

**CHAPTER 1****Descriptive Statistics**

- 1.1 Introduction**
  - 1.2 Basic concepts**
  - 1.3 Sampling schemes**
  - 1.4 Graphical representation of data**
  - 1.5 Numerical description of data**
  - 1.6 Computers and statistics**
  - 1.7 Chapter summary**
  - 1.8 Computer examples**
- Projects for Chapter 1**

Statistical software R is used for this book. All outputs and codes given are in R. R is a free statistical software, and it can be downloaded from the website: <http://www.r-project.org>

---

## Exercises 1.2

### 1.2.1

The suggested solutions:

For qualitative data we can have color, sex, race, Zip code and so on. For quantitative data we can have age, temperature, time, height, weight and so on. For cross section data we can have school funding for each department in 2000. For time series data we can have the crude oil price from 1995 to 2008.

### 1.2.2

The suggested solutions:

For qualitative data we collect the frequency information of the data and we want to see the comparison by either bar chart or pie chart.

For quantitative data we collect the numerical information of the data and we want to see the comparison by histogram distribution.

For cross section data we collect different section data on the same time and we want to make comparison between them.

For time series data we collect same type of data on different time spot and we want to see if there is any trend or pattern of this data with time shifting.

### 1.2.3

The suggested questions can be:

1. What types of data the amounts are?
2. Do these Federal Agency receive the same amount of funding? If not, why?
3. Which Federal Agency should receive more funding? Why?

The suggested inferences we can make are:

1. These Federal Agency get different amount of money.
2. There are big differences between funding the Agencies receive.

### 1.2.4

The suggested questions can be

1. How does the funding changes for each agency through time?
2. Should we change the proportion between the Agencies or not?
3. Should we increase the total amount or not?

The suggested inferences we can make is

1. The total money tends to be the same.
2. The proportion between the Agencies tends to be the same.

## Exercises 1.3

### 1.3.1

#### Simple Random Sample:

Say we have a population of 1,000 students, and we want a sample of 100 students. Using software or a random table, we randomly select 100 out of the 1,000 students. We want the selection probability for all the students to be equal. That is no student is more likely to be selected than any other student.

#### Systematic Sample:

Again, we have a population of 1,000 students, and we want a sample of 100 students. We need the sampling interval  $k = N/n = 10$ . Now, we need a random starting point between 1 and  $k$ . Let say, we randomly select 4. This gives us the sample: 4, 14, 24, ..., 74, 84, 94. This sample of numbers will correspond to ordered list of students.

#### Stratified Sample:

Suppose we decide to sample 100 college students from the population of 1000 ( that is 10% of the population). We know these 1000 students come from three different major, Math, Computer Science and Social Science. We have Math 200, CS 400 and SS 400 students. Then we choose 10% of each of them Math 20, CS 40 and SS 40 by using simple random sample within each major.

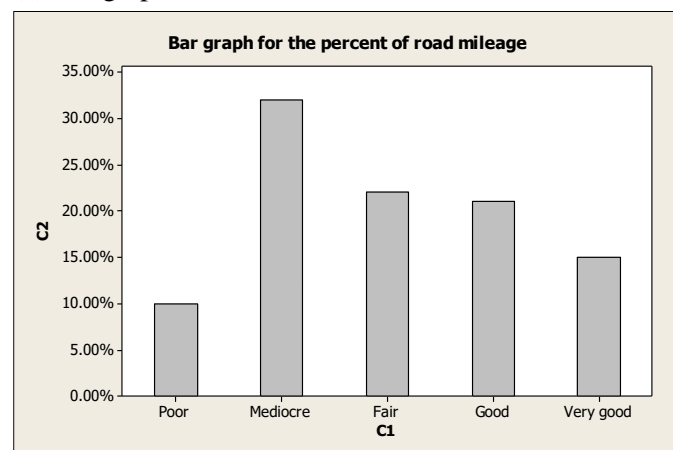
#### Cluster Sample:

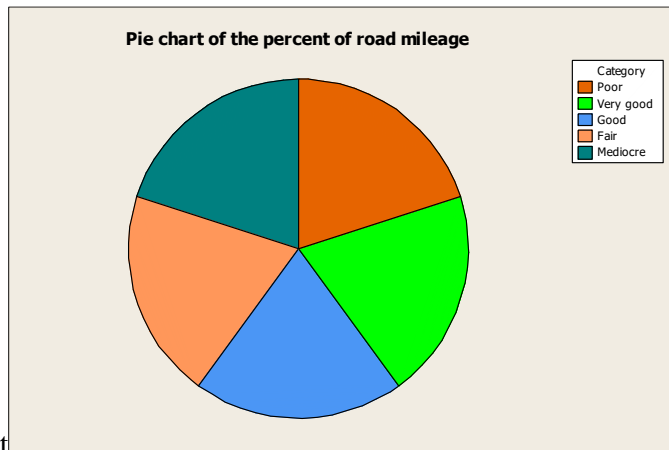
Presume we have a population of 1,000 students clustered into 10 departments. For our sample of students, we will randomly select a subset from the 10 departments. Let say we randomly select 3 out 10 departments. Now, all the students on those 3 department become the sample from the population of students.

## Exercises 1.4

### 1.4.1

#### (a) Bar graph

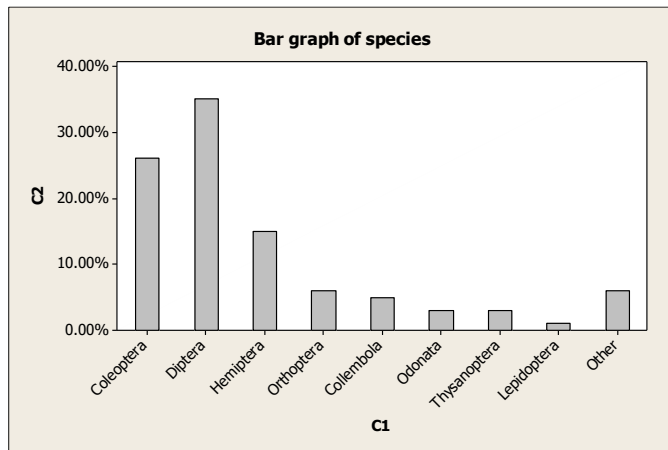




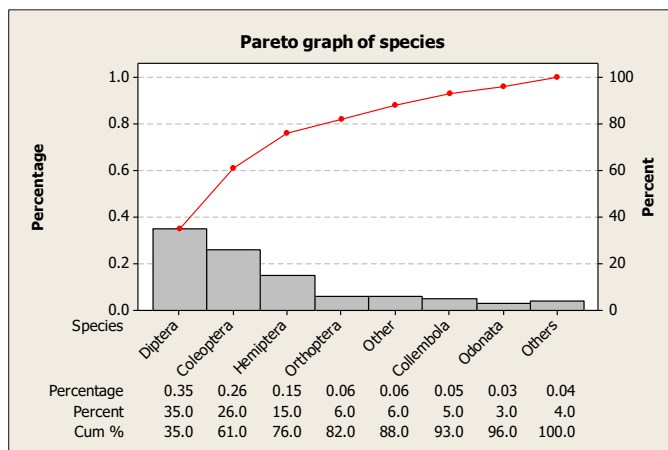
(b) Pie chart

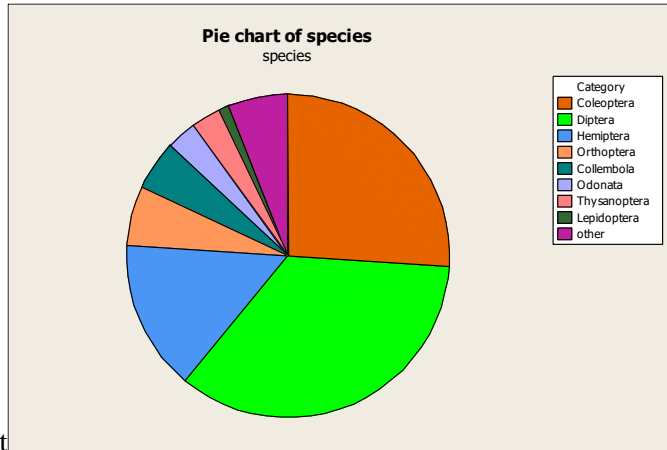
### 1.4.2

(a) Bar graph



(b) Pareto graph

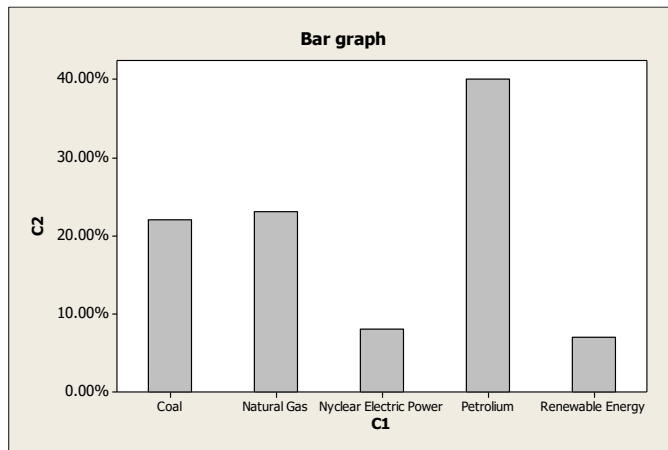




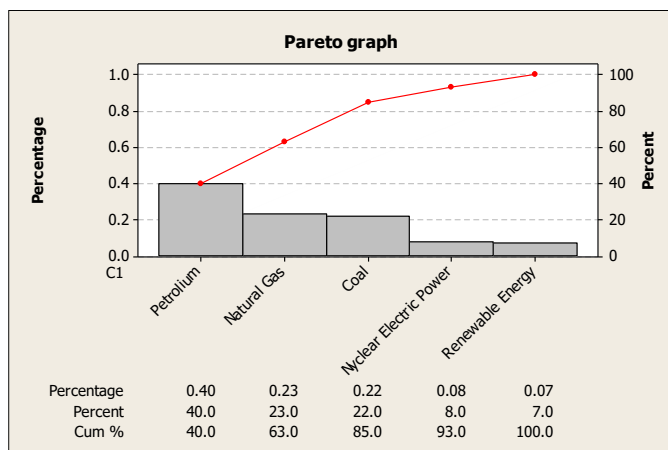
(c) Pie chart

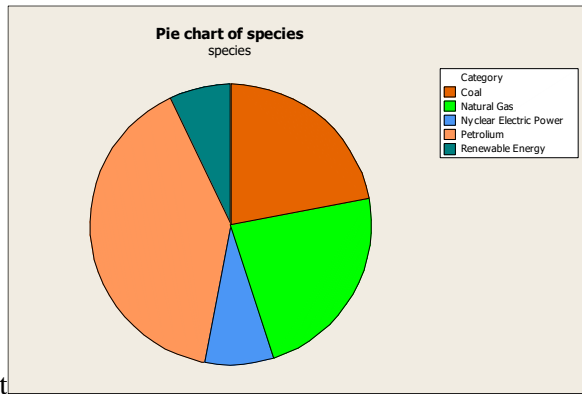
### 1.4.3

(a) Bar graph



(b) Pareto graph

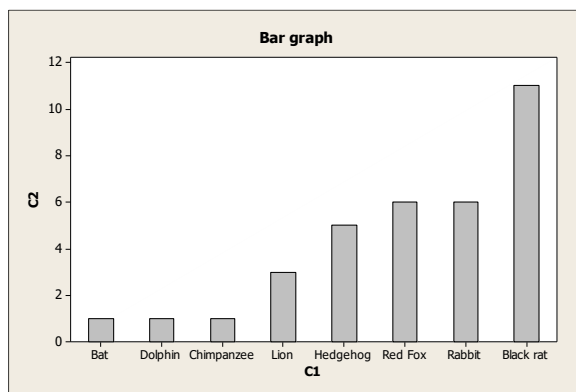




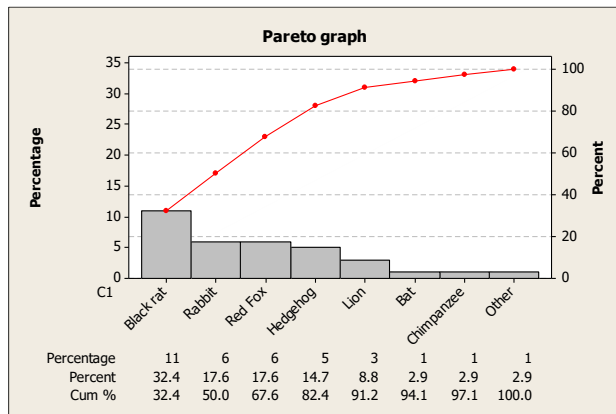
(c) Pie chart

1.4.4

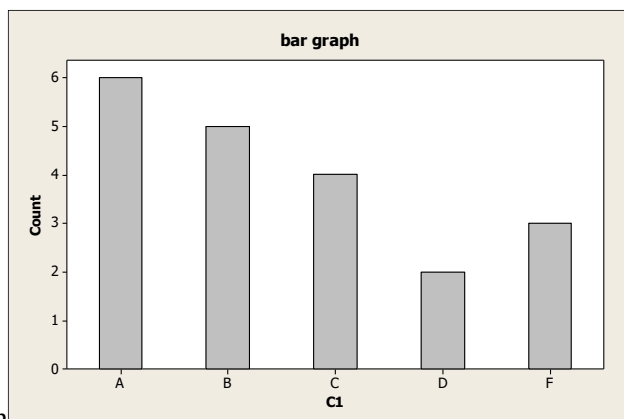
(a) Bar graph



(b) Pareto graph

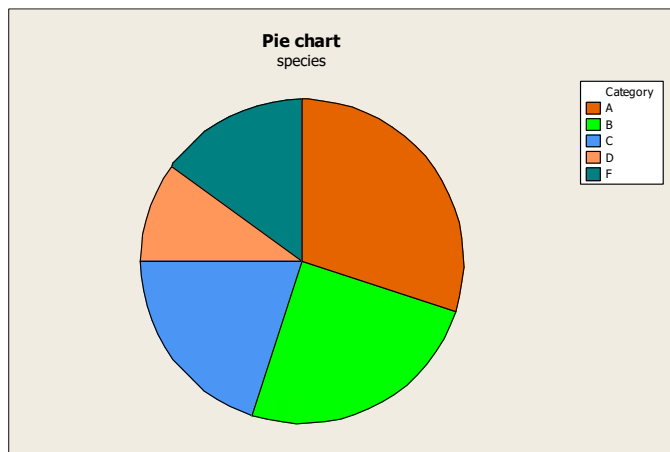


1.4.5



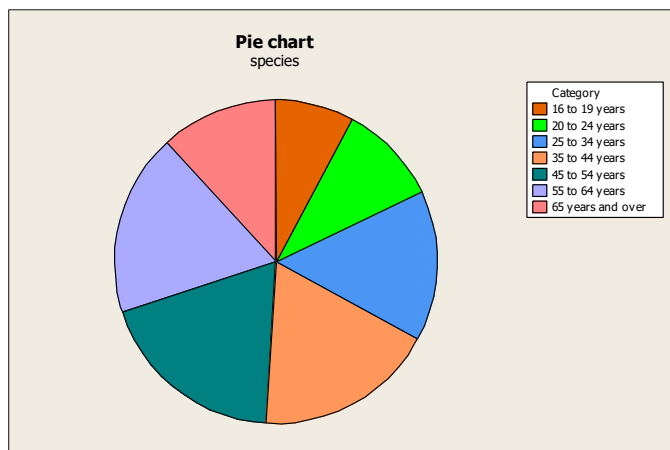
(a) Bar graph

(b) Pie chart

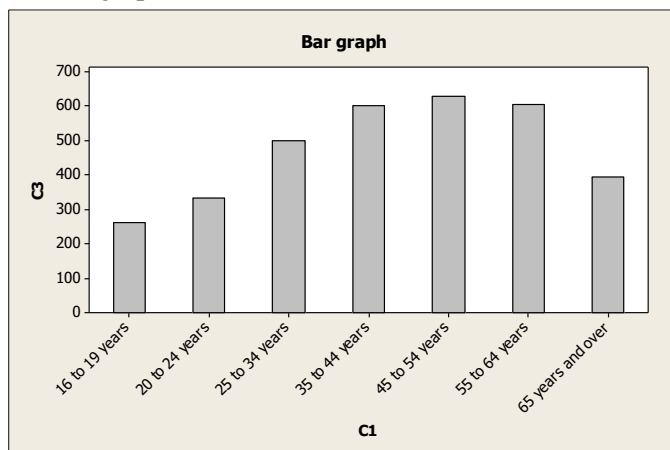


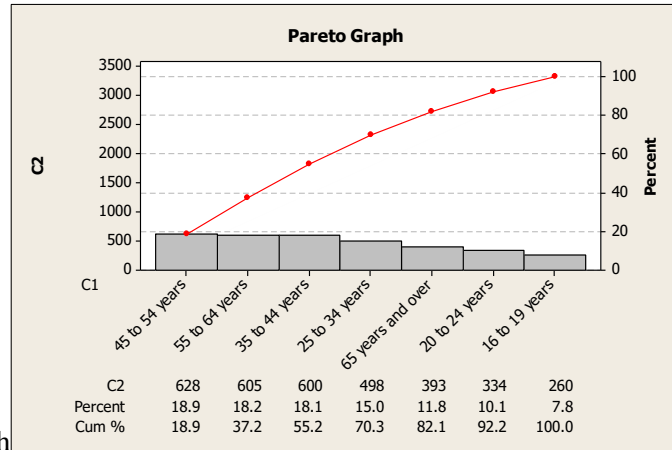
1.4.6

(a) Pie chart



(b) Bar graph

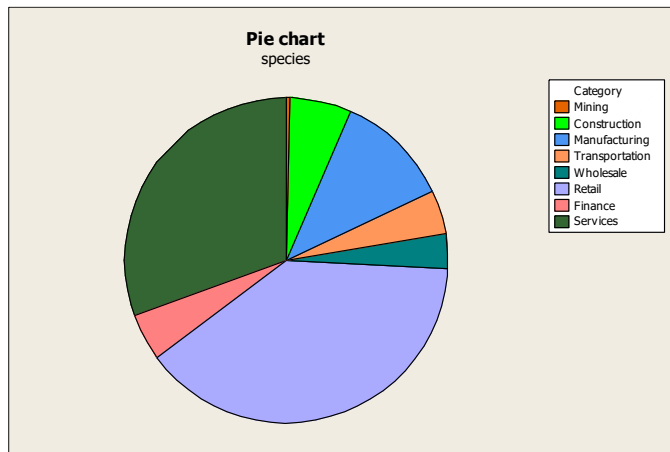




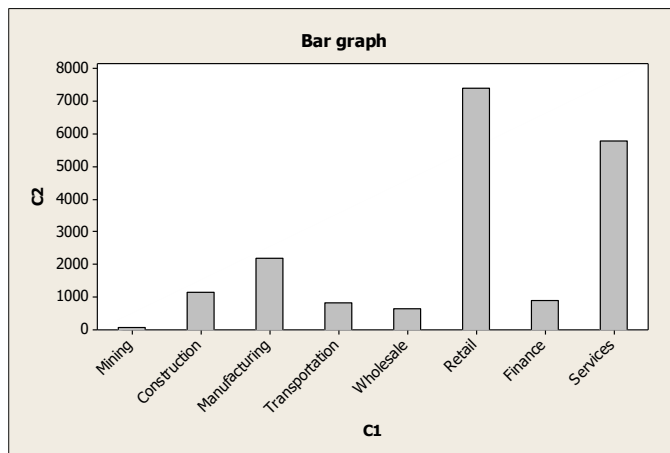
(c) Pareto graph

1.4.7

(a) Pie chart

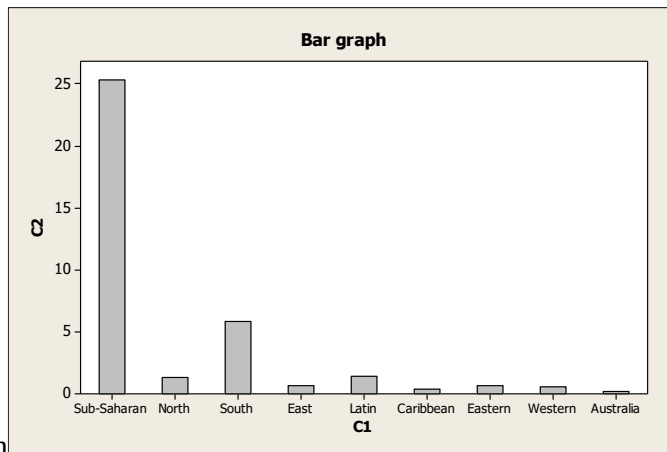


(b) Bar graph



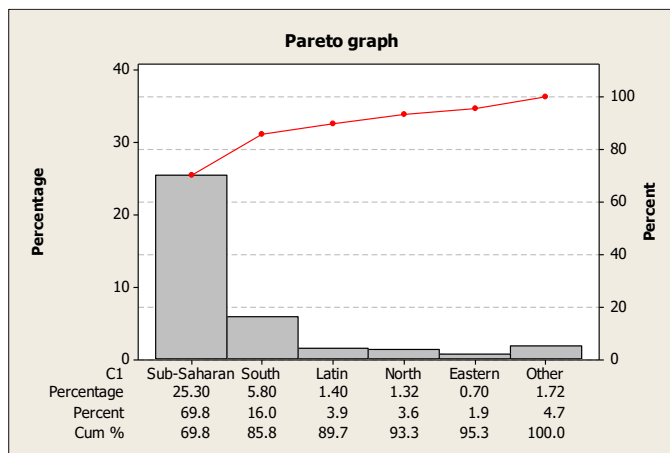


1.4.8



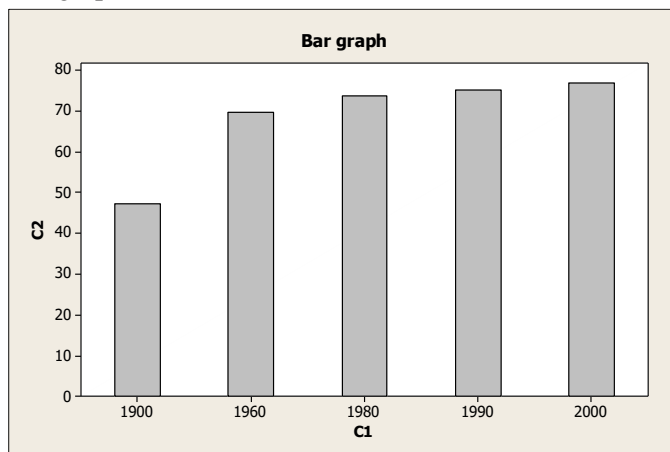
(a) Bar graph

(b) Pareto graph

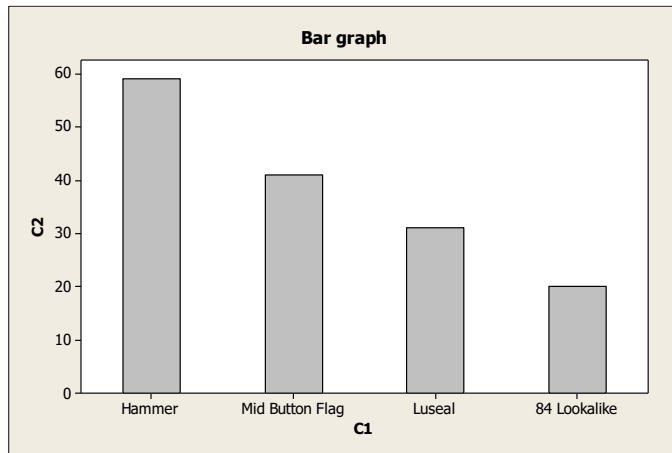


1.4.9

Bar graph

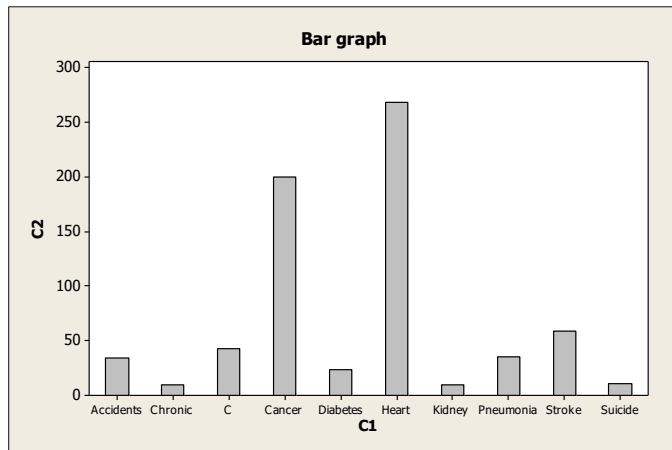


1.4.10

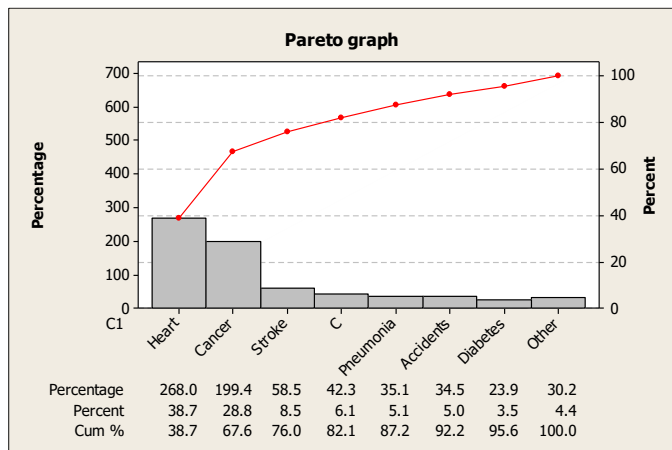


1.4.11

(a) Bar graph



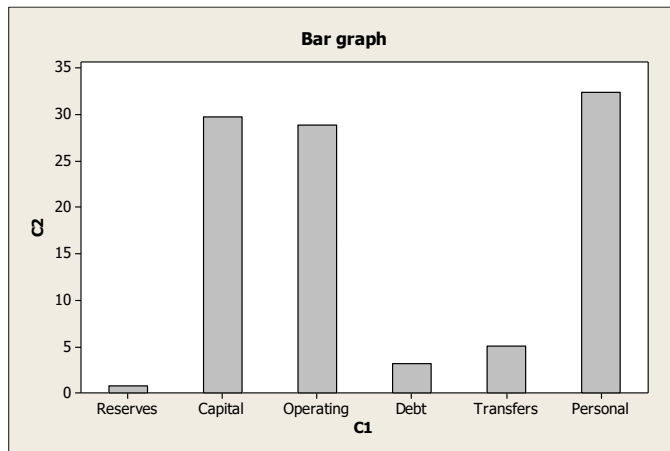
(b) Pareto graph



**1.4.12**

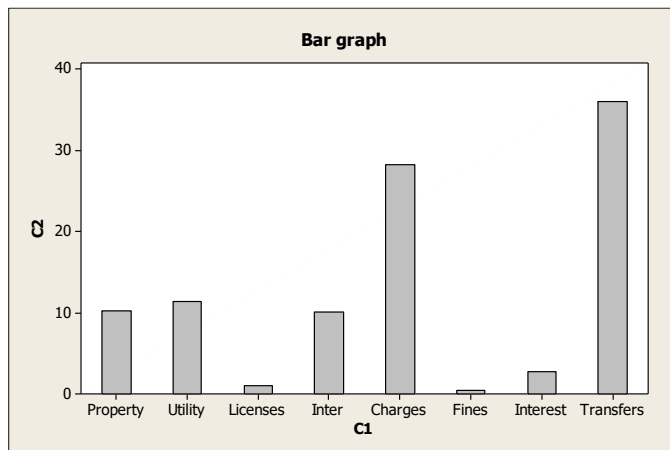
(a) Expenditure

Bar graph



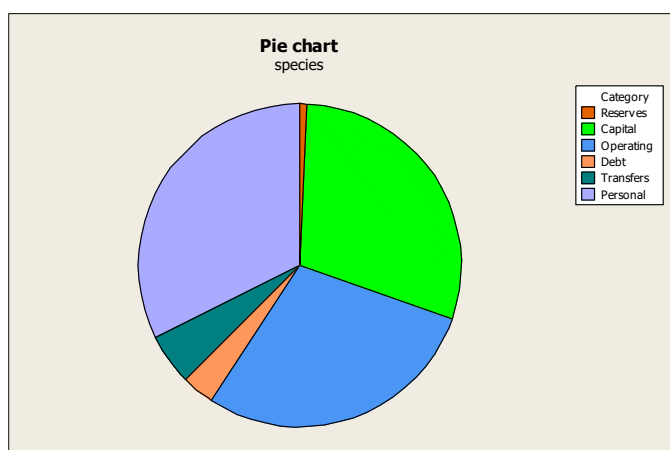
Revenues

Bar graph



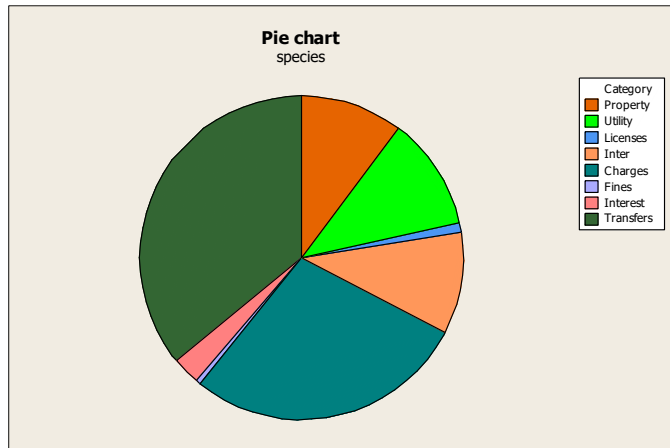
(b) Expenditure

Pie chart

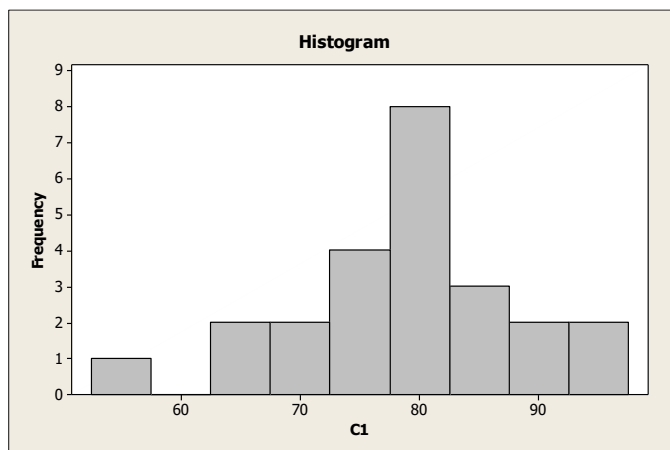


## Revenues

## Pie chart



## 1.4.13



## 1.4.14

(a) Stem and leaf

**Stem-and-Leaf Display: C1**

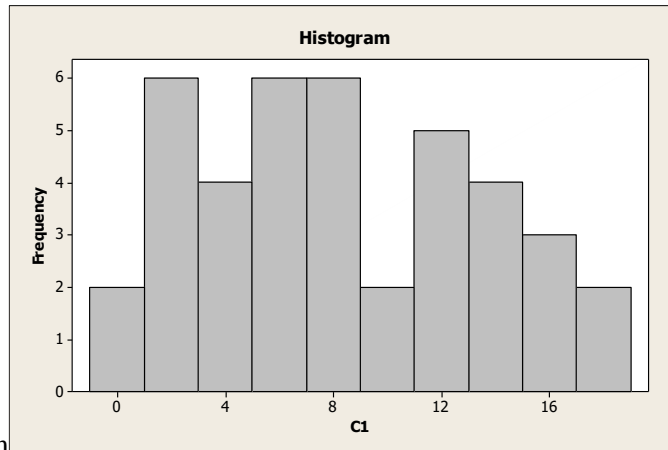
Stem-and-leaf of C1 N = 40

Leaf Unit = 1.0

```

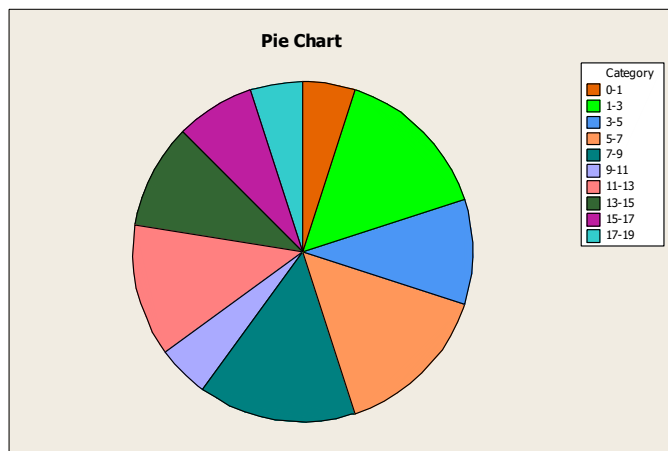
2  0 00
12 0 222223333
13 0 5
20 0 666677
20 0 888899
14 1 111
11 1 223333
5  1 55
3  1 677

```



(b) Histogram

(c) Pie chart

**1.4.15**

(a) Stem and leaf

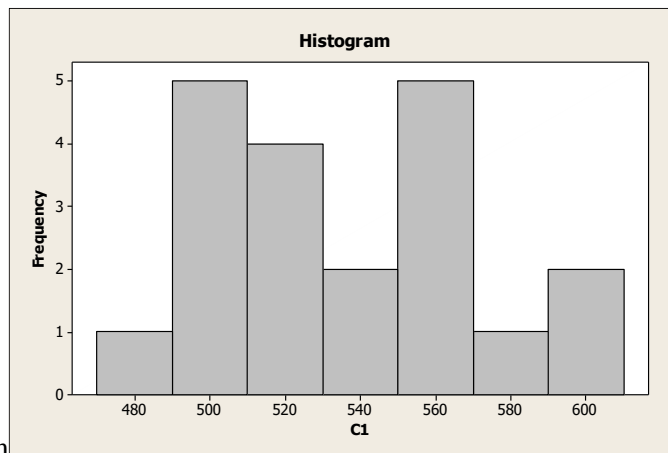
Stem-and-leaf of SAT Mathematics scores  $N = 20$ 

Leaf Unit = 10

```

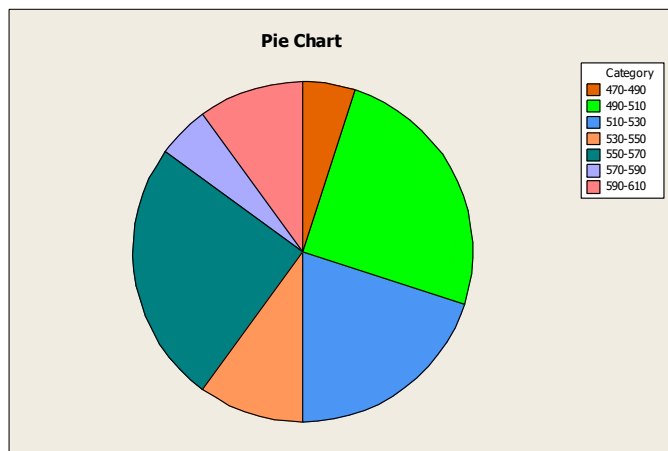
1  4 7
3  4 99
8  5 00011
10 5 22
10 5 4455
6  5 6667
2  5 9
1  6 0

```



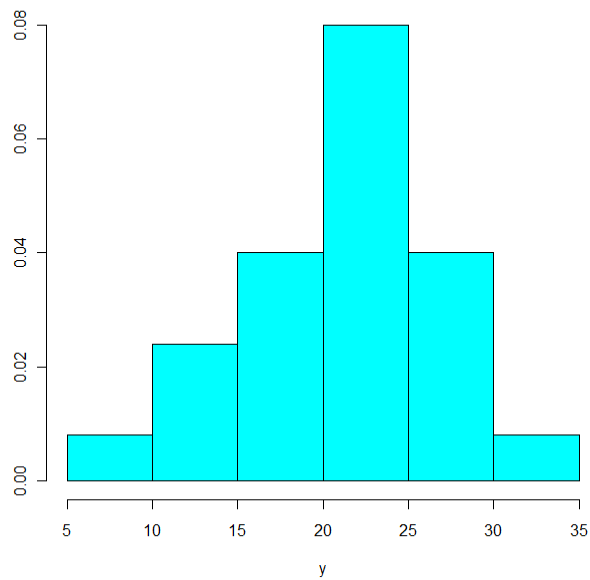
(b) Histogram

(c) Pie chart

**1.4.16**

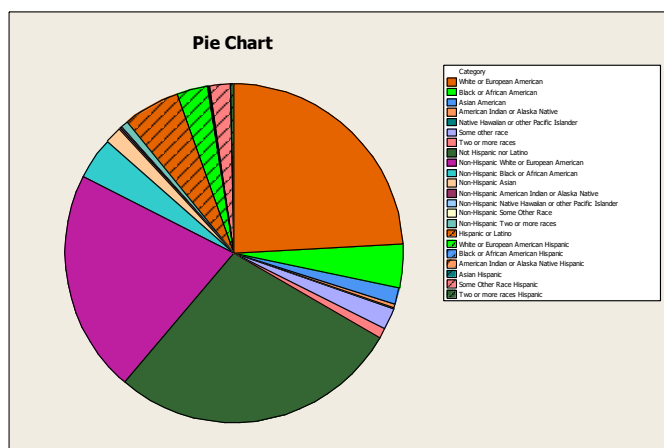
Frequency table

Interval	Frequency	Relative Freq	Percentage
5-9	1	.04	4
10-14	3	.12	12
15-19	5	.2	20
20-24	10	.4	40
25-29	5	.2	20
30-35	1	.04	4



Histogram

## 1.4.17



## Exercises 1.5

## 1.5.1

$$\bar{x} = \frac{1615 + 1780}{2} = 1597.5$$

$$s^2 = \frac{(1615 - 1597.5)^2 + (1780 - 1597.5)^2}{2} = 3882$$

$$s^2 = 3882$$

$$s = \sqrt{3882} = 62.3$$

## 1.5.2

(a)

$$\bar{x} = \frac{15 + 15 + 55 + 15}{4} = 25$$

$$s^2 = \frac{(15-25)^2 + (15-25)^2 + (55-25)^2 + (15-25)^2}{4}$$

$$s^2 = 58$$

$$s = \sqrt{58} \approx 7.6158$$

(b)

$$Q_1 = 6.625$$

$$Q_3 = \frac{7.5 + 7.625}{2} = 7.5625$$

$$M = 7.375$$

$$IQR = 7.5625 - 6.625 = .9375$$

$$LL = 6.625 - 1.5(.9375) = 5.21875$$

$$UL = 7.625 + 1.5(.9375) = 9.03125$$

There are no outliers.

## 1.5.3

Given information: mean=6, median = 4, mode = 3

We know that the value 3 can only be in the data twice. If not the median would be different than 4. This gives us the following: 3, 3, x, y. Where x and y are the missing values. We introduce a system of equations to solve for x and y.

$$\frac{3+3+x+y}{4} = 6 \qquad \frac{x+3}{2} = 4$$

$$x+y = 24-6 \qquad x = 8-3$$

$$x+y = 18 \qquad x = 5$$

$$y = 18-5$$

$$y = 13 \quad x = 5$$

Data: 3, 3, 5, 13

$$Var = \frac{1}{4} [(3-6)^2 + (3-6)^2 + (5-6)^2 + (13-6)^2]$$

$$= \frac{1}{4} (68)$$

$$= 17$$

$$Sd = \sqrt{Var}$$

$$= \sqrt{17}$$

$$= 4.123$$



## 1.5.4

$$\bar{x} = \frac{118 + 150 + 158 + 21}{4} = 125$$

$$s^2 = \frac{(118 - 125)^2 + (150 - 125)^2 + (158 - 125)^2 + (21 - 125)^2}{4-1}$$

$$(a) \quad s^2 = 7269.1$$

$$s = \sqrt{7269.1} = 85.25$$

$$Range = 21 - 21$$

$$Q_1 = 537$$

$$Q_3 = 1578$$

$$M = \frac{1117 + 1050}{2} = 1083.5$$

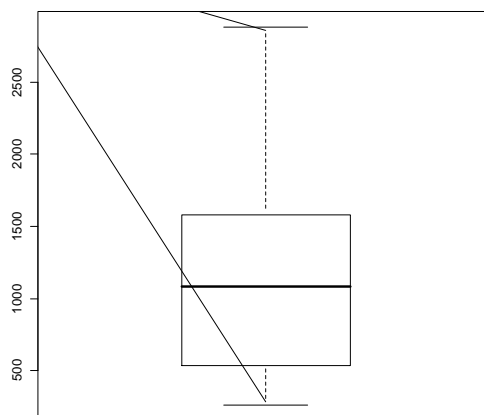
$$(b) \quad IQR = 1578 - 537 = 1041$$

$$LL = 537 - 1.5(1041) = -1024.5$$

$$UL = 1578 + 1.5(1041) = 3139.5$$

There are no outliers.

(c)



## 1.5.5

$$Q_1 = 80$$

$$Q_3 = 115$$

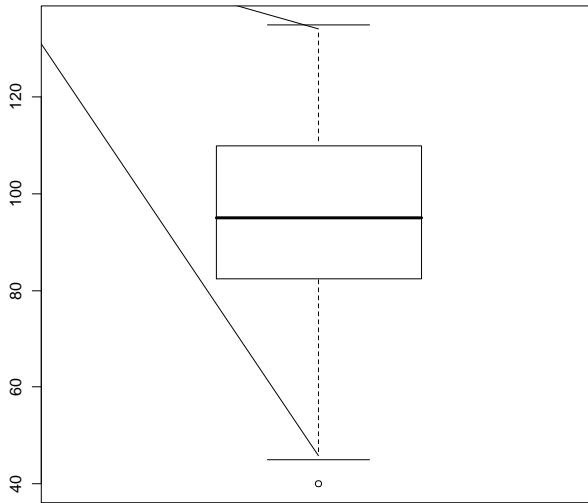
$$(a) \quad M = 95$$

$$IQR = 115 - 80 = 35$$

$$LL = 80 - 1.5(35) = 27.5$$

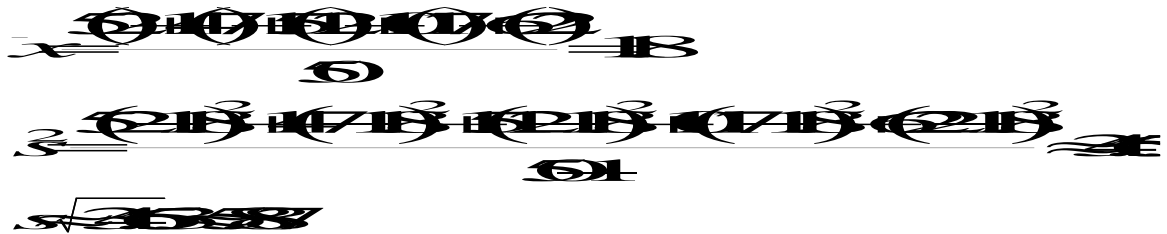
$$UL = 115 + 1.5(35) = 167.5$$

(b)

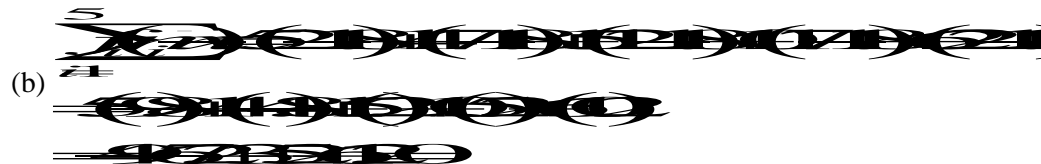


(c) There are no outliers.

1.5.6



1.5.7



## 1.5.8

$$\begin{aligned} \sum_{i=1}^n (x_i - \bar{x})^2 &= \sum_{i=1}^n (x_i^2 - 2x_i\bar{x} + \bar{x}^2) = \sum_{i=1}^n x_i^2 - 2\bar{x} \sum_{i=1}^n x_i + \sum_{i=1}^n \bar{x}^2 \\ &= \sum_{i=1}^n x_i^2 - 2\bar{x}^2 + n\bar{x}^2 = \sum_{i=1}^n x_i^2 - n\bar{x}^2 = \sum_{i=1}^n x_i^2 - n \left( \frac{\sum_{i=1}^n x_i}{n} \right)^2 \end{aligned}$$

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

(b)

~~$$\sum_{i=1}^n (x_i - \bar{x})^2 = \sum_{i=1}^n x_i^2 - 2\bar{x} \sum_{i=1}^n x_i + \sum_{i=1}^n \bar{x}^2$$

$$= \sum_{i=1}^n x_i^2 - 2\bar{x}^2 + n\bar{x}^2 = \sum_{i=1}^n x_i^2 - n\bar{x}^2 = \sum_{i=1}^n x_i^2 - n \left( \frac{\sum_{i=1}^n x_i}{n} \right)^2$$

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

## 1.5.9

~~$$\bar{x} = \frac{\sum_{i=1}^3 x_i}{3} = \frac{9526}{3} = 3175$$

$$s^2 = \frac{1}{3} \sum_{i=1}^3 (x_i - 3175)^2 = \frac{54832}{3} \approx 17944$$

$$s = \sqrt{17944} \approx 134$$~~

$$Q_1 = \frac{24.75 + 25.44}{2} = 25.095$$

$$Q_3 = \frac{42.19 + 43.25}{2} = 42.72$$

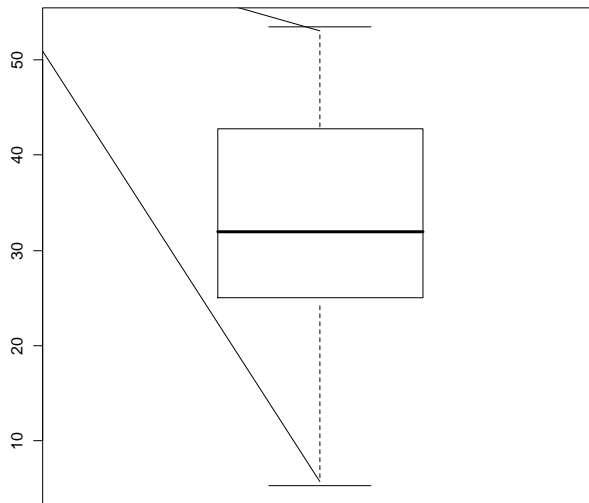
$$(b) M = \frac{32 + 32}{2} = 32$$

$$IQR = 42.72 - 25.095 = 17.625$$

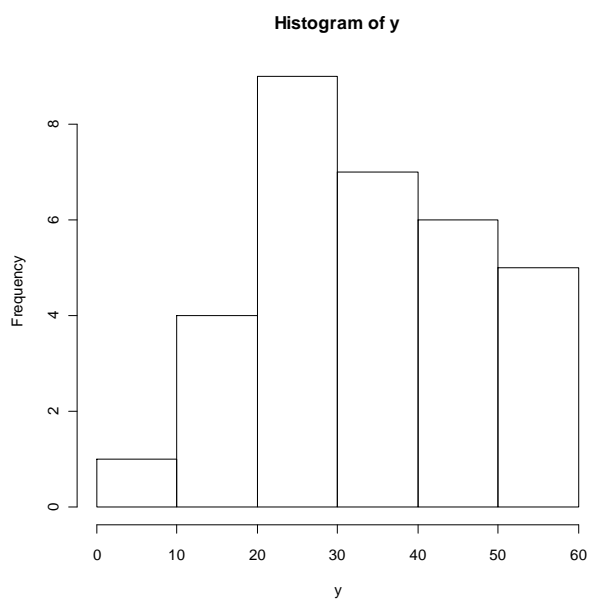
$$LL = 25.095 - 1.5(17.625) = -1.3425$$

$$UL = 42.72 + 1.5(17.625) = 69.1575$$

There are no outliers.



(c)



(d)

(e)

$$\bar{x} = 33.105$$

~~$\bar{x} \pm 1s = [19.954, 46.256]$~~  21 data point (65.625%) fall within 1 SD, empirical rule = 68%

~~$\bar{x} \pm 2s = [6.452, 60.758]$~~  31 data point (96.875%) fall within 2 SD, empirical rule = 95%

~~$\bar{x} \pm 3s = [-6.850, 67.055]$~~  32 data point (100%) fall within 3 SD, empirical rule = 99.7%

## 1.5.10

$$\bar{x} = \frac{\sum_{i=1}^n x_i}{n} = \frac{336}{40} = 8.4$$

$$(a) \quad s^2 = \frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})^2 \approx \frac{94376}{39} \approx 2420$$

$$\text{range} = 12.5 - 3.6 = 8.9$$

$$Q_1 = \frac{3.7 + 3.6}{2} = 3.65$$

$$Q_3 = \frac{12.8 + 12.3}{2} = 12.55$$

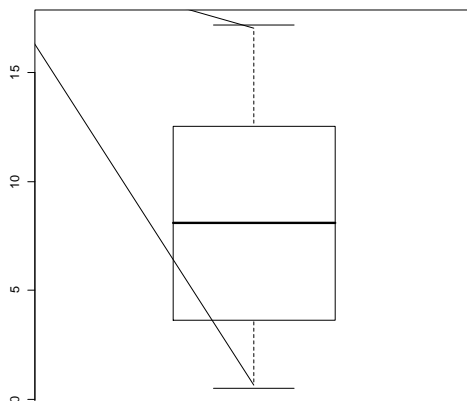
$$(b) \quad M = \frac{8.3 + 7.9}{2} = 8.1$$

$$IQR = 12.55 - 3.65 = 8.9$$

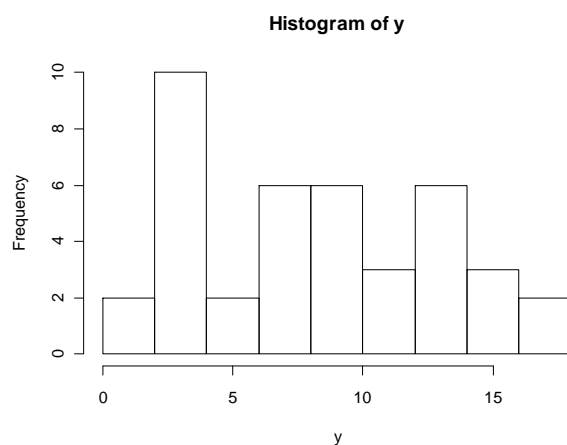
$$LL = 3.65 - 1.5(8.9) = -9.7$$

$$UL = 12.55 + 1.5(8.9) = 25.9$$

There are no outliers.



(c)



(d)

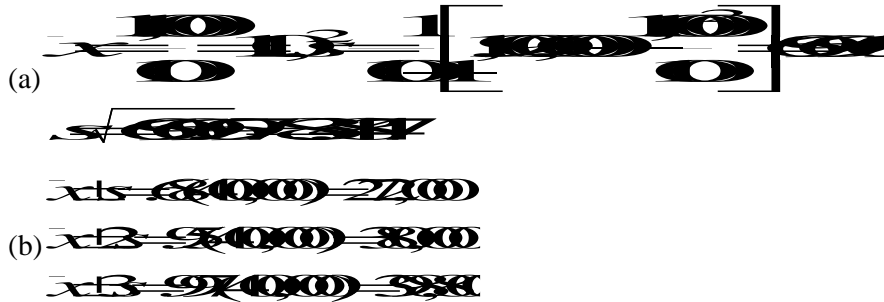
(e)  $\bar{x} = 8.34$

~~$\bar{x} \pm s = [3.2, 13.6]$~~  24 data point (60%) fall within 1 SD, empirical rule = 68%

~~$\bar{x} \pm 2s = [-1.5, 18.8]$~~  40 data point (100%) fall within 2 SD, empirical rule = 95%

~~$\bar{x} \pm 3s = [-6.2, 23.1]$~~  40 data point (100%) fall within 3 SD, empirical rule = 99.7%

### 1.5.11



### 1.5.12

$$Q_1 = 39$$

$$Q_3 = 46$$

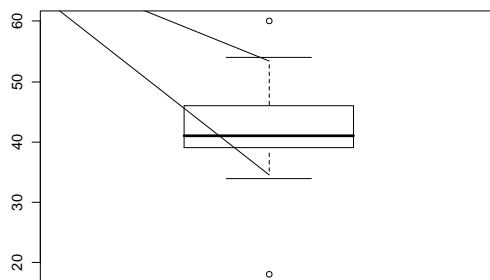
$$M = \frac{42 + 40}{2} = 41$$

$$IQR = 46 - 39 = 7$$

(a) 
$$\bar{x} = \frac{\sum_{i=1}^{10} x_i}{10} = \frac{418}{10} = 41.8$$

$$s^2 = \frac{1}{9} \sum_{i=1}^{10} (x_i - 8.34)^2 = \frac{1149.6}{9} \approx 127.733$$

$$s = \sqrt{127.733} \approx 11.302$$



(c)

$$KR=7$$

$$IL=3+15(7)=285$$

$$(d) \quad IL=46+15(7)=55$$

*Teaetwootias Bark6*

### 1.5.13

(a)

$$\bar{x} = \frac{\sum_{i=1}^n x_i}{n} = \frac{1123}{30} = 37.433$$

$$s^2 = \frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})^2 \approx 3.502$$

$$s = \sqrt{3.502} \approx 1.871$$

(b) Frequency table

Class	Interval	Frequency	$M_i$	$M_i \cdot f_i$
1	0-1.6	4	.8	3.2
2	1.7-3.3	10	2.5	25
3	3.4-5	9	4.2	37.8
4	5.1-6.7	5	5.9	29.5
5	6.8-8.4	2	7.6	15.2

(c) Grouped data:



The results from the grouped data are similar to the actual data.

### 1.5.14

(a)

$$\bar{x} = \frac{\sum_{i=1}^n x_i}{n} = \frac{184}{30} \approx 6.133$$

$$s^2 = \frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})^2 \approx 65.085$$

$$s \approx \sqrt{65.085} \approx 8.074$$

(b) Frequency table

Class	Interval	Frequency	$M_i$	$M_i \cdot f_i$
1	0-20	1	10	10
2	20-40	8	30	240
3	40-60	6	50	300
4	60-80	5	70	350
5	80-100	10	90	900

(c)

Grouped data:

$$\bar{x} = \frac{141610}{6}$$

$$s^2 = \frac{[(109)(39) + (139)(59) + (129)(79) + (109)(99)]}{5}$$

$$s = \sqrt{6523}$$

The results from the grouped data are similar to the actual data.

1.5.15

$$L = 25, f_m = 1396, w = 4$$

$$F_b = 1788, n = 5146$$

$$M_{f_m} = 227$$

1.5.16

(a)

$$\bar{x} = \frac{111111111}{175}$$

$$s^2 = \frac{[(1517)(17) + (1517)(17) + (1517)(17) + (1517)(17) + (1517)(17)]}{174}$$

$$s = \sqrt{13116}$$

(b)  $L = 175, f_m = 18, w = 9$

$$F_b = 19, n = 50$$

$$M_{f_m} = 1$$

1.5.17

$$\bar{x} = \frac{311111111}{40}$$

$$s^2 = \frac{[(19)(41) + (19)(41) + (19)(41) + (19)(41) + (19)(41)]}{39}$$

$$s = \sqrt{5125}$$

(b)  $L = 40, f_m = 59, w = 19, F_b = 69$