

Data: Types and Presentation

-
- 1.1 TYPES OF BIOLOGICAL DATA
 - 1.2 ACCURACY AND SIGNIFICANT FIGURES
 - 1.3 FREQUENCY DISTRIBUTIONS
 - 1.4 CUMULATIVE FREQUENCY DISTRIBUTIONS
-

Scientific study involves the systematic collection, organization, analysis, and presentation of knowledge. Many investigations in the biological sciences are quantitative, where knowledge is in the form of numerical observations called *data*. (One numerical observation is a *datum*.*) In order for the presentation and analysis of data to be valid and useful, we must use methods appropriate to the type of data obtained, to the design of the data collection, and to the questions asked of the data; and the limitations of the data, of the data collection, and of the data analysis should be appreciated when formulating conclusions. This chapter, and those that follow, will introduce many concepts relevant to this goal.

The word *statistics* is derived from the Latin for “state,” indicating the historical importance of governmental data gathering, which related principally to demographic information (including census data and “vital statistics”) and often to their use in military recruitment and tax collecting.[†]

The term *statistics* is often encountered as a synonym for *data*: One hears of college enrollment statistics (such as the numbers of newly admitted students, numbers of senior students, numbers of students from various geographic locations), statistics of a basketball game (such as how many points were scored by each player, how many fouls were committed), labor statistics (such as numbers of workers unemployed, numbers employed in various occupations), and so on. Hereafter, this use of the word *statistics* will not appear in this book. Instead, it will be used in its other common manner: to refer to the *orderly collection, analysis, and interpretation of data with a view to objective evaluation of conclusions based on the data*. (Section 2.4 will introduce another fundamentally important use of the term *statistic*.)

Statistics applied to biological problems is simply called *biostatistics* or, sometimes, *biometry*[‡] (the latter term literally meaning “biological measurement”). Although

*The term *data* is sometimes seen as a singular noun meaning “numerical information.” This book refrains from that use.

[†]Peters (1987: 79) and Walker (1929: 32) attribute the first use of the term *statistics* to a German professor, Gottfried Achenwall (1719–1772), who used the German word *Statistik* in 1749, and the first published use of the English word to John Sinclair (1754–1835) in 1791.

[‡]The word *biometry*, which literally means “biological measurement,” had, since the nineteenth century, been found in several contexts (such as demographics and, later, quantitative genetics; Armitage, 1985; Stigler, 2000), but using it to mean the application of statistical methods to biological information apparently was conceived between 1892 and 1901 by Karl Pearson, along with the name *Biometrika* for the still-important English journal he helped found; and it was first published in the inaugural issue of this journal in 1901 (Snedecor, 1954). The Biometrics Section of the American

their magnitudes relative to each other; or success in learning to run a maze may be recorded as *A*, *B*, or *C*.

It is often true that biological data expressed on the ordinal scale could have been expressed on the interval or ratio scale had exact measurements been obtained (or obtainable). Sometimes data that were originally on interval or ratio scales will be changed to ranks; for example, examination grades of 99, 85, 73, and 66% (ratio scale) might be recorded as *A*, *B*, *C*, and *D* (ordinal scale), respectively.

Ordinal-scale data contain and convey less information than ratio or interval data, for only relative magnitudes are known. Consequently, quantitative comparisons are impossible (e.g., we cannot speak of a grade of *C* being half as good as a grade of *A*, or of the difference between cell sizes 1 and 2 being the same as the difference between sizes 3 and 4). However, we will see that many useful statistical procedures are, in fact, applicable to ordinal data.

(d) Data in Nominal Categories. Sometimes the variable being studied is classified by some qualitative measure it possesses rather than by a numerical measurement. In such cases the variable may be called an *attribute*, and we are said to be dealing with *nominal*, or *categorical*, data. Genetic phenotypes are commonly encountered biological attributes: The possible manifestations of an animal's eye color might be brown or blue; and if human hair color were the attribute of interest, we might record black, brown, blond, or red. As other examples of nominal data (*nominal* is from the Latin word for "name"), people might be classified as male or female, or right-handed or left-handed. Or, plants might be classified as dead or alive, or as with or without fertilizer application. Taxonomic categories also form a nominal classification scheme (for example, plants in a study might be classified as pine, spruce, or fir).

Sometimes, data that might have been expressed on an ordinal, interval, or ratio scale of measurement may be recorded in nominal categories. For example, heights might be recorded as tall or short, or performance on an examination as pass or fail, where there is an arbitrary cut-off point on the measurement scale to separate tall from short and pass from fail.

As will be seen, statistical methods useful with ratio, interval, or ordinal data generally are not applicable to nominal data, and we must, therefore, be able to identify such situations when they occur.

(e) Continuous and Discrete Data. When we spoke previously of plant heights, we were dealing with a variable that could be any conceivable value within any observed range; this is referred to as a *continuous variable*. That is, if we measure a height of 35 cm and a height of 36 cm, an infinite number of heights is possible in the range from 35 to 36 cm: a plant might be 35.07 cm tall or 35.988 cm tall, or 35.3263 cm tall, and so on, although, of course, we do not have devices sensitive enough to detect this infinity of heights. A continuous variable is one for which there is a possible value between any other two values.

However, when speaking of the number of leaves on a plant, we are dealing with a variable that can take on only certain values. It might be possible to observe 27 leaves, or 28 leaves, but 27.43 leaves and 27.9 leaves are values of the variable that are impossible to obtain. Such a variable is termed a *discrete* or *discontinuous variable* (also known as a *meristic variable*). The number of white blood cells in 1 mm³ of blood, the number of giraffes visiting a water hole, and the number of eggs laid by a grasshopper are all discrete variables. The possible values of a discrete variable generally are consecutive integers, but this is not necessarily so. If the leaves on our

plants are always formed in pairs, then only even integers are possible values of the variable. And the ratio of number of wings to number of legs of insects is a discrete variable that may only have the value of 0, 0.3333 . . . , or 0.6666 . . . (i.e., $\frac{0}{6}$, $\frac{2}{6}$, or $\frac{4}{6}$, respectively).*

Ratio-, interval-, and ordinal-scale data may be either continuous or discrete. Nominal-scale data by their nature are discrete.

1.2 ACCURACY AND SIGNIFICANT FIGURES

Accuracy is the nearness of a measurement to the true value of the variable being measured. *Precision* is not a synonymous term but refers to the closeness to each other of repeated measurements of the same quantity. Figure 1.1 illustrates the difference between accuracy and precision of measurements.

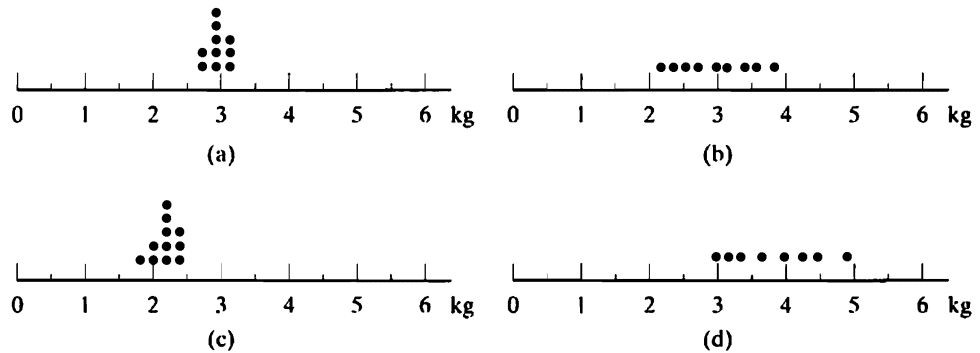


FIGURE 1.1: Accuracy and precision of measurements. A 3-kilogram animal is weighed 10 times. The 10 measurements shown in sample (a) are relatively accurate and precise; those in sample (b) are relatively accurate but not precise; those of sample (c) are relatively precise but not accurate; and those of sample (d) are relatively inaccurate and imprecise.

Human error may exist in the recording of data. For example, a person may miscount the number of birds in a tract of land or misread the numbers on a heart-rate monitor. Or, a person might obtain correct data but record them in such a way (perhaps with poor handwriting) that a subsequent data analyst makes an error in reading them. We shall assume that such errors have not occurred, but there are other aspects of accuracy that should be considered.

Accuracy of measurement can be expressed in numerical reporting. If we report that the hind leg of a frog is 8 cm long, we are stating the number 8 (a value of a continuous variable) as an estimate of the frog's true leg length. This estimate was made using some sort of a measuring device. Had the device been capable of more accuracy, we might have declared that the leg was 8.3 cm long, or perhaps 8.32 cm long. When recording values of continuous variables, it is important to designate the accuracy with which the measurements have been made. By convention, the value 8 denotes a measurement in the range of 7.50000 . . . to 8.49999 . . . , the value 8.3 designates a range of 8.25000 . . . to 8.34999 . . . , and the value 8.32 implies that the true value lies within the range of 8.31500 . . . to 8.32499 That is, the reported value is the midpoint of the implied range, and the size of this range is designated by the last decimal place in the measurement. The value of 8 cm implies an ability to

*The ellipsis marks (...) may be read as "and so on." Here, they indicate that $\frac{2}{6}$ and $\frac{4}{6}$ are repeating decimal fractions, which could just as well have been written as 0.333333333333 . . . and 0.666666666666 . . . , respectively.

determine length within a range of 1 cm, 8.3 cm implies a range of 0.1 cm, and 8.32 cm implies a range of 0.01 cm. Thus, to record a value of 8.0 implies greater accuracy of measurement than does the recording of a value of 8, for in the first instance the true value is said to lie between 7.95000 ... and 8.049999 ... (i.e., within a range of 0.1 cm), whereas 8 implies a value between 7.50000 ... and 8.49999 ... (i.e., within a range of 1 cm). To state 8.00 cm implies a measurement that ascertains the frog's limb length to be between 7.99500 ... and 8.00499 ... cm (i.e., within a range of 0.01 cm). Those digits in a number that denote the accuracy of the measurement are referred to as *significant figures*. Thus, 8 has one significant figure, 8.0 and 8.3 each have two significant figures, and 8.00 and 8.32 each have three.

In working with exact values of discrete variables, the preceding considerations do not apply. That is, it is sufficient to state that our frog has four limbs or that its left lung contains thirteen flukes. The use of 4.0 or 13.00 would be inappropriate, for as the numbers involved are exactly 4 and 13, there is no question of accuracy or significant figures.

But there are instances where significant figures and implied accuracy come into play with discrete data. An entomologist may report that there are 72,000 moths in a particular forest area. In doing so, it is probably not being claimed that this is the exact number but an estimate of the exact number, perhaps accurate to two significant figures. In such a case, 72,000 would imply a range of accuracy of 1000, so that the true value might lie anywhere from 71,500 to 72,500. If the entomologist wished to convey the fact that this estimate is believed to be accurate to the nearest 100 (i.e., to three significant figures), rather than to the nearest 1000, it would be better to present the data in the form of *scientific notation*,* as follows: If the number 7.2×10^4 ($= 72,000$) is written, a range of accuracy of 0.1×10^4 ($= 1000$) is implied, and the true value is assumed to lie between 71,500 and 72,500. But if 7.20×10^4 were written, a range of accuracy of 0.01×10^4 ($= 100$) would be implied, and the true value would be assumed to be in the range of 71,950 to 72,050. Thus, the accuracy of large values (and this applies to continuous as well as discrete variables) can be expressed succinctly using scientific notation.

Calculators and computers typically yield results with more significant figures than are justified by the data. However, it is good practice—to avoid rounding error—to retain many significant figures until the last step in a sequence of calculations, and on attaining the result of the final step to round off to the appropriate number of figures. A suggestion for the number of figures to report is given at the end of Section 6.2.

1.3 FREQUENCY DISTRIBUTIONS

When collecting and summarizing large amounts of data, it is often helpful to record the data in the form of a *frequency table*. Such a table simply involves a listing of all the observed values of the variable being studied and how many times each value is observed. Consider the tabulation of the frequency of occurrence of sparrow nests in each of several different locations. This is illustrated in Example 1.1, where the observed kinds of nest sites are listed, and for each kind the number of nests observed is recorded. The distribution of the total number of observations among the various categories is termed a *frequency distribution*. Example 1.1 is a frequency table for nominal data, and these data may also be presented graphically by means of a *bar graph* (Figure 1.2), where the height of each bar is proportional to the frequency in the class represented. The widths of all bars in a bar graph should be equal so

*The use of scientific notation—by physicists—can be traced back to at least the 1860s (Miller, 2004b).

EXAMPLE 1.1 The Location of Sparrow Nests: A Frequency Table of Nominal Data

The variable is nest site, and there are four recorded categories of this variable. The numbers recorded in these categories constitute the frequency distribution.

<i>Nest Site</i>	<i>Number of Nests Observed</i>
A. Vines	56
B. Building eaves	60
C. Low tree branches	46
D. Tree and building cavities	49

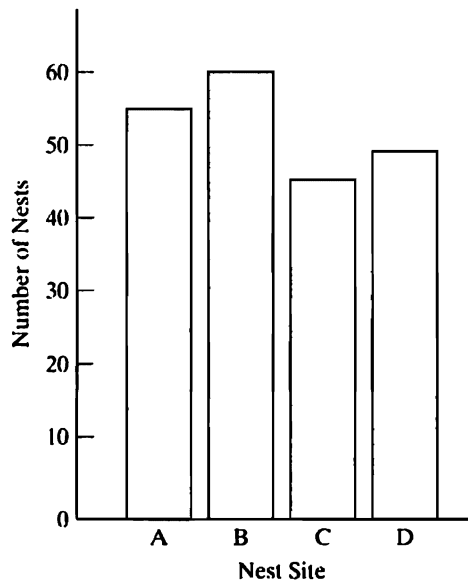


FIGURE 1.2: A bar graph of the sparrow nest data of Example 1.1. An example of a bar graph for nominal data.

that the eye of the reader is not distracted from the differences in bar heights; this also makes the area of each bar proportional to the frequency it represents. Also, the frequency scale on the vertical axis should begin at zero to avoid the apparent differences among bars. If, for example, a bar graph of the data of Example 1.1 were constructed with the vertical axis representing frequencies of 45 to 60 rather than 0 to 60, the results would appear as in Figure 1.3. Huff (1954) illustrates other techniques that can mislead the readers of graphs. It is good practice to leave space between the bars of a bar graph of nominal data, to emphasize the distinctness among the categories represented.

A frequency tabulation of ordinal data might appear as in Example 1.2, which presents the observed numbers of sunfish collected in each of five categories, each category being a degree of skin pigmentation. A bar graph (Figure 1.4) can be prepared for this frequency distribution just as for nominal data.

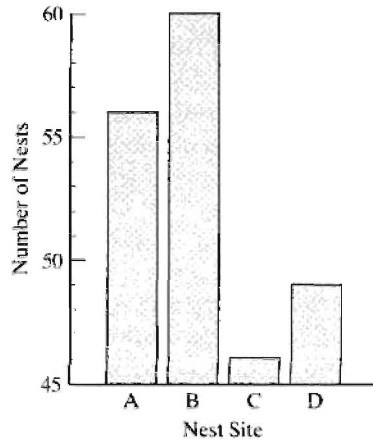


FIGURE 1.3: A bar graph of the sparrow nest data of Example 1.1, drawn with the vertical axis starting at 45. Compare this with Figure 1.1, where the axis starts at 0.

EXAMPLE 1.2 Numbers of Sunfish, Tabulated According to Amount of Black Pigmentation: A Frequency Table of Ordinal Data

The variable is amount of pigmentation, which is expressed by numerically ordered classes. The numbers recorded for the five pigmentation classes compose the frequency distribution.

<i>Pigmentation Class</i>	<i>Amount of Pigmentation</i>	<i>Number of Fish</i>
0	No black pigmentation	13
1	Faintly speckled	68
2	Moderately speckled	44
3	Heavily speckled	21
4	Solid black pigmentation	8

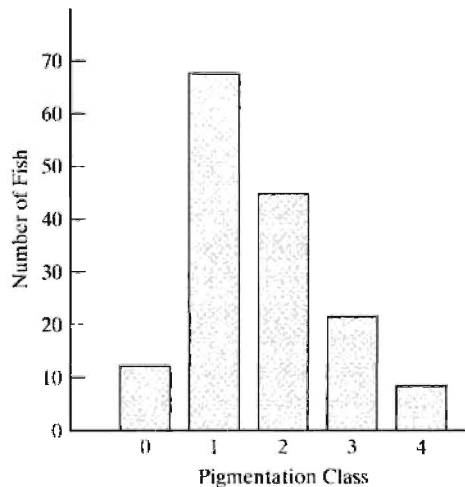


FIGURE 1.4: A bar graph of the sunfish pigmentation data of Example 1.2. An example of a bar graph for ordinal data.

In preparing frequency tables of interval- and ratio-scale data, we can make a procedural distinction between discrete and continuous data. Example 1.3 shows discrete data that are frequencies of litter sizes in foxes, and Figure 1.5 presents this frequency distribution graphically.

EXAMPLE 1.3 Frequency of Occurrence of Various Litter Sizes in Foxes: A Frequency Table of Discrete, Ratio-Scale Data

The variable is litter size, and the numbers recorded for the five litter sizes make up frequency distribution.

<i>Litter Size</i>	<i>Frequency</i>
3	10
4	27
5	22
6	4
7	1

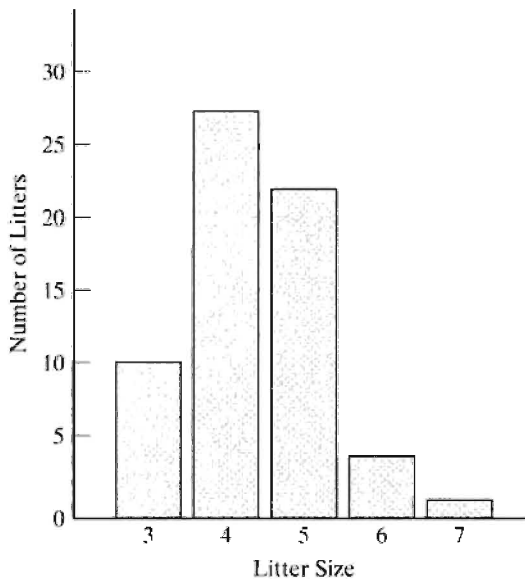


FIGURE 1.5: A bar graph of the fox litter data of Example 1.3. An example of a bar graph for discrete, ratio-scale data.

Example 1.4a shows discrete data that are the numbers of aphids found per clover plant. These data create quite a lengthy frequency table, and it is not difficult to imagine sets of data whose tabulation would result in an even longer list of frequencies. Thus, for purposes of preparing bar graphs, we often cast data into a frequency table by grouping them.

Example 1.4b is a table of the data from Example 1.4a arranged by grouping the data into size classes. The bar graph for this distribution appears as Figure 1.6. Such grouping results in the loss of some information and is generally utilized only to make frequency tables and bar graphs easier to read, and not for calculations performed on

the data. There have been several “rules of thumb” proposed to aid in deciding into how many classes data might reasonably be grouped, for the use of too few groups will obscure the general shape of the distribution. But such “rules” or recommendations are only rough guides, and the choice is generally left to good judgment, bearing in mind that from 10 to 20 groups are useful for most biological work. (See also Doane, 1976.) In general, groups should be established that are equal in the size interval of the variable being measured. (For example, the group size interval in Example 1.4b is four aphids per plant.)

EXAMPLE 1.4a Number of Aphids Observed per Clover Plant: A Frequency Table of Discrete, Ratio-Scale Data

<i>Number of Aphids on a Plant</i>	<i>Number of Plants Observed</i>	<i>Number of Aphids on a Plant</i>	<i>Number of Plants Observed</i>
0	3	20	17
1	1	21	18
2	1	22	23
3	1	23	17
4	2	24	19
5	3	25	18
6	5	26	19
7	7	27	21
8	8	28	18
9	11	29	13
10	10	30	10
11	11	31	14
12	13	32	9
13	12	33	10
14	16	34	8
15	13	35	5
16	14	36	4
17	16	37	1
18	15	38	2
19	14	39	1
		40	0
		41	1

Total number of observations = 424

Because continuous data, contrary to discrete data, can take on an infinity of values, one is essentially always dealing with a frequency distribution tabulated by groups. If the variable of interest were a weight, measured to the nearest 0.1 mg, a frequency table entry of the number of weights measured to be 48.6 mg would be interpreted to mean the number of weights grouped between 48.5500... and 48.6499... mg (although in a frequency table this class interval is usually written as 48.55–48.65). Example 1.5 presents a tabulation of 130 determinations of the amount of phosphorus, in milligrams per gram, in dried leaves. (Ignore the last two columns of this table until Section 1.4.)

EXAMPLE 1.4b Number of Aphids Observed per Clover Plant: A Frequency Table Grouping the Discrete, Ratio-Scale Data of Example 1.4a

<i>Number of Aphids on a Plant</i>	<i>Number of Plants Observed</i>
0–3	6
4–7	17
8–11	40
12–15	54
16–19	59
20–23	75
24–27	77
28–31	55
32–35	32
36–39	8
40–43	1

Total number of observations = 424

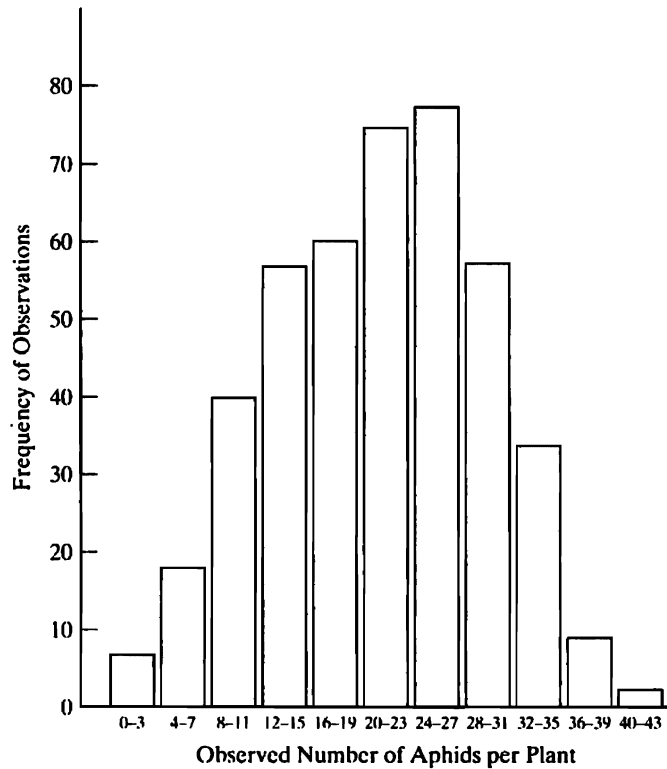


FIGURE 1.6: A bar graph of the aphid data of Example 1.4b. An example of a bar graph for grouped discrete, ratio-scale data.

EXAMPLE 1.5 Determinations of the Amount of Phosphorus in Leaves: A Frequency Table of Continuous Data

<i>Phosphorus (mg/g of leaf)</i>	<i>Frequency (i.e., number of determinations)</i>	Cumulative frequency	
		<i>Starting with Low Values</i>	<i>Starting with High Values</i>
8.15–8.25	2	2	130
8.25–8.35	6	8	128
8.35–8.45	8	16	122
8.45–8.55	11	27	114
8.55–8.65	17	44	103
8.65–8.75	17	61	86
8.75–8.85	24	85	69
8.85–8.95	18	103	45
8.95–9.05	13	116	27
9.05–9.15	10	126	14
9.15–9.25	4	130	4

Total frequency = 130 = n

In presenting this frequency distribution graphically, one can prepare a *histogram*,* which is the name given to a bar graph based on continuous data. This is done in Figure 1.7: note that rather than indicating the range on the horizontal axis, we indicate only the midpoint of the range, a procedure that results in less crowded printing on the graph. Note also that adjacent bars in a histogram are often drawn touching each other, to emphasize the continuity of the scale of measurement, whereas in the other bar graphs discussed they generally are not.

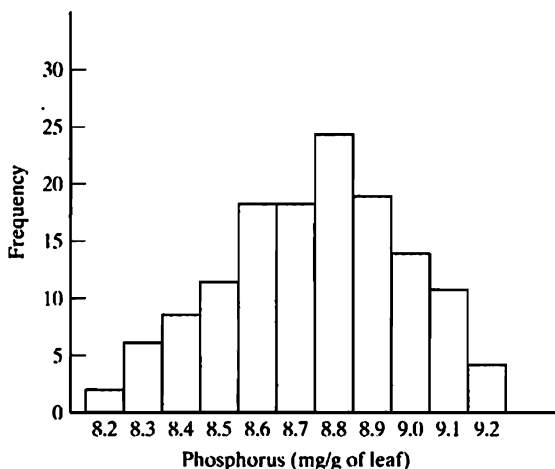


FIGURE 1.7: A histogram of the leaf phosphorus data of Example 1.5. An example of a histogram for continuous data.

*The term *histogram* is from Greek roots (referring to a pole-shaped drawing) and was first published by Karl Pearson in 1895 (David 1995).

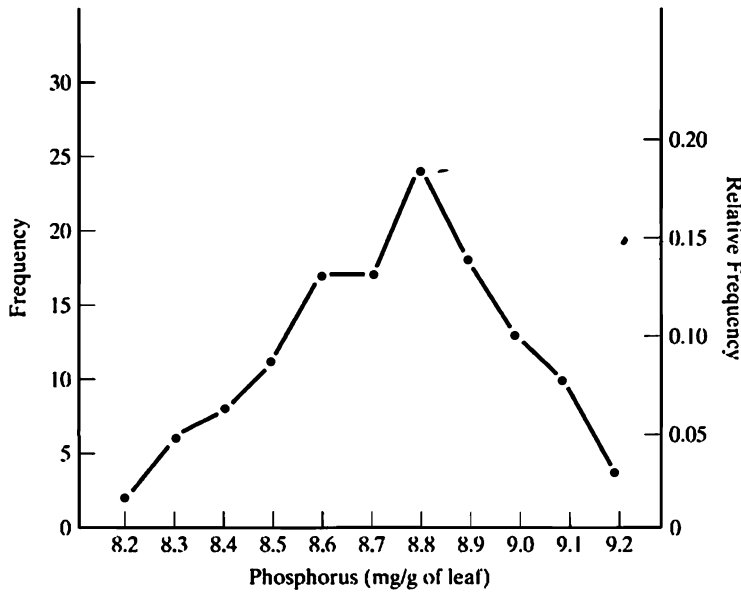


FIGURE 1.8: A frequency polygon for the leaf phosphorus data of Example 1.5.

Often a *frequency polygon* is drawn instead of a histogram. This is done by plotting the frequency of each class as a dot (or other symbol) at the class midpoint and then connecting each adjacent pair of dots by a straight line (Figure 1.8). It is, of course, the same as if the midpoints of the tops of the histogram bars were connected by straight lines. Instead of plotting frequencies on the vertical axis, one can plot *relative frequencies*, or proportions of the total frequency. This enables different distributions to be readily compared and even plotted on the same axes. Sometimes, as in Figure 1.8, frequency is indicated on one vertical axis and the corresponding relative frequency on the other. (Using the data of Example 1.5, the relative frequency for 8.2 mg/g is $2/130 = 0.015$, that for 8.3 mg/g is $6/130 = 0.046$, that for 9.2 mg/g is $4/130 = 0.030$, and so on. The total of all the frequencies is n , and the total of all the relative frequencies is 1.)

Frequency polygons are also commonly used for discrete distributions, but one can argue against their use when dealing with ordinal data, as the polygon implies to the reader a constant size interval horizontally between points on the polygon. Frequency polygons should not be employed for nominal-scale data.

If we have a frequency distribution of values of a continuous variable that falls into a large number of class intervals, the data may be grouped as was demonstrated with discrete variables. This results in fewer intervals, but each interval is, of course, larger. The midpoints of these intervals may then be used in the preparation of a histogram or frequency polygon. The user of frequency polygons is cautioned that such a graph is simply an aid to the eye in following trends in frequency distributions, and one should not attempt to read frequencies between points on the polygon. Also note that the method presented for the construction of histograms and frequency polygons requires that the class intervals be equal. Lastly, the vertical axis (e.g., the frequency scale) on frequency polygons and bar graphs generally should begin with zero, especially if graphs are to be compared with one another. If this is not done, the eye may be misled by the appearance of the graph (as shown for nominal-scale data in Figures 1.2 and 1.3).

1.4 CUMULATIVE FREQUENCY DISTRIBUTIONS

A frequency distribution informs us how many observations occurred for each value (or group of values) of a variable. That is, examination of the frequency table of Example 1.3 (or its corresponding bar graph or frequency polygon) would yield information such as, “How many fox litters of four were observed?”, the answer being 27. But if it is desired to ask questions such as, “How many litters of four or more were observed?”, or “How many fox litters of five or fewer were observed?”, we are speaking of *cumulative frequencies*. To answer the first question, we sum all frequencies for litter sizes four and up, and for the second question, we sum all frequencies from the smallest litter size up through a size of five. We arrive at answers of 54 and 59, respectively.

In Example 1.5, the phosphorus concentration data are cast into two cumulative frequency distributions, one with cumulation commencing at the low end of the measurement scale and one with cumulation being performed from the high values toward the low values. The choice of the direction of cumulation is immaterial, as can be demonstrated. If one desired to calculate the number of phosphorus determinations less than 8.55 mg/g, namely 27, a cumulation starting at the low end might be used, whereas the knowledge of the frequency of determinations greater than 8.55 mg/g, namely 103, can be readily obtained from the cumulation commencing from the high end of the scale. But one can easily calculate any frequency from a low-to-high cumulation (e.g., 27) from its complementary frequency from a high-to-low cumulation (e.g., 103), simply by knowing that the sum of these two frequencies is the total frequency (i.e., $n = 130$); therefore, in practice it is not necessary to calculate both sets of cumulations.

Cumulative frequency distributions are useful in determining medians, percentiles, and other quantiles, as discussed in Sections 3.2 and 4.2. They are not often presented in bar graphs, but *cumulative frequency polygons* (sometimes called *ogives*) are not

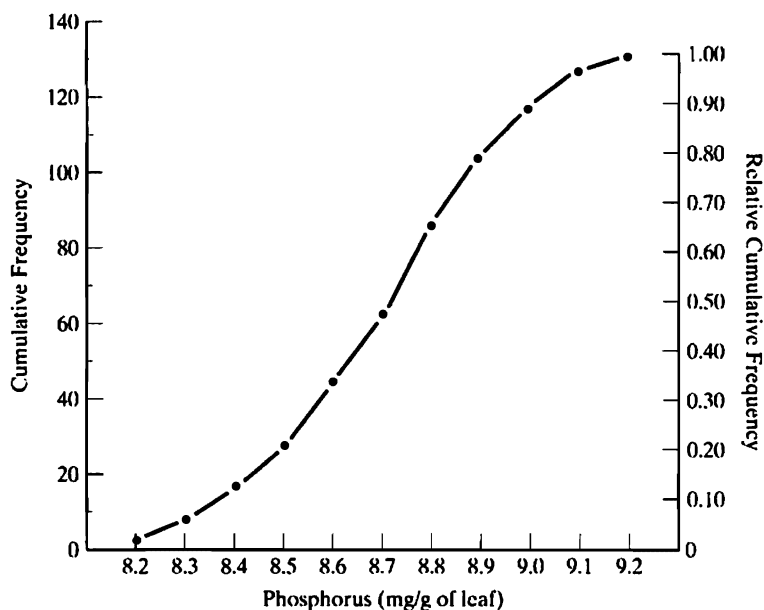


FIGURE 1.9: Cumulative frequency polygon of the leaf phosphorus data of Example 1.5, with cumulation commencing from the lowest to the highest values of the variable.

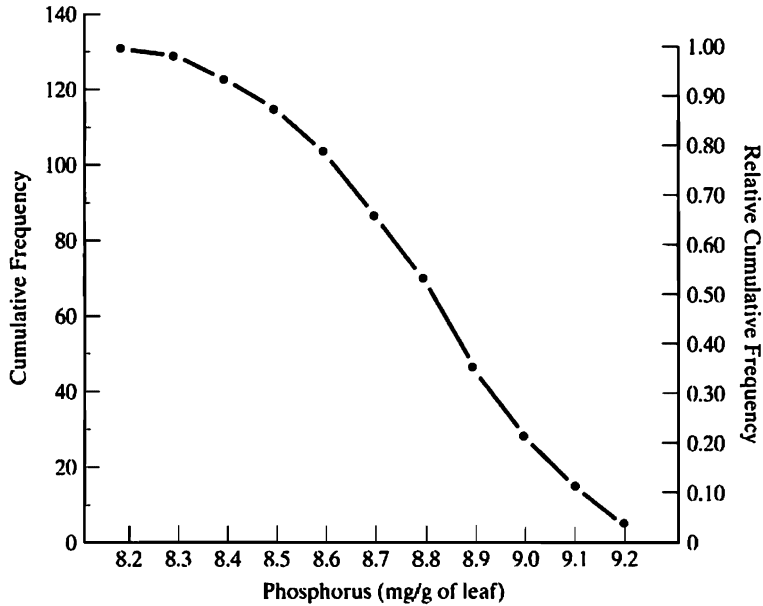


FIGURE 1.10: Cumulative frequency polygon of the leaf phosphorus data of Example 1.5, with cumulation commencing from the highest to the lowest values of the variable.

uncommon. (See Figures 1.9 and 1.10.) Relative frequencies (proportions of the total frequency) can be plotted instead of (or, as in Figures 1.9 and 1.10, in addition to) frequencies on the vertical axis of a cumulative frequency polygon. This enables different distributions to be readily compared and even plotted on the same axes. (Using the data of Example 1.5 for Figure 1.9, the relative cumulative frequency for 8.2 mg/g is $2/130 = 0.015$, that for 8.3 mg/g is $8/130 = 0.062$, and so on. For Figure 1.10, the relative cumulative frequency for 8.2 mg/g is $130/130 = 1.000$, that for 8.3 mg/g is $128/130 = 0.985$, and so on.)

Populations and Samples

-
- 2.1 POPULATIONS
 - 2.2 SAMPLES FROM POPULATIONS
 - 2.3 RANDOM SAMPLING
 - 2.4 PARAMETERS AND STATISTICS
 - 2.5 OUTLIERS
-

The primary objective of a statistical analysis is to infer characteristics of a group of data by analyzing the characteristics of a small sampling of the group. This generalization from the part to the whole requires the consideration of such important concepts as population, sample, parameter, statistic, and random sampling. These topics are discussed in this chapter.

2.1 POPULATIONS

Basic to statistical analysis is the desire to draw conclusions about a group of measurements of a variable being studied. Biologists often speak of a “population” as a defined group of humans or of another species of organisms. Statisticians speak of a *population* (also called a *universe*) as a group of measurements (not organisms) about which one wishes to draw conclusions. It is the latter definition, the statistical definition of *population*, that will be used throughout this book. For example, an investigator may desire to draw conclusions about the tail lengths of bobcats in Montana. All Montana bobcat tail lengths are, therefore, the population under consideration. If a study is concerned with the blood-glucose concentration in three-year-old children, then the blood-glucose levels in all children of that age are the population of interest.

Populations are often very large, such as the body weights of all grasshoppers in Kansas or the eye colors of all female New Zealanders, but occasionally populations of interest may be relatively small, such as the ages of men who have traveled to the moon or the heights of women who have swum the English Channel.

2.2 SAMPLES FROM POPULATIONS

If the population under study is very small, it might be practical to obtain all the measurements in the population. If one wishes to draw conclusions about the ages of all men who have traveled to the moon, it would not be unreasonable to attempt to collect all the ages of the small number of individuals under consideration. Generally, however, populations of interest are so large that obtaining all the measurements is unfeasible. For example, we could not reasonably expect to determine the body weight of every grasshopper in Kansas. What can be done in such cases is to obtain a subset of all the measurements in the population. This subset of measurements constitutes a *sample*, and from the characteristics of samples we can

draw conclusions about the characteristics of the populations from which the samples came.*

Biologists may sample a population that does not physically exist. Suppose an experiment is performed in which a food supplement is administered to 40 guinea pigs, and the sample data consist of the growth rates of these 40 animals. Then the population about which conclusions might be drawn is the growth rates of all the guinea pigs that conceivably might have been administered the same food supplement under identical conditions. Such a population is said to be “imaginary” and is also referred to as “hypothetical” or “potential.”

2.3 RANDOM SAMPLING

Samples from populations can be obtained in a number of ways; however, for a sample to be representative of the population from which it came, and to reach valid conclusions about populations by induction from samples, statistical procedures typically assume that the samples are obtained in a *random* fashion. To sample a population randomly requires that each member of the population has an equal and independent chance of being selected. That is, not only must each measurement in the population have an equal chance of being chosen as a member of the sample, but the selection of any member of the population must in no way influence the selection of any other member. Throughout this book, “sample” will always imply “random sample.”†

It is sometimes possible to assign each member of a population a unique number and to draw a sample by choosing a set of such numbers at random. This is equivalent to having all members of a population in a hat and drawing a sample from them while blindfolded. Appendix Table B.41 provides 10,000 random digits for this purpose. In this table, each digit from 0 to 9 has an equal and independent chance of appearing anywhere in the table. Similarly, each combination of two digits, from 00 to 99, is found at random in the table, as is each three-digit combination, from 000 to 999, and so on.

Assume that a random sample of 200 names is desired from a telephone directory having 274 pages, three columns of names per page, and 98 names per column. Entering Table B.41 at random (i.e., do not always enter the table at the same place), one might decide first to arrive at a random combination of three digits. If this three-digit number is 001 to 274, it can be taken as a randomly chosen page number (if it is 000 or larger than 274, simply skip it and choose another three-digit number, e.g., the next one on the table). Then one might examine the next digit in the table; if it is a 1, 2, or 3, let it denote a page column (if a digit other than 1, 2, or 3 is encountered, it is ignored, passing to the next digit that is 1, 2, or 3). Then one could look at the next two-digit number in the table; if it is from 01 to 98, let it represent a randomly selected name within that column. This three-step procedure would be performed a total of 200 times to obtain the desired random sample. One can proceed in any direction in the random number table: left to right, right to left, upward, downward, or diagonally; but the direction should be decided on before looking at the table. Computers are capable of quickly generating random numbers (sometimes called “pseudorandom” numbers because the number generation is not perfectly random), and this is how Table B.41 was derived.

*This use of the terms *population* and *sample* was established by Karl Pearson (1903).

†This concept of random sampling was established by Karl Pearson between 1897 and 1903 (Miller, 2004a).

Very often it is not possible to assign a number to each member of a population, and random sampling then involves biological, rather than simply mathematical, considerations. That is, the techniques for sampling Montana bobcats or Kansas grasshoppers require knowledge about the particular organism to ensure that the sampling is random. Researchers consult relevant books, periodical articles, or reports that address the specific kind of biological measurement to be obtained.

2.4 PARAMETERS AND STATISTICS

Several measures help to describe or characterize a population. For example, generally a preponderance of measurements occurs somewhere around the middle of the range of a population of measurements. Thus, some indication of a population “average” would express a useful bit of descriptive information. Such information is called a *measure of central tendency* (also called a *measure of location*), and several such measures (e.g., the mean and the median) will be discussed in Chapter 3.

It is also important to describe how dispersed the measurements are around the “average.” That is, we can ask whether there is a wide spread of values in the population or whether the values are rather concentrated around the middle. Such a descriptive property is called a *measure of variability* (or a *measure of dispersion*), and several such measures (e.g., the range and the standard deviation) will be discussed in Chapter 4.

A quantity such as a measure of central tendency or a measure of dispersion is called a *parameter* when it describes or characterizes a population, and we shall be very interested in discussing parameters and drawing conclusions about them. Section 2.2 pointed out, however, that one seldom has data for entire populations, but nearly always has to rely on samples to arrive at conclusions about populations. Thus, one rarely is able to calculate parameters. However, by random sampling of populations, parameters can be estimated well, as we shall see throughout this book. An estimate of a population parameter is called a *statistic*.^{*} It is statistical convention to represent population parameters by Greek letters and sample statistics by Latin letters; the following chapters will demonstrate this custom for specific examples.

The statistics one calculates will vary from sample to sample for samples taken from the same population. Because one uses sample statistics as estimates of population parameters, it behooves the researcher to arrive at the “best” estimates possible. As for what properties to desire in a “good” estimate, consider the following.

First, it is desirable that if we take an indefinitely large number of samples from a population, the long-run average of the statistics obtained will equal the parameter being estimated. That is, for some samples a statistic may underestimate the parameter of interest, and for others it may overestimate that parameter; but in the long run the estimates that are too low and those that are too high will “average out.” If such a property is exhibited by a statistic, we say that we have an *unbiased* statistic or an unbiased estimator.

Second, it is desirable that a statistic obtained from any single sample from a population be very close to the value of the parameter being estimated. This property of a statistic is referred to as *precision*,[†] *efficiency*, or *reliability*. As we commonly secure only one sample from a population, it is important to arrive at a close estimate of a parameter from a single sample.

^{*}This use of the terms *parameter* and *statistic* was defined by R. A. Fisher as early as 1922 (Miller, 2004a; Savage, 1976).

[†]The precision of a sample statistic, as defined here, should not be confused with the precision of a measurement, defined in Section 1.2.

Third, consider that one can take larger and larger samples from a population (the largest sample being the entire population). As the sample size increases, a *consistent* statistic will become a better estimate of the parameter it is estimating. Indeed, if the sample were the size of the population, then the best estimate would be obtained: the parameter itself.

In the chapters that follow, the statistics recommended as estimates of parameters are “good” estimates in the sense that they possess a desirable combination of unbiasedness, efficiency, and consistency.

2.5 OUTLIERS

Occasionally, a set of data will have one or more observations that are so different, relative to the other data in the sample, that we doubt they should be part of the sample. For example, suppose a researcher collected a sample consisting of the body weights of nineteen 20-week-old mallard ducks raised in individual laboratory cages, for which the following 19 data were recorded:

1.87, 3.75, 3.79, 3.82, 3.85, 3.87, 3.90, 3.94, 3.96, 3.99,
3.99, 4.00, 4.03, 4.04, 4.05, 4.06, 4.09, 8.97, and 39.8 kilograms.

Visual inspection of these 19 recorded data casts doubt upon the smallest datum (1.87 kg) and the two largest data (8.97 kg and 39.8 kg) because they differ so greatly from the rest of the weights in the sample. Data in striking disagreement with nearly all the other data in a sample are often called *outliers* or *discordant data*, and the occurrence of such observations generally calls for closer examination.

Sometimes it is clear that an outlier is the result of incorrect recording of data. In the preceding example, a mallard duck weight of 39.8 kg is highly unlikely (to say the least!), for that is about the weight of a 12-year-old boy or girl (and such a duck would probably not fit in one of the laboratory cages). In this case, inspection of the data records might lead us to conclude that this body weight was recorded with a careless placement of the decimal point and should have been 3.98 kg instead of 39.8 kg. And, upon interrogation, the research assistant may admit to weighing the eighteenth duck with the scale set to pounds instead of kilograms, so the metric weight of that animal should have been recorded as 4.07 (not 8.97) kg.

Also, upon further examination of the data-collection process, we may find that the 1.87-kg duck was taken from a wrong cage and was, in fact, only 4 weeks old, not 20 weeks old, and therefore did not belong in this sample. Or, perhaps we find that it was not a mallard duck, but some other bird species (and, therefore, did not belong in this sample). Statisticians say a sample is *contaminated* if it contains a datum that does not conform to the characteristics of the population being sampled. So the weight of a 4-week-old duck, or of a bird of a different species, would be a statistical contaminant and should be deleted from this sample.

There are also instances where it is known that a measurement was faulty—for example, when a laboratory technician spills coffee onto an electronic measuring device or into a blood sample to be analyzed. In such a case, the measurements known to be erroneous should be eliminated from the sample.

However, outlying data can also be correct observations taken from an intended population, collected purely by chance. As we shall see in Section 6.1, when drawing a random sample from a population, it is relatively likely that a datum in the sample will be around the average of the population and very unlikely that a sample datum will be dramatically far from the average. But sample data very far from the average still may be possible.

It should also be noted that in some situations the examination of an outlier may reveal the effect of a previously unsuspected factor. For example, the 1.87-kg duck might, indeed, have been a 20-week-old mallard but suffering from a genetic mutation or a growth-impeding disease deserving of further consideration in additional research.

In summary, it is not appropriate to discard data simply because they appear (to someone) to be unreasonably extreme. However, if there is a very obvious reason for correcting or eliminating a datum, such as the situations described previously, the incorrect data should be corrected or eliminated. In some other cases questionable data can be *accommodated* in statistical analysis, perhaps by employing statistical procedures that give them less weight or analytical techniques that are *robust* in that they are resistant to effects of discrepant data. And in situations when this cannot be done, dubious data will have to remain in the sample (perhaps encouraging the researcher to repeat the experiment with a new set of data).

The idea of rejecting erroneous data dates back over 200 years; and recommendations for formal, objective methods for such rejection began to appear about 150 years ago. Major discussions of outliers, their origin, and treatment (rejection or accommodation) are those of Barnett and Lewis (1994), Beckman and Cook (1983), and Thode (2002: 123–142).

Measures of Central Tendency

- 3.1 THE ARITHMETIC MEAN
- 3.2 THE MEDIAN
- 3.3 THE MODE
- 3.4 OTHER MEASURES OF CENTRAL TENDENCY
- 3.5 CODING DATA

In samples, as well as in populations, one generally finds a preponderance of values somewhere around the middle of the range of observed values. The description of this concentration near the middle is an *average*, or a *measure of central tendency* to the statistician. It is also termed a *measure of location*, for it indicates where, along the measurement scale, the sample or population is located. Various measures of central tendency are useful population parameters, in that they describe an important property of populations. This chapter discusses the characteristics of these parameters and the sample statistics that are good estimates of them.

3.1 THE ARITHMETIC MEAN

The most widely used measure of central tendency is the *arithmetic mean*,* usually referred to simply as the *mean*,† which is the measure most commonly called an “average.”

Each measurement in a population may be referred to as an X_i (read “ X sub i ”) value. Thus, one measurement might be denoted as X_1 , another as X_2 , another as X_3 , and so on. The subscript i might be any integer value up through N , the total number of X values in the population.‡ The mean of the population is denoted by the Greek letter μ (lowercase mu) and is calculated as the sum of all the X_i values divided by the size of the population.

The calculation of the population mean can be abbreviated concisely by the formula

$$\mu = \frac{\sum_{i=1}^N X_i}{N}. \quad (3.1)$$

*As an adjective, *arithmetic* is pronounced with the accent on the third syllable. In early literature on the subject, the adjective *arithmetical* was employed.

†The term *mean* (as applied to the arithmetic mean, as well as to the geometric and harmonic means of Section 3.4) dates from ancient Greece (Walker, 1929: 183), with its current statistical meaning in use by 1755 (Miller, 2004a; Walker, 1929: 176); *central tendency* appeared by the late 1920s (Miller, 2004a).

‡Charles Babbage (1791–1871) (O’Connor and Robertson, 1998) was an English mathematician and inventor who conceived principles used by modern computers—well before the advent of electronics—and who, in 1832, proposed the modern convention of italicizing Latin (also called Roman) letters to denote quantities; nonitalicized letters had already been employed for this purpose for more than six centuries (Miller, 2001).

The Greek letter Σ (capital sigma) means “summation”^{*} and $\sum_{i=1}^N X$ means “summation of all X_i values from X_1 through X_N .” Thus, for example, $\sum_{i=1}^4 X_i = X_1 + X_2 + X_3 + X_4$ and $\sum_{i=3}^5 X_i = X_3 + X_4 + X_5$. Since, in statistical computations, summations are nearly always performed over the entire set of X_i values, this book will assume $\sum X_i$ to mean “sum X_i ’s over all values of i ,” simply as a matter of printing convenience, and $\mu = \sum X_i/N$ would therefore designate the same calculation as would $\mu = \sum_{i=1}^N X_i/N$.

The most efficient, unbiased, and consistent estimate of the population mean, μ , is the sample mean, denoted as \bar{X} (read as “ X bar”). Whereas the size of the population (which we generally do not know) is denoted as N , the size of a sample is indicated by n , and \bar{X} is calculated as

$$\bar{X} = \frac{\sum_{i=1}^n X_i}{n} \quad \text{or} \quad \bar{X} = \frac{\sum X_i}{n}, \quad (3.2)$$

which is read “the sample mean equals the sum of all measurements in the sample divided by the number of measurements in the sample.”[†] Example 3.1 demonstrates the calculation of the sample mean. Note that the mean has the same units of measurement as do the individual observations. The question of how many decimal places should be reported for the mean will be answered at the end of Section 6.2; until then we shall simply record the mean with one more decimal place than the data.

EXAMPLE 3.1 A Sample of 24 from a Population of Butterfly Wing Lengths

X_i (in centimeters): 3.3, 3.5, 3.6, 3.6, 3.7, 3.8, 3.8, 3.8, 3.9, 3.9, 3.9, 4.0, 4.0, 4.0, 4.0, 4.1, 4.1, 4.1, 4.2, 4.2, 4.3, 4.3, 4.4, 4.5.

$$\begin{aligned} \sum X_i &= 95.0 \text{ cm} \\ n &= 24 \\ \bar{X} &= \frac{\sum X_i}{n} = \frac{95.0 \text{ cm}}{24} = 3.96 \text{ cm} \end{aligned}$$

^{*}Mathematician Leonhard Euler (1707–1783; born in Switzerland, worked mostly in Russia), in 1755, was the first to use Σ to denote summation (Cajori, 1928/9, Vol. II: 61).

[†]The modern symbols for plus and minus (“+” and “−”) appear to have first appeared in a 1456 unpublished manuscript by German mathematician and astronomer Regiomontanus (Johannes Müller, 1436–1476), with Bohemia-born Johann (Johannes) Widman (1562–1498) the first, in 1489, to use them in print (Cajori, 1928/9, Vol. I: 128, 231–232). The modern equal sign (“=”) was invented by Welsh physician and mathematician Robert Recorde (1510–1558), who published it in 1557 (though its use then disappeared in print until 1618), and it was well recognized starting in 1631 (Cajori, *ibid.*: 298; Gullberg, 1997: 107). Recorde also was the first to use the plus and minus symbols in an English work (Miller, 2004b). Using a horizontal line to express division derives from its use, in denoting fractions, by Arabic author Al-Ḥaṣṣār in the twelfth century, though it was not consistently employed for several more centuries (Cajori, *ibid.* I: 269, 310). The slash mark (“/”; also known as a solidus, virgule, or diagonal) was recommended to denote division by the English logician and mathematician Augustus De Morgan (1806–1871) in 1845 (*ibid.* I: 312–313), and the India-born Swiss author Johann Heinrich Rahn (1622–1676) proposed, in 1659, denoting division by the symbol “÷”, which previously was often used by authors as a minus sign (*ibid.*: 211, 270; Gullberg, 1997: 105). Many other symbols were used for mathematical operations, before and after these introductions (e.g., Cajori, *ibid.*: 229–245).

If, as in Example 3.1, a sample contains multiple identical data for several values of the variable, then it may be convenient to record the data in the form of a frequency table, as in Example 3.2. Then X_i can be said to denote each of k different measurements and f_i can denote the frequency with which that X_i occurs in the sample. The sample mean may then be calculated, using the sums of the products of f_i and X_i , as*

$$\bar{X} = \frac{\sum_{i=1}^k f_i X_i}{n} \tag{3.3}$$

Example 3.2 demonstrates this calculation for the same data as in Example 3.1.

EXAMPLE 3.2 The Data from Example 3.1 Recorded as a Frequency Table		
X_i (cm)	f_i	$f_i X_i$ (cm)
3.3	1	3.3
3.4	0	0
3.5	1	3.5
3.6	2	7.2
3.7	1	3.7
3.8	3	11.4
3.9	3	11.7
4.0	4	16.0
4.1	3	12.3
4.2	2	8.4
4.3	2	8.6
4.4	1	4.4
4.5	1	4.5
<hr/>		
$\sum f_i = 24$		$\sum f_i X_i = 95.0$ cm

$k = 13$

$\sum_{i=1}^k f_i = n = 24$

$\bar{X} = \frac{\sum_{i=1}^k f_i X_i}{n} = \frac{95.0 \text{ cm}}{24} = 3.96 \text{ cm}$

median = $3.95 \text{ cm} + \left(\frac{1}{4}\right)(0.1 \text{ cm})$

= $3.95 \text{ cm} + 0.025 \text{ cm}$

= 3.975 cm

A similar procedure is computing what is called a *weighted mean*, an expression of the average of several means. For example, we may wish to combine the mean of 3.96 cm from the sample of 24 measurements in Example 3.1 with a mean of 3.78 cm from a sample of 30 measurements and a mean of 4.02 cm from a sample of 15. These three means would be from a total of $24 + 30 + 15 = 69$ data; and if we had all 69 of the data we could sum them and divide the sum by 69 to obtain the overall mean length. However, that overall mean can be obtained without knowing the 69

*Denoting the multiplication of two quantities (e.g., a and b) by their adjacent placement (i.e., ab) derives from practices in Hindu manuscripts of the seventh century (Cajori, 1928/9, Vol. I: 77, 250). Modern multiplication symbols include a raised dot (as in $a \cdot b$), which was suggested in a 1631 posthumous publication of Thomas Harriot (1560?–1621) and prominently adopted in 1698 by the outstanding mathematician Gottfried Wilhelm Leibniz (1646–1716, in what is now Germany); the St. Andrew’s cross (as in $a \times b$), which was used in 1631 by English mathematician William Oughtred (1574–1660) though it was not in general use until more than 200 years later; and the letter X, which was used, perhaps by Oughtred, as early as 1618 (Cajori, *ibid.*: 251; Gullberg, 1997: 104; Miller 2004b). Johann Rahn’s 1659 use of an asterisk-like symbol (as in $a * b$) (Cajori, *ibid.*: 212–213) did not persist but resurfaced in electronic computer languages of the latter half of the twentieth century.

individual measurements, by employing Equation 3.3 with $f_1 = 24$, $X_1 = 3.96$ cm, $f_2 = 30$, $X_2 = 3.78$ cm, $f_3 = 15$, $X_3 = 4.02$ cm, and $n = 69$. This would yield a weighted mean of $\bar{X} = [(24)(3.96 \text{ cm}) + (30)(3.78 \text{ cm}) + (15)(4.02 \text{ cm})]/69 = (268.74 \text{ cm})/69 = 3.89$ cm.

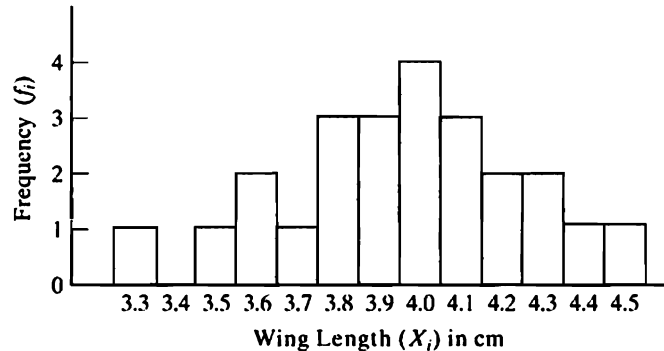


FIGURE 3.1: A histogram of the data in Example 3.2. The mean (3.96 cm) is the center of gravity of the histogram, and the median (3.975 cm) divides the histogram into two equal areas.

If data are plotted as a histogram (Figure 3.1), the mean is the *center of gravity* of the histogram.* That is, if the histogram were made of a solid material, it would balance horizontally with the fulcrum at \bar{X} . The mean is applicable to both ratio- and interval-scale data; it should not be used for ordinal data and cannot be used for nominal data.

3.2 THE MEDIAN

The median is typically defined as the middle measurement in an ordered set of data.† That is, there are just as many observations larger than the median as there are smaller. The sample median is the best estimate of the population median. In a symmetrical distribution (such as Figures 3.2a and 3.2b) the sample median is also an unbiased and consistent estimate of μ , but it is not as efficient a statistic as \bar{X} and should not be used as a substitute for \bar{X} . If the frequency distribution is asymmetrical, the median is a poor estimate of the mean.

The median of a sample of data may be found by first arranging the measurements in order of magnitude. The order may be either ascending or descending, but ascending order is most commonly used as is done with the samples in Examples 3.1, 3.2, and 3.3. Then, we define the sample median as

$$\text{sample median} = X_{(n+1)/2}. \quad (3.4)$$

*The concept of the mean as the center of gravity was used by L. A. J. Quetelet in 1846 (Walker, 1929: 73).

†The concept of the median was conceived as early as 1816, by K. F. Gauss; enunciated and reinforced by others, including F. Galton in 1869 and 1874; and independently discovered and promoted by G. T. Fechner beginning in 1874 (Walker, 1929: 83–88, 184). It received its name, in English, from F. Galton in 1882 (David, 1995) and, in French, from A. A. Cournot in 1843 (David, 1998a).

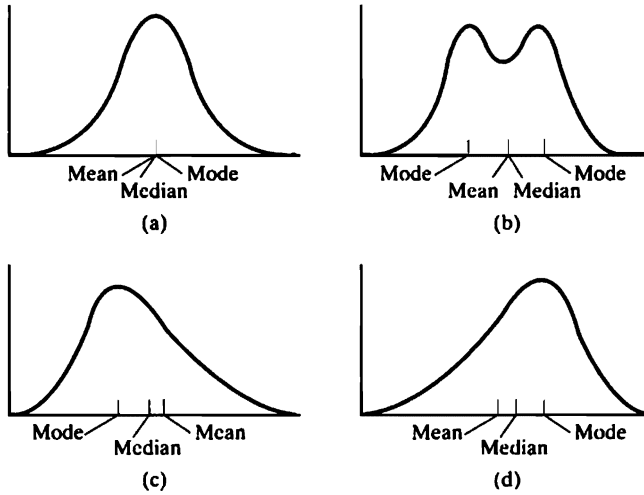


FIGURE 3.2: Frequency distributions showing measures of central tendency. Values of the variable are along the abscissa (horizontal axis), and the frequencies are along the ordinate (vertical axis). Distributions (a) and (b) are symmetrical, (c) is asymmetrical and said to be positively skewed, and (d) is asymmetrical and said to be negatively skewed. Distributions (a), (c), and (d) are unimodal, and distribution b is bimodal. In a unimodal asymmetric distribution, the median lies about one-third the distance between the mean and the mode.*

EXAMPLE 3.3 Life Span for Two Species of Birds in Captivity

The data for each species are arranged in order of magnitude

<i>Species A</i> X_i (mo)	<i>Species B</i> X_i (mo)
16	34
32	36
37	38
39	45
40	50
41	54
42	56
50	59
82	69
	91

$n = 9$	$n = 10$
median = $X_{(n+1)/2} = X_{(9+1)/2}$	median = $X_{(n+1)/2} = X_{(10+1)/2}$
= $X_5 = 40$ mo	= $X_{5.5} = 52$ mo
$\bar{X} = 42.11$ mo	$\bar{X} = 53.20$ mo

*An interesting relationship among the mean, median, and standard deviation is shown in Equation 4.21.

If the sample size (n) is odd, then the subscript in Equation 3.4 will be an integer and will indicate which datum is the middle measurement in the ordered sample. For the data of species *A* in Example 3.3, $n = 9$ and the sample median is $X_{(n+1)/2} = X_{(9+1)/2} = X_5 = 40$ mo. If n is even, then the subscript in Equation 3.4 will be a number midway between two integers. This indicates that there is not a middle value in the ordered list of data; instead, there are two middle values, and the median is defined as the midpoint between them. For the species *B* data in Example 3.3, $n = 10$ and $X_{(n+1)/2} = X_{(10+1)/2} = X_{5.5}$, which signifies that the median is midway between X_5 and X_6 , namely a median of $(50 \text{ mo} + 54 \text{ mo})/2 = 52$ mo.

Note that the median has the same units as each individual measurement. If data are plotted as a frequency histogram (e.g., Figure 3.1), the median is the value of X that divides the area of the histogram into two equal parts. In general, the sample median is a more efficient estimate of the population median when the sample size is large.

If we find the middle value(s) in an ordered set of data to be among identical observations (referred to as *tied* values), as in Example 3.1 or 3.2, a difficulty arises. If we apply Equation 3.4 to these 24 data, then we conclude the median to be $X_{12.5} = 4.0$ cm. But four data are tied at 4.0 cm, and eleven measurements are less than 4.0 cm and nine are greater. Thus, 4.0 cm does not fit the definition above of the median as that value for which there is the same number of data larger and smaller. Therefore, a better definition of the median of a set of data is that value for which no more than half the data are smaller and no more than half are larger.

When the sample median falls among tied observations, we may interpolate to better estimate the population median. Using the data of Example 3.2, we desire to estimate a value below which 50% of the observations in the population lie. Fifty percent of the observations in the sample would be 12 observations. As the first 7 classes in the frequency table include 11 observations and 4 observations are in class 4.0 cm, we know that the desired sample median lies within the range of 3.95 to 4.05 cm. Assuming that the four observations in class 4.0 cm are distributed evenly within the 0.1-cm range of 3.95 to 4.05 cm, then the median will be $\left(\frac{1}{4}\right)(0.1 \text{ cm}) = 0.025$ cm into this class. Thus, the median = 3.95 cm + 0.025 cm = 3.975 cm. In general, for the sample median within a class interval containing tied observations,

$$\text{median} = \left(\begin{array}{c} \text{lower limit} \\ \text{of interval} \end{array} \right) + \left(\frac{0.5n - \text{cum. freq.}}{\text{no. of observations in interval}} \right) \left(\begin{array}{c} \text{interval} \\ \text{size} \end{array} \right), \quad (3.5)$$

where “cum. freq.” refers to the cumulative frequency of the previous classes.* By using this procedure, the calculated median will be the value of X that divides the area of the histogram of the sample into two equal parts. As another example, refer back to Example 1.5, where, by Equation 3.5, median = 8.75 mg/g + $\{[(0.5)(130) - 61]/24\} \{0.10 \text{ mg/g}\} = 8.75 \text{ mg/g} + 0.02 \text{ mg/g} = 8.77 \text{ mg/g}$.

The median expresses less information than does the mean, for it does not take into account the actual value of each measurement, but only considers the rank of each measurement. Still, it offers advantages in some situations. For example, extremely high or extremely low measurements (“outliers”; Section 2.5) do not affect the median as much as they affect the mean (causing the sample median to be called a “resistant” statistic). Distributions that are not symmetrical around the mean (such as in Figures 3.2c and 3.2d) are said to be *skewed*.† When we deal with skewed

*This procedure was enunciated in 1878 by the German psychologist Gustav Theodor Fechner (1801–1887) (Walker, 1929: 86).

†This term, applied to a distribution and to a curve, was used as early as 1895 by Karl Pearson (Miller, 2004a).

populations and do not want the strong influence of outliers, we may prefer the median to the mean to express central tendency.

Note that in Example 3.3 the researcher would have to wait 82 months to compute a mean life expectancy for species *A* and 91 months for species *B*, whereas the median for species *A* could be determined in only 40 months and in only 52 months for species *B*. Also, to calculate a median one does not need to have accurate data for all members of the sample. If, for example, we did not have the first three data for species *A* accurately recorded, but could state them as “less than 39 months,” then the median could have been determined just as readily as if we had all 9 data fully recorded, while calculation of the mean would not have been possible.

The expression “LD fifty” (LD_{50}), used in some areas of biological research, is simply the median lethal dose (and is so named because the median is the 50th percentile, as we shall see in Section 4.2).

The median can be determined not only for interval-scale and ratio-scale data, but also for data on an ordinal scale, data for which the use of the mean usually would not be considered appropriate. But neither the median nor the mean is applicable to nominal data.

3.3 THE MODE

The *mode* is commonly defined as the most frequently occurring measurement in a set of data.* In Example 3.2, the mode is 4.0 cm. But it is perhaps better to define a mode as a measurement of relatively great concentration, for some frequency distributions may have more than one such point of concentration, even though these concentrations might not contain precisely the same frequencies. Thus, a sample consisting of the data 6, 7, 7, 8, 8, 8, 8, 8, 8, 9, 9, 10, 11, 12, 12, 12, 12, 12, 13, 13, and 14 mm would be said to have two modes: at 8 mm and 12 mm. (Some authors would refer to 8 mm as the “major mode” and call 12 mm the “minor mode.”) A distribution in which each different measurement occurs with equal frequency is said to have no mode. If two consecutive values of *X* have frequencies great enough to declare the *X* values modes, the mode of the distribution may be said to be the midpoint of these two *X*'s; for example, the mode of 3, 5, 7, 7, 7, 8, 8, 8, and 10 liters is 7.5 liters. A distribution with two modes is said to be *bimodal* (e.g., Figure 3.2b) and may indicate a combination of two distributions with different modes (e.g., heights of men and women). Modes are often discerned from histograms or frequency polygons; but we should be aware that the shape of such graphs (such as Figures 1.6, 1.7, and 1.8), and therefore the appearance of modes, may be influenced by the measurement intervals on the horizontal axis.

The sample mode is the best estimate of the population mode. When we sample a symmetrical unimodal population, the mode is an unbiased and consistent estimate of the mean and median (Figure 3.2a), but it is relatively inefficient and should not be so used. As a measure of central tendency, the mode is affected by skewness less than is the mean or the median, but it is more affected by sampling and grouping than these other two measures. The mode, but neither the median nor the mean, may be used for data on the nominal, as well as the ordinal, interval, and ratio scales of measurement. In a unimodal asymmetric distribution (Figures 3.2c and 3.2d), the median lies about one-third the distance between the mean and the mode.

The mode is not often used in biological research, although it is often interesting to report the number of modes detected in a population, if there are more than one.

*The term *mode* was introduced by Karl Pearson in 1895 (David, 1995).

3.4 OTHER MEASURES OF CENTRAL TENDENCY

(a) The Geometric Mean. The *geometric mean* is the n th root* of the product of the n data:

$$\bar{X}_G = \sqrt[n]{X_1 X_2 X_3 \dots X_n} = \sqrt[n]{\prod_{i=1}^n X_i}. \quad (3.6)$$

Capital Greek pi, Π , means “take the product”[†] in an analogous fashion as Σ indicates “take the sum.” The geometric mean may also be calculated as the antilogarithm of the arithmetic mean of the logarithms of the data (where the logarithms may be in any base); this is often more feasible computationally:

$$\bar{X}_G = \text{antilog} \left(\frac{\log X_1 + \log X_2 + \dots + \log X_n}{n} \right) = \text{antilog} \frac{\sum_{i=1}^n \log X_i}{n}. \quad (3.7)$$

The geometric mean is appropriate to use only for ratio-scale data and only when all of the data are positive (that is, greater than zero). If the data are all equal, then the geometric mean, \bar{X}_G , is equal to the arithmetic mean, \bar{X} (and also equal to the harmonic mean described below); if the data are not all equal, then[‡] $\bar{X}_G < \bar{X}$.

\bar{X}_G is sometimes used as a measure of location when the data are highly skewed to the right (i.e., when there are many more data larger than the arithmetic mean than there are data smaller than the arithmetic mean).

\bar{X}_G is also useful when dealing with data that represent ratios of change. As an illustration of this, Example 3.4 considers changes in the size of a population of organisms over four decades. Each of the original data (population size at the end of a decade) is expressed as a ratio, X_i , of the population size to the population size of the previous decade. The geometric mean of those ratios is computed and may be thought of as representing the average rate of growth per decade (which is the same as a constant rate of compound interest). This example demonstrates that the arithmetic mean of those ratios is $\bar{X} = 1.1650$ (i.e., 16.50% growth) per decade. But over the four decades of population change, this mean would have us calculate a final population size of $(10,000)(1.1650)(1.1650)(1.1650)(1.1650) = 18,421$, which is *not* the population size recorded at the end of the fourth decade. However, using the geometric mean, \bar{X}_G , to indicate the average rate of growth, the final population size would be computed to be $(10,000)(1.608)(1.608)(1.608)(1.608) = 18,156$, which is the fourth-decade population size that was observed.

*The second footnote in Section 4.5 outlines the origin of the square-root symbol, $\sqrt{\quad}$; indicating the cube root as $\sqrt[3]{\quad}$ was suggested by Albert Girard (1595–1632, French-born but studied and worked in the Netherlands) as early as 1629, but this symbol was not generally used until well into the eighteenth century (Cajori, 1928/9, Vol. I: 371–372). The cube-root symbol eventually was expanded to $\sqrt[n]{\quad}$ to denote the n th root.

[†]Use of this symbol to indicate taking the product was introduced by René Descartes (Gullberg, 1997: 105).

[‡]The symbols “<” and “>” (meaning “less than” and “greater than”) were inserted by someone else into a 1631 posthumous publication by the English mathematician and astronomer Thomas Harriot (1560?–1621), (Cajori, 1928/9, Vol. I: 199; Gullberg, 1997: 109; Miller, 2004b). The symbols for “less than or equal to” (\leq) and “greater than or equal to” (\geq) were written as \leq and \geq when introduced by the French scientist Pierre Bouguere (1698–1758) in 1734. (Gullberg, 1997: 109).

EXAMPLE 3.4 The Geometric Mean of Ratios of Change

Decade	Population Size	Ratio of Change X_i
0	10,000	
1	10,500	$\frac{10,500}{10,000} = 1.05$
2	11,550	$\frac{11,550}{10,500} = 1.10$
3	13,860	$\frac{13,860}{11,550} = 1.20$
4	18,156	$\frac{18,156}{13,860} = 1.31$

$$\bar{X} = \frac{1.05 + 1.10 + 1.20 + 1.31}{4} = \frac{4.66}{4} = 1.1650$$

$$\text{and } (10,000)(0.1650)(1.650)(1.650)(1.650) = 18,421$$

But,

$$\bar{X}_G = \sqrt[4]{(1.05)(1.10)(1.20)(1.31)} = \sqrt[4]{1.8157} = 1.1608$$

or

$$\begin{aligned} \bar{X}_G &= \text{antilog} \left[\frac{\log(1.05) + \log(1.10) + \log(1.20) + \log(1.31)}{4} \right] \\ &= \frac{\text{antilog}(0.0212 + 0.0414 + 0.0792 + 0.1173)}{4} = \frac{\text{antilog}(0.2591)}{4} \\ &= \text{antilog } 0.0648 = 1.1608 \end{aligned}$$

$$\text{and } (10,000)(1.1608)(1.1608)(1.1608)(1.1608) = 18,156$$

(b) The Harmonic Mean. The *harmonic mean* is the reciprocal of the arithmetic mean of the reciprocals of the data:

$$\bar{X}_H = \frac{1}{\frac{1}{n} \sum \frac{1}{X_i}} = \frac{n}{\sum \frac{1}{X_i}} \quad (3.8)$$

It may be used for ratio-scale data when no datum is zero. If all of the data are identical, then the harmonic mean, \bar{X}_H , is equal to the arithmetic mean, \bar{X} (and equal to the geometric mean, \bar{X}_G). If the data are all positive and not identical, then $\bar{X}_H < \bar{X}_G < \bar{X}$.

\bar{X}_H finds use when desiring an average of rates, as described by Croxton, Cowden, and Klein (1967: 182–188). For example, consider that a flock of birds flies from a roosting area to a feeding area 20 km away, flying at a speed of 40 km/hr (which

takes 0.5 hr). The flock returns to the roosting area along the same route (20 km), flying at 20 km/hr (requiring 1 hr of flying time). To ask what the average flying speed was, we might employ Equation 3.2 and calculate the arithmetic mean as $\bar{X} = (40 \text{ km/hr} + 20 \text{ km/hr})/2 = 30 \text{ km/hr}$. However, this answer may not be satisfying, because a total of 40 km was traveled in 1.5 hr, indicating a speed of $(40 \text{ km})/(1.5 \text{ hr}) = 26.7 \text{ km/hr}$. Example 3.5 shows that the harmonic mean (\bar{X}_H) is 26.7 km/hr.

EXAMPLE 3.5 The Harmonic Mean of Rates

$$X_1 = 40 \text{ km/hr}, X_2 = 20 \text{ km/hr}$$

$$\bar{X} = \frac{40 \text{ km/hr} + 20 \text{ km/hr}}{2} = \frac{60 \text{ km/hr}}{2} = 30 \text{ km/hr}$$

But

$$\begin{aligned} \bar{X}_H &= \frac{2}{\frac{1}{40 \text{ km/hr}} + \frac{1}{20 \text{ km/hr}}} = \frac{2}{0.0250 \text{ hr/km} + 0.0500 \text{ hr/km}} \\ &= \frac{2}{0.075 \text{ hr/km}} = 26.67 \text{ km/hr} \end{aligned}$$

(c) The Range Midpoint. The *range midpoint*, or *midrange*, is a measure of location defined as the point halfway between the minimum and the maximum values in the set of data. It may be used with data measured on the ratio, interval, or ordinal scale; but it is not generally a good estimate of location, for it utilizes relatively little information from the data. (However, the so-called mean daily temperature is often reported as the mean of the minimum and maximum and is, therefore, a range midpoint.)

The midpoint of any two symmetrically located percentiles (see Section 4.2), such as the point midway between the first and third quartiles (i.e., the 25th and 75th percentiles), may be used as a location measure in the same fashion as the range midpoint is used (see Dixon and Massey, 1969: 133–134). Such measures are not as adversely affected by aberrantly extreme values as is the range midpoint, and they may be applied to ratio or interval data. If used with ordinal data, they (and the range midpoint) would be the same as the median.

3.5 CODING DATA

Often in the manipulation of data, considerable time and effort can be saved if *coding* is employed. Coding is the conversion of the original measurements into easier-to-work-with values by simple arithmetic operations. Generally coding employs a *linear transformation* of the data, such as multiplying (or dividing) or adding (or subtracting) a constant. The addition or subtraction of a constant is sometimes termed a translation of the data (i.e., changing the origin), whereas the multiplication or division by a constant causes an expansion or contraction of the scale of measurement.

EXAMPLE 3.6 Coding Data to Facilitate Calculations**Sample 1 (Coding by Subtraction:
 $A = -840$ g)****Sample 2 (Coding by Division:
 $M = 0.001$ liters/ml)**

X_i (g)	coded $X_i = X_i - 840$ g	X_i (ml)	coded $X_i = (X_i)(0.001$ liters/ml) $= X_i$ liters
842	2	8,000	8.000
844	4	9,000	9.000
846	6	9,500	9.500
846	6	11,000	11.000
847	7	12,500	12.500
848	8	13,000	13.000
849	9		
$\sum X_i = 5922$ g		$\sum X_i = 63,000$ ml	
coded $\sum X_i = 42$ g		coded $\sum X_i = 63.000$ liters	
$\bar{X} = \frac{5922 \text{ g}}{7}$	coded $\bar{X} = \frac{42 \text{ g}}{7}$	$\bar{X} = 10,500$ ml	coded $\bar{X} = 10.500$ liters
$= 846$ g	$= 6$ g		
$\bar{X} = \text{coded } \bar{X} - A$	$\bar{X} = \text{coded } \frac{\bar{X}}{M}$		
$= 6 \text{ g} - (-840 \text{ g})$	$= \frac{10.500 \text{ liters}}{0.001 \text{ liters/ml}}$		
$= 846$ g	$= 10,500$ ml		

The first set of data in Example 3.6 are coded by subtracting a constant value of 840 g. Not only is each coded value equal to $X_i - 840$ g, but the mean of the coded values is equal to $\bar{X} - 840$ g. Thus, the easier-to-work-with coded values may be used to calculate a mean that then is readily converted to the mean of the original data, simply by adding back the coding constant.

In Sample 2 of Example 3.6, the observed data are coded by dividing each observation by 1000 (i.e., by multiplying by 0.001).^{*} The resultant mean only needs to be multiplied by the coding factor of 1000 (i.e., divided by 0.001) to arrive at the mean of the original data. As the other measures of central tendency have the same units as the mean, they are affected by coding in exactly the same fashion.

Coding affects the median and mode in the same way as the mean is affected. The widespread use of computers has greatly diminished the need for researchers to

^{*}In 1593, mathematician Christopher Clavius (1538–1612, born in what is now Germany but spent most of his life in what is now Italy; also credited with proposing the currently used Gregorian calendar rules regarding leap years: O'Connor and Robertson, 1996) became the first to use a decimal point to separate units from tenths; in 1617, the Scottish mathematician John Napier (1550–1617) used both points and commas for this purpose (Cajori, 1928/9, Vol. I: 322–323), and the comma is still so used in some parts of the world. In some countries a raised dot has been used—a symbol Americans sometimes employ to denote multiplication.

utilize coding (although computer software may use it). Appendix C presents coding for a variety of statistics.

EXERCISES

3.1. If $X_1 = 3.1$ kg, $X_2 = 3.4$ kg, $X_3 = 3.6$ kg, $X_4 = 3.7$ kg, and $X_5 = 4.0$ kg, calculate the value of

(a) $\sum_{i=1}^4 X_i$.

(b) $\sum_{i=2}^4 X_i$.

(c) $\sum_{i=1}^5 X_i$.

(d) $\sum X_i$.

3.2. (a) Calculate the mean of the five weights in Exercise 3.1.

(b) Calculate the median of those weights.

3.3. The ages, in years, of the faculty members of a university biology department are 32.2, 37.5, 41.7, 53.8, 50.2, 48.2, 46.3, 65.0, and 44.8.

(a) Calculate the mean age of these nine faculty members.

(b) Calculate the median of the ages.

(c) If the person 65.0 years of age retires and is replaced on the faculty with a person 46.5 years old, what is the new mean age?

(d) What is the new median age?

3.4. Consider the following frequency tabulation of leaf weights (in grams):

X_i	f_i
1.85–1.95	2
1.95–2.05	1
2.05–2.15	2
2.15–2.25	3
2.25–2.35	5
2.35–2.45	6
2.45–2.55	4
2.55–2.65	3
2.65–2.75	1

Using the midpoints of the indicated ranges of X_i ,

(a) Calculate the mean leaf weight using Equation 3.2, and

(b) Calculate the mean leaf weight using Equation 3.3.

(c) Calculate the median leaf weight using Equation 3.4, and

(d) Calculate the median using Equation 3.5.

(e) Determine the mode of the frequency distribution.

3.5. A fruit was collected from each of eight lemon trees, with the intent of measuring the calcium concentration in the rind (grams of calcium per 100 grams of dry rind). The analytical method used could only detect a concentration of at least 0.80 g/100 g of dry weight. Six of the eight concentrations were measured to be 1.02, 0.98, 0.91, 0.84, 0.87, 1.04 g/100 g of dry weight, and two of the concentrations were known to be less than 0.80 g/100 g of dry weight. What is the median of this sample of eight data?

Measures of Variability and Dispersion

-
- 4.1 THE RANGE
 - 4.2 DISPERSION MEASURED WITH QUANTILES
 - 4.3 THE MEAN DEVIATION
 - 4.4 THE VARIANCE
 - 4.5 THE STANDARD DEVIATION
 - 4.6 THE COEFFICIENT OF VARIATION
 - 4.7 INDICES OF DIVERSITY
 - 4.8 CODING DATA
-

In addition to a description of the central tendency of a set of data, it is generally desirable to have a description of the *variability*, or of the *dispersion*,* of the data. A measure of variability (or measure of dispersion, as it is often called) is an indication of the spread of measurements around the center of the distribution. Measurements that are concentrated around the center of a distribution of data have low variability (low dispersion), whereas data that are very spread out along the measurement scale have high variability (high dispersion). Measures of variability of a population are population parameters, and sample measures of variability are statistics that estimate those parameters.

4.1 THE RANGE

The difference between the highest and lowest measurements in a group of data is termed the *range*.† If sample measurements are arranged in increasing order of magnitude, as if the median were about to be determined, then

$$\text{sample range} = X_n - X_1, \quad (4.1)$$

which is

$$\text{sample range} = \text{largest } X - \text{smallest } X.$$

Sample 1 in Example 4.1 is a hypothetical set of ordered data in which $X_1 = 1.2$ g and $X_n = 2.4$ g. Thus, the range may be expressed as 1.2 to 2.4 g, or as $2.4 \text{ g} - 1.2 \text{ g} = 1.2 \text{ g}$. Note that the range has the same units as the individual measurements. Sample 2 in Example 4.1 has the same range as Sample 1.

*The statistical use of this term first appeared in an 1876 publication by Francis Galton (David, 1998a).

†This statistical term dates from an 1848 paper by H. Lloyd (David, 1995). It was already used by the Greek astronomer Hipparchus as a measure of dispersion in the second century B.C.E. (David, 1998b).

EXAMPLE 4.1 Calculation of Measures of Dispersion for Two Hypothetical Samples of 7 Insect Body Weights

Sample 1

X_i (g)	$X_i - \bar{X}$ (g)	$ X_i - \bar{X} $ (g)	$(X_i - \bar{X})^2$ (g ²)
1.2	-0.6	0.6	0.36
1.4	-0.4	0.4	0.16
1.6	-0.2	0.2	0.04
1.8	0.0	0.0	0.00
2.0	0.2	0.2	0.04
2.2	0.4	0.4	0.16
2.4	0.6	0.6	0.36

$$\begin{aligned} \sum X_i &= 12.6 \text{ g} & \sum (X_i - \bar{X}) &= 0.0 \text{ g} & \sum |X_i - \bar{X}| &= 2.4 \text{ g} & \sum (X_i - \bar{X})^2 &= 1.12 \text{ g}^2 \end{aligned}$$

= sum of squared deviations from the mean
= "sum of squares"

$$n = 7; \bar{X} = \frac{\sum X_i}{n} = \frac{12.6 \text{ g}}{7} = 1.8 \text{ g}$$

$$\text{range} = X_7 - X_1 = 2.4 \text{ g} - 1.2 \text{ g} = 1.2 \text{ g}$$

$$\text{interquartile range} = Q_3 - Q_1 = 2.2 \text{ g} - 1.4 \text{ g} = 0.8 \text{ g}$$

$$\text{mean deviation} = \frac{\sum |X_i - \bar{X}|}{n} = \frac{2.4 \text{ g}}{7} = 0.34 \text{ g}$$

$$\text{variance} = s^2 = \frac{\sum (X_i - \bar{X})^2}{n - 1} = \frac{1.12 \text{ g}^2}{6} = 0.1867 \text{ g}^2$$

$$\text{standard deviation} = s = \sqrt{0.1867 \text{ g}^2} = 0.43 \text{ g}$$

Sample 2

X_i (g)	$X_i - \bar{X}$ (g)	$ X_i - \bar{X} $ (g)	$(X_i - \bar{X})^2$ (g ²)
1.2	-0.6	0.6	0.36
1.6	-0.2	0.2	0.04
1.7	-0.1	0.1	0.01
1.8	0.0	0.0	0.00
1.9	0.1	0.1	0.01
2.0	0.2	0.2	0.04
2.4	0.6	0.6	0.36

$$\begin{aligned} \sum X_i &= 12.6 \text{ g} & \sum (X_i - \bar{X}) &= 0.0 \text{ g} & \sum |X_i - \bar{X}| &= 1.8 \text{ g} & \sum (X_i - \bar{X})^2 &= 0.82 \text{ g}^2 \end{aligned}$$

= sum of squared deviations from the mean
= "sum of squares"

$$n = 7; \bar{X} = \frac{\sum X_i}{n} = \frac{12.6 \text{ g}}{7} = 1.8 \text{ g}$$

$$\text{range} = X_7 - X_1 = 2.4 \text{ g} - 1.2 \text{ g} = 1.2 \text{ g}$$

$$\begin{aligned} \text{interquartile range} &= Q_3 - Q_1 = 2.0 \text{ g} - 1.6 \text{ g} = 0.4 \text{ g} \\ \text{mean deviation} &= \frac{\sum |X_i - \bar{X}|}{n} = \frac{1.8 \text{ g}}{7} = 0.26 \text{ g} \\ \text{variance} = s^2 &= \frac{\sum (X_i - \bar{X})^2}{n - 1} = \frac{0.82 \text{ g}^2}{6} = 0.1367 \text{ g}^2 \\ \text{standard deviation} = s &= \sqrt{0.1367 \text{ g}^2} = 0.37 \text{ g} \end{aligned}$$

The range is a relatively crude measure of dispersion, inasmuch as it does not take into account any measurements except the highest and the lowest. Furthermore, it is unlikely that a sample will contain both the highest and lowest values in the population, so the sample range usually underestimates the population range; therefore, it is a biased and inefficient estimator. Nonetheless, it is considered useful by some to present the sample range as an estimate (although a poor one) of the population range. For example, taxonomists are often concerned with having an estimate of what the highest and lowest values in a population are expected to be. Whenever the range is specified in reporting data, however, it is usually a good practice to report another measure of dispersion as well. The range is applicable to ordinal-, interval-, and ratio-scale data.

4.2 DISPERSION MEASURED WITH QUANTILES

Because the sample range is a biased and inefficient estimate of the population range, being sensitive to extremely large and small measurements, alternative measures of dispersion may be desired. Just as the median (Section 3.2) is the value above and below which lies half the set of data, one can define measures, called *quantiles*, above or below which lie other fractional portions of the data.

For example, if the data are divided into four equal parts, we speak of *quartiles*. One-fourth of all the ranked observations are smaller than the first quartile, one-fourth lie between the first and second quartiles, one-fourth lie between the second and third quartiles, and one-fourth are larger than the third quartile. The second quartile is identical to the median. As with the median, the first and third quartiles might be one of the data or the midpoint between two of the data. The first quartile, Q_1 , is

$$Q_1 = X_{(n+1)/4}; \quad (4.2)$$

if the subscript, $(n + 1)/4$, is not an integer or half-integer, then it is rounded up to the nearest integer or half-integer. The second quartile is the median, and the subscript on X for the third quartile, Q_3 , is

$$n + 1 - (\text{subscript on } X \text{ for } Q_1, \text{ after any rounding}). \quad (4.3)$$

Examining the data in Example 3.3: For species *A*, $n = 9$, $(n + 1)/4 = 2.5$, and $Q_1 = X_{2.5} = 34.5$ mo; and $Q_3 = X_{10-2.5} = X_{7.5} = 46$ mo. For species *B*, $n = 10$, $(n + 1)/4 = 2.75$ (which we round up to 3), and $Q_1 = X_3 = 38$ mo, and $Q_3 = X_{11-3} = X_8 = 59$ mo.

The distance between Q_1 and Q_3 , the first and third quartiles (i.e., the 25th and 75th percentiles), is known as the *interquartile range* (or *semiquartile range*):

$$\text{interquartile range} = Q_3 - Q_1. \quad (4.4)$$

One may also encounter the *semi-interquartile range*:

$$\text{semi-interquartile range} = \frac{Q_3 - Q_1}{2}, \quad (4.5)$$

also known as the *quartile deviation*.*

If the distribution of data is symmetrical, then 50% of the measurements lie within one quartile deviation above and below the median. For Sample 1 in Example 4.1, $Q_1 = 1.4$ g, $Q_3 = 2.2$ g, and the interquartile range is $2.2 \text{ g} - 1.4 \text{ g} = 0.8 \text{ g}$. And for Sample 2, $Q_1 = 1.6$ g, $Q_3 = 2.0$ g, and the interquartile range is $2.0 \text{ g} - 1.6 \text{ g} = 0.4 \text{ g}$.

Similarly, values that partition the ordered data set into eight equal parts (or as equal as n will allow) are called *octiles*. The first octile, \mathcal{O}_1 , is

$$\mathcal{O}_1 = X_{(n+1)/8}; \quad (4.6)$$

and if the subscript, $(n + 1)/8$, is not an integer or half-integer, then it is rounded up to the nearest integer or half-integer. The second, fourth, and sixth octiles are the same as quartiles; that is, $\mathcal{O}_2 = Q_1$, $\mathcal{O}_4 = Q_2 = \text{median}$ and $\mathcal{O}_6 = Q_3$. The subscript on X for the third octile, \mathcal{O}_3 , is

$$2(\text{subscript on } X \text{ for } Q_1) - \text{subscript on } X \text{ for } \mathcal{O}_1; \quad (4.7)$$

the subscript on X for the fifth octile, \mathcal{O}_5 , is

$$n + 1 - \text{subscript on } X \text{ for } \mathcal{O}_3; \quad (4.8)$$

and the subscript on X for the seventh octile, \mathcal{O}_7 , is

$$n + 1 - \text{subscript on } X \text{ for } \mathcal{O}_1. \quad (4.9)$$

Thus, for the data of Example 3.3: For species *A*, $n = 9$, $(n + 1)/8 = 1.5$ and $\mathcal{O}_1 = X_{1.5} = 35$ mo; $2(2.5) - 1.5 = 3.5$, so $\mathcal{O}_3 = X_{3.5} = 38$ mo; $n + 1 - 3.5 = 6.5$, so $\mathcal{O}_5 = X_{6.5} = 41.5$ mo; and $n + 1 - 1.5 = 8.5$, so $\mathcal{O}_7 = 61$. For species *B*, $n = 10$, $(n + 1)/8 = 1.25$ (which we round up to 1.5) and $\mathcal{O}_1 = X_{1.5} = 35$ mo; $2(3) - 1.5 = 4.5$, so $\mathcal{O}_3 = X_{4.5} = 39.5$ mo; $n + 1 - 4.5 = 6.5$, so $\mathcal{O}_5 = X_{6.5} = 41.5$ mo; and $n + 1 - 1.5 = 9.5$, so $\mathcal{O}_7 = 44.5$ mo.

Besides the median, quartiles, and octiles, ordered data may be divided into fifths, tenths, or hundredths by quantities that are respectively called *quintiles*, *deciles*, and *centiles* (the latter also called *percentiles*). Measures that divide a group of ordered data into equal parts are collectively termed *quantiles*.[†] The expression “LD₅₀,” used in some areas of biological research, is simply the 50th percentile of the lethal doses, or the median lethal dose. That is, 50% of the experimental subjects survived this dose, whereas 50% did not. Likewise, “LC₅₀” is the median lethal concentration, or the 50th percentile of the lethal concentrations.

Instead of distance between the 25th and 75th percentiles, distances between other quantiles (e.g., 10th and 90th percentiles) may be used as a dispersion measure. Quantile-based measures of dispersion are valid for ordinal-, interval-, or ratio-scale data, and they do not exhibit the bias and inefficiency of the range.

*This measure was proposed in 1846 by L. A. J. Quetelet (1796–1874); Sir Francis Galton (1822–1911) later called it the “quartile deviation” (Walker, 1929: 84) and, in 1882, used the terms “quartile” and “interquartile range” (David, 1995).

[†]Sir Francis Galton developed the concept of percentiles, quartiles, deciles, and other quantiles in writings from 1869 to 1885 (Walker, 1929: 86–87, 177, 179). The term *quantile* was introduced in 1940 by M. G. Kendall (David, 1995).

4.3 THE MEAN DEVIATION

As is evident from the two samples in Example 4.1, the range conveys no information about how clustered about the middle of the distribution the measurements are. As the mean is so useful a measure of central tendency, one might express dispersion in terms of deviations from the mean. The sum of all deviations from the mean, that is, $\sum(X_i - \bar{X})$, will always equal zero, however, so such a summation would be useless as a measure of dispersion (as seen in Example 4.1).

Using the absolute values of the deviations from the mean eliminates the negative signs of the deviations, and summing those absolute values results in a quantity that is an expression of dispersion about the mean. Dividing this quantity by n yields a measure known as the *mean deviation*, or *mean absolute deviation*,* of the sample; this measure has the same units as do the data. In Example 4.1, Sample 1 is more variable (or more dispersed, or less concentrated) than Sample 2. Although the two samples have the same range, the mean deviations, calculated as

$$\text{sample mean deviation} = \frac{\sum |X_i - \bar{X}|}{n}, \quad (4.10)$$

express the differences in dispersion.† A different kind of mean deviation can be defined by using the sum of the absolute deviations from the median instead of from the mean.

Mean deviations are seldom encountered, because their utility is far less than that of the statistics in Sections 4.4 and 4.5.

4.4 THE VARIANCE

Another method of eliminating the negative signs of deviations from the mean is to square the deviations. The sum of the squares of the deviations from the mean is often simply called the *sum of squares*, abbreviated SS, and is defined as follows:‡

$$\text{population SS} = \sum (X_i - \mu)^2 \quad (4.11)$$

$$\text{sample SS} = \sum (X_i - \bar{X})^2. \quad (4.12)$$

It can be seen from the above two equations that as a measure of variability, or dispersion, the sum of squares considers how far the X_i 's deviate from the mean. In

*The term *mean deviation* is apparently due to Karl Pearson (1857–1936) (Walker, 1929: 55) and *mean absolute deviation*, in 1972, to D. F. Andrews, P. J. Bickel, F. R. Hampel, P. J. Huber, W. H. Rogers, and J. W. Tukey (David, 1995).

†Karl Weierstrass, in 1841, was the first to denote the absolute value of a quantity by enclosing it within two vertical lines (Cajori, 1928/9, Vol. II: p. 123); that is, $|a| = a$ and $|-a| = a$.

‡The modern notation using raised numerals as exponents was introduced by René Descartes in 1637, and many other kinds of notation for exponents were employed before and after that (Cajori, 1928/9, Vol. I: 358; Gullberg, 1997: 134). An 1845 notation of Augustus De Morgan, $a \wedge b$ to indicate a^b (Cajori, *ibid.*: 358), has reemerged in modern computer use. Nicolas Chuquet (1445–1488) was the first to use negative exponents, and Nicole (also known as Nicolaus) Oresme (1323–1382) was the first to use fractional exponents, though neither of these French mathematicians employed the modern notation of Isaac Newton (1642–1727), the colossal English mathematician, physicist, and astronomer (Cajori, *ibid.*: 91, 102, 354–355):

$$x^{-a} = \frac{1}{x^a}; \quad x^{\frac{1}{a}} = \sqrt[a]{x}.$$

Using parentheses or brackets to group quantities dates from the mid-sixteenth century, though it was not common mathematical notation until more than two centuries later (*ibid.*: 392).

Sample 1 of Example 4.1, the sample mean is 1.8 g and it is seen (in the last column) that

$$\begin{aligned}\text{Sample SS} &= (1.2 - 1.8)^2 + (1.4 - 1.8)^2 + (1.6 - 1.8)^2 + (1.8 - 1.8)^2 \\ &\quad + (2.0 - 1.8)^2 + (2.2 - 1.8)^2 + (2.4 - 1.8)^2 \\ &= 0.36 + 0.16 + 0.04 + 0.00 + 0.04 + 0.16 + 0.36 \\ &= 1.12\end{aligned}$$

(where the units are grams²).^{*} The sum of squares may also be visualized as a measure of the average extent to which the data deviate from each other, for (using the same seven data from Sample 1 in Example 4.1):

$$\begin{aligned}\text{SS} &= [(1.2 - 1.4)^2 + (1.2 - 1.6)^2 + (1.2 - 1.8)^2 + (1.2 - 2.0)^2 \\ &\quad + (1.2 - 2.2)^2 + (1.2 - 2.4)^2 + (1.4 - 1.6)^2 + (1.4 - 1.8)^2 \\ &\quad + (1.4 - 2.0)^2 + (1.4 - 2.2)^2 + (1.4 - 2.4)^2 + (1.6 - 1.8)^2 \\ &\quad + (1.6 - 2.0)^2 + (1.6 - 2.2)^2 + (1.6 - 2.4)^2 + (1.8 - 2.0)^2 \\ &\quad + (1.8 - 2.2)^2 + (1.8 - 2.4)^2 + (2.0 - 2.2)^2 + (2.0 - 2.4)^2 \\ &\quad + (2.2 - 2.4)^2]/7 \\ &= [0.04 + 0.16 + 0.36 + 0.64 + 1.00 + 1.44 + 0.04 + \cdots + 0.04 + 0.16 \\ &\quad + 0.04]/7 \\ &= 7.84/7 = 1.12\end{aligned}$$

(again in grams²).

The mean sum of squares is called the *variance* (or *mean square*,[†] the latter being short for *mean squared deviation*), and for a population is denoted by σ^2 (“sigma squared,” using the lowercase Greek letter):

$$\sigma^2 = \frac{\sum(X_i - \mu)^2}{N}. \quad (4.14)$$

The best estimate of the population variance, σ^2 , is the sample variance, s^2 :

$$s^2 = \frac{\sum(X_i - \bar{X})^2}{n - 1}. \quad (4.15)$$

If, in Equation 4.14, we replace μ by \bar{X} and N by n , the result is a quantity that is a biased estimate of σ^2 in that it underestimates σ^2 . Dividing the sample sum of squares

^{*}Owing to an important concept in statistics, known as *least squares*, the sum of squared deviations from the mean is smaller than the sum of squared deviations from any other quantity (e.g., the median). Indeed, if Equation 4.12 is applied using some quantity in place of the mean, the resultant “sum of squares” would be

$$SS + nd^2, \quad (4.13)$$

where d is the difference between the mean and the quantity used. For the population sum of squares (defined in Equation 4.11), the relationship would be $SS + Nd^2$.

[†]The term *mean square* dates back at least to an 1875 publication of Sir George Biddell Airy (1801–1892), Astronomer Royal of England (Walker, 1929: 54). The term *variance* was introduced in 1918 by English statistician Sir Ronald Aylmer Fisher (1890–1962) (*ibid.*: 189; David, 1995).

by $n - 1$ (called the *degrees of freedom*,* often abbreviated DF), rather than by n , yields an unbiased estimate, and it is Equation 4.15 that should be used to calculate the sample variance.

If all observations in a sample are equal, then there is no variability (that is, no dispersion) and $s^2 = 0$. And s^2 becomes increasingly large as the amount of variability, or dispersion, increases. Because s^2 is a mean sum of squares, it can never be a negative quantity.

The variance expresses the same type of information as does the mean deviation, but it has certain very important mathematical properties relative to probability and hypothesis testing that make it superior. Thus, the mean deviation is very seldom encountered in biostatistical analysis.

The calculation of s^2 can be tedious for large samples, but it can be facilitated by the use of the equality

$$\text{sample SS} = \sum X_i^2 - \frac{(\sum X_i)^2}{n}. \quad (4.16)$$

This formula is equivalent to Equation 4.12 but is much simpler to work with. Example 4.2 demonstrates its use to obtain a sample sum of squares.

Because the sample variance equals the sample SS divided by DF,

$$s^2 = \frac{\sum X_i^2 - \frac{(\sum X_i)^2}{n}}{n - 1}. \quad (4.17)$$

This last formula is often referred to as a “working formula,” or “machine formula,” because of its computational advantages. There are, in fact, two major advantages in calculating SS by Equation 4.16 rather than by Equation 4.12. First, fewer computational steps are involved, a fact that decreases chance of error. On many calculators the summed quantities, $\sum X_i$ and $\sum X_i^2$, can both be obtained with only one pass through the data, whereas Equation 4.12 requires one pass through the data to calculate \bar{X} and at least one more pass to calculate and sum the squares of the deviations, $X_i - \bar{X}$. Second, there may be a good deal of rounding error in calculating each $X_i - \bar{X}$, a situation that leads to decreased accuracy in computation, but that is avoided by the use of Equation 4.16.†

For data recorded in frequency tables,

$$\text{sample SS} = \sum f_i X_i^2 - \frac{(\sum f_i X_i)^2}{n}, \quad (4.18)$$

*Given the sample mean (\bar{X}) and sample size (n) in Example 4.1, *degrees of freedom* means that the data could have been weights different from those shown, but when any six (i.e., $n - 1$) of the seven weights are specified, then the seventh weight is also known. The term was first used, though in a different context, by Ronald Aylmer Fisher in 1922 (David, 1955).

†Computational formulas advantageous on calculators may not prove accurate on computers (Wilkinson and Dallal, 1977), largely because computers may use fewer significant figures. (Also see Ling, 1974.) Good computer programs use calculation techniques designed to help avoid rounding errors.

where f_i is the frequency of observations with magnitude X_i . But with a calculator or computer it is often faster to use Equation 4.18 for the individual observations, disregarding the class groupings.

The variance has square units. If measurements are in grams, their variance will be in grams squared, or if the measurements are in cubic centimeters, their variance will be in terms of cubic centimeters squared, even though such squared units have no physical interpretation. The question of how many decimal places to report for the variance will be considered at the end of Section 6.2.

EXAMPLE 4.2 "Machine Formula" Calculation of Variance, Standard Deviation, and Coefficient of Variation (These are the data of Example 4.1)

Sample 1		Sample 2	
X_i (g)	X_i^2 (g ²)	X_i (g)	X_i^2 (g ²)
1.2	1.44	1.2	1.44
1.4	1.96	1.6	2.56
1.6	2.56	1.7	2.89
1.8	3.24	1.8	3.24
2.0	4.00	1.9	3.61
2.2	4.84	2.0	4.00
2.4	5.76	2.4	5.76

$\sum X_i = 12.6 \text{ g}$	$\sum X_i^2 = 23.80 \text{ g}^2$	$\sum X_i = 12.6 \text{ g}$	$\sum X_i^2 = 23.50 \text{ g}^2$
$n = 7$		$n = 7$	
$\bar{X} = \frac{12.6 \text{ g}}{7} = 1.8 \text{ g}$		$\bar{X} = \frac{12.6 \text{ g}}{7} = 1.8 \text{ g}$	
$SS = \sum X_i^2 - \frac{(\sum X_i)^2}{n}$		$SS = 23.50 \text{ g}^2 - \frac{(12.6 \text{ g})^2}{7}$	
$= 23.80 \text{ g}^2 - \frac{(12.6 \text{ g})^2}{7}$		$= 0.82 \text{ g}^2$	
$= 23.80 \text{ g}^2 - 22.68 \text{ g}^2$		$s^2 = \frac{0.82 \text{ g}^2}{6} = 0.1367 \text{ g}^2$	
$= 1.12 \text{ g}^2$		$s = \sqrt{0.1367 \text{ g}^2} = 0.37 \text{ g}$	
$s^2 = \frac{SS}{n - 1}$		$V = \frac{0.37 \text{ g}}{1.8 \text{ g}} = 0.21 = 21\%$	
$= \frac{1.12 \text{ g}^2}{6} = 0.1867 \text{ g}^2$			
$s = \sqrt{0.1867 \text{ g}^2} = 0.43 \text{ g}$			
$V = \frac{s}{\bar{X}} = \frac{0.43 \text{ g}}{1.8 \text{ g}} = 0.24 = 24\%$			

4.5 THE STANDARD DEVIATION

The *standard deviation** is the positive square root[†] of the variance; therefore, it has the same units as the original measurements. Thus, for a population,

$$\sigma = \sqrt{\frac{\sum X_i^2 - \frac{(\sum X_i)^2}{N}}{N}}. \quad (4.19)$$

And for a sample,[‡]

$$s = \sqrt{\frac{\sum X_i^2 - \frac{(\sum X_i)^2}{n}}{n - 1}}. \quad (4.20)$$

Examples 4.1 and 4.2 demonstrate the calculation of s . This quantity frequently is abbreviated SD, and on rare occasions is called the *root mean square deviation* or *root mean square*. Remember that the standard deviation is, by definition, always a nonnegative quantity.[§] The end of Section 6.2 will explain how to determine

*It was the great English statistician Karl Pearson (1857–1936) who coined the term *standard deviation* and its symbol, σ , in 1893, prior to which this quantity was called the *mean error* (Eells, 1926; Walker, 1929: 54–55, 183, 188). In early literature (e.g., by G. U. Yule in 1919), it was termed *root mean square deviation* and acquired the symbol s , and (particularly in the fields of education and psychology) it was occasionally computed using deviations from the median (or even the mode) instead of from the mean (Eells, 1926).

[†]The square root sign ($\sqrt{\quad}$) was introduced by Silesian-born Austrian mathematician Christoff Rudolf (1499–1545) in 1525; by 1637 René Descartes (1596–1650) combined this with a vinculum (a horizontal bar placed above quantities to group them as is done with parentheses or brackets) to obtain the symbol $\sqrt{\quad}$, but Gottfried Wilhelm Leibniz (1646–1716) preferred $\sqrt{(\quad)}$, which is still occasionally seen (Cajori, 1928/9, Vol. I: 135, 208, 368, 372, 375). The first footnote in Section 3.4 speaks to the origin of the cube root symbol ($\sqrt[3]{\quad}$).

[‡]The sample s is actually a slightly biased estimate of the population σ , in that on the average it is a slightly low estimate, especially in small samples. But this fact is generally considered to be offset by the statistic's usefulness. Correction for this bias is sometimes possible (e.g., Bliss, 1967: 131; Dixon and Massey, 1969: 136; Gurland and Tripathi, 1971; Tolman, 1971), but it is rarely employed.

[§]It can be shown that the median of a distribution is never more than one standard deviation away from the mean (μ); that is,

$$|\text{median} - \mu| \leq \sigma \quad (4.21)$$

(Hotelling and Solomon, 1932; O'Connell, 1990; Page and Murty, 1982; Watson, 1994). This is a special case, where $p = 50$, of the relationship

$$\mu - \sigma \sqrt{\frac{1 - p/100}{p/100}} \leq X_p \leq \mu + \sigma \sqrt{\frac{p/100}{1 - p/100}}, \quad (4.22)$$

where X_p is the p th percentile of the distribution (Dharmadhikari, 1991). Also, Page and Murty (1982) have shown these population-parameter relationships between the standard deviation and the range and between the standard deviation and the mean, median, and mode:

$$\text{range}/\sqrt{2n} \leq \sigma \leq \text{range}/2; \quad (4.22a)$$

$$|\text{mode} - \mu| \leq \sigma\sqrt{n/m} \text{ and } |\text{mode} - \text{median}| \leq \sigma(n/m), \quad (4.22b)$$

where m is the number of data at the modal value.

the number of decimal places that may appropriately be recorded for the standard deviation.

4.6 THE COEFFICIENT OF VARIATION

The *coefficient of variation** or *coefficient of variability*, is defined as

$$V = \frac{s}{\bar{X}} \quad \text{or} \quad V = \frac{s}{\bar{X}} \cdot 100\%. \quad (4.23)$$

As s/\bar{X} is generally a small quantity, it is frequently multiplied by 100% in order to express V as a percentage. (The coefficient of variation is often abbreviated as CV.)

As a measure of variability, the variance and standard deviation have magnitudes that are dependent on the magnitude of the data. Elephants have ears that are perhaps 100 times larger than those of mice. If elephant ears were no more variable, relative to their size, than mouse ears, relative to their size, the standard deviation of elephant ear lengths would be 100 times as great as the standard deviation of mouse ear lengths (and the variance of the former would be $100^2 = 10,000$ times the variance of the latter). The sample coefficient of variation expresses sample variability relative to the mean of the sample (and is on rare occasion referred to as the “relative standard deviation”). It is called a measure of *relative variability* or *relative dispersion*.

Because s and \bar{X} have identical units, V has no units at all, a fact emphasizing that it is a relative measure, divorced from the actual magnitude or units of measurement of the data. Thus, had the data in Example 4.2 been measured in pounds, kilograms, or tons, instead of grams, the calculated V would have been the same. The coefficient of variation of a sample, namely V , is an estimate of the coefficient of variation of the population from which the sample came (i.e., an estimate of σ/μ). The coefficient of variation may be calculated only for ratio scale data; it is, for example, not valid to calculate coefficients of variation of temperature data measured on the Celsius or Fahrenheit temperature scales. Simpson, Roe, and Lewontin (1960: 89–95) present a good discussion of V and its biological application, especially with regard to zoomorphological measurements.

4.7 INDICES OF DIVERSITY

For nominal-scale data there is no mean or median or ordered measurements to serve as a reference for discussion of dispersion. Instead, we can invoke the concept of *diversity*, the distribution of observations among categories. Consider that sparrows are found to nest in four different types of location (vines, eaves, branches, and cavities). If, out of twenty nests observed, five are found at each of the four locations, then we would say that there was great diversity in nesting sites. If, however, seventeen nests were found in cavities and only one in each of the other three locations, then we would consider the situation to be one of very low nest-site diversity. In other words, observations distributed evenly among categories display high diversity, whereas a set of observations where most of the data occur in very few of the categories is one exhibiting low diversity.

A large number of diversity measures have been introduced, especially for ecological data (e.g., Brower, Zar, and von Ende, 1998: 177–184; Magurran, 2004), a few of which are presented here.

*The term *coefficient of variation* was introduced by the statistical giant Karl Pearson (1857–1936) in 1896 (David, 1995). In early literature the term was variously applied to the ratios of different measures of dispersion and different measures of central tendency (Eells, 1926).

Among the quantitative descriptions of diversity available are those based on a field known as *information theory*.^{*} The underlying considerations of these measures can be visualized by considering *uncertainty* to be synonymous with diversity. If seventeen out of twenty nest sites were to be found in cavities, then one would be relatively certain of being able to predict the location of a randomly encountered nest site. However, if nests were found to be distributed evenly among the various locations (a situation of high nest-site diversity), then there would be a good deal of uncertainty involved in predicting the location of a nest site selected at random. If a set of nominal scale data may be considered to be a random sample, then a quantitative expression appropriate as a measure of diversity is that of Shannon (1948):

$$H' = - \sum_{i=1}^k p_i \log p_i \quad (4.24)$$

(often referred to as the Shannon-Wiener diversity index or the Shannon-Weaver index). Here, k is the number of categories and p_i is the proportion of the observations found in category i . Denoting n to be sample size and f_i to be the number of observations in category i , then $p_i = f_i/n$; and an equivalent equation for H' is

$$H' = \frac{n \log n - \sum_{i=1}^k f_i \log f_i}{n}, \quad (4.25)$$

a formula that is easier to use than Equation 4.24 because it eliminates the necessity of calculating the proportions (p_i). Published tables of $n \log n$ and $f_i \log f_i$ are available (e.g., Brower, Zar, and von Ende, 1998: 181; Lloyd, Zar, and Karr, 1968). Any logarithmic base may be used to compute H' ; bases 10, e , and 2 (in that order of commonness) are the most frequently encountered. A value of H' (or of any other measure of this section except evenness measures) calculated using one logarithmic base may be converted to that of another base; Table 4.1 gives factors for doing this for bases 10, e , and 2. Unfortunately, H' is known to be an underestimate of the diversity in the sampled population (Bowman et al., 1971). However, this bias decreases with increasing sample size. Ghent (1991) demonstrated a relationship between H' and testing hypotheses for equal abundance among the k categories.

The magnitude of H' is affected not only by the distribution of the data but also by the number of categories, for, theoretically, the maximum possible diversity for a set of data consisting of k categories is

$$H'_{\max} = \log k. \quad (4.26)$$

Therefore, some users of Shannon's index prefer to calculate

$$J' = \frac{H'}{H'_{\max}} \quad (4.27)$$

instead of (or in addition to) H' , thus expressing the observed diversity as a proportion of the maximum possible diversity. The quantity J' has been termed *evenness* (Pielou, 1966) and may also be referred to as *homogeneity* or *relative diversity*. The measure

^{*}Claude Elwood Shannon (1916–2001) founded what he first called “a mathematical theory of communication” and has become known as “information theory.”

TABLE 4.1: Multiplication Factors for Converting among Diversity Measures (H , H' , H_{\max} , or H'_{\max}) Calculated Using Different Logarithmic Bases*

To convert to:	To convert from:		
	Base 2	Base e	Base 10
Base 2	1.0000	1.4427	3.3219
Base e	0.6931	1.0000	2.3026
Base 10	0.3010	0.4343	1.0000

For example, if $H' = 0.255$ using base 10; H' would be $(0.255)(3.3219) = 0.847$ using base 2.

*The measures J and J' are unaffected by change in logarithmic base.

1 - J' may then be viewed as a measure of *heterogeneity*; it may also be considered a measure of *dominance*, for it reflects the extent to which frequencies are concentrated in a small number of categories. The number of categories in a sample (k) is typically an underestimate of the number of categories in the population from which the sample came, because some categories (especially the rarer ones) are likely to be missed in collecting the sample. Therefore, the sample evenness, J' , is typically an overestimate of the population evenness. (That is, J' is a biased statistic.) Example 4.3 demonstrates the calculation of H' and J' .

If a set of data may not be considered a random sample, then Equation 4.24 (or 4.25) is not an appropriate diversity measure (Pielou, 1966). Examples of such

EXAMPLE 4.3 Indices of Diversity for Nominal Scale Data: The Nesting Sites of Sparrows

Category (i)	Observed Frequencies (f_i)
	<i>Sample 1</i>
Vines	5
Eaves	5
Branches	5
Cavities	5

$$\begin{aligned}
 H' &= \frac{n \log n - \sum f_i \log f_i}{n} = [20 \log 20 - (5 \log 5 + 5 \log 5 + 5 \log 5 + 5 \log 5)]/20 \\
 &= [26.0206 - (3.4949 + 3.4949 + 3.4949 + 3.4949)]/20 \\
 &= 12.0410/20 = 0.602
 \end{aligned}$$

$$H'_{\max} = \log 4 = 0.602$$

$$J' = \frac{0.602}{0.602} = 1.00$$

Sample 2

<i>Vines</i>	1
<i>Eaves</i>	1
<i>Branches</i>	1
<i>Cavities</i>	17

$$\begin{aligned}
 H' &= \frac{n \log n - \sum f_i \log f_i}{n} = [20 \log 20 - (1 \log 1 + 1 \log 1 + 1 \log 1 \\
 &\quad + 17 \log 17)]/20 \\
 &= [26.0206 - (0 + 0 + 0 + 20.9176)]/20 \\
 &= 5.1030/20 = 0.255 \\
 H'_{\max} &= \log 4 = 0.602 \\
 J' &= \frac{0.255}{0.602} = 0.42
 \end{aligned}$$

Sample 3

<i>Vines</i>	2
<i>Eaves</i>	2
<i>Branches</i>	2
<i>Cavities</i>	34

$$\begin{aligned}
 H' &= \frac{n \log n - \sum f_i \log f_i}{n} = [40 \log 40 - (2 \log 2 + 2 \log 2 + 2 \log 2 \\
 &\quad + 34 \log 34)]/40 \\
 &= [64.0824 - (0.6021 + 0.6021 + 0.6021 \\
 &\quad + 52.0703)]/40 \\
 &= 10.2058/40 = 0.255 \\
 H'_{\max} &= \log 4 = 0.602 \\
 J' &= \frac{0.255}{0.602} = 0.42
 \end{aligned}$$

situations may be when we have, in fact, data composing an entire population, or data that are a sample obtained nonrandomly from a population. In such a case, one may use the information-theoretic diversity measure of Brillouin (1962: 7–8):*

$$H = \frac{\log \left(\frac{n!}{\prod_{i=1}^k f_i!} \right)}{n}, \quad (4.28)$$

*The notation $n!$ is read as “ n factorial” and signifies the product $(n)(n - 1)(n - 2) \cdots (2)(1)$. It was proposed by French physician and mathematician Christian Kramp (1760–1826) around 1798; he originally called this function *faculty* (“*facultés*” in French) but in 1808 accepted the term *factorial* (“*factorielle*” in French) used by Alsatian mathematician Louis François Antoine Arbogast (1759–1803) (Cajori, 1928/9, Vol. II: 72; Gullberg, 1997: 106; Miller, 2004a; O’Connor and Robertson, 1997). English mathematician Augustus De Morgan (1806–1871) decried the adoption of this symbol as a “barbarism” because it introduced into mathematics a symbol that already had an established meaning in written language, thus giving “the appearance of expressing surprise or admiration” in a mathematical result (Cajori, *ibid.*: 328).

where Π (capital Greek pi) means to take the product, just as Σ means to take the sum. Equation 4.28 may be written, equivalently, as

$$H = \frac{\log \frac{n!}{f_1!f_2! \dots f_k!}}{n} \quad (4.29)$$

or as

$$H = \frac{(\log n! - \sum \log f_i!)}{n}. \quad (4.30)$$

Table B.40 gives logarithms of factorials to ease this calculation. Other such tables are available, as well (e.g., Brower, Zar, and von Ende 1998: 183; Lloyd, Zar, and Karr, 1968; Pearson and Hartly, 1966: Table 51).* Ghent (1991) discussed the relationship between H and the test of hypotheses about equal abundance among k categories.

The maximum possible Brillouin diversity for a set of n observations distributed among k categories is

$$H_{\max} = \frac{\log n! - (k - d) \log c! - d \log (c + 1)!}{n}, \quad (4.35)$$

where c is the integer portion of n/k , and d is the remainder. (For example, if $n = 17$ and $k = 4$, then $n/k = 17/4 = 4.25$ and $c = 4$ and $d = 0.25$.) The Brillouin-based evenness measure is, therefore,

$$J = \frac{H}{H_{\max}}, \quad (4.36)$$

with $1 - J$ being a dominance measure. When we consider that we have data from an entire population, k is a population measurement, rather than an estimate of one, and J is not a biased estimate as is J' .

For further considerations of these and other diversity measures, see Brower, Zar, and von Ende (1998: Chapter 5B) and Magguran (2004: 100–121).

4.8 CODING DATA

Section 3.5 showed how coding data may facilitate statistical computations of measures of central tendency. Such benefits are even more apparent when calculating SS , s^2 ,

*For moderate to large n (or f_i), “Stirling’s approximation” is excellent (see note after Table B.40):

$$n! = \sqrt{2\pi n} (n/e)^n = \sqrt{2\pi} \sqrt{n} e^{-n} n^n, \quad (4.31)$$

of which this is an easily usable derivation:

$$\log n! = (n + 0.5) \log n - 0.434294n + 0.399090. \quad (4.32)$$

An approximation with only half the error of the above is

$$n! = \sqrt{2\pi} \left(\frac{n + 0.5}{e} \right)^{n+0.5} \quad (4.33)$$

and

$$\log n! = (n + 0.5) \log(n + 0.5) - 0.434294(n + 0.5) + 0.399090. \quad (4.34)$$

This is named for James Stirling, who published something similar to the latter approximation formula in 1730, making an arithmetic improvement in the approximation earlier known by Abraham de Moivre (Kemp, 1989; Pearson, 1924; Walker, 1929: 16).

and s , because of the labor, and concomitant chances of error, involved in the unwieldy squaring of large or small numbers.

When data are coded by adding or subtracting a constant (call it A), the measures of dispersion of Sections 4.1 through 4.5 are not changed from what they were for the data before coding. This is because these measures are based upon deviations, and deviations are not changed by moving the data along the measurement scale (e.g., the deviation between 1 and 10 is the same as the deviation between 11 and 20). Sample 1 in Example 4.4 demonstrates this.

However, when coding by multiplying by a constant (call it M), the measures of dispersion are affected, for the magnitudes of the deviations will be changed. With such coding, the range, mean deviation, and standard deviation are changed by a factor of M , in the same manner as the arithmetic mean and the median are, whereas the sum of squares and variance are changed in accordance with the square of the coding constant (i.e., M^2), and the coefficient of variance is not affected. This is demonstrated in Sample 2 of Example 4.4.

Appendix C presents the results of coding these and many other statistics, where a coded datum is described as

$$[X_i] = MX_i + A. \quad (4.37)$$

EXAMPLE 4.4 Coding Data to Facilitate the Calculation of Measures of Dispersion			
Sample 1 (Coding by Subtraction: $A = -840$ g)			
<i>Without Coding X_i</i>		<i>Using Coding $[X_i]$</i>	
X_i (g)	X_i^2 (g ²)	$[X_i]$ (g)	$[X_i]^2$ (g ²)
842	708,964	2	4
843	710,649	3	9
844	712,336	4	16
846	715,716	6	36
846	715,716	6	36
847	717,409	7	49
848	719,104	8	64
849	720,801	9	81
$\sum X_i = 6765$ g		$\sum [X_i] = 45$ g	
$\sum X_i^2 = 5,720,695$ g ²		$\sum [X_i]^2 = 295$ g ²	
$s^2 = \frac{5720695 \text{ g}^2 - \frac{(6765 \text{ g})^2}{8}}{7}$		$[s^2] = \frac{295 \text{ g}^2 - \frac{(45 \text{ g})^2}{8}}{7}$	
$= 5.98 \text{ g}^2$		$= 5.98 \text{ g}^2$	
$s = 2.45$ g		$[s] = 2.44$ g	
$\bar{X} = 845.6$ g		$[\bar{X}] = 5.6$ g	
$V = \frac{s}{\bar{X}} = \frac{2.45 \text{ g}}{845.6 \text{ g}}$			
$= 0.0029 = 0.29\%$			