

# PROBABILITY & STATISTICS

---

BS 1402

# Contents

- Organization of Data
  - Frequency Distribution of Data
- Graphical Representation of Qualitative Data
  - Histograms
  - Frequency Polygons
  - Ogives
  - Stem and Leaf Plots

# Frequency Distribution of Data

When the data are in original form, they are called **raw data**

- Since little information can be obtained from looking at raw data, the researcher organizes the data into what is called a *frequency distribution*.
- A frequency distribution consists of *classes* and their corresponding *frequencies*. Each raw data value is placed into a quantitative or qualitative category called a **class**.
- The **frequency** of a class then is the number of data values contained in a specific class.

# Frequency Distribution of Data

A **frequency distribution** is the organization of raw data in table form, using classes and frequencies.

## Example 1:

Suppose a researcher wished to do a study on the ages of the top 50 wealthiest people in the world. The researcher first would have to get the data on the ages of the people. The raw data is listed below:

49	57	38	73	81
74	59	76	65	69
54	56	69	68	78
65	85	49	69	61
48	81	68	37	43
78	82	43	64	67
52	56	81	77	79
85	40	85	59	80
60	71	57	61	69
61	83	90	87	74

## Example 1 (cont.):

The frequency distribution for the preceding data can be shown as:

Class limits	Tally	Frequency
35–41	///	3
42–48	///	3
49–55	////	4
56–62	 	10
63–69	 	10
70–76		5
77–83	 	10
84–90		5
		<hr/>
		Total 50

Now some general observations can be made from looking at the frequency distribution. For example, it can be stated that the majority of the wealthy people in the study are over 55 years old.

# Frequency Distribution of Data

Two types of frequency distributions that are most often used are the *categorical frequency distribution* and the *grouped frequency distribution*.

- **Categorical Frequency Distribution:**

The categorical frequency distribution is used for data that can be placed in specific categories, such as nominal- or ordinal-level data. For example, data such as political affiliation, religious affiliation, or major field of study would use categorical frequency distributions.

# Frequency Distribution of Data

- **Grouped Frequency Distribution:**

When the range of the data is large, the data must be grouped into classes that are more than one unit in width, in what is called a grouped frequency distribution.



## Example 2:

### Categorical Frequency Distribution

Twenty-five army inductees were given a blood test to determine their blood type. The data set is

A	B	B	AB	O
O	O	B	AB	B
B	B	O	A	O
A	O	O	O	AB
AB	A	O	B	A

Construct a frequency distribution for the data.

#### **Solution:**

Since the data are categorical, discrete classes can be used. There are four blood types:

A, B, O, and AB.

## Example 2: (cont.)

**Step 1** Make a table as shown.

A Class	B Tally	C Frequency	D Percent
A			
B			
O			
AB			

**Step 2** Tally the data and place the results in column B.

**Step 3** Count the tallies and place the results in column C.

**Step 4** Find the percentage of values in each class by using the formula

$$\% = \frac{f}{n} \cdot 100\%$$

where  $f$  = frequency of the class and  $n$  = total number of values. For example, in the class of type A blood, the percentage is

$$\% = \frac{5}{25} \cdot 100\% = 20\%$$

Percentages are not normally part of a frequency distribution, but they can be added since they are used in certain types of graphs such as pie graphs. Also, the decimal equivalent of a percent is called a *relative frequency*.

## Example 2: (cont.)

**Step 5** Find the totals for columns C (frequency) and D (percent). The completed table is shown.

A Class	B Tally	C Frequency	D Percent
A		5	20
B	//	7	28
O	//	9	36
AB		<u>4</u>	<u>16</u>
		Total 25	100

For the sample, more people have type O blood than any other type.

## Example 3:

### Grouped Frequency Distribution

A distribution of the number of hours that boat batteries lasted is the following.

Class limits	Class boundaries	Tally	Frequency
24–30	23.5–30.5	///	3
31–37	30.5–37.5	/	1
38–44	37.5–44.5	////	5
45–51	44.5–51.5	//// //	9
52–58	51.5–58.5	//// /	6
59–65	58.5–65.5	/	1
			<hr/> 25

In this distribution, the values 24 and 30 of the first class are called class limits. The lower class limit is 24; it represents the smallest data value that can be included in the class. The upper class limit is 30; it represents the largest data value that can be included in the class. The numbers in the second column are called class boundaries. These numbers are used to separate the classes so that there are no gaps in the frequency distribution. The gaps are due to the limits; for example, there is a gap between 30 and 31.

## Example 3 (cont.):

$$\text{Lower limit} - 0.5 = 31 - 0.5 = 30.5 = \text{lower boundary}$$

$$\text{Upper limit} + 0.5 = 37 + 0.5 = 37.5 = \text{upper boundary}$$

If the data are in tenths, such as 6.2, 7.8, and 12.6, the limits for a class hypothetically might be 7.8–8.8, and the boundaries for that class would be 7.75–8.85.

Finally, the class width for a class in a frequency distribution is found by subtracting the lower (or upper) class limit of one class from the lower (or upper) class limit of the next class. For example, the class width in the preceding distribution on the duration of boat batteries is 7, found from  $31 - 24 = 7$ . The class width can also be found by subtracting the lower boundary from the upper boundary for any given class. In this case,  $30.5 - 23.5 = 7$

The researcher must decide how many classes to use and the width of each class. To construct a frequency distribution, follow these rules:

1. *There should be between 5 and 20 classes*
2. *It is preferable but not absolutely necessary that the class width be an odd number*

This ensures that the midpoint of each class has the same place value as the data. The class midpoint  $X_m$  is obtained by adding the lower and upper boundaries and dividing by 2, or adding the lower and upper limits and dividing by 2:

$$X_m = \frac{\text{lower boundary} + \text{upper boundary}}{2}$$

or

$$X_m = \frac{\text{lower limit} + \text{upper limit}}{2}$$

For example, the midpoint of the first class in the example with boat batteries is

$$\frac{24 + 30}{2} = 27 \quad \text{or} \quad \frac{23.5 + 30.5}{2} = 27$$

3. *The classes must be mutually exclusive.* Mutually exclusive classes have nonoverlapping class limits so that data cannot be placed into two classes. Manytimes, frequency distributions such as

<u>Age</u>
10–20
20–30
30–40
40–50

are found in the literature or in surveys. If a person is 40 years old, into which class should she or he be placed? A better way to construct a frequency distribution is to use classes such as

<u>Age</u>
10–20
21–31
32–42
43–53

4. *The classes must be continuous.* Even if there are no values in a class, the class must be included in the frequency distribution. There should be no gaps in a frequency distribution. The only exception occurs when the class with a zero frequency is the first or last class. A class with a zero frequency at either end can be omitted without affecting the distribution.

5. *The classes must be exhaustive.* There should be enough classes to accommodate all the data.

6. *The classes must be equal in width.* This avoids a distorted view of the data.

One exception occurs when a distribution has a class that is open-ended. That is, the class has no specific beginning value or no specific ending value. A frequency distribution with an open-ended class is called an **open-ended distribution**. Here are two examples of distributions with open-ended classes.

<u>Age</u>	<u>Frequency</u>	<u>Minutes</u>	<u>Frequency</u>
10–20	3	Below 110	16
21–31	6	110–114	24
32–42	4	115–119	38
43–53	10	120–124	14
54 and above	8	125–129	5



# Graphical Representation of Data

The three most commonly used graphs in research are

1. The histogram.
2. The frequency polygon.
3. The cumulative frequency graph, or ogive

# Histograms

The **histogram** is a graph that displays the data by using contiguous vertical bars (unless the frequency of a class is 0) of various heights to represent the frequencies of the classes.

## Example 4:

### Record High Temperatures

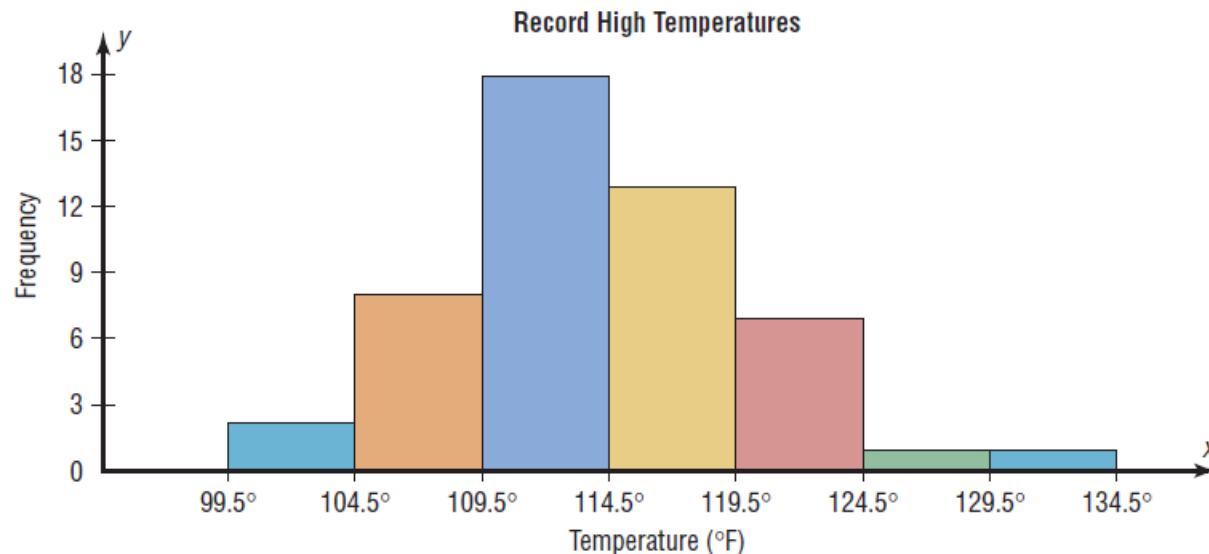
Construct a histogram to represent the data shown for the record high temperatures for each of the 50 states (see Example 2–2).

Class boundaries	Frequency
99.5–104.5	2
104.5–109.5	8
109.5–114.5	18
114.5–119.5	13
119.5–124.5	7
124.5–129.5	1
129.5–134.5	1

## Example 4 (cont.):

### Solution

- Step 1** Draw and label the  $x$  and  $y$  axes. The  $x$  axis is always the horizontal axis, and the  $y$  axis is always the vertical axis.
- Step 2** Represent the frequency on the  $y$  axis and the class boundaries on the  $x$  axis.
- Step 3** Using the frequencies as the heights, draw vertical bars for each class. See Figure 2–2.



# Frequency Polygons

The **frequency polygon** is a graph that displays the data by using lines that connect points plotted for the frequencies at the midpoints of the classes. The frequencies are represented by the heights of the points.

## Example 5:

Construct a frequency polygon to represent the data for the record high temperatures given in previous example (Example 4)

### Solution

**Step 1** Find the midpoints of each class. Recall that midpoints are found by adding the upper and lower boundaries and dividing by 2:

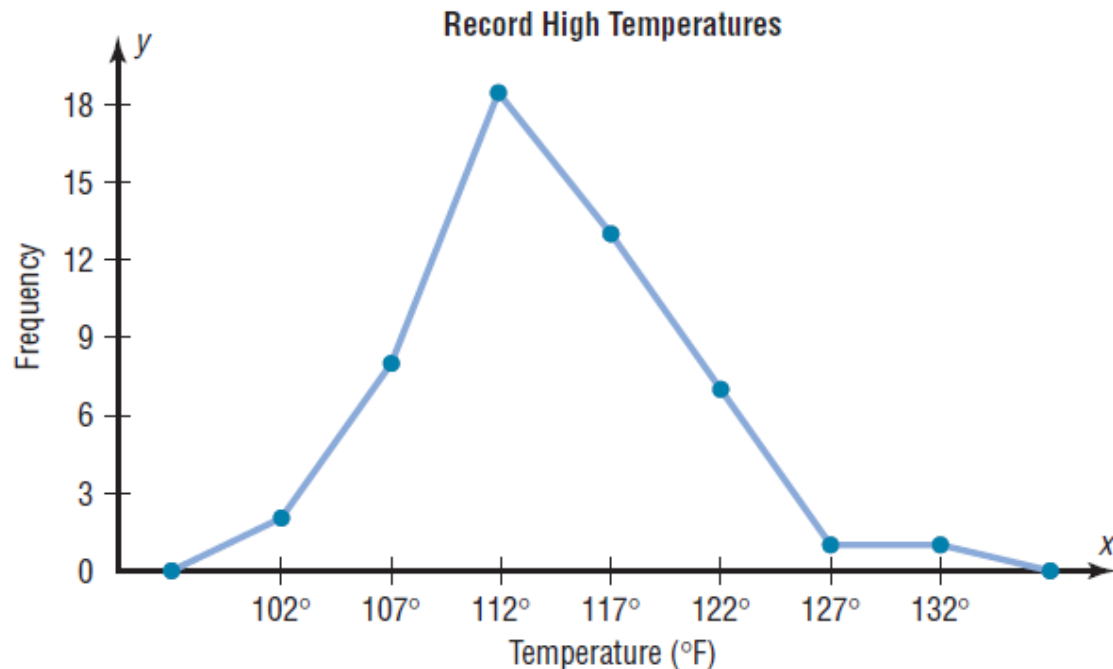
$$\frac{99.5 + 104.5}{2} = 102 \quad \frac{104.5 + 109.5}{2} = 107$$

and so on. The midpoints are

Class boundaries	Midpoints	Frequency
99.5–104.5	102	2
104.5–109.5	107	8
109.5–114.5	112	18
114.5–119.5	117	13
119.5–124.5	122	7
124.5–129.5	127	1
129.5–134.5	132	1

## Example 5 (cont.):

- Step 2** Draw the  $x$  and  $y$  axes. Label the  $x$  axis with the midpoint of each class, and then use a suitable scale on the  $y$  axis for the frequencies.
- Step 3** Using the midpoints for the  $x$  values and the frequencies as the  $y$  values, plot the points.
- Step 4** Connect adjacent points with line segments. Draw a line back to the  $x$  axis at the beginning and end of the graph, at the same distance that the previous and next midpoints would be located, as shown in Figure 2–3.



# Ogives

The third type of graph that can be used represents the cumulative frequencies for the classes. This type of graph is called the *cumulative frequency graph*, or *ogive*. The **cumulative frequency** is the sum of the frequencies accumulated up to the upper boundary of a class in the distribution.

The **ogive** is a graph that represents the cumulative frequencies for the classes in a frequency distribution.



## Example 6:

Construct an ogive to represent the data for the record high temperatures given in Example 4

### Solution

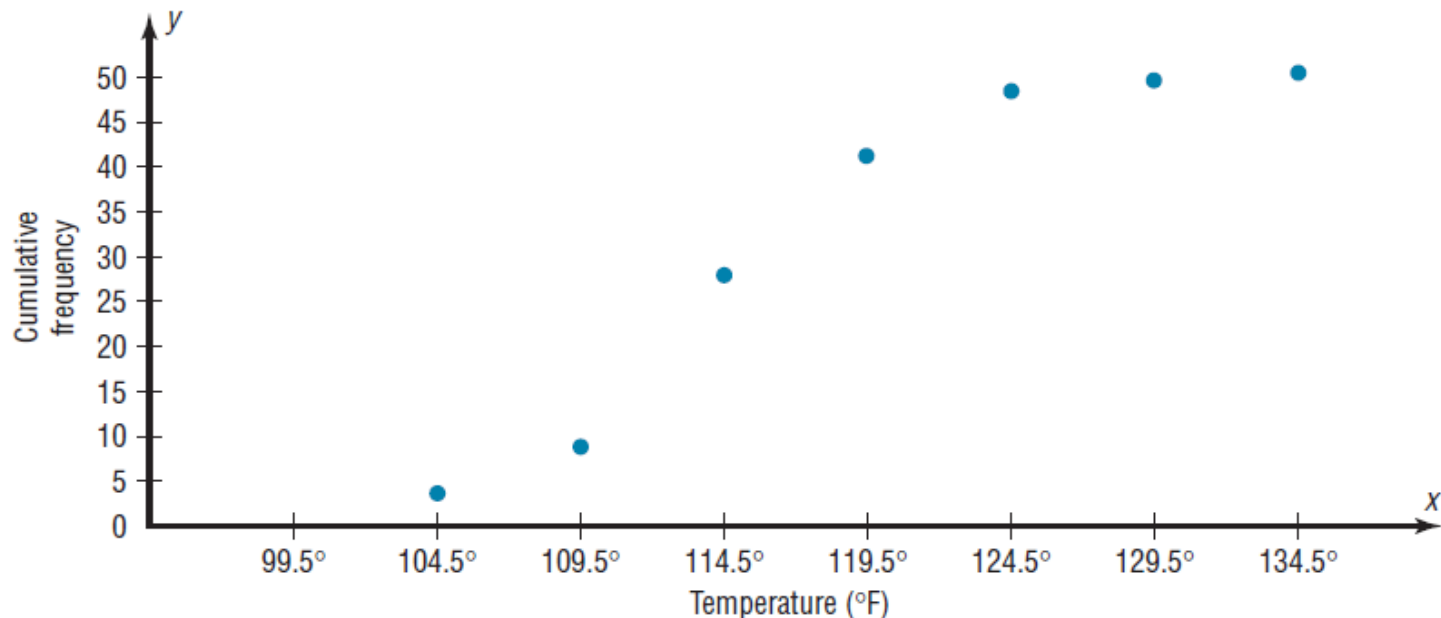
**Step 1** Find the cumulative frequency for each class.

	Cumulative frequency
Less than 99.5	0
Less than 104.5	2
Less than 109.5	10
Less than 114.5	28
Less than 119.5	41
Less than 124.5	48
Less than 129.5	49
Less than 134.5	50

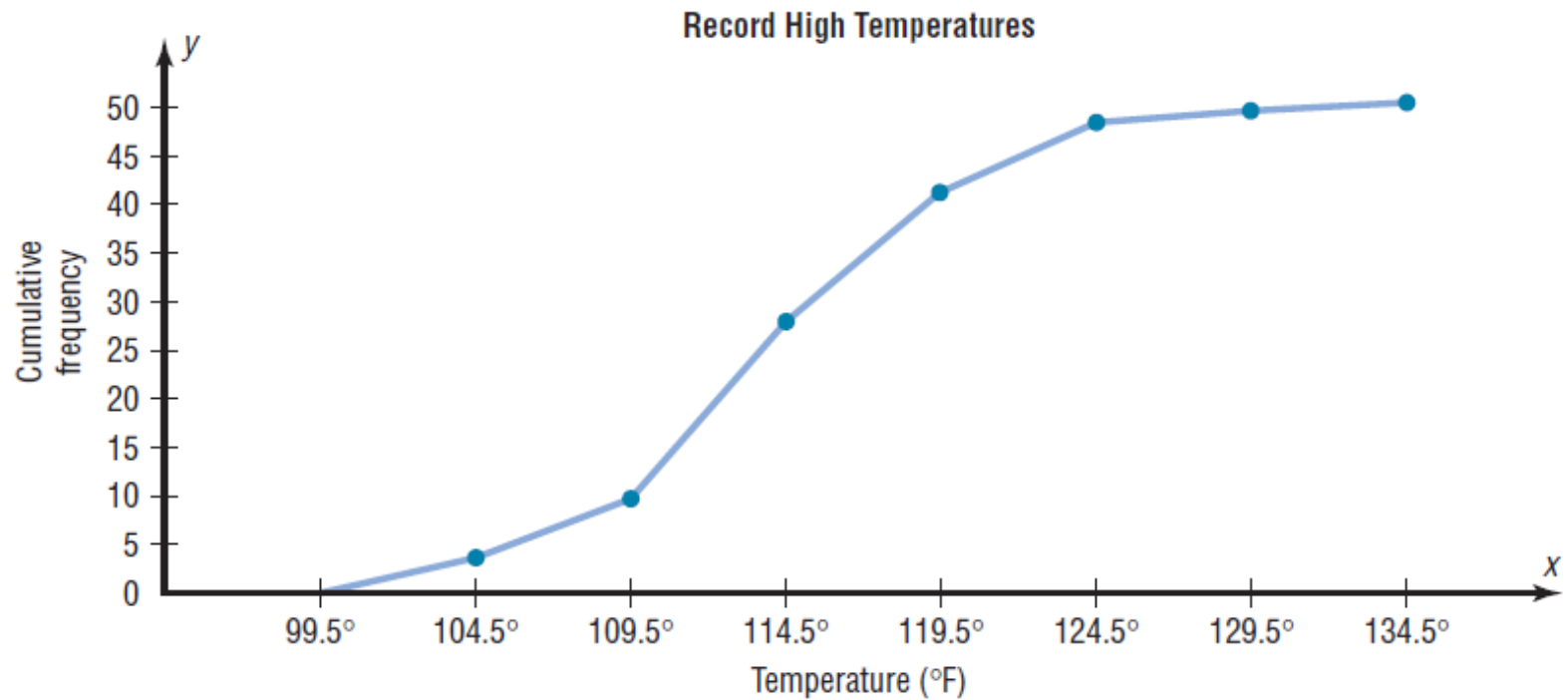
**Step 2** Draw the  $x$  and  $y$  axes. Label the  $x$  axis with the class boundaries. Use an appropriate scale for the  $y$  axis to represent the cumulative frequencies. (Depending on the numbers in the cumulative frequency columns, scales such as 0, 1, 2, 3, . . . , or 5, 10, 15, 20, . . . , or 1000, 2000, 3000, . . . can be used. Do *not* label the  $y$  axis with the numbers in the cumulative frequency column.) In this example, a scale of 0, 5, 10, 15, . . . will be used.

## Example 6 (cont.):

- Step 3** Plot the cumulative frequency at each upper class boundary, as shown in Figure 2–4. Upper boundaries are used since the cumulative frequencies represent the number of data values accumulated up to the upper boundary of each class.
- Step 4** Starting with the first upper class boundary, 104.5, connect adjacent points with line segments, as shown in Figure 2–5. Then extend the graph to the first lower class boundary, 99.5, on the  $x$  axis.

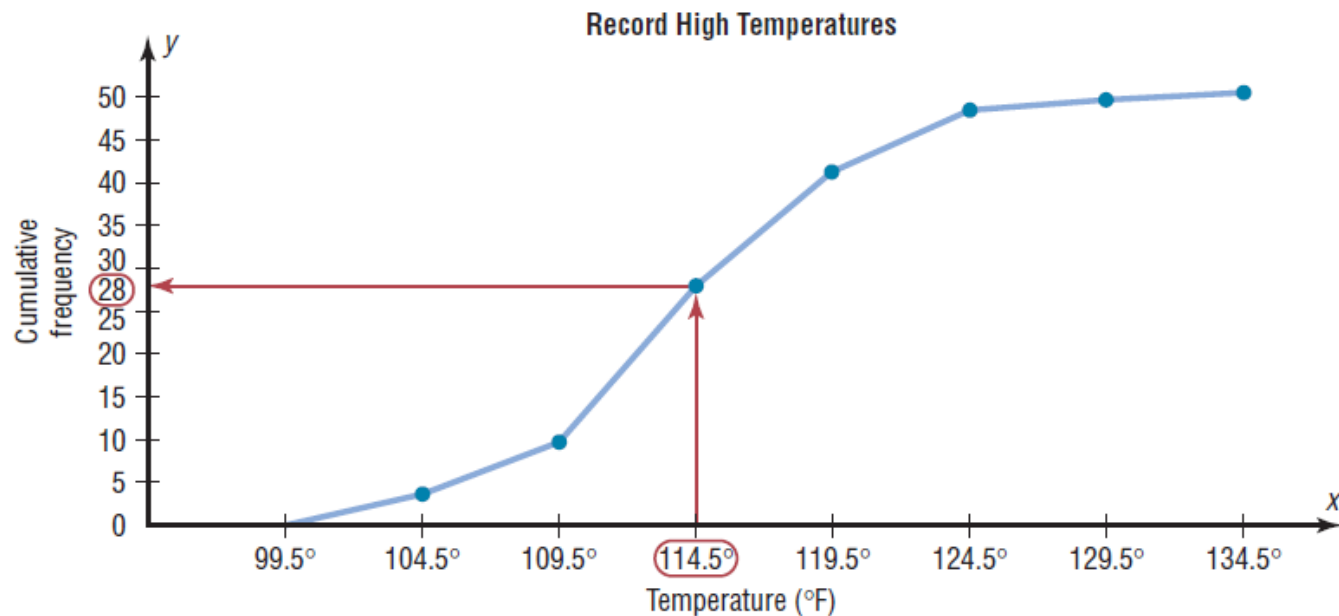


## Example 6 (cont.):



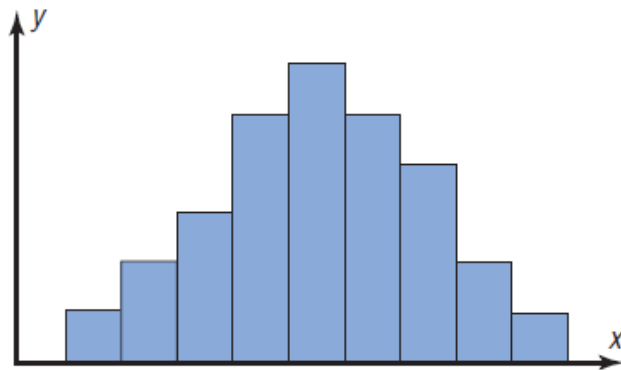
## Example 6 (cont.):

Cumulative frequency graphs are used to visually represent how many values are below a certain upper class boundary. For example, to find out how many record high temperatures are less than 114.5F, locate 114.5F on the x axis, draw a vertical line up until it intersects the graph, and then draw a horizontal line at that point to the y axis.

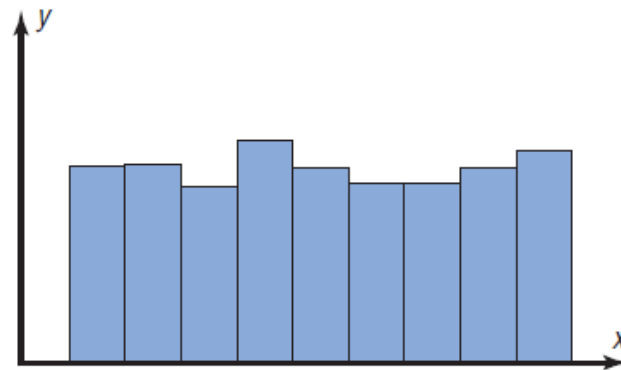


# Distribution Shapes

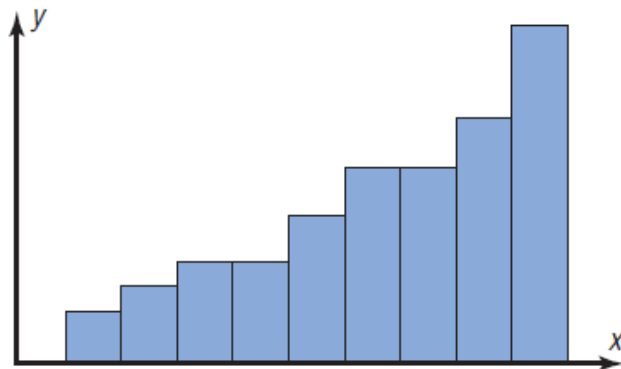
A distribution can have many shapes, and one method of analyzing a distribution is to draw a histogram or frequency polygon for the distribution. Several of the most common shapes are shown in Figure below:



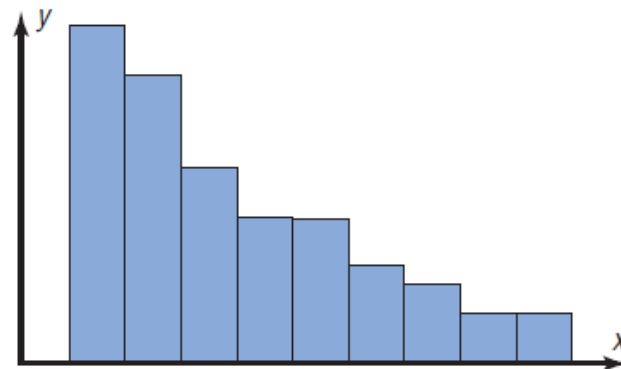
(a) Bell-shaped



(b) Uniform

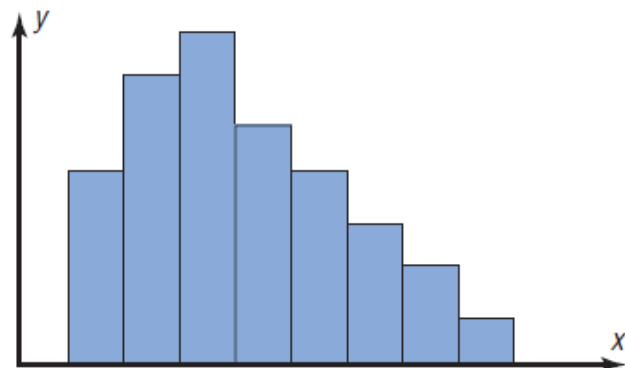


(c) J-shaped

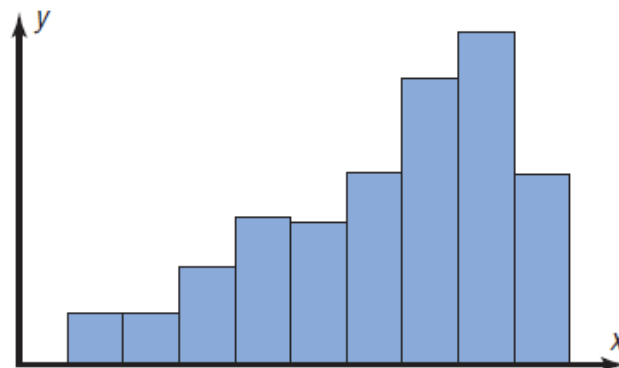


(d) Reverse J-shaped

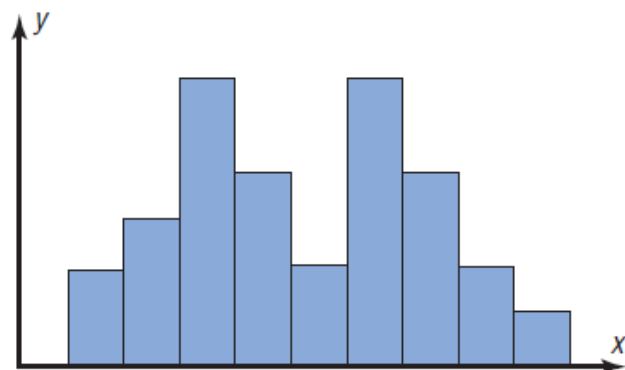
# Distribution Shapes



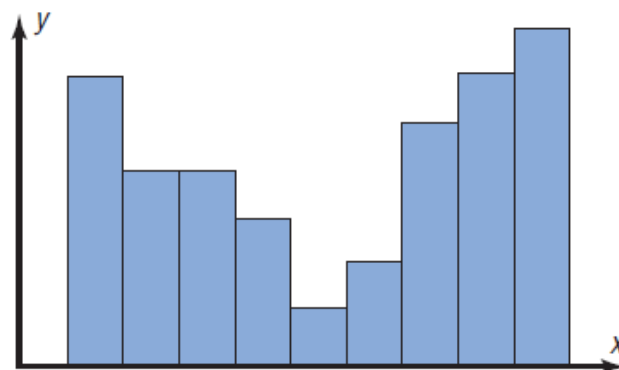
(e) Right-skewed



(f) Left-skewed



(a) Bimodal



(h) U-shaped

# Stem and Leaf Plots

The stem and leaf plot is a method of organizing data and is a combination of sorting and graphing. It has the advantage over a grouped frequency distribution of retaining the actual data while showing them in graphical form.

A **stem and leaf plot** is a data plot that uses part of the data value as the stem and part of the data value as the leaf to form groups or classes.

## Example 7:

At an outpatient testing center, the number of cardiograms performed each day for 20 days is shown. Construct a stem and leaf plot for the data.

25	31	20	32	13
14	43	02	57	23
36	32	33	32	44
32	52	44	51	45

### Solution

**Step 1** Arrange the data in order:

02, 13, 14, 20, 23, 25, 31, 32, 32, 32,  
32, 33, 36, 43, 44, 44, 45, 51, 52, 57

*Note:* Arranging the data in order is not essential and can be cumbersome when the data set is large; however, it is helpful in constructing a stem and leaf plot. The leaves in the final stem and leaf plot should be arranged in order.



## Example 7(cont.):

**Step 2** Separate the data according to the first digit, as shown.

02	13, 14	20, 23, 25	31, 32, 32, 32, 32, 33, 36
43, 44, 44, 45	51, 52, 57		

**Step 3** A display can be made by using the leading digit as the *stem* and the trailing digit as the *leaf*. For example, for the value 32, the leading digit, 3, is the stem and the trailing digit, 2, is the leaf. For the value 14, the 1 is the stem and the 4 is the leaf. Now a plot can be constructed as shown in Figure 2–22.

Leading digit (stem)	Trailing digit (leaf)
0	2
1	3 4
2	0 3 5
3	1 2 2 2 2 3 6
4	3 4 4 5
5	1 2 7

## Example 7 (cont.):

**Figure 2–22**

Stem and Leaf Plot for  
Example 2–13

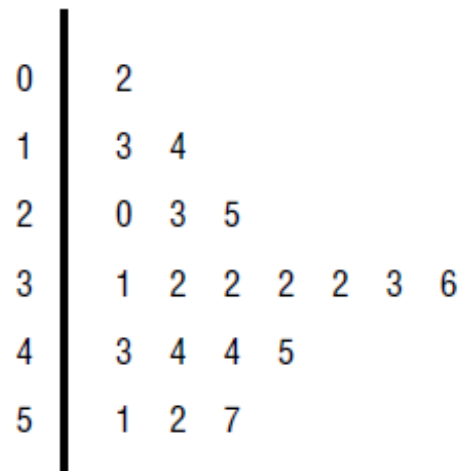


Figure 2–22 shows that the distribution peaks in the center and that there are no gaps in the data. For 7 of the 20 days, the number of patients receiving cardiograms was between 31 and 36. The plot also shows that the testing center treated from a minimum of 2 patients to a maximum of 57 patients in any one day.

If there are no data values in a class, you should write the stem number and leave the leaf row blank. Do not put a zero in the leaf row.

## Example 8:

An insurance company researcher conducted a survey on the number of car thefts in a large city for a period of 30 days last summer. The raw data are shown. Construct a stem and leaf plot by using classes 50–54, 55–59, 60–64, 65–69, 70–74, and 75–79.

52	62	51	50	69
58	77	66	53	57
75	56	55	67	73
79	59	68	65	72
57	51	63	69	75
65	53	78	66	55

## Example 8(cont.):

### Solution

**Step 1** Arrange the data in order.

50, 51, 51, 52, 53, 53, 55, 55, 56, 57, 57, 58, 59, 62, 63,  
65, 65, 66, 66, 67, 68, 69, 69, 72, 73, 75, 75, 77, 78, 79

**Step 2** Separate the data according to the classes.

50, 51, 51, 52, 53, 53      55, 55, 56, 57, 57, 58, 59  
62, 63      65, 65, 66, 66, 67, 68, 69, 69      72, 73  
75, 75, 77, 78, 79

**Step 3** Plot the data as shown here.

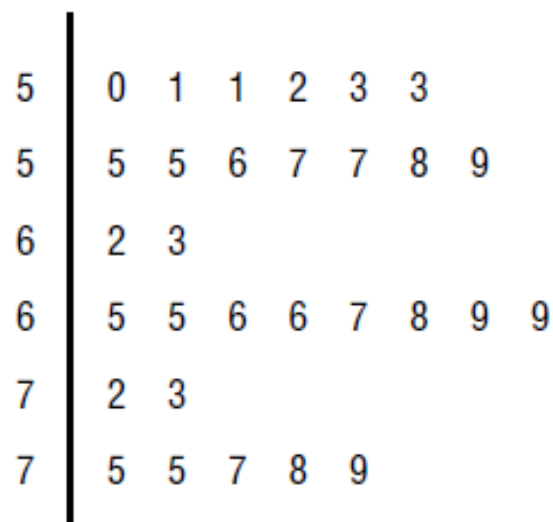
Leading digit (stem)	Trailing digit (leaf)
5	0 1 1 2 3 3
5	5 5 6 7 7 8 9
6	2 3
6	5 5 6 6 7 8 9 9
7	2 3
7	5 5 7 8 9

The graph for this plot is shown in Figure 2–23.

## Example 8(cont.):

**Figure 2-23**

Stem and Leaf Plot for  
Example 2-14



Thank You.