

# I *Graphical Excellence*

Excellence in statistical graphics consists of complex ideas communicated with clarity, precision, and efficiency. Graphical displays should

- show the data
- induce the viewer to think about the substance rather than about methodology, graphic design, the technology of graphic production, or something else
- avoid distorting what the data have to say
- present many numbers in a small space
- make large data sets coherent
- encourage the eye to compare different pieces of data
- reveal the data at several levels of detail, from a broad overview to the fine structure
- serve a reasonably clear purpose: description, exploration, tabulation, or decoration
- be closely integrated with the statistical and verbal descriptions of a data set.

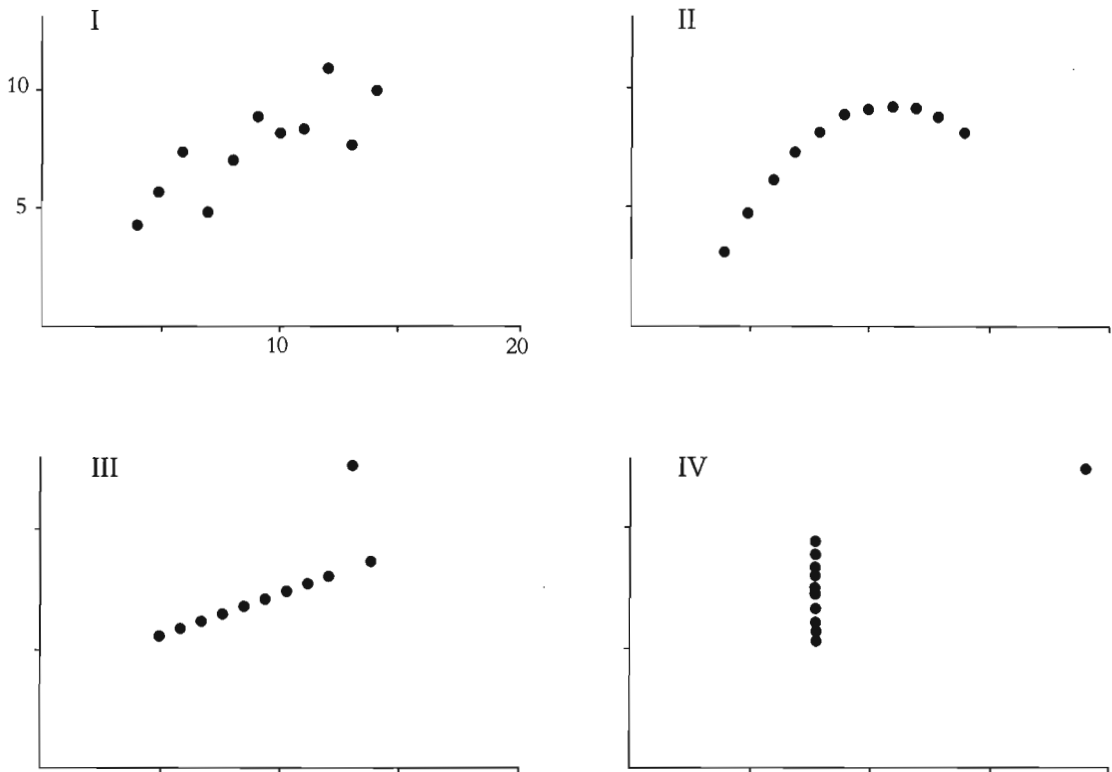
Graphics *reveal* data. Indeed graphics can be more precise and revealing than conventional statistical computations. Consider Anscombe's quartet: all four of these data sets are described by exactly the same linear model (at least until the residuals are examined).

I		II		III		IV	
X	Y	X	Y	X	Y	X	Y
10.0	8.04	10.0	9.14	10.0	7.46	8.0	6.58
8.0	6.95	8.0	8.14	8.0	6.77	8.0	5.76
13.0	7.58	13.0	8.74	13.0	12.74	8.0	7.71
9.0	8.81	9.0	8.77	9.0	7.11	8.0	8.84
11.0	8.33	11.0	9.26	11.0	7.81	8.0	8.47
14.0	9.96	14.0	8.10	14.0	8.84	8.0	7.04
6.0	7.24	6.0	6.13	6.0	6.08	8.0	5.25
4.0	4.26	4.0	3.10	4.0	5.39	19.0	12.50
12.0	10.84	12.0	9.13	12.0	8.15	8.0	5.56
7.0	4.82	7.0	7.26	7.0	6.42	8.0	7.91
5.0	5.68	5.0	4.74	5.0	5.73	8.0	6.89

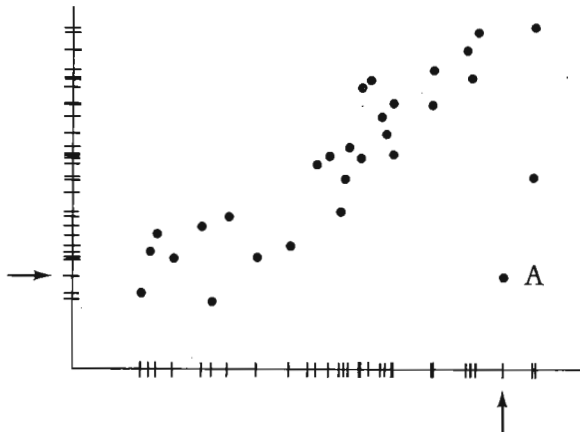
$N = 11$   
 mean of X's = 9.0  
 mean of Y's = 7.5  
 equation of regression line:  $Y = 3 + 0.5X$   
 standard error of estimate of slope = 0.118  
 $t = 4.24$   
 sum of squares  $X - \bar{X} = 110.0$   
 regression sum of squares = 27.50  
 residual sum of squares of Y = 13.75  
 correlation coefficient = .82  
 $r^2 = .67$

And yet how they differ, as the graphical display of the data makes vividly clear:

F. J. Anscombe, "Graphs in Statistical Analysis," *American Statistician*, 27 (February 1973), 17-21.



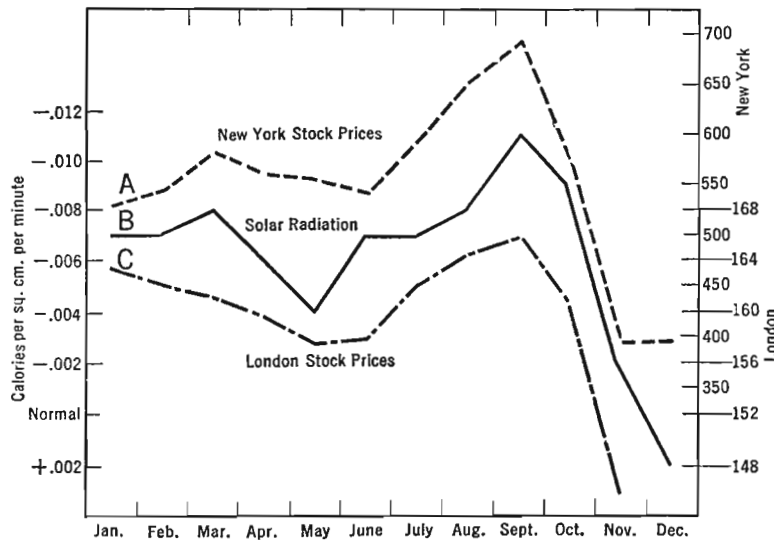
And likewise a graphic easily reveals point A, a wildshot observation that will dominate standard statistical calculations. Note that point A hides in the marginal distribution but shows up as clearly exceptional in the bivariate scatter.



Stephen S. Brier and Stephen E. Fienberg, "Recent Econometric Modelling of Crime and Punishment: Support for the Deterrence Hypothesis?" in Stephen E. Fienberg and Albert J. Reiss, Jr., eds., *Indicators of Crime and Criminal Justice: Quantitative Studies* (Washington, D.C., 1980), p. 89.

Of course, statistical graphics, just like statistical calculations, are only as good as what goes into them. An ill-specified or preposterous model or a puny data set cannot be rescued by a graphic (or by calculation), no matter how clever or fancy. A silly theory means a silly graphic:

Edward R. Dewey and Edwin F. Dakin,  
*Cycles: The Science of Prediction* (New  
York, 1947), p. 144.



SOLAR RADIATION AND STOCK PRICES

A. New York stock prices (Barron's average). B. Solar Radiation, inverted, and C. London stock prices, all by months, 1929 (after Garcia-Mata and Shaffner).

Let us turn to the practice of graphical excellence, the efficient communication of complex quantitative ideas. Excellence, nearly always of a multivariate sort, is illustrated here for fundamental graphical designs: data maps, time-series, space-time narrative designs, and relational graphics. These examples serve several purposes, providing a set of high-quality graphics that can be discussed (and sometimes even redrawn) in constructing a theory of data graphics, helping to demonstrate a descriptive terminology, and telling in brief about the history of graphical development. Most of all, we will be able to see just how good statistical graphics can be.

## Data Maps

These six maps report the age-adjusted death rate from various types of cancer for the 3,056 counties of the United States. Each map portrays some 21,000 numbers.<sup>1</sup> Only a picture can carry such a volume of data in such a small space. Furthermore, all that data, thanks to the graphic, can be thought about in many different ways at many different levels of analysis—ranging from the contemplation of general overall patterns to the detection of very fine county-by-county detail. To take just a few examples, look at the

- high death rates from cancer in the northeast part of the country and around the Great Lakes
- low rates in an east-west band across the middle of the country
- higher rates for men than for women in the south, particularly Louisiana (cancers probably caused by occupational exposure, from working with asbestos in shipyards)
- unusual hot spots, including northern Minnesota and a few counties in Iowa and Nebraska along the Missouri River
- differences in types of cancer by region (for example, the high rates of stomach cancer in the north-central part of the country—probably the result of the consumption of smoked fish by Scandinavians)
- rates in areas where you have lived.

The maps provide many leads into the causes—and avoidance—of cancer. For example, the authors report:

In certain situations . . . the unusual experience of a county warrants further investigation. For example, Salem County, New Jersey, leads the nation in bladder cancer mortality among white men. We attribute this excess risk to occupational exposures, since about 25 percent of the employed persons in this county work in the chemical industry, particularly the manufacturing of organic chemicals, which may cause bladder tumors. After the finding was communicated to New Jersey health officials, a company in the area reported that at least 330 workers in a single plant had developed bladder cancer during the last 50 years. It is urgent that surveys of cancer risk and programs in cancer control be initiated among workers and former workers in this area.<sup>2</sup>

<sup>1</sup>Each county's rate is located in two dimensions and, further, at least four numbers would be necessary to reconstruct the size and shape of each county. This yields  $7 \times 3,056$  entries in a data matrix sufficient to reproduce a map.

In highest decile,  
statistically significant



Significantly high, but  
not in highest decile



In highest decile, but not  
statistically significant



Not significantly different  
from U.S. as a whole

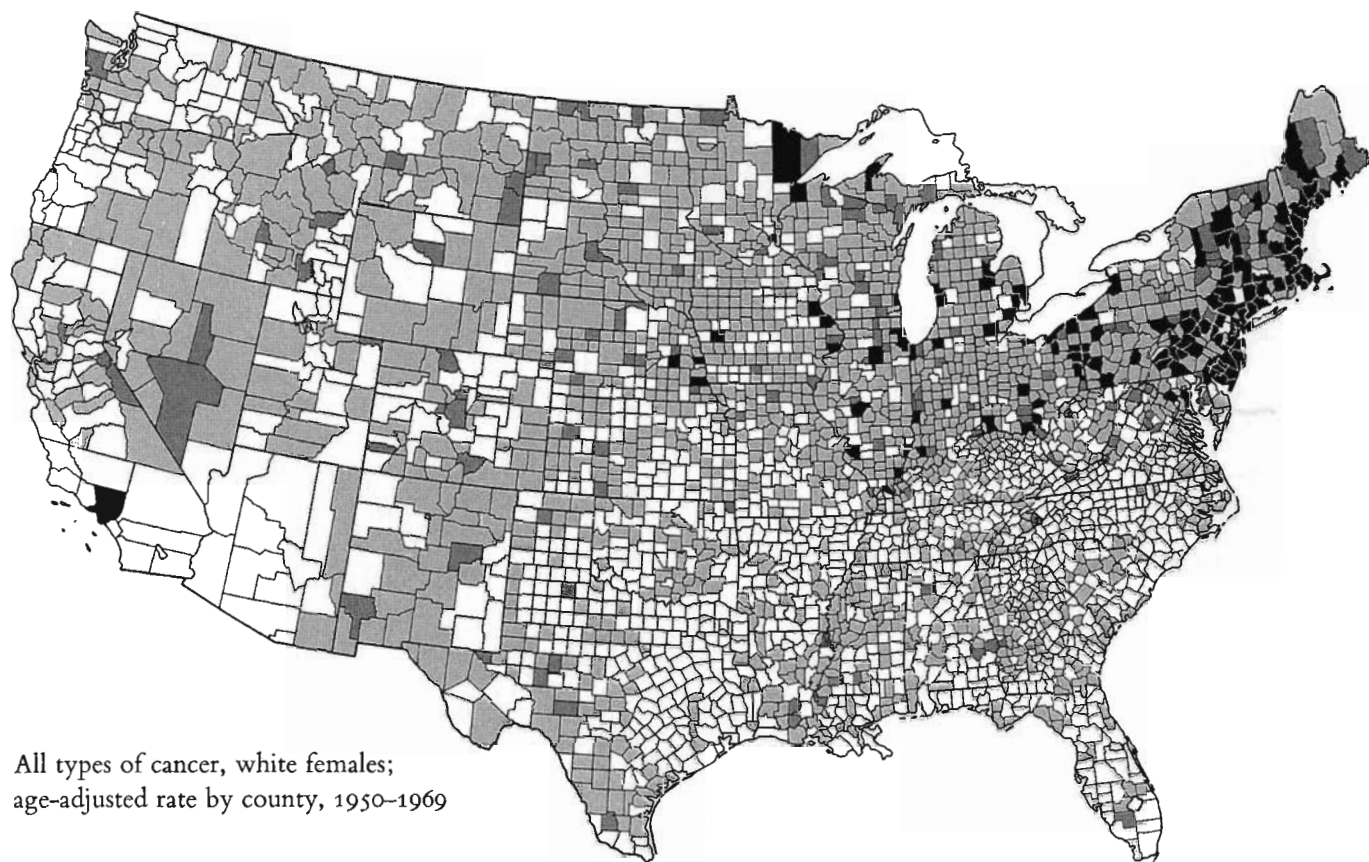


Significantly lower than  
U.S. as a whole

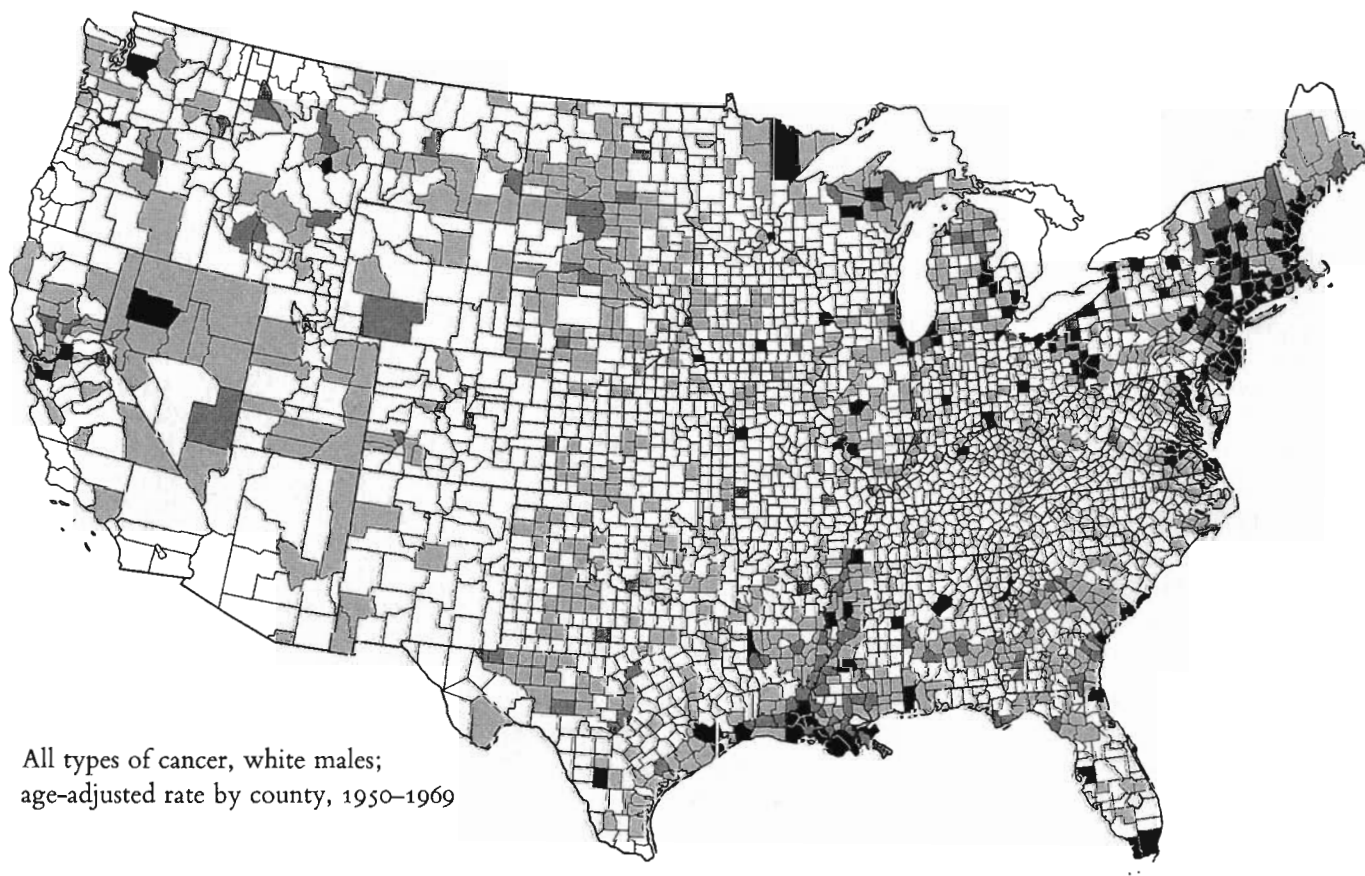


<sup>2</sup>Robert Hoover, Thomas J. Mason, Frank W. McKay, and Joseph F. Fraumeni, Jr., "Cancer by County: New Resource for Etiologic Clues," *Science*, 189 (September 19, 1975), 1006.

Maps from *Atlas of Cancer Mortality for U.S. Counties: 1950–1969*, by Thomas J. Mason, Frank W. McKay, Robert Hoover, William J. Blot, and Joseph F. Fraumeni, Jr. (Washington, D.C.: Public Health Service, National Institutes of Health, 1975). The six maps shown here were redesigned and redrawn by Lawrence Fahey and Edward Tufte.

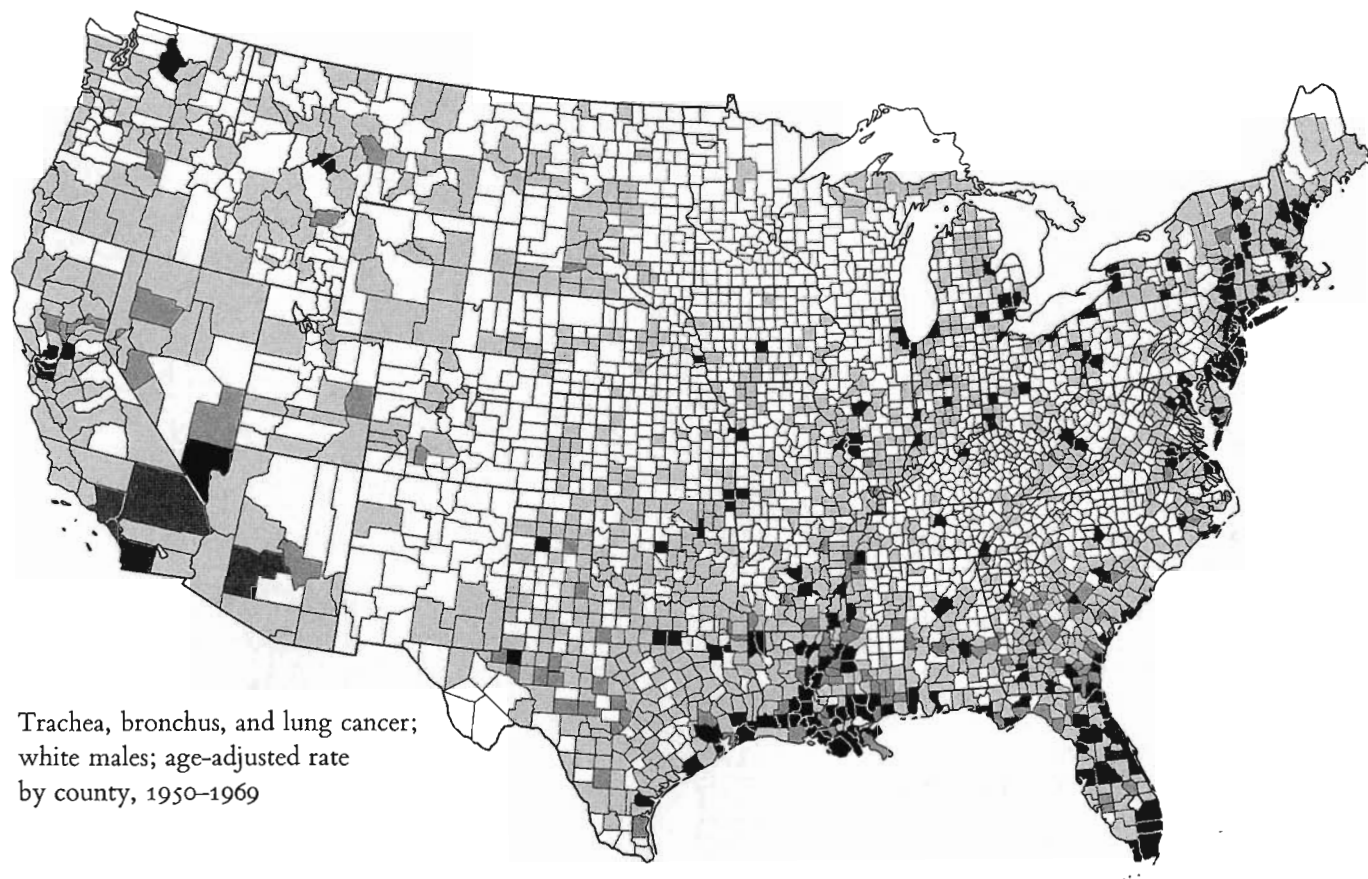
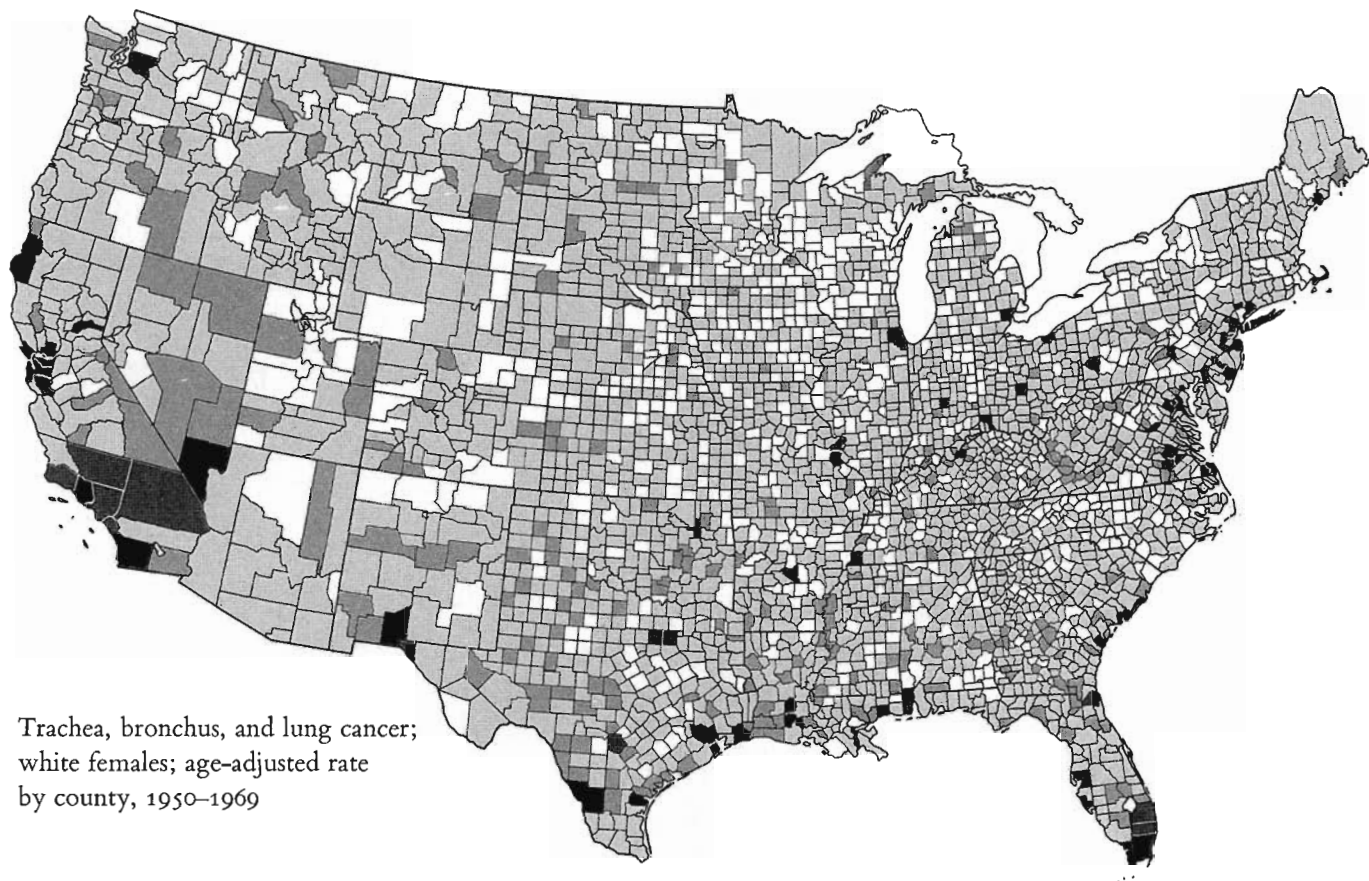


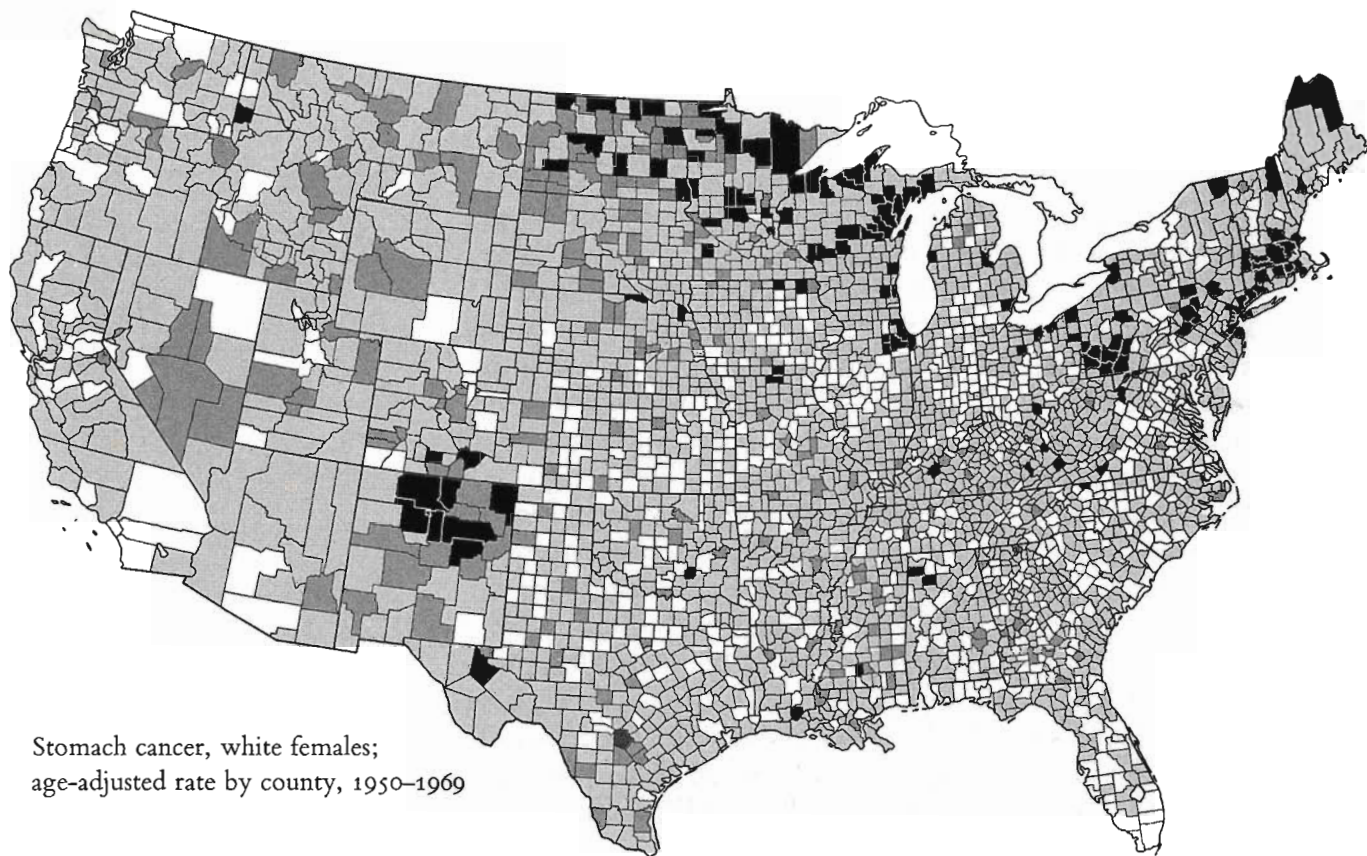
All types of cancer, white females;  
age-adjusted rate by county, 1950-1969



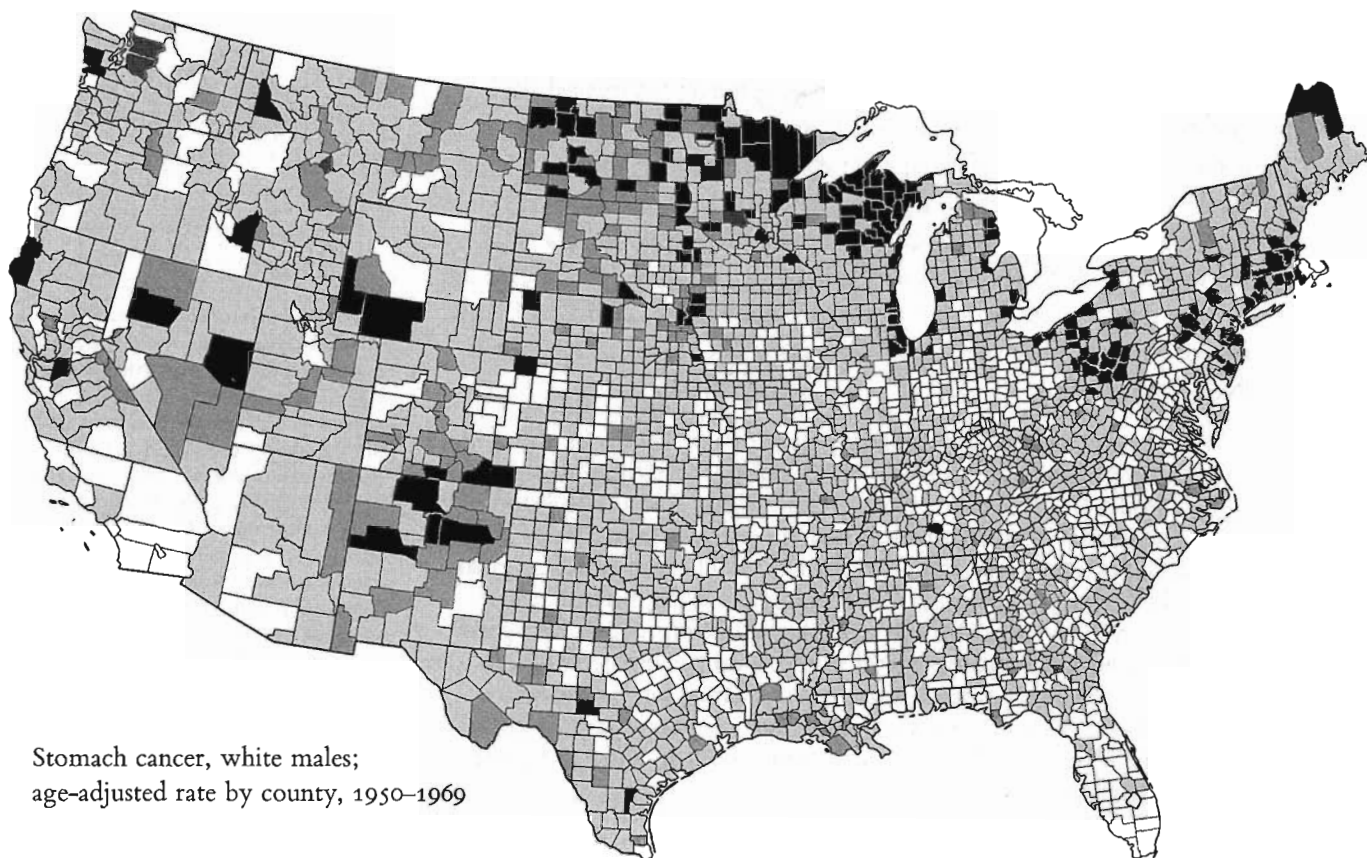
All types of cancer, white males;  
age-adjusted rate by county, 1950-1969







Stomach cancer, white females;  
age-adjusted rate by county, 1950-1969



Stomach cancer, white males;  
age-adjusted rate by county, 1950-1969

The maps repay careful study. Notice how quickly and naturally our attention has been directed toward exploring the substantive content of the data rather than toward questions of methodology and technique. Nonetheless the maps do have their flaws. They wrongly equate the visual importance of each county with its geographic area rather than with the number of people living in the county (or the number of cancer deaths). Our visual impression of the data is entangled with the circumstance of geographic boundaries, shapes, and areas—the chronic problem afflicting shaded-in-area designs of such “blot maps” or “patch maps.”

A further shortcoming, a defect of data rather than graphical composition, is that the maps are founded on a suspect data source, death certificate reports on the cause of death. These reports fall under the influence of diagnostic fashions prevailing among doctors and coroners in particular places and times, a troublesome adulterant of the evidence purporting to describe the already sometimes ambiguous matter of the exact bodily site of the primary cancer. Thus part of the regional clustering seen on the maps, as well as some of the hot spots, may reflect varying diagnostic customs and fads along with the actual differences in cancer rates between areas.

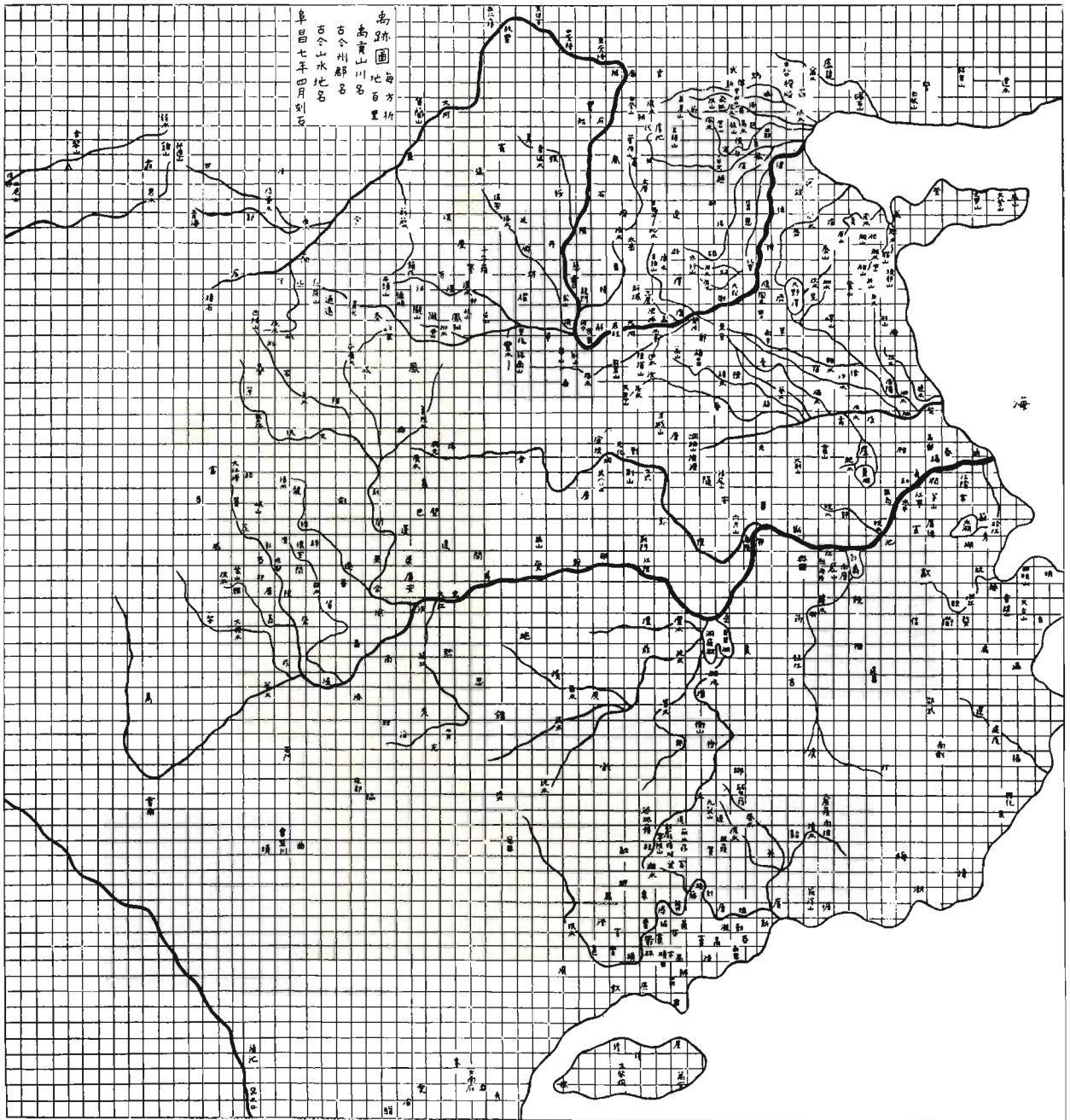
Data maps have a curious history. It was not until the seventeenth century that the combination of cartographic and statistical skills required to construct the data map came together, fully 5,000 years after the first geographic maps were drawn on clay tablets. And many highly sophisticated geographic maps were produced centuries before the first map containing any statistical material was drawn.<sup>3</sup> For example, a detailed map with a full grid was engraved during the eleventh century A.D. in China. The Yü Chi Thu (Map of the Tracks of Yü the Great) shown here is described by Joseph Needham as the

... most remarkable cartographic work of its age in any culture, carved in stone in +1137 but probably dating from before +1100. The scale of the grid is 100 *li* to the division. The coastal outline is relatively firm and the precision of the network of river systems extraordinary. The size of the original, which is now in the Pei Lin Museum at Sian, is about 3 feet square. The name of the geographer is not known. . . . Anyone who compares this map with the contemporary productions of European religious cosmography cannot but be amazed at the extent to which Chinese geography was at that time ahead of the West. . . . There was nothing like it in Europe till the Escorial MS. map of about +1550. . . .<sup>4</sup>

<sup>3</sup>Data maps are usually described as “thematic maps” in cartography. For a thorough account, see Arthur H. Robinson, *Early Thematic Mapping in the History of Cartography* (Chicago, 1982). On the history of statistical graphics, see H. Gray Funkhouser, “Historical Development of the Graphical Representation of Statistical Data,” *Osiris*, 3 (November 1937), 269–404; and James R. Beniger and Dorothy L. Robyn, “Quantitative Graphics in Statistics: A Brief History,” *American Statistician*, 32 (February 1978), 1–11.

<sup>4</sup>Joseph Needham, *Science and Civilisation in China* (Cambridge, 1959), vol. 3, 546–547.

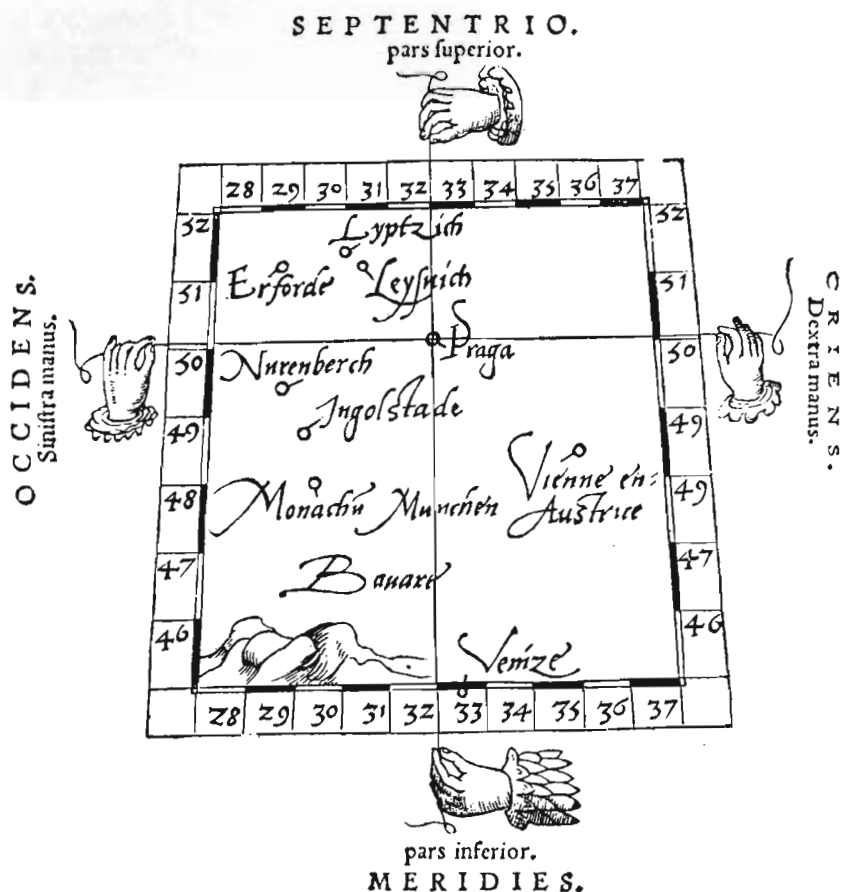




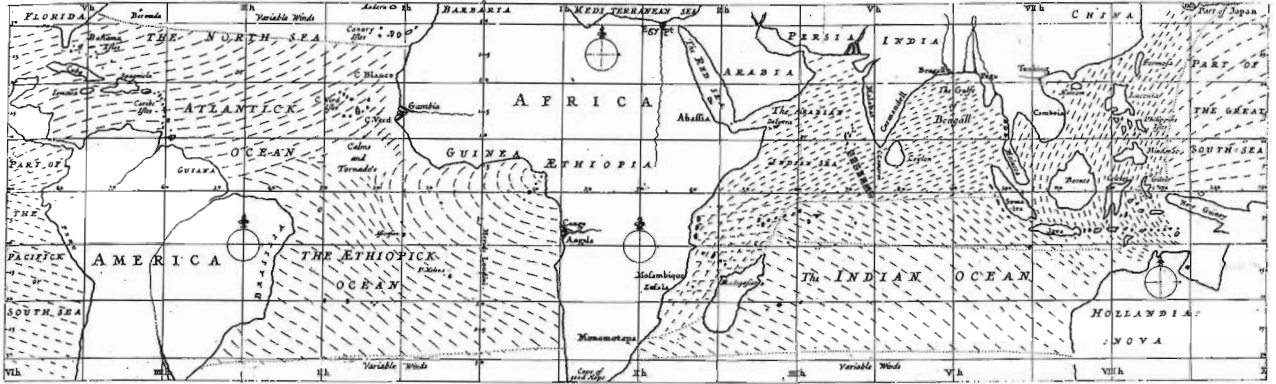
E. Chavannes, "Les Deux Plus Anciens  
 Spécimens de la Cartographie Chinoise,"  
*Bulletin de l'École Française de l'Extrême  
 Orient*, 3 (1903), 1-35, Carte B.

# Ecce formulam, vsum, atque

structuram Tabularum Ptolomæ, cum quibusdam locis, in quibus studiosus Geographiæ se satis exercere potest.

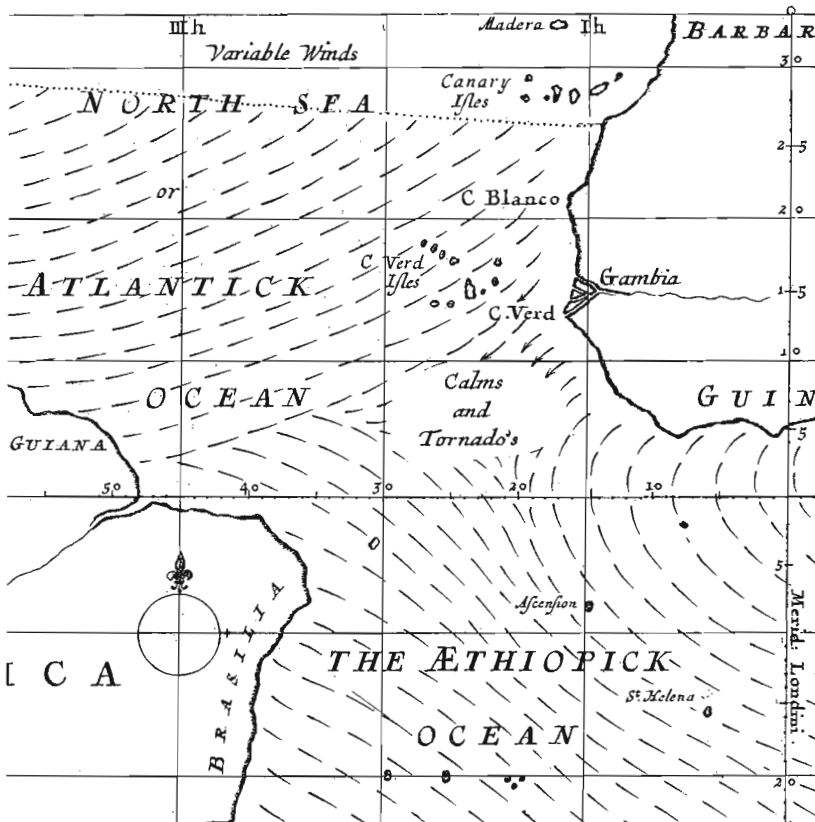


The 1546 edition of *Cosmographia* by Petrus Apianus contained examples of map design that show how very close European cartography by that time had come to achieving statistical graphicacy, even approaching the bivariate scatterplot. But, according to the historical record, no one had yet made the quantitative abstraction of placing a measured quantity on the map's surface at the intersection of the two threads instead of the name of a city, let alone the more difficult abstraction of replacing latitude and longitude with some other dimensions, such as time and money. Indeed, it was not until 1786 that the first economic time-series was plotted.



One of the first data maps was Edmond Halley's 1686 chart showing trade winds and monsoons on a world map.<sup>5</sup> The detailed section below shows the cartographic symbolization; with, as Halley wrote, "... the sharp end of each little stroak pointing out that part of the Horizon, from whence the wind continually comes; and where there are Monsoons the rows of stroaks run alternately backwards and forwards, by which means they are thicker [denser] than elsewhere."

<sup>5</sup>Norman J. W. Thrower, "Edmond Halley as a Thematic Geo-Cartographer," *Annals of the Association of American Geographers*, 59 (December 1969), 652-676.



Edmond Halley, "An Historical Account of the Trade Winds, and Monsoons, Observable in the Seas Between and Near the Tropicks; With an Attempt to Assign the Physical Cause of Said Winds," *Philosophical Transactions*, 183 (1686), 153-168.

An early and most worthy use of a map to chart patterns of disease was the famous dot map of Dr. John Snow, who plotted the location of deaths from cholera in central London for September 1854. Deaths were marked by dots and, in addition, the area's eleven water pumps were located by crosses. Examining the scatter over the surface of the map, Snow observed that cholera occurred almost entirely among those who lived near (and drank from) the Broad Street water pump. He had the handle of the contaminated pump removed, ending the neighborhood epidemic which had taken more than 500 lives.<sup>6</sup> The pump is located at the center of the map, just to the right of the D in BROAD STREET. Of course the link between the pump and the disease might have been revealed by computation and analysis without graphics, with some good luck and hard work. But, here at least, graphical analysis testifies about the data far more efficiently than calculation.

<sup>6</sup> E. W. Gilbert, "Pioneer Maps of Health and Disease in England," *Geographical Journal*, 124 (1958), 172-183. Shown here is a redrawing of John Snow's map. For a reproduction and detailed analysis of the original map, see Edward Tufte, *Visual Explanations: Images and Quantities, Evidence and Narrative* (Cheshire, Connecticut, 1997), Chapter 2. Ideally, see John Snow, *On the Mode of Communication of Cholera* (London, 1855).



Charles Joseph Minard gave quantity as well as direction to the data measures located on the world map in his portrayal of the 1864 exports of French wine:

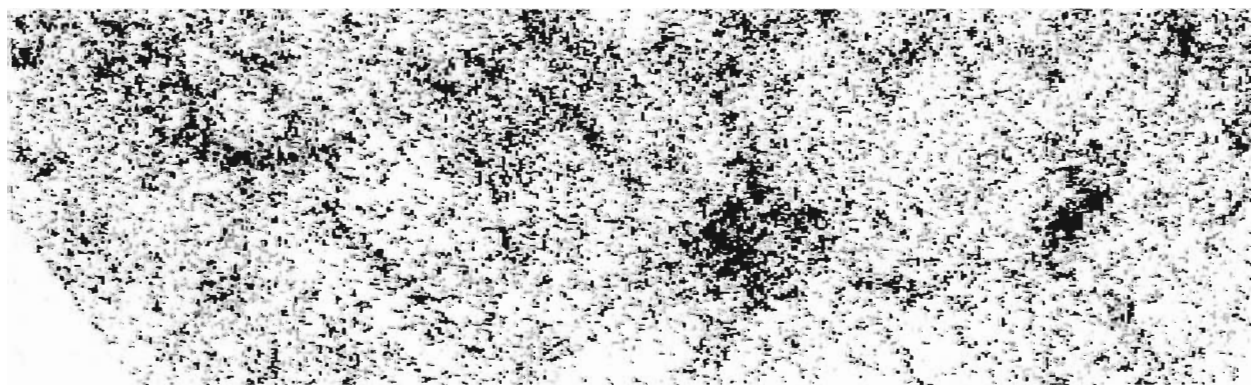




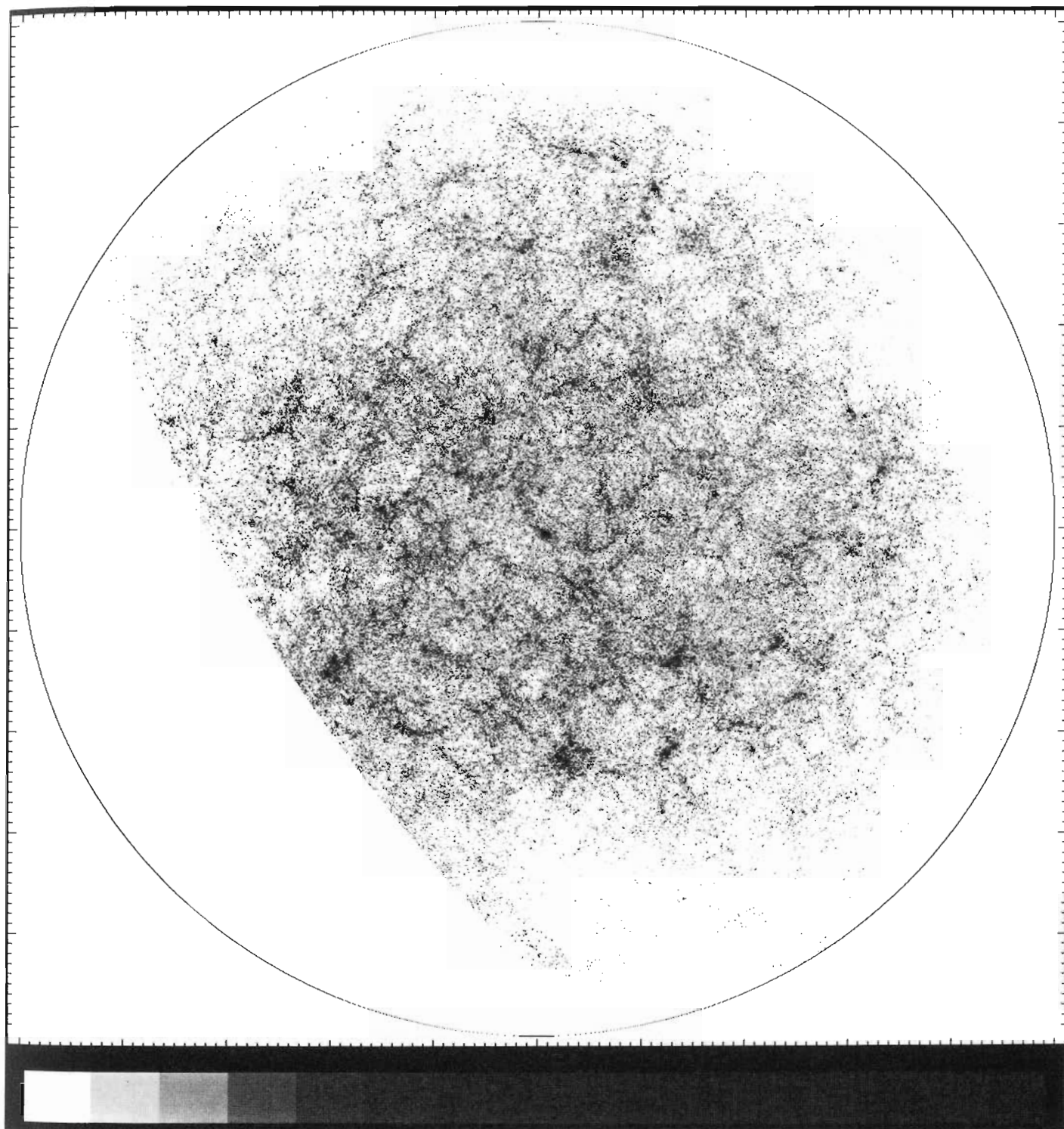
Computerized cartography and modern photographic techniques have increased the density of information some 5,000-fold in the best of current data maps compared to Halley's pioneering effort. This map shows the distribution of 1.3 million galaxies (including some overlaps) in the northern galactic hemisphere. The map divides the sky into  $1,024 \times 2,222$  rectangles. The number of galaxies counted in each of the 2,275,328 rectangles is represented by ten gray tones; the darker the tone, the greater the number of galaxies counted. The north galactic pole is at the center. The sharp edge on the left results from the earth blocking the view from the observatory. In the area near the perimeter of the map, the view is obscured by the interstellar dust of the galaxy in which we live (the Milky Way) as the line of sight passes through the flattened disk of our galaxy. The curious texture of local clusters of galaxies seen in this truly new view of the universe was not anticipated by students of galaxies, who had, of course, microscopically examined millions of photographs of galaxies before seeing this macroscopic view. Although the clusters are clearly evident (and accounted for by a theory of galactic origins), the seemingly random filaments may be happenstance. The producers of the map note the "strong temptation to conclude that the galaxies are arranged in a remarkable filamentary pattern on scales of approximately  $5^\circ$  to  $15^\circ$ , but we caution that this visual impression may be misleading because the eye tends to pick out linear patterns even in random noise. Indeed, roughly similar patterns are seen on maps constructed from simulated catalogs where no linear structure has been built in. . . ."<sup>7</sup>

<sup>7</sup> Michael Seldner, B. H. Siebers, Edward J. Groth and P. James E. Peebles, "New Reduction of the Lick Catalog of Galaxies," *Astronomical Journal*, 82 (April 1977), 249-314. See Gillian R. Knapp, "Mining the Heavens: The Sloan Digital Sky Survey," *Sky & Telescope* (August 1997), 40-48; Margaret J. Geller and John P. Huchra, "Mapping the Universe," *Sky & Telescope* (August 1991), 134-139.

The most extensive data maps, such as the cancer atlas and the count of the galaxies, place millions of bits of information on a single page before our eyes. No other method for the display of statistical information is so powerful.







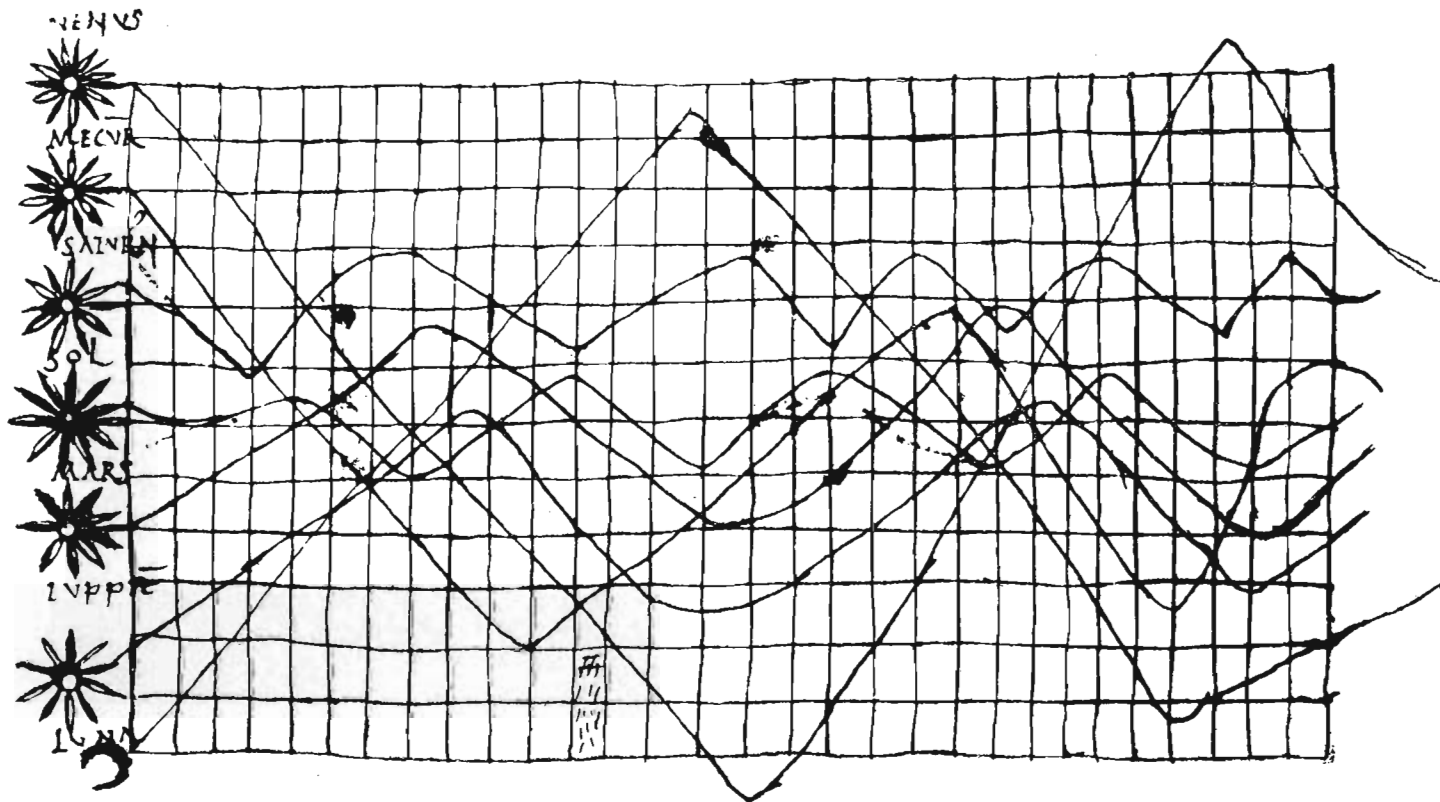
## Time-Series

The time-series plot is the most frequently used form of graphic design.<sup>8</sup> With one dimension marching along to the regular rhythm of seconds, minutes, hours, days, weeks, months, years, centuries, or millennia, the natural ordering of the time scale gives this design a strength and efficiency of interpretation found in no other graphic arrangement.

This reputed tenth- (or possibly eleventh-) century illustration of the inclinations of the planetary orbits as a function of time, apparently part of a text for monastery schools, is the oldest known example of an attempt to show changing values graphically. It appears as a mysterious and isolated wonder in the history of data graphics, since the next extant graphic of a plotted time-series shows up some 800 years later. According to Funkhouser, the astronomical content is confused and there are difficulties in reconciling the graph and its accompanying text with the actual movements of the planets. Particularly disconcerting is the wavy path ascribed to the sun.<sup>9</sup> An erasure and correction of a curve occur near the middle of the graph.

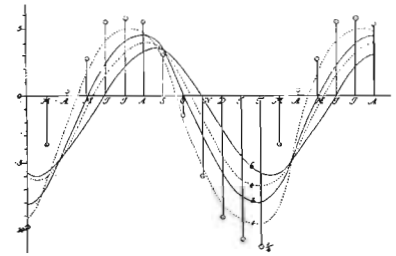
<sup>8</sup> A random sample of 4,000 graphics drawn from 15 of the world's newspapers and magazines published from 1974 to 1980 found that more than 75 percent of all the graphics published were time-series. Chapter 3 reports more on this.

<sup>9</sup> H. Gray Funkhouser, "A Note on a Tenth Century Graph," *Osiris*, 1 (January 1936), 260-262.

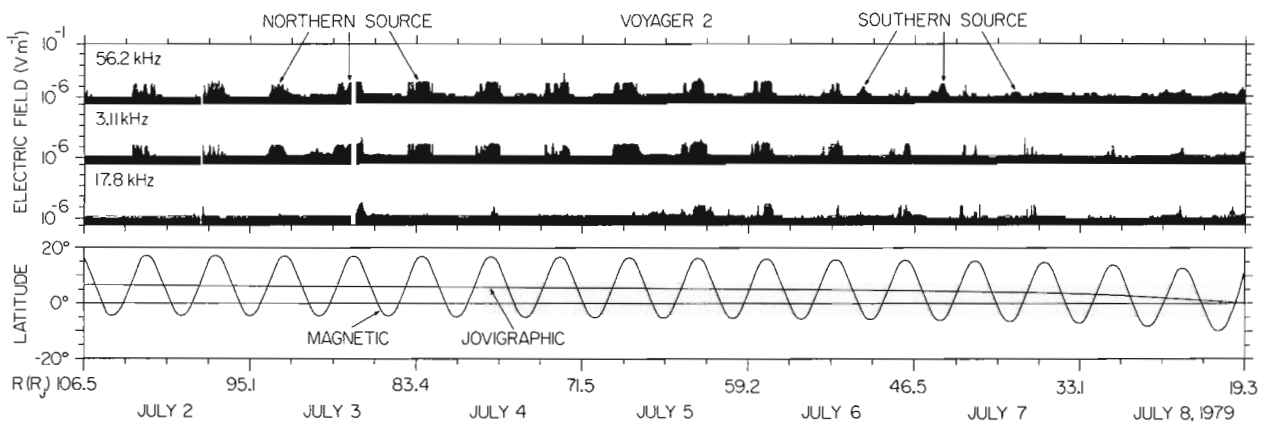




It was not until the late 1700s that time-series charts began to appear in scientific writings. This drawing of Johann Heinrich Lambert, one of a long series, shows the periodic variation in soil temperature in relation to the depth under the surface. The greater the depth, the greater the time-lag in temperature responsiveness. Modern graphic designs showing time-series periodicities differ little from those of Lambert, although the data bases are far larger.



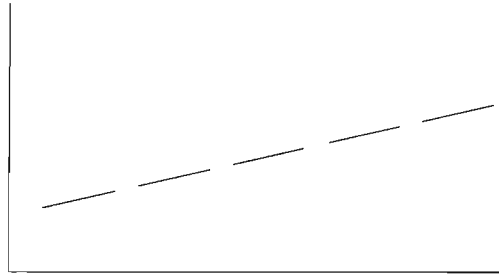
J. H. Lambert, *Pyrometrie* (Berlin, 1779).



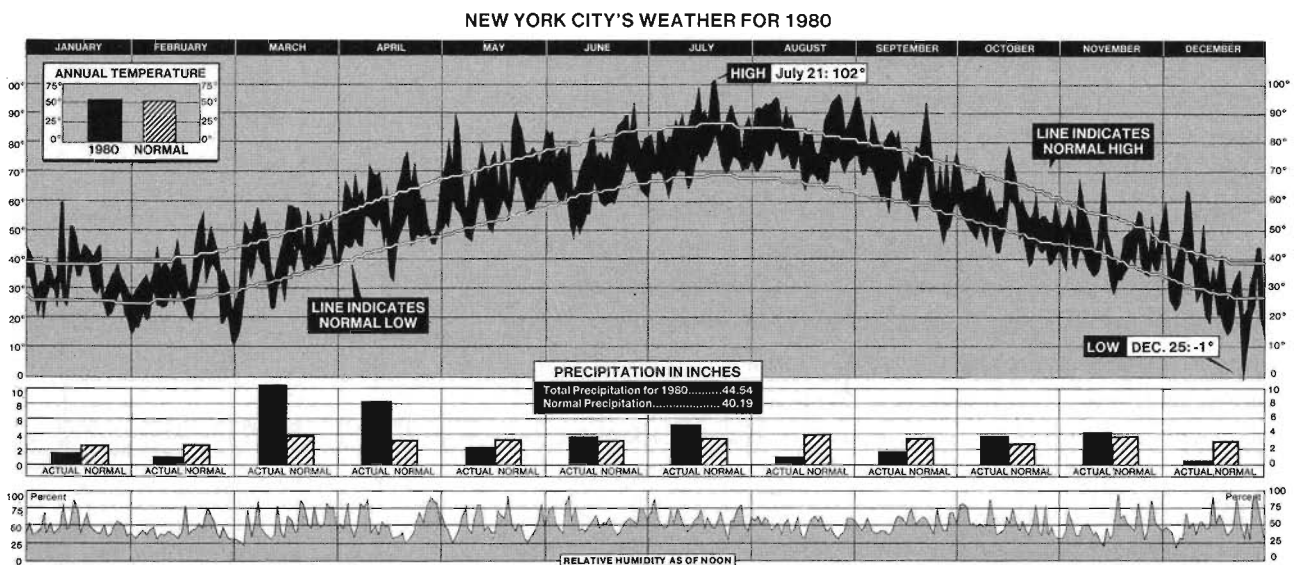
This plot of radio emissions from Jupiter is based on data collected by Voyager 2 in its pass close by the planet in July 1979. The radio intensity increases and decreases in a ten-hour cycle as Jupiter rotates. Maximum intensity occurs when the Jovian north magnetic pole is tipped toward the spacecraft, indicating a northern hemisphere source. A southern source was detected on July 7, as the spacecraft neared the equatorial plane. The horizontal scale shows the distance of the spacecraft from the planet measured in terms of Jupiter radii (R). Note the use of dual labels on the horizontal to indicate both the date and distance from Jupiter. The entire bottom panel also serves to label the horizontal scale, describing the changing orientation of the spacecraft relative to Jupiter as the planet is approached. The multiple time-series enforce not only comparisons within each series over time (as do all time-series plots) but also comparisons between the three different sampled radio bands shown. This richly multivariate display is based on 453,600 instrument samples of eight bits each. The resulting 3.6 million bits were reduced by peak and average processing to the 18,900 points actually plotted on the graphic.

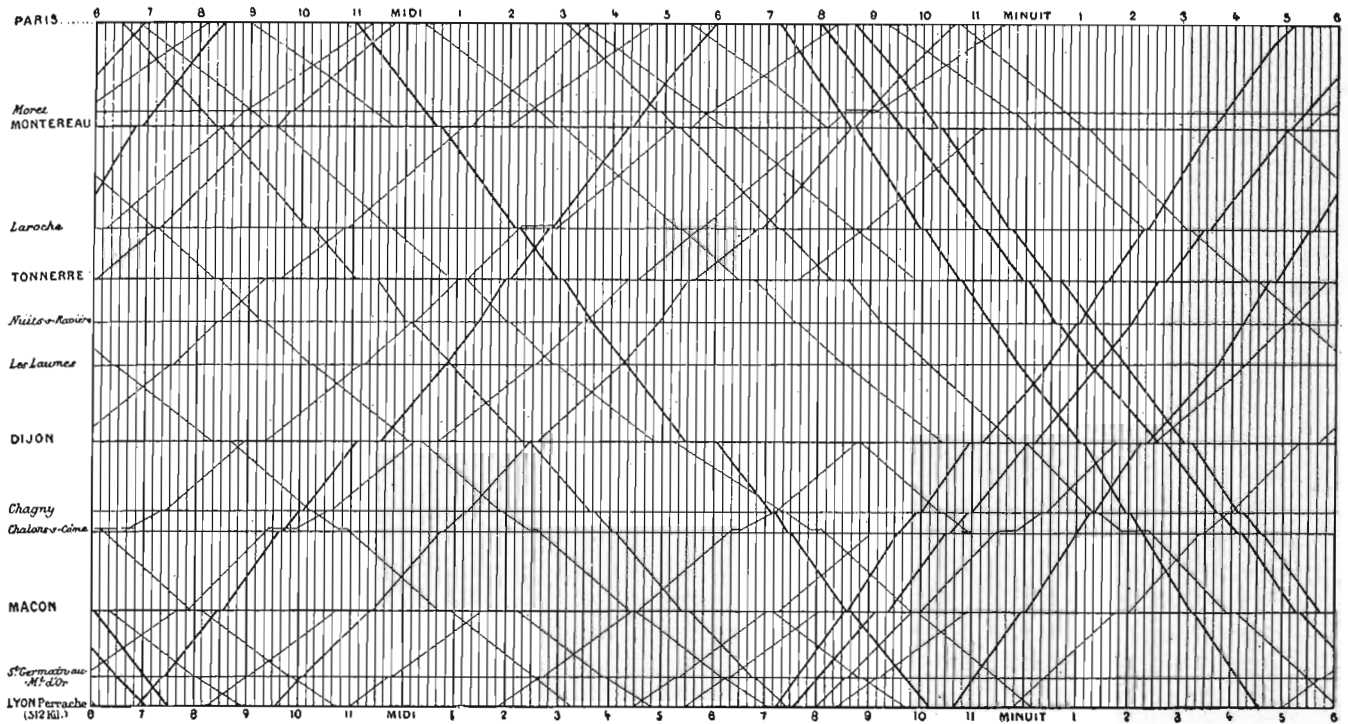
D. A. Gurnett, W. S. Kurth, and F. L. Scarf, "Plasma Wave Observations Near Jupiter: Initial Results from Voyager 2," *Science* 206 (November 23, 1979), 987-991; and letter from Donald A. Gurnett to Edward R. Tufte, June 27, 1980.

Time-series displays are at their best for big data sets with real variability. Why waste the power of data graphics on simple linear changes,



which can usually be better summarized in one or two numbers? Instead, graphics should be reserved for the richer, more complex, more difficult statistical material. This New York City weather summary for 1980 depicts 1,888 numbers. The daily high and low temperatures are shown in relation to the long-run average. The path of the normal temperatures also provides a forecast of expected change over the year; in the middle of February, for instance, New York City residents can look forward to warming at the rate of about 1.5 degrees per week all the way to July, the yearly peak. This distinguished graphic successfully organizes a large collection of numbers, makes comparisons between different parts of the data, and tells a story.





E. J. Marey, *La méthode graphique* (Paris, 1885), p. 20. The method is attributed to the French engineer, Ibry.

A design with similar strengths is Marey's graphical train schedule for Paris to Lyon in the 1880s. Arrivals and departures from a station are located along the horizontal; length of stop at a station is indicated by the length of the horizontal line. The stations are separated in proportion to their actual distance apart. The slope of the line reflects the speed of the train: the more nearly vertical the line, the faster the train. The intersection of two lines locates the time and place that trains going in opposite directions pass each other.

In 1981 a new express train from Paris to Lyon cut the trip to under three hours, compared to more than nine hours when Marey published the graphical train schedule. The path of the modern TGV (*train à grande vitesse*) is shown, overlaid on the schedule of 100 years before:

