some without indicating nonnormality. If a sample of 30 observations contains 4 outliers, two of which are extreme, would it be reasonable to assume the population from which the data were collected has a normal frequency curve? Probably not.



```
Stem-and-Leaf Display: C1

Stem-and-leaf of C1 N = 250

Leaf Unit = 1.0

1 1 8

5 2 1378

9 3 3379

17 4 41002567
```

17 2385 9820 (320 904 3242 6 9242 16	45678911123456	$\begin{array}{c} 11223567\\ 13455789\\ 2222444458899\\ 11222233455555667777888889\\ 0000111122334455555666666678888889\\ 1111222334445555555667888899999\\ 00012333444444455566667788889\\ 00111112223344455566668899\\ 000111112223344455556668899\\ 0001111122234446579\\ 011136677778\\ 001133\\ 04669 \end{array}$
6 1	17	04669

The boxplot is better at highlighting outliers than are other displays. The histogram and stem and leaf displays below appear to have the same basic shape as a normal curve (unimodal, symmetric). However, the boxplot shows that we have a dozen outliers in a sample of 250 observations. We would only expect about two outliers in 250 observations when sampling from a population with a normal frequency curve. The frequency curve is best described as unimodal, symmetric, and **heavy-tailed**.



```
Stem-and-Leaf Display: C2
Stem-and-leaf of C2
                Ν
                  = 250
Leaf Unit = 1.0
1
11
45
124
(84)
42
16
7
2
2
1
     6
7
8
9
10
 1
        5
        2578899999
        00011222333333344567777778888889999
        11
        00000011122222333345556689
     12
13
        000113567
        12345
     14
15
16
        3
1
```

Not all symmetric distributions are mound-shaped, as the display below suggests. The boxplot shows symmetry, but the tails of the distribution are shorter (lighter) than in the normal distribution. Note that the distance between quartiles is roughly constant here.



Stem-and-Leaf Display: C3

Stem-and-leaf of C3 $\,\mathbb{N}$ = 250 Leaf Unit = 1.0

29	5	001111223345556666677777899999
56	6	0000011112233345666666777889
82	7	000111233344455555666678889
108	8	111222223445566677788888889
(18)	9	001113334466788889
124	10	0000012234555666667788899999
97	11	000001112233345666688899
73	12	00113334444445555666678899
48	13	00000111233344456777888999
22	14	00012444555666666777999

The mean and median are identical in a population with a (exact) symmetric frequency curve. The histogram and stem and leaf displays for a sample selected from a symmetric population will tend to be fairly symmetric. Further, the sample means and medians will likely be close.

The distribution below is unimodal, and asymmetric or **skewed**. The distribution is said to be **skewed to the right**, or upper end, because the right tail is much longer than the left tail. The boxplot also shows the skewness - the region between the minimum observation and the median contains half the data in less than 1/5 the range of values. In addition, the upper tail contains several outliers.



The distribution below is unimodal and **skewed to the left**. The two examples show that extremely skewed distributions often contain outliers in the longer tail of the distribution.



Not all distributions are unimodal. The distribution below has two modes or peaks, and is said to be **bimodal**. Distributions with three or more peaks are called **multi-modal**.



Stem-and-Leaf Display: C6

Stem-and-leaf of C6 $\,$ N $\,$ = 250 Leaf Unit = 10 $\,$

4	0	2233
12	Ó	44455555
32	Ó	6666677777777777777777
64	0	8888888888888888889999999999999999999
95	1	0000000000000111111111111111111
115	1	2222222222233333333
(15)	1	44444445555555
120	1	6666677777
110	1	88899999999
99	2	0000000000011111111111111111
71	2	222222222222222333333333333333333333333
38	2	444444445555555555555555555555555555555

The boxplot and histogram or stem and leaf display (or dotplot) are used **together** to describe the distribution. The boxplot does not provide information about modality - it only tells you about skewness and the presence of outliers.

As noted earlier, many statistical methods assume the population frequency curve is normal. Small deviations from normality usually do not dramatically influence the operating characteristics of these methods. We worry most when the deviations from normality are severe, such as extreme skewness or heavy tails containing multiple outliers.

Interpretations for Examples

The **MAO** samples are fairly symmetric, unimodal (?), and have no outliers. The distributions do not deviate substantially from normality. The various measures of central location (\overline{Y}, M) are fairly close, which is common with reasonably symmetric distributions containing no outliers.

The **SIDS** sample is unimodal, and skewed to the right due to the presence of four outliers in the upper tail. Although not given, we expect the mean to be noticeably higher than the median (Why?). A normality assumption here is unrealistic.

Example: Length of Stay in a Psychiatric Unit

Data on all 58 persons committed voluntarily to the acute psychiatric unit of a health care canter in Wisconsin during the first six months of a year are stored in the worksheet HCC that installs with Minitab. Two of the variables are Length (of stay, in number of days), and Reason (for discharge, 1=normal, 2=other). It is of interest to see if the length of stay differs for the two types of discharge. The main parts of the boxplots comparing the groups are rather compressed (and not very useful) because outliers are using up all the scale.



The solution in a case like this is to zoom in using the Data Options in the boxplot display. In this case let us exclude rows where Length > 30. Now we have a little more basis for comparison.



Does it look like there is a really large difference between the groups? What would you say about the shape of the distributions? Does it look like these are normally distributed values?

Examine the descriptive statistics. What is a reasonable summary here and what probably is pretty distorted? What is your summary of the data based upon the boxplots and numerical summaries?

Descriptive Statistics: Length

Variable Length	Reason 1 2	N 42 16	N* 0 0	Mean 11.55 6.44	SE Mean 2.25 1.78	StDev 14.60 7.11	Minimum 0.00 0.00	Q1 1.75 1.25	Median 7.50 4.50	Q3 13.50 8.25
Variable Length	Reason 1 2	Maximum 75.00 25.00								