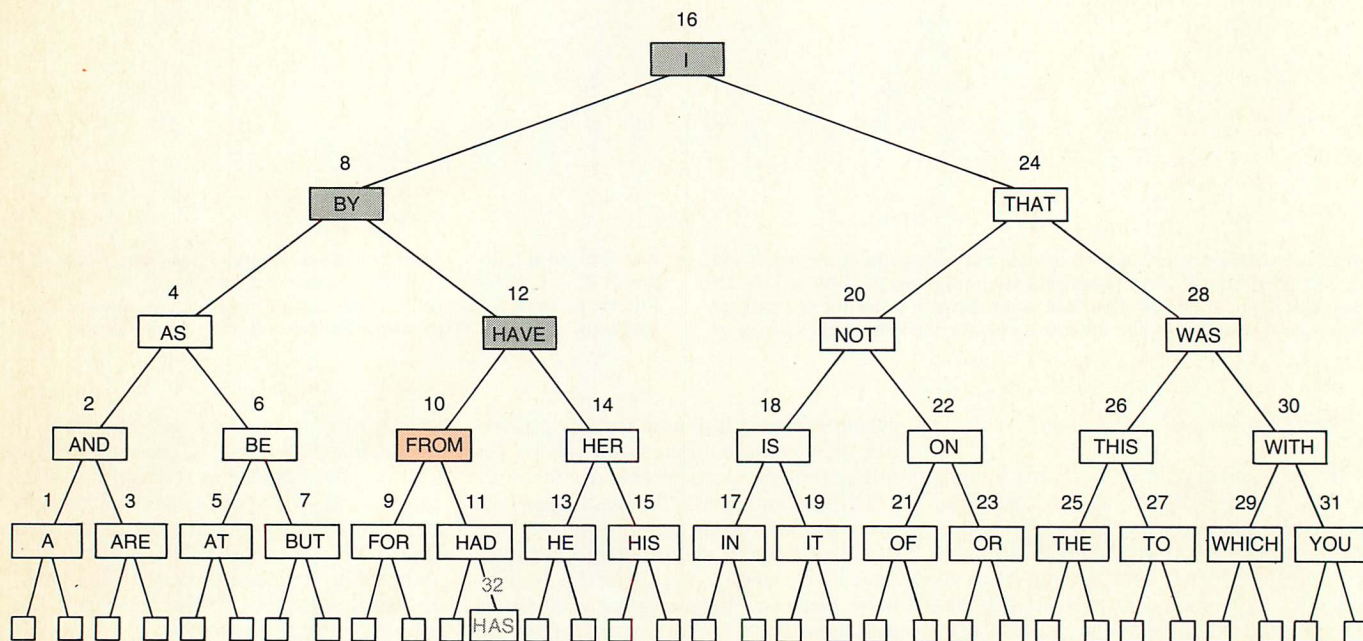
Nel caso delle 31 parole inglesi più comuni potremmo convertire ogni chiave in un numero compreso tra 1 e 32, sommando il valore numerico delle rispettive lettere e prendendo il resto della divisione per 32. Così, l'hash-indirizzo di «the» sarebbe $20 + 8 + 5 - 32 = 1$, quello di «of» risulterebbe essere $15 + 6 = 21$, e così via per il resto della lista. Se si è fortunati,

ogni parola disporrà di un differente hash-indirizzo e le ricerche diventeranno veramente spedite.

In generale, supponiamo di voler memorizzare n chiavi nelle m posizioni di memoria di un elaboratore, dove m è maggiore di n . Poiché $n = 31$, diciamo che m è uguale a 32. Supponiamo inoltre che esista una hash-funzione $h(x)$, che

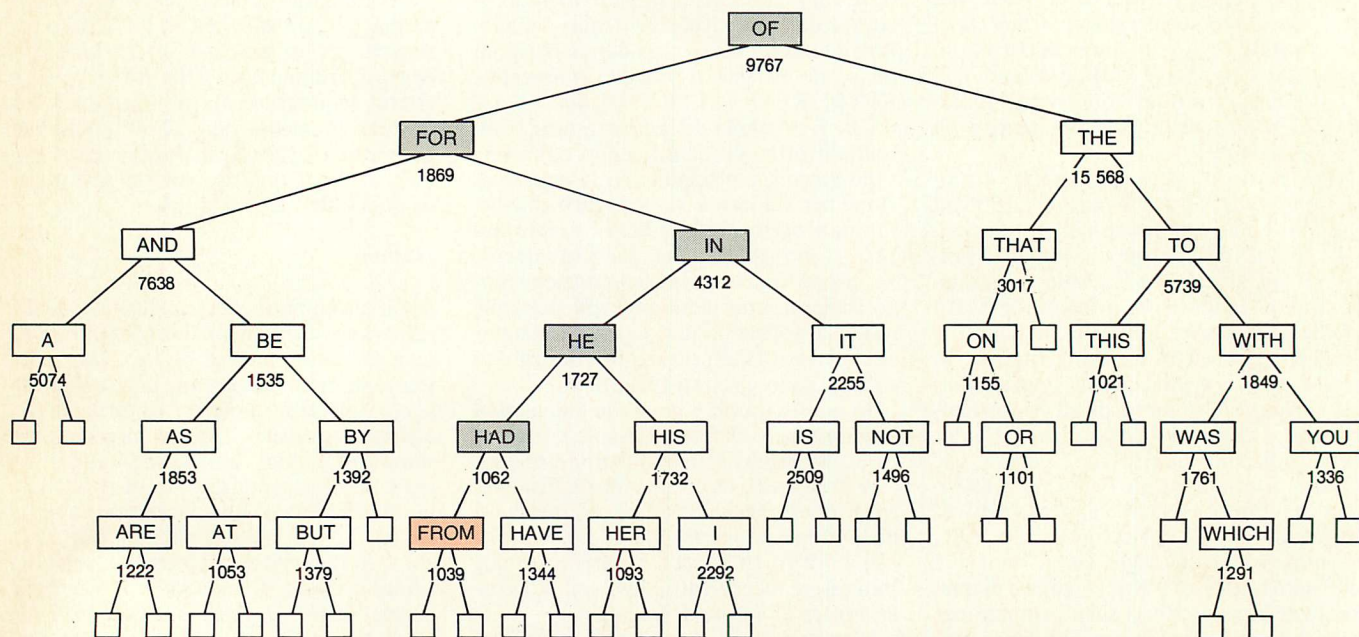
assegna a ogni parola x un numero tra 1 e m . Una hash-funzione è adeguata se è improbabile che $h(x)$ e $h(y)$ siano uguali quando x e y sono parole diverse.

A meno che m non sia molto più grande di n , quasi tutte le hash-funzioni portano, però, a qualche «collisione» tra i valori $h(x)$ e $h(y)$. È estremamente improbabile che n numeri casuali indipen-



Nell'algoritmo B è implicito un albero di ricerca binario; quello qui raffigurato illustra graficamente l'ordine in cui l'algoritmo B esamina una tabella alfabetica composta dalle 31 parole inglesi più frequenti. Partendo dalla radice, o sommità, dell'albero, la parola in ingresso x viene confrontata con il punto centrale della tabella, la parola «I». Se x precede alfabeticamente «I», la ricerca prosegue lungo la ramificazione di sinistra; se x è invece successiva, lungo quella di destra. Così, se x è la parola «from», si constata dapprima che x precede «I»,

quindi che segue «by», precede «have» e infine coincide con «from». Se x non fosse in tabella, la ricerca si fermerebbe a uno dei 32 zeri (quadrati) sotto l'albero. Quando le ramificazioni vengono rappresentate esplicitamente nella memoria del calcolatore, piuttosto che implicitamente, come nell'algoritmo B (il che richiede il calcolo del punto centrale), la ricerca prosegue più spedita e l'inserimento di nuove informazioni viene agevolato: così, la parola «had», (in grigio) può essere inserita, in ordine alfabetico, al posto di uno degli zeri.



L'albero di ricerca binario ottimale, che mostra la migliore disposizione possibile delle 31 parole, si basa sulle frequenze relative (riportate sotto ogni parola). L'albero non è perfettamente bilanciato come nel caso precedente; di conseguenza la ricerca sarà in qualche caso più lunga. Per identificare la parola «from» saranno necessari, con

questo albero, sei passi invece di quattro (sentieri in grigio e in colore). In media però l'albero ottimale consente a un calcolatore migliori prestazioni, poiché le parole più comuni sono esaminate per prime. Si noti che la parola «the», pur essendo la più frequente, non si trova alla radice, perché troppo lontana dal centro dell'alfabeto.

denti tra 1 e m siano tutti differenti. Consideriamo un esempio banale: è noto che quando 23 o più persone siano presenti in una stanza, vi sono buone probabilità che due di queste siano nate nello stesso giorno dell'anno. D'altra parte, in un gruppo di 88 persone è verosimile che ci siano tre individui nati nello stesso giorno dell'anno. Per quanto il fenomeno sembri paradossale a molta gente, il controllo matematico è semplice, e molte altre coincidenze apparentemente impossibili possono essere spiegate allo stesso modo.

Il paradosso del compleanno può essere riformulato dicendo che una hash-funzione con $m=365$ e $n=23$ porterà ad almeno una collisione, più spesso che no. Ogni procedura di ricerca basata su hash-funzioni deve essere dunque in grado di far fronte al problema della collisione.

Supponiamo di voler consultare una tabella per x , ma che l'indirizzo $h(x)$ sia già occupato dalla parola y . Il modo più semplice di trattare la collisione consiste nell'esaminare le posizioni $h(x)$, $h(x)-1$, $h(x)-2$, fino a che non si trovi x o si giunga a una posizione vuota. Se la ricerca giunge al termine della tabella prima di averla consultata tutta, ricomincia dall'altro capo. Questa procedura, detta a scandaglio lineare, può essere espressa sotto forma di algoritmo:

Algoritmo D; hashing con scandaglio lineare.

D1. [Inizio.] Poni $j \leftarrow h(x)$.

D2. [Esito negativo?] Se l'ingresso j della tabella è vuoto, 0 in uscita e l'algoritmo termina.

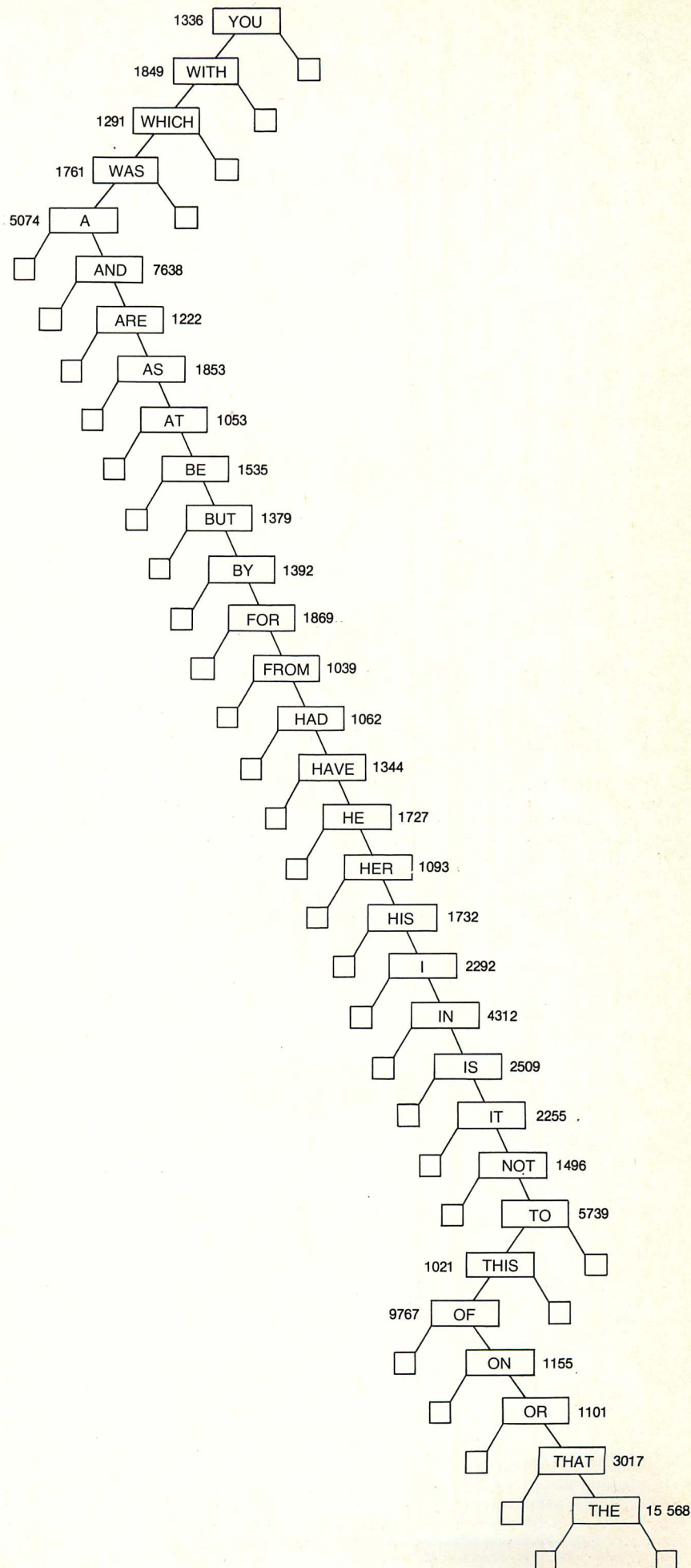
D3. [Esito positivo?] Se $x = \text{KEY}[j]$, j in uscita e l'algoritmo termina.

D4. [Passaggio alla posizione successiva.] Poni $j \leftarrow j-1$; se $j=0$ poni $j \leftarrow m$ (La posizione m è considerata vicina alla posizione 1.) Ritorna a D2.

Se x non è in tabella e l'algoritmo termina senza successo al passo D2 perché l'ingresso j della tabella è vuoto, potremo porre $\text{KEY}[j] \leftarrow x$, utilizzando il valore di j in questione. Una successiva ricerca di x , seguendo esattamente la stessa strada, a partire da $h(x)$ per muoversi poi in $h(x)-1$, ecc, troverà x nella posizione j . In questo modo la ricerca procederà correttamente anche in caso di collisioni.

Ritornando all'esempio delle 31 parole, supponiamo che queste vengano inserite una a una in ordine decrescente di frequenza («the» per prima, «of» per seconda, ecc.) a partire da una tabella vuota. Il risultato è la hash-tabella a sinistra nell'illustrazione della pagina seguente. La maggior parte delle parole appare al proprio hash-indirizzo o nelle immediate vicinanze; fanno eccezione

L'albero «pessimale» è il peggiore tra gli alberi di ricerca binari per le 31 parole più frequenti della lingua inglese. In questo caso vanno persi i vantaggi della struttura ad albero, poiché per ogni confronto abbiamo un ramo «morto».



quelle inserite per ultime. La parola meno frequente, «this», si trova nella posizione 8, nonostante che il suo indirizzo sia 24, perché le posizioni dal 9 al 24 erano già occupate quando è stata inserita. Malgrado queste anomalie il numero medio di unità di tempo necessarie all'al-

goritmo D per trovare una parola si aggira intorno a 1,666: meno della metà rispetto ai confronti richiesti per individuare la parola utilizzando l'albero di ricerca binario ottimale. Al tempo di consultazione bisogna aggiungere quello necessario a calcolare $h(x)$ al passo D1;

ma anche così, per una grande quantità di dati, il metodo dello *hashing* supera decisamente qualsiasi algoritmo di confronto binario.

In pratica è preferibile evitare di riempire la tabella come nell'esempio. Il numero m delle posizioni tabellari è solitamente scelto sufficientemente alto, in modo da saturare la tabella all'80-90 per cento. Si può dimostrare che il numero medio di prove necessarie a trovare una parola tra n ugualmente probabili, che siano state inserite in modo casuale in una tabella di lunghezza m , è $1 + [(n-1)/m + (n-1)(n-2)/m^2 + (n-1)(n-2)(n-3)/m^3 + \dots]/2$.

Sia α la ragione di riempimento o «fattore di carico» n/m della tabella. Per n tendente all'infinito, il numero medio di prove richieste per trovare una qualsiasi parola in un elenco approssima il valore $1 + (\alpha + \alpha^2 + \alpha^3 + \dots)/2$, che equivale a $[1 + 1/(1-\alpha)]/2$. Peraltro, l'effettivo numero medio di prove sarà sempre minore di questo valore limite. Se poi la tabella in esame è carica all'80 per cento, l'algoritmo D eseguirà in media 3 consultazioni per ricerca riuscita.

È da notare che il limite superiore stabilito vale per tutte le tabelle a parità di carico, indipendentemente dalla loro lunghezza. Lo stesso non può dirsi nel caso dell'algoritmo a confronto binario, in cui il tempo di ricerca medio aumenta indefinitamente al crescere del numero n delle parole da esaminare.

Indagini senza esito

Le osservazioni del paragrafo precedente circa l'esiguo numero di prove richieste dall'algoritmo D, valgono solo quando x si trovi effettivamente nella tabella. Se x non c'è, il numero medio di confronti necessari ad accertare il fatto è maggiore, e precisamente $1 + [2n/m + 3n(n-1)/m^2 + 4n(n-1)(n-2)/m^3 + \dots]/2$, che al crescere di n si approssima a $[1 + 1/(1-\alpha)^2]/2$. In altre parole, una ricerca a vuoto in un hash-tabella riempita all'80 per cento, richiede in media 13 prove. D'altronde nell'esempio delle 31 parole immagazzinate in 32 posizioni tabellari, si può notare come tutte le ricerche a vuoto terminino all'unica posizione vuota, 5, indipendentemente da quella dell'indirizzo di partenza $h(x)$. Una situazione analoga si presenta con l'algoritmo a ricerca sequenziale A, le cui ricerche a vuoto terminano tutte alla posizione 0.

Nel 1973 O. Amble dell'Università di Oslo ha osservato che il problema della ricerca a vuoto può essere sdrammatizzato, combinando i due concetti di hashing e di ordinamento alfabetico. Inseriamo le 31 parole inglesi più comuni in una tabella in ordine alfabetico inverso, invece che in ordine di frequenza decrescente. Poiché la tabella viene consultata a partire da $h(x)$, passando a $h(x)-1$, ecc, tutte le parole collocate tra l'indirizzo $h(x)$ e l'effettiva posizione di x devono precedere alfabeticamente x , nell'evenienza di collisioni. La ricerca di x può

1	THE	(1)
2	HAVE	(4)
3	TO	(3)
4	HIS	(4)
5		
6	BE	(7)
7	FOR	(7)
8	THIS	(24)
9	I	(9)
10	BUT	(11)
11	WAS	(11)
12	HAD	(13)
13	HE	(13)
14	FROM	(20)
15	AT	(21)
16	NOT	(17)
17	THAT	(17)
18	WHICH	(19)
19	AND	(19)
20	AS	(20)
21	OF	(21)
22	ON	(29)
23	IN	(23)
24	ARE	(24)
25	YOU	(29)
26	BY	(27)
27	WITH	(28)
28	IS	(28)
29	IT	(29)
30	HER	(31)
31	OR	(1)
32	A	(1)

1	THE	(1)
2	HAVE	(4)
3	TO	(3)
4	HIS	(4)
5		
6	BE	(7)
7	FOR	(7)
8	AND	(19)
9	I	(9)
10	BUT	(11)
11	WAS	(11)
12	HAD	(13)
13	HE	(13)
14	ARE	(24)
15	AS	(20)
16	NOT	(17)
17	THAT	(17)
18	AT	(21)
19	WHICH	(19)
20	FROM	(20)
21	OF	(21)
22	BY	(27)
23	IN	(23)
24	THIS	(24)
25	IS	(28)
26	IT	(29)
27	ON	(29)
28	WITH	(28)
29	YOU	(29)
30	A	(1)
31	HER	(31)
32	OR	(1)

La hash-tabella costituisce una buona soluzione per esaminare grandi quantità di dati con un calcolatore; sfruttandone infatti la capacità di operare calcoli ad alta velocità, possiamo ottenere, per ogni parola x , il rispettivo indirizzo $h(x)$ già all'inizio della ricerca (il numero in parentesi accanto alla parola). Nell'esempio l'indirizzo è ottenuto sommando i valori numerici di ogni lettera ($A=1, b=2...Z=26$), dividendo per 32 il risultato e considerando il resto. Talvolta a due parole x e y corrisponde uno stesso indirizzo $h(x)$, avviene così una «collisione». Se x non si trova in $h(x)$, la ricerca prosegue esaminando $h(x)-1$, $h(x)-2$, ecc. Così, l'indirizzo di «his» è $h+i+s$, cioè $8+9+19=32=4$. L'indirizzo di «have» è anch'esso 4; per trovarne la posizione effettiva l'algoritmo esamina quindi la posizione 4 (grigio chiaro), quindi la 3 (grigio scuro) e infine la 2 (colore), dove appunto si trova «have». Se x non fosse in tabella, la ricerca terminerebbe alla posizione vuota 5. La hash-tabella ordinata (a destra) combina il concetto di hash con i vantaggi dell'ordinamento alfabetico e si rivela particolarmente efficiente quando la parola x non figura nella tabella. Tutte le parole tra $h(x)$ e la posizione effettiva di x sono alfabeticamente successive a x . Una ricerca senza esito non si fermerà quindi necessariamente alla posizione vuota 5, ma non appena incontri una parola alfabeticamente precedente a x . Per x uguale a «has», il cui indirizzo è 28 (grigio chiaro), la ricerca si ferma alla posizione 22 corrispondente a «by» (grigio scuro).

quindi terminare (a vuoto), quando si incontra una parola alfabeticamente seguente a x . Ciò equivale al seguente:

Algoritmo E; scandaglio lineare in una hash-tabella ordinata. Questo algoritmo assume che $KEY[j]$ sia 0, quando l'ingresso j è vuoto, e che tutte le parole x assumano un valore numerico maggiore di 0.

E1. [Inizio.] Poni $j \leftarrow h(x)$.

E2. [Esito negativo?] Se $KEY[j] < x$, 0 in uscita e l'algoritmo termina.

E3. [Esito positivo?] Se $KEY[j] = x$, j in uscita e l'algoritmo termina.

E4. [Passaggio al successivo.] Poni $j \leftarrow j+1$; se $j=0$, poni $j \leftarrow m$. Torna al passo E2.

L'utilità dell'algoritmo E è rilevabile dalla hash-tabella ordinata mostrata nella parte destra dell'illustrazione della pagina a fronte. Supponiamo di voler determinare se «has» sia una delle 31 parole più usate in inglese; il suo hash-indirizzo è $8+1+19=28$. L'algoritmo E termina la ricerca in sei passi, quando arriva a $j=22$ («by»), senza percorrere la tabella fino all'ingresso vuoto a $j=5$.

In una hash-tabella ordinata il numero medio di sondaggi per ricerca a vuoto si riduce a $1 + [n/m + n(n-1)/m^2 + n(n-1)(n-2)/m^3 + \dots]/2$, un numero quasi sempre minore di $[1 + 1/(1-\alpha)]/2$. Quindi, sia per una ricerca con esito positivo sia per una a esito negativo, il limite è lo stesso. Per una tabella carica all'80 per cento l'algoritmo E effettuerà in media tre sondaggi, quale che sia n .

Ciò vale quando l'hash-tabella sia stata costruita inserendo le chiavi in ordine alfabetico decrescente, nel modo descritto. In pratica, tuttavia, non è sempre possibile assumere una tale struttura della tabella. Spesso le tabelle variano dinamicamente con l'uso, venendo introdotte in ordine casuale sempre nuove parole. Mentre la struttura di un albero di ricerca binario (algoritmo C) e di una hash-tabella non ordinata (algoritmo D) possono adattarsi senza difficoltà a tali ac-

crescenti dinamici, lo stesso non può dirsi della struttura di una hash-tabella ordinata. Fortunatamente disponiamo di un algoritmo molto semplice, che consente l'inserimento di nuove parole in una simile tabella.

Algoritmo F; inserimento di una nuova parola in una hash-tabella ordinata. Questo algoritmo introduce la parola x e riordina gli altri ingressi, in modo da preservare la validità della ricerca dell'algoritmo E.

F1. [Inizio.] Poni $j \leftarrow h(x)$.

F2. [Confronto.] Se $KEY[j] < x$, scambia i valori di $KEY[j]$ e x (Vale a dire: poni x al precedente valore di $KEY[j]$ e poni $KEY[j]$ al precedente valore di x).

F3. [Termine?] Se $x=0$ l'algoritmo si arresta.

F4. [Passaggio al successivo.] Poni $j \leftarrow j+1$; se $j=0$ poni nuovamente $j \leftarrow m$. Torna al passo F2.

Se decidessimo di inserire con l'algoritmo F la parola «has» nella hash-tabella delle 31 parole inglesi di uso più comune, la procedura collocherebbe «has» nella posizione 22, grazie al trasferimento di «by» dalla posizione 22 alla 18, di «at» dalla 18 alla 15, di «as» dalla 15 alla 14, di «are» dalla 14 alla 8, e infine di «and» dalla 8 alla 5. Il lavoro può sembrare molto, ma il tempo impiegato è appena superiore a quello necessario a inserire «has» in una hash-tabella non ordinata, utilizzando l'algoritmo D. Generalmente l'inserimento di una parola in una hash-tabella ordinata comporta lo stesso numero di iterazioni dell'inserzione della stessa parola in una hash-tabella non ordinata. Inoltre il numero medio di parole che deve essere spostato al passo F2, per far posto alla nuova parola, è $(n-1)/2m + 2(n-1)(n-2)/3m^2 + 3(n-1)(n-2)/4m^3 + \dots$ approssimativamente uguale a $1/(1-\alpha) + [\log_e(1-\alpha)]/\alpha$, dove e sta per 2,71828. L'impegno richiesto all'algoritmo F per l'inserimento di una parola è quindi ragionevole.

Nel caso specifico non abbiamo potuto inserire effettivamente «has» nella tabella, poiché occorre mantenerla vuota almeno una posizione. Nella nostra tabella completa è presente la prima parola dell'ordinamento alfabetico, «a», quindi lo scandaglio lineare dell'algoritmo F funzionerà in tutti i casi. Se «a» non comparisse nella tabella, sarebbe necessaria una posizione vuota, per evitare una ricerca senza termine qualora la parola in entrata fosse proprio «a».

Una delle più sorprendenti caratteristiche delle hash-tabelle ordinate è la loro unicità. Dato un qualsiasi insieme di parole, se utilizziamo l'algoritmo F per costruirne una, otterremo sempre la medesima tabella quale che sia l'ordine di inserimento delle parole.

Conclusioni

La mia discussione sulle modalità di ricerca delle informazioni immagazzinate nella memoria di un calcolatore era intesa a illustrare alcuni punti di importanza generale concernenti gli algoritmi. Abbiamo visto che un algoritmo deve essere descritto puntualmente, e ciò è più difficile di quanto non sembri. Quando si tenta di risolvere un problema con l'aiuto di un calcolatore, il primo algoritmo che viene in mente può spesso venire perfezionato. La strutturazione dei dati (si vedano gli alberi binari) è un fattore determinante per la costruzione di un algoritmo efficiente. Quando si comincia a studiare la velocità di un algoritmo o si cerca il più adatto a una particolare applicazione, si giunge spesso a risultati interessanti e si può constatare che le risposte ai quesiti che man mano sorgono, sono molto sottili. Talvolta proprio il «migliore» degli algoritmi può essere perfezionato modificando le regole di base. Poiché i calcolatori «pensano» in modo differente dalle persone, i canoni operativi adatti alla mente umana non sono necessariamente i più efficienti quando siano trasferiti sulle macchine.

La teoria delle catastrofi

Le trasformazioni che avvengono in modo improvviso hanno resistito a lungo a un'analisi matematica; oggi è possibile descrivere con metodi topologici tali fenomeni in termini di sette catastrofi elementari

di E.C. Zeeman

I ricercatori si servono spesso della costruzione di modelli matematici per descrivere i fenomeni. Anzi, quando questo tipo di modelli si rivela particolarmente efficace e fruttuoso, si suol dire che esso non solo descrive i fenomeni, ma anche «li spiega»; se il modello può essere ricondotto a una equazione semplice può anche essere considerato come una legge naturale. Il metodo più diffuso nella costruzione di modelli di questo tipo è da 300 anni in qua il calcolo differenziale ideato da Newton e Leibniz. Newton stesso espresse le sue leggi del moto e della gravitazione come equazioni differenziali, e James Clerk Maxwell se ne servì per la sua teoria dell'elettromagnetismo. Anche la teoria generale della relatività di Einstein si può esprimere in una serie di equazioni differenziali. E a questi esempi se ne potrebbero aggiungere molti altri meno noti. Tuttavia, le equazioni differenziali usate come linguaggio descrittivo presentano un limite intrinseco: possono descrivere soltanto quei fenomeni di trasformazione in cui il cambiamento avviene in maniera graduale e continua. In termini matematici le soluzioni di un'equazione differenziale devono essere funzioni differenziabili. Sono relativamente pochi i fenomeni che hanno un andamento regolare e ordinato; al contrario, il mondo è pieno di trasformazioni improvvise e di imprevedibili discontinuità che richiedono l'uso di funzioni che non sono differenziabili.

Soltanto recentemente è stato elaborato un metodo matematico per descrivere i fenomeni discontinui e divergenti. Questo metodo è potenzialmente in grado di descrivere l'evoluzione dei sistemi in ogni manifestazione naturale, e perciò comporta una vasta generalità; può essere applicato con particolare efficacia in quelle situazioni in cui forze e motivazioni che cambiano gradualmente portano a cambiamenti improvvisi del comportamento. Per questo motivo il metodo è stato chiamato «teoria delle catastrofi». Molti fatti nel campo della fisica possono ora essere classificati come esempi di catastrofi. Tuttavia, le applicazioni più rilevanti di questa teoria pos-

sono essere effettuate nel campo della biologia e delle scienze sociali, dove i fenomeni discontinui e divergenti sono i più diffusi, e in cui altre tecniche matematiche si sono dimostrate così largamente insufficienti. La teoria delle catastrofi potrebbe perciò fornire un linguaggio matematico a quelle scienze che sono dette ancor oggi «inesatte».

La teoria della catastrofe è stata elaborata da René Thom dell'Istituto di studi scientifici superiori di Bures-sur-Yvette in Francia. Egli espone la sua teoria in un libro pubblicato nel 1972, *Stabilité Structurelle et Morphogénèse*. La teoria fa largo uso della topologia, il settore della matematica in cui si studiano tra l'altro le proprietà delle superfici a più dimensioni. Il coinvolgimento della topologia è dovuto al fatto che le forze che sono presenti nella natura possono essere descritte mediante superfici di equilibrio uniformi; ed è quando l'equilibrio si spezza che avviene la catastrofe. Il problema dunque per la teoria delle catastrofi è quello di descrivere la forma di tutte le possibili superfici di equilibrio. Thom ha risolto questo problema mediante alcune strutture basilari che egli chiama catastrofi elementari. Per processi determinati da non più di quattro fattori, Thom ha dimostrato che si danno esattamente sette catastrofi elementari. La dimostrazione del teorema di Thom è piuttosto difficile, ma è abbastanza facile comprenderne i risultati. È possibile comprendere le catastrofi elementari stesse e applicarle ai problemi senza bisogno di riferirsi alla dimostrazione.

Un modello dell'aggressione

Per chiarire la natura dei modelli derivati dalla teoria della catastrofe prenderemo in esame una serie di esempi. Comincerò con il considerare un modello di aggressione nel cane. Konrad Z. Lorenz ha notato che il comportamento aggressivo è influenzato da due forze in conflitto tra loro: collera e paura, e ha ipotizzato che nel cane queste componenti possono essere misurate con una certa attendibilità. La collera di un cane è

correlata al grado di apertura della bocca e di esposizione della dentatura; la paura si manifesta da quanto appiattisce le orecchie verso la nuca. Prendendo l'espressione del muso come indicatore dello stato emotivo del cane, possiamo cercare di stabilire in che misura il suo comportamento varia in funzione dell'umore.

Se è presente soltanto una delle componenti emozionali in conflitto, la risposta del cane è relativamente facile da prevedere. Se il cane è adirato allora ci si può attendere un'azione di tipo aggressivo, come l'attacco. Quando il cane è spaventato ma non si provoca la sua collera, l'attacco diventa improbabile e il cane molto verosimilmente batterà in ritirata. La previsione è univoca anche quando non sia presente nessuno stimolo; in questo caso è probabile che il cane manifesti un tipo di comportamento neutrale, non connesso né all'aggressione né alla sottomissione.

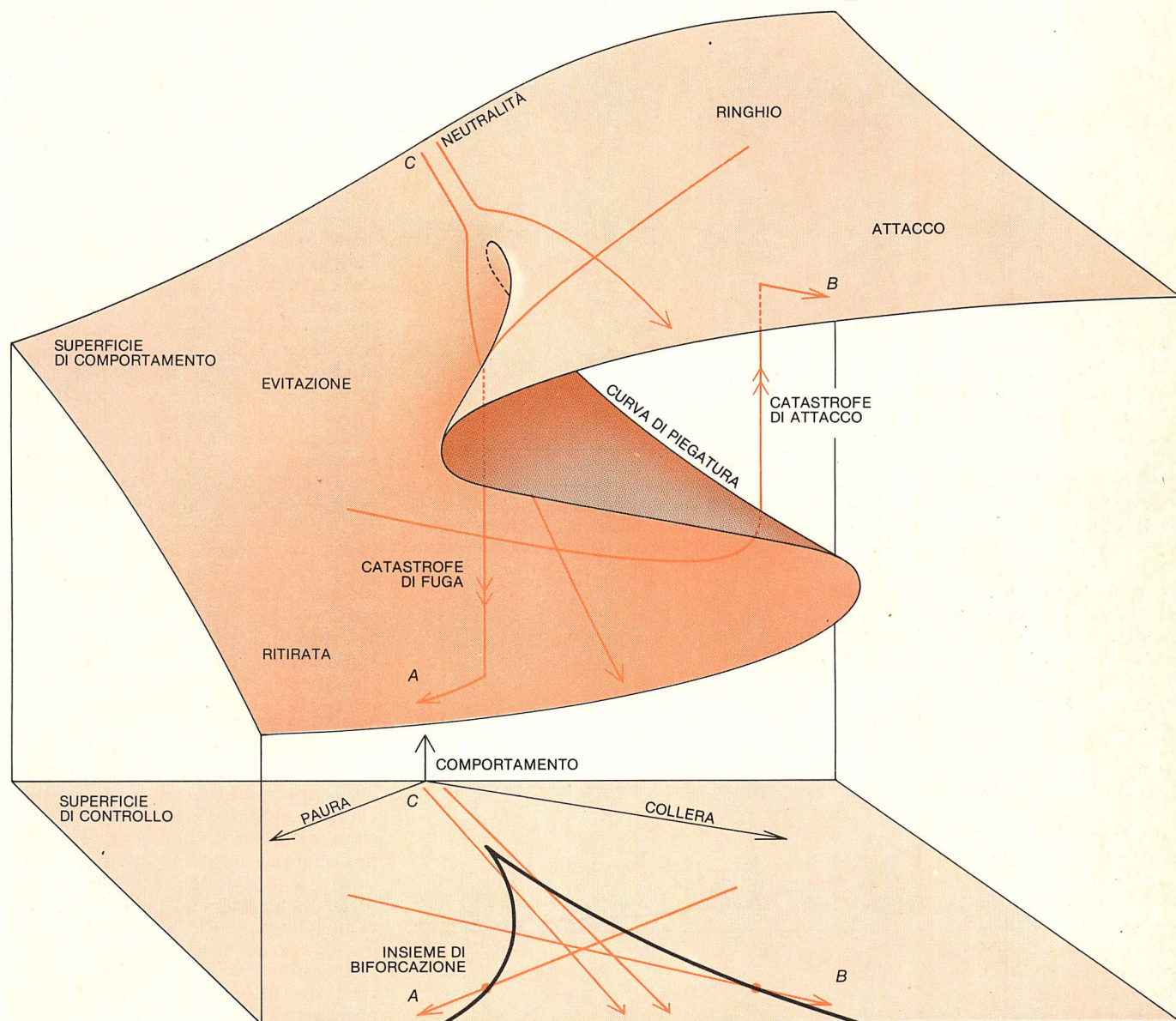
E che cosa succede se il cane si trova a provare nello stesso istante collera e paura? In questo caso le due componenti che lo influenzano sono direttamente in conflitto. I modelli semplici che non sono in grado di esprimere la discontinuità potrebbero in questo caso prevedere che i due stimoli si annullino a vicenda, dando come risultato un comportamento neutrale. Questa previsione dimostra appieno l'insufficienza esplicativa di tali modelli semplicistici, poiché la neutralità è di fatto l'ultimo dei comportamenti probabili. Quando un cane è nello stesso tempo in preda alla collera e alla paura sono alte le probabilità di un comportamento fortemente polarizzato in un senso o nell'altro: il cane può attaccare o abbandonare, ma non rimarrà indifferente. Il punto di forza del modello derivato dalla teoria delle catastrofi è che si può rendere conto di questa distribuzione bimodale delle probabilità. E ancora, il modello fornisce una base per prevedere, in particolari circostanze, quale comportamento sceglierà il cane.

Per costruire il modello, anzitutto rappresentiamo i due parametri di controllo, collera e paura, su due assi in un piano orizzontale, che chiameremo su-

perficie di controllo. Il comportamento del cane viene poi misurato su un terzo asse, l'asse del comportamento, che è perpendicolare ai primi due. Dovremmo assumere che esiste una serie continua di possibili comportamenti: dalla ritirata completa, all'atto di acquattarsi, di schermirsi, alla neutralità, al ringhiare e

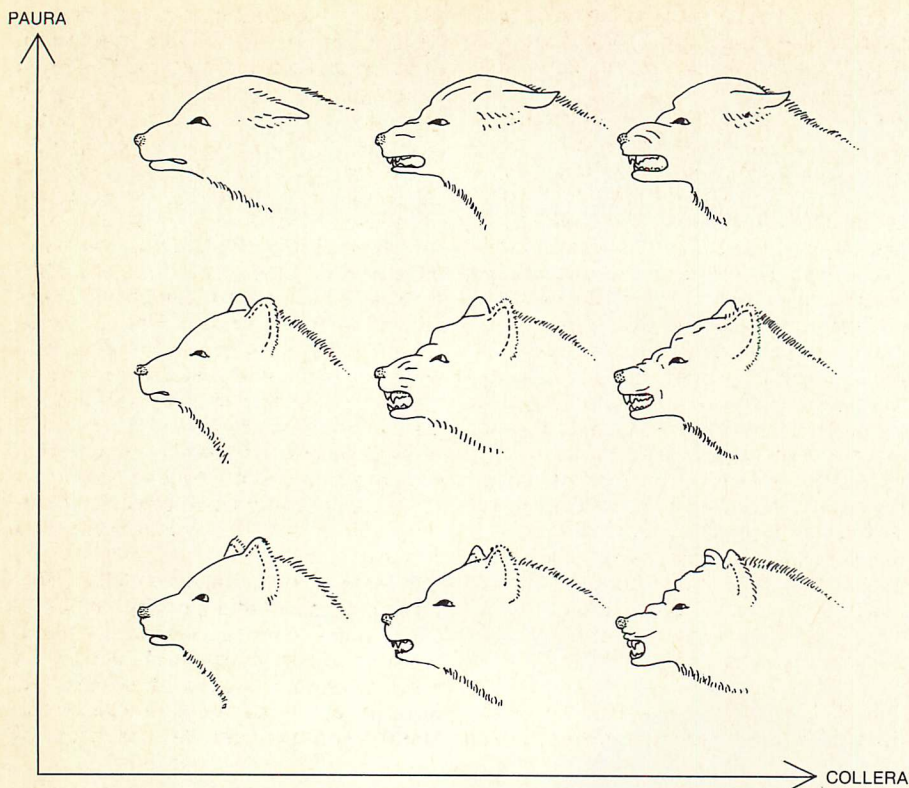
arruffare il pelo fino all'attaccare. I modi di comportamento più aggressivi sono indicati come valori alti sull'asse del comportamento, quelli meno aggressivi con i valori bassi. Per ogni punto sulla superficie di controllo (cioè, per ogni combinazione di ira e paura) esiste almeno un comportamento più probabile, che rap-

presentiamo con un punto direttamente superiore al punto della superficie di controllo collocato a un'altezza corrispondente al comportamento. Ci sarà così un solo punto per il comportamento in corrispondenza di molti punti sulla superficie di controllo, nei quali predomini la collera o la paura. Verso il centro del

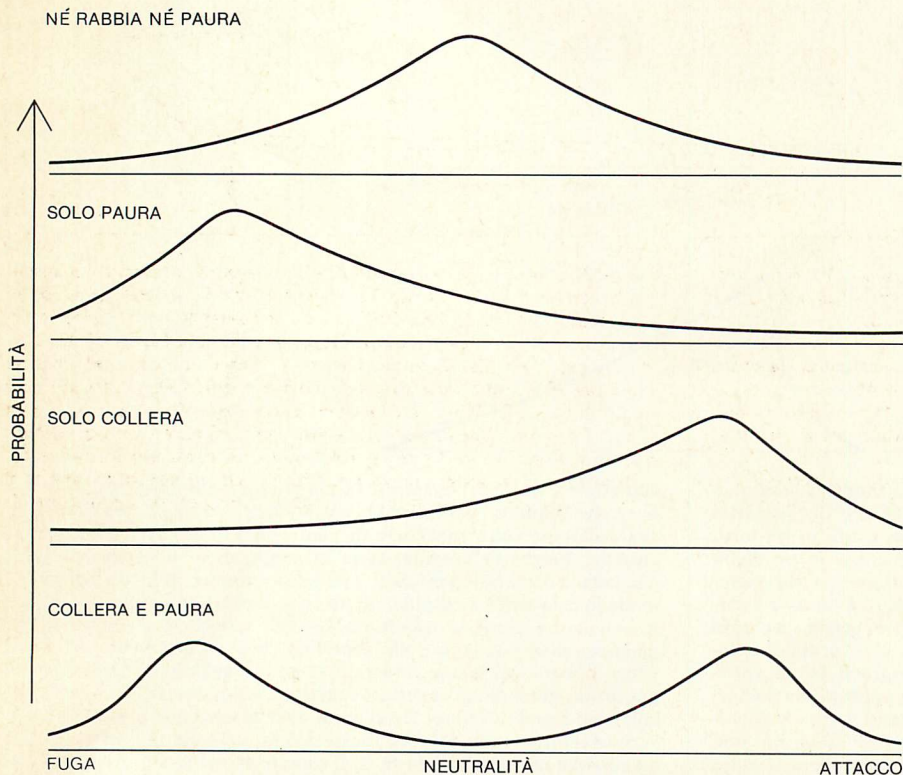


L'aggressività nei cani può essere descritta con un modello basato su una delle catastrofi elementari. Il modello assume che il comportamento aggressivo sia determinato da due fattori in conflitto tra loro, collera e paura, che vengono rappresentate su due assi in un piano orizzontale: la superficie di controllo. Il comportamento del cane, che va dall'attacco alla ritirata, viene rappresentato su un asse verticale. Per ogni combinazione di collera e paura, e perciò per ogni punto della superficie di controllo, esiste almeno una forma di comportamento probabile, indicata con un punto sopra il punto corrispondente sulla superficie di controllo, ad altezza appropriata sull'asse del comportamento. L'insieme di tutti questi punti forma la superficie di comportamento. Nella maggioranza dei casi esiste un solo modo di comportarsi probabile, ma nelle zone in cui collera e paura approssimativamente si equivalgono, ci sono due modi possibili: un cane che sia nello stesso tempo incollerito e spaventato può sia attaccare sia ritirarsi, quindi ci sono al centro del grafico due piani che rappresentano il comportamento possibile, ed essi sono collegati da un terzo piano formando una superficie continua e ripiegata. Il terzo dei piani, al centro, disegnato in grigio, presenta una significativa diversità rispetto agli altri due: rappresenta il comportamento meno probabile, in questo caso la neutralità. Verso l'origine la piega sulla

superficie di comportamento si fa più stretta e alla fine scompare. La linea che delimita i confini della piega è detta curva di piegatura, e la sua posizione sulla superficie di controllo è una curva a forma di cuspid. Poiché la cuspid segna la soglia dove il comportamento comincia a diventare bimodale, è chiamata insieme di biforcazione e il modello catastrofe a cuspid. Se un cane incollerito viene spaventato, il suo umore segue la traiettoria A sulla superficie di controllo. Il percorso corrispondente sulla superficie di comportamento si sposta verso sinistra sul piano superiore fino a raggiungere la curva di piegatura; qui il piano superiore scompare e il percorso va a cadere di colpo sul piano inferiore. E qui che il cane interrompe il suo attacco e si ritira improvvisamente. Allo stesso modo, un cane spaventato che si incollerisce segue la traiettoria B. Il cane resta nel piano inferiore fino a che questo piano scompare, e poi non appena salta al piano superiore, smette di ritirarsi e attacca improvvisamente. Un cane che è nello stesso tempo incollerito e spaventato deve seguire una delle due traiettorie in C. Il fatto che si sposti sul piano superiore e diventi più aggressivo o che si sposti sul piano inferiore e diventi più sottomesso dipende in maniera critica dai valori della collera e della paura. Un lieve cambiamento degli stimoli può essere all'origine di un cambiamento assai rilevante del comportamento: il fenomeno è divergente.



I fattori di controllo nel modello dell'aggressività sono la collera e la paura, che nel cane si rilevano dall'espressione del muso: la collera dal grado di apertura della bocca, e la paura dal grado di appiattimento delle orecchie verso la nuca. Da questi indicatori è possibile rendersi conto dello stato emozionale del cane, e per mezzo del modello prevedere il suo comportamento.



Una funzione di probabilità determina il comportamento del cane sotto l'influenza della collera e della paura in conflitto. Quando non è presente nessuno stimolo, il comportamento più probabile è quello di neutralità; quando è presente solo la collera, è l'aggressione, quando c'è solo la paura, è la resa. Quando il cane è incollerito e spaventato, il grafico della probabilità diventa bimodale: l'attacco e la fuga sono entrambe scelte possibili, mentre la neutralità è la risposta meno probabile. La bimodalità è rispecchiata nella superficie di comportamento della catastrofe a cuspide, che presenta due piani che rappresentano il comportamento più probabile.

grafico tuttavia, laddove la collera e la paura approssimativamente si equivalgono, ogni punto appartenente alla superficie di controllo ha due punti di comportamento, uno con valori alti sull'asse del comportamento, che rappresenta l'azione di aggressione, l'altro con valori bassi che rappresenta l'atto di sottomissione. Possiamo inoltre segnare un terzo punto, che verrà sempre a cadere tra questi due, che rappresenta il comportamento più neutrale meno probabile.

Se individuiamo i punti di comportamento per l'intera superficie di controllo e li uniamo tra loro, essi formeranno una superficie continua: la superficie di comportamento. La superficie presenta nel suo complesso una pendenza, dai valori alti, dove è predominante la collera, fino ai valori bassi nella regione in cui lo stato prevalente è la paura; ma la pendenza non è la caratteristica più tipica della superficie. La teoria delle catastrofi rivela che al centro della superficie deve esserci una piega doppia e continua, che dà luogo a una rientranza senza grinze verso la parte posteriore della superficie che alla fine scompare in un punto in cui i tre piani della piega si fondono (si veda l'illustrazione della pagina precedente). È la piega che conferisce a questo modello le sue caratteristiche più interessanti. Ciascun punto della superficie di comportamento rappresenta il comportamento più probabile del cane, con l'eccezione di quelli che vengono a cadere nel mezzo della piega, che rappresentano il comportamento meno probabile. Con la teoria delle catastrofi possiamo ricavare la forma dell'intera superficie dal fatto che il comportamento è bimodale per alcuni dei punti di controllo.

Per comprendere come questo modello può prevedere il comportamento, dobbiamo considerare la reazione del cane a degli stimoli che cambiano. Supponiamo che la condizione emozionale del cane sia di neutralità e che possa essere rappresentata da un punto all'origine della superficie di controllo. Il comportamento del cane, dato dal punto corrispondente sulla superficie di comportamento, è anch'esso di neutralità. Se in seguito alcuni stimoli fanno crescere l'ira del cane senza influenzare la paura, il comportamento cambia gradualmente, spostandosi sulla superficie di comportamento a quote superiori, verso posizioni di maggiore aggressività; se la collera è aumentata abbastanza, il cane attacca. Se ora la paura del cane comincia ad aumentare mentre la collera rimane a un livello elevato, il punto che rappresenta lo stato emozionale sulla superficie di controllo si sposterà verso il centro del grafico. Il punto che rappresenta il comportamento deve evidentemente seguirlo, ma poiché la pendenza della superficie di comportamento in questa regione non è accentuata, il comportamento cambia soltanto leggermente, e il cane continua a essere aggressivo. Man mano che la paura continua a crescere, però, il punto del comportamento raggiunge infine il bordo della piega. E qui si evidenziano le originali e illuminanti proprietà di que-

sto modello. Al bordo della piega il piano lungo il quale si è spostato il punto del comportamento si piega all'ingiù, e viene perciò ad annullarsi; con un ulteriore aumento della paura, la superficie scompare. La situazione comportamentale deve perciò cadere perpendicolarmente sulla superficie inferiore del grafico, che rappresenta modi di comportarsi assai diversi. Gli stati aggressivi del piano superiore non sono più possibili: si verifica un cambiamento improvviso, anzi catastrofico, un atteggiamento più mite. Il modello dunque prevede che se un cane adirato si impaurisce progressivamente, alla fine interromperà il suo attacco e si ritirerà. Il cambiamento improvviso del comportamento potrebbe essere chiamato una catastrofe di fuga. Il grafico permette inoltre di individuare un modello opposto di comportamento: una catastrofe di attacco. In uno stato iniziale dominato dalla paura il comportamento del cane permane stabilmente sul piano inferiore, ma in seguito a un adeguato aumento di collera passa dalla parte opposta della piega saltando di colpo al piano superiore, corrispondente a uno stato emozionale più aggressivo.

In altre parole, un cane spaventato, messo in una situazione in cui la rabbia aumenta costantemente, può attaccare improvvisamente.

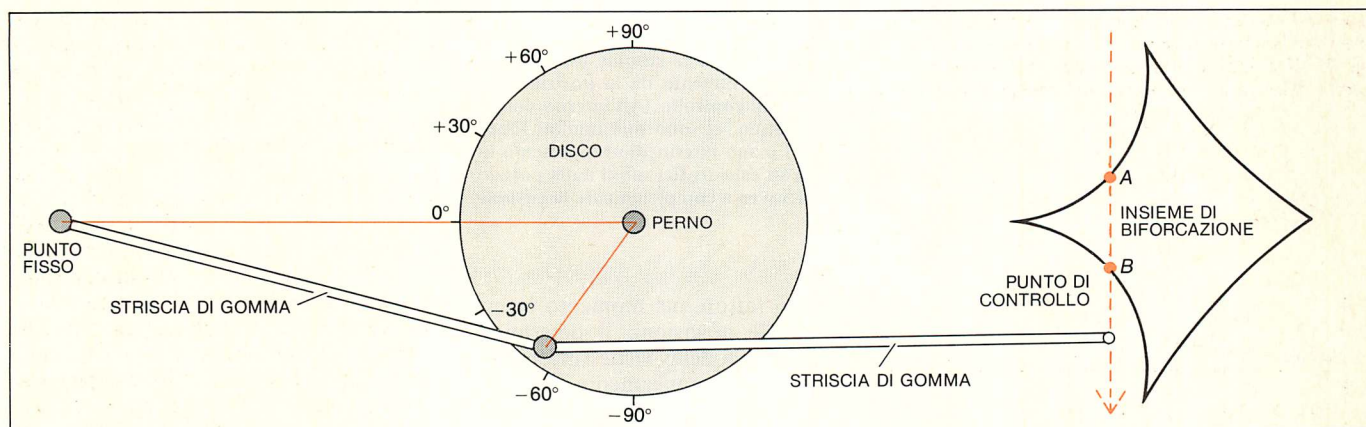
Consideriamo infine il comportamento di un cane, il cui umore iniziale è di neutralità, quando la collera e la paura aumentano contemporaneamente. Il punto di comportamento si trova all'inizio nell'origine, e sotto l'influenza di stimoli in conflitto tra loro si sposta su una linea retta verso la parte anteriore del grafico. Arrivato però alla singolarità che la superficie presenta dove inizia a piegarsi, il punto si deve spostare sul piano superiore quando il cane diventa più aggressivo o sul piano inferiore quando il cane diventa meno aggressivo. Il piano che viene scelto dipende in modo critico dall'umore appena prima che raggiunga la singolarità. Il grafico è detto divergente: un cambiamento lieve nelle condizioni iniziali dà luogo, nella condizione finale, a un cambiamento molto rilevante.

La catastrofe a cuspide

La linea che segna i bordi della piega nella superficie di comportamento, dove

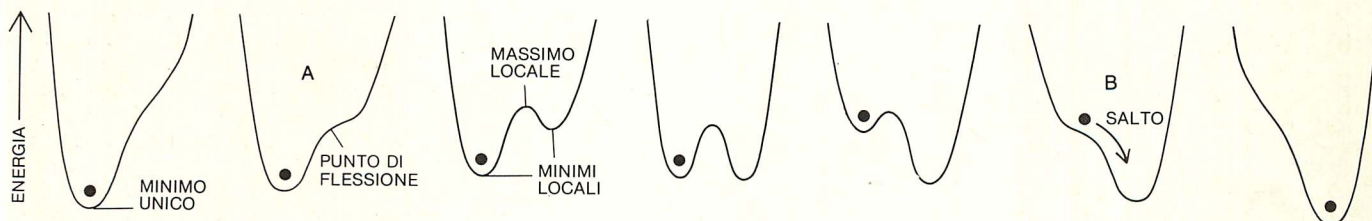
il piano superiore e il piano inferiore si flettono per formare il piano intermedio, è chiamata curva di piegatura, Proiettata sulla superficie di controllo, dà luogo a una curva a forma di cuspide; per questo il modello è chiamato catastrofe a cuspide. È una delle più semplici catastrofi elementari, e si è dimostrata di gran lunga la più utilizzabile. La cuspide sulla superficie di controllo è detta insieme di biforcazione della catastrofe a cuspide e definisce le soglie entro cui si verificano i cambiamenti improvvisi. Fino a che lo stato del sistema si mantiene all'esterno della cuspide, il comportamento varia in maniera graduale e continua in funzione dei parametri di controllo. Anche sulla soglia della cuspide non si osserva nessun cambiamento brusco. Quando il punto di controllo continua il percorso attraversando la cuspide, è inevitabile una catastrofe.

In ogni punto all'interno della biforcazione ci sono due possibili modi di comportarsi; al di fuori di essa, c'è soltanto un modo possibile. E ancora, ci sono all'interno della cuspide solo due modi di comportamento, anche se nel punto corrispondente la superficie di



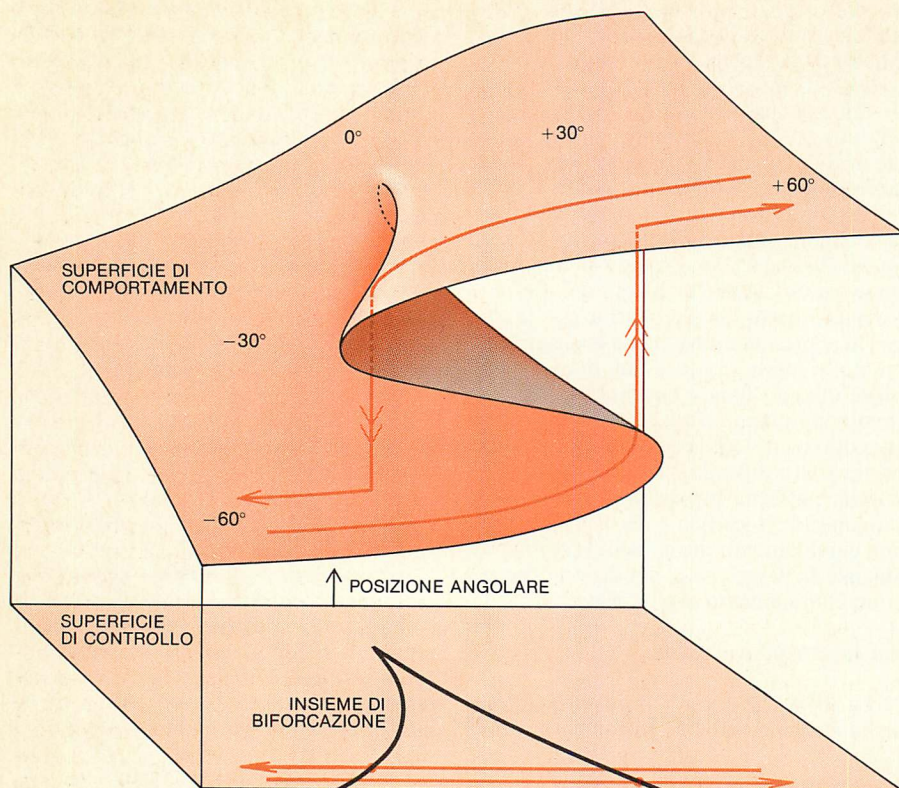
La macchina della catastrofe ideata dall'autore presenta un comportamento discontinuo che può essere descritto per mezzo di una catastrofe a cuspide. La macchina consiste in un disco di cartone fissato con un perno nel centro, con due strisce di gomma fissate in un punto vicino al perimetro. La lunghezza di ogni striscia, non sottoposta a tensione, è approssimativamente uguale al diametro del disco. L'estremità libera di una delle strisce è fissata a lato dell'apparecchiatura, e la macchina si mette in funzione spostando l'altra striscia di gomma, l'estremità libera della quale è detta punto di controllo. Il comportamento misurato è l'angolo formato dal punto fisso, dal perno e dal punto in cui le due strisce sono fissate al disco. Molti spostamenti del

punto di controllo provocano solo delle leggere rotazioni del disco, ma in certi casi il disco scatta improvvisamente da un lato all'altro. Se si segnano tutte le posizioni del punto di controllo in cui avvengono questi movimenti improvvisi, ne risulta una curva a forma di diamante. La curva consta di quattro cuspidi, ciascuna delle quali forma l'insieme di biforcazione di una catastrofe a cuspide. Spostando il punto di controllo lungo la traiettoria in colore, nella posizione A non si verifica nessun movimento, e il disco gira di colpo fino a portarsi nella posizione B. Se il percorso del punto di controllo è rovesciato, esso attraversa B senza che si verifichi nessun movimento, mentre il disco si sposta quando il punto di controllo raggiunge la posizione A.



Una funzione dell'energia regola il comportamento della macchina della catastrofe. La macchina tende sempre ad assumere una posizione di energia minima, vale a dire, il disco ruota fino a rendere minima la tensione delle strisce di gomma. Quando il punto di controllo si trova al di fuori dell'insieme di biforcazione, esiste una sola posizione di energia minima, che corrisponde all'unica posizione stabile del disco. Quando il punto di controllo viene spostato attraverso l'insieme

di biforcazione, in A si sviluppa un secondo minimo locale di energia che alla fine diventa più profondo del primo. La macchina tuttavia non può raggiungere il nuovo minimo locale, perché ne è separata da un massimo locale. Soltanto quando il punto di controllo attraversa la seconda linea della cuspide in B il massimo locale viene eliminato; a questo punto l'equilibrio si spezza e la macchina si sposta in modo improvviso nella nuova posizione caratterizzata da energia minima.



Il modello a cuspide della macchina della catastrofe dà luogo a una superficie di comportamento piegata, su una parte dell'insieme di biforcazione, esattamente come la cuspide più vicina al disco. Ogni punto del piano superiore della superficie di comportamento dà la posizione del disco che ha il minimo di energia per quella posizione del punto di controllo. All'interno dell'insieme di biforcazione, dove sono due le posizioni stabili del disco, ci sono due minimi locali, uno collocato sul piano superiore e uno sul piano inferiore. Il piano intermedio rappresenta invece il massimo locale nella funzione dell'energia. I cambiamenti catastrofici subiti dalla posizione angolare si verificano ogni volta che il punto di controllo attraversa completamente la cuspide.

comportamento presenta tre piani. Ciò in quanto il piano che si trova nella regione centrale della piega è formato da punti che rappresentano il comportamento meno probabile. Il piano intermedio è inserito nel grafico soprattutto affinché la superficie di comportamento risulti omogenea e continua; il punto di comportamento non occupa mai il piano intermedio. Anzi, non esiste alcun percorso sulla superficie di controllo, che potrebbe condurre il punto di comportamento sul piano intermedio. Ogniqualvolta la curva di piegatura viene attraversata, il punto «salta» dal piano superiore al piano inferiore, e viceversa.

La costruzione di questo modello si basava in un primo tempo su un'ipotesi sostanzialmente deterministica: che il comportamento del cane avrebbe potuto essere previsto dal suo umore così come veniva esternato dall'atteggiamento del muso. La bimodalità del grafico che ne risulta potrebbe sembrare in un primo tempo un elemento in grado di minare questa ipotesi, poiché l'esistenza di due possibili maniere di comportarsi per uno stato emozionale dato rende impossibile una previsione univoca. Effettivamente se conosciamo soltanto lo stato emozionale attuale (e se questo stato viene a cadere all'interno della regione bimodale del grafico) non possiamo prevedere che cosa farà il cane. Tuttavia se considera-

mo un altro fattore nel momento in cui facciamo delle previsioni, il determinismo di questo modello viene rinforzato, rendendo contemporaneamente il modello più complesso e sofisticato. Il comportamento del cane può essere previsto se noi conosciamo, oltre al suo stato emozionale attuale, anche la storia recente delle sue emozioni. Non sarebbe sorprendente che gli effetti ottenuti spaventando un cane adirato siano diversi da quelli ottenuti facendo adirare un cane spaventato.

Modelli di comportamento umano

La catastrofe a cuspide fornisce anche un'interpretazione di certi tipi di comportamento umano. Per esempio, una discussione spesso comporta la manifestazione di aggressività, e il suo sviluppo è fortemente determinato dalla collera e dalla paura. Si può costruire una catastrofe a cuspide prendendo questi stati emozionali come parametri di controllo, considerando l'intensità del conflitto come asse del comportamento.

All'inizio, sulla superficie di comportamento si trova il comportamento più neutrale: la discussione razionale. Con l'aumento della collera e della paura, il punto di comportamento si sposta in avanti sul grafico, verso quella parte di superficie che è piegata, dove il compor-

tamento è bimodale. I contendenti vedono chiudersi la possibilità di condurre un discorso equilibrato, e sono costretti o a fare asserzioni sempre più recise, o a fare concessioni. Con un aumento ulteriore di emotività, le scelte di comportamento possibili divergono ulteriormente, e le alternative sono l'invettiva o le scuse; alla fine gli avversari devono scegliere tra la sfuriata e le lacrime. Una volta che la discussione è divenuta calda, un leggero aumento sia di collera che di paura può causare un cambiamento subitaneo nel comportamento. Chi sostiene con aggressività una tesi, quando comincia a oscillare nella sua sicurezza, può abbandonare la posizione e cominciare a scusarsi; un oppositore timido, costretto a fare ripetute concessioni può improvvisamente perdere la calma e diventare furibondo. Il modello suggerisce anche una strategia per una persuasione efficace. Se è probabile che una discussione sollevi sia collera sia paura, allora è molto meglio esporre i fatti senza cercare di imporre subito il proprio punto di vista, lasciando sbollire le emozioni e mettendo in grado l'avversario di riguadagnare la via del pensiero razionale.

Un altro schema di comportamento umano che può essere descritto dal modello a cuspide è l'autocommiserazione e la catarsi che a volte interviene ad alleviarla. In questo caso i parametri di controllo, analoghi alla collera e alla paura, sono le emozioni meno violente della frustrazione e dell'ansia. E ancora, l'asse del comportamento non misura un comportamento manifesto, che è negli animali l'unico che può essere osservato, ma gli umori sottostanti che nell'uomo possono essere identificati direttamente. Una tipica serie di stati d'animo potrebbe andare dalla collera alla noia, attraverso stati d'animo di neutralità, arrivando fino all'abbattimento e all'autocommiserazione.

Nel modello, un aumento notevole dell'ansia provoca uno stato d'animo persistente di autocommiserazione; il punto che rappresenta lo stato d'animo finisce nel piano inferiore della superficie di comportamento piegata (si veda l'illustrazione in alto a pagina 164). L'autocommiserazione è un atteggiamento difensivo adottato assai frequentemente dai bambini e sembra spesso che un atteggiamento di solidarietà e di comprensione risulti inefficace ad attenuarla. Eppure una battuta spiritosa può provocare una subitanea diminuzione della collera e liberando la tensione può aprire la possibilità di un ritorno a uno stato meno dominato dall'emozione. È spiacevole che sia l'ironia ad aver successo e la comprensione a fallire, ma la causa di ciò appare chiara nel modello. Il sarcasmo porta a un aumento della frustrazione, e il risultato è che il punto che rappresenta lo stato d'animo attraversa la superficie di comportamento fino alla curva di piegatura; dopo aver raggiunto l'estremità del piano inferiore, non può che compiere un salto catastrofico fino al piano superiore, e l'autocommiserazione si trasforma in collera.

Questi esempi della catastrofe a cuspidi offrono un modello interessante e apparentemente efficace di certe modalità di comportamento animale e umano, ma si tratta solo di un modello fenomenologico; non si può ancora dire che sia in grado di spiegare il comportamento. Il problema della motivazione per cui un cane fa quello che fa non è stato risolto ma semplicemente spostato a un livello di maggiore astrazione. Dobbiamo ora chiederci perché il modello funziona. In particolare, perché il punto del comportamento si muove sulla superficie fino al bordo di una piega, e poi compie un balzo catastrofico in un altro piano? Perché non passa gradualmente da una superficie all'altra, in modo che si verifichi una transizione continua? Qual è il meccanismo che regola lo stato del sistema sulla superficie di comportamento? La risposta a queste domande può essere ottenuta per mezzo di un altro esempio di catastrofe a cuspidi, esempio che riguarda il comportamento di un sistema assai più semplice di quello di un cane o di un uomo.

Una macchina della catastrofe

È possibile generare artificialmente catastrofi elementari per mezzo di un semplice dispositivo composto da un cartoncino rigido, strisce di gomma e pochi altri oggetti. La parte principale della macchina consiste in un disco di cartoncino con un perno al centro, e di due strisce di gomma attaccate in un punto vicino al perimetro. Il capo libero di una delle due strisce è fissato a un punto che si trova al di fuori del disco; l'altra striscia di gomma serve a controllare i movimenti del disco, e si indica la posizione del suo capo libero come il punto di controllo (*si veda l'illustrazione in alto a pagina 161*). La macchina della catastrofe si mette in moto spostando il punto di controllo nel piano del disco. In molti casi il risultato è una leggera rotazione del disco; tuttavia, in una regione, prossima al punto diametralmente opposto al punto in cui è agganciata la striscia fissa, un leggero movimento del punto di controllo può provocare un movimento improvviso del disco. Segnando la posizione del punto di controllo ogni volta che il disco scatta, ne risulta una curva concava, a forma di diamante. Questa curva è formata da quattro cuspidi connesse, gli insiemi di biforcazione di quattro catastrofi a cuspidi.

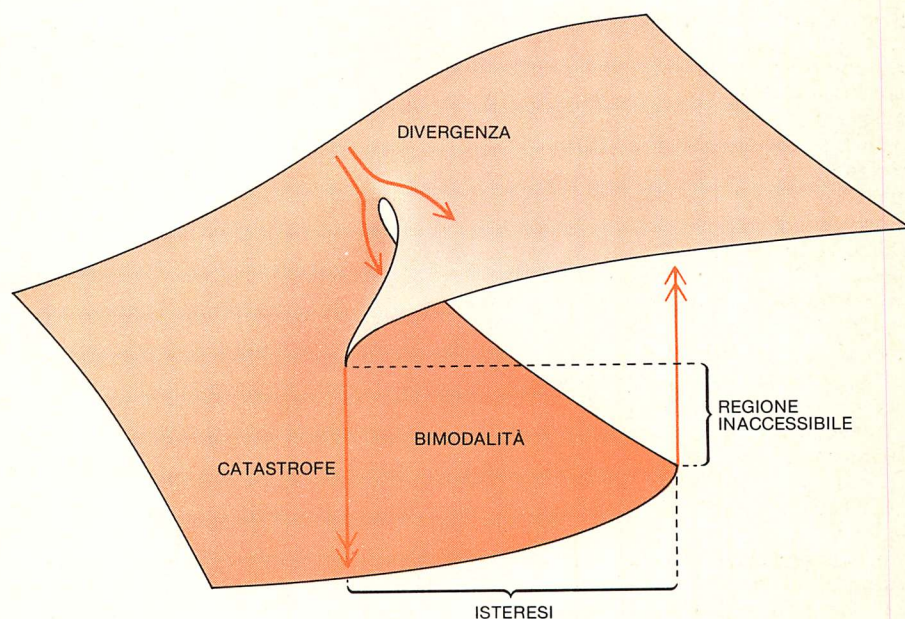
Se consideriamo solo una di queste quattro cuspidi, quella più vicina al disco, possiamo costruire la corrispondente superficie di comportamento sistemando la curva di piegatura in modo che venga a trovarsi proprio sopra la cuspidi. Per ogni posizione del punto di controllo al di fuori della cuspidi, la superficie di comportamento ha un solo piano e il disco ha una sola posizione stabile. Se il punto di controllo si trova all'interno della cuspidi, la superficie del comportamento ha tre piani, ma il piano di mezzo deve essere escluso perché corrisponde a un equilibrio instabile. Come

risultato, dunque, ci sono due posizioni stabili del disco. Il fatto che il comportamento della macchina risponda a questo modello può essere verificato spostando il punto di controllo attraverso il grafico da destra a sinistra. Il disco si muove lievemente e solo in misura trascurabile fino a quando il punto di controllo tocca il bordo di destra della cuspidi; il disco gira poi di colpo non appena il punto di comportamento cade al di fuori dell'estremità del piano inferiore e balza su quello superiore. Quando il percorso del punto di controllo va in senso contrario, il punto attraversa il bordo destro della cuspidi senza che si verifichi l'evento, e il disco continua a muoversi lievemente fino a che il punto raggiunge il bordo sinistro della cuspidi; qui, il punto di comportamento salta dal piano superiore a quello inferiore.

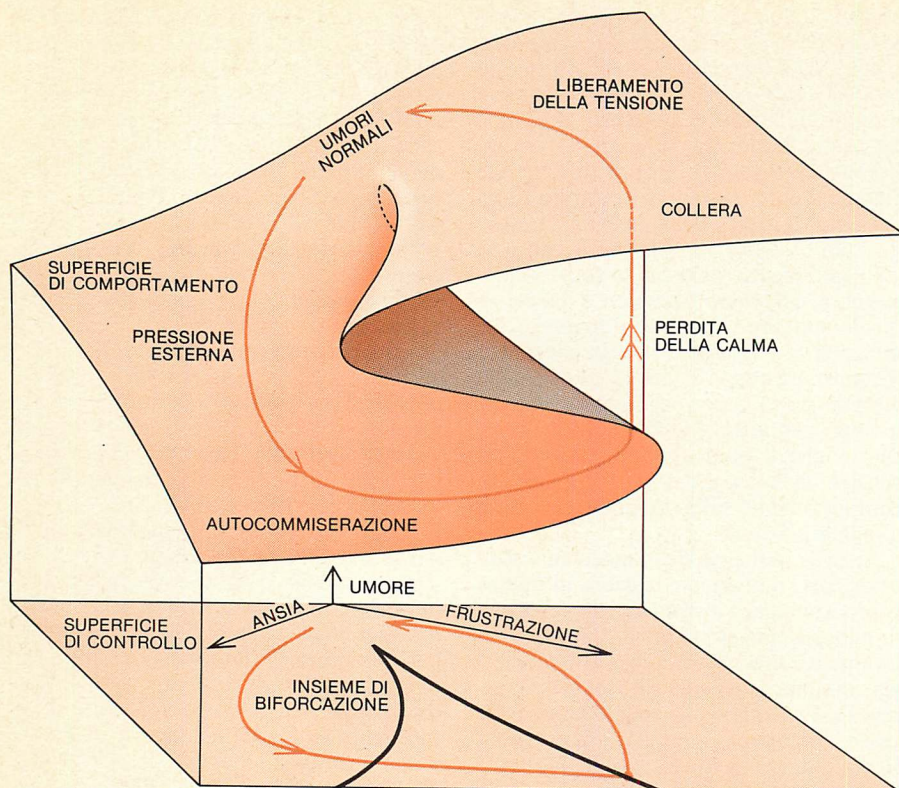
È facile, nel caso della macchina della catastrofe, individuare la causa di questo comportamento: si tratta della tendenza di ogni sistema fisico in cui l'attrito è rilevante a raggiungere uno stato di energia minima. L'energia da ridurre al minimo è l'energia potenziale situata nelle strisce di gomma, e infatti il disco ruota fino a quando la tensione sulle due fasce non è al minimo. In questa posizione, la macchina è in equilibrio stabile. A meno che si aggiunga energia al sistema, la macchina rimane in uno stato di equilibrio. Il processo che la conduce a quel punto è detto «dinamica».

Il funzionamento della dinamica può essere dimostrato con una serie di grafici, ciascuno dei quali mostra per ogni singola posizione del punto di controllo l'energia della macchina per tutte le possibili rotazioni del disco (*si veda l'illustra-*

zione in basso a pagina 161). Fino a che il punto di controllo resta al di fuori della cuspidi, il grafico consiste in una curva continua con una sola depressione, o minimo, e lo stato della macchina si sposterà rapidamente alla condizione di energia minima nella parte più bassa della depressione. Non appena il punto di controllo cade all'interno della cuspidi, si sviluppa, vicino alla prima, una seconda depressione, che rappresenta l'energia minima locale. Questa seconda depressione diviene sempre più profonda, ma la macchina non può entrare nello stato corrispondente a essa poiché le due depressioni sono separate da un picco, che rappresenta un massimo locale dell'energia. Lo stato della macchina non cambia fino a che la seconda depressione non si unisce al valore massimo locale. Ciò avviene quando il punto di controllo oltrepassa la seconda linea della cuspidi, e lo stato della macchina di conseguenza viene portato di colpo dalla dinamica in una nuova e unica posizione di energia minima. Si chiarisce allora il significato della dinamica nel momento in cui si scopre che la superficie di comportamento della catastrofe a cuspidi è il grafico di tutti i minimi e di tutti i massimi della funzione dell'energia. Al di fuori della cuspidi esiste un solo minimo di energia, e non esistono massimi; la superficie di comportamento perciò presenta un solo piano. In corrispondenza della cuspidi si formano un nuovo minimo e un nuovo massimo locale e di conseguenza si formano due nuovi piani sulla superficie di comportamento. Lo stato della macchina non può mai rimanere stabile nel piano intermedio poiché tale posizione corrisponde al massimo di energia.



Cinque proprietà caratterizzano i fenomeni che possono essere descritti per mezzo di una catastrofe a cuspidi. Il comportamento è sempre bimodale in qualche parte del dominio e si osservano sbalzi improvvisi tra un modo di comportarsi e l'altro. Il salto dal piano superiore al piano inferiore della superficie di comportamento non avviene nella stessa posizione in cui avviene il salto dal piano inferiore a quello superiore, e questo effetto è chiamato isteresi. Tra il piano superiore e quello inferiore esiste una zona inaccessibile sull'asse di comportamento; il piano intermedio, che rappresenta il comportamento meno probabile, è stato omissso per maggiore chiarezza. La catastrofe a cuspidi implica la possibilità di un comportamento divergente.



La liberazione catartica dall'autocommisurazione è descritta per mezzo di una catastrofe a cuspide in cui l'ansia e la frustrazione sono i fattori in conflitto che influenzano lo stato d'animo. L'autocommisurazione viene indotta da un aumento di ansietà; vi si può porre fine con una battuta ironica, che provoca un aumento di frustrazione. Non appena il punto di controllo attraversa la cuspide, lo stato d'animo cambia in maniera catastrofica passando alla collera; l'abbassamento della tensione che ne segue riapre la via a stati emozionali più sereni.

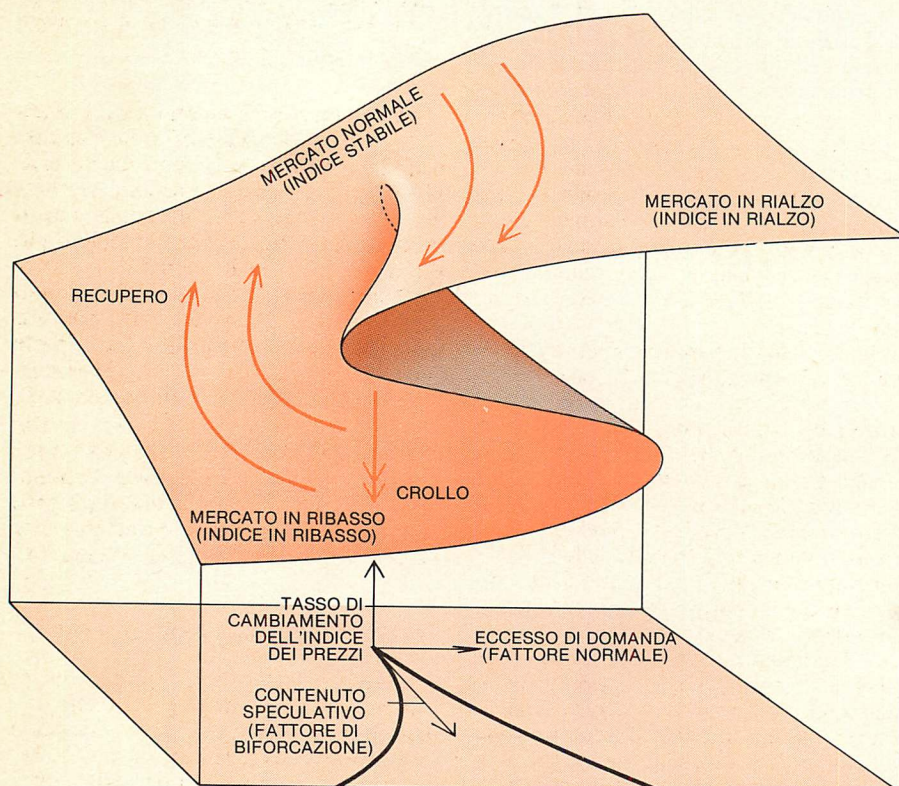
Il procedimento matematico per delineare la superficie di comportamento è un problema di analisi elementare: la superficie di comportamento è un grafico comprendente tutti i punti in cui la derivata prima della funzione dell'energia è uguale a zero. Non è indispensabile comprendere come si effettua questa operazione: è sufficiente sapere che la derivata prima è uguale a zero in ogni parte del grafico in cui la funzione dell'energia è orizzontale (dove la sua pendenza è uguale a zero). La funzione è orizzontale solo ai minimi e ai massimi e ai punti di flessione. I minimi danno luogo ai piani stabili, quello superiore e quello inferiore, i massimi danno luogo al piano intermedio, instabile, e i punti di flessione danno luogo alla curva di piegatura che segna i confini dei piani.

Il comportamento associato alla cuspide più lontana dal disco potrebbe essere analizzato nello stesso modo con cui abbiamo descritto la cuspide più vicina. Le due cuspidi laterali, tuttavia, differiscono per un aspetto fondamentale: la loro funzione dell'energia è rovesciata, in modo che all'interno della cuspide ci sono due punti - che corrispondono a due posizioni del disco - con massimo di energia e solo un punto con minimo di energia. La dinamica fa restare la macchina nell'unica posizione stabile di energia minima. Anche sulla superficie di comportamento la posizione dei minimi e dei massimi si trova rovesciata: il piano intermedio rappresenta minimi stabili di energia, il piano superiore e quello inferiore rappresentano massimi instabili. Il punto di comportamento perciò può trovarsi soltanto sul piano intermedio. Il grafico è chiamato catastrofe a cuspide duale.

Il ruolo della «dinamica»

Ora è possibile spiegare il successo della teoria della catastrofe nel rendere conto del comportamento della macchina delle catastrofi. Il concetto centrale è quello di dinamica, che ha due funzioni. Primo, tiene fermo il punto di comportamento sul piano superiore o sul piano inferiore della superficie di comportamento. Se si gira a mano il disco in senso contrario a quello della tensione delle strisce di gomma, e poi lo si lascia andare, la dinamica lo riporta subito indietro, nella posizione di equilibrio, vale a dire, riporta il punto di comportamento sulla superficie di comportamento. Secondo, quando il punto di comportamento attraversa la curva di piegatura, è la dinamica che provoca il salto catastrofico da un piano all'altro.

Gli stessi principi possono essere applicati ai modelli psicologici considerati più sopra. Le funzioni di probabilità in quei modelli sono analoghe alla funzione dell'energia nella macchina della catastrofe, con la sola eccezione del fatto che le posizioni dei minimi e dei massimi sono ribaltate. Il piano superiore e quello inferiore della superficie di comportamento sono formati da tutti i punti che rappresentano la massima probabilità e



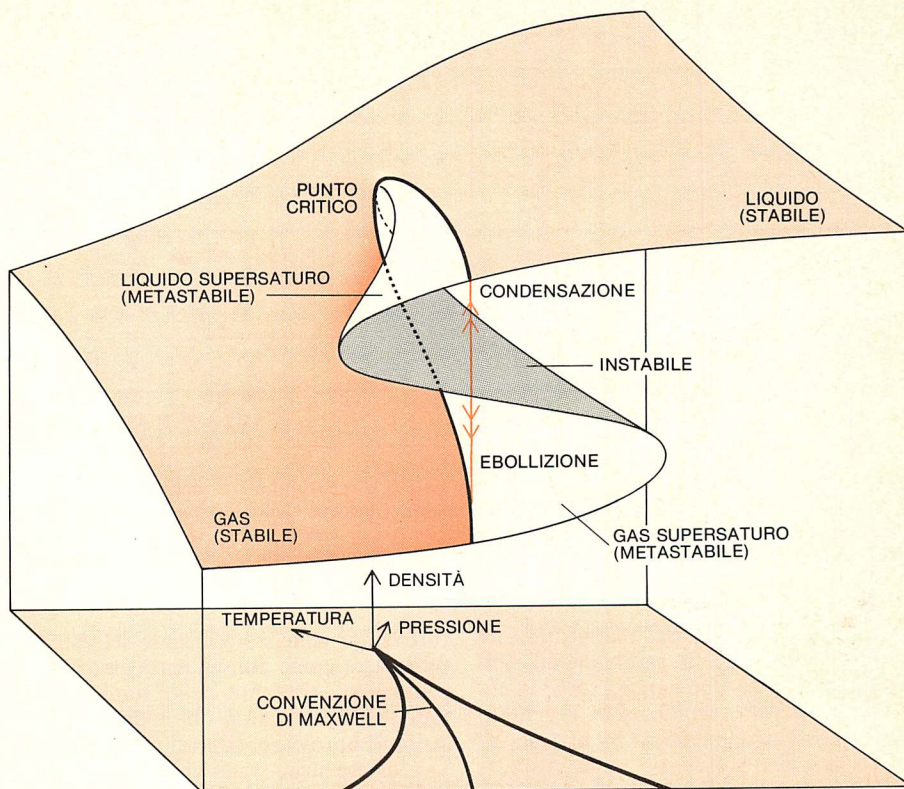
Il comportamento del mercato finanziario viene descritto da un modello in cui i parametri di controllo sono una domanda superiore all'offerta e il rapporto tra il mercato controllato dagli speculatori e quello controllato da chi investe a lungo termine. Il comportamento stesso viene misurato dal tasso con cui l'indice dei prezzi sale o scende. I fattori di controllo sono orientati nel grafico non come fattori in conflitto, ma come fattore normale e fattore di biforcazione. Una caduta dal piano superiore a quello inferiore rappresenta un crollo finanziario; il lento recupero ha luogo per effetto della retroazione dell'indice del prezzo sui parametri di controllo.

il piano intermedio è formato da tutti quelli che rappresentano il minimo della probabilità. Resta una domanda importante: in che cosa consiste la dinamica? Nel modello dell'aggressione che cosa spinge il cane a manifestare il comportamento più probabile, e nel modello dell'autocommiserazione perché è proprio lo stato d'animo «più probabile» quello che viene adottato?

Un minimo di energia in un sistema fisico come quello della macchina della catastrofe è un esempio speciale di un concetto chiamato «attrattore». In questo caso si tratta del tipo più semplice di attrattore, uno stato unico e stabile, e il suo effetto è simile a quello di un magnete: ogni cosa nel suo raggio di influenza è trascinata verso di esso. Sotto l'influenza dell'attrattore il sistema assume uno stato di equilibrio stabile. Anche nei modelli psicologici devono esistere degli attrattori, anche se non è detto che debbano essere semplici come questo. L'attrattore di un sistema che sia in equilibrio dinamico consiste nell'intero ciclo stabile degli stati attraverso i quali passa il sistema. Per esempio una corda di violino strofinata dall'archetto ripete lo stesso ciclo di posizioni più e più volte, alla propria frequenza di risonanza, e questo ciclo di posizioni rappresenta un attrattore per la corda strofinata.

Nei modelli psicologici, l'ovvia sede in cui ricercare gli attrattori è il funzionamento neurologico del cervello. Evidentemente il cervello è di gran lunga più complicato e misterioso che non una corda di violino, ma si sa che i miliardi di neuroni che lo compongono sono organizzati in vaste reti interconnesse che formano un sistema dinamico; gli stati di equilibrio di ogni sistema dinamico possono essere rappresentati come degli attrattori. Alcuni di questi attrattori possono consistere in singoli stati, ma altri sono più simili a cicli stabili di stati o ad analoghi in dimensioni superiori di cicli stabili. Poiché le varie parti del cervello si influenzano a vicenda, gli attrattori compaiono e scompaiono, in certi casi rapidamente, in altri più gradualmente. Non appena un attrattore lascia il posto a un altro, la stabilità del sistema può essere conservata; ma spesso ciò non accade; allora si verifica un salto catastrofico nello stato del cervello.

La teoria di Thom stabilisce che tutti i possibili sbalzi improvvisi tra gli attrattori più semplici - i punti di equilibrio stabile - sono determinati dalle catastrofi elementari. Quindi se la dinamica del cervello avesse solo questo tipo di attrattori, si potrebbero verificare solo catastrofi elementari. Ma non è così; un'indicazione ovvia dell'esistenza di attrattori più complessi è il ritmo alfa delle onde cerebrali: un attrattore ciclico. Non si conoscono ancora le regole che controllano gli sbalzi tra gli attrattori ciclici e quelli di dimensioni superiori; è probabile che comprendano non solo catastrofi elementari, ma anche catastrofi generalizzate, e il loro studio costituisce oggi un'area assai attiva di ricerca matematica. Non esiste oggi una teoria soddisfa-



La fase di transizione dallo stato liquido della materia a quello gassoso è rispecchiata in un modello a cuspide in cui i fattori di controllo sono la temperatura e la pressione. In condizioni normali, ebollizione e condensazione si verificano agli stessi valori di temperatura e pressione, quindi si hanno cambiamenti catastrofici, ma non isteresi. In circostanze speciali, tuttavia, il vapore può venire raffreddato al di sotto del punto di condensazione, e un liquido può essere riscaldato al di sopra del suo punto di ebollizione, così che si segue la superficie del comportamento per tutto il percorso fino alla curva di piegatura. Il punto critico, in cui il liquido e il gas esistono contemporaneamente, è rappresentato dalla singolarità in cui la piega scompare.

cente che descriva tutte le dinamiche del cervello, tuttavia le catastrofi elementari forniscono dei modelli significativi di alcune attività cerebrali. I modelli sono espliciti e a volte di una semplicità cristallina, ma la potenza della teoria matematica che sta alla loro base rende ragione della complessità delle reti nervose sottostanti.

Il concetto di attrattore nella dinamica del cervello fornisce ciò che occorre nei nostri modelli di comportamento umano e animale. Non si conosce bene quale sia il meccanismo responsabile di uno stato d'animo come l'autocommiserazione, ma l'esistenza di questo stato d'animo come stato stabile implica che quel meccanismo è un attrattore. Nel modello dell'autocommiserazione ogni punto della superficie di comportamento corrisponde a un attrattore per il sistema del cervello che determina questo stato d'animo. Se quel sistema neurologico viene in qualche modo disturbato, esso torna rapidamente, sotto l'influenza di un attrattore, alla superficie di comportamento, proprio come la macchina della catastrofe ritorna nel suo stato di equilibrio. Ci si imbatte in cambiamenti bruschi di umore quando si interrompe la stabilità di un attrattore, permettendo al sistema determinante lo stato d'animo di cadere sotto l'influenza di un altro attrattore, verso il quale esso si sposta immediatamente.

Attraverso questo meccanismo ipotetico la teoria delle catastrofi non solo fornisce un modello del comportamento manifestato, ma anche dell'attività cerebrale che lo produce. È probabile che il modello si adatti meglio alle regioni più primitive del cervello, come il mesencefalo, dove le reti nervose hanno un alto grado di integrazione e quindi possono agire come un tutto unico. Nella parte filogeneticamente più recente, la corteccia cerebrale, gli schemi di attività sono assai più complessi. I modelli psicologici che abbiamo considerato riguardano soprattutto l'emozione e lo stato d'animo, e si ritiene che la regolazione degli stati emozionali sia dovuta soprattutto alla parte del mesencefalo detta sistema limbico.

Caratteristiche della catastrofe a cuspide

L'analisi degli esempi precedenti fa pensare che sono molti gli elementi comuni a tutte le catastrofi a cuspide. Una caratteristica costante è che il comportamento è bimodale su una parte del dominio e che si osservano cambiamenti improvvisi nel passaggio da un modo di comportarsi all'altro. Inoltre, il cambiamento improvviso presenta costantemente l'effetto chiamato isteresi, e cioè la transizione dal piano superiore a quello inferiore non avviene nello stesso punto

CATASTROFE		DIMENSIONI DI CONTROLLO	DIMENSIONI DI COMPORTAMENTO	FUNZIONE	DERIVATA PRIMA
CUSPIDI	PIEGA	1	1	$\frac{1}{3}x^3 - ax$	$x^2 - a$
	CUSPIDE	2	1	$\frac{1}{4}x^4 - ax - \frac{1}{2}bx^2$	$x^3 - a - bx$
	CODA DI RONDINE	3	1	$\frac{1}{5}x^5 - ax - \frac{1}{2}bx^2 - \frac{1}{3}cx^3$	$x^4 - a - bx - cx^2$
	FARFALLA	4	1	$\frac{1}{6}x^6 - ax - \frac{1}{2}bx^2 - \frac{1}{3}cx^3 - \frac{1}{4}dx^4$	$x^5 - a - bx - cx^2 - dx^3$
OMBELICHI	IPERBOLICO	3	2	$x^3 + y^3 + ax + by + cxy$	$3x^2 + a + cy$ $3y^2 + b + cx$
	ELLITTICO	3	2	$x^3 - xy^2 + ax + by + cx^2 + cy^2$	$3x^2 - y^2 + a + 2cx$ $-2xy + b + 2cy$
	PARABOLICO	4	2	$x^2y + y^4 + ax + by + cx^2 + dy^2$	$2xy + a + 2cx$ $x^2 + 4y^3 + b + 2dy$

Sette catastrofi elementari descrivono tutte le possibili discontinuità nei fenomeni controllati da non più di quattro fattori. Ciascuna delle catastrofi è associata a una funzione in cui i parametri di controllo sono rappresentati come coefficienti (a, b, c, d) e il comportamento

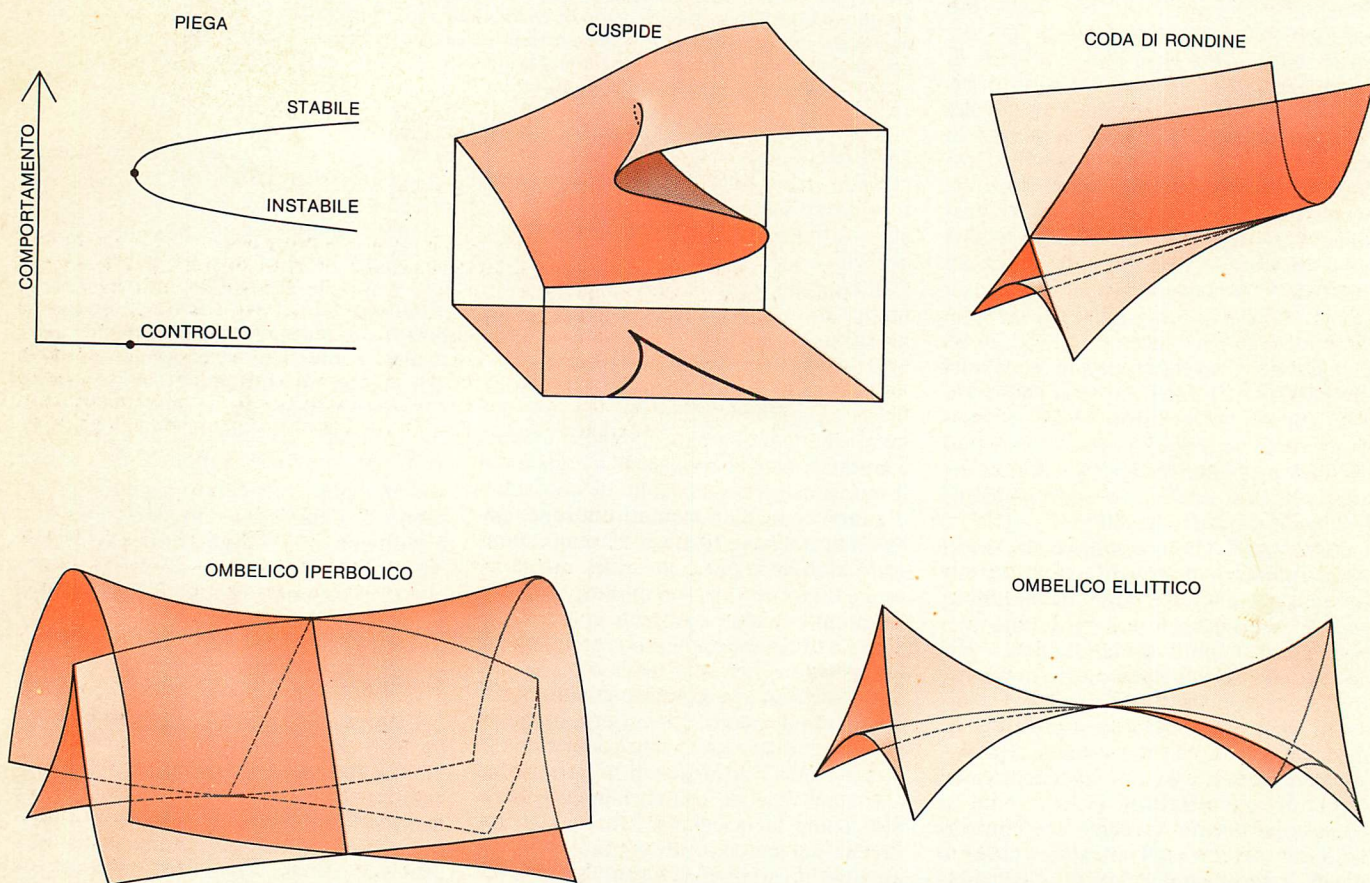
del sistema è determinato dalle variabili (x, y). In ogni modello di catastrofe, la superficie di comportamento è il grafico di tutti i punti in cui la derivata prima di questa funzione è uguale a zero o, se vi sono due derivate prime, i punti in cui sono entrambe uguali a zero.

in cui avviene quella dal piano inferiore a quello superiore. Il cambiamento improvviso non avviene nel mezzo della cuspide, bensì solo quando viene raggiunto l'insieme di biforcazione. Un'altra caratteristica è che all'interno della cuspide, dove il comportamento si fa bimodale, la zona centrale sull'asse del comportamento diventa inaccessibile. Infine, il modello implica la possibilità di una divergenza, in modo che una leggera per-

turbazione nello stato iniziale del sistema può avere come conseguenza una differenza assai rilevante nello stato finale. Queste cinque qualità - bimodalità, transizioni improvvise, isteresi, inaccessibilità e divergenza - sono correlate tra loro dal modello stesso. Se una qualsiasi di queste qualità si evidenzia in un processo, bisognerebbe ricercare le altre quattro, e se se ne trova più di una, allora il processo dovrebbe essere considerato

passibile di una descrizione per mezzo della catastrofe a cuspide.

Un processo in cui possono essere individuate le cinque qualità è quello dello sviluppo delle ostilità tra nazioni, situazione che presenta delle ovvie analogie con i modelli del litigio e della aggressione. In quei modelli i parametri erano la collera e la paura; qui vi sostituiamo la minaccia e il costo. L'asse del comportamento descrive le azioni possibili della



I grafici di cinque catastrofi elementari chiariscono le loro possibilità geometriche. La catastrofe a piega è una sezione trasversale di una curva di piegatura della catastrofe a cuspide, e il suo insieme di biforcazione consiste di un solo punto. La catastrofe a cuspide è tra le catastrofi

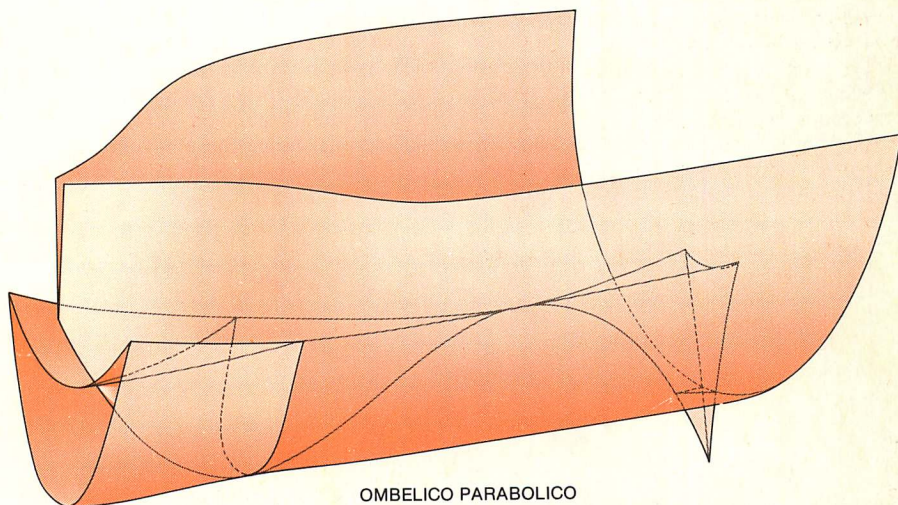
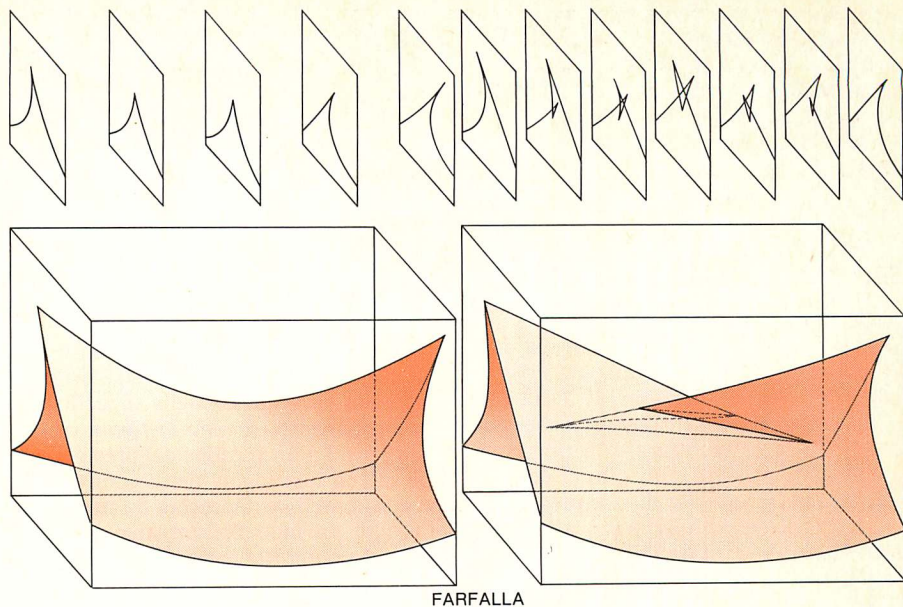
disegnabili quella di ordine massimo. La coda di rondine è a quattro dimensioni e l'ombelico iperbolico e l'ombelico ellittico sono a cinque dimensioni. Per queste è possibile disegnare solo gli insiemi di biforcazione a tre dimensioni; le superfici di comportamento non compaiono.

nazione, che variano dall'attacco generale, a risposte militari meno vistose, come i blocchi militari, fino alla sospensione delle ostilità e alla resa. In una situazione in cui sono alti sia la minaccia sia il costo, l'opinione pubblica diventa bimodale, poiché la nazione è divisa in «colombe», favorevoli alla resa, e in «falchi», favorevoli all'attacco. Nel modello, la dinamica è la sensibilità del governo nei confronti dei propri elettori; il governo cerca sempre di modificare la propria politica al fine di aumentare il consenso, e perciò esso resta sulla superficie di comportamento. Possiamo dedurre e riconoscere per mezzo del modello le catastrofi possibili. Una nazione minacciata può fare ripetute concessioni, ma c'è un limite al di là del quale un'ulteriore intimidazione provoca un'improvvisa dichiarazione di guerra. Viceversa, con l'aumento del costo una nazione può continuare a intensificare l'azione di guerra, ma c'è un limite al di là del quale un'ulteriore aumento del costo può portare a una resa improvvisa. L'isteresi può essere osservata nel ritardo con cui può avvenire la dichiarazione di guerra o la resa. La regione inaccessibile dell'asse del comportamento è la zona intermedia che rappresenta il negoziato o il compromesso. Infine, la divergenza si può osservare in un conflitto tra due nazioni di uguale potenza, in cui l'andamento dell'opinione pubblica è simile, ma la risposta del governo è completamente diversa: uno diventa sempre più aggressivo, l'altro sempre più remissivo.

Un altro fenomeno che può essere analizzato mediante una catastrofe a cuspidè è il mercato finanziario, dove i termini di «mercato in rialzo» e «mercato in ribasso» suggeriscono un'ovvia bimodalità. E ancora, un crollo o un collasso del mercato viene prontamente illustrato dal grafico con un salto catastrofico da un piano all'altro della superficie di comportamento.

Per costruire questo modello è necessaria una piccola modifica. In questo caso gli assi di controllo non divergono su ciascun lato della cuspidè, come avviene nei casi precedenti; invece un asse viene in avanti sul grafico, bisecando la cuspidè, e l'altro è perpendicolare alla cuspidè (si veda la figura in basso a pagina 164). Il parametro che biseca la cuspidè è chiamato fattore di biforcazione poiché un suo aumento provoca una divergenza sempre crescente tra il piano superiore e quello inferiore. L'altro fattore è il fattore normale, poiché nella parte posteriore della superficie di comportamento il comportamento si sviluppa in modo continuo rispetto a esso.

Nel modello del mercato finanziario il fattore normale è una domanda eccedente rispetto all'offerta; il fattore di biforcazione è forse più difficile da identificare, ma potrebbe essere messo in relazione all'ammontare delle azioni controllate dagli speculatori rapportato a quello in possesso di coloro che investono a lungo termine. Un mercato con un indice che sale è un mercato in rialzo, e il suo punto di comportamento si trova sul pia-



Le sezioni sono l'unico espediente che permette di rappresentare graficamente le ultime due catastrofi, poiché anche i loro insiemi di biforcazione hanno più di tre dimensioni. L'insieme di biforcazione a quattro dimensioni della catastrofe a farfalla è mostrato in sezioni a tre dimensioni; la quarta dimensione è il fattore farfalla, e se è il tempo, una configurazione si trasforma nell'altra. Il movimento da sinistra a destra in ogni disegno equivale al cambiamento del fattore inclinazione; le sezioni a due dimensioni illustrano assai più chiaramente l'effetto di questo fattore. L'insieme di biforcazione a quattro dimensioni della catastrofe a ombelico parabolico è anch'esso visibile solo in una sezione tridimensionale ed è il risultato di un tracciato elaborato con l'aiuto di un calcolatore da A.N. Godwin del Lancaster Polytechnic in Inghilterra.

no superiore; un indice in diminuzione, ossia un mercato in ribasso, vede il punto di comportamento sul piano inferiore.

Si può ora comprendere il meccanismo del crollo. Un mercato con una domanda in eccesso e una grande proporzione di speculatori è un mercato in rialzo sulla superficie di comportamento superiore. Un crollo può essere provocato da qualsiasi evento che riduca la domanda fino a spingere il punto di comportamento al di là della curva di piegatura. Più grande è la parte del mercato controllata dagli speculatori, più grave sarà il crollo. Si potrebbe subito chiedere come mai la ripresa che segue è generalmente lenta, e come mai non c'è un «crollo di risalita», da un mercato in ribasso a uno in rialzo. La risposta più plausibile a questo interrogativo è che

l'asse del comportamento (il tasso di cambiamento dell'indice) ha un'influenza sui parametri di controllo per mezzo di un ciclo di retroazione. Un mercato in ribasso scoraggia la speculazione, ma dopo un certo tempo la svalutazione dei titoli che ne consegue incoraggia gli investimenti a lungo termine e di conseguenza, dopo un crollo, il fattore di biforcazione si riduce e il mercato ritorna a spostarsi sul grafico verso una regione in cui la superficie di comportamento non è più bimodale. Con il crescere della fiducia e con il prodursi di un eccesso di domanda l'indice si rialza, ma a poco a poco e in maniera continua, senza catastrofi. Ora è la speculazione a essere incoraggiata e gli investimenti a essere scoraggiati, ed esistono di nuovo le condizioni ideali per la formazione di un nuo-

vo ciclo di espansione e di depressione.

Come si è precisato più sopra, la maggior parte delle applicazioni della teoria delle catastrofi sono state fatte in biologia e nelle scienze sociali, laddove altre tecniche di costruzione di modelli risultano inefficaci ai fini dell'informazione; ma si danno molte situazioni nel campo della fisica (scienza che dispone di un linguaggio matematico altamente sviluppato), in cui questa teoria può offrire un utile contributo alla comprensione degli eventi. Una di queste situazioni è ad esempio la transizione della materia dallo stato liquido a quello gassoso. Possiamo riscrivere le equazioni di J.D. van der Waals come una catastrofe a cuspidi, in cui la temperatura e la pressione sono i due fattori di controllo in conflitto e la densità è l'asse di comportamento. La superficie superiore è allora lo stato liquido e la superficie inferiore è lo stato gassoso; le due catastrofi rappresentano l'ebollizione e la condensazione. Il vertice della cuspidi è il punto critico, in cui il liquido e il gas esistono nello stesso tempo. Girando attorno alla parte posteriore della cuspidi, un liquido può diventare gas senza passare attraverso l'ebollizione.

In circostanze particolari il sistema fisico può seguire perfettamente la superficie di comportamento fino al bordo di

entrambi i piani. È possibile, con le dovute precauzioni, per esempio, scaldare l'acqua molto al di là del suo normale punto di ebollizione, e il vapore acqueo può essere raffreddato molto al di sotto del punto di condensazione prima che inizi la fase di transizione. Questo tipo di riscaldamento e di raffreddamento oltre al punto critico è sfruttato nelle camere a bolle e nelle camere a nebbia usate per individuare particelle subatomiche. Normalmente, tuttavia, una sostanza bolle e si condensa a un medesimo valore di temperatura e pressione, così che sezionando la parte intermedia della piega si forma un «precipizio» nella superficie di comportamento (si veda la figura di pagina 165). La formazione del precipizio si spiega con una regola detta «convenzione di Maxwell» e rispecchia il fatto che si tratta di un modello statistico; che descrive il comportamento medio di molte particelle.

Un altro splendido esempio nel campo dell'ottica è offerto dalla caustica che si forma quando la luce si riflette o si rifrange su una superficie curva. Una caustica assai nota è la curva a forma di cuspidi che talvolta compare sulla superficie di una tazzina da caffè che brilla alla luce del Sole, che è formata dal riflesso dei raggi del Sole sulla superficie interna della tazza.

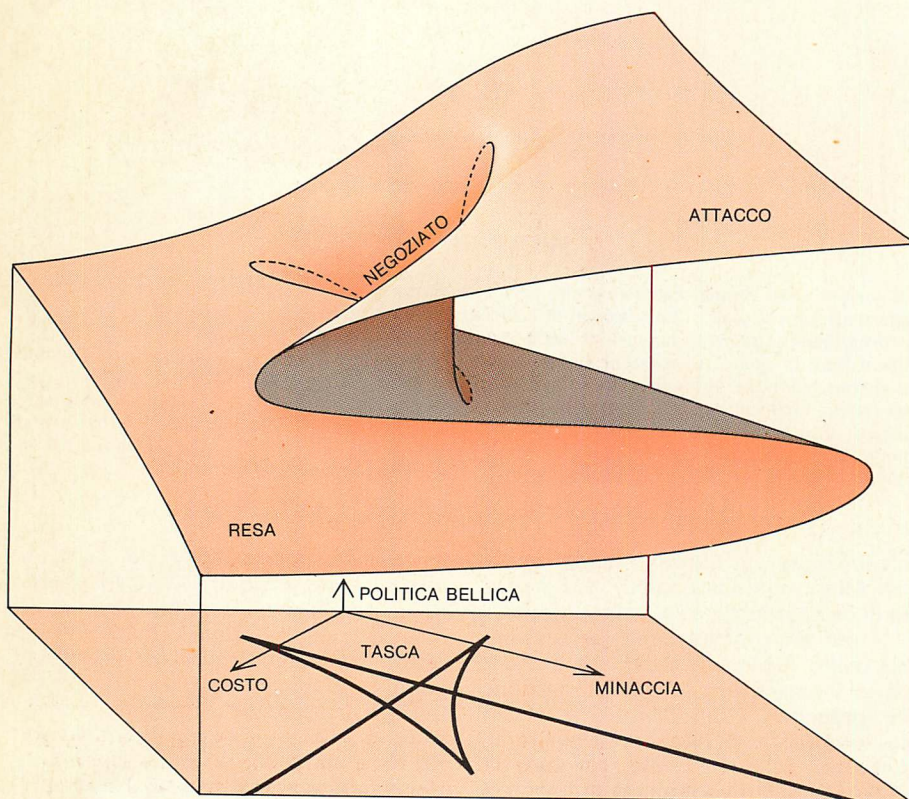
Un altro esempio altrettanto comune di caustica, che presenta discontinuità sia temporali sia spaziali della luminosità, è il disegno mutevole prodotto dai raggi di Sole sul fondo di una piscina. L'arcobaleno è un insieme di conici caustici colorati. Altre caustiche più complesse possono essere prodotte facendo riflettere un raggio di sole attraverso uno specchio concavo o attraverso lenti sferiche o cilindriche (come una lampadina o una coppa piena d'acqua). In questa applicazione la teoria della catastrofe ha permesso una migliore comprensione del fenomeno; Thom ha dimostrato che le caustiche stazionarie possono avere solo tre tipi di punti singolari. Una sottigliezza matematica dell'analisi delle caustiche per mezzo della teoria della catastrofe è che non esiste dinamico; esiste invece un principio variazionale che dà uguale importanza sia ai minimi sia ai massimi.

La catastrofe a cuspidi è una figura tridimensionale: due dimensioni servono per i due parametri di controllo e una terza serve per l'asse di comportamento. Di fatto l'asse di comportamento non è detto che rappresenti la variazione di un singolo comportamento; nei modelli del funzionamento del cervello, per esempio, può rappresentare gli stati di miliardi di neuroni, che variano tutti contemporaneamente. Tuttavia la teoria delle catastrofi dimostra che è sempre possibile selezionare una singola variabile e tracciare la superficie di comportamento soltanto rispetto a quell'asse, così da ottenere il familiare grafico tridimensionale. Se si riduce il grafico a due dimensioni, ne risulta un modello ancora più semplice: la catastrofe a piega. Nella catastrofe a piega c'è solo un parametro di controllo; lo spazio di controllo è una linea retta e l'insieme di biforcazione è un unico punto su quella linea. Lo spazio di comportamento è una parabola, metà della quale rappresenta gli stati stabili, e l'altra metà quelli instabili. Le due regioni sono separate da un punto di piegatura direttamente sopra il punto di biforcazione.

Il teorema di classificazione

La catastrofe a piega può essere considerata come una sezione trasversale della curva di piegatura della catastrofe a cuspidi. A sua volta, la cuspidi può essere considerata come una quantità di catastrofi a piega, con un nuovo punto di singolarità all'origine. Sullo stesso schema si possono costruire catastrofi più complicate, di dimensioni superiori: ognuna di queste è formata da tutte le catastrofi di ordine inferiore, più una nuova singolarità all'origine.

Se lo spazio di controllo diventa tridimensionale mentre lo spazio di comportamento resta a una dimensione, si può costruire un'unica catastrofe a quattro dimensioni. La superficie di comportamento diventa un'ipersuperficie a tre dimensioni, e anziché essere piegata lungo delle curve, come nella catastrofe a cuspidi, è piegata lungo delle intere superfici, configurazione questa che non può



La catastrofe a farfalla permette di prevedere l'emergere di un'opinione di compromesso in un modello che descrive lo sviluppo di una politica bellica. Nella catastrofe a farfalla sono necessari quattro fattori di controllo, ma qui ne sono visibili solo due (la minaccia e il costo) mentre gli altri due si assumono costanti. L'insieme di biforcazione è una delle sezioni mostrate nella pagina precedente; si tratta di una curva complessa, con tre cuspidi e una «tasca» nel mezzo. Sulla superficie di comportamento in corrispondenza della tasca c'è un nuovo piano, che serve a rappresentare un nuovo modo di comportamento, intermedio rispetto agli altri. Se sono alti sia il valore della minaccia sia quello del costo, il modello a cuspidi terrà conto solo delle posizioni estreme che si pronunciano per l'attacco o per la resa. Il nuovo piano nel modello a farfalla rappresenta l'emergere di un'opinione favorevole al compromesso che si pronuncia per il negoziato.

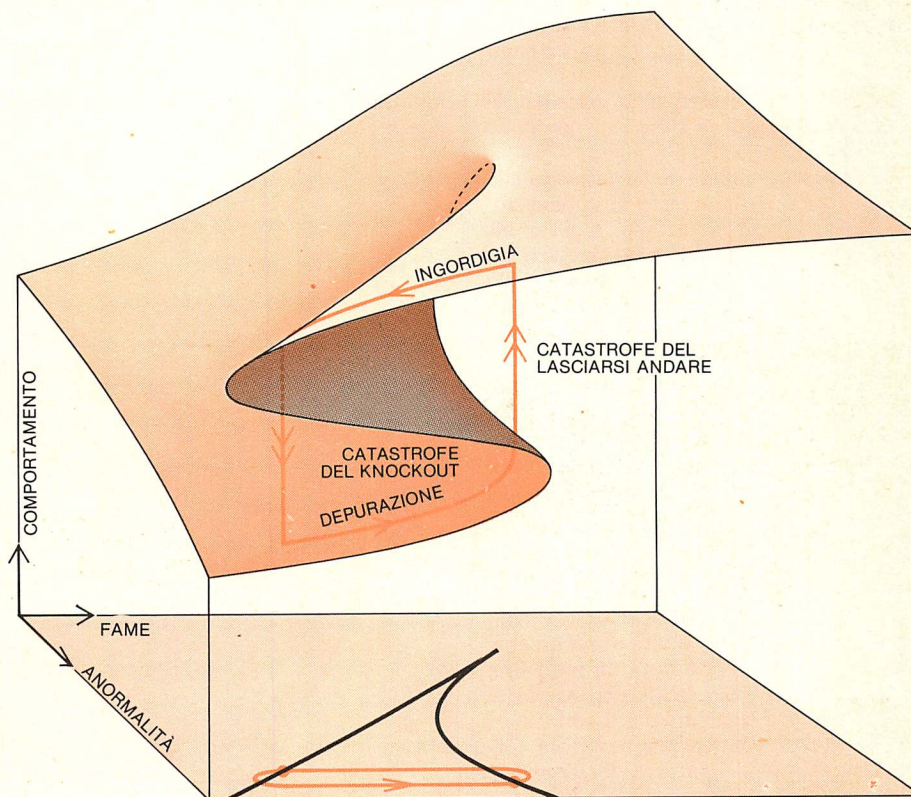
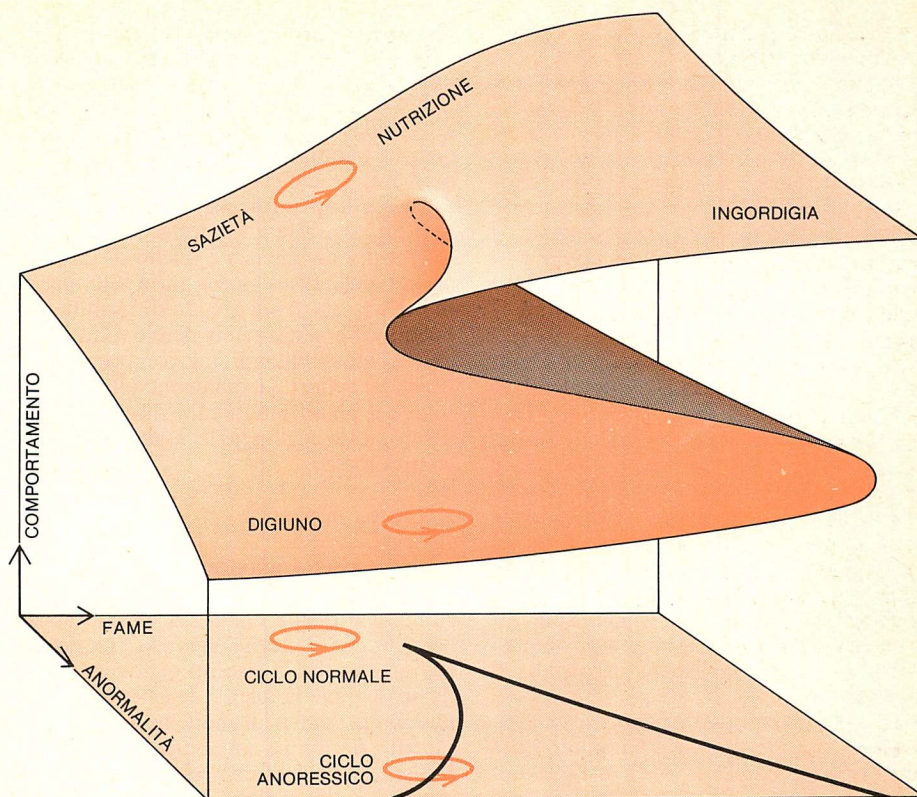
facilmente essere visualizzata. L'insieme di biforcazione non consiste più di curve che formano una cuspide a due dimensioni, ma è costituito da superfici a tre dimensioni che si intersecano dando luogo a cuspidi ai loro bordi. All'origine compare una nuova singolarità, detta catastrofe a coda di rondine. È impossibile rappresentare la catastrofe a coda di rondine in modo completo, per l'impossibilità di disegnare immagini a 4 dimensioni. Possiamo tuttavia rappresentare graficamente il suo insieme di biforcazione, che è tridimensionale, e da questo disegno è possibile derivare alcune intuizioni geometriche circa la coda di rondine, allo stesso modo in cui è possibile descrivere la catastrofe a cuspide disegnando il suo insieme di biforcazione (la cuspide) in due dimensioni, e ricordando che la superficie di comportamento è bimodale sopra la parte interna della cuspide.

Se aggiungiamo ancora un altro parametro di controllo, si forma una catastrofe a 5 dimensioni. La piega, la cuspide e la coda di rondine appaiono anche qui come sezioni, e una nuova singolarità viene associata a una «tasca» formata dall'interpenetrazione di parecchie superfici. La forma di questa tasca, o delle sue sezioni, ha suggerito il nome di catastrofe a farfalla. Nella catastrofe a farfalla anche l'insieme di biforcazione è a quattro dimensioni e perciò non può essere disegnato. Lo si può illustrare solo mediante sezioni a due o a tre dimensioni (si veda l'illustrazione di pagina 167).

Ci sono ancora due catastrofi a cinque dimensioni, che si formano quando lo spazio di controllo ha tre dimensioni e lo spazio di comportamento ha due dimensioni. Sono chiamate la catastrofe a ombelico iperbolico e la catastrofe a ombelico ellittico. Come nel caso della coda di rondine, i loro insiemi di biforcazione consistono di superfici con i bordi a cuspide, e possono essere disegnate dal momento che sono a tre dimensioni. Infine, la catastrofe a sei dimensioni generata da uno spazio di controllo a quattro dimensioni e da uno spazio di comportamento a due dimensioni è chiamata ombelico parabolico. La sua geometria è complessa, e ancora una volta si può rappresentare graficamente solo il suo insieme di biforcazione.

Col crescere del numero delle dimensioni dello spazio di controllo e dello spazio di comportamento si può costruire una serie infinita di catastrofi. Il matematico russo V.I. Arnold le ha classificate fino ad almeno 25 dimensioni. Per modelli che riguardano fenomeni del mondo reale, tuttavia, le 7 descritte sopra sono probabilmente le più importanti, poiché sono le uniche che hanno uno spazio di controllo che non ha più di quattro dimensioni. Una classe particolarmente comune di processi, quelli determinati dalla posizione nello spazio e nel tempo, non possono richiedere uno spazio di controllo con più di quattro dimensioni, poiché il nostro mondo ha solo tre dimensioni spaziali e una temporale.

Anche le catastrofi che non possono



L'anoressia nervosa, un disturbo nervoso che si manifesta in adolescenti e giovani donne che si sottopongono a digiuni ossessivi, può essere descritto per mezzo di una catastrofe a farfalla. Due dei fattori di controllo sono la fame e un atteggiamento anormale verso il cibo. Negli individui normali la fame conduce a un ciclo di comportamento che oscilla tra il nutrirsi e la sazietà; nell'individuo anoressico, con atteggiamenti anormali, il medesimo ciclo di fame conduce a comportamenti assai differenti. Nella prima fase della malattia (*in alto*) il ciclo resta nel piano inferiore della superficie di comportamento, e l'individuo anoressico rimane costantemente in uno stato mentale concentrato sul digiuno. La seconda fase (*in basso*) viene determinata dal cambiamento di un terzo fattore, il controllo cosciente. Quando la paziente perde il controllo di se stessa dopo un periodo di due anni o più, l'insieme di biforcazione viene deviato a sinistra fino a quando il ciclo della fame attraversa il lato destro della cuspide. A questo punto la persona entra in un ciclo di isteresi; digiuna fino a quando la fame la costringe a «lasciarsi andare» in maniera catastrofica, e allora si abbuffa fino a che, dopo una catastrofe di «knock-out», ricomincia a digiunare e a depurarsi da quello che essa vive come una contaminazione.

essere rappresentate possono essere impiegate per descrivere i fenomeni per mezzo di modelli. La loro geometria è completamente determinata, e il movimento di un punto sulla superficie di comportamento può essere studiato analiticamente anche se non può essere rappresentato graficamente. Ogni catastrofe è definita da una funzione potenziale, e in ogni caso la superficie di comportamento è il grafico di tutti i punti in cui la derivata prima di quella funzione è uguale a zero.

Il punto di forza della teoria di Thom sta nella sua generalità e completezza. Essa stabilisce che se un processo viene determinato dalla minimizzazione o dalla massimizzazione di alcune funzioni, e se è controllato da non più di 4 fattori, allora ogni e qualsiasi singolarità della superficie di comportamento che ne risulta deve essere analoga a una delle sette catastrofi che abbiamo descritto. Se il processo è controllato da due soli fattori, allora la superficie di comportamento può avere solo pieghe e cuspidi. Il teorema stabilisce, insomma, che in ogni processo regolato da due fattori, la catastrofe a cuspidi è la cosa più complicata che può verificarsi nel grafico. La dimostrazione del teorema è troppo tecnica e troppo lunga per poter essere presentata qui, ma le sue conseguenze sono molto chiare: ogni volta che una forza che cambia in maniera continua produce come effetto un cambiamento improvviso, il processo deve essere descritto mediante una catastrofe.

Dopo la catastrofe a cuspidi, quella che ha lo spettro più vasto di applicazioni è la catastrofe a farfalla. Così come un comportamento bimodale determina il modello a cuspidi, allo stesso modo un comportamento trimodale determina quello a farfalla. Nel modello a cuspidi della politica bellica, per esempio, in cui l'opinione pubblica è divisa tra «colombe» e «falchi», il modello a farfalla contempla l'emergere di un'opinione di compromesso favorevole al negoziato. Questa nuova alternativa si presenta come un nuovo piano della superficie di comportamento, che sale in maniera continua nella parte posteriore della piega (si veda l'illustrazione a pagina 168).

La geometria della catastrofe a farfalla è controllata da quattro parametri. Due di questi sono già noti dai modelli a cuspidi: il fattore normale e il fattore di biforcazione. I due rimanenti sono nuovi: il fattore divergenza e il fattore farfalla. L'effetto del fattore divergenza è quello di alterare la posizione e la forma della cuspidi: esso sposta la parte principale della cuspidi a destra o a sinistra, mentre il vertice della cuspidi si piega nella direzione opposta. Nello stesso tempo, il fattore divergenza sposta la superficie di comportamento in su e in giù.

L'effetto del fattore farfalla è quello di dar luogo al terzo modo stabile di comportamento. Con l'aumento del fattore farfalla, la cuspidi sulla superficie di controllo si sviluppa in tre cuspidi, che formano una «tasca» triangolare. Sopra la tasca c'è la nuova superficie

triangolare, sulla superficie di comportamento, tra il piano superiore e quello inferiore.

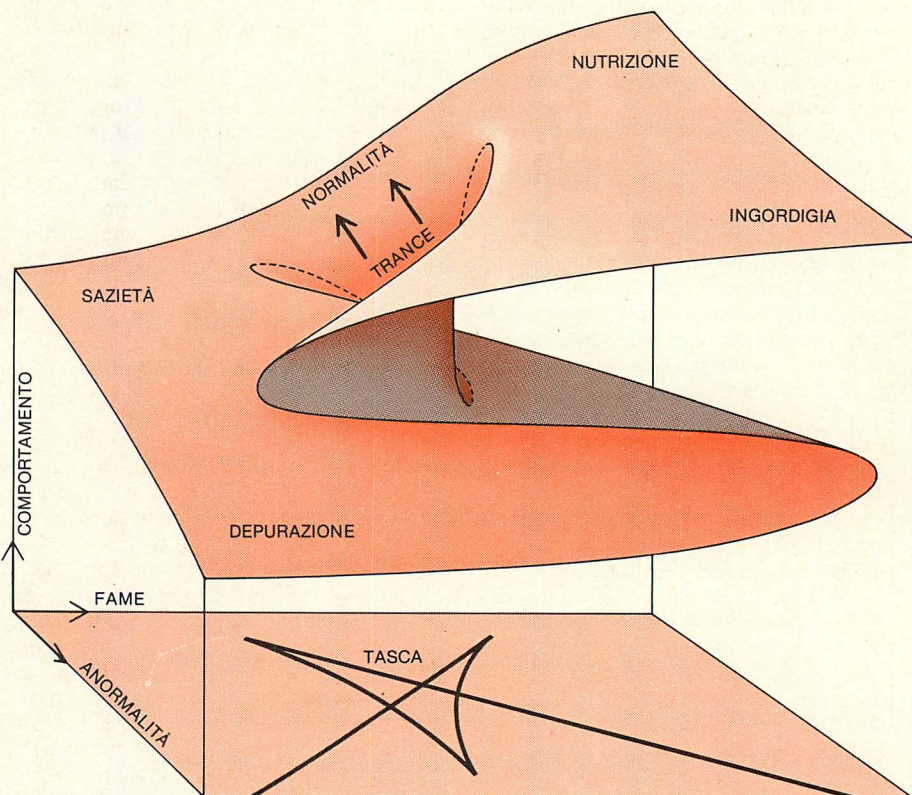
Per tracciare il grafico della catastrofe a farfalla occorre sopprimere due dei quattro parametri, e di solito si eliminano il fattore inclinazione e il fattore farfalla. La loro influenza sul grafico non può tuttavia essere ignorata. Un effetto del fattore divergenza è quello di ridurre un lato della tasca fino a farla scomparire in una catastrofe a coda di rondine; il fattore divergenza tende perciò a distruggere la possibilità di un compromesso. Poiché il fattore farfalla controlla la crescita del piano intermedio di comportamento, esso sviluppa la stabilità di un compromesso.

Anoressia nervosa

Una seconda applicazione della catastrofe a farfalla, di una straordinaria fecondità, è quella che si fa all'anoressia nervosa, un disturbo nervoso frequente soprattutto in adolescenti e giovani donne per le quali la dieta è degenerata in un digiuno ossessivo. Il modello è stato sviluppato da me in collaborazione con J. Hevesi, uno psicoterapista inglese che ha introdotto la terapia per mezzo della *trance* nel trattamento dell'anoressia. In una recente indagine, risulta che su 1000 pazienti sofferenti di anoressia, quelli che erano stati curati da lui erano i soli dei quali si potesse affermare che fossero guariti completamente.

Nella fase iniziale dell'anoressia il digiuno ossessivo può condurre all'inedia, e in qualche caso anche alla morte. Col passare del tempo, gli atteggiamenti della paziente verso il cibo diventano progressivamente anormali. Dopo circa due anni si sviluppa di solito una seconda fase, chiamata bulimia, durante la quale la malata alternativamente digiuna e si rimpinzia di cibo. Il comportamento bimodale di questa seconda fase suggerisce immediatamente una catastrofe a cuspidi. La malata di anoressia precipita in un ciclo di isteresi, saltando in maniera catastrofica tra due estremi, e non riesce ad approdare al comportamento normale che sta tra questi due. La teoria delle catastrofi suggerisce anche una cura teorica: se si potesse indurre una nuova biforcazione, come nella catastrofe a farfalla, si potrebbe riaprire la via alla normalità.

La superficie di comportamento in questo modello rappresenta il comportamento manifesto della paziente, che varia dal riempirsi esageratamente di cibo, attraverso il cibarsi normalmente a sazietà, fino al digiuno ossessivo. Essa fornisce anche alcune indicazioni circa gli stati cerebrali sottostanti al comportamento; come nel modello dell'aggressione, abbiamo a che fare con degli stati emozionali che hanno probabilmente origine nel sistema limbico. I dati raccolti dagli psicologi suggeriscono che la variabile del comportamento può costituire effettivamente una misura dell'importanza re-



Il trattamento dell'anoressia si fonda sulla creazione di un terzo modo di comportarsi intermedio. Il nuovo comportamento è reso possibile aumentando il quarto parametro di controllo della catastrofe a farfalla: la rassicurazione. Ciò crea una tasca nell'insieme di biforcazione e quindi un piano intermedio nella superficie di comportamento. In un primo tempo la paziente entra ed esce dalla trance con salti catastrofici dal piano intermedio al piano superiore e a quello

lativa che il sistema limbico assegna agli stimoli provenienti dal corpo in contrapposizione a quelli provenienti dalla corteccia cerebrale. In una persona normale questi stimoli possono essere in un certo senso bilanciati, ma nella malata anoressica l'uno o l'altro tendono a essere dominanti.

Tra i parametri di controllo, il fattore normale è la fame, che nelle persone normali regola il ciclo ritmico di nutrizione e sazietà. Il fattore di biforcazione è il grado di anormalità degli atteggiamenti della paziente anoressica verso il cibo; l'anormalità aumenta man mano che la sua condizione peggiora.

Il fattore divergenza nel grafico a farfalla è la perdita dell'autocontrollo, che può essere misurata con la perdita di peso. Nella prima fase del disturbo gli atteggiamenti della malata di anoressia sono già anormali, ma essa mantiene l'autocontrollo. Il risultato è che essa rimane relegata sul piano inferiore della superficie di comportamento; per tutto il tempo della veglia, il sistema limbico resta in uno stato corrispondente a uno stato d'animo concentrato sul digiuno, anche quando prende i suoi pasti ridottissimi.

Man mano che la malata anoressica diminuisce di peso, perde anche l'autocontrollo, e il fattore divergenza aumenta gradualmente. Il risultato è che la cuspidi oscilla verso la sinistra del grafico (si veda l'illustrazione in alto a pagina 169); se si sposta abbastanza lontano, la parte destra della cuspidi interseca il ciclo anoressico, provocando l'inizio improvviso della seconda fase. Ora la malata non è più imprigionata in un ciclo di digiuno costante, ma cade in un ciclo di isteresi, saltando avanti e indietro dal piano inferiore a quello superiore. Con le parole di un'anoressica tipica, il salto catastrofico dal digiuno all'ingordigia avviene quando la persona «si lascia andare» e osserva indifesa il «mostro dentro di lei» che si abbuffa di cibo per diverse ore, a volte anche vomitando. Il ritorno

catastrofico al digiuno avviene quando la stanchezza, il disgusto e l'umiliazione la travolgono, esperienza questa che molte malate chiamano *knockout*.

Il periodo di digiuno che segue il *knockout* nel ciclo di isteresi è diverso dal digiuno regolare della prima fase. Si trova in una posizione diversa sull'asse di comportamento e si potrebbe chiamare più appropriatamente fase di depurazione. Lo stato limbico che è associato al primo digiuno è dominato dagli impulsi provenienti dalla corteccia cerebrale ed è finalizzato a impedire l'entrata del cibo. Durante il periodo di ingordigia il sistema limbico è dominato da impulsi provenienti dal corpo. Lo stato limbico sottostante il periodo seguente di depurazione è ancora dominato da impulsi cerebrali, ma di questi una parte proviene dal corpo ed è diretta a liberarlo dalla contaminazione.

La terapia della trance usata da Hevesi rassicura la paziente, riduce la sua incertezza e la mette perciò in grado di riguadagnare la via verso il comportamento normale. Le malate anoressiche tendono a dormire irregolarmente, e quando sono sveglie vivono esperienze spontanee simili alla trance; è su questi periodi che interviene il terapeuta. La trance può rappresentare il terzo stato del sistema limbico, nella zona altrimenti inaccessibile tra gli stati di ingordigia e di depurazione. Quando la paziente digiuna vive con grande ansietà il mondo esterno, e quando si abbuffa se ne sente sopraffatta; ma durante la trance essa è isolata, la sua mente è libera sia dal cibo sia dalla preoccupazione di evitarlo. È solo a questo punto che è possibile intervenire con la rassicurazione.

La rassicurazione diviene il fattore farfalla nel modello. Essa dà luogo a un nuovo piano sulla superficie di comportamento, piano che si trova tra gli altri due e che eventualmente può dare accesso alla regione stabile, della normalità, dietro la cuspidi. Poiché la terapia generalmente avviene durante il momento del digiuno del ciclo, l'entrata nella trance equivale a un salto catastrofico dal piano inferiore a quello intermedio. L'uscita dalla trance è un'altra catastrofe, che può condurre la paziente sia al piano inferiore sia a quello superiore.

Dopo circa due settimane di trattamento e circa nella settima seduta di trance, l'atteggiamento anormale della paziente nei confronti del cibo si esaurisce in maniera catastrofica, e la personalità si riunifica di nuovo in un tutto integrato. Quando la paziente si risveglia dalla trance, ne può parlare come di un «momento di rinascita» e scopre che può mangiare di nuovo senza paura di abbuffarsi. La trance ha apparentemente riaperto una via di ritorno del cervello a stati limbici più equilibrati, così che la paziente può riacquistare un comportamento normale. La seduta di trance seguente rinforza l'esperienza.

Uno dei punti di forza del modello della teoria della catastrofe sull'anoressia nervosa è che esso rende conto della descrizione che la paziente fa del proprio

stato. I termini quasi incomprensibili con i quali alcune malate di anoressia descrivono il loro malessere si rivelano piuttosto ovvi considerati nel quadro delle superfici della catastrofe. Il vantaggio di un linguaggio matematico in questo tipo di applicazioni è che è psicologicamente neutro. Esso permette una sintesi di osservazioni che potrebbero altrimenti sembrare sconnesse.

Il futuro della teoria delle catastrofi

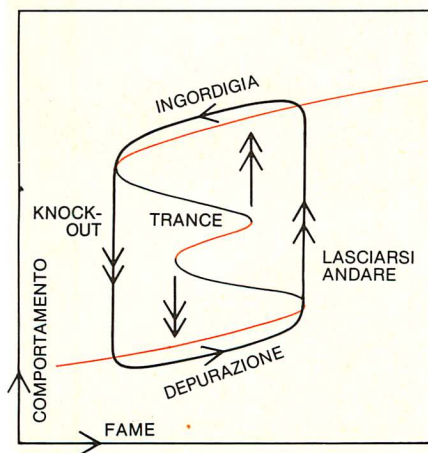
La teoria della catastrofe è assai recente: Thom ha pubblicato il primo articolo nel 1968. Finora il suo impatto maggiore è avvenuto nei confronti della matematica stessa; in particolare ha stimolato altri settori della matematica che erano necessari per dimostrare i suoi teoremi. I principali problemi insoluti di questa teoria riguardano la comprensione e la classificazione delle catastrofi generalizzate e le catastrofi più sofisticate che si presentano quando si impongono condizioni di simmetria. Inoltre esistono problemi connessi al modo in cui la teoria delle catastrofi può essere usata insieme ad altri metodi e concetti matematici.

Si stanno studiando nuove applicazioni di questa teoria in molti campi. In fisica e in ingegneria sono stati elaborati modelli della propagazione di onde sismiche, sull'area minima delle superfici, le oscillazioni non lineari, la dispersione e l'elasticità. Michael V. Berry dell'Università di Bristol ha recentemente utilizzato le catastrofi a ombelico per ottenere nuovi risultati nella fisica delle caustiche e nella meccanica dei fluidi, e ha confermato sperimentalmente questi risultati.

Structural Stability and Morphogenesis di Thom, ispirato dal lavoro di D'Arcy Wentworth Thompson e C.H. Waddington, era largamente connesso all'embriologia, ma per ora soltanto pochi biologi hanno tentato un'applicazione sperimentale delle sue ipotesi in laboratorio. Io ho costruito dei modelli a catastrofe del battito cardiaco, della propagazione degli impulsi nervosi e della formazione della gastrula e dei somiti nell'embrione. Esperimenti recenti condotti da J. Cooke del laboratorio del Medical Research Council di Edimburgo sembrano confermare alcune delle mie previsioni.

Gran parte del mio stesso lavoro, tuttavia, si è svolto nel campo delle scienze umane, come si vede anche dai modelli descritti in questo articolo. Un numero sempre maggiore di ricercatori sta sviluppando modelli derivati dalla teoria della catastrofe, e nei prossimi dieci anni ritengo che questi modelli saranno verificati sperimentalmente. Solo allora potremo valutare il valore reale di questo metodo.

Thom ha usato la sua teoria nel tentativo di capire come si forma il linguaggio. Ed è affascinante pensare che lo stesso apparato matematico possa spiegare non solo come il codice genetico controlla lo sviluppo dell'embrione ma anche come la parola scritta faccia schiudere la nostra immaginazione.



inferiore, come indicato nella sezione a destra. Quando la terapia ha un effetto positivo, tuttavia, la paziente compie una graduale transizione dal piano intermedio ai modi di comportamento normali rappresentati dietro la tasca.

NOTE BIOGRAFICHE E BIBLIOGRAFICHE

1 Sui paradossi di Zenone

Autore

MARTIN GARDNER ha creato la rubrica *Mathematical Games* che appare su «Scientific American» dal 1956. Gardner si è laureato all'Università di Chicago nel 1936; prima di prestare servizio nella marina durante la seconda guerra mondiale è stato giornalista per il «Tulsa Tribune» e ha collaborato allo University of Chicago Press Relations Department, per diventare in seguito giornalista *freelance*. Oltre a essere autore di vari libri di giochi matematici, oramai classici, tratti dal materiale presentato nella sua rubrica mensile, Gardner ha pubblicato diversi volumi al di fuori di questo ambito. Tra questi ricordiamo *The Ambidextrous Universe* (1964), *The Annotated Ancient Mariner* (1965) e *The Philosophic Foundations of Physics* di R. Carnap come curatore.

Bibliografia

GRÜNBAUM ADOLF, *Modern Science and Zeno's Paradoxes*, Wesleyan University Press, 1967.

2 I fondamenti della matematica

Autore

W.V. QUINE è Edgar Pierce Professor of Philosophy della Harvard University. Dopo essersi diplomato in matematica nel 1930 all'Oberlin College si laureò due anni dopo a Harvard con una dissertazione di logica sotto la guida di Alfred North Whitehead. Dopo un anno trascorso tra le università di Vienna, Praga e Varsavia fu eletto alla Harvard's Society of Fellows nel 1933. Nel 1936 iniziò a insegnare a Harvard e nel 1948 divenne professore di filosofia. In seguito ha lavorato presso l'Università di Oxford, l'Institute for Advanced Study di Princeton e il Centre for Advanced Study in the Behavioral Sciences di Palo Alto. Tra le sue opere ricordiamo l'articolo *New Foundations for Mathematical Logic* pubblicato nel 1937 nella «Ameri-

can Mathematical Monthly» e i libri *Mathematical Logic* e *Set Theory and its Logic* del 1940 e del 1963.

Bibliografia

BERNAYS PAUL e FRAENKEL A.A., *Axiomatic Set Theory*, North-Holland Publishing Company, 1958.

FREGE G., *Die Grundlagen der Arithmetik: Eine logisch-matematische Untersuchung über den Begriff der Zahl*, Köbner, Breslavia 1884; traduzione italiana in *Aritmetica e Logica*, a cura di Corrado Mangione.

RUSSELL BERTRAND, *Introduction to Mathematical Philosophy*, George Allen & Unwin, Ltd., 1919.

QUINE WILLARD VAN ORMAN, *Set Theory and Its Logic*, Harvard University Press, 1963.

HATCHER WILLIAM S., *Fondamenti della matematica*, Boringhieri, Torino 1973.

3 La teoria non cantoriana degli insiemi

Autori

PAUL J. COHEN e REUBEN HERSH sono professori di matematica rispettivamente alla Stanford University e all'Università del Nuovo Messico. Cohen ha studiato al Brooklyn College, si è laureato all'Università di Chicago nel 1958 e ha insegnato all'Università di Rochester, al Massachusetts Institute of Technology, alla Harvard University e ha tenuto dei corsi all'Institute for Advanced Study di Princeton. Quale riconoscimento dei suoi studi, e in particolare di quello qui pubblicato, ha ricevuto la Fields Medal al Congresso internazionale di matematica del 1966; assieme ad altri ha vinto il premio di 10 000 dollari che la Research Corporation of America assegna ogni anno ad uno scienziato americano; nel 1963 aveva anche vinto il premio Böcher assegnato annualmente dalla American Mathematical Society per eminenti ricerche di analisi. Hersh ha insegnato all'Università di New York, alla Fairleigh Dickinson University e a Stanford. Prima di laurearsi in matematica a New York nel 1962 aveva conseguito il diploma *magna cum laude* in letteratura inglese a Harvard, era stato militare in Corea, aveva lavorato per quattro anni nella redazione di *Scientific American* e per altri

quattro anni come tornitore e collaudatore di macchine utensili.

Bibliografia

NAGEL E., JAMES R., NEWMAN R., *La prova di Gödel*, Boringhieri, 1961.

CASARI E., *Questioni di filosofia della matematica*, Feltrinelli, 1964.

EVES H. e NEWSON C.V., *An Introduction to the Foundations and Fundamental Concepts of Mathematics*, Holt, Rinehart and Winston, 1965.

SCOTT D., *Proof of the Independence of the Continuum Hypothesis*, in «Mathematical Systems Theory», vol. 1, n. 2, maggio 1967.

Che cos'è il problema del continuo di Cantor? di K. Gödel in *La filosofia della matematica*, a cura di Carlo Cellucci, Laterza, 1967.

COHEN P.J., *Set Theory and the Continuum Hypothesis*, W.A. Benjamin Inc., 1966 traduzione italiana *La teoria degli insiemi e l'ipotesi del continuo*, Feltrinelli, Milano 1973.

LOLLI GABRIELE, *Teoria assiomatica degli insiemi*, Boringhieri, Torino 1974.

4 L'analisi non-standard

Autori

MARTIN DAVIS e REUBEN HERSH sono rispettivamente professore di matematica all'Istituto Courant di scienze matematiche dell'Università di New York e professore di matematica all'Università del Nuovo Messico. Davis si è laureato alla Princeton University nel 1950. Hersh, dopo essersi diplomato in letteratura inglese all'Harvard College, ha lavorato per quattro anni nella redazione di Scientific American. Successivamente cominciò a interessarsi di matematica e si è laureato in questa materia all'Università di New York.

Bibliografia

BOYER CARL B., *The History of the Calculus and Its Conceptual Development*, Dover Publications, 1959.

ROBINSON ABRAHAM, *Non-Standard Analysis*, North-Holland Publishing Co., 1966.

LOLLI GABRIELE, *Nuovi modelli del sistema dei numeri reali*, in questo volume a pagina 42.

LIGHTSTONE A.H., *Infinitesimals* in «The American Mathematical Monthly», 79, n. 3, 1972.

ROBINSON A., *Teoria dei modelli e metamatematica dell'algebra*, Boringhieri, Torino 1974.

5 Nuovi modelli del sistema dei numeri reali

Autore

GABRIELE LOLLI è professore straordinario di logica presso l'Università di Genova. Laureato in matematica nel 1965 all'Università di Torino, si è specializzato in ricerche di logica matematica e di teoria degli insiemi. Nel 1970-71 è stato *visiting fellow* presso il Dipartimento di matematica della Yale University. Tra le sue ultime pubblicazioni ricordiamo *Categorie, universi e principi di riflessione*, Boringhieri, Torino 1977.

Bibliografia

FEFERMAN SOLOMON, *The Number Systems*, Addison-Wesley, 1964.

CHOQUET, G., *La retta numerica in Strutture algebriche e strutture topologiche*, Feltrinelli, 1963.

HALMOS, P.R., *Measure Theory*, Van Nostrand, 1950.

RASIOWA H., SIKORSKI R., *The Mathematics of Metamathematics*, P.W.N. Warszawa, 1963.

SCOTT D., *A Proof of the Independence of the Continuum Hypothesis*, in «Mathematical System Theory», n. 1, 1967.

LOLLI GABRIELE, *Teoria assiomatica degli insiemi*, Boringhieri, Torino 1974.

6 Tre personaggi della matematica

Autore

BRUNO DE FINETTI è professore di calcolo delle probabilità nella Facoltà di scienze dell'Università di Roma. Appena laureatosi in matematica alla Università di Milano entrò all'Istituto centrale di statistica, e poi in una compagnia di assicurazione (con mansioni statistico-attuariali e di razionalizzazione e meccanizzazione di procedure tecnico-contabili), dedicandosi contemporaneamente, nel tempo libero, a ricerche matematiche, prevalentemente nel campo della teoria delle probabilità. Lasciò tali attività nel 1946 per entrare nella carriera universitaria (a Trieste e poi a Roma), dapprima alla cattedra di matematica attuariale. Si interessa fattivamente anche ai problemi dell'insegnamento della matematica (di cui propugna un profondo rinnovamento), in particolare operando in seno alla Commissione italiana per l'insegnamento matematico e alla società *Mathesis*, di cui dal 1970 è presidente. Fa anche parte del Comitato tecnico per la programmazione del Ministero della pubblica istruzione che ha pubblicato le «Proposte per il nuovo piano della scuola».

Bibliografia

KLEIN FELIX, *Elementarmathematik vom Höheren Standpunkte aus*, (3 voll.), Springer, Berlino, 1924.

TRICOMI FRANCESCO, *Funzioni analitiche*, Zanichelli, Bologna, 1946.

DE FINETTI BRUNO, *Matematica logico-intuitiva*, Cremonese, Roma, 1960.

DE FINETTI BRUNO, *Il « saper vedere » in matematica*, Loescher, Torino, 1967.

DE FINETTI BRUNO, *Le proposte per la matematica nei nuovi licei: informazioni, commenti critici, suggerimenti*, in «Periodico di Matematiche», 45, 2, 1967.

BRUNER JEROME S., *Il conoscere: Saggi per la mano sinistra*, Armando, Roma, 1968.

POLYA GEORGE, *Come risolvere i problemi di matematica; La scoperta matematica (I e II)*, Feltrinelli, Milano, 1967, 1971-1974.

CAMPEDELLI LUIGI, *La geometria dei parallelogrammi*, Le Monnier, Firenze, 1970.

DIEUDONNÉ JEAN, *Algebra lineare e geometria elementare*, Feltrinelli, Milano, 1970.

7 Algebre di Boole, diagrammi di Venn e calcolo proposizionale

Autore

MARTIN GARDNER, per la biografia si rimanda all'articolo precedente.

Bibliografia

WHITESITT ELTON J., *Boolean Algebra and Its Applications*, Addison Wesley Publishing Company, 1961.

ARNOLD B.H., *Logic and Boolean Algebra*, Prentice-Hall, Inc. 1962.

MENDELSON ELLIOTT, *Algebra di Boole*, Etas Libri, Milano, 1977.

8 Verità è dimostrazione

Autore

ALFRED TARSKI è professore di matematica all'Università di California a Berkeley. Nato in Polonia, si è laureato alla Università di Varsavia nel 1924 ove è restato fino al 1939, allorché si è trasferito negli USA. E' stato professore associato alla Harvard University per due anni e quindi membro dell'Institute for Advanced Study per un anno, prima di passare a Berkeley. Tra i suoi vari scritti, citiamo *Introduction to Logic and to the Methodology of Deductive Sciences* che è stato pubblicato in 10 lingue (traduzione italiana, *Introduzione alla logica*, Bompiani, Milano 1969). L'articolo qui pubblicato è basato su una conferenza tenuta a Berkeley nel 1963 e su un discorso tenuto all'Università di Londra nel 1966.

Bibliografia

TARSKI ALFRED, *The Semantic Conception of Truth in Semantics and the Philosophy of Language*, a cura di Leonard Linsky, University of Illinois Press, 1952.

TARSKI ALFRED, *The Concept of Truth in Formalized Languages in Logic, Semantics, Metamathematics*, Clarendon Press, 1956.

TARSKI ALFRED, *Introduction to Logic and to the Methodology of Deductive Science*, Oxford University Press, 1965 (traduzione italiana, *Introduzione alla logica*, Bompiani, Milano, 1969).

DALLA CHIARA SCABIA MARIA LUISA, *Modelli sintattici e semantici delle teorie elementari*, Feltrinelli, Milano 1969.

MENDELSON ELLIOTT, *Introduzione alla logica matematica*, Boringhieri, Torino 1972.

DALLA CHIARA SCABIA MARIA LUISA, *La logica*, ISEDI, Milano 1973.

9 Un libro di logica smarrito di Lewis Carroll

Autore

W.W. BARTLEY III è professore di filosofia e di storia della filosofia e della scienza all'Università di Pittsburg, dove inoltre lavora come Direttore aggregato del Centro per la filosofia della scienza. E' anche docente di filosofia all'Università statale della California ad Hayward. Bartley si è diplomato ad Harvard nel 1956 e si è laureato in logica e metodo scientifico alla London School of Economics nel 1962.

Bibliografia

CARROLL LEWIS, *Il gioco della logica*, Astrolabio, 1969.

CARROLL LEWIS, *Una storia ingarbugliata*, Astrolabio, 1969.

MARTIN ROBERT L. (a cura di), *The Paradox of the Liar*, Yale University Press, 1970.

COPI IRVING M. e GOULD JAMES A. (a cura di), *Readings on Logic*, The Macmillan Company, 1972.

MANGIONE CORRADO, *La svolta della logica nell'Ottocento; Logica e problemi dei fondamenti nella seconda metà dell'Ottocento; La logica nel XX secolo in Storia del pensiero filosofico e scientifico* di Ludovico Geymonat (vol. 5 e 6), Garzanti, 1970-72.

10 Un nuovo livello di astrazione: la teoria delle categorie

Autore

LUCIO LOMBARDO RADICE è professore ordinario di algebra all'Università di Roma. Nato a Catania si è laureato in matematica nel 1938. Fu costretto a interrompere l'attività accademica perché incarcerato per attività antifascista. Nel 1945 riprese il suo posto di assistente all'Università di Roma dove restò fino al 1956, quando divenne professore di geometria all'Università di Palermo. Nel 1960 ritornò a Roma come professore prima di geometria e poi di algebra quando questa cattedra fu istituita nell'università italiana. Oltre che di vari scritti e memorie scientifiche, è autore del volume *Istituzioni di algebra astratta*.

Bibliografia

FREYD P.M., *Abelian Categories*, Columbia University Press, 1962.

MITCHELL B., *Theory of Categories*, Academic Press, New York, 1965.

COHN P.M., *Universal Algebra*, Harper and Row, New York, 1965.

HASSE M. e MICHLER L., *Theorie der Kategorien*, VEB Deutscher Verlag der Wissenschaften, 1965.

BRINKMANN H.B. e PUPPE D., *Kategorien und Funktoren*, Springer, Berlino, 1966.

MAC LANE S. e BIRKHOFF G., *Algebra*, McMillan, Londra, 1967, traduzione italiana *Algebra*, Mursia, Milano 1975.

BALDASSARRI GHEZZO S., MARGAGLIO C. e MILLEVOI T., *Introduzione ai metodi della geometria algebrica*, Cremonese, Roma, 1967.

BUCUR I. e DELEANU A., *Introduction to the Theory of Categories and Functors*, Interscience Publications New York, 1968.

PAIREIGIS BODO, *Categories and Functors*, Academic Press, New York, 1970.

MAC LANE SAUNDERS, *Categories for the Working Mathematician*, Springer, New York, Heidelberg, Berlino, 1971; traduzione italiana *Categorie nella pratica matematica*, Boringhieri, Torino 1977.

HIGGINS PHILIP J., *Notes on Categories and Groupoids*, van Nostrand, Londra, 1971.

HORST SCHUBERT, *Categories*, Springer, Berlino, 1972.

11 Induzione e probabilità

Autori

DOMENICO COSTANTINI e MARCO MONDADORI si occupano ormai da anni di logica induttiva. Costantini si è laureato in scienze statistiche all'Università di Roma e quindi specializzato in filosofia della scienza all'Università di Milano; attualmente è docente di calcolo delle probabilità all'Università di Bologna. Mondadori si è laureato in filosofia all'Università di Milano e attualmente insegna all'Università di Bologna.

Bibliografia

CARNAP R. e JEFFREY R.C., *Studies in Inductive logic and Probability*, University of California Press, 1971.

CARNAP R., *Analiticità, Significanza, Induzione* a cura di A. Meotti e M. Mondadori, Il Mulino, 1971.

COSTANTINI D., *Fondamenti del calcolo delle probabilità*, Feltrinelli, 1970.

ESSLER W.K., *Induktive Logik, Grundlagen und Voraussetzung*, Verlag K. Alber, 1970.

COSTANTINI D., *Introduzione alla probabilità*, Boringhieri, Torino 1977.

12 I problemi della conferma

Autore

WESLEY C. SALMON dopo essere stato professore di filosofia della scienza all'Università dell'Indiana è passato all'Università dell'Arizona come professore di filosofia. Si è diplomato all'Università di Chicago e si è laureato in filosofia all'Università della California a Los Angeles. Ha lasciato recentemente la carica di presidente dell'Associazione americana per la filosofia della scienza.

Bibliografia

HUME DAVID, *Ricerche sull'intelletto umano*, Laterza, 1968.

CARNAP RUDOLF, *Logical Foundations of Probability*, The University of Chicago Press, 1962.

SKYRMS BRIAN, *Choice and Chance: An Introduction to Inductive Logic*, Dickenson Publishing Co., 1966.

SALMON WESLEY C., *The Foundations of Scientific Inference*, University of Pittsburgh Press, 1967.

LUCKENBACH SIDNEY A. (a cura di), *Probabilities, Problems and Paradoxes: Readings in Inductive Logic*, Dickenson Publishing Co., 1971.

SALMON WESLEY C., *Confirmation and Relevance in «Minnesota Studies in the Philosophy of Science»*, 1973.

CARNAP RUDOLF, *Analiticità, significanza, induzione*, Il Mulino, 1971.

13 Problemi non risolti dell'aritmetica

Autore

HOWARD DELONG è professore incaricato di filosofia al Trinity College di Hartford. All'università si è specializzato in matematica, laureandosi al Williams College nel 1957. Si è poi dedicato alla filosofia, laureandosi nel 1960 presso la Princeton University. Insegna al Trinity College dal 1960. «Tra i miei interessi esterni – scrive – c'è lo studio del mercato economico da un punto di vista logico. Per me presenta un fascino intellettuale del tutto indipendente da quello monetario corrente».

Bibliografia

MYHILL JOHN, *Some Philosophical Implications of Mathematical Logic, I: Three classes of Ideas*, in «The Review of Metaphysics», 6, n. 2, 1952.

KEMENY JOHN G., *Undecidable Problems of Elementary Number Theory*, in «Mathematische Annalen», 135, n. 2, 1958.

REID CONSTANCE, *Hilbert*, Springer-Verlag, 1970.

DELONG HOWARD, *A profile of Mathematical Logic*, Addison-Wesley, 1970.

14 La macchina di Turing e la questione da essa sollevata: può una macchina pensare?

Autore

MARTIN GARDNER, per la biografia si veda l'articolo precedente.

Bibliografia

BERNSTEIN JEREMY, *The Analytical Engine: Computers-Past, Present and Future*, Random House, 1964.

ROSS ANDERSON ALAN (a cura) *Minds and Machines*, Prentice-Hall, Inc. 1964.

LOEHLIN JOHN C., *Computer Models of Personality*, Randon House, 1968.

MINSKY MARVIN (a cura), *Semantic Information Processing*, The MIT Press, 1968.

MELTZER BERNARD e MICHIE DONALD, *Machine Intelligence 4*, Edinburgh University Press, 1969.

15 Giochi, logica e calcolatori

Autore

HAO WANG (*Giochi, logica e calcolatori*) insegna logica matematica alla Harvard University. Wang è nato in Cina e si è diplomato in matematica alla National Associated University di Kunming nel 1943. Giunto negli Stati Uniti nel 1946 si laureò in logica a Harvard, rimanendo presso questa università per diversi anni come ricercatore e professore assistente. Nel 1954 andò in Inghilterra come ricercatore della Rockefeller Foundation. Nel 1955 fu John Locke Lecturer in Philosophy all'Università di Oxford e lettore di filosofia della matematica dal 1956 al 1961, anno in cui assunse la sua carica attuale a Harvard. Wang è autore di *A Survey of Mathematical Logic*.

Bibliografia

WANG HAO, *Dominoes and the AEA Case of the Decision Problem* in «Proceedings of the Symposium on Mathematical Theory of Automata: MRI Symposia Series, vol. XII», Polytechnic Press of the Polytechnic Institute of Brooklyn, 1963.

HILBERT D. e ACKERMANN W., *Principles of Mathematical Logic*, Chelsea Publishing Company, 1950.

The Undecidable: Basic Papers on Undecidable Propositions, Unsolvability Problems and Computable Functions, Raven Press, 1965.

HERMES HANS, *Enumerabilità, decidibilità, computabilità*, Boringhieri, Torino 1975.

16 Il decimo problema di Hilbert

Autori

MARTIN DAVIS e REUBEN HERSH per la biografia si rimanda all'articolo precedente.

Bibliografia

DAVIS MARTIN, *Computability & Unsolvability*, McGraw-Hill Book Company, 1958.

WANG HAO, *Giochi, logica e calcolatori*, in questo volume a pagina 130.

REID CONSTANCE, *Hilbert*, Springer-Verlag, 1970.

DAVIS MARTIN, *Hilbert's Tenth Problem is Unsolvability* in «American Mathematical Monthly», 80, n. 3, 1973.

17 Gli algoritmi

Autore

DONALD E. KNUTH insegna scienza dei calcolatori alla Stanford University. Diplomatosi in matematica al Case Institute of Technology si è laureato in matematica al California Institute of Technology nel 1963, dove è rimasto per cinque anni dopo la laurea diventando professore di matematica. Nel 1968 è entrato a far parte del Dipartimento

di scienza dei calcolatori di Stanford. Knuth ha ricevuto diversi riconoscimenti per le sue ricerche, compreso lo A.M. Turing Award dell'Association for Computing Machinery nel 1974. Il Turing Award faceva menzione della serie di sette volumi che sta scrivendo intitolata *The art of Computer Programming*, di cui conta di completare i primi cinque entro il 1980.

Bibliografia

KNUTH DONALD E., *Ancient Babylonian Algorithms* in «Communications of the ACM», 15, n. 7, luglio 1972.

KNUTH DONALD E., *The Art of Computer Programming: Vols. I-III*, Addison-Wesley Publishing Company 1973.

KNUTH DONALD E., *Computer Science and Its Relation to Mathematics* in «The American Mathematical Monthly», 81, n. 4, aprile 1974.

AMBLE O. e KNUTH D.E., *Ordered Hash Tables* in «The Computer Journal», 17, maggio 1974.

NIEVERGELT J., *Binary Search Trees and File Organization* in «Computing Surveys», 6, n. 3, settembre 1974.

YAO ANDREW C. e YAO FRANCES F., *The Complexity of Searching an Ordered Random Table* in «Symposium on Foundations of Computer Science», IEEE Computer Society's Technical Committee, IEEE Computer Society, 1976.

18 La teoria delle catastrofi

Autore

E.C. ZEEMAN insegna matematica all'Università di Warwick dove ha fondato e dirige il Mathematics Research Centre. Zeeman ha iniziato i suoi studi di matematica su-

periore nel 1947, dopo quattro anni di servizio come ufficiale d'aviazione nella Royal Air Force. Dopo essersi diplomato nel 1948 e laureato nel 1954 all'Università di Cambridge, ha insegnato in questa università per 10 anni prima di passare a Warwick nel 1964. Fino al 1970 i suoi interessi erano soprattutto per la topologia e in questo campo ha prodotto una cinquantina di articoli. Nel frattempo, intorno al 1960, ha scritto tre articoli sul cervello che, stando alle sue parole, non erano nulla di eccezionale ma hanno avuto il merito di indirizzare verso la biologia gli interessi di René Thom, contribuendo quindi in modo indiretto all'elaborazione da parte di quest'ultimo della teoria delle catastrofi. Negli ultimi cinque anni gli interessi di Zeeman sono stati concentrati soprattutto sullo sviluppo della parte matematica della teoria delle catastrofi e sulle applicazioni di tale teoria a diversi settori della scienza.

Bibliografia

ZEEMAN E.C., *Primary and Secondary Waves in Development Biology* in «Lectures on Mathematics in the Life Sciences», 7: Some Mathematical Questions in Biology, a cura di Simon A. Levin, American Mathematical Society, 1974.

THOM R. e ZEEMAN E.C., *Catastrophe Theory; Its Present State and Future Perspectives* in «Dynamical Systems-Warwick 1974: Proceedings of a Symposium Held at the University of Warwick 1973/74» a cura di Anthony Manning, Springer-Verlag, 1975.

THOM RENÉ, *Structural Stability and Morphogenesis*, W.A. Benjamin, Inc., 1975.

TROTMAN D.J.A. e ZEEMAN E.C., *The Classification of Elementary Catastrophes of Codimension 5* in «Symposium on Catastrophe Theory; Seattle 1976» Springer-Verlag (in stampa).

INDICE ANALITICO

Abel N. E., 90
 algebra di Boole, 47, 48, 49, 66-69
 algoritmo, 118, 121, 139, 130-137, 144, 147-157
 algoritmo di Euclide, 141
 analisi non standard, 34, 39, 41
 anello, 89-94
 antinomia del mentitore, 73, 78, 87, 118
 antinomie, 26
 Archimede, 34, 35, 89
 archimedeica proprietà, 34, 40
 Aristotele, 34, 37, 70, 96, 111, 114
 assioma di continuità, 45, 48
 assioma di scelta, 27-30, 32, 33
 assiomi, 27

Bacone F., 96
 Baker A., 141
 Berkeley J., 37, 38, 41
 Bernoulli Jacopo e Giovanni, 37, 38, 98, 101
 Berry M., 171
 Beth E.W., 84
 Bolyai J., 30, 33
 Boole G., 49, 84, 89
 Brouwer L.E.J., 29
 Bruner J.S., 50

calcolo proposizionale, 66-69
 campo, 91
 Cantor G., 25, 26, 28, 115
 cardinalità, 25, 47
 Carnap R., 102, 112
 Carroll Lewis (Charles Lutwidge Dogson), 81, 86-88
 categoria, 89, 93-95
 Cauchy A. L., 17, 43, 44
 Cavendish H., 104
 Church A., 118, 120-123, 139
 classe di equivalenza, 44
 Cohen P., 29, 30, 42, 45
 completezza, 123
 completo (sistema formale), 118
 computabile, 142
 conforme (proiezione), 58, 59
 connettivi logici, 47, 77
 consistenza, 29, 30, 115
 consistenza relativa, 32

continuo (insieme), 26
 convergente, 43
 convergere, 43
 coordinate polari, 57, 59, 60, 61
 corpo commutativo, 45
 corpo di insiemi, 46
 corpo ordinato, 43, 45, 47
 corpo ordinato continuo, 44, 45, 48
 corrispondenza biunivoca, 25, 27, 47, 89
 costrizione (*forcing*), 33
 costruibili (insiemi), 30
 crivello di Eratostene, 145
 curva esponenziale, 52, 58, 60
 Cusano N., 36, 41

D'Alembert J., 41
 D'Arcy Wentworth Th., 171
 Davis M., 139, 145, 146
 Dedekind R., 44
 De Finetti R., 99, 100
 de l'Hôpital, 37, 41
 Democrito, 34
 diagramma commutativo, 92
 diagrammi di Venn, 66-69
 dimostrazione, 75-79
 Diofanto di Alessandria, 139, 141
 Duhem P., 109, 111

Eilenberg S., 89
 Einstein A., 33, 110, 158
 elemento neutro, 45, 47, 49, 90, 93
 equazione diofantea, 139, 141, 144-146
 equipotenza, 25, 26
 Euatlo, 113
 Eubulide, 73
 Euclide, 30, 34, 41, 76, 89, 111, 114, 117, 121, 122
 Eudosso, 42
 Eulero L., 37, 38, 59

Fermat P., 117, 120, 141
 fluenti, 37
 flussioni, 37, 38
 Fraenkel A., 28
 Frege G., 20, 26, 76, 111
 frequenza relativa, 99
 funtore, 93
 funzione derivata, 58, 59
 funzione esponenziale, 53, 55, 59, 62



funzione inversa, 53, 54
funzione logaritmo, 53, 55

Gauss K. F., 30, 33, 115, 123, 141
geometria euclidea 27-30, 32, 115
geometria non euclidea, 28-30, 32, 42, 115
geometria riemanniana, 33
Gödel K., 24, 28, 30, 32, 33, 40, 78, 111, 118-123, 139
Goldbach C., 117, 119, 120, 123
Goodman N., 105
gruppo, 90, 92, 93
gruppo commutativo (abeliano), 90, 91, 93, 94

Hafele J.C., 109
Halley E., 38
Hamilton W.R., 43
hashing, 153, 156
Hempel C.G., 105
Henkin L., 17, 40
Hilbert D., 26, 28, 100, 138, 139, 141, 142, 145, 146

implicazione logica, 106
indivisibili, 34
induzione, 96
inferenza logica, 27
infinitesimo, 34, 35, 37-41
infinito (insieme), 25
iniezione, 91
insieme, 93
insieme potenza, 26, 28
insiemi di Borel, 46, 48
inverso, 45, 93
ipotesi del continuo, 28-30, 32, 33, 49
isomorfo (modello), 44, 45, 47, 49
Keynes M., 99
Keplero J., 36
Kleene S.C., 139
Kolmogorov A.N., 100, 101
Kotarbinski T., 73
Laplace P.S., 96, 98, 99
Leibniz G.W., 17, 34, 37, 38, 41, 96, 158

Leonardo Pisano (il Fibonacci), 144, 146
limitato superiormente (insieme), 44, 45, 48
limite, 43
linguaggio formale, 39, 74, 77
linguaggio oggetto, 75, 77, 78, 118, 119
linguaggio simbolico, 98
Lobacevskij N.I., 30, 33, 115, 123
Löwenheim L., 83
logaritmi decimali, 54
logaritmi naturali, 53, 54
logica, 114-117
logica del primo ordine, 45
logica del secondo ordine, 45
Lorenz K., 158
Lukasiewicz J., 73

macchina di Turing, 126-129, 132, 136, 137
Mac Lane S., 89

Malcev A., 40
Matyasevich Y., 138, 139, 141, 144-146
Maxwell J.C., 158
metalinguaggio, 75, 77, 78, 118-119
metodo epsilon-delta, 36, 41
metodo di esaustione, 36, 41
Mill. J.S., 96
minimo confine superiore, 44, 45
modello, 28, 29
modello non standard, 40, 41, 46, 49
monoide, 93
morfismo, 92, 93-95
Myhill J., 123

Napier J., 53
Newton Isaac, 17, 34, 37, 38, 104, 158
non archimedeo (numero), 34
nucleo, 91, 92, 95
numerabile (insieme), 25, 26
numeri complessi, 54, 56
numeri di Fibonacci, 144, 146
numeri interi relativi, 42
numeri irrazionali, 42, 43
numeri naturali, 25, 40, 42, 89
numeri razionali, 42, 43
numeri reali, 40, 42, 44

oggetto finale, 94
oggetto iniziale, 94
omomorfismo, 90-94
omotetia, 55, 56, 57

paradossi di Zenone 13-15 87, 116
paradossi logici, 105
paradosso dei corvi, 105
paradosso del barbiere, 84, 111
paradosso della conferma, 102, 103
paradosso del verdlù-blerde, 105
paradosso di Russell, 23, 26, 73, 116
Pascal B., 36
piano affine, 57, 58
Pitagora, 39, 43
Platone, 89
Playfair J., 30, 32
Pioncaré H., 29, 138
Post E., 121, 139
postulati, 27
postulato della parallela, 17, 30
probabilità, 42, 96, 98, 99, 100, 103, 106-112
procedura di decisione, 123, 147
programma, 147
proposizioni elementari, 47
Protagora, 113
pseudosfera, 29, 32
punto singolare 61, 62
Putnam H., 139, 145, 146

quantificatori, 44, 77
quoziente (struttura), 43, 44

radiante, 59
Ramsey F.P. 99, 100
rappresentazione su (suriezione), 91
regole di inferenza, 76, 82
relazione di equivalenza, 44
Riemann B., 29, 30, 33
Robinson A., 34, 37, 40, 41

Robinson J., 139, 145, 146
Ross A., 88
Russell B., 23, 26, 81, 105, 111, 116

Savage L.J., 112
Scott D., 48
serie di Leibniz, 60
serie di potenze, 53, 54
sezioni di Dedekind, 44
sillogismo, 114
sistema formale, 117, 118
Shannon C., 54
Skolem T., 40

spazio di probabilità, 46, 47, 49
spazio topologico, 93
stereografica (proiezione), 58, 59
successione fondamentale (o di Cauchy), 43, 44, 49
sviluppo in serie, 60, 61

Talete, 27, 111
Tarski A., 87
teorema cinese del resto, 145
teorema di Bayes, 112
teorema di Cauchy-Hadamard, 61
teorema di compattezza, 40
teorema di completezza, 40
teorema di incompletezza (di Gödel), 78, 111, 118, 119
teoremi, 87

teoria degli insiemi, 23, 25-33, 42, 45, 46, 115
teoria dei modelli, 41
teoria delle catastrofi, 158-171
teoria formale, 77
teoria non cantoriana degli insiemi, 30
teoria non standard, 42

Thom R. 158, 165, 168, 170, 171
totalmente ordinato (insieme), 48
Turing A.M., 121, 139

unità immaginaria, 54
universo non standard, 39
universo standard, 39

variabili casuali, 46-49
velocità istantanea, 37
verità, 70-79
von Mises R., 98-100

Waddington C.H., 171
Weierstrass K., 34, 36-38, 41
Weyl H., 29
Whitehead A.N., 81

Yao F.F., 151
Yao A.C., 151

Zenone, 27
Zermelo E., 27, 28

VERITA' E DIMOSTRAZIONE

Questioni di matematica

Gli articoli presentati in questo volume, alcuni dei quali inediti in Italia, offrono una panoramica dello stato della ricerca contemporanea in matematica, sia per quanto riguarda gli studi fondazionali dalla portata più nettamente filosofica, sia per quanto riguarda il lato applicativo, legato alla soluzione di problemi concreti. Caratteristica di questi articoli è il fatto che gli autori in gran parte coincidono con i matematici e i filosofi che hanno dato contributi fondamentali alla ricerca in questo secolo; basti pensare a Quine, a Tarski e a Cohen. Nella prima parte del volume è possibile trovare una chiara esposizione del problema dei fondamenti, e a questo riguardo sono accostati articoli su argomenti specifici, come l'indipendenza dell'ipotesi del continuo o l'analisi non-standard, ad articoli di carattere più generale sull'importanza fondazionale della teoria degli insiemi e sui fondamenti dell'analisi. Nella seconda parte è l'aspetto logico matematico che prende il sopravvento. Il lettore è introdotto ai risultati fondamentali riguardanti la decidibilità e la completezza delle teorie e agli strumenti essenziali per poter comprendere, almeno nelle linee generali, i procedimenti con cui tali risultati sono stati ottenuti, primo fra tutti il concetto di algoritmo o di procedura effettiva, codificato nelle macchine di Turing, argomento con cui si apre la terza parte del volume. A questo punto il problema dei fondamenti della matematica viene espresso in una delle sue formulazioni più mature, come ricerca di una adeguatezza tra semantica, ossia teoria degli insiemi, e sintassi, ossia teorie formali, di una corrispondenza fra verità e dimostrazione. Il fatto che tale corrispondenza sia solo parziale (è questa una delle conseguenze del famoso teorema di Gödel) costituisce uno dei risultati più interessanti tra quelli conseguiti dalla ricerca nel campo della logica matematica:

Articoli

- M. GARDNER
 SUI PARADOSSI DI ZENONE
- W.V. QUINE
 I FONDAMENTI DELLA MATEMATICA
- P.J. COHEN e R. HERSCH
 LA TEORIA NON CANTORIANA DEGLI INSIEMI
- M. DAVIS e R. HERSCH
 L'ANALISI NON-STANDARD
- G. LOLLI
 NUOVI MODELLI DEL SISTEMA DEI NUMERI REALI
- B. DE FINETTI
 TRE PERSONAGGI DELLA MATEMATICA: π , e , i
- M. GARDNER
 ALGEBRA DI BOOLE, DIAGRAMMI DI VENN
 E CALCOLO PROPOSIZIONALE
- A. TARSKI
 VERITÀ E DIMOSTRAZIONE
- W.W. BARTLEY III
 UN LIBRO DI LOGICA SMARRITO DI LEWIS CARROLL
- L. LOMBARDO-RADICE
 UN NUOVO LIVELLO DI ASTRAZIONE:
 LA TEORIA DELLE CATEGORIE
- D. COSTANTINI e M. MONDADORI
 INDUZIONE E PROBABILITÀ
- W.C. SALMON
 I PROBLEMI DELLA CONFERMA
- H. DELONG
 I PROBLEMI NON RISOLTI DELL'ARITMETICA
- M. GARDNER
 LA MACCHINA DI TURING E LA QUESTIONE DA ESSA
 SOLLEVATA: PUÒ UNA MACCHINA PENSARE?
- H. WANG
 GIOCHI, LOGICA E CALCOLATORI
- M. DAVIS e R. HERSCH
 IL X PROBLEMA DI HILBERT
- D.E. KNUTH
 GLI ALGORITMI
- E.C. ZEEMAN
 LA TEORIA DELLE CATASTROFI