



Development Team

Module 2

Analysing Data Numerically

Measures of Central Tendency

- Mean
- Median
- Mode

Measures of Spread

- Range
- Standard Deviation
- Inter-Quartile Range

Mean & Standard Deviation using a Calculator

Calculate the mean and standard deviation of the following 10 students heights by:

- using the data unsorted
- creating a frequency table

Unsorted Data

	Gender	Height /cm
1.	Boy	165
2.	Girl	165
3.	Boy	150
4.	Boy	171
5.	Girl	153
6.	Boy	171
7.	Girl	153
8.	Girl	153
9.	Boy	166
10.	Boy	179

Frequency Table

Height /cm	150	153	165	166	171	179
Frequency	1	3	2	1	2	1

Mean = 162.6 cm

S.D. = 9.32cm

Central Tendency: The Mean

Advantages:

- Mathematical centre of a distribution
- Does not ignore any information

Disadvantages:

- Influenced by extreme scores and skewed distributions
- May not exist in the data

Central Tendency: The Mode

Advantages:

- Good with nominal data
- Easy to work out and understand
- The score exists in the data set

Disadvantages:

- Small samples may not have a mode
- More than one mode might exist

Central Tendency: The Median

Advantages:

- Not influenced by extreme scores or skewed distribution
- Good with ordinal data
- Easier to calculate than the mean
- Considered as the typical observation

Disadvantages:

- May not exist in the data
- Does not take actual data into account only its (ordered) position
- Difficult to handle theoretically

Summary: Relationship between the 3 M's

Characteristics	Mean	Median	Mode
Consider all the data in calculation	Yes	No	No
Easily affected by extreme data	Yes	No	No
Can be obtained from a graph	No	Yes	Yes
Should be one of the data	No	No	Yes
Need to arrange the data in ascending order	No	Yes	No

Using Appropriate Averages

Example

There are 10 houses in Pennylane Close.

On Monday, the numbers of letters delivered to the houses are:

0 2 5 3 34 4 0 1 0 2

Calculate the mean, mode and the median of the number of letters.

Comment on your results.

Solution

$$\begin{aligned}\text{Mean} &= \frac{0+2+5+3+34+4+0+1+0+2}{10} \\ &= 5.1\end{aligned}$$

$$\text{Mode} = 0$$

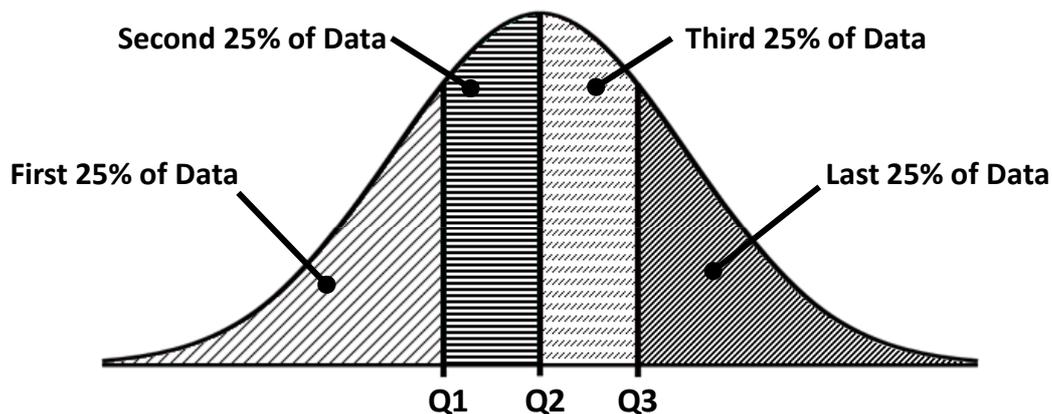
$$\text{Median} = 2$$

In this case the **mean** (5.1) has been distorted by the large number of letters delivered to one of the houses. It is, therefore, not a good measure of a 'typical' number of letters delivered to any house in Pennylane Close.

The **mode** (0) is also not a good measure of a 'typical' number of letters delivered to a house, since 7 out of the 10 houses do actually receive some letters.

The **median** (2) is perhaps the best measure of the 'typical' number of letters delivered to each house, since half of the houses received 2 or more letters and the other half received 2 or fewer letters.

Quartiles



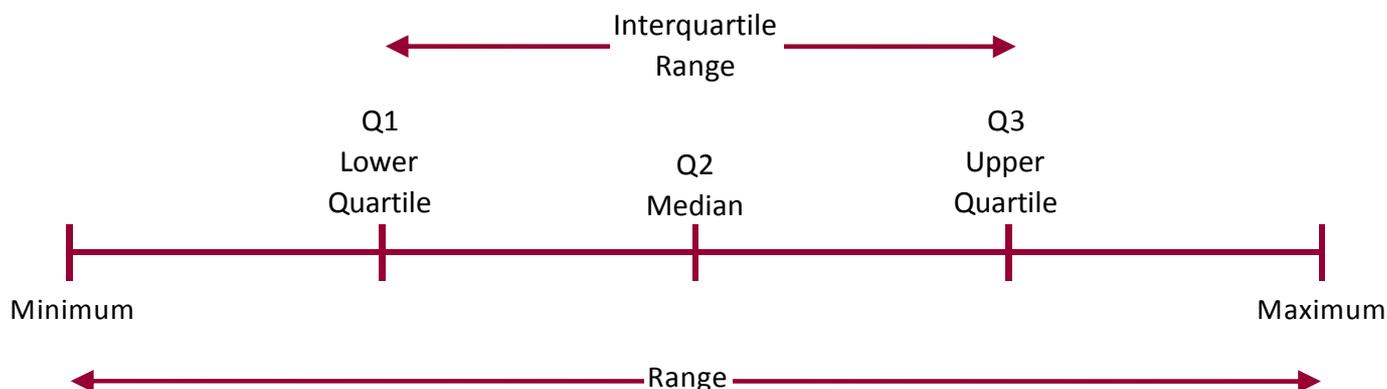
Quartiles

When we arrange the data in ascending order of magnitude and divide them into four equal parts, the values which divide the data into four equal parts are called **quartiles**.

They are usually denoted by the symbols Q_1 is the lowest quartile (or first quartile) where 25% of the data lie below it; Q_2 is the middle quartile (or second quartile or median) where 50% of the data lie below it; Q_3 is the upper quartile (or third quartile) where 75% of the data lie below it.

Interquartile Range & Range

Data in ascending order of magnitude			
First 25% of data	Second 25% of data	Third 25% of data	Last 25% of data



Standard Deviation

Example

Two machines A and B are used to measure the diameter of a washer. 50 measurements of a washer are taken by each machine. If the standard deviations of measurements taken by machine A and B are 0.4mm and 0.15mm respectively, which instrument gives more consistent measurements?

Solution

Standard deviation of A = 0.4 mm

Standard deviation of B = 0.15 mm

The smaller the standard deviation, the less widely dispersed the data is. This means that more measurements are closer to the mean. Therefore, the measurements taken by instrument B are more consistent.

Guide to Distributions

Example 1

Seven teenagers at a youth club were asked their age.
They gave the following ages:

16, 14, 19, 16, 13, 18, 16

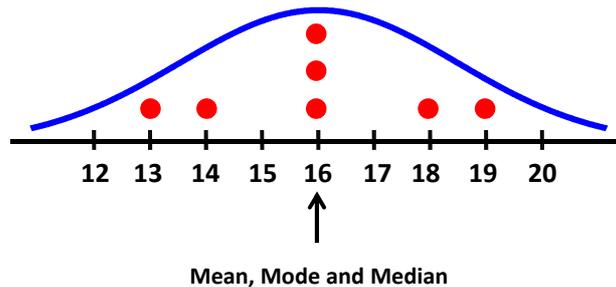
The mean, mode and the median of their ages are as follows:

Mean = 16

Mode = 16

Median = 16

If a line plot of their ages is drawn we get the following.



When the **Mean** and **Median** are the same value the plot is symmetrical (Bell Shaped).
The **Mode** affects the height of the of the Bell Shaped curve.

Example 2:

Seven people were asked how many text messages they send (on average) every week. The results were as follows:

3, 25, 30, 30, 30, 33, 40

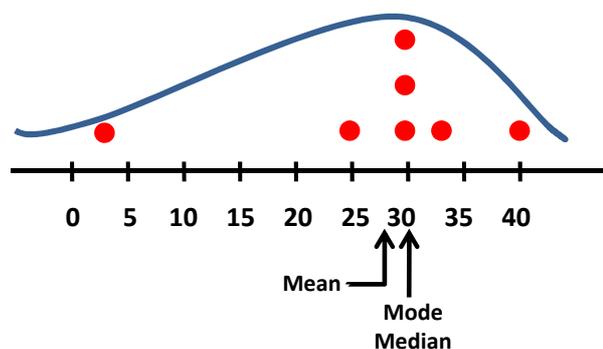
The mean mode and the median of text messages sent are as follows:

Mean = 27.29

Mode = 30

Median = 30

If a line plot of the number of texts sent is drawn we get the following:



When the **Mean** is to the left of the **Median** the data is said to be skewed to the **left** or **negatively skewed**.

The **Mode** affects the height of the curve.

Example 3:

Eight factory workers were asked to give their annual salary.

The results are as follows: (figures in thousands of Euro)

20, 22, 25, 26, 27, 27, 70

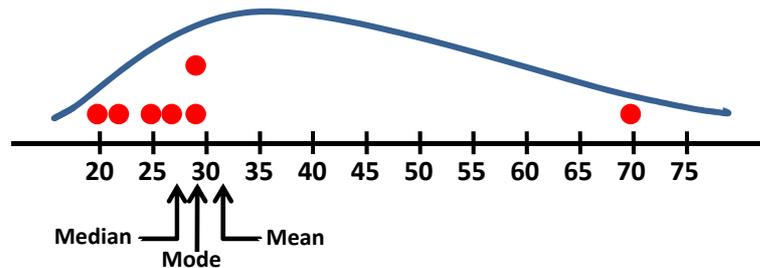
The mean mode and the median of their annual salaries are as follows:

Mean = 31

Mode = 27

Median = 26

If a line plot of their salaries is drawn we get the following.



When the **Mean** is to the **right** of the **Median** the data is said to be skewed to the **right** or **positively skewed**.

The **Mode** affects the height of the curve.

Bivariate Data

1. Involves **2 variables**
2. Deals with causes or relationships
3. The major purpose of bivariate analysis is to determine whether relationships exist

We will look at the following:

- Scatter plots
- Correlation
- Correlation coefficient
- Correlation & causality
- Line of Best Fit
- Correlation coefficient not equal to slope

Sample question:

Is there a relationship between the scores of students who study Physics and their scores in Mathematics?

Univariate Data versus Bivariate Data

Univariate data: Only one item of data is collected e.g. height

Bivariate Data: Data collected in pairs to see if there is a relationship between the variables e.g. height and arm span, mobile phone bill and age etc.

Examples:

Categorical paired data: Colour of eyes and gender

Discrete paired: Number of bars eaten per week and number of tooth fillings

Continuous paired: Height and weight

Category and discrete paired: Type of dwelling and number of occupants etc.
[Look at C@S questionnaire]

Correlation

- Correlation: is about assessing the strength of the relationship between pairs of data. The first step in determining the relationship between 2 variables is to draw a Scatter Plot.
- After establishing if a Linear Relationship (Line of Best Fit) exists between 2 variables X and Y, the strength of the relationship can be measured and is known as the correlation coefficient (r).
- *Correlation* is a precise term describing the strength and direction of the **linear** relationship between quantitative variables.

Correlation Coefficient (r)

$$-1 \leq r \leq 1$$

$r = +1$

Corresponds to a perfect positive linear correlation where the points lie exactly on a straight line
[The line will have positive slope]

$r = 0$:

Corresponds to little or no correlation i.e. as x increases there is no definite tendency for the values of y to increase or decrease in a straight line

$r = -1$:

Corresponds to a perfect negative linear correlation where the points lie exactly on a straight line
[The line will have negative slope]

r close to $+1$:

Indicates a strong positive linear correlation, i.e. y tends to increase as x increases

r close to -1 :

Indicates a strong negative linear correlation, i.e. y tends to decrease as x increases

The correlation coefficient (r) is a numerical measure of the direction and strength of a linear association.

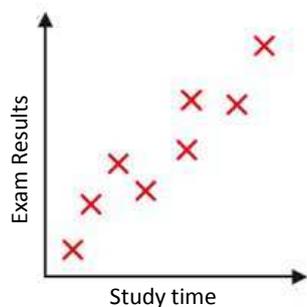
Scatter Plots

- Can show the relationship between 2 variables using ordered pairs plotted on a coordinate plane
- The data points are not joined
- The resulting pattern shows the type and strength of the relationship between the two variables
- Where a relationship exists, a line of best fit can be drawn (by eye) between the points
- Scatter plots can show positive or negative correlation, weak or strong correlation, outliers and spread of data

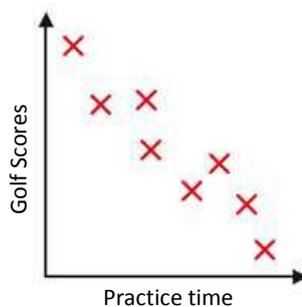
An outlier is a data point that does not fit the pattern of the rest of the data. There can be several reasons for an outlier including mistakes made in the data entry or simply an unusual value.

Describing Correlation

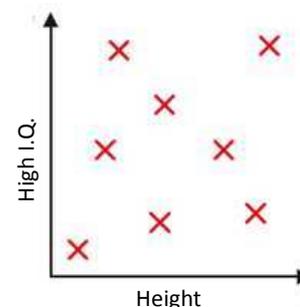
- **Form:** Straight, Curved, No pattern
- **Direction:** Positive, Negative, Neither
- **Strength:** Weak, Moderate, Strong
- **Unusual Features:** Outliers, Subgroups



Positive correlation
As one quantity increases so does the other



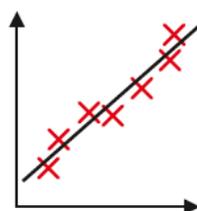
Negative correlation
As one quantity increases the other decreases.



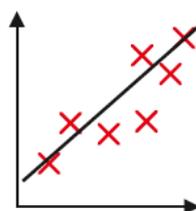
No correlation
Both quantities vary with no clear relationship

Line of Best Fit

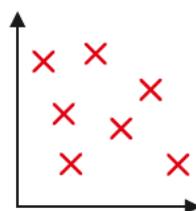
- Roughly goes through the middle of the scatter of the points
- To describe it generally: it has about as many points on one side of the line as the other, and it doesn't have to go through any of the points
- It can go through some, all or none of the points
- Strong correlation is when the scatter points lie very close to the line
- It also depends on the size of the sample from which the data was chosen



Strong positive correlation



Moderate positive correlation



No correlation – no linear relationship



Moderate negative correlation



Strong negative correlation

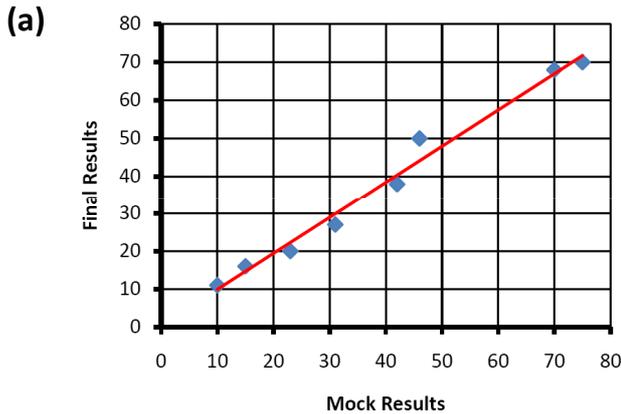
Example

A set of students sat their mock exam in English, they sat their final exam in English at a later date. The marks obtained by the students in both examinations were as follows:

Students	A	B	C	D	E	F	G	H
Mock Results	10	15	23	31	42	46	70	75
Final Results	11	16	20	27	38	50	68	70

- (a) Draw a Scatter Plot for this data and draw a line of best fit.
(b) Is there a correlation between mock results and final results?

Solution



- (b) There is a positive correlation between mock results and final results.

Correlation Coefficient by Calculator

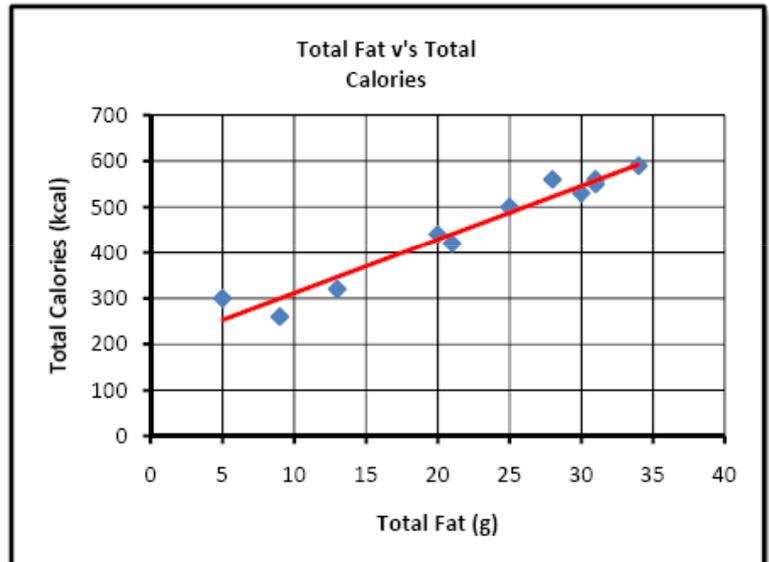
Food	Total Fat (g)	Total Calories (kcal)
Hamburger	9	260
Cheeseburger	13	320
Quarter Pounder	21	420
Quarter Pounder with	30	530
Big Mac	31	560
Special	31	550
Special with Bacon	34	590
Crispy Chicken	25	500
Fish Fillet	28	560
Grilled Chicken	20	440
Grilled Chicken Light	5	300

$$r = 0.9746$$

$$a = 193.85$$

$$b = 11.73$$

$$y = 193.85 + 11.73x$$



Before doing this on the calculator, the class should do a scatter plot using the data in the table. Discuss the relationship between the data (i.e. grams of fat v calories).

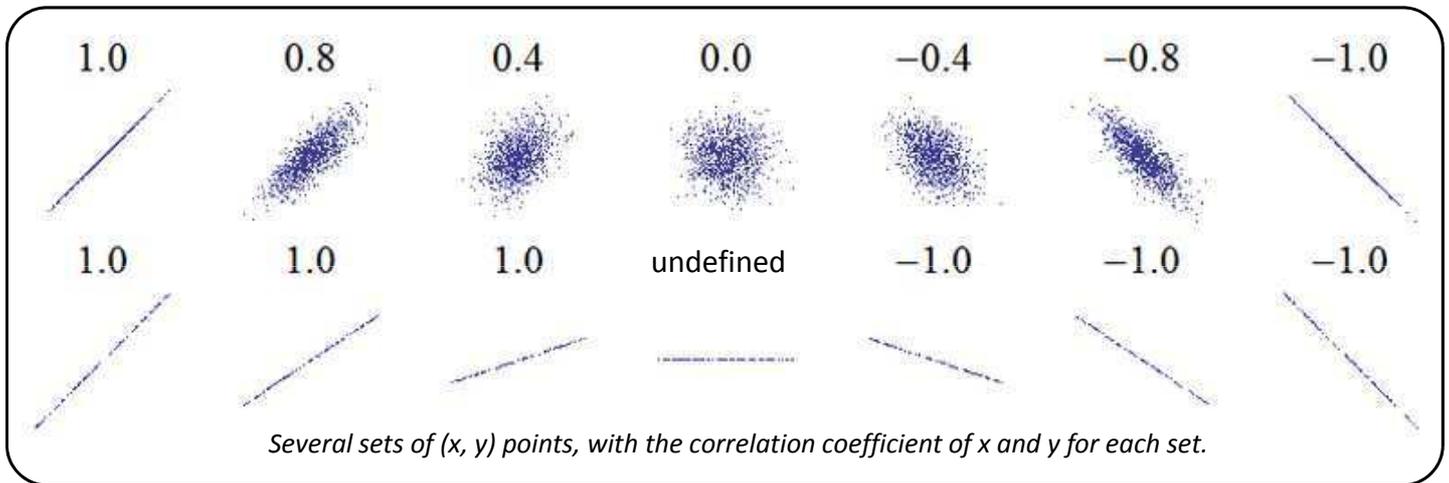
Correlation versus Causation

- Correlation is a mathematical relationship between 2 variables which are measured
- A Correlation of 0, means that there is no linear relationship between the 2 variables and knowing one does not allow prediction of the other
- Strong Correlation may be no more than a statistical association and does not imply causality
- Just because there is a strong correlation between 2 things does not mean that one causes the other. A consistently strong correlation may suggest causation but does not prove it.
- Look at these examples which show a strong correlation but do not prove causality:
- **E.g. 1** With the decrease in the number of pirates we have seen an increase in global warming over the same time period. Does this mean global warming is caused by the decrease in pirates?
- **E.g. 2** With the increase in the number of television sets sold an electrical shop has seen an increase in the number of calculators sold over the same time period. Does this mean that buying a television causes you to buy a calculator?

Criteria for Establishing Causation

- There has to be a strong consistent association found in repeated studies
- The cause has to be plausible and precede the effect in time
- Higher doses will result in stronger responses

Correlation & Slope



Note that the correlation reflects the spread and direction of a linear relationship but not the gradient (slope) of that relationship, **N.B.:** the figure in the centre of the second line has a slope of 0 but in that case the correlation coefficient is undefined because the variance of Y is zero.

The gradient (slope) of the line of best fit is not important when dealing with correlation, except that a vertical or horizontal line of best fit means that the variables are not connected. *[The sign of the slope of the line of best fit will be the same as that of the correlation coefficient because both will be in the same direction.]*

Notes