

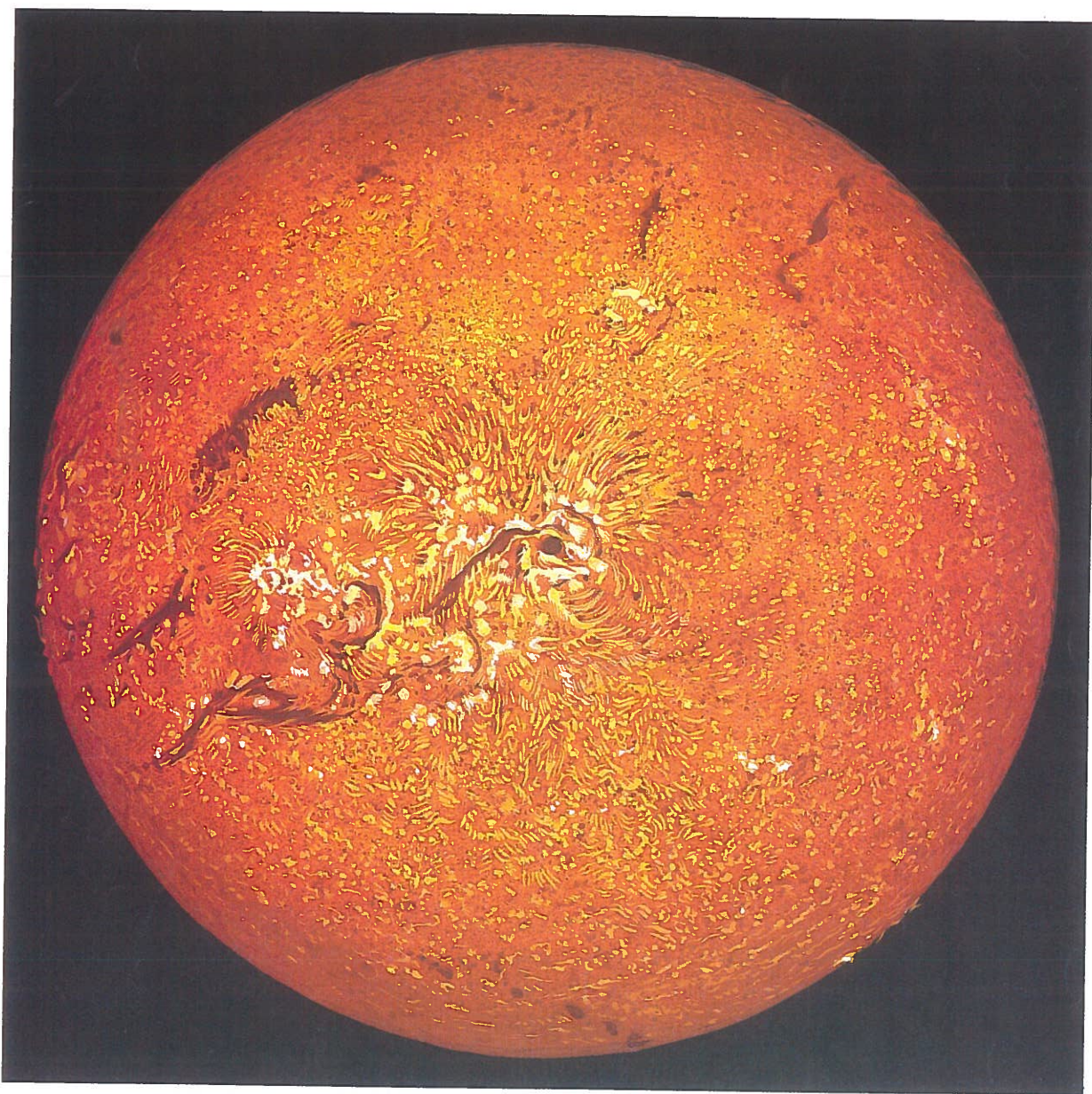
SCIENTIFIC AMERICAN

FEBRUARY 1990
\$2.95

Why morphine taken for pain is not addictive.

Gallium arsenide: a dynamic technology comes of age.

Altruism in—would you believe it?—vampire bats.



Variable inferno: the churning sun grows brighter and more turbulent as it approaches the peak of its 11-year activity cycle.

Positive Feedbacks in the Economy

A new economic theory elucidates mechanisms whereby small chance events early in the history of an industry or technology can tilt the competitive balance

by W. Brian Arthur

Conventional economic theory is built on the assumption of diminishing returns. Economic actions engender a negative feedback that leads to a predictable equilibrium for prices and market shares. Such feedback tends to stabilize the economy because any major changes will be offset by the very reactions they generate. The high oil prices of the 1970's encouraged energy conservation and increased oil exploration, precipitating a predictable drop in prices by the early 1980's. According to conventional theory, the equilibrium marks the "best" outcome possible under the circumstances: the most efficient use and allocation of resources.

Such an agreeable picture often does violence to reality. In many parts of the economy, stabilizing forces appear not to operate. Instead positive feedback magnifies the effects of small economic shifts; the economic models that describe such effects differ vastly from the conventional ones. Diminishing returns imply a single equilibrium point for the economy, but positive feedback—increasing returns—makes for many possible equilibrium points. There is no guarantee that the particular economic outcome selected from among the many alter-

natives will be the "best" one. Furthermore, once random economic events select a particular path, the choice may become locked-in regardless of the advantages of the alternatives. If one product or nation in a competitive marketplace gets ahead by "chance," it tends to stay ahead and even increase its lead. Predictable, shared markets are no longer guaranteed.

During the past few years I and other economic theorists at Stanford University, the Santa Fe Institute in New Mexico and elsewhere have been developing a view of the economy based on positive feedback. Increasing-returns economics has roots that go back 70 years or more, but its application to the economy as a whole is largely new. The theory has strong parallels with modern nonlinear physics (instead of the pre-20th-century physical models that underlie conventional economics), it requires new and challenging mathematical techniques and it appears to be the appropriate theory for understanding modern high-technology economies.

The history of the videocassette recorder furnishes a simple example of positive feedback. The VCR market started out with two competing formats selling at about the same price: VHS and Beta. Each format could realize increasing returns as its market share increased: large numbers of VHS recorders would encourage video outlets to stock more prerecorded tapes in VHS format, thereby enhancing the value of owning a VHS recorder and leading more people to buy one. (The same would, of course, be true for Beta-format players.) In this way, a small gain in market share would improve the competitive position of one system and help it further increase its lead.

Such a market is initially unstable. Both systems were introduced at about the same time and so began with roughly equal market shares; those shares fluctuated early on because of external circumstance, "luck" and corporate maneuvering. Increasing returns on early gains eventually tilted the competition toward VHS: it accumulated enough of an advantage to take virtually the entire VCR market. Yet it would have been impossible at the outset of the competition to say which system would win, which of the two possible equilibria would be selected. Furthermore, if the claim that Beta was technically superior is true, then the market's choice did not represent the best economic outcome.

Conventional economic theory offers a different view of competition between two technologies or products performing the same function. An example is the competition between water and coal to generate electricity. As hydroelectric plants take more of the market, engineers must exploit more costly dam sites, thereby increasing the chance that a coal-fired plant will be cheaper. As coal plants take more of the market, they bid up the price of coal (or trigger the imposition of costly pollution controls) and so tip the balance toward hydropower. The two technologies end up sharing the market in a predictable proportion that best exploits the potentials of each, in contrast to what happened to the two video-recorder systems.

The evolution of the VCR market would not have surprised the great Victorian economist Alfred Marshall, one of the founders of today's conventional economics. In his 1890 *Principles of Economics*, he noted that if firms' production costs fall as their market shares increase, a firm that simply by good fortune gained a high

W. BRIAN ARTHUR is Morrison Professor of Population Studies and Economics at Stanford University. He obtained his Ph.D. from the University of California, Berkeley, in 1973 and holds graduate degrees in operations research, economics and mathematics. Until recently Arthur was on leave at the Santa Fe Institute, a research institute dedicated to the study of complex systems. There he directed a team of economists, physicists, biologists and others investigating behavior of the economy as an evolving, complex system.

proportion of the market early on would be able to best its rivals; "whatever firm first gets a good start" would corner the market. Marshall did not follow up this observation, however, and theoretical economics has until recently largely ignored it.

Marshall did not believe that increasing returns applied everywhere; agriculture and mining—the mainstays of the economies of his time—were subject to diminishing returns caused by limited amounts of fertile land or high-quality ore deposits. Manufacturing, on the other hand, enjoyed increasing returns because large plants allowed improved organization. Modern economists do not see economies of scale as a reliable source of increasing returns. Sometimes large plants have proved more economical; often they have not.

I would update Marshall's insight by observing that the parts of the economy that are resource-based (agriculture, bulk-goods production, mining) are still for the most part subject to diminishing returns. Here conventional economics rightly holds sway. The parts of the economy that are knowledge-based, on the other hand, are largely subject to increasing returns. Products such as computers, pharmaceuticals, missiles, aircraft, automobiles, software, telecommunications equipment or fiber optics are complicated to design and to manufacture.

They require large initial investments in research, development and tooling, but once sales begin, incremental production is relatively cheap. A new airframe or aircraft engine, for example, typically costs between \$2 and \$3 billion to design, develop, certify and put into production. Each copy thereafter costs perhaps \$50 to \$100 million. As more units are built, unit costs continue to fall and profits increase.

Increased production brings additional benefits: producing more units means gaining more experience in the manufacturing process and achieving greater understanding of how to produce additional units even more cheaply. Moreover, experience gained with one product or technology can make it easier to produce new products incorporating similar or related technologies. Japan, for example, leveraged an initial investment in building precision instruments into a capacity for building consumer electronics products and then the integrated circuits that went into them.

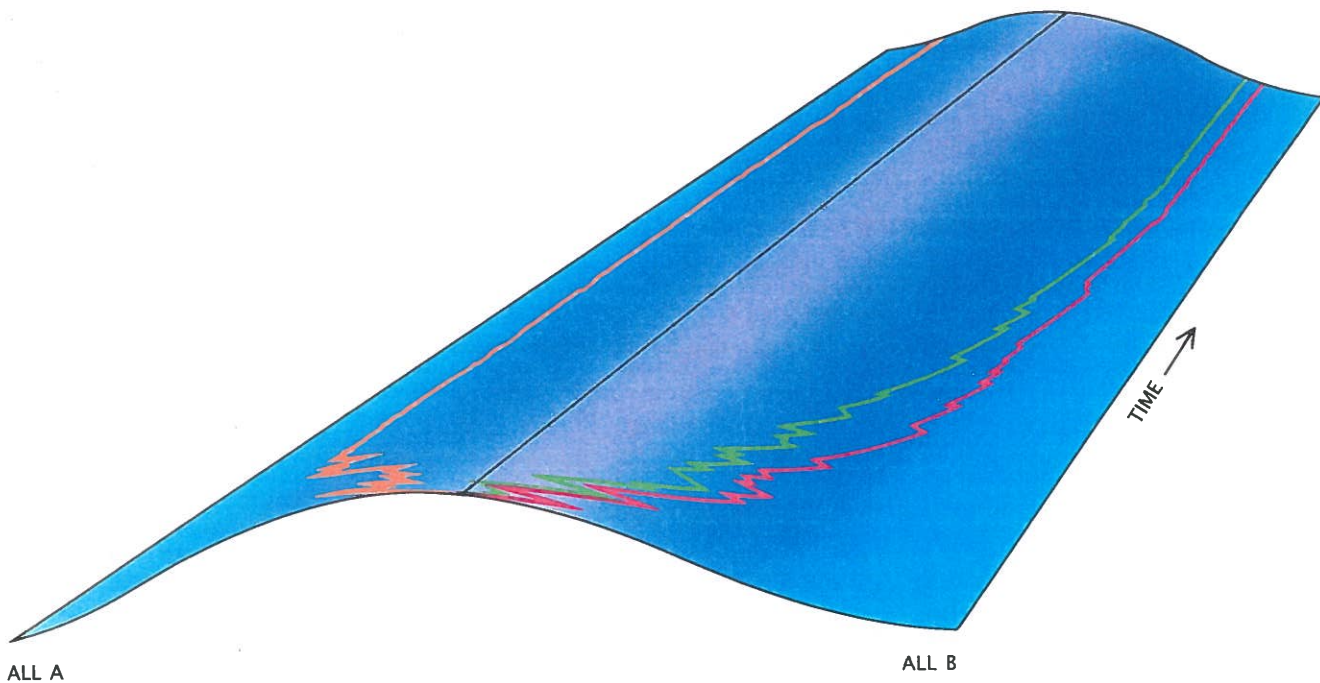
Not only do the costs of producing high-technology products fall as a company makes more of them, but the benefits of using them increase. Many items such as computers or telecommunications equipment work in networks that require compatibility; when one brand gains a significant market share, people have a strong incentive to buy more of the same prod-

uct so as to be able to exchange information with those using it already.

If increasing returns are important, why were they largely ignored until recently? Some would say that complicated products—high technology—for which increasing returns are so important, are themselves a recent phenomenon. This is true but is only part of the answer. After all, in the 1940's and 1950's, economists such as Gunnar K. Myrdal and Nicholas Kaldor identified positive-feedback mechanisms that did not involve technology. Orthodox economists avoided increasing returns for deeper reasons.

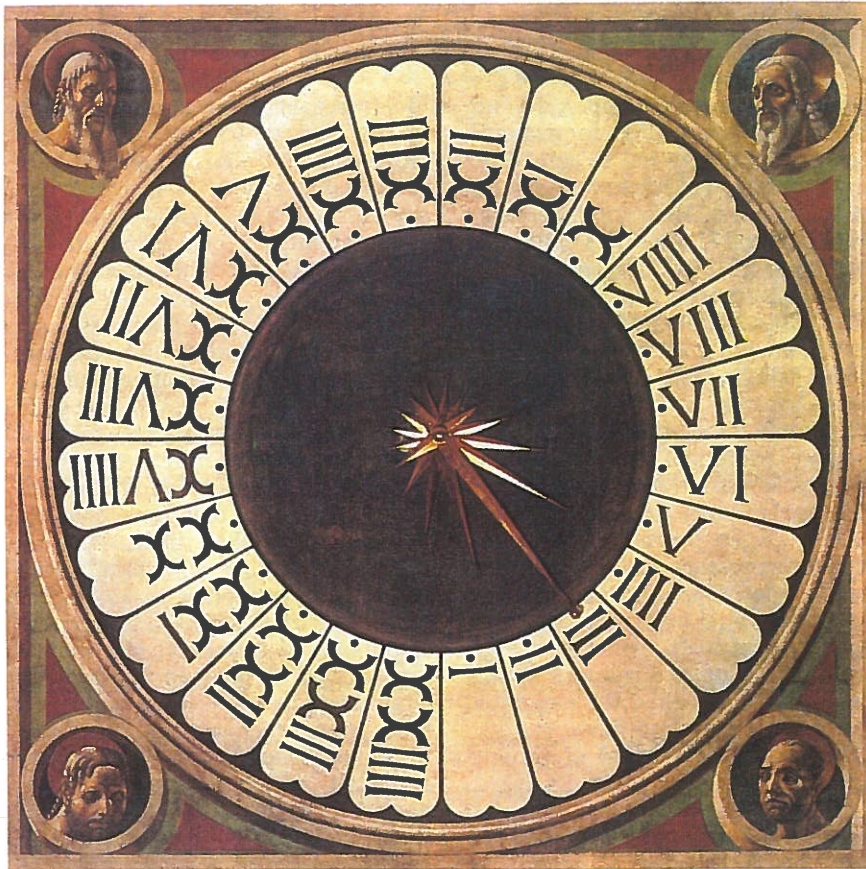
Some economists found the existence of more than one solution to the same problem distasteful—unscientific. "Multiple equilibria," wrote Joseph A. Schumpeter in 1954, "are not necessarily useless, but from the standpoint of any exact science the existence of a uniquely determined equilibrium is, of course, of the utmost importance, even if proof has to be purchased at the price of very restrictive assumptions; without any possibility of proving the existence of [a] uniquely determined equilibrium—or at all events, of a small number of possible equilibria—at however high a level of abstraction, a field of phenomena is really a chaos that is not under analytical control."

Other economists could see that



RANDOM WALK on a convex surface illustrates increasing-returns competition between two technologies. Chance determines early patterns of adoption and so influences how

fast each competitor improves. As one technology gains more adherents (corresponding to motion downhill toward either edge of the surface), further adoption is increasingly likely.



FLORENCE CATHEDRAL CLOCK has hands that move "counterclockwise" around its 24-hour dial. When Paolo Uccello designed the clock in 1443, a convention for clockfaces had not emerged. Competing designs were subject to increasing returns: the more clockfaces of one kind were built, the more people became used to reading them. Hence, it was more likely that future clockfaces would be of the same kind. After 1550, "clockwise" designs displaying only 12 hours had crowded out other designs. The author argues that chance events coupled with positive feedback, rather than technological superiority, will often determine economic developments.

theories incorporating increasing returns would destroy their familiar world of unique, predictable equilibria and the notion that the market's choice was always best. Moreover, if one or a few firms came to dominate a market, the assumption that no firm is large enough to affect market prices on its own (which makes economic problems easy to analyze) would also collapse. When John R. Hicks surveyed these possibilities in 1939 he drew back in alarm. "The threatened wreckage," he wrote, "is that of the greater part of economic theory." Economists restricted themselves to diminishing returns, which presented no anomalies and could be analyzed completely.

Still others were perplexed by the question of how a market could select one among several possible solutions. In Marshall's example, the firm that is the largest at the outset has the lowest production costs and must inevitably win in the market. In that case, why would smaller firms compete at all?

On the other hand, if by some chance a market started with several identical firms, their market shares would remain poised in an unstable equilibrium forever.

Studying such problems in 1979, I believed I could see a way out of many of these difficulties. In the real world, if several similar-size firms entered a market at the same time, small fortuitous events—unexpected orders, chance meetings with buyers, managerial whims—would help determine which ones achieved early sales and, over time, which firm dominated. Economic activity is quantized by individual transactions that are too small to observe, and these small "random" events can accumulate and become magnified by positive feedbacks so as to determine the eventual outcome. These facts suggested that situations dominated by increasing returns should be modeled not as static, deterministic problems

but as dynamic processes based on random events and natural positive feedbacks, or nonlinearities.

With this strategy an increasing-returns market could be re-created in a theoretical model and watched as its corresponding process unfolded again and again. Sometimes one solution would emerge, sometimes (under identical conditions) another. It would be impossible to know in advance which of the many solutions would emerge in any given run. Still, it would be possible to record the particular set of random events leading to each solution and to study the probability that a particular solution would emerge under a certain set of initial conditions. The idea was simple, and it may well have occurred to economists in the past. But making it work called for nonlinear random-process theory that did not exist in their day.

Every increasing-returns problem need not be studied in isolation; many turn out to fit a general nonlinear probability schema. It can be pictured by imagining a table to which balls are added one at a time; they can be of several possible colors—white, red, green or blue. The color of the ball to be added next is unknown, but the probability of a given color depends on the current proportions of colors on the table. If an increasing proportion of balls of a given color increases the probability of adding another ball of the same color, the system can demonstrate positive feedback. The question is, Given the function that maps current proportions to probabilities, what will be the proportions of each color on the table after many balls have been added?

In 1931 the mathematician George Polya solved a very particular version of this problem in which the probability of adding a color always equaled its current proportion. Three U.S. probability theorists, Bruce M. Hill of the University of Michigan at Ann Arbor and David A. Lane and William D. Suderth of the University of Minnesota at Minneapolis, solved a more general, nonlinear version in 1980. In 1983 two Soviet probability theorists, Yuri M. Ermoliev and Yuri M. Kaniovski, both of the Glushkov Institute of Cybernetics in Kiev, and I found the solution to a very general version. As balls continue to be added, we proved, the proportions of each color must settle down to a "fixed point" of the probability function—a set of values where the probability of adding each color is equal to the proportion of that color on the table. Increasing returns allow several such sets of fixed points.

This means that we can determine the possible patterns or solutions of an increasing-returns problem by solving the much easier challenge of finding the sets of fixed points of its probability function. With such tools economists can now define increasing-returns problems precisely, identify their possible solutions and study the process by which a solution is reached. Increasing returns are no longer "a chaos that is not under analytical control."

In the real world, the balls might be represented by companies and their colors by the regions where they decide to settle. Suppose that firms enter an industry one by one and choose their locations so as to maximize profit. The geographic preference of each firm (the intrinsic benefits it gains from being in a particular region) varies; chance determines the preference of the next firm to enter the industry. Also suppose, however, that firms' profits increase if they are near other firms (their suppliers or customers). The first firm to enter the industry picks a location based purely on geographic preference. The second firm decides based on preference modified by the benefits gained by locating near the first firm. The third firm is influenced by the positions of the first two firms, and so on. If some location by good fortune attracts more firms than the others in the early stages of this evolution, the probability that it will attract more firms increases. Industrial concentration becomes self-reinforcing.

The random historical sequence of firms entering the industry determines which pattern of regional settlement results, but the theory shows that not all patterns are possible. If the attractiveness exerted by the presence of other firms always rises as more firms are added, some region will always dominate and shut out all others. If the attractiveness levels off, other solutions, in which regions share the industry, become possible. Our new tools tell us which types of solutions can occur under which conditions.

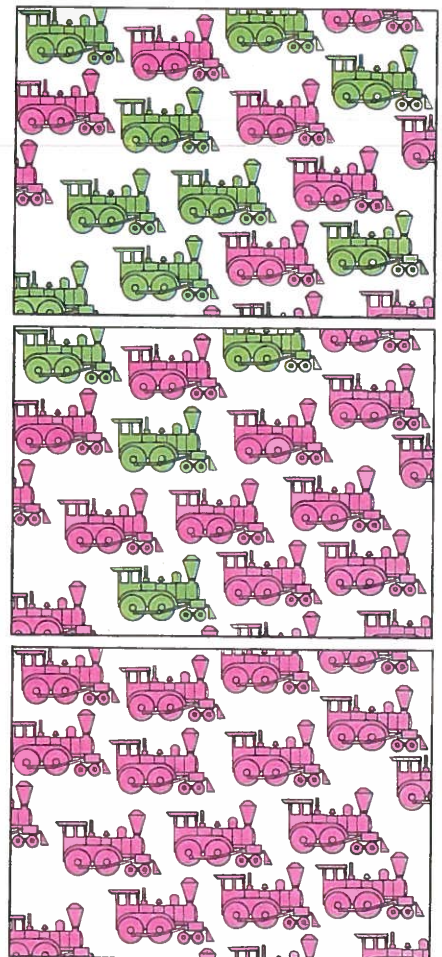
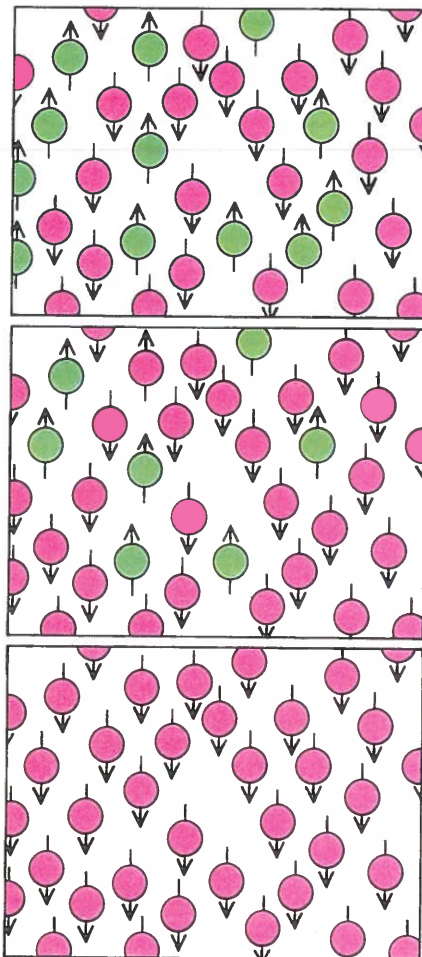
Do some regions in fact amass a large proportion of an industry because of historical chance rather than geographic superiority? Santa Clara County in California (Silicon Valley) is a likely example. In the 1940's and early 1950's certain key people in the U.S. electronics industry—the Varian brothers, William Hewlett and David Packard, William Shockley—set up shop near Stanford University; the local availability of engineers, supplies

and components that these early firms helped to create made Santa Clara County extremely attractive to the 900 or so firms that followed. If these early entrepreneurs had preferred other places, the densest concentration of electronics in the country might well be somewhere else.

On a grander scale, if small events in history had been different, would the location of cities themselves be different? I believe the answer is yes. To the degree that certain locations are natural harbors or junction points on rivers or lakes, the pattern of cities today reflects not chance but geography. To the degree that industry and people are attracted to places where such resources are already gathered, small, early chance concentrations may have been the seeds of today's configuration of urban centers. "Chance and necessity," to use Jacques Monod's

phrase, interact. Both have played crucial roles in the development of urban centers in the U.S. and elsewhere.

Self-reinforcing mechanisms other than these regional ones work in international high-tech manufacturing and trade. Countries that gain high volume and experience in a high-technology industry can reap advantages of lower cost and higher quality that may make it possible for them to shut out other countries. For example, in the early 1970's, Japanese automobile makers began to sell significant numbers of small cars in the U.S. As Japan gained market volume without much opposition from Detroit, its engineers and production workers gained experience, its costs fell and its products improved. These factors, together with improved sales networks, allowed Japan to increase



FERROMAGNETS AND REGIONAL RAIL GAUGES become ordered in much the same way. As a disordered magnetic material is cooled (left), the atomic dipoles inside it exert forces on one another, causing neighboring dipoles to align. Eventually all the dipoles in a sample line up, but the direction they all take (up or down) cannot be predicted beforehand. Similarly, as Douglas Puffert of Swarthmore College has shown, neighboring private railroads (right) in the past century adopted the same gauge to extend their range more easily. Eventually all (or most) railroads used the same gauge. Similar equations describe the behavior of these two systems.

its share of the U.S. market; as a result, workers gained still more experience, costs fell further and quality improved again. Before Detroit responded seriously, this positive-feedback loop had helped Japanese companies to make serious inroads into the U.S. market for small cars. Similar sequences of events have taken place in the markets for television sets, integrated circuits and other products.

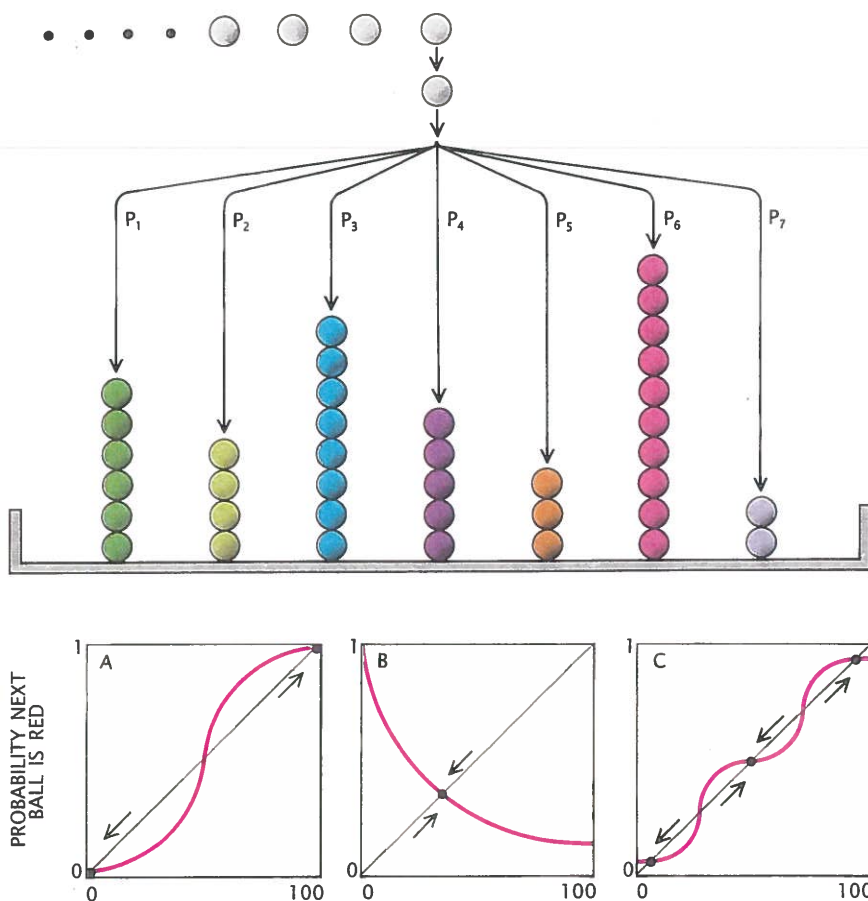
How should countries respond to a world economy where such rules apply? Conventional recommendations for trade policy based on constant or diminishing returns tend toward low-profile approaches. They rely on the open market, discourage monopolies and leave issues such as R&D spending to companies. Their underlying assumption is that there is a fixed world price at which producers load goods onto the market, and so inter-

ference with local costs and prices by means of subsidies or tariffs is unproductive. These policies are appropriate for the diminishing-returns parts of the economy, not for the technology-based parts where increasing returns dominate.

Policies that are appropriate to success in high-tech production and international trade would encourage industries to be aggressive in seeking out product and process improvements. They would strengthen the national research base on which high-tech advantages are built. They would encourage firms in a single industry to pool their resources in joint ventures that share up-front costs, marketing networks, technical knowledge and standards. They might even foster strategic alliances, enabling companies in several countries to enter a complex industry that none could

tackle alone. Increasing-returns theory also points to the importance of timing when undertaking research initiatives in new industries. There is little sense in entering a market that is already close to being locked-in or that otherwise offers little chance of success. Such policies are slowly being advocated and adopted in the U.S.

The value of other policies, such as subsidizing and protecting new industries—bioengineering, for example—to capture foreign markets, is debatable. Dubious feedback benefits have sometimes been cited to justify government-sponsored white elephants. Furthermore, as Paul R. Krugman of the Massachusetts Institute of Technology and several other economists have pointed out, if one country pursues such policies, others will retaliate by subsidizing their own high-technology industries. Nobody gains. The question of optimal industrial and trade policy based on increasing returns is currently being studied intensely. The policies countries choose will determine not only the shape of the global economy in the 1990's but also its winners and its losers.

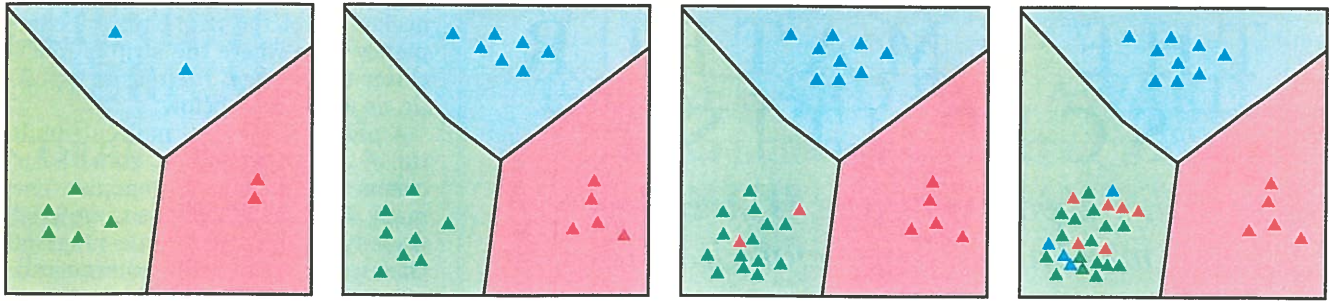


NONLINEAR PROBABILITY THEORY can predict the behavior of systems subject to increasing returns. In this model, balls of different colors are added to a table; the probability that the next ball will have a specific color depends on the current proportions of colors (top). Increasing returns occur in A (the graph shows the two-color case; arrows indicate likely directions of motion): a red ball is more likely to be added when there is already a high proportion of red balls. This case has two equilibrium points: one at which almost all balls are red; the other at which very few are red. Diminishing returns occur in B: a higher proportion of red balls lowers the probability of adding another. There is a single equilibrium point. A combination of increasing and diminishing returns (C) yields many equilibrium points.

Increasing-returns mechanisms do not merely tilt competitive balances among nations; they can also cause economies—even such successful ones as those of the U.S. and Japan—to become locked into inferior paths of development. A technology that improves slowly at first but has enormous long-term potential could easily be shut out, locking an economy into a path that is both inferior and difficult to escape.

Technologies typically improve as more people adopt them and firms gain experience that guides further development. This link is a positive-feedback loop: the more people adopt a technology, the more it improves and the more attractive it is for further adoption. When two or more technologies (like two or more products) compete, positive feedbacks make the market for them unstable. If one pulls ahead in the market, perhaps by chance, its development may accelerate enough for it to corner the market. A technology that improves more rapidly as more people adopt it stands a better chance of surviving—it has a “selectional advantage.” Early superiority, however, is no guarantee of long-term fitness.

In 1956, for example, when the U.S. embarked on its nuclear-power program, a number of designs were proposed: reactors cooled by gas, light water, heavy water, even liquid sodi-



COMPANIES CHOOSE LOCATIONS to maximize profits, which are determined by intrinsic geographic preference (shown by color) and by the presence of other companies. In this computer-generated example, most of the first few companies set-

tle in the green region, and so all new companies eventually settle there. Such clustering might appear to imply that the green region is somehow superior. In other runs of the program, however, the red and blue regions dominate instead.

um. Robin Cowan of New York University has shown that a series of trivial circumstances locked virtually the entire U.S. nuclear industry into light water. Light-water reactors were originally adapted from highly compact units designed to propel nuclear submarines. The role of the U.S. Navy in early reactor-construction contracts, efforts by the National Security Council to get a reactor—any reactor—working on land in the wake of the 1957 *Sputnik* launch as well as the predilections of some key officials all acted to favor the early development of light-water reactors. Construction experience led to improved light-water designs and, by the mid-1960's, fixed the industry's path. Whether other designs would, in fact, have been superior in the long run is open to question, but much of the engineering literature suggests that high-temperature, gas-cooled reactors would have been better.

Technological conventions or standards, as well as particular technologies, tend to become locked-in by positive feedback, as my colleague Paul A. David of Stanford has documented in several historical instances. Although a standard itself may not improve with time, widespread adoption makes it advantageous for newcomers to a field—who must exchange information or products with those already working there—to fall in with the standard, be it the English language, a high-definition television system, a screw thread or a typewriter keyboard. Standards that are established early (such as the 1950's-vintage computer language FORTRAN) can be hard for later ones to dislodge, no matter how superior would-be successors may be.

Until recently conventional economics texts have tended to portray the economy as something akin to a large Newtonian system, with a unique equilibrium solu-

tion preordained by patterns of mineral resources, geography, population, consumer tastes and technological possibilities. In this view, perturbations or temporary shifts—such as the oil shock of 1973 or the stock-market crash of 1987—are quickly negated by the opposing forces they elicit. Given future technological possibilities, one should in theory be able to forecast accurately the path of the economy as a smoothly shifting solution to the analytical equations governing prices and quantities of goods. History, in this view, is not terribly important; it merely delivers the economy to its inevitable equilibrium.

Positive-feedback economics, on the other hand, finds its parallels in modern nonlinear physics. Ferromagnetic materials, spin glasses, solid-state lasers and other physical systems that consist of mutually reinforcing elements show the same properties as the economic examples I have given. They “phase lock” into one of many possible configurations; small perturbations at critical times influence which outcome is selected, and the chosen outcome may have higher energy (that is, be less favorable) than other possible end states.

This kind of economics also finds parallels in the evolutionary theory of punctuated equilibrium. Small events (the mutations of history) are often averaged away, but once in a while they become all-important in tilting parts of the economy into new structures and patterns that are then preserved and built on in a fresh layer of development.

In this new view, initially identical economies with significant increasing-returns sectors do not necessarily select the same paths. Instead they eventually diverge. To the extent that small events determining the overall path always remain beneath the resolution of the economist's lens, accurate forecasting of an economy's future may be

theoretically, not just practically, impossible. Steering an economy with positive feedbacks into the best of its many possible equilibrium states requires good fortune and good timing—a feel for the moments when beneficial change from one pattern to another is most possible. Theory can help identify these states and times, and it can guide policymakers in applying the right amount of effort (not too little but not too much) to dislodge locked-in structures.

The English philosopher of science Jacob Bronowski once remarked that economics has long suffered from a fatally simple structure imposed on it in the 18th century. I find it exciting that this is now changing. With the acceptance of positive feedbacks, economists' theories are beginning to portray the economy not as simple but as complex, not as deterministic, predictable and mechanistic but as process-dependent, organic and always evolving.

FURTHER READING

- MARKET STRUCTURE AND FOREIGN TRADE. Elhanan Helpman and Paul Krugman. The MIT Press, 1985.
- PATH-DEPENDENT PROCESSES AND THE EMERGENCE OF MACRO-STRUCTURE. W. Brian Arthur, Yu M. Ermoliev and Yu M. Kaniowski in *European Journal of Operational Research*, Vol. 30, pages 294-303; 1987.
- SELF-REINFORCING MECHANISMS IN ECONOMICS. W. Brian Arthur in *The Economy as an Evolving Complex System*. Edited by Philip W. Anderson, Kenneth J. Arrow and David Pines. Addison-Wesley Publishing Co., 1988.
- PATH-DEPENDENCE: PUTTING THE PAST INTO THE FUTURE OF ECONOMICS. Paul David. I.M.S.S.S. Tech Report No. 533, Stanford University; November, 1988.
- COMPETING TECHNOLOGIES, INCREASING RETURNS, AND LOCK-IN BY HISTORICAL EVENTS. W. Brian Arthur in *The Economic Journal*, Vol. 99, No. 394, pages 116-131; March, 1989.