Random Sampling and Data Description

CHAPTER OUTLINE

- 6-1 NUMERICAL SUMMARIES
- 6-2 STEM-AND-LEAF DIAGRAMS
- 6-3 FREQUENCY DISTRIBUTIONS AND HISTOGRAMS

- 6-4 BOX PLOTS
- 6-5 TIME SEQUENCE PLOTS
- 6-6 PROBABILITY PLOTS

LEARNING OBJECTIVES

After careful study of this chapter you should be able to do the following:

- Compute and interpret the sample mean, sample variance, sample standard deviation, sample median, and sample range
- 2. Explain the concepts of sample mean, sample variance, population mean, and population variance
- Construct and interpret visual data displays, including the stem-and-leaf display, the histogram, and the box plot
- 4. Explain the concept of random sampling
- 5. Construct and interpret normal probability plots
- Explain how to use box plots and other data displays to visually compare two or more samples of data
- 7. Know how to use simple time series plots to visually display the important features of timeoriented data.

Definition: Sample Mean

If the *n* observations in a sample are denoted by x_1, x_2, \ldots, x_n , the sample mean is

$$\overline{x} = \frac{x_1 + x_2 + \dots + x_n}{n} = \frac{\sum_{i=1}^n x_i}{n}$$
(6-1)

Example 6-1

Let's consider the eight observations collected from the prototype engine connectors from Chapter 1. The eight observations are $x_1 = 12.6$, $x_2 = 12.9$, $x_3 = 13.4$, $x_4 = 12.3$, $x_5 = 13.6$, $x_6 = 13.5$, $x_7 = 12.6$, and $x_8 = 13.1$. The sample mean is

$$\overline{x} = \frac{x_1 + x_2 + \dots + x_n}{n} = \frac{\sum_{i=1}^{8} x_i}{8} = \frac{12.6 + 12.9 + \dots + 13.1}{8}$$
$$= \frac{104}{8} = 13.0 \text{ pounds}$$

A physical interpretation of the sample mean as a measure of location is shown in the dot diagram of the pull-off force data. See Figure 6-1. Notice that the sample mean $\overline{x} = 13.0$ can be thought of as a "balance point." That is, if each observation represents 1 pound of mass placed at the point on the *x*-axis, a fulcrum located at \overline{x} would exactly balance this system of weights. Dr. Saed TARAPIAH Descriptive Statistics 4



Figure 6-1 The sample mean as a balance point for a system of weights.

Population Mean

For a finite population with N measurements, the mean is

$$\mu = \sum_{i=1}^{N} x_i f(x_i) = \frac{\sum_{i=1}^{N} x_i}{\frac{N}{N}}$$

$$(6-2)$$

The sample mean is a reasonable estimate of the population mean.

Definition: Sample Variance

If x_1, x_2, \ldots, x_n is a sample of *n* observations, the sample variance is

$$s^{2} = \frac{\sum_{i=1}^{n} (x_{i} - \overline{x})^{2}}{n - 1}$$
(6-3)

The sample standard deviation, s, is the positive square root of the sample variance.

How Does the Sample Variance Measure Variability?



Figure 6-2 How the sample variance measures variability through the deviations $x_i - \overline{x}$.

Dr. Saed TARAPIAH

Example 6-2

Table 6-1 displays the quantities needed for calculating the sample variance and sample standard deviation for the pull-off force data. These data are plotted in Fig. 6-2. The numerator of s^2 is

$$\sum_{i=1}^{8} (x_i - \overline{x})^2 = 1.60$$

so the sample variance is

$$s^2 = \frac{1.60}{8-1} = \frac{1.60}{7} = 0.2286 \text{ (pounds)}^2$$

and the sample standard deviation is

$$s = \sqrt{0.2286} = 0.48$$
 pounds

Dr. Saed TARAPIAH

Descriptive Statistics

9

Table 6-1 Calculation of Terms for the Sample Variance and Sample Standard Deviation

i	x_i	$x_i - \overline{x}$	$(x_i - \overline{x})^2$
1	12.6	-0.4	0.16
2	12.9	-0.1	0.01
3	13.4	0.4	0.16
4	12.3	-0.7	0.49
5	13.6	0.6	0.36
6	13.5	0.5	0.25
7	12.6	-0.4	0.16
8	13.1	0.1	0.01
	104.0	0.0	1.60

Computation of s²

$$s^{2} = \frac{\sum_{i=1}^{n} x_{i}^{2} - \frac{\left(\sum_{i=1}^{n} x_{i}\right)^{2}}{n}}{n-1}$$

Dr. Saed TARAPIAH

Descriptive Statistics

(6-4)

Population Variance

When the population is finite and consists of N values, we may define the population variance as

$$\sigma^{2} = \frac{\sum_{i=1}^{N} (x_{i} - \mu)^{2}}{N}$$
(6-5)

The sample variance is a reasonable estimate of the population variance.

Definition

If the *n* observations in a sample are denoted by x_1, x_2, \ldots, x_n , the sample range is

$$r = \max(x_i) - \min(x_i) \tag{6-6}$$

A stem-and-leaf diagram is a good way to obtain an informative visual display of a data set x_1, x_2, \ldots, x_n , where each number x_i consists of at least two digits. To construct a stem-and-leaf diagram, use the following steps.

Steps for Constructing a Stem-and-Leaf Diagram

- Divide each number x_t into two parts: a stem, consisting of one or more of the leading digits and a leaf, consisting of the remaining digit.
- (2) List the stem values in a vertical column.
- (3) Record the leaf for each observation beside its stem.
- (4) Write the units for stems and leaves on the display.

Example 6-4

To illustrate the construction of a stem-and-leaf diagram, consider the alloy compressive strength data in Table 6-2. We will select as stem values the numbers 7, 8, 9, ..., 24. The resulting stem-and-leaf diagram is presented in Fig. 6-4. The last column in the diagram is a frequency count of the number of leaves associated with each stem. Inspection of this display immediately reveals that most of the compressive strengths lie between 110 and 200 psi and that a central value is somewhere between 150 and 160 psi. Furthermore, the strengths are distributed approximately symmetrically about the central value. The stem-and-leaf diagram enables us to determine quickly some important features of the data that were not immediately obvious in the original display in Table 6-2.

Table 6-2	Comp	ressive Streng	gth (in psi)	of 80 Alumi	inum-Lithiu	m Alloy Spe	cimens
105	221	183	186	121	181	180	143
97	154	153	174	120	168	167	141
245	228	174	199	181	158	176	110
163	131	154	115	160	208	158	133
207	180	190	193	194	133	156	123
134	178	76	167	184	135	229	146
218	157	101	171	165	172	158	169
199	151	142	163	145	171	148	158
160	175	149	87	160	237	150	135
196	201	200	176	150	170	118	149

...

Dr. Saed TARAPIAH

Figure 6-4 Stem-andleaf diagram for the compressive strength data in Table 6-2.

Stem	Leaf	Frequency
7	6	1
8	7	1
9	7	1
10	5 1	2
11	580	3
12	103	3
13	413535	6
14	29583169	8
15	471340886808	12
16	3073050879	10
17	8544162106	10
18	0361410	7
19	960934	6
20	7108	4
21	8	1
22	189	3
23	7	1
24	5	1

Stem : Tens and hundreds digits (psi); Leaf: Ones digits (psi)

Example 6-5

Figure 6-5 illustrates the stem-and-leaf diagram for 25 observations on batch yields from a chemical process. In Fig. 6-5(a) we have used 6, 7, 8, and 9 as the stems. This results in too few stems, and the stem-and-leaf diagram does not provide much information about the data. In Fig. 6-5(b) we have divided each stem into two parts, resulting in a display that more

	Stem	Leaf	Stem	Leaf	Stem	Leaf	
	6	134556	6L	134	6z	1	
	7	011357889	6U	556	6t	3	
	8	1344788	7L	0113	6f	455	
	9	235	7U	57889	6s	6	
	(;	a)	8L	1344	6e		
			8U	788	7z	011	
Figure 6-5			9L	23	7t	3	
Stom-and-			9U	5	7f	5	
Stemanu-			(b)	7s	7	
leaf displays					7e	889	
for Example					8z	1	
					8t	3	
6-5. Stem:					8f	44	
					8s	7	
iens algits.					8e	88	
Leaf [.] Ones					9z		
					9t	23	
digits.					91	5	
	r				9s	10	
Dr. Saed TARAPIAH	L	Descriptive	Statistics		9e	119	
					(.,	

Figure 6-6 Stemand-leaf diagram from Minitab.

Character Stem-and-Leaf Display

Stem-and-leaf of Strength

N = 80	Leaf Unit = 1.0		
1	7	6	
2	8	7	
3	9	7	
5	10	15	
8	11	058	
11	12	013	
17	13	1 3 3 4 5 5	
25	14	12356899	
37	15	001344678888	
(10)	16	0 0 0 3 3 5 7 7 8 9	
33	17	0112445668	
23	18	0 0 1 1 3 4 6	
16	19	034699	
10	20	0178	
6	21	8	
5	22	189	
2	23	7	
1	24	5	

Data Features

• The **median** is a measure of central tendency that divides the data into two equal parts, half below the median and half above. If the number of observations is even, the median is halfway between the two central values.

From Fig. 6-6, the 40th and 41st values of strength as 160 and 163, so the median is (160 + 163)/2 = 161.5. If the number of observations is odd, the median is the *central* value.

The **range** is a measure of variability that can be easily computed from the ordered stem-and-leaf display. It is the maximum minus the minimum measurement. From Fig.6-6 the range is 245 - 76 =169. Dr. Saed TARAPIAH **Descriptive Statistics** 21

Data Features

When an **ordered** set of data is divided into four equal parts, the division points are called **quartiles**.

The **first** or **lower quartile**, q_1 , is a value that has approximately one-fourth (25%) of the observations below it and approximately 75% of the observations above.

The **second quartile**, q_2 , has approximately one-half (50%) of the observations below its value. The second quartile is *exactly* equal to the **median**.

The **third** or **upper quartile**, q_3 , has approximately three-fourths (75%) of the observations below its value. As in the case of the median, the quartiles may not be unique.

Dr. Saed TARAPIAH

Data Features

• The compressive strength data in Figure 6-6 contains n = 80 observations. Minitab software calculates the first and third quartiles as the(n + 1)/4 and 3(n + 1)/4 ordered observations and interpolates as needed.

For example, (80 + 1)/4 = 20.25 and 3(80 + 1)/4 = 60.75.

Therefore, Minitab interpolates between the 20th and 21st ordered observation to obtain $q_1 = 143.50$ and between the 60th and 61st observation to obtain $q_3 = 181.00$.

Data Features

• The **interquartile range** is the difference between the upper and lower quartiles, and it is sometimes used as a measure of variability.

• In general, the 100*k*th **percentile** is a data value such that approximately 100k% of the observations are at or below this value and approximately 100(1 - k)% of them are above it.

- A **frequency distribution** is a more compact summary of data than a stem-and-leaf diagram.
- To construct a frequency distribution, we must divide the range of the data into intervals, which are usually called **class intervals**, **cells**, or **bins**. **Constructing a Histogram (Equal Bin Widths)**:
- (1) Label the bin (class interval) boundaries on a horizontal scale.
- (2) Mark and label the vertical scale with the frequencies or the relative frequencies.
- (3) Above each bin, draw a rectangle where height is equal to the frequency (or relative frequency) corresponding to that bin.





from Minitab with 17 bins.

Dr. Saed TARAPIAH



Figure 6-9 A histogram of the compressive strength data from Minitab with nine bins.

Figure 6-9 A histogram of the compressive strength data from Minitab with nine bins.

Dr. Saed TARAPIAH



Figure 6-10 A cumulative distribution plot of the compressive strength data from Minitab.



Figure 6-11 Histograms for symmetric and skewed distributions.

- The **box plot** is a graphical display that simultaneously describes several important features of a data set, such as center, spread, departure from symmetry, and identification of observations that lie unusually far from the bulk of the data.
- Whisker
- Outlier
- Extreme outlier

6-4 Box Plots



Figure 6-13 Description of a box plot.

6-4 Box Plots



Figure 6-14 Box plot for compressive strength data in Table 6-2. Dr. Saed TARAPIAH Descriptive Statistics

6-4 Box Plots

Figure 6-15

Comparative box plots of a quality index at three plants.



- A **time series** or **time sequence** is a data set in which the observations are recorded in the order in which they occur.
- A **time series plot** is a graph in which the vertical axis denotes the observed value of the variable (say *X*) and the horizontal axis denotes the time (which could be minutes, days, years, etc.).
- When measurements are plotted as a time series, we often see

•trends,

cycles, orother broad features of the data

Dr. Saed TARAPIAH



Figure 6-16 Company sales by year (a) and by quarter (b).



Figure 6-17 A digidot plot of the compressive strength data in Table 6-2.

Dr. Saed TARAPIAH



Figure 6-18 A digidot plot of chemical process concentration readings, observed hourly.

• **Probability plotting** is a graphical method for determining whether sample data conform to a hypothesized distribution based on a subjective visual examination of the data.

• Probability plotting typically uses special graph paper, known as **probability paper**, that has been designed for the hypothesized distribution. Probability paper is widely available for the normal, lognormal, Weibull, and various chi-square and gamma distributions.

Example 6-7

Ten observations on the effective service life in minutes of batteries used in a portable personal computer are as follows: 176, 191, 214, 220, 205, 192, 201, 190, 183, 185. We hypothesize that battery life is adequately modeled by a normal distribution. To use probability plotting to investigate this hypothesis, first arrange the observations in ascending order and calculate their cumulative frequencies (j - 0.5)/10 as shown in Table 6-6.

j	$x_{(j)}$	(j - 0.5)/10	z_j
1	176	0.05	-1.64
2	183	0.15	-1.04
3	185	0.25	-0.67
4	190	0.35	-0.39
5	191	0.45	-0.13
6	192	0.55	0.13
7	201	0.65	0.39
8	205	0.75	0.67
9	214	0.85	1.04
10	220	0.95	1.64
	Daga	mintizza Statistica	

Table 6-6 Calculation for Constructing a Normal Probability Plot

Example 6-7 (continued)

The pairs of values $x_{(j)}$ and (j - 0.5)/10 are now plotted on normal probability paper. This plot is shown in Fig. 6-19. Most normal probability paper plots 100(j - 0.5)/n on the left vertical scale and 100[1 - (j - 0.5)/n] on the right vertical scale, with the variable value plotted on the horizontal scale. A straight line, chosen subjectively, has been drawn through the plotted points. In drawing the straight line, you should be influenced more by the points near the middle of the plot than by the extreme points. A good rule of thumb is to draw the line approximately between the 25th and 75th percentile points. This is how the line in Fig. 6-19 was determined. In assessing the "closeness" of the points to the straight line, imagine a "fat pencil" lying along the line. If all the points are covered by this imaginary pencil, a normal distribution adequately describes the data. Since the points in Fig. 6-19 would pass the "fat pencil" test, we conclude that the normal distribution is an appropriate model.

Figure 6-19 Normal probability plot for battery life.



Figure 6-20 Normal probability plot obtained from standardized normal scores.





Figure 6-21 Normal probability plots indicating a nonnormal distribution. (a) Light-tailed distribution. (b) Heavy-tailed distribution. (c) A distribution with positive (or right) skew.

IMPORTANT TERMS AND CONCEPTS

Box plot Frequency distribution and histogram Median, quartiles and percentiles Multivariable Data Normal probability plot Pareto chart Population mean Population standard deviation Population variance Probability plot Relative Frequency Distribution Sample mean Sample standard deviation Sample variance Stem-and-leaf diagram Time series plots