

This tutorial explores sigma notation, numerical summaries and the boxplot.

### Numerical Summaries

Given a sample  $\{x_i\}$  and ordered data  $\{x_{(i)}\}$  for  $i = 1, \dots, n$

sample mean	$\bar{x} = \frac{1}{n} \sum_{i=1}^n x_i$
sample variance	$s^2 = \frac{1}{n-1} \sum_{i=1}^n (x_i - \bar{x})^2 = \frac{1}{n-1} \left( \sum_{i=1}^n x_i^2 - n\bar{x}^2 \right)$
sample standard deviation	$s = \sqrt{\frac{1}{n-1} \sum_{i=1}^n (x_i - \bar{x})^2} = \sqrt{\frac{1}{n-1} \left( \sum_{i=1}^n x_i^2 - n\bar{x}^2 \right)}$
1st quartile	$Q_1 = \frac{x_{(\lceil \frac{n}{4} \rceil)} + x_{(\lfloor \frac{n}{4} + 1 \rfloor)}}{2}$
2nd quartile (or median)	$Q_2 = \tilde{x} = \frac{x_{(\lceil \frac{n+1}{2} \rceil)} + x_{(\lfloor \frac{n+1}{2} \rfloor)}}{2}$
3rd quartile	$Q_3 = \frac{x_{(\lceil \frac{3n}{4} \rceil)} + x_{(\lfloor \frac{3n}{4} + 1 \rfloor)}}{2}$
five number summary	$(x_{(1)}, Q_1, Q_2, Q_3, x_{(n)})$
Interquartile Range	$IQR = Q_3 - Q_1$
Boxplot Thresholds (for outliers)	$LT = Q_1 - 1.5IQR, UT = Q_3 + 1.5IQR$

Note: There are 2 formulae for the variance: the 1st one is the definition formula and the 2nd one is the calculation formula. The 2nd formula can also be written as  $s^2 = \frac{1}{n-1} \left( \sum_{i=1}^n x_i^2 - \frac{1}{n} \left( \sum_{i=1}^n x_i \right)^2 \right)$ .

### 1. Sigma Notation, Ceiling and Floor Functions

(a) Given the data  $x = \{1, 2, 3, 6, 7, 9\}$  and  $y = \{1, 1, 2, 3, 4, 4\}$ , by hand calculate

$$\sum_{i=1}^6 x_i \quad \sum_{i=1}^6 x_i^2 \quad \sum_{i=1}^6 x_i y_i \quad \sum_{i=1}^3 (x_i - 5)^2 \quad \sum_{i=2}^3 y_{(i)}^2$$

(b) Calculate the mean and standard deviation of  $x$ .

(c) If each data point in  $x$  is increased by 1, how would the mean and standard deviation change? Why? Check numerically.

(d) Find the 1st and 3rd quartiles of  $x$ , by first calculating

$$x_{(\lceil \frac{n}{4} \rceil)} \quad x_{(\lfloor \frac{n}{4} + 1 \rfloor)} \quad x_{(\lceil \frac{3n}{4} \rceil)} \quad x_{(\lfloor \frac{3n}{4} + 1 \rfloor)}$$

(e) Check your working in (a) to (d) using R.

> x=c(1,2,3,6,7,9)	> mean(x)
> y=c(1,1,2,3,4,4)	> sd(x)
> sum(x)	> quantile(x,type=2)
> sum(x^2)	
> sum(x*y)	> m=x+1
> sum((x-5)^2)[1:3])	> mean(m)
> sum((sort(y)^2)[2:3])	> sd(m)

- (f) Extension (This is not examinable. Just for students who want to challenge themselves.)  
 Given  $m_i = x_i + 1$ , show algebraically that  $\bar{m} = \bar{x} + 1$  and  $s_m^2 = s_x^2$ .

**Solution**

(a)

$$\sum_{i=1}^6 x_i = 1 + 2 + \dots + 9 = 28.$$

$$\sum_{i=1}^6 x_i^2 = 1^2 + 2^2 + \dots + 9^2 = 180.$$

$$\sum_{i=1}^6 x_i y_i = 1 \times 1 + 2 \times 1 + \dots + 9 \times 4 = 91.$$

$$\sum_{i=1}^3 (x_i - 5)^2 = (1 - 5)^2 + (2 - 5)^2 + (3 - 5)^2 = 29.$$

$$\sum_{i=2}^3 y_{(i)}^2 = 1^2 + 2^2 = 5, \text{ as } \{y_{(i)}\} = \{y_i\} \text{ as the data is already sorted.}$$

Note: If say  $\{y_i\} = \{8, 1, 4\}$ , then  $\{y_{(i)}\} = \{1, 4, 8\}$ , so  $\sum_{i=2}^3 y_{(i)}^2 = 4^2 + 8^2 = 80$ .

(b)

$$\bar{x} = \frac{1}{6} \sum_{i=1}^6 x_i = \frac{28}{6} \approx 4.67$$

$$s = \sqrt{\frac{1}{5} \left( \sum_{i=1}^6 x_i^2 - 6\bar{x}^2 \right)} = \sqrt{\frac{1}{5} \left( 180 - 6\left(\frac{28}{6}\right)^2 \right)} \approx 3.14$$

(c)

Mean: Increases by 1 (as the full data set shifts up by 1, the average shifts up by 1).

Standard deviation: Unchanged ((as the full data set shifts up by 1, the spread is unchanged).

Check numerically:

Given  $\{m_i\} = \{x_i + 1\} = \{2, 3, 4, 7, 8, 10\}$ , then  $\bar{m} \approx 5.67$  and  $s_m^2 \approx 3.14$ .

(d)

$$x(\lceil \frac{n}{4} \rceil) = x(\lceil \frac{6}{4} \rceil) = x_{(2)} = 2$$

$$x(\lfloor \frac{n}{4} + 1 \rfloor) = x(\lfloor \frac{6}{4} + 1 \rfloor) = x_{(2)} = 2$$

$$x(\lceil \frac{3n}{4} \rceil) = x(\lceil \frac{3 \times 6}{4} \rceil) = x_{(5)} = 7$$

$$x(\lfloor \frac{3n}{4} + 1 \rfloor) = x(\lfloor \frac{3 \times 6}{4} + 1 \rfloor) = x_{(5)} = 7$$

$$\text{Hence } Q_1 = \frac{x_{(2)} + x_{(2)}}{2} = 2 \text{ and } Q_3 = \frac{x_{(5)} + x_{(5)}}{2} = 7$$

(e)

Commands given.

The standard deviation could also be found (the long way) in R by using the formula:

```
> sqrt( 1/5*(sum(x^2) - 6*mean(x)^2) )
```

(f)

If  $m_i = x_i + 1$  then:

$$\bar{m} = \frac{1}{n} \sum_{i=1}^n (x_i + 1) = \frac{1}{n} \sum_{i=1}^n x_i + \frac{1}{n} \sum_{i=1}^n 1 = \bar{x} + 1.$$

$$s_m^2 = \frac{1}{n-1} \sum_{i=1}^n (m_i - \bar{m})^2 = \frac{1}{n-1} \sum_{i=1}^n ((x_i + 1) - (\bar{x} + 1))^2 = \frac{1}{n-1} \sum_{i=1}^n (x_i - \bar{x})^2 = s_x^2.$$

## 2. Numerical Summaries

A sample of 36 mice was used to investigate the use of iron in  $\text{Fe}^+$  form as a dietary supplement. The iron was given orally and was radioactively labelled so that the exact percentage of iron retained could be measured accurately. The measurements were

7.6	1.2	4.9	5.7	13.0	1.0	3.4	0.2	10.8	1.0	2.4	12.3
0.7	1.1	0.7	0.9	6.5	1.6	4.0	29.1	0.2	0.1	9.2	11.9
0.3	14.4	1.8	9.9	3.4	3.8	9.9	4.1	4.1	24.0	21.0	11.9

(a) Using the the following R output, fill out the table.

Size of data	Mean	Median	Standard deviation	Variance	1st Quartile	3rd Quartile	IQR

```
> sum(x)
[1] 238.1
> sum(x^2)
[1] 3333.85
> sort(x)
 [1]  0.1  0.2  0.2  0.3  0.7  0.7  0.9  1.0  1.0  1.1  1.2  1.6  1.8  2.4
[15]  3.4  3.4  3.8  4.0  4.1  4.1  4.9  5.7  6.5  7.6  9.2  9.9  9.9 10.8
[29] 11.9 11.9 12.3 13.0 14.4 21.0 24.0 29.1
```

(b) What is the five number summary of  $x$ ?

(c) Construct a boxplot by hand.

(d) In order to compare the sensitivities to outliers of the mean, median, standard deviation and IQR, remove the largest value and recompute these four numerical summaries. Compute the relative change in each, as a percentage.

	Mean	Median	Standard deviation	IQR
Data without largest value				
Relative Change (%)				

Hint: the sum and sum of squares become 209 and 2487.04.

(e) Comment on your findings.

(f) Check your answers in (a) to (c) using R.

```
> x=scan(file=url("http://www.maths.usyd.edu.au/MATH1005/r/wk3q1.txt"))
> length(x)
> mean(x)
> median(x)
```

```

> sd(x)
> var(x)
> quantile(x,type=2)
> iqr=quantile(x,type=2)[4]-quantile(x,type=2)[2]
> boxplot(x)

```

### Solution

(a)

Size of data	Mean	Median	Standard deviation	Variance	1st Quartile	3rd Quartile	IQR
36	6.61	4.05	7.09	50.26	1.05	10.35	9.3

Calculations:

$$\bar{x} = \frac{238.1}{36} = 6.61$$

$$\tilde{x} = \frac{x_{(18)} + x_{(19)}}{2} = \frac{4.0 + 4.1}{2} = 4.05$$

$$s = \sqrt{\frac{1}{35} \left( 3333.85 - 36 \left( \frac{238.1}{36} \right)^2 \right)} \approx 7.09$$

$$s^2 \approx 50.26$$

$$Q_1 = \frac{x_{(\lceil \frac{36}{4} \rceil)} + x_{(\lfloor \frac{36}{4} + 1 \rfloor)}}{2} = \frac{x_{(9)} + x_{(10)}}{2} = \frac{1 + 1.1}{2} = 1.05$$

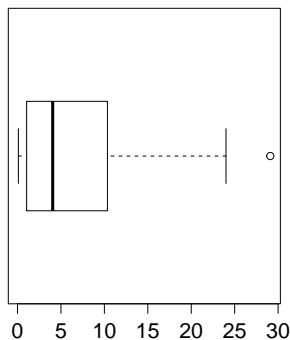
$$Q_3 = \frac{x_{(\lceil \frac{3 \times 36}{4} \rceil)} + x_{(\lfloor \frac{3 \times 36}{4} + 1 \rfloor)}}{2} = \frac{x_{(27)} + x_{(28)}}{2} = \frac{9.9 + 10.8}{2} = 10.35$$

$$IQR = 10.35 - 1.05 = 9.3.$$

(b) 5 number summary is: (0.1, 1.05, 4.05, 10.35, 29.1)

(c) Construction of Boxplot:

1. Draw a box from  $Q_1 = 1.05$  to  $Q_3 = 10.35$ , with a line within the box for  $Q_2 = 4.05$ .
2. Calculate  $LT = Q_1 - 1.5 \times IQR = -12.9$  and  $UT = Q_3 + 1.5 \times IQR = 24.3$ , which determines the maximum possible length of the whiskers.
3. Draw a whisker from the box to the nearest points within LT and UT: this is 0.1 and 24.0
4. Any points lying outside the thresholds are outliers (designated by circles): 29.1.



(d) If the new data set (without the outlier) is  $\{y\}$ , then we have  $\sum_{i=1}^{35} y_i = 209$  and  $\sum_{i=1}^{35} y_i^2 = 2487.04$ .

	Mean	Median	Standard deviation	IQR
Data without largest value	5.97	4	6.04	8.9
Relative Change (%)	$(6.61-5.97)/6.61$ = 9.6	$(4.05-4)/4.05$ = 1.2	$(7.09-6.04)/7.09$ = 14.8	$(9.3-8.9)/9.3$ = 4.3

Calculations:

$$\bar{y} = \frac{209}{35} = 5.97$$

$$\tilde{y} = x_{(18)} = 4$$

$$s_y = \sqrt{\frac{1}{34} \left( 2487.04 - 35 \left( \frac{209}{35} \right)^2 \right)} \approx 6.04$$

$$s_y^2 \approx 36.44$$

$$Q_1 = \frac{x_{(\lceil \frac{35}{4} \rceil)} + x_{(\lfloor \frac{35}{4} + 1 \rfloor)}}{2} = \frac{x_{(9)} + x_{(9)}}{2} = 1$$

$$Q_3 = \frac{x_{(\lceil \frac{3 \times 35}{4} \rceil)} + x_{(\lfloor \frac{3 \times 35}{4} + 1 \rfloor)}}{2} = \frac{x_{(27)} + x_{(27)}}{2} = 9.9$$

$$IQR = 9.9 - 1 = 8.9.$$

(e) Comment: Notice that the median and IQR have small relative change compared to the mean and sd, as they are robust.

(f) Commands given.

Note: R calculates quantiles using a few different commands. For our definition of quartiles, use the `quantile` command. Don't use the `IQR` command or the `summary` command.

```
> quantile(y,type=2)
 0%  25%  50%  75% 100%
0.1  1.0  4.0  9.9 24.0
```

### 3. Comparison of Boxplots

Students completed an online quiz consisting of 20 questions, resulting in the following marks.

Students who had Studied (A): 9 10 11 12 12 13 14 15 15 16 17 17 18

Students who had not studied (B): 1 3 5 8 9 9 10 10 12 12 14 15 16

(a) By hand, produce boxplots for A and B.

(b) Scan the data into R and produce the boxplots.

```
>boxplot(a,b)
>boxplot(a,b, horizontal=TRUE,col=c("green","blue"))    # More colourful version!
```

(c) Comment on your findings.

### **Solution**

(a)

For A:

$$\tilde{x} = x_{(7)} = 12$$

$$Q_1 = \frac{x_{(\lceil \frac{13}{4} \rceil)} + x_{(\lfloor \frac{13}{4} + 1 \rfloor)}}{2} = \frac{x_{(4)} + x_{(4)}}{2} = 12$$

$$Q_3 = \frac{x_{(\lceil \frac{3 \times 13}{4} \rceil)} + x_{(\lfloor \frac{3 \times 13}{4} + 1 \rfloor)}}{2} = \frac{x_{(10)} + x_{(10)}}{2} = 16$$

$$IQR = 16 - 12 = 4$$

$$LT = Q_1 - 1.5 * IQR = 12 - 6 = 6$$

$$UT = Q_3 + 1.5 * IQR = 16 + 6 = 22$$

Comparing to sorted data, we see there are no data points outside LT and UT, hence there are no outliers.

For B:

$$\tilde{x} = x_{(7)} = 10$$

$$Q_1 = \frac{x_{(\lceil \frac{13}{4} \rceil)} + x_{(\lfloor \frac{13}{4} + 1 \rfloor)}}{2} = \frac{x_{(4)} + x_{(4)}}{2} = 8$$

$$Q_3 = \frac{x_{(\lceil \frac{3 \times 13}{4} \rceil)} + x_{(\lfloor \frac{3 \times 13}{4} + 1 \rfloor)}}{2} = \frac{x_{(10)} + x_{(10)}}{2} = 12$$

$$IQR = 12 - 8 = 4$$

$$LT = Q_1 - 1.5 * IQR = 8 - 6 = 2$$

$$UT = Q_3 + 1.5 * IQR = 12 + 6 = 18$$

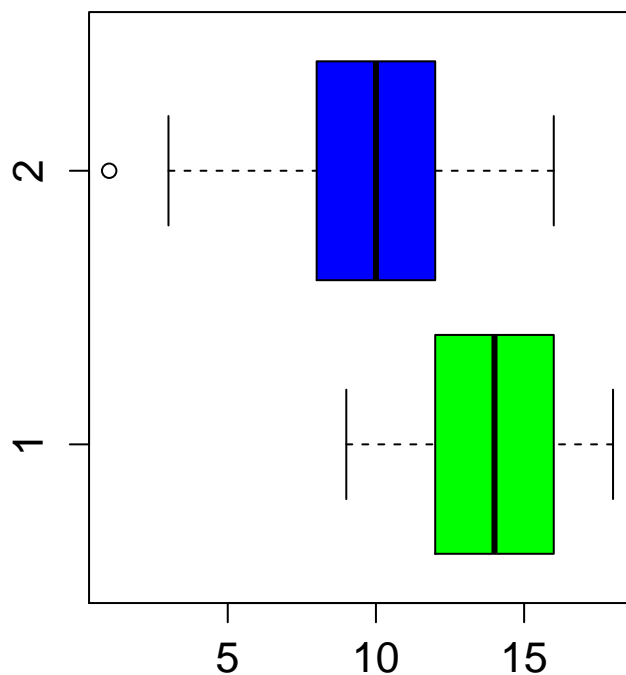
Comparing to sorted data, we see there is one data point lower than LT, hence there is one outlier at 1.

Check your 2 boxplots against output in (b).

(b) Check quantiles in R:

```
> quantile(a,type=2)
 0%  25%  50%  75% 100%
  9   12   14   16   18
> quantile(b,type=2)
 0%  25%  50%  75% 100%
  1    8   10   12   16
```

```
>boxplot(a,b, horizontal=TRUE,col=c("green","blue"))
```



(c) Comments: The students who hadn't studied (B) have a lower median and a much bigger spread of marks than the students who had studied (A). One student in B has a mark (1) which is much lower than the whole rest of cohort.

#### 4. Mean and median

- The sample average age of 5 people in a room is 30 years. A 36 year old person walks into the room. Now what is the average age of the people in the room?
- Suppose the median age is 30 years and a 36 year old person enters the room. Can you find the new median age from this information?

**Solution**

Let  $x_i$  denote the  $i^{th}$  person in the room. Initially there are 5 people in the room.

$$\begin{aligned}\bar{x} &= \frac{1}{5} \sum_{i=1}^5 x_i \\ &= 30\end{aligned}$$

Hence the sum of the ages of the 5 people is  $\sum_{i=1}^5 x_i = 150$ . When person 6 enters the room the sum of the ages is  $\sum_{i=1}^6 x_i = 150 + 36$ . Thus the new mean age will be

$$\begin{aligned}\bar{x} &= \frac{1}{6} \sum_{i=1}^6 x_i \\ &= 31\end{aligned}$$

Order the individuals by age, initially we can deduce that  $x_{(3)} = 30$  by the median. When person 6 enters the room the median will now be  $Q_2 = \frac{x_{(3)} + x_{(4)}}{2}$ , as there is no way to determine  $x_{(4)}$  we can not deduce the new median age.

**5. Extension (Sigma Notation)**

Show that the 2 formulae for variance are equal.

**Solution**

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