

- The “outer fences” are located three hinge-spreads beyond the hinges:²¹

$$OF_L = H_L - 3 \times H\text{-spread}$$

$$OF_U = H_U + 3 \times H\text{-spread}$$

Observations beyond the outer fences are termed “far outside” and are represented by filled circles. There are no far-outside observations in the infant mortality data.

- The “whisker” growing from each end of the central box extends either to the extreme observation on its side of the distribution (as at the low end of the infant mortality data) or to the most extreme nonoutlying observation, called the “adjacent value” (as at the high end of the infant mortality distribution).²²

The boxplot of infant mortality in Figure 3.11 clearly reveals the skewness of the distribution: The lower whisker is much shorter than the upper whisker; the median is closer to the lower hinge than to the upper hinge; and there are several outside observations at the upper end of the infant mortality distribution, but none at the lower end. The apparent bimodality of the infant mortality data is not captured by the boxplot, however.

There are many useful univariate displays, including the traditional histogram. The stem-and-leaf plot is a modern variant of the histogram for small data sets, constructed directly from numerical data. Nonparametric density estimation may be employed to smooth a histogram. Quantile comparison plots are useful for comparing data with a theoretical probability distribution. Boxplots summarize some of the most important characteristics of a distribution, including center, spread, skewness, and outliers.

3.2 Plotting Bivariate Data

The *scatterplot*—a direct geometric representation of observations on two quantitative variables (generically, Y and X)—is the most useful of all statistical graphs. The scatterplot is a natural representation of data partly because the media on which we draw plots—paper, computer screens—are intrinsically two dimensional. Scatterplots are as familiar and essentially simple as they are useful; I will therefore limit this presentation to a few points. There are many examples of bivariate scatterplots in this book, including in the preceding chapter.

- In analyzing data, it is convenient to work in a computing environment that permits the interactive identification of observations in a scatterplot.
- Because relationships between variables in the social sciences are often weak, scatterplots can be dominated visually by “noise.” It often helps, therefore, to plot a nonparametric regression of Y on X .

²¹Here is a rough justification for the fences: In a normal population, the hinge-spread is 1.349 standard deviations, and so $1.5 \times H\text{-spread} = 1.5 \times 1.349 \times \sigma \approx 2\sigma$. The hinges are located $1.349/2 \approx 0.7$ standard deviations above and below the mean. The inner fences are, therefore, approximately at $\mu \pm 2.7\sigma$, and the outer fences at $\mu \pm 4.7\sigma$. From the standard normal table, $\Pr(Z > 2.7) \approx .003$, so we expect slightly less than 1% of the observations beyond the inner fences ($2 \times .003 = .006$); likewise, because $\Pr(Z > 4.7) \approx 1.3 \times 10^{-6}$, we expect less than one observation in 100,000 beyond the outer fences.

²²All of the folksy terminology—“hinges,” “fences,” “whiskers,” and so on—originates with Tukey (1977).

150
100
50
0

Infant Mortality Rate (per 1000)

Figure 3.13 Scatterplot for lowest smooth levels of GDP

- Scatterplots in which the bulk of the data is concentrated in a small region, as in Figure 3.13. It often helps to jitter the data between Y and X .²³
- Scatterplots in which the data points are widely scattered, as in Figure 3.14. In this instance of this phenomenon, the vocabulary test again. Surveys conducted in the United States and include in total 100 discrete—is to focus on a single variable, for example, can be discussed below). A random quantity to a random variable on a Paradoxically, the test randomly “jittered”

The bivariate scatterplot of two quantitative variables. The nonparametric regression line is shown. Scatterplots of the residuals are also shown, jittering the data.

²³See Chapter 4.

²⁴The idea of jittering a scatterplot

e hinges:²¹

e” and are represented by
unt mortality data.

ends either to the extreme
f the infant mortality data)
“adjacent value” (as at the

ewness of the distribution:
is closer to the lower hinge
the upper end of the infant
mortality of the infant mortality

l histogram. The stem-
ts, constructed directly
employed to smooth a
data with a theoretical
important characteristics

n two quantitative variables
. The scatterplot is a natu-
aw plots—paper, computer
ar and essentially simple as
. There are many examples
ter.

environment that permits the

are often weak, scatterplots
re, to plot a nonparametric

id is 1.349 standard deviations, and
0.7 standard deviations above and
the outer fences at $\mu \pm 4.7\sigma$. From
the observations beyond the inner
ss than one observation in 100,000

es with Tukey (1977).

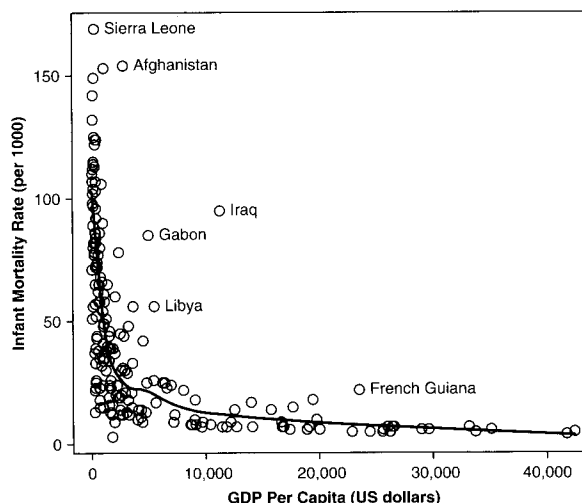


Figure 3.13 Scatterplot for infant mortality and GDP per capita for 193 nations. The line is for a lowess smooth with a span of 1/2. Several nations with high infant mortality for their levels of GDP are identified.

- Scatterplots in which one or both variables are highly skewed are difficult to examine, because the bulk of the data congregate in a small part of the display. Consider, for example, the scatterplot for infant mortality and gross domestic product (GDP) per capita in Figure 3.13. It often helps to “correct” substantial skewness prior to examining the relationship between Y and X .²³
- Scatterplots in which the variables are discrete can also be difficult to examine. An extreme instance of this phenomenon is shown in Figure 3.14, which plots scores on a 10-item vocabulary test against years of education. The data are from 16 of the U.S. General Social Surveys conducted by the National Opinion Research Center between 1974 and 2004, and include in total 21,638 observations. One solution—especially useful when only X is discrete—is to focus on the conditional distribution of Y for each value of X . Boxplots, for example, can be employed to represent the conditional distributions (see Figure 3.16, discussed below). Another solution is to separate overlapping points by adding a small random quantity to the discrete scores. In Figure 3.15, for example, I have added a uniform random variable on the interval $[-0.4, +0.4]$ to each value of vocabulary and education. Paradoxically, the tendency for vocabulary to increase with education is much clearer in the randomly “jittered” display.²⁴

The bivariate scatterplot is a natural graphical display of the relationship between two quantitative variables. Interpretation of a scatterplot can often be assisted by graphing a nonparametric regression, which summarizes the relationship between the two variables. Scatterplots of the relationship between discrete variables can be enhanced by randomly jittering the data.

²³See Chapter 4.

²⁴The idea of jittering a scatterplot, as well as the terminology, is due to Cleveland (1994).

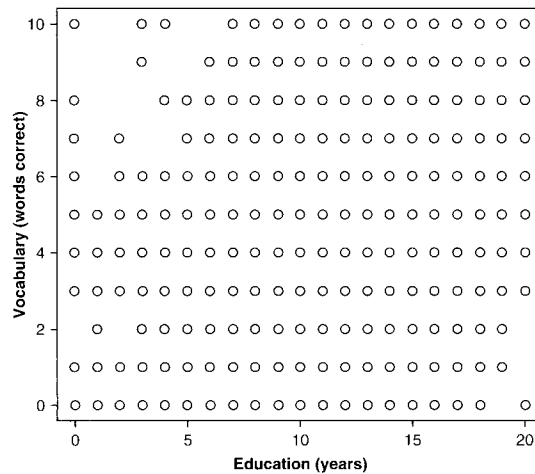


Figure 3.14 Scatterplot of scores on a 10-item vocabulary test versus years of education. Although there are nearly 22,000 observations in the data set, most of the plotted points fall on top of one another.

SOURCE: National Opinion Research Center (2005).

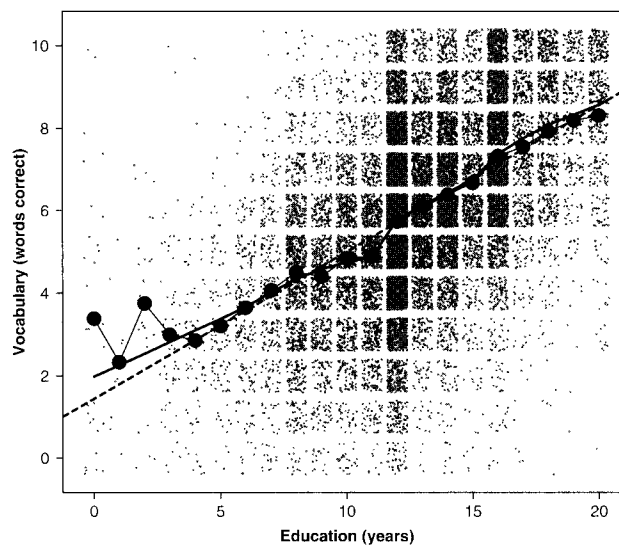


Figure 3.15 Jittered scatterplot for vocabulary score versus years of education. A uniformly distributed random quantity between -0.4 and $+0.4$ was added to each score for both variables. The heavier solid line is for a lowess fit to the data, with a span of 0.2 ; the broken line is the linear least-squares fit; the conditional means for vocabulary given education are represented by the dots, connected by the lighter solid line.

As mentioned, when the explanatory variable is discrete, parallel boxplots can be used to display the conditional distributions of Y . One common case occurs when the explanatory variable is a qualitative/categorical variable. An example is shown in Figure 3.16, using data collected by Michael Ornstein (1976) on interlocking directorates among the 248 largest Canadian firms. The response

Figure 3.16 Number of interlocking directorates among the 248 largest Canadian firms. SOURCE: Personal communication with Michael Ornstein.

variable in this graph is the number of interlocking directorates by each firm with other firms in the sample. The corporation is controlled by each firm with other firms in the sample. It is apparent from the graph that the conditional distribution of the number of interlocking directorates is different for foreign and Canadian corporations and the United States. The conditional distribution of the number of interlocking directorates is different for the conditional distribution of the number of interlocking directorates among firms than among U.K. firms.

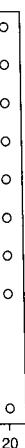
Parallel boxplots display the conditional distributions of Y for discrete (categorical) explanatory variables.

3.3 Plotting Multivariate Data

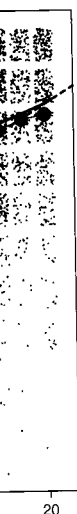
Because paper and computer graphics are limited, data is intrinsically difficult to visualize. A two-dimensional “point cloud” is a higher-dimensional case (see below). The essential character of the data is revealed on the basis of a statistical analysis.

²⁵We will revisit this example in Chapter 4, where we will discuss outliers in the plot.

²⁶We will apply these powerful techniques in Chapter 4.



years of education.
set, most of the plotted



education. A uniformly
s added to each score for
the data, with a span of 0.2;
nal means for vocabulary
by the lighter solid line.

xplots can be used to display
the explanatory variable is a
ing data collected by Michael
Canadian firms. The response

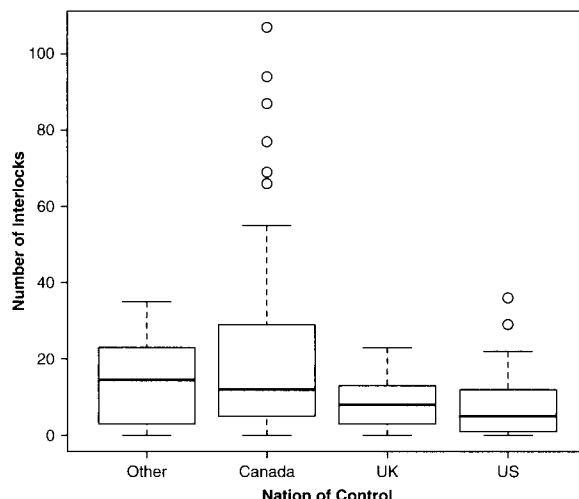


Figure 3.16 Number of interlocking directorate and executive positions by nation of control, for 248 dominant Canadian firms.

SOURCE: Personal communication from Michael Ornstein.

variable in this graph is the number of interlocking directorships and executive positions maintained by each firm with others in the group of 248. The explanatory variable is the nation in which the corporation is controlled, coded as Canada, United Kingdom, United States, and other foreign.

It is apparent from the graph that the average level of interlocking is greater among other-foreign and Canadian corporations than among corporations controlled in the United Kingdom and the United States. It is relatively difficult to discern detail in this display: first, because the conditional distributions of interlocks are positively skewed; and, second, because there is an association between level and spread—variation is also greater among other-foreign and Canadian firms than among U.K. and U.S. firms.²⁵

Parallel boxplots display the relationship between a quantitative response variable and a discrete (categorical or quantitative) explanatory variable.

3.3 Plotting Multivariate Data

Because paper and computer screens are two dimensional, graphical display of multivariate data is intrinsically difficult. Multivariate displays for quantitative data often project the higher-dimensional “point cloud” of the data onto a two-dimensional space. It is, of course, impossible to view a higher-dimensional scatterplot directly (but see the discussion of the three-dimensional case below). The essential trick of effective multidimensional display is to select projections that reveal important characteristics of the data. In certain circumstances, projections can be selected on the basis of a statistical model fit to the data or on the basis of explicitly stated criteria.²⁶

²⁵We will revisit this example in Section 4.4. Because the names of the firms are unavailable, I have not identified the outliers in the plot.

²⁶We will apply these powerful ideas in Chapters 11 and 12.

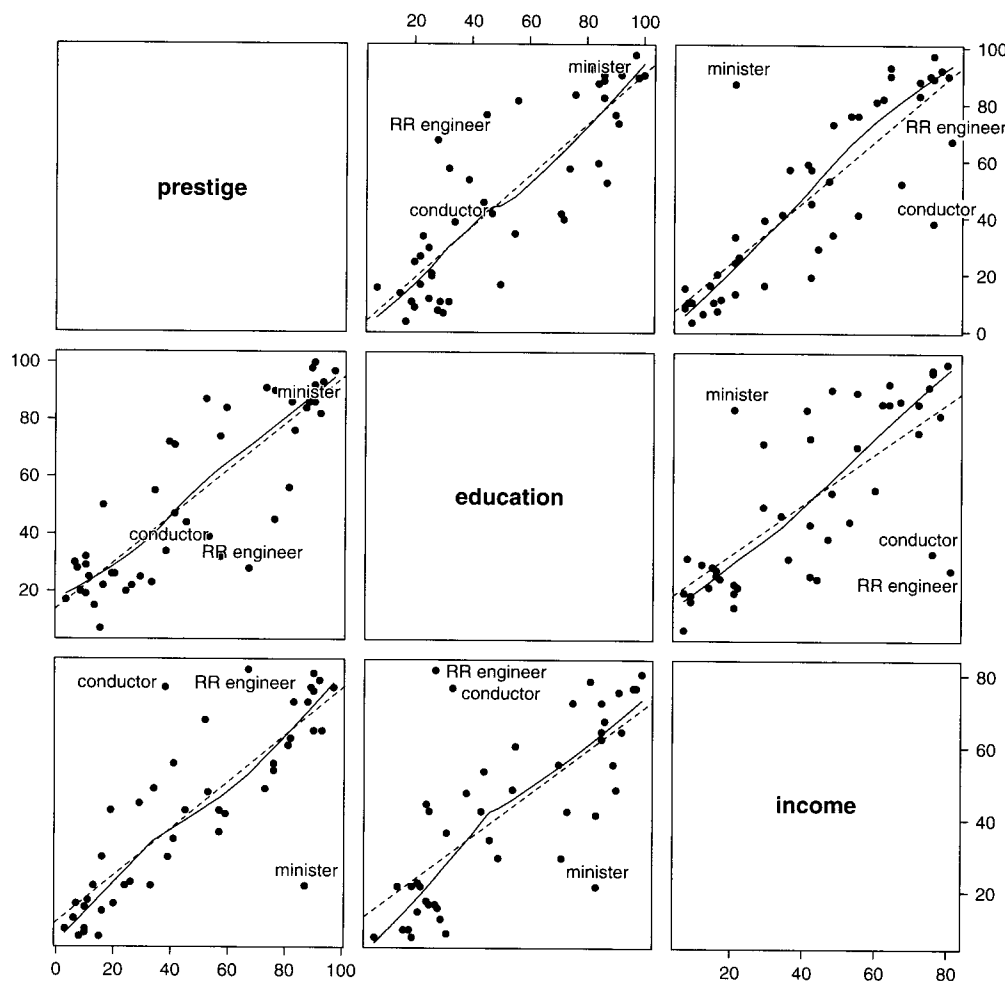


Figure 3.17 Scatterplot matrix for occupational prestige, level of education, and level of income, for 45 U.S. occupations in 1950. The least-squares regression line (broken line) and lowess smooth (for a span of 0.6, solid line) are shown on each plot. Three unusual observations are identified.

SOURCE: Duncan (1961).

3.3.1 Scatterplot Matrices

A simple approach to multivariate data, which does not require a statistical model, is to examine bivariate scatterplots for all pairs of variables. Arraying these plots in a *scatterplot matrix* produces a graphical analog to the correlation matrix.

An illustrative scatterplot matrix, for data on the prestige, education, and income levels of 45 U.S. occupations, appears in Figure 3.17. In this data set, first analyzed by Duncan (1961), “prestige” represents the percentage of respondents in a survey who rated an occupation as “good” or “excellent” in prestige; “education” represents the percentage of incumbents in the occupation in the 1950 U.S. Census who were high-school graduates; and “income” represents the percentage of occupational incumbents who earned incomes in excess of \$3500. Duncan’s purpose was to use a regression analysis of prestige on income and education to predict the prestige levels of

3.3. Plotting Multivariate Data

other occupations, for which there were no direct prestige ratings.

The variable names on the columns of the display: For the first column, the variable is “prestige”; the horizontal axis is “education”. Thus, the scatterplot in the first row shows the relationship between education (on the horizontal axis) and prestige (on the vertical axis).

It is important to understand that when analyzing multivariate data, the plot focuses on the *marginal* relationship between the object of data analysis for scatterplots (between pairs of variables). For example, the relationship between prestige and education ignores the relationship between prestige and income.

The response variable Y is prestige. The partial relationship between prestige and education, controlling for income, is a partial association between Y and X themselves are nonlinear. The relationship can be nonlinear even when the relationship between Y and X is linear.

Despite this intrinsic limitation, scatterplots are indeed the best way to visualize multivariate data, and this is indeed the case for the data on prestige, education, and income. *Ministers* have relatively high prestige for their level of education. *RR engineers* have relatively high income for their level of education. *Conductors* also have relatively high income for their level of education. This is all for the least-squares line.

3.3.2 Coded Scatterplots

Information about a categorical variable can be added to the plotting symbols. The number of degrees of fill, distinguishing between different categories.

Figure 3.18 shows a scatterplot of prestige versus education for men. The data are displayed in 12 categories, which, recall, represent different levels of prestige. The data are near the line $Y = X$; it is also as one would expect, and the data are heavier than everyone else.

3.3.3 Three-Dimensional Scatterplots

Another useful multivariate plot is the *three-dimensional scatterplot*. Most

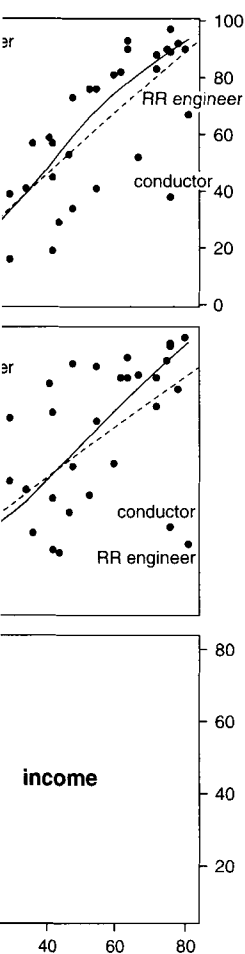
²⁷We will return to this regression analysis.

²⁸These ideas are explored in Chapter 4.

²⁹See the discussion of Duncan’s data in Chapter 4.

³⁰See Spence and Lewandowsky (1968) on coded scatterplots.

³¹Davis’s data were introduced in Chapter 4.



ion, and level of income,
on line (broken line) and
each plot. Three unusual

tical model, is to examine
scatterplot matrix produces

on, and income levels of
lyzed by Duncan (1961),
d an occupation as “good”
umbents in the occupation
represents the percentage
Duncan’s purpose was to
dict the prestige levels of

other occupations, for which data on income and education were available, but for which there were no direct prestige ratings.²⁷

The variable names on the diagonal of the scatterplot matrix in Figure 3.17 label the rows and columns of the display: For example, the vertical axis for the two plots in the first row of the display is “prestige”; the horizontal axis for the two plots in the second column is “education.” Thus, the scatterplot in the first row, second column is for prestige (on the vertical axis) versus education (on the horizontal axis).

It is important to understand an essential limitation of the scatterplot matrix as a device for analyzing multivariate data: By projecting the multidimensional point cloud onto pairs of axes, the plot focuses on the *marginal* relationships between the corresponding pairs of variables. The object of data analysis for several variables, however, is typically to investigate *partial* relationships (between pairs of variables, “controlling” statistically for other variables), not marginal associations. For example, in the Duncan data set, we are more interested in the partial relationship of prestige to education holding income constant than in the marginal relationship between prestige and education ignoring income.

The response variable Y can be related marginally to a particular X , even when there is no partial relationship between the two variables controlling for other X s. It is also possible for there to be a partial association between Y and an X but no marginal association. Furthermore, if the X s themselves are nonlinearly related, then the marginal relationship between Y and a specific X can be nonlinear even when their partial relationship is linear.²⁸

Despite this intrinsic limitation, scatterplot matrices often uncover interesting features of the data, and this is indeed the case in Figure 3.17, where the display reveals three unusual observations: *Ministers* have relatively low income for their relatively high level of education, and relatively high prestige for their relatively low income; *railroad conductors* and *railroad engineers* have relatively high incomes for their more-or-less average levels of education; *railroad conductors* also have relatively low prestige for their relatively high incomes. This pattern bodes ill for the least-squares linear regression of prestige on income and education.²⁹

3.3.2 Coded Scatterplots

Information about a categorical third variable can be entered on a bivariate scatterplot by coding the plotting symbols. The most effective codes use different colors to represent categories, but degrees of fill, distinguishable shapes, and distinguishable letters can also be effective.³⁰

Figure 3.18 shows a scatterplot of Davis’s data on measured and reported weight.³¹ Observations for men are displayed as Ms, for women as Fs. Except for the outlying point (number 12—which, recall, represents an error in the data), the points both for men and for women cluster near the line $Y = X$; it is also clear from the display that most men are heavier than most women, as one would expect, and that, discounting the bad data point, one man (number 21) is quite a bit heavier than everyone else.

3.3.3 Three-Dimensional Scatterplots

Another useful multivariate display, directly applicable to three variables at a time, is the *three-dimensional scatterplot*. Moreover, just as data can be projected onto a judiciously chosen plane

²⁷We will return to this regression problem in Chapter 5.

²⁸These ideas are explored in Chapter 12.

²⁹See the discussion of Duncan’s occupational-prestige regression in Chapter 11.

³⁰See Spence and Lewandowsky (1990) for a fine review of the literature on graphical perception, including information on coded scatterplots.

³¹Davis’s data were introduced in Chapter 2, where only the data for women were presented.

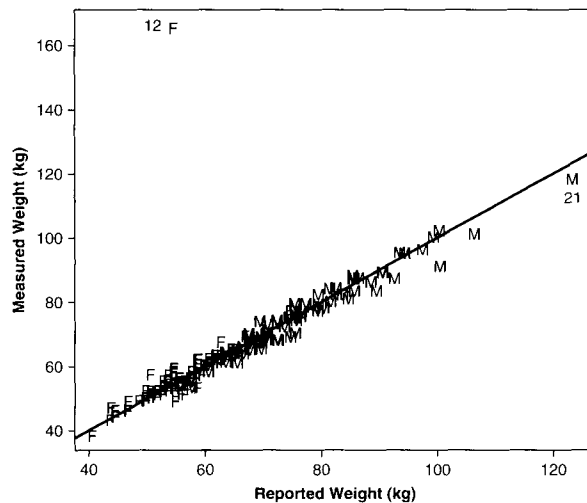


Figure 3.18 Davis's data on measured and reported weight, by gender. Data points for men are represented by Ms, for women by Fs, and are jittered slightly to reduce overplotting. The line on the graph is $Y = X$. In the combined data set for men and women, the outlying observation is number 12.

in a two-dimensional plot, higher-dimensional data can be projected onto a three-dimensional space, expanding the range of application of three-dimensional scatterplots.³²

Barring the use of a true stereoscopic display, the three-dimensional scatterplot is an illusion produced by modern statistical software: The graph represents a projection of a three-dimensional space onto a two-dimensional computer screen. Nevertheless, motion (e.g., rotation) and the ability to interact with the display—possibly combined with the effective use of perspective, color, depth cueing, and other visual devices—can produce a vivid impression of directly examining objects in three-dimensional space.

It is literally impossible to convey this impression adequately on the static, two-dimensional page of a book, but Figure 3.19 shows Duncan's prestige data rotated interactively into a revealing orientation: Looking down the cigar-shaped scatter of most of the data, the three unusual observations stand out very clearly.

3.3.4 Conditioning Plots

Conditioning plots (or *coplots*), described in Cleveland (1993), are another graphical device for examining multidimensional data. The essential idea of the coplot is to focus on the relationship between the response variable and a particular explanatory variable, dividing the data into groups based on the values of other explanatory variables—the *conditioning variables*. If the conditioning variables are discrete, then this division is straightforward and natural. If a conditioning variable is continuous, it can be binned: Cleveland suggests using overlapping bins, which are called “shingles.”

An illustrative coplot, for the General Social Survey vocabulary data, is shown in Figure 3.20. This graph displays the relationship between vocabulary score and education “controlling for”

³²For example, there are three-dimensional versions of the added-variable and component-plus-residual plots discussed in Chapters 11 and 12. See, for example, Cook and Weisberg (1989).

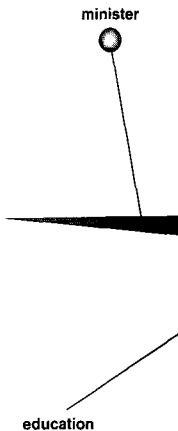


Figure 3.19 Three-dimensional scatterplot of Duncan's prestige data, rotated into an orientation that reveals the cigar-shaped scatter of the data. The least-squares regression plane is shown.

gender and the year of the survey. The three unusual observations are the three panels of the coplot; that is, the three panels show the relationship between vocabulary score and education, controlling for gender and the year of the survey. In a few panels, the lowest line of the scatter is quite sparse.

Although they can be effective, three-dimensional plots are two, or perhaps three, conditions away from revealing the relationship between the response variable and the explanatory variables. Second, because coplots are often used for large data sets, an issue that

Visualizing multivariate data is a challenge. Higher-dimensional scatterplots can be projected onto two or three dimensions, but the resulting dynamic three-dimensional plots are often difficult to interpret.

Summary

- Statistical graphs are central to data analysis and in statistical inference.



Data points for men are
tly to reduce overplotting.
or men and women, the

l onto a three-dimensional
plots.³²

al scatterplot is an illusion
tion of a three-dimensional
e.g., rotation) and the ability
of perspective, color, depth
directly examining objects

the static, two-dimensional
interactively into a reveal-
the data, the three unusual

another graphical device for
to focus on the relationship
dividing the data into groups
variables. If the conditioning
1. If a conditioning variable
ing bins, which are called

ata, is shown in Figure 3.20.
education “controlling for”

onent-plus-residual plots discussed

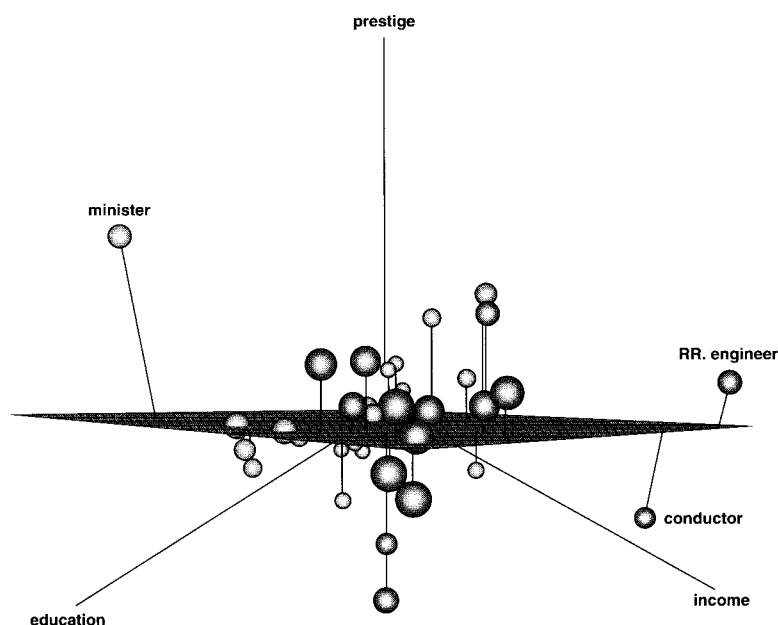


Figure 3.19 Three-dimensional scatterplot for Duncan's occupational prestige data, rotated into an orientation that reveals three unusual observations. From this orientation, the least-squares regression plane, also shown in the plot, is viewed nearly edge on.

gender and the year of the survey. The partial relationships are remarkably similar in the different panels of the coplot; that is, gender and year appear to make little difference to the relationship between vocabulary score and education. The relationships also appear to be very close to linear: In a few panels, the lowess line departs from the linear least-square line at the far left, but data in this region are quite sparse.

Although they can be effective graphs, coplots have limitations: First, if there are more than two, or perhaps three, conditioning variables, it becomes difficult to perceive how the partial relationship between the response and the focal explanatory variable changes with the conditioning variables. Second, because coplots require the division of the data into groups, they are most useful for large data sets, an issue that grows more acute as the number of conditioning variables increases.

Visualizing multivariate data is intrinsically difficult because we cannot directly examine higher-dimensional scatterplots. Effective displays project the higher-dimensional point cloud onto two or three dimensions; these displays include the scatterplot matrix, the dynamic three-dimensional scatterplot, and the conditioning plot.

Summary

- Statistical graphs are central to effective data analysis, both in the early stages of an investigation and in statistical modeling.

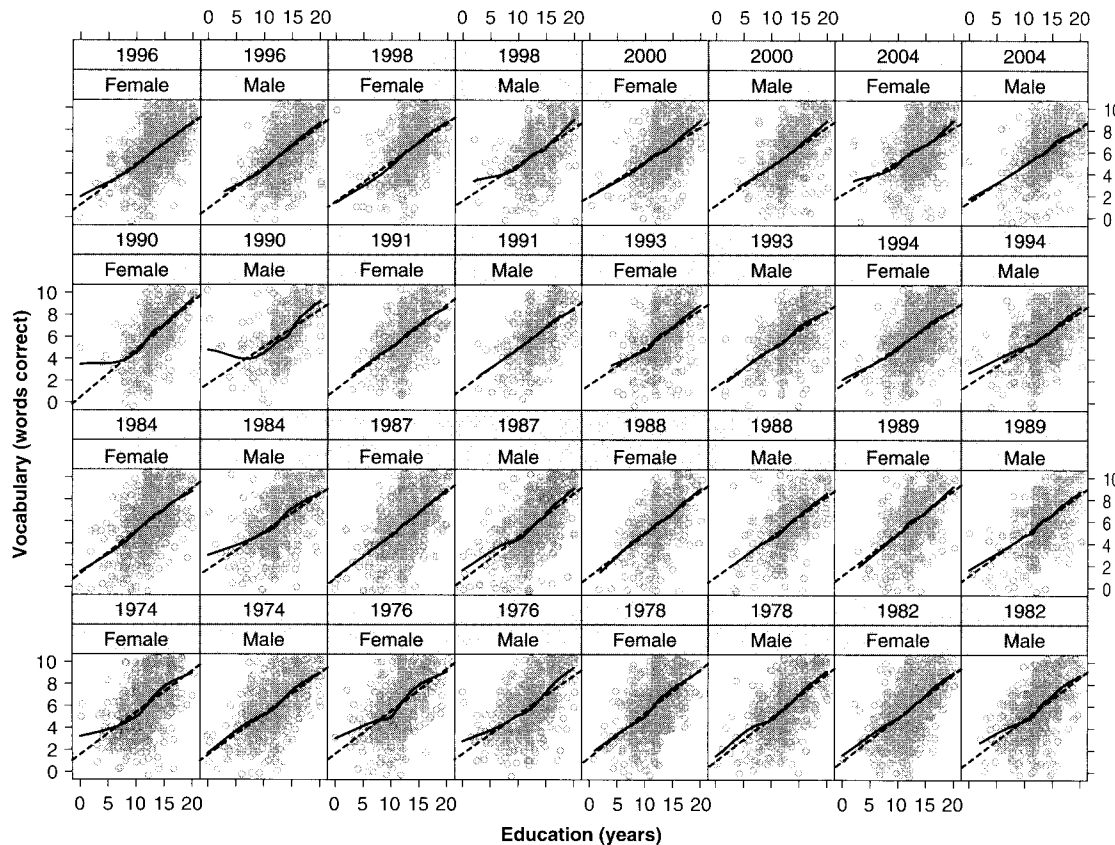


Figure 3.20 Coplot showing the relationship between vocabulary score and education controlling for year and gender. The points in each panel are jittered to reduce overplotting. The broken line shows the linear least-square fit, while the solid line gives the lowess fit for a span of 0.6.

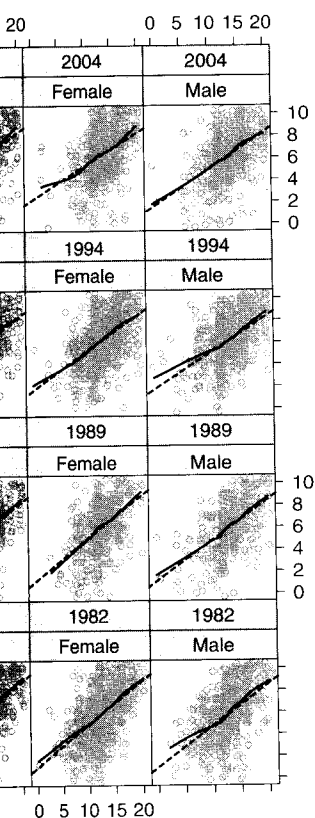
- There are many useful univariate displays, including the traditional histogram. The stem-and-leaf plot is a modern variant of the histogram for small data sets, constructed directly from numerical data. Nonparametric density estimation may be employed to smooth a histogram. Quantile-comparison plots are useful for comparing data with a theoretical probability distribution. Boxplots summarize some of the most important characteristics of a distribution, including center, spread, skewness, and outliers.
- The bivariate scatterplot is a natural graphical display of the relationship between two quantitative variables. Interpretation of a scatterplot can often be assisted by graphing a nonparametric regression, which summarizes the relationship between the two variables. Scatterplots of the relationship between discrete variables can be enhanced by randomly jittering the data.
- Parallel boxplots display the relationship between a quantitative response variable and a discrete explanatory variable.
- Visualizing multivariate data is intrinsically difficult because we cannot directly examine higher-dimensional scatterplots. Effective displays project the higher-dimensional point

cloud onto two or three dynamic three-dimension

Recommended Reading

The literature—especially the furnish only the briefest of bit

- Fox (2000c) presents a history of the subject. J. social scientists.
- Tufte's (1983) influential opinionated but well w graphics, broadly constr
- Modern interest in statist data analysis; unfortun (Tukey, 1977) difficult introduction to the topic exploratory data analys 1985).
- Tukey's influence made described in two access in Chambers, Cleveland start.
- Modern statistical grap The S statistical comp 1998; Chambers & Ha graphical capabilities. F ace. Cook and Weisber (Tierney, 1990) to prod a variety of statistical g of the methods describ ern statistical graphs u statistical computing e
- Atkinson (1985) prese as do Cook (1998) and



and education
 are jittered to reduce
 overlap, while the solid line

on a histogram. The stem-
 plots, constructed directly
 employed to smooth a his-
 togram with a theoretical prob-
 ability distribution, are im-
 portant characteristics of a

relationship between two
 variables. Graphing a
 scatterplot between the two variables
 can be enhanced by randomly

response variable and a

we cannot directly examine
 the higher-dimensional point

cloud onto two or three dimensions; these displays include the scatterplot matrix, the dynamic three-dimensional scatterplot, and the conditioning plot.

Recommended Reading

The literature—especially the recent literature—on statistical graphics is truly voluminous. I will furnish only the briefest of bibliographies:

- Fox (2000c) presents a brief overview of statistical graphics, including information on the history of the subject. Jacoby (1997, 1998) gives a more extended overview addressed to social scientists.
- Tufte's (1983) influential book on graphical presentation of quantitative information is opinionated but well worth reading. (Tufte has since published several other books on graphics, broadly construed, but I prefer his first book.)
- Modern interest in statistical graphics is the direct result of John Tukey's work on exploratory data analysis; unfortunately, Tukey's idiosyncratic writing style makes his seminal book (Tukey, 1977) difficult to read. Velleman and Hoaglin (1981) provide a more digestible introduction to the topic. There is interesting information on the statistical theory underlying exploratory data analysis in two volumes edited by Hoaglin, Mosteller, and Tukey (1983, 1985).
- Tukey's influence made Bell Labs a center of work on statistical graphics, much of which is described in two accessible and interesting books by William Cleveland (1993, 1994) and in Chambers, Cleveland, Kleiner, and Tukey (1983). Cleveland (1994) is a good place to start.
- Modern statistical graphics is closely associated with advances in statistical computing: The S statistical computing environment (Becker, Chambers, & Wilks, 1988; Chambers, 1998; Chambers & Hastie, 1992), also a product of Bell Labs, is particularly strong in its graphical capabilities. R, a free, open-source implementation of S, was mentioned in the preface. Cook and Weisberg (1994, 1999) use the Lisp-Stat statistical computing environment (Tierney, 1990) to produce an impressive statistical package, called Arc, which incorporates a variety of statistical graphics of particular relevance to regression analysis (including many of the methods described later in this text). Friendly (1991) describes how to construct modern statistical graphs using the SAS/Graph system. Brief presentations of these and other statistical computing environments appear in a book edited by Stine and Fox (1996).
- Atkinson (1985) presents a variety of innovative graphs in support of regression analysis, as do Cook (1998) and Cook and Weisberg (1994, 1999).