

EPILOGUE

The Art and Science of Cause and Effect

*A public lecture delivered November 1996 as part of
the UCLA Faculty Research Lectureship Program*

The topic of this lecture is causality – namely, our awareness of what causes what in the world and why it matters.

Though it is basic to human thought, causality is a notion shrouded in mystery, controversy, and caution, because scientists and philosophers have had difficulties defining when one event *truly causes* another.

We all understand that the rooster’s crow does not cause the sun to rise, but even this simple fact cannot easily be translated into a mathematical equation.

Today, I would like to share with you a set of ideas which I have found very useful in studying phenomena of this kind. These ideas have led to practical tools that I hope you will find useful on your next encounter with cause and effect.

It is hard to imagine anyone here who is *not* dealing with cause and effect.

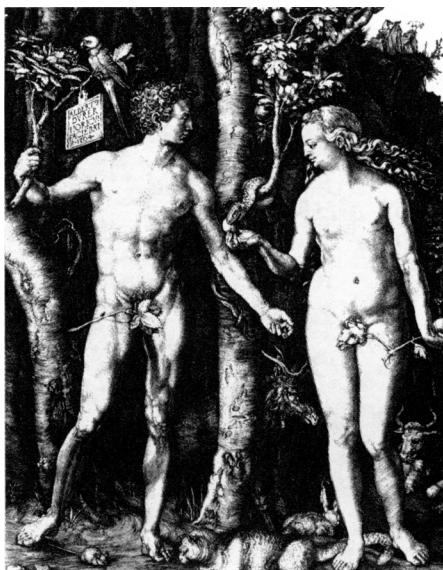
Whether you are evaluating the impact of bilingual education programs or running an experiment on how mice distinguish food from danger or speculating about why Julius Caesar crossed the Rubicon or diagnosing a patient or predicting who will win the presidential election, you are dealing with a tangled web of cause–effect considerations.

The story that I am about to tell is aimed at helping researchers deal with the complexities of such considerations, and to clarify their meaning.

This lecture is divided into three parts.

I begin with a brief historical sketch of the difficulties that various disciplines have had with causation.

Next I outline the ideas that reduce or eliminate several of these historical difficulties.





Eve is just as skillful: “The serpent deceived me, and I ate.”

The thing to notice about this story is that God did not ask for explanation, only for the facts – it was Adam who felt the need to explain. The message is clear: causal explanation is a man-made concept.

Another interesting point about the story: explanations are used exclusively for passing responsibilities.

Indeed, for thousands of years explanations had no other function. Therefore, only Gods, people, and animals could cause things to happen, not objects, events, or physical processes.

Natural events entered into causal explanations much later because, in the ancient world, events were simply *predetermined*.



Storms and earthquakes were *controlled* by the angry gods [slide 2] and could not in themselves assume causal responsibility for the consequences.

Even an erratic and unpredictable event such as the roll of a die [3] was not considered a *chance* event but rather a divine message demanding proper interpretation.

One such message gave the prophet Jonah the scare of his life when he was identified as God’s renegade and was thrown overboard [4].

Quoting from the book of Jonah: “And the sailors said: ‘Come and let us cast lots to find out who is to blame for this ordeal.’ So they cast lots and the lot fell on Jonah.”

Obviously, on this luxury Phoenician cruiser, “casting lots” was used not for recreation but for communication – a one-way modem for processing messages of vital importance.

In summary, the agents of causal forces in the ancient world were either deities, who cause things to happen for a purpose, or human beings and animals, who possess free will, for which they are punished and rewarded.

This notion of causation was naive, but clear and unproblematic.

The problems began, as usual, with engineering; when machines had to be constructed to do useful jobs [5].

As engineers grew ambitious, they decided that the earth, too, can be moved [6], but not with a single lever.

Systems consisting of many pulleys and wheels [7], one driving another, were needed for projects of such magnitude.

And, once people started building multistage systems, an interesting thing happened to causality – *physical objects began acquiring causal character*.

When a system like that broke down, it was futile to blame God or the operator – instead, a broken rope or a rusty pulley were more useful explanations, simply because these could be replaced easily and make the system work.

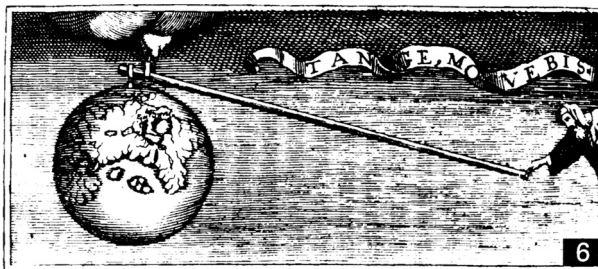
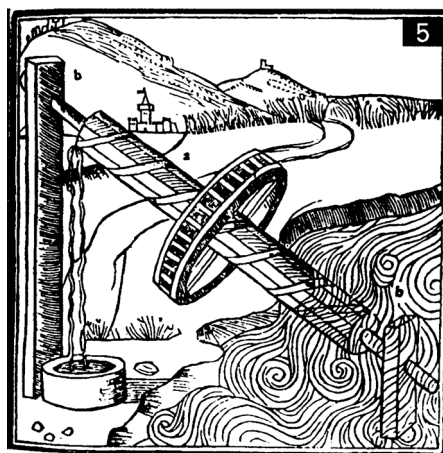
At that point in history, Gods and humans ceased to be the sole agents of causal forces – lifeless objects and processes became partners in responsibility.

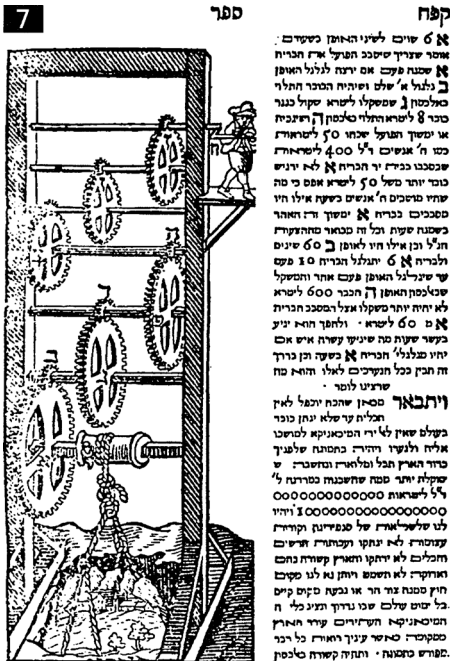
A wheel turned and stopped *because* the wheel preceding it turned and stopped – the human operator became secondary.

Not surprisingly, these new agents of causation took on some of the characteristics of their predecessors – Gods and humans.

Natural objects became not only carriers of credit and blame but also carriers of force, will, and even purpose.

Aristotle regarded explanation in terms of a *purpose* to be the only complete and satisfactory explanation for why a thing is what it is.





He even called it a *final cause* – namely, the final aim of scientific inquiry.

From that point on, causality served a dual role: *causes* were the targets of credit and blame on one hand and the carriers of physical flow of control on the other.

This duality survived in relative tranquility [8] until about the time of the Renaissance, when it encountered conceptual difficulties.

What happened can be seen on the title page [9] of Recorde's book "The Castle of Knowledge," the first science book in English, published in 1575.

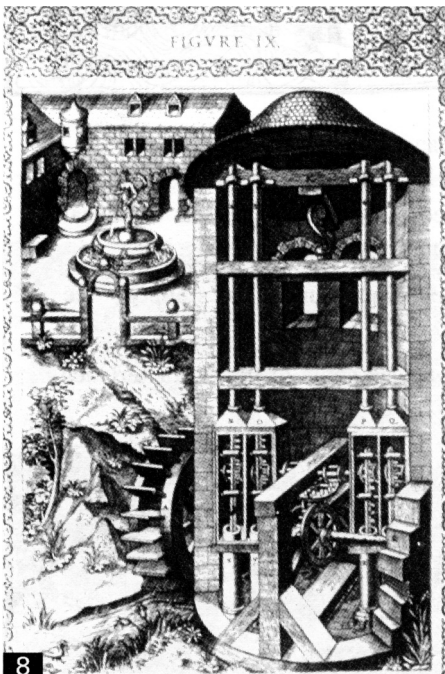
The wheel of fortune is turned, not by the wisdom of God, but by the ignorance of man.

And, as God's role as the final cause was taken over by human knowledge, the whole notion of causal explanation came under attack.

The erosion started with the work of Galileo [10].

Most of us know Galileo as the man who was brought before the inquisition and imprisoned [11] for defending the heliocentric theory of the world.

But while all that was going on, Galileo also managed to quietly engineer the most profound revolution that science has ever known.



This revolution, expounded in his 1638 book "Discorsi" [12], published in Leyden, far from Rome, consists of two maxims:

One, description first, explanation second – that is, the "how" precedes the "why"; and

Two, description is carried out in the language of mathematics; namely, equations.

Ask not, said Galileo, whether an object falls because it is pulled from below or pushed from above.

Ask how well you can predict the time it takes for the object to travel a certain distance, and how that time will vary from object to object and as the angle of the track changes.

Moreover, said Galileo, do not attempt to answer such questions in the qualitative and slippery nuances of human language; say it in the form of mathematical equations [13].

It is hard for us to appreciate today how strange that idea sounded in 1638, barely 50 years after the introduction of algebraic notation by Vieta. To proclaim algebra the *universal* language of science would sound today like proclaiming Esperanto the language of economics.

Why would Nature agree to speak algebra? Of all languages?

But you can't argue with success.

The distance traveled by an object turned out indeed to be proportional to the square of the time.

Even more successful than predicting outcomes of experiments were the computational aspects of algebraic equations.

They enabled engineers, for the first time in history, to ask “how to” questions in addition to “what if” questions.

In addition to asking: “What if we narrow the beam, will it carry the load?”, they began to ask more difficult questions: “How to shape the beam so that it *will* carry the load?” [14]

This was made possible by the availability of methods for solving equations.

The algebraic machinery does not discriminate among variables; instead of predicting behavior in terms of parameters, we can turn things around and solve for the parameters in terms of the desired behavior.

Let us concentrate now on Galileo's first maxim – “description first, explanation second” – because that idea was taken very seriously by the scientists and changed the character of science from speculative to empirical.

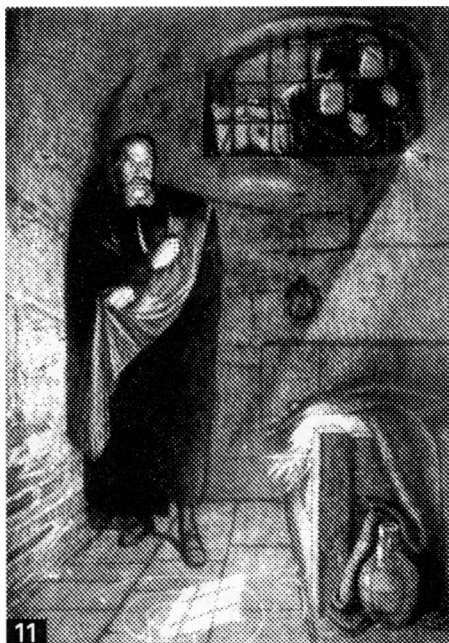
Physics became flooded with empirical laws that were extremely useful.

Snell's law [15], Hooke's law, Ohm's law, and Joule's law are examples of purely empirical generalizations that were discovered and used long before they were explained by more fundamental principles.

Philosophers, however, were reluctant to give up the idea of causal explanation and continued to search for the origin and justification of those successful Galilean equations.

For example, Descartes ascribed cause to *eternal truth*.





cause and the other *effect*, and infer the existence of the one from that of the other.”

Thus, causal connections according to Hume are the product of observations. Causation is a learnable habit of the mind, almost as fictional as optical illusions and as transitory as Pavlov’s conditioning.

**DISCORSI
E
DIMOSTRAZIONI
MATEMATICHE,
intorno à due nuoue scienze**

Attenenti alla
MECANICA & i MOVIMENTI LOCALI;

del Signor
GALILEO GALILEI LINCEO,
Filosofo e Matematico primario del Serenissimo
Grand Duca di Toscana.

Con una Appendice del centro di gravità à d'alcuni Solidi.



IN LEIDA,
Appresso gli Elſevirii. M. D. C. XXXVIII. **12**

Liebniz made cause a self-evident logical law.

Finally, about one hundred years after Galileo, a Scottish philosopher by the name of David Hume [16] carried Galileo’s first maxim to an extreme [17].

Hume argued convincingly that the *why* is not merely second to the *how*, but that the *why* is totally superfluous as it is subsumed by the *how*.

On page 156 of Hume’s “Treatise of Human Nature” [18], we find the paragraph that shook up causation so thoroughly that it has not recovered to this day.

I always get a kick reading it: “Thus we remember to have seen that species of object we call *flame*, and to have felt that species of sensation we call *heat*. We likewise call to mind their constant conjunction in all past instances. Without any farther ceremony, we call the one

It is hard to believe that Hume was not aware of the difficulties inherent in his proposed recipe.

He knew quite well that the rooster crow *stands* in constant conjunction to the sunrise, yet it does not *cause* the sun to rise.

He knew that the barometer reading *stands* in constant conjunction to the rain but does not *cause* the rain.

Today these difficulties fall under the rubric of *spurious correlations*, namely “correlations that do not imply causation.”

Now, taking Hume’s dictum that all knowledge comes from experience encoded in the mind as correlation, and our observation that correlation does not imply causation, we are led into our first riddle of causation: How do people *ever* acquire knowledge of *causation*?

We saw in the rooster example that regularity of succession is not sufficient; what *would* be sufficient?

What patterns of experience would justify calling a connection “causal”?

Moreover: What patterns of experience *convince* people that a connection is “causal”?

If the first riddle concerns the *learning* of causal connection, the second concerns its usage: What *difference* would it make if I told you that a certain connection is or is not causal?

Continuing our example, what difference would it make if I told you that the rooster does cause the sun to rise?

This may sound trivial.

The obvious answer is that knowing “what causes what” makes a big difference in how we act.

If the rooster’s crow causes the sun to rise, we could make the night shorter by waking up our rooster earlier and making him crow – say, by telling him the latest rooster joke.

But this riddle is *not* as trivial as it seems.

If causal information has an empirical meaning beyond regularity of succession, then that information should show up in the laws of physics.

But it does not!

The philosopher Bertrand Russell made this argument [19] in 1913:

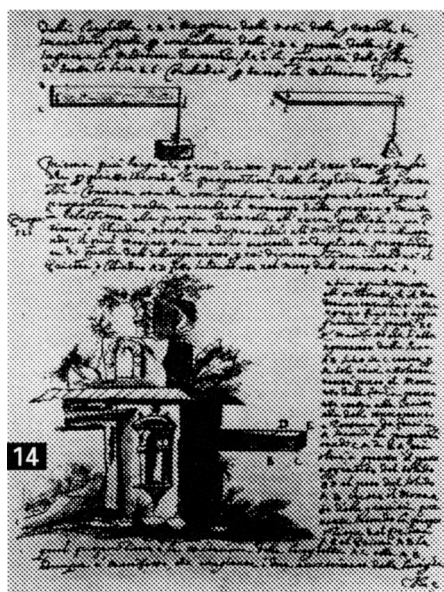
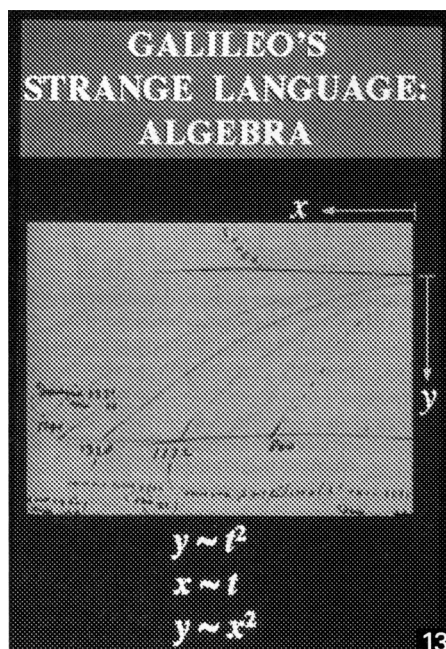
“All philosophers,” says Russell, “imagine that causation is one of the fundamental axioms of science, yet oddly enough, in advanced sciences, the word ‘cause’ never occurs.... The law of causality, I believe, is a relic of bygone age, surviving, like the monarchy, only because it is erroneously supposed to do no harm.”

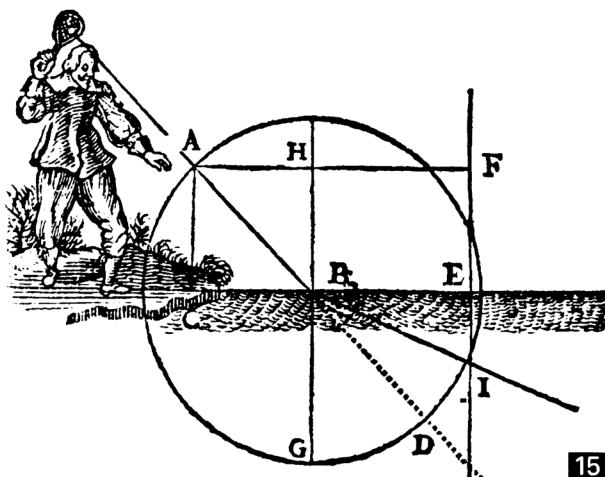
Another philosopher, Patrick Suppes, who argued for the importance of causality, noted that:

“There is scarcely an issue of ‘Physical Review’ that does not contain at least one article using either ‘cause’ or ‘causality’ in its title.”

What we conclude from this exchange is that physicists talk, write, and think one way and formulate physics in another.

Such bilingual activity would be forgiven if causality was used merely as a convenient communication device – a shorthand for expressing complex patterns of physical relationships that would otherwise take many equations to write.





Take, for instance, Newton's law:

$$f = ma.$$

The rules of algebra permit us to write this law in a wild variety of syntactic forms, all meaning the same thing – that if we know any two of the three quantities, the third is determined.

Yet, in ordinary discourse we say that force causes acceleration – not that acceleration causes force, and we feel very strongly about this distinction.

Likewise, we say that the ratio f/a helps us *determine* the mass, not that it *causes* the mass.

Such distinctions are not supported by the equations of physics, and this leads us to ask whether the whole causal vocabulary is purely metaphysical, “surviving, like the monarchy ...”.

Fortunately, very few physicists paid attention to Russell's enigma. They continued to write equations in the office and talk cause–effect in the cafeteria; with astonishing success they smashed the atom, invented the transistor and the laser.

The same is true for engineering.

But in another arena the tension could not go unnoticed, because in that arena the demand for distinguishing causal from other relationships was very explicit.

This arena is statistics.

The story begins with the discovery of correlation, about one hundred years ago.



Francis Galton [20], inventor of fingerprinting and cousin of Charles Darwin, quite understandably set out to prove that talent and virtue run in families.

Galton's investigations drove him to consider various ways of measuring how properties of one class of individuals or objects are related to those of another class.

In 1888, he measured the length of a person's forearm and the size of that person's head and asked to what degree one of these quantities can predict the other [21].

He stumbled upon the following discovery: If you plot one quantity against the other and scale the two axes properly, then the slope of the best-fit line has some nice mathematical properties. The slope is 1 only when one quantity can predict the other precisely; it is zero whenever the prediction is no better than a random guess; and, most remarkably, the slope is the same no matter if you plot X against Y or Y against X .

"It is easy to see," said Galton, "that co-relation must be the consequence of the variations of the two organs being partly due to common causes."

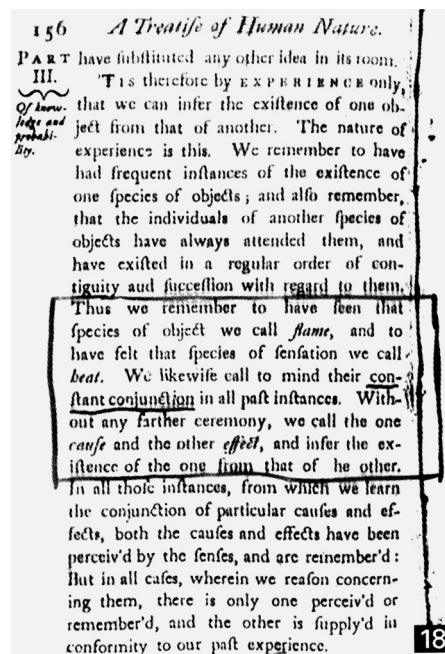
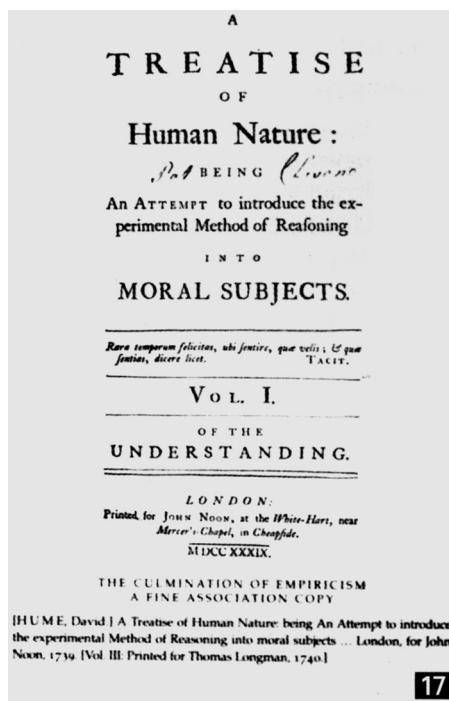
Here we have, for the first time, an objective measure of how two variables are "related" to each other, based strictly on the data, clear of human judgment or opinion.

Galton's discovery dazzled one of his disciples, Karl Pearson [22], now considered to be one of the founders of modern statistics.

Pearson was 30 years old at the time, an accomplished physicist and philosopher about to turn lawyer, and this is how he describes, 45 years later [23], his initial reaction to Galton's discovery:

"I felt like a buccaneer of Drake's days. . . .

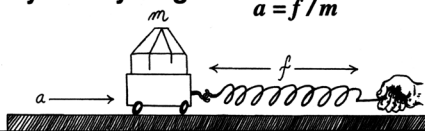
"I interpreted . . . Galton to mean that there was a category broader than causation, namely correlation, of which causation was only the limit, and that this new conception of correlation brought psychology, anthropology, medicine, and sociology in large parts into the field of mathematical treatment."



PURGING CAUSALITY FROM PHYSICS?

- **BERTRAND RUSSELL (1913):**
In advanced sciences the word “cause” never occurs. Causality is a relic of bygone ago.
- **PATRICK SUPPES (1970):**
“Causality” is commonly used by physicists

The symmetry enigma: $f = m a$
 $a = f / m$



Now, Pearson has been described as a man “with the kind of drive and determination that took Hannibal over the Alps and Marco Polo to China.”

When Pearson felt like a buccaneer, you can be sure he gets his bounty.

The year 1911 saw the publication of the third edition of his book “The Grammar of Science.” It contained a new chapter titled “Contingency and Correlation – The Insufficiency of Causation,” and this is

what Pearson says in that chapter:

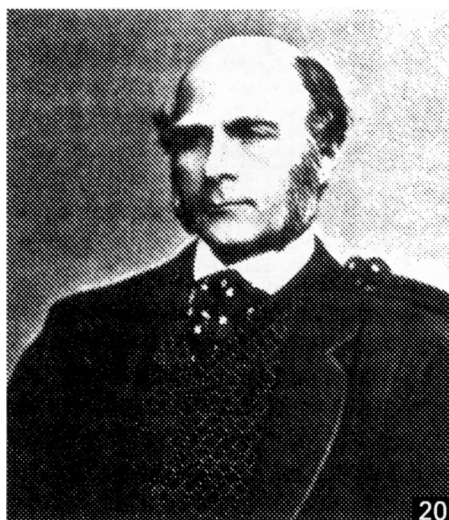
“Beyond such discarded fundamentals as ‘matter’ and ‘force’ lies still another fetish amidst the inscrutable arcana of modern science, namely, the category of cause and effect.”

And what does Pearson substitute for the archaic category of cause and effect? You wouldn’t believe your ears: *contingency tables* [24].

“Such a table is termed a contingency table, and the ultimate scientific statement of description of the relation between two things can always be thrown back upon such a contingency table. . . .

“Once the reader realizes the nature of such a table, he will have grasped the essence of the conception of association between cause and effect.”

Thus, Pearson categorically denies the need for an independent concept of causal relation beyond correlation.



He held this view throughout his life and, accordingly, did not mention causation in *any* of his technical papers.

His crusade against animistic concepts such as “will” and “force” was so fierce and his rejection of determinism so absolute that he *exterminated* causation from statistics before it had a chance to take root.

It took another 25 years and another strong-willed person, Sir Ronald Fisher [25], for statisticians to formulate the randomized experiment – the only scientifically proven method of testing causal relations from data, and to this day, the one and only causal concept permitted in mainstream statistics.

And that is roughly where things stand today.

If we count the number of doctoral theses, research papers, or textbook pages written on causation, we get the impression that Pearson still rules statistics.

The “Encyclopedia of Statistical Science” devotes twelve pages to correlation but only two pages to causation – and spends one of those pages demonstrating that “correlation does not imply causation.”

Let us hear what modern statisticians say about causality.

Philip Dawid, the current editor of “Biometrika” (the journal founded by Pearson), admits: “Causal inference is one of the most important, most subtle, and most neglected of all the problems of statistics.”

Terry Speed, former president of the Biometric Society (whom you might remember as an expert witness at the O. J. Simpson murder trial), declares: “Considerations of causality should be treated as they have always been treated in statistics: preferably not at all but, if necessary, then with very great care.”

Sir David Cox and Nanny Wermuth, in a book published just a few months ago, apologize as follows: “We did not in this book use the words *causal* or *causality*. . . . Our reason for caution is that it is rare that firm conclusions about causality can be drawn from one study.”

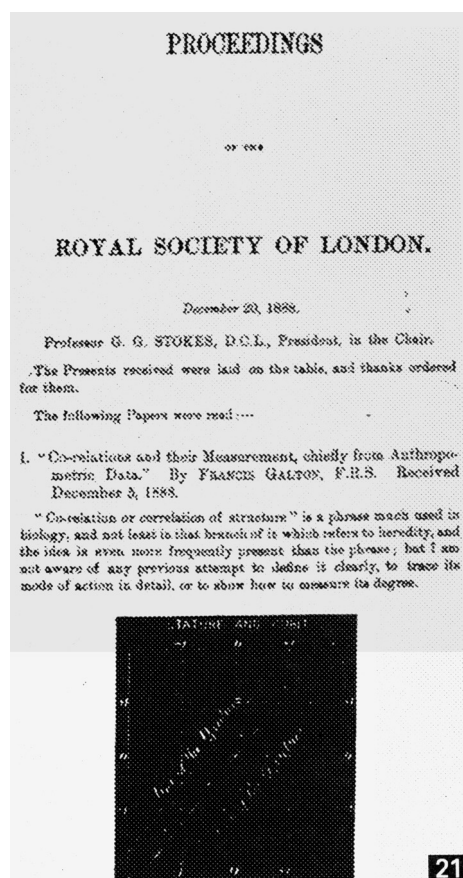
This position of caution and avoidance has paralyzed many fields that look to statistics for guidance, especially economics and social science.

A leading social scientist stated in 1987: “It would be very healthy if more researchers abandon thinking of and using terms such as cause and effect.”

Can this state of affairs be the work of just one person? Even a buccaneer like Pearson?

I doubt it.

But how else can we explain why statistics, the field that has given the world such powerful concepts as the testing of hypothesis and the





design of experiment, would give up so early on causation?

One obvious explanation is, of course, that causation is much harder to measure than correlation.

Correlations can be estimated directly in a single uncontrolled study, while causal conclusions require controlled experiments.

But this is too simplistic; statisticians are not easily deterred by difficulties, and children manage to learn cause–effect relations *without* running controlled experiments.

The answer, I believe, lies deeper, and it has to do with the official language of statistics – namely, the language of probability.

This may come as a surprise to some of you, but the word *cause* is not in the vocabulary of probability theory; we cannot express in the language of probabilities the sentence, *mud does not cause*

rain – all we can say is that the two are mutually correlated or dependent – meaning that if we find one, we can expect the other.

Naturally, if we lack a language to express a certain concept explicitly, we can't expect to develop scientific activity around that concept.

Scientific development requires that knowledge be transferred reliably from one study to another and, as Galileo showed 350 years ago, such transference requires the precision and computational benefits of a formal language.

I will soon come back to discuss the importance of language and notation, but first I

wish to conclude this historical survey with a tale from another field in which causation has had its share of difficulty.

This time it is computer science – the science of symbols – a field that is relatively new yet one that has placed a tremendous emphasis on language and notation and therefore may offer a useful perspective on the problem.

When researchers began to encode causal relationships using computers, the two riddles of causation were awakened with renewed vigor.

24 CONTINGENCY AND CORRELATION 159

B_1 occurs n_{p1} , B_2 occurs n_{p2} times, and so on. We thus are able to obtain a general distribution of B's for each class of A that we can form, and were we to go through the whole population, N, of A's in this manner we should obtain a table of the following kind :—

		TYPE OF A OBSERVED										Total.
TYPE OF B OBSERVED		A_1	A_2	A_3	A_p	
	B_1	n_{11}	n_{21}	n_{31}	n_{p1}	$n_{.1}$
	B_2	n_{12}	n_{22}	n_{32}	n_{p2}	$n_{.2}$
	B_3	n_{13}	n_{23}	n_{33}	n_{p3}	$n_{.3}$

	B_q	n_{1q}	n_{2q}	n_{3q}	n_{pq}	$n_{.q}$

	Total	$n_{1.}$	$n_{2.}$	$n_{3.}$	$n_{p.}$	N

Put yourself in the shoes of this robot [26] who is trying to make sense of what is going on in a kitchen or a laboratory.

Conceptually, the robot's problems are the same as those faced by an economist seeking to model the national debt or an epidemiologist attempting to understand the spread of a disease.

Our robot, economist, and epidemiologist all need to track down cause–effect relations from the environment, using limited actions and noisy observations.

This puts them right at Hume's first riddle of causation: *how?*

The second riddle of causation also plays a role in the robot's world.

Assume we wish to take a shortcut and teach our robot all we know about cause and effect in this room [27].

How should the robot organize and make use of this information?

Thus, the two philosophical riddles of causation are now translated into concrete and practical questions:

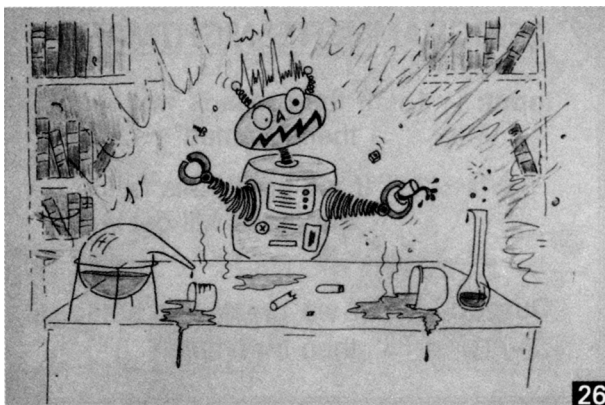
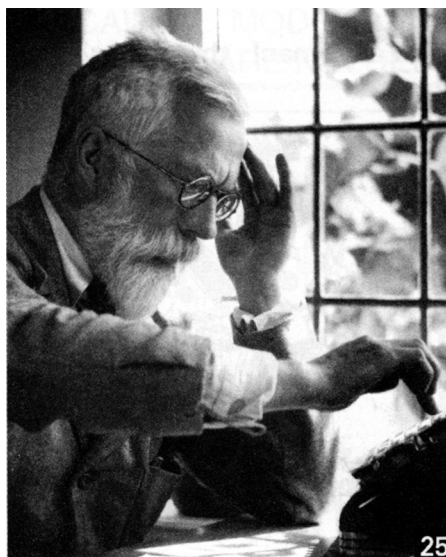
How should a robot acquire causal information through interaction with its environment? How should a robot process causal information received from its creator–programmer?

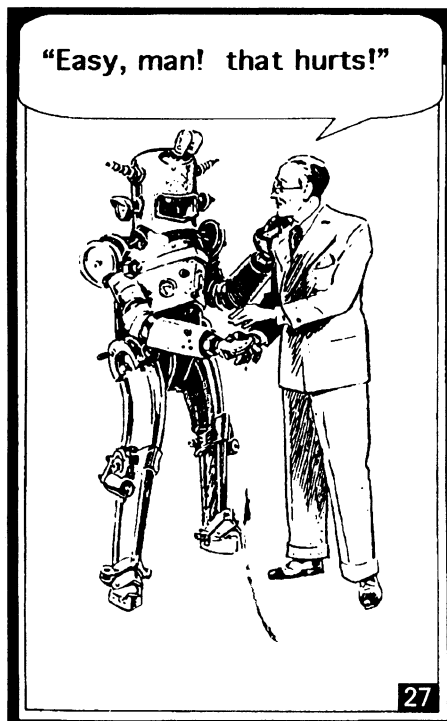
Again, the second riddle is not as trivial as it might seem. Lord Russell's warning that causal relations and physical equations are incompatible now surfaces as an apparent flaw in logic.

For example, when given the information, “If the grass is wet, then it rained” and “If we break this bottle, the grass will get wet,” the computer will conclude “If we break this bottle, then it rained” [28].

The swiftness and specificity with which such programming bugs surface have made Artificial Intelligence programs an ideal laboratory for studying the fine print of causation.

This brings us to the second part of the lecture: how the second riddle of causation can be solved by combining equations with graphs, and how this solution makes the first riddle less formidable.





The overriding ideas in this solution are:

First – treating causation as a summary of behavior under interventions; and

Second – using equations and graphs as a mathematical language within which causal thoughts can be represented and manipulated.

And to put the two together, we need a *third* concept: treating interventions as a surgery over equations.

Let us start with an area that uses causation extensively and never had any trouble with it: engineering.

Here is an engineering drawing [29] of a circuit diagram that shows cause–effect relations among the signals in the circuit. The circuit consists of *and* gates and *or* gates, each performing some logical function between input and output. Let us examine this diagram closely, since its simplicity and familiarity are very deceiving. This diagram is, in fact, one of the greatest marvels of science. It is capable of conveying more information than mil-

lions of algebraic equations or probability functions or logical expressions. What makes this diagram so much more powerful is the ability to predict not merely how the circuit behaves under normal conditions but also how the circuit will behave under millions of *abnormal* conditions. For example, given this circuit diagram, we can easily tell what the output will be if some input changes from 0 to 1. This is normal and can easily be expressed by a simple input–output equation. Now comes the abnormal part. We can also tell what the output will be when we set *Y* to 0 (zero), or tie it to *X*, or change this *and* gate to an *or* gate, or when we perform any of the millions of combinations of these

operations. The designer of this circuit did not anticipate or even consider such weird interventions, yet, miraculously, we can predict their consequences. How? Where does this representational power come from?

It comes from what early economists called *autonomy*. Namely, the gates in this diagram represent independent mechanisms – it is easy to change one without changing the other. The diagram takes advantage of this independence and

CAUSATION AS A PROGRAMMER'S NIGHTMARE

- Input:**
1. “If the grass is wet, then it rained”
 2. “If we break this bottle, the grass will get wet”
- Output:** “If we break this bottle, then it rained”

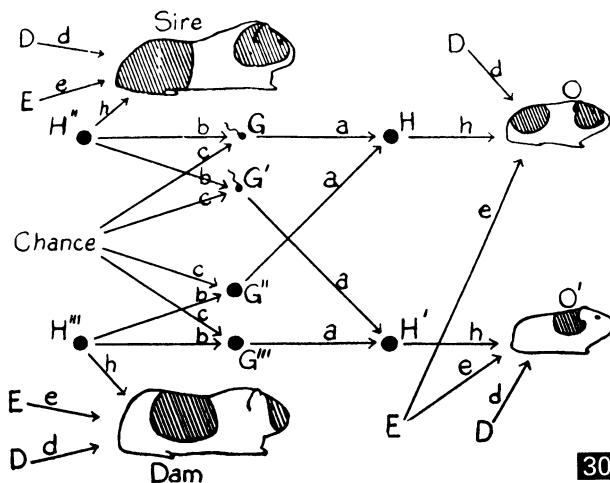
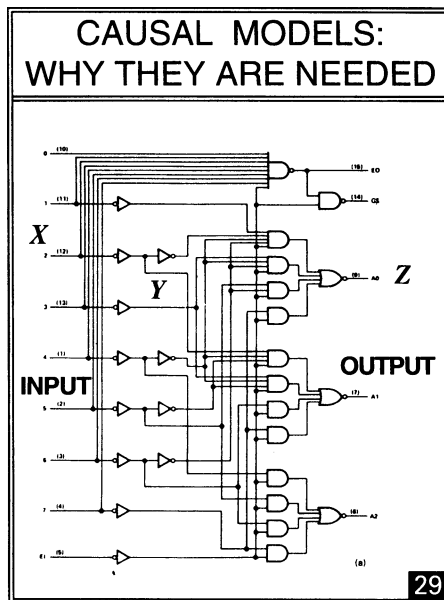
describes the normal functioning of the circuit *using precisely those building blocks that will remain unaltered under intervention*.

My colleagues from Boelter Hall are surely wondering why I stand here before you blathering about an engineering triviality as if it were the eighth wonder of the world. I have three reasons for doing this. First, I will try to show that there is a lot of unexploited wisdom in practices that engineers take for granted.

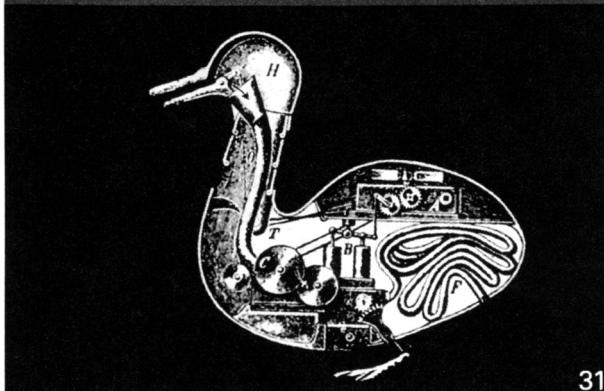
Second, I am trying to remind economists and social scientists of the benefits of this diagrammatic method. They have been using a similar method on and off for over 75 years, called structural equation modeling and path diagrams, but in recent years they have allowed algebraic convenience to suppress the diagrammatic representation, together with its benefits. Finally, these diagrams capture, in my opinion, the very essence of causation – the ability to predict the consequences of abnormal eventualities and new manipulations. In S.Wright's diagram [30], for example, it is possible to predict what coat pattern the guinea-pig litter is likely to have if we change environmental factors, shown here as input (E), or even genetic factors, shown as intermediate nodes between parents and offsprings (H). Such predictions cannot be made on the basis of algebraic or correlational analysis.

Viewing causality this way explains why scientists pursue causal explanations with such zeal and why attaining a causal model is accompanied by a sense of gaining “deep understanding” and “being in control.”

Deep understanding [31] means knowing not merely how things behaved yesterday but also how things will behave under new hypothetical circumstances, control being one such circumstance. Interestingly, when we have such understanding we feel “in control” even if we have no practical way of controlling things. For example, we have no practical way to control celestial motion, and still the theory of gravitation gives us a feeling of understanding and control, because it provides a blueprint for hypothetical control. We can predict the effect on tidal waves of unexpected new events – say, the moon being hit by a meteor or the gravitational constant suddenly diminishing by a



DEEP UNDERSTANDING = HOW THINGS WORK WHEN TAKEN APART



factor of 2 – and, just as important, the gravitational theory gives us the assurance that ordinary manipulation of earthly things will *not* control tidal waves. It is not surprising that causal models are viewed as the litmus test for distinguishing deliberate reasoning from reactive or instinctive response. Birds and monkeys may possibly be trained to perform complex tasks such as fixing a broken wire, but that requires trial-and-error training. Deliberate reasoners, on the other hand, can anticipate the consequences of new manipulations *without ever trying* those manipulations.

Let us magnify [32] a portion of the circuit diagram so that we can understand why the diagram can predict outcomes that equations cannot. Let us also switch from logical gates to linear equations (to make everyone here more comfortable), and assume we are dealing with a system containing just two components: a multiplier and an adder. The *multiplier* takes the input and multiplies it by a factor of 2; the *adder* takes its input and adds a 1 to it. The equations describing these two components are given here on the left.

But are these equations *equivalent* to the diagram on the right? Obviously not! If they were, then let us switch the variables around, and the resulting two equations should be equivalent to the circuit shown below. But these two circuits are different. The top one tells us that if we physically manipulate Y it will affect Z , while the bottom one shows that manipulating Y will affect X and will have no effect on Z . Moreover, performing some additional algebraic operations on our equations, we can obtain two

new equations, shown at the bottom, which point to no structure *at all*; they simply represent two constraints on three variables without telling us how they influence each other.

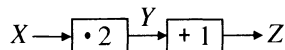
Let us examine more closely the mental process by which we determine the effect of physically manipulating Y – say, setting Y to 0 [33].

Clearly, when we set Y to 0, the relation between X and Y is no longer given by the multiplier – a

EQUATIONS VS. DIAGRAMS

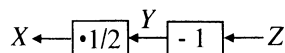
$$Y = 2X$$

$$Z = Y + 1$$



$$X = Y/2$$

$$Y = Z - 1$$

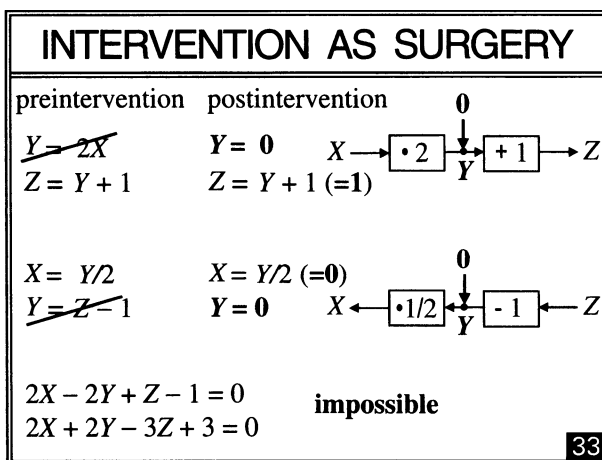


$$2X - 2Y + Z - 1 = 0$$

$$2X + 2Y - 3Z + 3 = 0$$

32

new mechanism now controls Y , in which X has no say. In the equational representation, this amounts to replacing the equation $Y = 2X$ by a new equation $Y = 0$ and solving a new set of equations, which gives $Z = 1$. If we perform this surgery on the lower pair of equations, representing the lower model, we get of course a different solution. The second equation will need to be replaced, which will yield $X = 0$ and leave Z unconstrained.

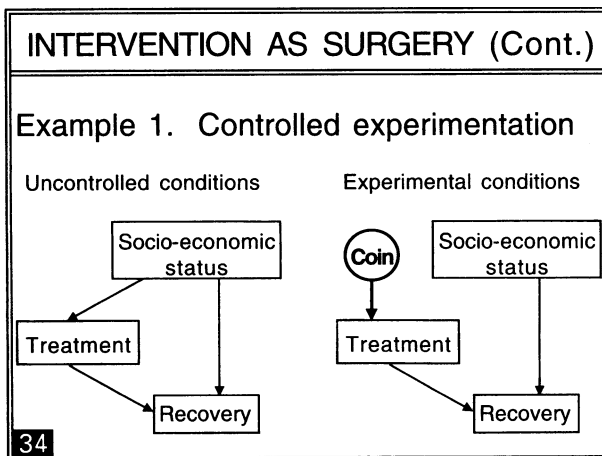


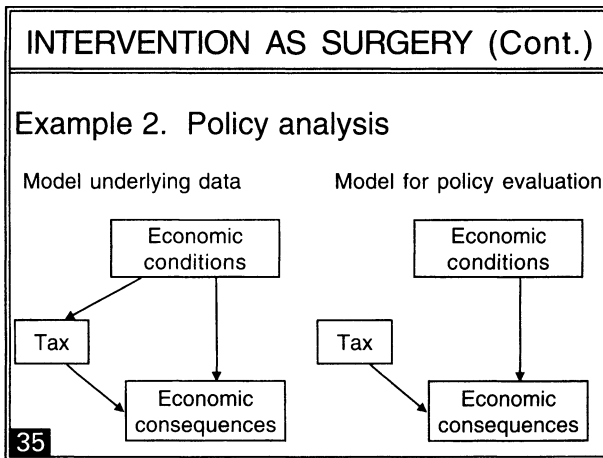
We now see how this model of intervention leads to a formal definition of causation: “ Y is a cause of Z if we can change Z by manipulating Y , namely, if after surgically removing the equation for Y , the solution for Z will depend on the new value we substitute for Y .” We also see how vital the diagram is in this process. *The diagram tells us which equation is to be deleted when we manipulate Y .* That information is totally washed out when we transform the equations into algebraically equivalent form, as shown at the bottom of the screen. From this pair of equations alone, it is impossible to predict the result of setting Y to 0, because we do not know what surgery to perform – there is no such thing as “the equation for Y .”

In summary, *intervention amounts to a surgery on equations* (guided by a diagram) and *causation means predicting the consequences of such a surgery*.

This is a universal theme that goes beyond physical systems. In fact, the idea of modeling interventions by “wiping out” equations was first proposed in 1960 by an *economist*, Herman Wold, but his teachings have all but disappeared from the economics literature. History books attribute this mysterious disappearance to Wold’s personality, but I tend to believe that the reason goes deeper: Early econometricians were very careful mathematicians; they fought hard to keep their algebra clean and formal, and they could not agree to have it contaminated by gimmicks such as diagrams. And as we see on the screen, the surgery operation makes no mathematical sense without the diagram, as it is sensitive to the way we write the equations.

Before expounding on the properties of this new mathematical operation, let me demonstrate how useful it is for clarifying concepts in statistics and economics.





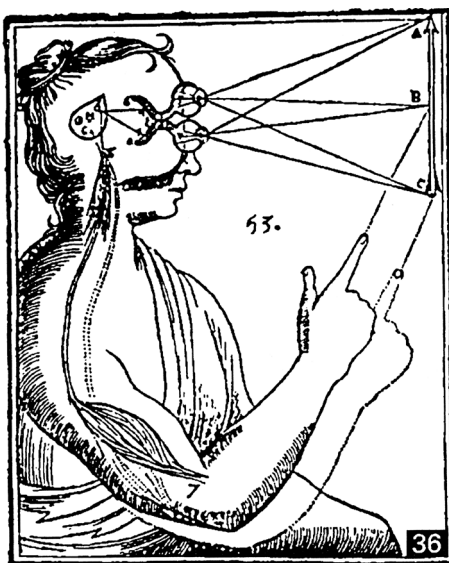
Under uncontrolled conditions, the choice of treatment is up to the patients and may depend on the patients' socioeconomic backgrounds. This creates a problem, because we can't tell if changes in recovery rates are due to treatment or to those background factors. What we wish to do is compare patients of like backgrounds, and that is precisely what Fisher's *randomized experiment* accomplishes. How? It actually consists of two parts, randomization and *intervention*.

Intervention means that we change the natural behavior of the individual: we separate subjects into two groups, called treatment and control, and we convince the subjects to obey the experimental policy. We assign treatment to some patients who, under normal circumstances, will not seek treatment, and we give a placebo to patients who otherwise would receive treatment. That, in our new vocabulary, means *surgery* – we are severing one functional link and replacing it with another. Fisher's great insight was that connecting the new link to a random coin flip *guarantees* that the link we wish to break

is actually broken. The reason is that a random coin is assumed to be unaffected by anything we can measure on a macroscopic level – including, of course, a patient's socioeconomic background.

This picture provides a meaningful and formal rationale for the universally accepted procedure of randomized trials. In contrast, our next example uses the surgery idea to point out inadequacies in widely accepted procedures.

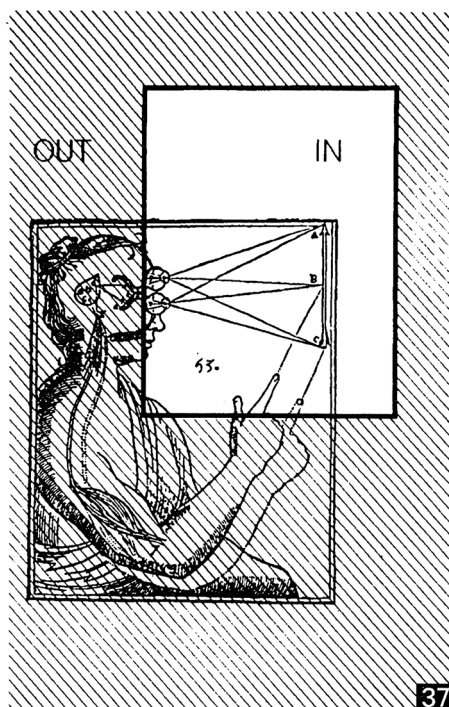
The example [35] involves a government official trying to evaluate the economic consequences of some policy – say, taxation. A deliberate decision to raise or lower taxes is a surgery on the model of the economy because it modifies the conditions prevailing when the model was built. Economic models are built on the basis of data taken over some period of time, and during this period



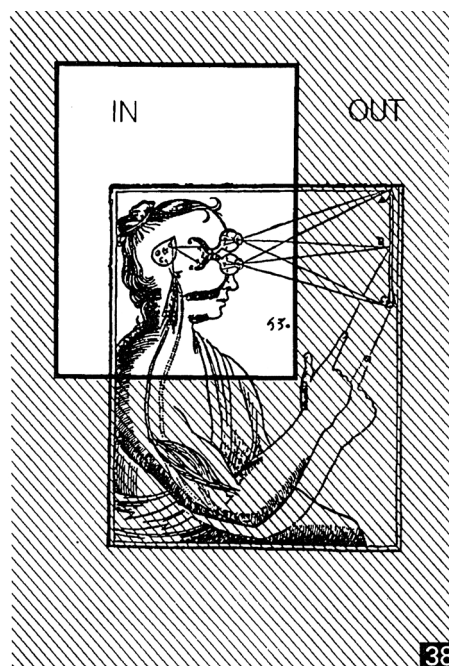
of time taxes were lowered and raised in response to some economic conditions or political pressure. However, when we *evaluate* a policy, we wish to compare alternative policies under the *same* economic conditions – namely, we wish to sever this link that, in the past, has tied policies to those conditions. In this setup, it is of course impossible to connect our policy to a coin toss and run a controlled experiment; we do not have the time for that, and we might ruin the economy before the experiment is over. Nevertheless the analysis that we *should conduct* is to infer the behavior of this mutilated model from data governed by a nonmutilated model.

I said *should conduct* because you will not find such analysis in any economics textbook. As I mentioned earlier, the surgery idea of Herman Wold was stamped out of the economics literature in the 1970s, and all discussions on policy analysis that I could find assume that the mutilated model prevails throughout. That taxation is under government control at the time of evaluation is assumed to be sufficient for treating taxation as an exogenous variable throughout, when in fact taxation is an endogenous variable during the model-building phase and turns exogenous only when evaluated. Of course, I am not claiming that reinstating the surgery model would enable the government to balance its budget overnight, but it is certainly something worth trying.

Let us now examine how the surgery interpretation resolves Russell's enigma concerning the clash between the directionality of causal relations and the symmetry of physical equations. The equations of physics are indeed symmetrical, but when we compare the phrases "*A causes B*" versus "*B causes A*," we are not talking about a single set of equations. Rather, we are comparing two world models, represented by two different sets of equations: one in which the equation for *A* is surgically removed; the other where the equation for *B* is removed. Russell would probably stop us at this point and ask: "How can you talk about *two* world models when in fact there is only one world model, given by all the equations of physics put together?" The answer is: *yes*. If you wish to



37



38

FROM PHYSICS TO CAUSALITY

Physics:

Symmetric equations of motion

Causal models:

Symmetric equations of motion

Circumscription (in vs. out)

Locality (autonomy of mechanisms)

Intervention = surgery on mechanisms

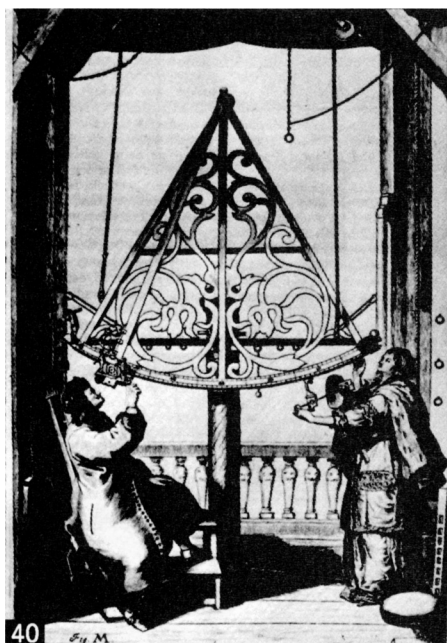
39

include the entire universe in the model, causality disappears because interventions disappear – the manipulator and the manipulated lose their distinction. However, scientists rarely consider the entirety of the universe as an object of investigation. In most cases the scientist carves a piece from the universe and proclaims that piece *in* – namely, the *focus* of investigation. The rest of the universe is then considered *out* or *background* and is summarized by what we call *boundary conditions*. This choice of *ins*

and *outs* creates asymmetry in the way we look at things, and it is this asymmetry that permits us to talk about “outside intervention” and hence about causality and cause–effect directionality.

This can be illustrated quite nicely using Descartes’ classical drawing [36]. As a whole, this hand–eye system knows nothing about causation. It is merely a messy plasma of particles and photons trying their very best to obey Schroedinger’s equation, which is symmetric.

However, carve a chunk from it – say, the object part [37] – and we can talk about the motion of the hand *causing* this light ray to change angle.



Carve it another way, focusing on the brain part [38], and lo and behold it is now the light ray that causes the hand to move – precisely the opposite direction. The lesson is that it is the way we carve up the universe that determines the directionality we associate with cause and effect. Such carving is tacitly assumed in every scientific investigation. In artificial intelligence it was called “circumscription” by J. McCarthy. In economics, circumscription amounts to deciding which variables are deemed endogenous and which exogenous, *in* the model or *external* to the model.

Let us summarize the essential differences between equational and causal models [39]. Both use a set of symmetric equations to describe normal conditions. The causal model, however, contains three additional ingredients: (i) a distinction between the *in* and the *out*; (ii) an assumption that each equation corresponds to an independent mechanism and hence must be preserved as a

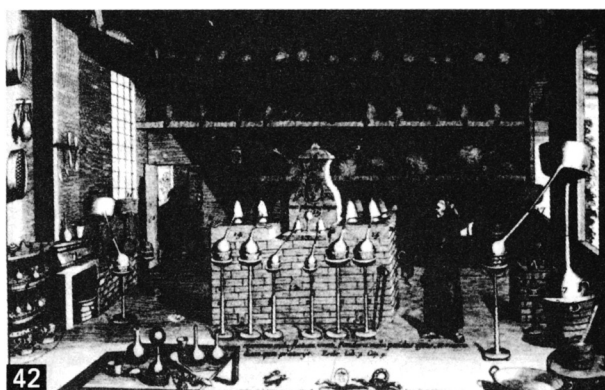
separate mathematical sentence; and (iii) interventions that are interpreted as surgeries over those mechanisms. This brings us closer to realizing the dream of making causality a friendly part of physics. But one ingredient is missing: *the algebra*. We discussed earlier how important the computational facility of algebra was to scientists and engineers in the Galilean era. Can we expect such algebraic facility to serve causality as well? Let me rephrase it differently: Scientific activity, as we know it, consists of two basic components:

Observations [40] and interventions [41].

The combination of the two is what we call a *laboratory* [42], a place where we control some of the conditions and observe others. It so happened that standard algebras have served the observational component very well but thus far have not benefitted the interventional component. This is true for the algebra of equations, Boolean algebra, and probability calculus – all are geared to serve observational sentences but not interventional sentences.

Take, for example, probability theory. If we wish to find the chance that it rained, given that we see the grass wet, we can express our question in a formal sentence written like that: $P(\text{Rain} \mid \text{Wet})$, to be read: the probability of Rain, given Wet [43]. The vertical bar stands for the phrase: “given that we see.” Not only can we express this question in a formal sentence, we can also use the machinery of probability theory and transform the sentence into other expressions. In our example, the sentence on the left can be transformed to the one on the right, if we find it more convenient or informative.

But suppose we ask a different question: “What is the chance it rained if we *make* the grass wet?” We cannot even express our query in the syntax of probability, because the vertical bar is already taken to mean “given that we see.” We can invent a new symbol *do*, and each time we see a *do* after the bar we read it *given that we do* – but this does not help us compute the answer to our question, because the rules of probability do not apply to this new reading. We know intuitively what the answer should be: $P(\text{Rain})$, because making



NEEDED: ALGEBRA OF DOING

Available: algebra of seeing

e.g., What is the chance it rained
if we **see** the grass wet?

$$P(\text{rain} \mid \text{wet}) = ? \quad \{ = P(\text{wet} \mid \text{rain}) \frac{P(\text{rain})}{P(\text{wet})} \}$$

Needed: algebra of doing

e.g., What is the chance it rained
if we **make** the grass wet?

$$P(\text{rain} \mid \text{do}(\text{wet})) = ? \quad \{ = P(\text{rain}) \}$$

43

the grass wet does not change the chance of rain. But can this intuitive answer, and others like it, be derived mechanically, so as to comfort our thoughts when intuition fails?

The answer is *yes*, and it takes a new algebra. First, we assign a symbol to the new operator “given that we do.” Second, we find the rules for manipulating sentences containing this new symbol. We do that by a process analogous to the way mathematicians found the rules of standard algebra.

Imagine that you are a mathematician in the sixteenth century, you are now an expert in the algebra of *addition*, and you feel an urgent need to introduce a new operator, *multiplication*, because you are tired of adding a number to itself all day long [44]. The first thing you do is assign the new operator a symbol: *multiply*. Then you go down to the meaning of the operator, from which you can deduce its rules of transformations. For example: the commutative law of multiplication can be deduced that way, the associative law, and so on. We now learn all this in high school.

In exactly the same fashion, we can deduce the rules that govern our new symbol: *do* (\cdot). We have an algebra for seeing – namely, probability theory. We have a new operator, with a brand new outfit and a very clear meaning, given to us by the surgery procedure. The door is open for deduction, and the result is given in the next slide [45].

Please do not get alarmed, I do not expect you to read these equations right now, but I think you can still get the flavor of this new calculus. It consists of three rules that permit us to transform expressions involving actions and observations into other expressions of this type. The first allows us to ignore an irrelevant observation, the third to

ignore an irrelevant action; the second allows us to exchange an action with an observation of the same fact. What are those symbols on the right? They are the “green lights” that the diagram gives us whenever the transformation is legal. We will see them in action on our next example.

This brings us to part three of the lecture, where I will demonstrate how the ideas presented thus far can be used to solve new problems of practical importance.

NEEDED: ALGEBRA OF DOING (Cont.)

Algebra of Multiplication

Available: algebra of addition

e.g., $a + b = b + a$,
 $a + (b + c) = (a + b) + c$

New operation: $a \times b$

Meaning: add a to itself b times

New rules:

$$\begin{aligned} a \times b &= b \times a, \\ a \times (b \times c) &= (a \times b) \times c \\ a \times (b + c) &= a \times b + a \times c \end{aligned}$$

By Analogy

Available: algebra of seeing

e.g., $P(x \mid y) = \frac{P(x, y)}{P(y)}$

New operation: $\text{do}(z)$

Meaning: surgery + substitution

New rules: $P(x \mid y, \text{do}(z)) = ?$

44

Consider the century-old debate concerning the effect of smoking on lung cancer [46]. In 1964, the Surgeon General issued a report linking cigarette smoking to death, cancer, and most particularly lung cancer. The report was based on nonexperimental studies in which a strong correlation was found between smoking and lung cancer, and the claim was that the correlation found is causal: If we ban smoking, then the rate of cancer cases will be roughly the same as the one we find today among non-smokers in the population.

These studies came under severe attacks from the tobacco industry, backed by some very prominent statisticians, among them Sir Ronald Fisher. The claim was that the observed correlations can also be explained by a model in which there is no causal connection between smoking and lung cancer. Instead, an unobserved genotype might exist that simultaneously causes cancer and produces an inborn craving for nicotine. Formally, this claim would be written in our notation as: $P(\text{Cancer} \mid do(\text{Smoke})) = P(\text{Cancer})$, meaning that making the population smoke or stop smoking would have no effect on the rate of cancer cases. Controlled experiments could decide between the two models, but these are impossible (and now also illegal) to conduct.

This is all history. Now we enter a hypothetical era where representatives of both sides decide to meet and iron out their differences. The tobacco industry concedes that there might be some weak causal link between smoking and cancer and representatives of the health group concede that there might be some weak links to genetic factors. Accordingly, they draw this combined model, and the question boils down to assessing, from the data, the strengths of the various links. They submit the query to a statistician and the answer comes back immediately: *impossible*. Meaning: there is no way to estimate the strength from the data, because any data whatsoever can perfectly fit either one of these two extreme models. So they give up and decide to continue the political battle as usual. Before parting, a suggestion comes up: perhaps we can resolve our differences if we measure some auxiliary factors. For example, since the

RULES OF CAUSAL CALCULUS

Rule 1: Ignoring observations

$$P(y \mid do\{x\}, z, w) = P(y \mid do\{x\}, w) \quad \text{if } (Y \perp\!\!\!\perp Z \mid X, W)_{G_{\overline{X}}}$$

Rule 2: Action/observation exchange

$$P(y \mid do\{x\}, do\{z\}, w) = P(y \mid do\{x\}, z, w) \quad \text{if } (Y \perp\!\!\!\perp Z \mid X, W)_{G_{\overline{X}, \overline{Z}}}$$

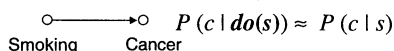
Rule 3: Ignoring actions

$$P(y \mid do\{x\}, do\{z\}, w) = P(y \mid do\{x\}, w) \quad \text{if } (Y \perp\!\!\!\perp Z \mid X, W)_{G_{\overline{X}, \overline{Z}, \overline{W}}}$$

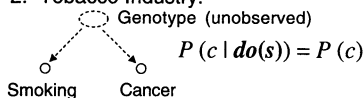
45

SMOKING AND CANCER: HANDLING COMPETING MODELS

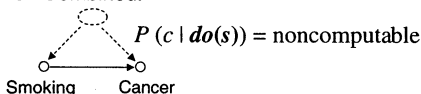
1. Surgeon General (1964):



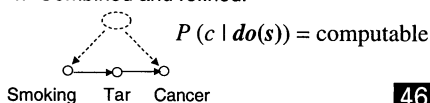
2. Tobacco Industry:



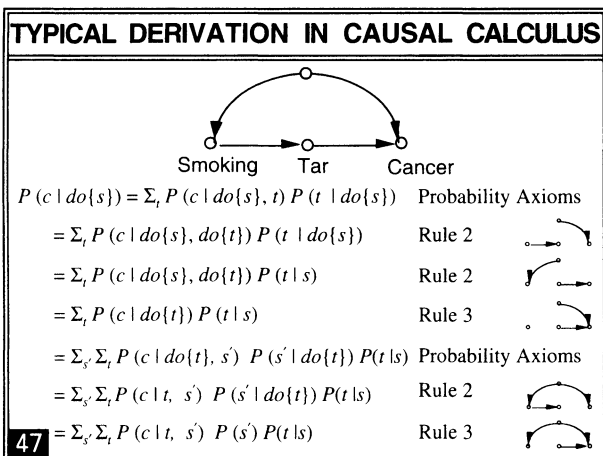
3. Combined:



4. Combined and refined:



46



deposits is available? The statistician comes back with good news: *it is computable* and, moreover, the solution is given in closed mathematical form. *How?*

SIMPSON'S PARADOX

(Pearson et al. 1899; Yule 1903; Simpson 1951)

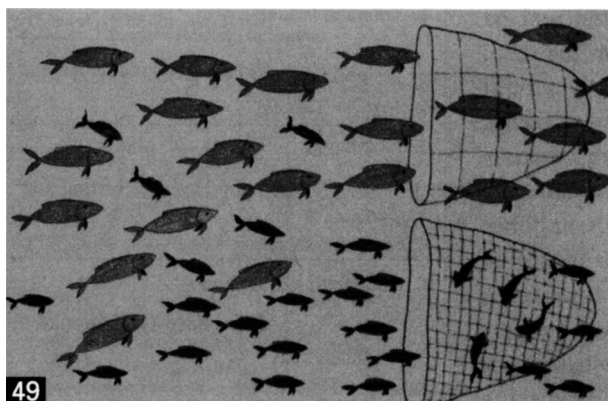
- Any statistical relationship between two variables may be **reversed** by including additional factors in the analysis.

Application: The adjustment problem

- Which factors **should** be included in the analysis.

48

to a formula involving no “do” symbols, which denotes an expression that is computable from nonexperimental data.



49

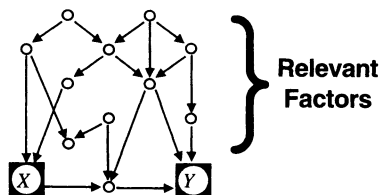
causal-link model is based on the understanding that smoking affects lung cancer through the accumulation of tar deposits in the lungs, perhaps we can measure the amount of tar deposits in the lungs of sampled individuals, and this might provide the necessary information for quantifying the links. Both sides agree that this is a reasonable suggestion, so they submit a new query to the statistician: Can we find the effect of smoking on cancer assuming that an intermediate measurement of tar

The statistician receives the problem and treats it as a problem in high school *algebra*: We need to compute $P(\text{Cancer})$, under hypothetical action, from nonexperimental data – namely, from expressions involving *no actions*. Or: We need to eliminate the “do” symbol from the initial expression. The elimination proceeds like ordinary solution of algebraic equations – in each stage [47], a new rule is applied, licensed by some subgraph of the diagram, eventually leading

You are probably wondering whether this derivation solves the smoking–cancer debate. The answer is *no*. Even if we could get the data on tar deposits, our model is quite simplistic, as it is based on certain assumptions that both parties might not agree to – for instance, that there is no direct link between smoking and lung cancer unmediated by tar deposits. The model would need to be refined

then, and we might end up with a graph containing twenty variables or more. There is no need to panic when someone tells us: “you did not take this or that factor into account.” On the contrary, the graph welcomes such new ideas, because it is so easy to add factors and measurements into the model. Simple tests are now available that permit an investigator to merely glance at the graph and decide if we can compute the effect of one variable on another.

THE ADJUSTMENT PROBLEM



Given: Causal graph
 Needed: Effect of X on Y
 Decide: Which measurements should be taken?

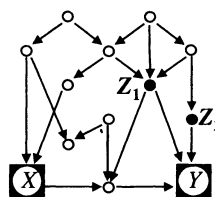
50

Our next example illustrates how a long-standing problem is solved by purely graphical means – proven by the new algebra. The problem is called *the adjustment problem* or “the covariate selection problem” and represents the practical side of Simpson’s paradox [48].

Simpson’s paradox, first noticed by Karl Pearson in 1899, concerns the disturbing observation that every statistical relationship between two variables may be *reversed* by including additional factors in the analysis. For example, you might run a study and find that students who smoke get higher grades; however, if you adjust for *age*, the opposite is true in every *age group*, that is, smoking predicts lower grades. If you further adjust for *parent income*, you find that smoking predicts higher grades again, in every *age–income* group, and so on.

Equally disturbing is the fact that no one has been able to tell us which factors *should* be included in the analysis. Such factors can now be identified by simple graphical means. The classical case demonstrating Simpson’s paradox took place in 1975, when UC-Berkeley

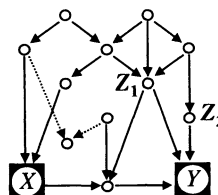
GRAPHICAL SOLUTION OF THE ADJUSTMENT PROBLEM



Subproblem:
 Test if Z_1 and Z_2 are sufficient measurements
STEP 1: Z_1 and Z_2 should not be descendants of X

51

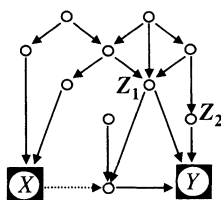
GRAPHICAL SOLUTION OF THE ADJUSTMENT PROBLEM (Cont.)



STEP 2: Delete all non-ancestors of $\{X, Y, Z\}$

52

GRAPHICAL SOLUTION OF THE ADJUSTMENT PROBLEM (Cont.)

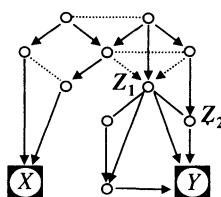


STEP 3: Delete all arcs emanating from X

53

ine a fishing boat with two different nets, a large mesh and a small net [49]. A school

GRAPHICAL SOLUTION OF THE ADJUSTMENT PROBLEM (Cont.)

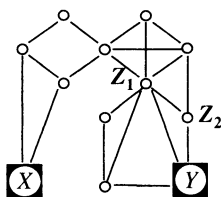


STEP 4: Connect any two parents sharing a common child

54

science literature in the 1970s. Should we, in salary discrimination cases, compare

GRAPHICAL SOLUTION OF THE ADJUSTMENT PROBLEM (Cont.)



STEP 5: Strip arrow-heads from all edges

55

was investigated for sex bias in graduate admission. In this study, overall data showed a higher rate of admission among male applicants; but, broken down by departments, data showed a slight bias in favor of admitting female applicants. The explanation is simple: female applicants tended to apply to more competitive departments than males, and in these departments, the rate of admission was low for both males and females.

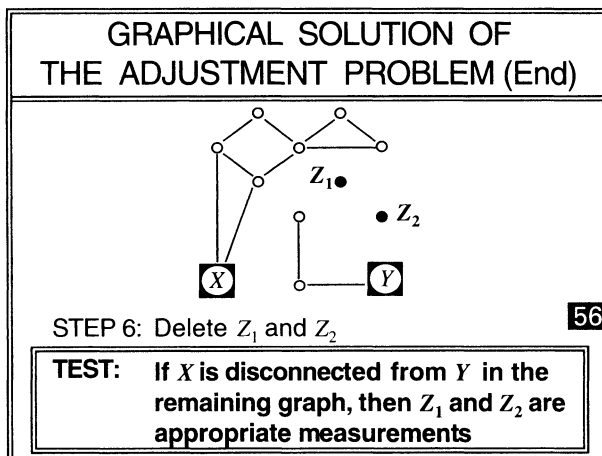
To illustrate this point, imagine a fishing boat with two different nets, a large mesh and a small net [49]. A school of fish swim toward the boat and seek to pass it. The female fish try for the small-mesh challenge, while the male fish try for the easy route. The males go through and only females are caught. Judging by the final catch, preference toward females is clearly evident. However, if analyzed separately, each individual net would surely trap males more easily than females.

Another example involves a controversy called “reverse regression,” which occupied the social science literature in the 1970s. Should we, in salary discrimination cases, compare salaries of equally qualified men and women or instead compare qualifications of equally paid men and women?

Remarkably, the two choices led to opposite conclusions. It turned out that men earned a higher salary than equally qualified women and, *simultaneously*, men were more qualified than equally paid women. The moral is that all conclusions are extremely sensitive to which variables we choose to hold constant when we are comparing,

and that is why the adjustment problem is so critical in the analysis of observational studies.

Consider an observational study where we wish to find the effect of X on Y , for example, treatment on response [50]. We can think of many factors that are relevant to the problem; some are affected by the treatment, some are affecting the treatment, and some are affecting both treatment and response. Some of these factors may be unmeasurable, such as genetic trait or life style; others are measurable, such as gender, age, and salary level. Our problem is to select a subset of these factors for measurement and adjustment so that, if we compare subjects under the same value of those measurements and average, we get the right result.



Let us follow together the steps that would be required to test if two candidate measurements, Z_1 and Z_2 , would be sufficient [51]. The steps are rather simple, and can be performed manually even on large graphs. However, to give you the feel of their mechanizability, I will go through them rather quickly. Here we go [52–56].

At the end of these manipulations, we end up with the answer to our question: “If X is disconnected from Y , then Z_1 and Z_2 are appropriate measurements.”

I now wish to summarize briefly the central message of this lecture. It is true that testing for cause and effect is difficult. Discovering causes of effects is even more difficult. But causality is not *mystical* or *metaphysical*. It can be understood in terms of simple processes, and it can be expressed in a friendly mathematical language, ready for computer analysis.

What I have presented to you today is a sort of pocket calculator, an *abacus* [57], to help us investigate certain problems of cause and effect with mathematical precision. This does not solve all the problems of causality, but the power of *symbols* and mathematics should not be underestimated [58].

Many scientific discoveries have been delayed over the centuries for the lack of a mathematical language that can amplify ideas and let scientists communicate results. I am convinced that many discoveries have been delayed in our century for lack of a mathematical language that can handle

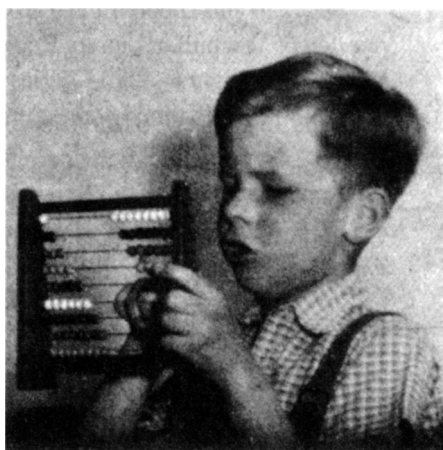


Fig. 155 Little Johnny and his “calculating machine.”



causation. For example, I am sure that Karl Pearson could have thought up the idea of *randomized experiment* in 1901 if he had allowed causal diagrams into his mathematics.

But the really challenging problems are still ahead: We still do not have a causal understanding of *poverty* and *cancer* and *intolerance*, and only the accumulation of data and the insight of great minds will eventually lead to such understanding.

The data is all over the place, the insight is yours, and now an abacus is at your disposal, too. I hope the combination amplifies each of these components.

Thank you.

Acknowledgments

Slide 1 (Dürer, *Adam and Eve*, 1504 engraving) courtesy of the Fogg Art Museum, Harvard University Art Museums, Gift of William Gray from the collection of Francis Calley Gray. Photo by Rick Stafford; image copyright © President and Fellows of Harvard College, Harvard University. Slide 2 (Doré, *The Flight of Lot*) copyright William H. Wise & Co. Slide 3 (Egyptian wall painting of Neferronpe playing a board game) courtesy of the Oriental Institute of the University of Chicago.

The following images were reproduced from antiquarian book catalogs, courtesy of Bernard Quaritch, Ltd. (London): slides 4, 5, 6, 7, 8, 9, 15, 27, 31, 36, 37, 38, 40, 42, and 58.

Slides 10 and 11 copyright The Courier Press. Slides 13 and 14 reprinted with the permission of Macmillan Library Reference USA, from *The Album of Science*, by I. Bernard Cohen. Copyright © 1980 Charles Scribner's Sons.

Slide 16 courtesy of the Library of California State University, Long Beach. Slides 20 and 22 reprinted with the permission of Cambridge University Press. Slide 25: copyright photograph by A. C. Barrington Brown, reproduced with permission.

Slide 30: from S. Wright (1920) in *Proceedings of the National Academy of Sciences*, vol. 6; reproduced with the permission of the American Philosophical Society and the University of Chicago Press. Slide 57 reprinted with the permission of Vandenhoeck & Ruprecht and The MIT Press.

NOTE: Color versions of slides 19, 26, 28–29, 32–35, and 43–56 may be downloaded from <http://www.cs.ucla.edu/~judea/>.