

Statistics 1

Revision Notes

November 2016

Statistics 1

1	Statistical modelling	5
	Statistical modelling	5
	Definition	5
	Advantages	5
	Disadvantages.....	5
2	Representation of sample data	6
	Variables	6
	Qualitative variables.....	6
	Quantitative variables.....	6
	Continuous variables	6
	Discrete variables	6
	Frequency distributions	6
	Frequency distribution.....	6
	Cumulative frequency	6
	Stem and leaf & back-to-back stem and leaf diagrams	7
	Comparing two distributions from a back to back stem and leaf diagram.	7
	Grouped frequency distributions	7
	Class boundaries and widths	7
	Cumulative frequency curves for grouped data.....	8
	Histograms	9
	Width and height of a bar in centimetres.....	11
3	Mode, mean (and median)	13
	Mode	13
	Mean	13
	Coding.....	14
	Coding and calculating the mean	14
	Median	15
	When to use mode, median and mean	15
	Mode	15
	Median.....	15
	Mean.....	15

4	Median (Q_2), quartiles (Q_1, Q_3) and percentiles	16
	Discrete lists and discrete frequency tables	16
	Interquartile range	16
	Range	16
	Discrete lists.....	16
	Discrete frequency tables	16
	Grouped frequency tables, continuous and discrete data.....	17
	Grouped frequency tables, continuous data	17
	Grouped frequency tables, discrete data	18
	Percentiles.....	19
	Box Plots.....	19
	Outliers.....	19
	Skewness.....	20
	Positive skew	20
	Negative skew.....	21
5	Measures of spread	22
	Range & interquartile range.....	22
	Range	22
	Interquartile range.....	22
	Variance and standard deviation.....	22
	Proof of the alternative formula for variance	22
	Rough checks, $m \pm s$, $m \pm 2s$	22
	Coding and variance	23
	Mean and variance of a combined group.....	24
6	Probability	25
	Relative frequency	25
	Sample spaces, events and equally likely outcomes.....	25
	Probability rules and Venn diagrams.....	25
	Diagrams for two dice etc.	26
	Tree diagrams.....	27
	Independent events.....	28
	To prove that A and B are independent	28
	Exclusive events.....	29
	Number of arrangements of n objects.....	29

7	Correlation	30
	Scatter diagrams.....	30
	Positive, negative, no correlation & line of best fit.....	30
	Product moment correlation coefficient, PMCC	30
	Formulae.....	30
	Coding and the PMCC	31
	Interpretation of the PMCC.....	31
8	Regression	33
	Explanatory and response variables.....	33
	Regression line.....	33
	Least squares regression line	33
	Interpretation	34
	Coding and Regression Lines	35
	Which formula for S_{xy} , S_{xx} , etc.....	36
9	Discrete Random Variables	39
	Random Variables	39
	Continuous and discrete random variables	39
	Continuous random variables.....	39
	Discrete random variables	39
	Probability distributions.....	39
	Cumulative probability distribution.....	40
	Expectation or expected values	40
	Expected mean or expected value of X	40
	Expected value of a function	40
	Expected variance.....	40
	Expectation algebra	40
	Interpretation of <i>expected</i> value	41
	The discrete uniform distribution	42
	Conditions for a discrete uniform distribution.....	42
	Expected mean and variance	43
	Non-standard uniform distribution	43
10	The Normal Distribution $N(\mu, \sigma^2)$	44
	The standard normal distribution $N(0, 1^2)$	44
	The general normal distribution $N(\mu, \sigma^2)$	44
	Use of tables	44

11	Context questions and answers	48
	Accuracy	48
	Statistical models	48
	Histograms	49
	Averages	49
	Skewness.....	50
	Correlation	51
	Discrete uniform distribution.....	53
	Normal distribution.....	53
12	Appendix	55
	$-1 \leq \text{P.M.C.C.} \leq 1$	55
	Cauchy-Schwartz inequality	55
	P.M.C.C. between -1 and $+1$	55
	Regression line and coding	56
	Proof	56
	Normal Distribution, $Z = \frac{X-\mu}{\sigma}$	57

1 Statistical modelling

Statistical modelling

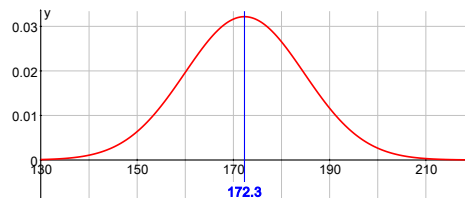
Example: When a die is rolled, we say that the probability of each number is $\frac{1}{6}$. This is a statistical model, but the assumption that each face is equally likely might not be true.

Suppose the die is weighted to increase the chance of a six. We might then find, after experimenting, that the probability of a *six* is $\frac{1}{4}$ and the probability of a *one* is $\frac{1}{12}$, with the probability of other faces remaining at $\frac{1}{6}$, all adding up to 1. In this case we have *refined*, or improved, the model to give a truer picture.

Example: The heights of a large group of adults are measured. The mean is 172.3 cm and the standard deviation is 12.4 cm.

It could be thought that the general shape of the histogram can be modelled by the curve

$$f(x) = \frac{1}{12.4\sqrt{2\pi}} e^{-\frac{1}{2}\left(\frac{x-172.3}{12.4}\right)^2}$$



This might not give a true picture, in which case we would have to change the equation, or *refine the model*.

Definition

A *statistical model* is a *simplification* of a real world situation, usually describing a real world situation using equations. It can be used to make *predictions* about a real world problem. By analysing and *refining* the model an *improved understanding* may be obtained.

Advantages

- the model is *quick* and *easy* to produce
- the model helps our *understanding* of the real world problem
- the model helps us to make *predictions*
- the model helps us to *control* a situation – e.g. railway timetables, air traffic control etc.

Disadvantages

- the model *simplifies* the situation and *only describes a part* of the real world problem.
- the model may *only work in certain situations*, or for a *particular range of values*.

2 Representation of sample data

Variables

Qualitative variables

Non-numerical - e.g. red, blue or long, short etc.

Quantitative variables

Numerical - e.g. length, age, time, number of coins in pocket, etc

Continuous variables

Can take **any** value within a given range - e.g. height, time, age, etc.

Discrete variables

Can only take certain values - e.g. shoe size, cost in £ and p, number of coins.

Frequency distributions

Frequency distribution

A *distribution* is best thought of as a *table*. Thus a *frequency distribution* can be thought of as a *frequency table*, i.e. a list of discrete values and their frequencies.

Example: The number of **M&Ms** is counted in several bags, and recorded in the *frequency distribution*/table below:

<i>number of M&Ms</i>	37	38	39	40	41	42	43
<i>frequency</i>	3	8	11	19	13	7	2

Cumulative frequency

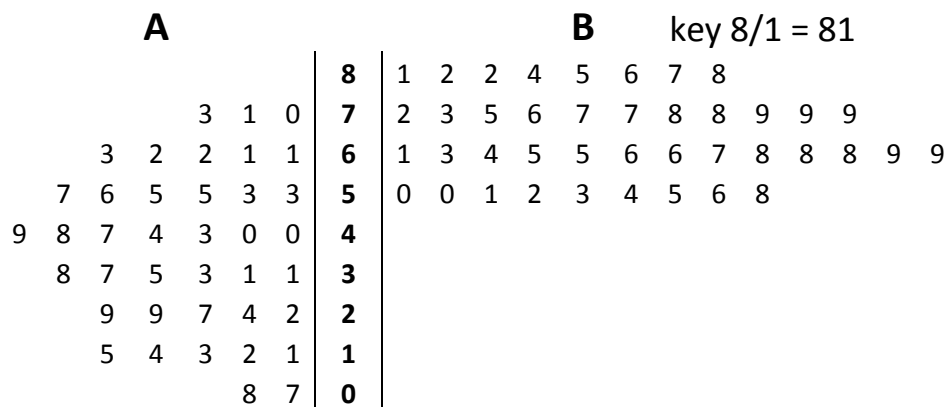
Add up the frequencies as you go down/along the list

<i>number of M&Ms</i>	37	38	39	40	41	42	43
<i>frequency</i>	3	8	11	19	13	7	2
<i>cumulative frequency</i>	3	11	22	41	54	61	63

Stem and leaf & back-to-back stem and leaf diagrams

Line up the digits as leaves on the branches so that it looks like a bar chart.
 The leaves on each branch must be in numerical order, starting from the stem.
 Add a key; e.g. 5|2 means 52, or 4|3 means 4.3 etc.

Comparing two distributions from a back to back stem and leaf diagram.



Comparison:

1. The values in **A** are on average smaller than those in **B**
2. The values in **A** are more spread out than those in **B**.

By looking at the diagram you should always be able to make at least two comparisons, usually one concerning the median and another concerning the spread, often the inter-quartile range.

Grouped frequency distributions

Class boundaries and widths

When deciding class *boundaries* you **must not leave a gap** between one class and another, whether dealing with *continuous* or *discrete* distributions.

For *discrete* distributions avoid leaving gaps between classes by using class boundaries as shown below: $X = 0, 1, 2, 3, 4, 5, 6, 7, \dots$ *discrete*.

Class interval – as given	Class contains	Class boundaries without gaps
0 – 4	0, 1, 2, 3, 4	0 – 4.5
5 – 9	5, 6, 7, 8, 9	4.5 – 9.5
10 – 12	10, 11, 12	9.5 – 12.5
etc		

For *continuous* distributions the class boundaries can be anywhere.

Example: $0 \leq x < 5$, $5 \leq x < 9$, $9 \leq x < 16$, etc.

Note that each interval starts at the point where the previous interval ended, but 5, for example, goes in the second interval $5 \leq x < 9$.

Cumulative frequency curves for grouped data

For a *discrete frequency distribution/table*.

class interval	class boundaries	frequency	class	cumulative frequency
0 - 4	0 to 4.5	27	≤ 4.5	27
5 - 9	4.5 to 9.5	36	≤ 9.5	63
10 - 19	9.5 to 19.5	54	≤ 19.5	117
20 - 29	19.5 to 29.5	49	≤ 29.5	166
30 - 59	29.5 to 59.5	24	≤ 59.5	190
60 - 99	59.5 to 99.5	10	≤ 99.5	200

Plot points at ends of intervals, **(4.5, 27)**, **(9.5, 63)**, **(19.5, 117)** etc. and join points with a smooth curve.

Histograms

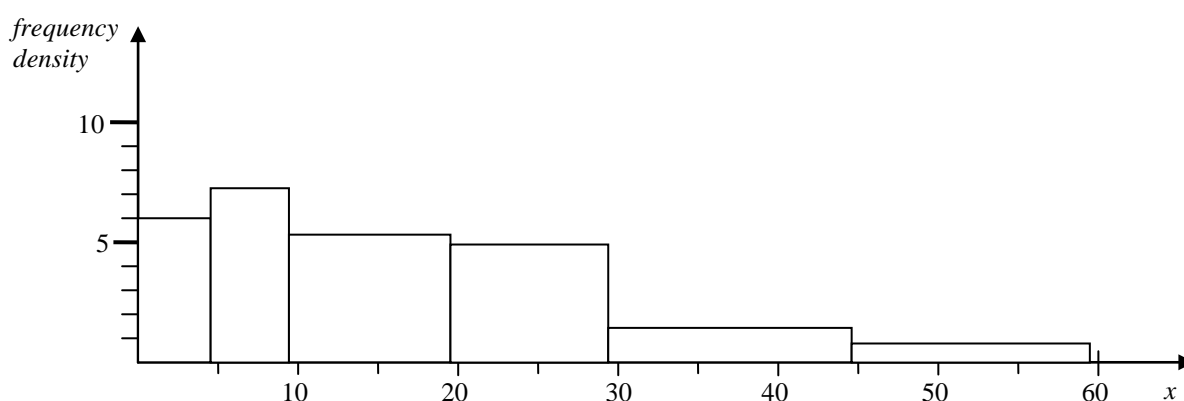
Histograms deal with continuous data. **NO SPACES BETWEEN THE BARS.**

Therefore when we have class intervals of

$3 - 4$, $5 - 7$, 8 , $9 - 12$ they must be thought of as
 $2.5 - 4.5$, $4.5 - 7.5$, $7.5 - 8.5$, $8.5 - 12.5$
 class width 2 class width 3 class width 1 class width 4

To draw a histogram, first draw up a table showing the class intervals, class boundaries, class widths, frequencies and then height of each bar is calculated as $\frac{\text{frequency}}{\text{class width}}$ – as shown below (the *height* is called the *frequency density*):

class interval	class boundaries	class width	frequency	frequency density
0 - 4	0 to 4.5	4.5	27	6
5 - 9	4.5 to 9.5	5	36	7.2
10 - 19	9.5 to 19.5	10	54	5.4
20 - 29	19.5 to 29.5	10	49	4.9
30 - 44	29.5 to 44.5	15	24	1.6
45 - 59	44.5 to 59.5	15	12	0.8

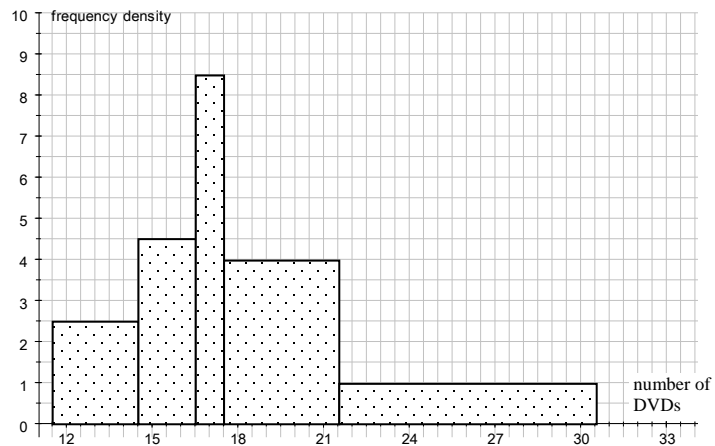


In the above diagram and table we see that the tallest bar is not necessarily the one with the greatest *frequency*, but the one with greatest *frequency density*. The class interval with tallest bar is 4.5 - 9.5, whereas the class interval with the greatest *frequency* is 9.5 - 19.5 (it is in fact the third tallest bar).

In this example we have taken $\text{frequency density} \times \text{class width} = \text{area}$. This is not always the case, but $\text{frequency density} \times \text{class width}$ is **always proportional to area**.

Example: 100 people were asked how many DVDs they owned. The results are shown below.

Number of DVDs	frequency
12 to 14	15
15 to 16	
17	
18 to 21	
22 to 30	



- What are the class boundaries for the interval 12 to 14?
- What is the area of the bar for the class 12 to 14? Compare this with the frequency for this class.
- Fill in the blanks in the frequency table. Check that your frequencies add up to 100.
- Find the number of people who own between 15 and 24 DVDs, inclusive.

Solution:

- Class boundaries for the interval 12 to 14 are 11.5 to 14.5
- Class width of bar for 12 to 14 is $14.5 - 11.5 = 3$,
height of bar = 2.5, \Rightarrow area = $3 \times 2.5 = 7.5$ which is half the frequency, 15.
Thus the frequency is always *twice the area*.
- 15 to 16 is 14.5 to 16.5 width 2, height 4.5 \Rightarrow area 9 \Rightarrow frequency = $2 \times 9 = 18$
17 is 16.5 to 17.5 width 1, height 8.5 \Rightarrow area 8.5 \Rightarrow frequency = $2 \times 8.5 = 17$
18 to 21 is 17.5 to 21.5 width 4, height 4 \Rightarrow area 16 \Rightarrow frequency = $2 \times 16 = 32$
22 to 30 is 21.5 to 30.5 width 9, height 1 \Rightarrow area 9 \Rightarrow frequency = $2 \times 9 = 18$

Check $15 + 18 + 17 + 32 + 18 = 100$, which is correct

- 15 to 24 is 14.5 to 24.5
14.5 to 16.5 and 16.5 to 17.5 and 17.5 to 21.5 gives $f = 18 + 17 + 32$
21.5 to 24.5 has width 3 and height 1 so area = 3 $\Rightarrow f = 2 \times 3 = 6$
 \Rightarrow number between 15 and 24 inclusive is $18 + 17 + 32 + 6 = 73$.

Width and height in centimetres

When the width and height are measured in cm , and the area in cm^2

The **class width** is proportional to the **width** in cm ,

the **frequency density** is proportional to the **height** in cm ,

and the **frequency** is proportional to the **area** in cm^2 of the bar.

The frequency might be *equal* to the area of the bar, or

the frequency might be *double* the area of the bar, or

the frequency might be one *third* of the area of the bar, etc.

Example: A histogram is to be drawn from the following data:

x	f
2 – 4	15
5 – 6	9
7 – 12	18
13 – 16	10

The width of the 2 – 4 interval is 6 cm , and its height is 8 cm .

Find the widths and heights, in cm , for the 5 – 6, the 7 – 12, and the 13 – 16 bars.

Solution:

(a) The 2 – 4 interval has width 6 cm , height 8 $cm \Rightarrow$ **area** = $6 \times 8 = 48 \text{ cm}^2$.

The 2 – 4 interval is 1.5 to 4.5, so has width $3 \propto 6 \text{ cm}$

and frequency $15 \propto 48 \text{ cm}^2$.

(b) The 5 – 6 interval is 4.5 – 6.5, so has width 2,

which is $\frac{2}{3}$ of the width of the 2 to 4 interval, so its width (in cm) is $\frac{2}{3} \times 6 = 4 \text{ cm}$.

The 5 – 6 bar has a **frequency** of 6, which is $\frac{6}{15}$ of the frequency of the 2 – 4 bar, so the **area** of the 5 – 6 bar is $\frac{6}{15} \times 48 = 19.2 \text{ cm}^2$,

\Rightarrow height is $\text{area} \div \text{width} = 19.2 \text{ cm}^2 \div 4 \text{ cm} = 4.8 \text{ cm}$.

(c) We can now fill in the table below, working from left to right.

x	boundaries	width	frequency	width, cm	area, cm^2	height, cm
2 – 4	1.5 – 4.5	3	15	6 cm	$6 \text{ cm} \times 8 \text{ cm} = 48 \text{ cm}^2$	8 cm
5 – 6	4.5 – 6.5	2	9	$\frac{2}{3} \times 6 = 4 \text{ cm}$	$\frac{9}{15} \times 48 = 28.8 \text{ cm}^2$	$28.8 \div 4 = 7.2 \text{ cm}$
7 – 12	6.5 – 12.5	6	18	$\frac{6}{3} \times 6 = 12 \text{ cm}$	$\frac{18}{15} \times 48 = 57.6 \text{ cm}^2$	$57.6 \div 12 = 4.8 \text{ cm}$
13 – 16	12.5 – 16.5	4	10	$\frac{4}{3} \times 6 = 8 \text{ cm}$	$\frac{10}{15} \times 48 = 32 \text{ cm}^2$	$32 \div 8 = 4 \text{ cm}$

Example: A grouped frequency table for the weights of adults has the following entries:

<i>weight kg</i>	50 – 60	60 – 70	70 – 85	...
<i>frequency</i>	60	35	20	...

In a histogram, the bar for the class 50 – 60 kg is 2 cm wide and 9 cm high.

Find the width and height of the bar for the 70 – 85 kg class.

Solution: 50 – 60 is usually taken to mean $50 \leq \text{weight} < 60$

The width of the 50 – 60 class is 10 kg \equiv 2 cm

$$\Rightarrow \text{width of the 70 – 85 class is } 15 \text{ kg} \equiv \frac{15}{10} \times 2 = 3 \text{ cm}$$

The *area* of the 50 – 60 bar is $2 \text{ cm} \times 9 \text{ cm} = 18 \text{ cm}^2 \Leftrightarrow$ *frequency* is 60

\Rightarrow the *frequency* of the 70 – 85 bar is 20

$$\Rightarrow \text{the } \mathbf{area} \text{ of the 70 – 85 bar is } \frac{20}{60} \times 18 = 6 \text{ cm}^2$$

$$\Rightarrow \text{the height of the 70 – 85 bar is } \mathbf{area} \div \mathbf{width} = 6 \text{ cm}^2 \div 3 \text{ cm} = 2 \text{ cm}.$$

Answer the width of 70 – 85 kg bar is 3 cm, and the height is 2 cm.

3 Mode, mean (and median)

Mode

The mode is the value, or class interval, which occurs most often.

Mean

The mean of the values x_1, x_2, \dots, x_n with frequencies f_1, f_2, \dots, f_n the mean is

$$m = \bar{x} = \frac{1}{N} \sum_{i=1}^n x_i f_i, \quad \text{where } N = \sum_{i=1}^n f_i$$

Example: Find the mean for the following table showing the number of children per family.

Solution:

Number of children	Frequency	
x	f	xf
0	5	0
1	8	8
2	12	24
3	18	54
4	9	36
5	4	20
	<hr/>	<hr/>
	56	142

$$\Sigma x_i f_i = 142, \text{ and } N = \Sigma f_i = 56$$

$$\Rightarrow \bar{x} = \frac{142}{56} = 2.54 \text{ to 3 S.F.}$$

Note that the mean m is often written as $m = \frac{1}{N} \sum_1^n x_i$, and the f_i is implied.

In a grouped frequency table you must use the mid-interval value.

Example: The table shows the numbers of children in prep school classes in a town.

Solution:

Number of children	Mid-interval value	Frequency	
	x	f	xf
1 - 10	5.5	5	27.5
11 - 15	13	8	104
16 - 20	18	12	216
21 - 30	25.5	18	459
31 - 40	35.5	<u>11</u>	<u>390.5</u>
		54	1197

$$\Sigma x_i f_i = 1197, \text{ and } N = \Sigma f_i = 54$$

$$\Rightarrow \bar{x} = \frac{1197}{54} = 22.2 \text{ to 3 S.F.}$$

Coding

The weights of a group of people are given as x_1, x_2, \dots, x_n in *kilograms*. These weights are now changed to *grammes* and given as t_1, t_2, \dots, t_n .

In this case $t_i = 1000 \times x_i$ – this is an example of *coding*.

Another example of coding could be $t_i = \frac{x_i - 20}{5}$.

Coding and calculating the mean

With the coding, $t_i = \frac{x_i - 20}{5}$, we are subtracting 20 from each x -value and then dividing the result by 5.

We first find the mean for t_i , and then we reverse the process to find the mean for x_i

\Rightarrow we find the mean for t_i , multiply by 5 and add 20, giving $\bar{x} = 5\bar{t} + 20$

Proof: $t_i = \frac{x_i - 20}{5} \Rightarrow x_i = 5t_i + 20$

$$\bar{x} = \frac{1}{N} \sum_{i=1}^n x_i f_i = \frac{1}{N} \sum_{i=1}^n (5t_i + 20) f_i$$

$$\Rightarrow \bar{x} = \frac{5}{N} \sum_{i=1}^n t_i f_i + \frac{20}{N} \sum_{i=1}^n f_i$$

$$\Rightarrow \bar{x} = 5\bar{t} + 20$$

$$\text{since } \bar{t} = \frac{1}{N} \sum_{i=1}^n t_i f_i \text{ and } N = \sum_{i=1}^n f_i$$

Example: Use the coding $t_i = \frac{x_i - 165}{10}$ to find the mean weight for the following distribution.

Weight, kg	Mid-interval x_i	Coded value $t_i = \frac{x_i - 165}{10}$	Frequency f_i	$t_i f_i$
140 - 150	145	-2	9	-18
150 - 160	155	-1	21	-21
160 - 170	165	0	37	0
170 - 180	175	1	28	28
180 - 190	185	2	11	22
			106	11

$$\Rightarrow \bar{t} = \frac{11}{106}$$

$$\text{and } t_i = \frac{x_i - 165}{10} \Rightarrow \bar{x} = 10\bar{t} + 165 = 10 \times \frac{11}{106} + 165 = 166.0377358$$

\Rightarrow mean weight is 166.04 kg to 2 D.P.

Here the coding simplified the arithmetic

for those who like to work without a calculator!

Median

The median is the middle number in an ordered list. Finding the median is explained in the next section.

When to use mode, median and mean

Mode

You should use the mode if the data is qualitative (colour etc.) or if quantitative (numbers) with a clearly defined mode (or bi-modal). It is not much use if the distribution is fairly even.

Median

You should use this for quantitative data (numbers), when the data is skewed, i.e. when the median, mean and mode are probably not equal, and when there might be extreme values (outliers).

Mean

This is for quantitative data (numbers), and uses all pieces of data. It gives a true measure, and should only be used if the data is fairly symmetrical (not skewed), i.e. the mean could not be affected by extreme values (outliers).

4 Median (Q_2), quartiles (Q_1, Q_3) and percentiles

Discrete lists and discrete frequency tables

To find medians and quartiles

1. Find $k = \frac{n}{2}$ (for Q_2), $\frac{n}{4}$ (for Q_1), $\frac{3n}{4}$ (for Q_3).
2. If k is an integer, use the mean of the k^{th} and $(k + 1)^{\text{th}}$ numbers in the list.
3. If k is not an integer, use the next integer **up**, and find the number with that position in the list.

Interquartile range

The interquartile range, I.Q.R., is $Q_3 - Q_1$.

Range

The range is the largest number minus the smallest (including outliers).

Discrete lists

A discrete list of 10 numbers is shown below:

x	11	13	17	25	33	34	42	49	51	52
-----	----	----	----	----	----	----	----	----	----	----

$n = 10$	for Q_1 ,	$\frac{n}{4} = 2.5$ so use 3 rd number,	\Rightarrow	$Q_1 = 17$	
	for Q_2 ,	$\frac{n}{2} = 5$ so use mean of 5 th and 6 th ,	\Rightarrow	$Q_2 = 33\frac{1}{2}$	median
	for Q_3 ,	$\frac{3n}{4} = 7.5$ so use 8 th number,	\Rightarrow	$Q_3 = 49$	

The interquartile range, I.Q.R., is $Q_3 - Q_1 = 49 - 17 = 32$, and the range is $52 - 11 = 41$.

Discrete frequency tables

x	5	6	7	8	9	10	11	12
f	3	6	8	10	9	8	6	4
<i>cum freq</i>	3	9	17	27	36	44	50	54

$n = 54$	for Q_1 ,	$\frac{n}{4} = 13\frac{1}{2}$ so use 14 th number,	\Rightarrow	$Q_1 = 7$	
	for Q_2 ,	$\frac{n}{2} = 27$ so use mean of 27 th and 28 th ,	\Rightarrow	$Q_2 = 8\frac{1}{2}$	median
	for Q_3 ,	$\frac{3n}{4} = 40\frac{1}{2}$ so use 41 st number,	\Rightarrow	$Q_3 = 10$	

The interquartile range, I.Q.R., is $Q_3 - Q_1 = 10 - 7 = 3$, and the range is $12 - 5 = 7$

Grouped frequency tables, continuous and discrete data

To find medians and quartiles

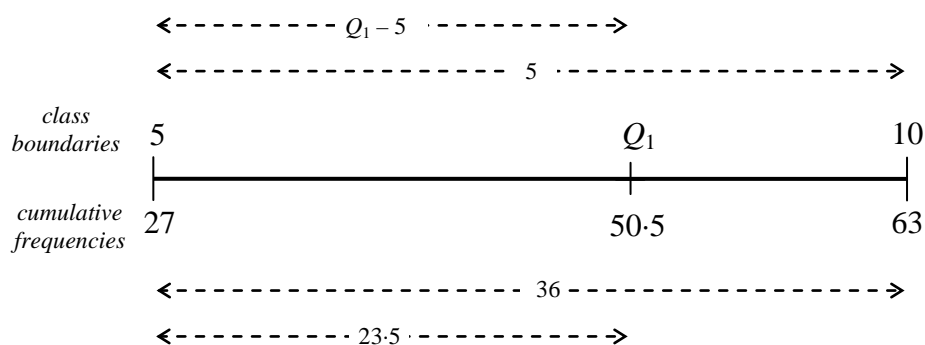
1. Find $k = \frac{n}{2}$ (for Q_2), $\frac{n}{4}$ (for Q_1), $\frac{3n}{4}$ (for Q_3).
2. **Do not round k up or change it in any way.**
3. Use linear interpolation to find median and quartiles – **note** that you must use the correct intervals for discrete data (start at the halves).

Grouped frequency tables, continuous data

class boundaries	frequency	cumulative frequency
$0 \leq x < 5$	27	27
5 to 10	36	63
10 to 20	54	117
20 to 30	49	166
30 to 60	24	190
60 to 100	12	202

With *continuous* data, the end of one interval is the same as the start of the next – no gaps.

To find Q_1 , $n = 202 \Rightarrow \frac{n}{4} = 50\frac{1}{2}$ **do not change it**



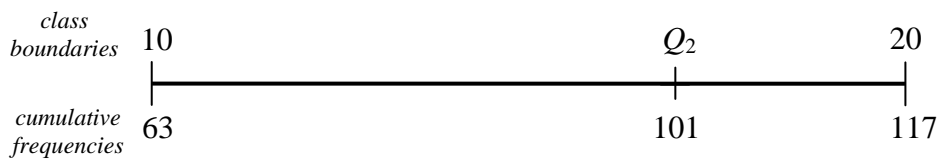
In the diagram we think of the 27th value as the end of the (0 – 5) interval, *and* the start of the (5 – 10) interval.

Also the 63rd value is taken as the end of (5 – 10) *and* the start of (10 – 20)

From the diagram $\frac{Q_1 - 5}{5} = \frac{23 \cdot 5}{36}$

$$\Rightarrow Q_1 = 5 + 5 \times \frac{23 \cdot 5}{36} = 8.263888889 = 8.26 \text{ to 3 s.f.}$$

To find Q_2 , $n = 202 \Rightarrow \frac{n}{2} = 101$ **do not change it**



From the diagram $\frac{Q_2-10}{20-10} = \frac{101-63}{117-63} \Rightarrow Q_2 = 10 + 10 \times \frac{38}{54} = 17.037\dots = 17.0$ to 3 S.F.

Similarly for Q_3 , $\frac{3n}{4} = 151.5$, so Q_3 lies in the interval (20, 30)

$\Rightarrow \frac{Q_3-20}{30-20} = \frac{151.5-117}{166-117} \Rightarrow Q_3 = 20 + 10 \times \frac{34.5}{49} = 27.0408\dots = 27.0$ to 3 S.F.

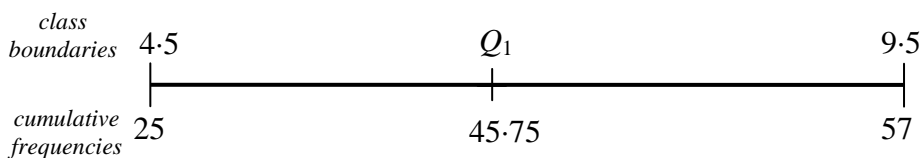
Grouped frequency tables, discrete data

The *discrete* data in grouped frequency tables is treated as *continuous*.

1. Change the class boundaries to the 4.5, 9.5 etc.
2. Proceed as for grouped frequency tables for continuous data.

class interval	class boundaries	frequency	cumulative frequency
0 – 4	0 to 4.5	25	25
5 – 9	4.5 to 9.5	32	57
10 – 19	9.5 to 19.5	51	108
20 – 29	19.5 to 29.5	47	155
30 – 59	29.5 to 59.5	20	175
60 – 99	59.5 to 99.5	8	183

To find Q_1 , $n = 183 \Rightarrow \frac{n}{4} = 45.75$



From the diagram $\frac{Q_1-4.5}{9.5-4.5} = \frac{45.75-25}{57-25}$

$\Rightarrow Q_1 = 4.5 + 5 \times \frac{20.75}{32} = 7.7421875\dots = 7.74$ to 3 S.F.

Q_2 and Q_3 can be found in a similar way.

Percentiles

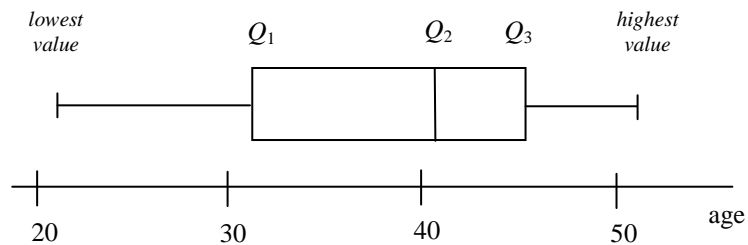
Percentiles are calculated in exactly the same way as quartiles.

Example: For the 90th percentile, find $\frac{90n}{100}$ and proceed as above.

Box Plots

In a group of people the youngest is 21 and the oldest is 52. The quartiles are 32 and 45, and the median age is 41.

We can illustrate this information with a box plot as below – remember to include a scale and label the axis.



Outliers

An outlier is an extreme value. You are not required to remember how to find an outlier – you will always be given a rule.

Example: The ages of 11 children are given below.

age 3 6 12 12 13 14 14 15 17 21 26

$Q_1 = 12$, $Q_2 = 14$ and $Q_3 = 17$.

Outliers are values outside the range $Q_1 - 1.5 \times (Q_3 - Q_1)$ to $Q_3 + 1.5 \times (Q_3 - Q_1)$.

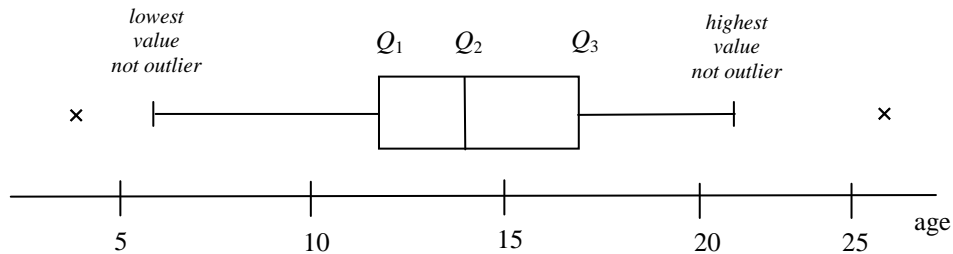
Find any outliers, and draw a box plot.

Solution: Lower boundary for outliers is $12 - 1.5 \times (17 - 12) = 4.5$

Upper boundary for outliers is $17 + 1.5 \times (17 - 12) = 24.5$

\Rightarrow 3 and 26 are the only outliers.

To draw a box plot, put crosses at 3 and 26 (outliers), and draw the lines to 6 (the lowest value which is *not* an outlier), and to 21 (the highest value which is *not* an outlier).



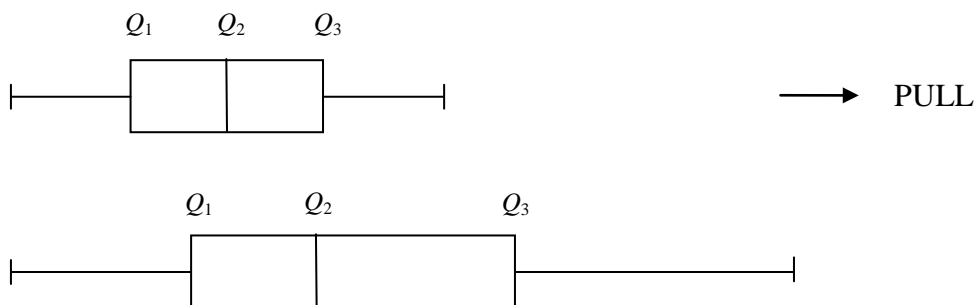
Note that there are other ways of drawing box plots with outliers, but this is the safest and will never be wrong – so why not use it.

Skewness

A distribution which is symmetrical is **not** skewed

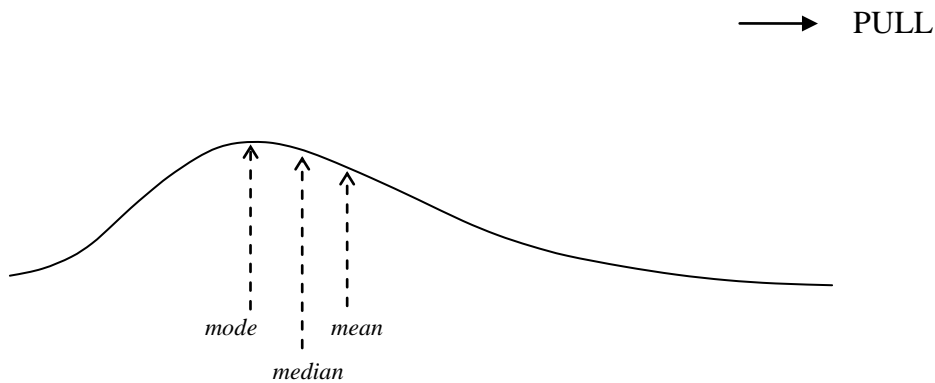
Positive skew

If a symmetrical box plot is stretched in the direction of the positive x -axis, then the resulting distribution has *positive* skew.



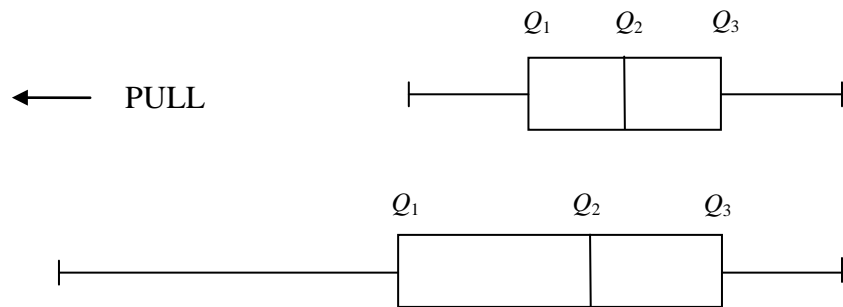
For positive skew the diagram shows that $Q_3 - Q_2 > Q_2 - Q_1$

The same ideas apply for a continuous distribution, and a little bit of thought should show that for *positive skew* $mean > median > mode$.



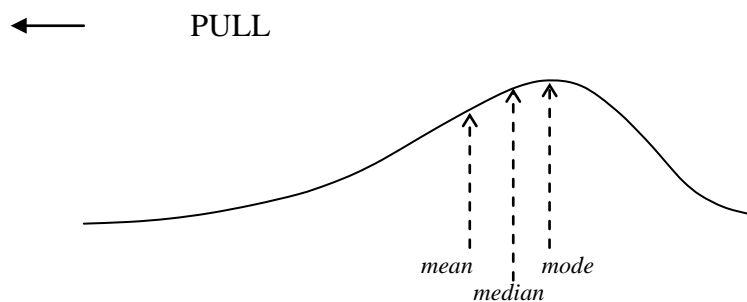
Negative skew

If a symmetrical box plot is stretched in the direction of the negative x -axis, then the resulting distribution has *negative skew*.



For negative skew the diagram shows that $Q_3 - Q_2 < Q_2 - Q_1$

The same ideas apply for a continuous distribution, and a little bit of thought should show that for *negative skew* $mean < median < mode$.



5 Measures of spread

Range & interquartile range

Range

The *range* is found by subtracting the smallest value from the largest value.

Interquartile range

The *interquartile range* is found by subtracting the lower quartile from the upper quartile, so I.Q.R. = $Q_3 - Q_1$.

Variance and standard deviation

Variance is the square of the standard deviation.

$$(\text{sd})_x^2 = \frac{1}{N} \sum (x_i - \bar{x})^2 f_i, \quad \text{or}$$

$$(\text{sd})_x^2 = \frac{1}{N} \sum x_i^2 f_i - \bar{x}^2.$$

When calculating the variance from a table, it is nearly always easier to use the *second* formula. To answer more general questions you should use the *first* formula.

Variance and standard deviation measure the *spread* of the distribution.

Proof of the alternative formula for variance

$$\begin{aligned} (\text{sd})_x^2 &= \frac{1}{N} \sum (x_i - \bar{x})^2 f_i = \frac{1}{N} \sum (x_i^2 - 2x_i\bar{x} + \bar{x}^2) f_i \\ &= \frac{1}{N} \sum x_i^2 f_i - \frac{1}{N} \sum 2x_i\bar{x} f_i + \frac{1}{N} \sum \bar{x}^2 f_i \\ &= \frac{1}{N} \sum x_i^2 f_i - \frac{2\bar{x}}{N} \sum x_i f_i + \frac{\bar{x}^2}{N} \sum f_i \\ &= \frac{1}{N} \sum x_i^2 f_i - 2\bar{x}^2 + \bar{x}^2 = \frac{1}{N} \sum x_i^2 f_i - \bar{x}^2 \end{aligned}$$

since $\bar{x} = \frac{1}{N} \sum x f$ and $N = \sum f$

Rough checks, $m \pm s$, $m \pm 2s$

When calculating a standard deviation, you should check that there is approximately 65 – 70% of the population within 1 s.d. of the mean and approximately 95% within 2 s.d. of the mean.

These approximations are best for a fairly symmetrical distribution.

Coding and variance

Using the coding $t_i = \frac{x_i - k}{a}$ we see that

$$x_i = at_i + k \quad \Rightarrow \quad \bar{x} = a\bar{t} + k \quad \text{as shown in the section on coding and the mean}$$

$$\begin{aligned} \Rightarrow (\text{sd})_x^2 &= \frac{1}{N} \sum (x_i - \bar{x})^2 f_i = \frac{1}{N} \sum ((at_i + k) - (a\bar{t} + k))^2 f_i \\ &= \frac{1}{N} \sum (at_i - a\bar{t})^2 f_i = \frac{a^2}{N} \sum (t_i - \bar{t})^2 f_i \\ \Rightarrow (\text{sd})_x^2 &= a^2 (\text{sd})_t^2 \quad \Rightarrow \quad (\text{sd})_x = a(\text{sd})_t \end{aligned}$$

Notice that subtracting k has no effect, since this is equivalent to translating the graph, and therefore does not change the *spread*, and if all the x -values are divided by a , then we need to multiply $(\text{sd})_t$ by a to find $(\text{sd})_x$.

Example: Find the mean and standard deviation for the following distribution.

Here the x -values are nasty, but if we change them to form $t_i = \frac{x_i - 210}{5}$ then the arithmetic in the last two columns becomes much easier.

x	$t_i = \frac{x_i - 210}{5}$	f	$t_i f_i$	$t_i^2 f_i$
200	-2	12	-24	48
205	-1	23	-23	23
210	0	42	0	0
215	1	30	30	30
220	2	10	20	40
		117	3	141

the mean of t is $\bar{t} = \frac{1}{N} \sum t_i f_i = \frac{3}{117} = \frac{1}{39}$

and the variance of t is $(\text{sd})_t^2 = \frac{1}{N} \sum t_i^2 f_i - \bar{t}^2 = \frac{141}{117} - \left(\frac{1}{39}\right)^2 = 1.204470743$

$(\text{sd})_t = \sqrt{1.204470743} = 1.097483824 = 1.10$ to 3 S.F.

To find \bar{x} , using $t_i = \frac{x_i - 210}{5}$, $\Rightarrow \bar{x} = 5\bar{t} + 210 = 5 \times \frac{1}{39} + 210 = 210.13$ to 2 D.P.

To find the standard deviation of x

$(\text{sd})_x = 5(\text{sd})_t = 5 \times 1.0974838... = 5.49$ to 3 S.F.

We would need to multiply the variance by $5^2 = 25$

$\Rightarrow (\text{sd})_x^2 = 25(\text{sd})_t^2 = 25 \times 1.204470743 = 30.1$ to 3 S.F.

Mean and variance of a combined group

Example: One class of 20 boys has a mean mass of 62 kg with a standard deviation of 3 kg.
A second class of 30 boys has a mean mass of 60 kg with a standard deviation of 4 kg.
Find the mean and standard deviation of the combined class of 50 boys.

Solution: Mean $\bar{x} = \frac{1}{n} \sum x_i$, and variance = $(sd)^2 = \frac{1}{n} \sum x_i^2 - \bar{x}^2$

1st class: mean = 62 = $\frac{1}{20} \sum x_i \Rightarrow \sum x_i = 20 \times 62 = 1240$

2nd class: mean = 60 = $\frac{1}{30} \sum x_i \Rightarrow \sum x_i = 30 \times 60 = 1800$

Combined: mean = $\bar{x} = \frac{1}{n} \sum x_i = \frac{1}{50} (1240 + 1800) = 60.8$ kg

1st class: variance = $3^2 = \frac{1}{20} \sum x_i^2 - 62^2 \Rightarrow \sum x_i^2 = (62^2 + 3^2) \times 20 = 77060$

2nd class: variance = $4^2 = \frac{1}{30} \sum x_i^2 - 60^2 \Rightarrow \sum x_i^2 = (60^2 + 4^2) \times 30 = 108480$

Combined: $\sum x_i^2 = 77060 + 108480 = 185540$, and the mean is 60.8, from above

\Rightarrow variance = $\frac{1}{50} \sum x_i^2 - \text{mean}^2$

= $\frac{1}{50} \times 185540 - 60.8^2 = 14.16$

\Rightarrow standard deviation = $\sqrt{14.16} = 3.76$ to 3 S.F.

Answer : mean = 60.8 kg and standard deviation = 3.76 kg, for combined group.

6 Probability

Relative frequency

After tossing a drawing pin a large number of times the *relative frequency* of it landing point up is $\frac{\text{number of times with point up}}{\text{total number of tosses}}$; this can be thought of as the experimental probability.

Sample spaces, events and equally likely outcomes

A *sample space* is the set of all possible outcomes, all *equally likely*.

An *event* is a set of possible outcomes.

$P(A) = \frac{\text{number of ways } A \text{ can happen}}{\text{total number in the sample space}} = \frac{n(A)}{N}$, where N is number in sample space.

Probability rules and Venn diagrams

All outcomes must be *equally likely* to happen.

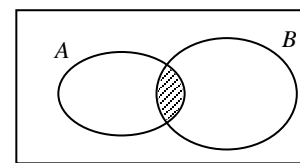
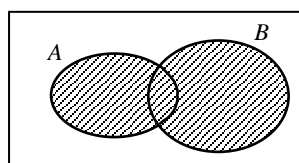
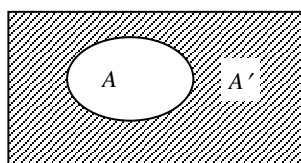
$$P(A) = \frac{n(A)}{N}$$

$$P(A') = P(\text{not } A) = 1 - P(A)$$

A' is the *complement* of A .

$A \cup B$ means *A or B or both*.

$A \cap B$ means *both A and B*,



$$P(A \cup B) = P(A) + P(B) - P(A \cap B)$$

$P(A | B)$ means the probability that A has occurred *given that* we know that B has already occurred, and should always be re-written as

$$P(A | B) = \frac{P(A \cap B)}{P(B)}$$

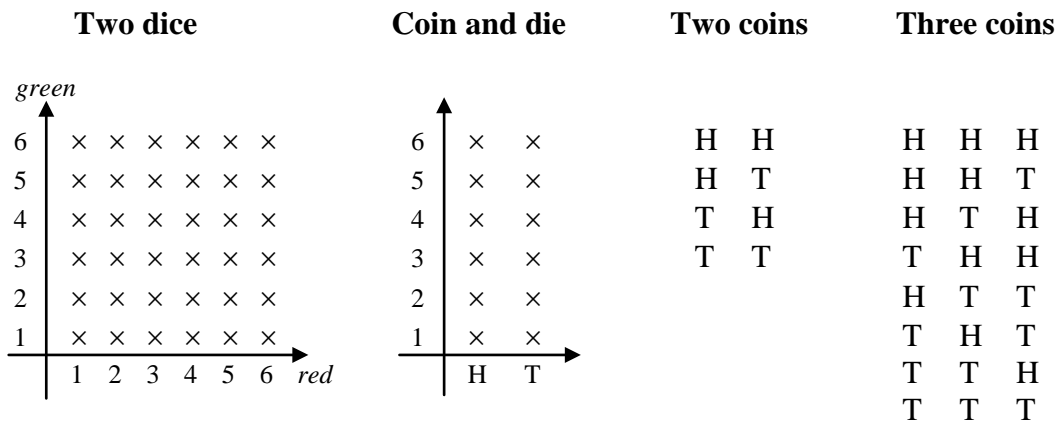
If we know that B has already happened, we can think of B as the new sample space with $n(B)$ elements.

Then the number of ways that A can now occur is $n(A \cap B)$

$$\Rightarrow P(A | B) = \frac{n(A \cap B)}{n(B)} = \frac{\frac{n(A \cap B)}{N}}{\frac{n(B)}{N}} = \frac{P(A \cap B)}{P(B)}$$

Diagrams for two dice etc.

When considering two dice, two spinners or a coin and a die, the following types of diagram are often useful – they ensure that all outcomes are *equally likely* to happen.



From these diagrams, in which all outcomes are *equally likely*, it should be easy to see that

For two dice: $P(\text{total } 10) = \frac{3}{36},$

$$P(\text{red} > \text{green}) = \frac{15}{36},$$

$$P(\text{total } 10 \mid 4 \text{ on green}) = \frac{P(\text{total } 10 \text{ and } 4 \text{ on green})}{P(4 \text{ on green})} = \frac{\frac{1}{36}}{\frac{6}{36}} = \frac{1}{6}.$$

For coin and die: $P(\text{Head and an even number}) = \frac{3}{12}.$

For three dice: $P(\text{exactly two Heads}) = \frac{3}{8}.$

Tree diagrams

The rules for tree diagrams are

Select which branches you need

Multiply along each branch

Add the results of each branch needed.

Make sure that you include enough working to show which branches you are using (method).

Be careful to allow for selection *with and without replacement*.

Example: In the launch of a rocket, the probability of an electrical fault is 0.2. If there is an electrical fault the probability that the rocket crashes is 0.4, and if there is no electrical fault the probability that the rocket crashes is 0.3.

Draw a tree diagram. The rocket takes off, and is seen to crash. What is the probability that there was an electrical fault?

Solution:

Let E be an electrical fault, and C be a crash.

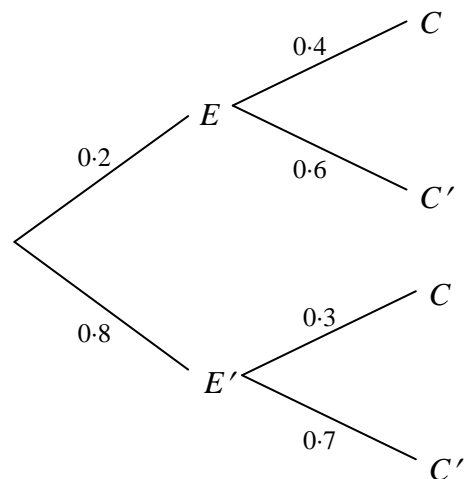
We want to find $P(E | C)$.

$$P(E | C) = \frac{P(E \cap C)}{P(C)}$$

$$P(E \cap C) = 0.2 \times 0.4 = 0.08$$

$$\text{and } P(C) = 0.2 \times 0.4 + 0.8 \times 0.3 = 0.32$$

$$\Rightarrow P(E | C) = \frac{0.08}{0.32} = 0.25$$



Independent events

Definition. **A and B are independent** $\Leftrightarrow P(A \cap B) = P(A) \times P(B)$

It is also true that $P(A | B) = P(A | B') = P(A)$.

A and B are not linked, they have no effect on each other.

To prove that A and B are independent

first find $P(A)$, $P(B)$ and $P(A \cap B)$ **without** assuming that $P(A \cap B) = P(A) \times P(B)$,

second show that $P(A \cap B) = P(A) \times P(B)$.

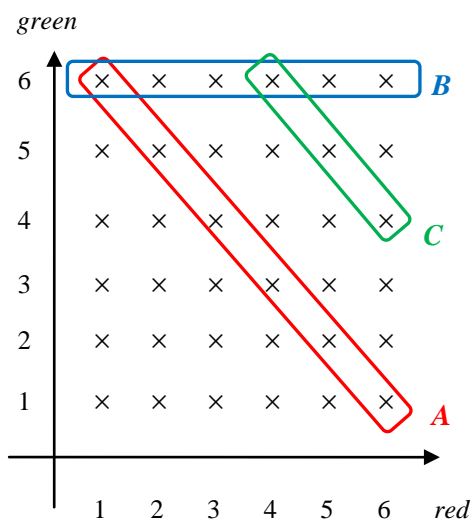
Note: If A and B are **not** independent then $P(A \cap B) \neq P(A) \times P(B)$, and must be found in another way, usually considering sample spaces and/or Venn diagrams.

Example: A red die and a green die are rolled and the total score recorded.

A is the event 'total score is 7', B is the event 'green score is 6' and C is the event 'total score is 10'.

Show that A and B are independent, but B and C are not independent.

Solution: The events A, B and C are shown on the diagram.



$$P(A) = \frac{6}{36} = \frac{1}{6}, \quad P(B) = \frac{6}{36} = \frac{1}{6}$$

$$\text{and } P(A \cap B) = \frac{1}{36} \quad \text{from diagram}$$

$$P(A) \times P(B) = \frac{1}{6} \times \frac{1}{6} = \frac{1}{36} = P(A \cap B)$$

\Rightarrow A and B are independent.

$$P(B) = \frac{6}{36} = \frac{1}{6}, \quad P(C) = \frac{3}{36} = \frac{1}{12}$$

$$\text{and } P(B \cap C) = \frac{1}{36} \quad \text{from diagram}$$

$$P(B) \times P(C) = \frac{1}{6} \times \frac{1}{12} = \frac{1}{72} \neq P(B \cap C)$$

\Rightarrow B and C are **not** independent

Example: A and B are independent events. $P(A) = 0.5$ and $P(A \cap B') = 0.3$. Find $P(B)$.

Solution: $P(A) = 0.5$ and $P(A \cap B') = 0.3$

$$\Rightarrow P(A \cap B) = P(A) - P(A \cap B') = 0.5 - 0.3$$

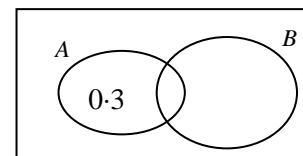
$$\Rightarrow P(A \cap B) = 0.2$$

But $P(A \cap B) = P(A) \times P(B)$ since A and B are independent

$$\Rightarrow P(A \cap B) = 0.5 \times P(B)$$

$$\Rightarrow 0.5 \times P(B) = 0.2$$

$$\Rightarrow P(B) = \frac{0.2}{0.5} = 0.4$$



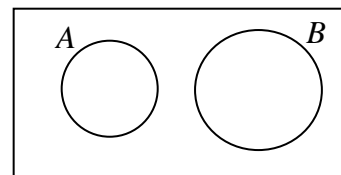
Exclusive events

Definition. **A and B are mutually exclusive**

$$\Leftrightarrow P(A \cap B) = 0$$

i.e. they cannot both occur at the same time

$$\Rightarrow P(A \cup B) = P(A) + P(B)$$



Note: If A and B are **not** exclusive then $P(A \cup B) \neq P(A) + P(B)$, and must be found in another way, usually considering sample spaces and/or Venn diagrams.

Example: $P(A) = 0.3$, $P(B) = 0.6$ and $P(A' \cap B') = 0.1$.

Prove that A and B are mutually exclusive.

Solution: $A' \cap B'$ is shaded in the diagram

$$\Rightarrow P(A' \cap B') = 1 - P(A \cup B)$$

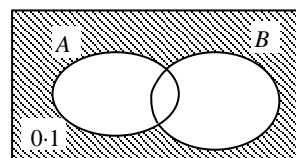
$$\Rightarrow P(A \cup B) = 1 - 0.1 = 0.9$$

$$P(A \cup B) = P(A) + P(B) - P(A \cap B)$$

$$\Rightarrow 0.9 = 0.3 + 0.6 - P(A \cap B)$$

$$\Rightarrow P(A \cap B) = 0$$

$$\Rightarrow A \text{ and } B \text{ are mutually exclusive.}$$



Number of arrangements

Example: A bag contains 5 Red beads, 7 Yellow beads, and 6 White beads. Three beads are drawn *without replacement* from the bag. Find the probability that there are 2 Red beads and 1 Yellow bead.

Solution: These beads can be drawn in any order, *RRY*, *RYR*, *YRR*

$$\Rightarrow P(RRY \text{ or } RYR \text{ or } YRR)$$

$$= P(RRY) + P(RYR) + P(YRR)$$

$$= \frac{5}{18} \times \frac{4}{17} \times \frac{7}{16} + \frac{5}{18} \times \frac{7}{17} \times \frac{4}{16} + \frac{7}{18} \times \frac{5}{17} \times \frac{4}{16} = \frac{35}{408}.$$

You must always remember the possibility of more than one order. In rolling four DICE, exactly TWO SIXES can occur in **six** ways:

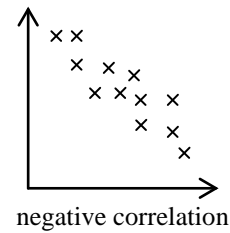
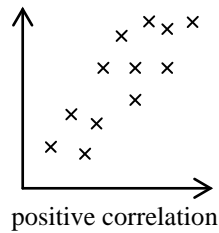
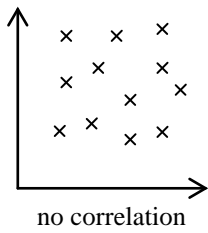
SSNN, SNSN, SNNS, NSSN, NSNS, NNSS, each of which would have the same probability $\left(\frac{1}{6}\right)^2 \times \left(\frac{5}{6}\right)^2 = \frac{25}{1296}$

and so the probability of exactly two sixes with four dice is $6 \times \frac{25}{1296} = \frac{25}{216}$.

7 Correlation

Scatter diagrams

Positive, negative, no correlation & line of best fit.



The pattern of a scatter diagram shows **linear** correlation in a general manner.
A line of best fit can be drawn by eye, **but only when the points nearly lie on a straight line.**

Product moment correlation coefficient, PMCC

Formulae

The S_{**} are all similar to each other and make other formulae simpler to learn and use:

$$S_{xy} = \sum (x_i - \bar{x})(y_i - \bar{y}) = \sum x_i y_i - \frac{1}{N} (\sum x_i) (\sum y_i)$$

$$S_{xx} = \sum (x_i - \bar{x})^2 = \sum x_i^2 - \frac{1}{N} (\sum x_i)^2$$

$$S_{yy} = \sum (y_i - \bar{y})^2 = \sum y_i^2 - \frac{1}{N} (\sum y_i)^2$$

The **PMCC** $r = \frac{S_{xy}}{\sqrt{S_{xx} S_{yy}}}$, $-1 \leq r \leq +1$ for proof, see appendix

To calculate the PMCC first calculate S_{xx} , S_{yy} and S_{xy} using the second formula on each line.

N.B. These formulae are all in the formula booklet.

Coding and the PMCC

To see the effect of coding on the PMCC, it is better to use the first formula on each line.

Example: Investigate the effect of the coding $t = \frac{x-k}{a}$ on the PMCC.

$$\text{Solution: } r_{xy} = \frac{\sum(x_i - \bar{x})(y - \bar{y})}{\sqrt{\sum(x_i - \bar{x})^2 \sum(y_i - \bar{y})^2}}$$

$$t_i = \frac{x_i - k}{a} \Rightarrow x_i = at_i + k \quad \text{and} \quad \bar{x} = a\bar{t} + k$$

$$\Rightarrow r_{xy} = \frac{\sum(at_i + k - (a\bar{t} + k))(y - \bar{y})}{\sqrt{\sum(at_i + k - (a\bar{t} + k))^2 \sum(y_i - \bar{y})^2}}$$

$$= \frac{\sum a(t_i - \bar{t})(y - \bar{y})}{\sqrt{\sum a^2(t_i - \bar{t})^2 \sum(y_i - \bar{y})^2}}$$

$$= \frac{a \sum(t_i - \bar{t})(y - \bar{y})}{\sqrt{a^2 \sum(t_i - \bar{t})^2 \sum(y_i - \bar{y})^2}}$$

$$= \frac{\sum(t_i - \bar{t})(y - \bar{y})}{\sqrt{\sum(t_i - \bar{t})^2 \sum(y_i - \bar{y})^2}}$$

$$= r_{ty}$$

In other words, the coding on x has had no effect on the PMCC. Similarly, coding on y has no effect on the PMCC.

\Rightarrow **Coding has no effect on the PMCC.**

Interpretation of the PMCC

It can be shown that $-1 \leq r \leq +1$

see appendix

if $r = +1$ there is perfect *positive linear correlation*,

if $r = -1$ there is perfect *negative linear correlation*,

if $r = 0$ (or close to zero) there is *no linear correlation*.

PMCC tests to see if there is a **linear connection** between the variables.

For strong correlation, the points on a scatter graph will lie very close to a straight line, and r will be close to 1 or -1 .

Example: Bleep tests are used to measure people's fitness. A higher score means a higher level of fitness. The heart rate, p beats per minute, and bleep score, s , for 12 people were recorded and coded, using $x = p - 60$ and $y = 10s - 50$.

x	0	-6	9	-1	5	8	30	19	28	20	36	23
y	55	62	38	-7	50	44	8	8	3	20	-14	3

$$\Sigma x = 171, \quad \Sigma y = 270, \quad \Sigma x^2 = 4477, \quad \Sigma y^2 = 13540, \quad \Sigma xy = 1020.$$

- Find the PMCC between x and y .
- Write down the PMCC between p and s .
- Explain why your answer to (b) might suggest that there is a linear relationship between p and s .
- Interpret the significance of the PMCC.

Solution: (a)
$$S_{xy} = \sum x_i y_i - \frac{1}{N} (\sum x_i) (\sum y_i) = 1020 - \frac{171 \times 270}{12} = -2827.5$$

$$S_{xx} = \sum x_i^2 - \frac{1}{N} (\sum x_i)^2 = 4477 - \frac{171^2}{12} = 2040.25$$

$$S_{yy} = \sum y_i^2 - \frac{1}{N} (\sum y_i)^2 = 13540 - \frac{270^2}{12} = 7465$$

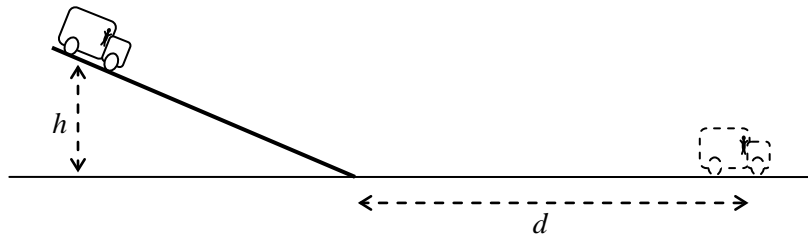
$$\Rightarrow r = \frac{S_{xy}}{\sqrt{S_{xx} S_{yy}}} = \frac{-2827.5}{\sqrt{2040.25 \times 7465}} = -0.7245127195 = -0.725 \text{ to 3 S.F.}$$

- As coding has no effect on the PMCC, the product moment correlation coefficient for p and s is also -0.725 , to 3 S.F.
- $r = -0.725$ is 'quite close' to -1 , and therefore the points on a scatter diagram would lie close to a straight line
 \Rightarrow there is evidence of a linear relation between p and s .
- Always give **both** answers – the technical describing the correlation, and the wordy one which copies out the question. Then, whatever the examiner wants, you should get the mark(s).

There is **negative correlation** between p and s , which means that **as heart rate increases, the bleep score decreases, or people with higher heart rate tend to have lower bleep scores.**

8 Regression

Explanatory and response variables



In an experiment a toy car is released from rest on a ramp from a height of h . The horizontal distance, d , is then measured. The experimenter can control the height, h , and the distance, d , depends on the height chosen.

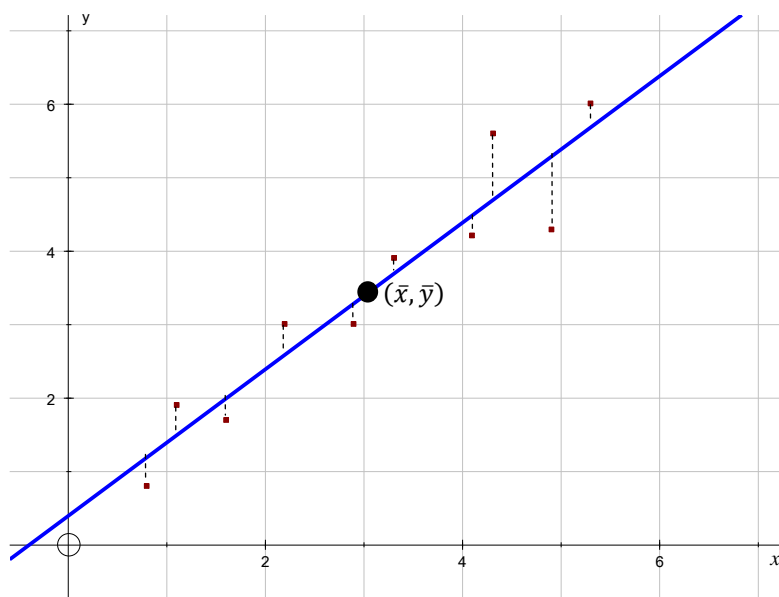
h is called the **explanatory** (or independent) variable and is plotted on the horizontal axis. d is called the **response** (or dependent) variable and is plotted on the vertical axis.

In some cases it may not be possible to control the explanatory variable. For example the temperature at a given time may affect the sales of ice cream; the researcher cannot control the temperature, but it is the temperature which affects the ice cream sales.

Therefore the temperature is the *explanatory* variable, and the ice cream sales is the *response* variable.

Regression line

Least squares regression line



The scatter diagram shows the regression line of y on x . The *regression* line is drawn to minimise the sum of the squares of the vertical distances between the line and the points.

It can be shown that the *regression* line has equation $y = a + bx$, where $b = \frac{S_{xy}}{S_{xx}}$,

also that the *regression* line passes through the ‘mean point’, (\bar{x}, \bar{y}) ,

and so we can find a from the equation $\bar{y} = a + b\bar{x} \Rightarrow a = \bar{y} - b\bar{x}$

Interpretation

In the equation $y = a + bx$

a is the value of y when x is zero (or when x is not present)

b is the amount by which y increases for an increase of 1 in x .

You must write your interpretation in context – copy out the question word for word.

Example: A local authority is investigating the cost of reconditioning its incinerators. Data from 10 randomly chosen incinerators were collected. The variables monitored were the operating time x (in thousands of hours) since last reconditioning and the reconditioning cost y (in £1000). None of the incinerators had been used for more than 3000 hours since last reconditioning.

The data are summarised below,

$$\Sigma x = 25.0, \Sigma x^2 = 65.68, \Sigma y = 50.0, \Sigma y^2 = 260.48, \Sigma xy = 130.64.$$

- (a) Find the equation of the regression line of y on x .
 (b) Give interpretations of a and b .

Solution:

$$(a) \quad \bar{x} = \frac{25.0}{10} = 2.50, \quad \bar{y} = \frac{50.0}{10} = 5.00,$$

$$S_{xy} = 130.64 - \frac{25.0 \times 50.0}{10} = 5.64, \quad S_{xx} = 65.68 - \frac{(25.0)^2}{10} = 3.18$$

$$\Rightarrow b = \frac{S_{xy}}{S_{xx}} = \frac{5.64}{3.18} = 1.773584906$$

$$\Rightarrow a = \bar{y} - b\bar{x} = 5.00 - 1.773584906 \times 2.50 = 0.5660377358$$

$$\Rightarrow \text{regression line equation is } y = 0.566 + 1.77x \quad \text{to 3 s.f.}$$

- (b) a is the cost in £1000 of reconditioning an incinerator which has not been used, so the cost of reconditioning an incinerator which has not been used is £566.

(a is the value of y when x is zero)

b is the increase in cost (in £1000) of reconditioning for every extra 1000 hours of use, so it costs an extra £1774 to recondition an incinerator for every 1000 hours of use. (b is the gradient of the line)

Coding and Regression Lines

If we have the regression line of y on x : $y = a + bx$,

and the coding $p = 2y$, $q = x - 5$

then the regression line of p on q can be found by re-arranging the coding to make x and y the subjects (if necessary), and substituting in the regression line of y on x .

$$p = 2y, \quad q = x - 5$$

$$\Rightarrow y = \frac{p}{2} \quad \text{and} \quad x = q + 5$$

Substitute in $y = a + bx$

$$\Rightarrow \frac{p}{2} = a + b(q + 5)$$

$$\Rightarrow \text{regression line of } p \text{ on } q \text{ is } p = (2a + 10b) + 2bq$$

Example: A drilling machine can run at various speeds, but in general the higher the speed the sooner the drill needs to be replaced.

Over several months, 15 pairs of observations of speed, s revolutions per minute, and the life of the drill, h hours, are collected.

For convenience the data are coded so that $s = x + 20$ and $h = 5y + 100$

The regression line of y on x is calculated as $y = 57.14 - 3.090x$.

Find the regression line of h on s .

Solution: First find x and y in terms of s and h .

$$s = x + 20 \quad \text{and} \quad h = 5y + 100$$

$$\Rightarrow x = s - 20 \quad \text{and} \quad y = \frac{h-100}{5}$$

Then substitute into the equation of the regression line, $y = 57.14 - 3.090x$.

$$\Rightarrow \frac{h-100}{5} = 57.14 - 3.090(s - 20)$$

$$\Rightarrow h - 100 = 285.7 - 15.45s + 309.0$$

$$\Rightarrow h = 694.7 - 15.45s \text{ is the regression line of } h \text{ on } s.$$

Which formula for S_{xy} , S_{xx} , etc.

Mean	$\frac{1}{N} \sum x_i$	
Variance = $(sd)_x^2$	$\frac{1}{N} \sum x_i^2 - \bar{x}^2$	$\frac{1}{N} \sum (x_i - \bar{x})^2$
S_{xy}	$\sum x_i y_i - \frac{\sum x_i \sum y_i}{N}$	$\sum (x_i - \bar{x})(y_i - \bar{y})$
S_{xx}	$\sum x_i^2 - \frac{(\sum x_i)^2}{N}$	$\sum (x_i - \bar{x})^2$
S_{yy}	$\sum y_i^2 - \frac{(\sum y_i)^2}{N}$	$\sum (y_i - \bar{y})^2$

Or we could use the equivalent formulae with $x_i f_i$, $x_i^2 f_i$, $(x_i - \bar{x})^2 f_i$, and $(x_i - \bar{x})(y_i - \bar{y}) f_i$.

In general, we use the formulae in the first column when we want to calculate, and the formulae in the second column for general problems about standard deviation, r , b , etc.

Example: The weights, w , and heights, h of four people are shown in the table below. (In practice there would be more people, but it is easier to show the method with a small number of pairs.)

- Find the mean weight, mean height and the PMCC.
- A fifth person joins the group with weight 77 kg and height 171 cm. Find the new value of S_{hh} .
- Given** that the addition of the fifth person increases the value of S_{wh} , what is the effect of the fifth person on the gradient of the regression line of w on h ?

Solution:

(a)

w	h	wh	w^2	h^2
75	173	12975	5625	29929
69	167	11523	4761	27889
80	174	13920	6400	30276
72	170	12240	5184	28900
296	684	50658	21970	116994

$$\bar{w} = \frac{296}{4} = 74 \quad \text{and} \quad \bar{h} = \frac{684}{4} = 171$$

$$S_{wh} = 50658 - \frac{296 \times 684}{4} = 42, \quad S_{ww} = 21970 - \frac{296^2}{4} = 66, \quad S_{hh} = 116994 - \frac{684^2}{4} = 30$$

$$\Rightarrow r = \frac{42}{\sqrt{66 \times 30}} = 0.943879807 = 0.944 \text{ to 3 s.f.}$$

(b) To answer part (b) we use the formulae in the second column.

$$\begin{aligned} S_{hh} &= \sum_1^4 (h_i - \bar{h})^2 \\ &= (173-171)^2 + (167-171)^2 + (174-171)^2 + (170-171)^2 \end{aligned}$$

The height, $h_5 = 171 \text{ cm}$, of the fifth person does not change the mean height, $\bar{h} = 171 \text{ cm}$. To find the new S_{hh} after the fifth person joins the group, we have to add the term

$$(h_5 - \bar{h})^2 = (171 - 171)^2 = 0.$$

and so zero is added to S_{hh} , and the new $S_{hh} =$ the old S_{hh} .

(c) The regression line of w on h is $w = a + bh$.

The gradient is $b = \frac{S_{wh}}{S_{hh}}$.

We are told that the addition of the fifth person increases the value of S_{wh} , and we have shown that S_{hh} is not changed

\Rightarrow the new gradient = $\frac{\text{bigger}}{\text{same}}$, and is therefore bigger than the original gradient for four people.

Example: A class of 25 students takes a physics exam and a mathematics exam. The equation of the regression line of physics marks, p , against mathematics marks, m , is found using

$$S_{pp} = 1241, S_{pm} = 649 \text{ and } S_{mm} = 847, \bar{p} = 57 \text{ and } \bar{m} = 65.$$

$$b = \frac{S_{pm}}{S_{mm}} = \frac{649}{847} = 0.766$$

- (a) Which is the explanatory variable?
- (b) A pupil who takes the exams at a later date scores 62 in physics and 65 marks in mathematics.
- (i) Find the new value of S_{mm} .
- (ii) What can you say, without detailed calculation, about the new values of S_{pm} and S_{pp} ?

Solution:

- (a) m is the explanatory variable (and p is the response variable).
- (b) (i) The addition of an extra mark, 65, which is equal to the mean, $\bar{m} = 65$, of the original group, will leave the mean unchanged.

$$S_{mm} = \sum(m_i - \bar{m})^2$$

The result, 65 in mathematics, will give an extra term in S_{mm} .

$$\text{The extra term} = (65 - \bar{m})^2 = (65 - 65)^2 = 0$$

Adding 0 to $S_{mm} = \sum(m_i - \bar{m})^2$ has no effect

$$\Rightarrow S_{mm} \text{ is unchanged, } \Rightarrow \text{ new } S_{mm} = 847$$

$$(ii) \text{ new } S_{pm} = \sum(p_i - \bar{p}_{new})(m_i - \bar{m}); \text{ and } \text{new } S_{pp} = \sum(p_i - \bar{p}_{new})^2$$

As we do not know the new value \bar{p}_{new} , we cannot, without detailed calculation, say anything about S_{pm} or S_{pp} .

9 Discrete Random Variables

Random Variables

A random variable must take a numerical value:

Examples: the number on a single throw of a die
the height of a person
the number of cars travelling past a fixed point in a certain time

But **not** the *colour* of hair as this is not a number

Continuous and discrete random variables

Continuous random variables

A continuous random variable is one which can take **any** value in a certain interval;

Examples: height, time, weight.

Discrete random variables

A discrete random variable can only take certain values in an interval

Examples: Score on die (1, 2, 3, 4, 5, 6)
Number of coins in pocket (0, 1, 2, ...)

Probability distributions

A probability *distribution* (or *table*) is the set of possible outcomes together with their probabilities, similar to a frequency *distribution* (or *table*).

Example:

score on two dice, X	2	3	4	5	6	7	8	9	10	11	12
probability, $f(x)$	$\frac{1}{36}$	$\frac{2}{36}$	$\frac{3}{36}$	$\frac{4}{36}$	$\frac{5}{36}$	$\frac{6}{36}$	$\frac{5}{36}$	$\frac{4}{36}$	$\frac{3}{36}$	$\frac{2}{36}$	$\frac{1}{36}$

is the probability *distribution* (*table*) for the random variable, X , the total score on two dice.

Note that the sum of the probabilities **must** be 1, i.e. $\sum_{x=2}^{12} P(X = x) = 1$.

Cumulative probability distribution

Just like cumulative frequencies, the cumulative probability, F , that the total score on two dice is less than or equal to 4 is $F(4) = P(X \leq 4) = P(X = 2, 3, 4) = \frac{1}{36} + \frac{2}{36} + \frac{3}{36} = \frac{6}{36} = \frac{1}{6}$.

Note that $F(4.3)$ means $P(X \leq 4.3)$ and seeing as there are no scores between 4 and 4.3 this is the same as $P(X \leq 4) = F(4)$.

You are expected to recognise that capital $F(X)$ means the cumulative probability.

Expectation or expected values

Expected mean or expected value of X .

For a discrete probability distribution the expected mean of X , or the expected value of X is

$$\mu = E[X] = \sum x_i p_i$$

Expected value of a function

The expected value of any function, $h(X)$, is defined as

$$E[h(X)] = \sum h(x_i) p_i$$

Note that for any constant, k , $E[k] = k$,

since $\sum k p_i = k \sum p_i = k \times 1 = k$

Expected variance

The expected variance of X is

$$\sigma^2 = \text{Var}[X] = \sum (x_i - \mu)^2 p_i = \sum x_i^2 p_i - \mu^2, \quad \text{or}$$

$$\sigma^2 = \text{Var}[X] = E[(X - \mu)^2] = E[X^2] - \mu^2 = E[X^2] - (E[X])^2$$

Expectation algebra

$$E[aX + b] = \sum (ax_i + b) p_i = \sum ax_i p_i + \sum b p_i = a \sum x_i p_i + b$$

$$= aE[X] + b$$

since $\sum x_i p_i = \mu$, and $\sum p_i = 1$

$$\text{Var}[aX + b] = E[(aX + b)^2] - (E[aX + b])^2$$

$$= E[(a^2 X^2 + 2abX + b^2)] - (aE[X] + b)^2$$

$$= \{a^2 E[X^2] + 2ab E[X] + E[b^2]\} - \{a^2 (E[X])^2 + 2ab E[X] + b^2\}$$

$$= a^2 E[X^2] - a^2 (E[X])^2 = a^2 \{E[X^2] - (E[X])^2\} = a^2 \text{Var}[X]$$

Thus we have two important results:

$$E[aX + b] = aE[X] + b$$

$\text{Var}[aX + b] = a^2 \text{Var}[X]$ which are equivalent to the results for coding done earlier.

Example: A fair die is rolled and the score recorded.

- (a) Find the expected mean and variance for the score, X .
- (b) A ‘prize’ is awarded which depends on the score on the die. The value of the prize is $\$Z = 3X - 6$. Find the expected mean and variance of Z .

Solution:

(a)	score	probability	$x_i p_i$	$x_i^2 p_i$
	1	$\frac{1}{6}$	$\frac{1}{6}$	$\frac{1}{6}$
	2	$\frac{1}{6}$	$\frac{2}{6}$	$\frac{4}{6}$
	3	$\frac{1}{6}$	$\frac{3}{6}$	$\frac{9}{6}$
	4	$\frac{1}{6}$	$\frac{4}{6}$	$\frac{16}{6}$
	5	$\frac{1}{6}$	$\frac{5}{6}$	$\frac{25}{6}$
	6	$\frac{1}{6}$	$\frac{6}{6}$	$\frac{36}{6}$
			$\frac{21}{6}$	$\frac{91}{6}$

$$\Rightarrow \mu = E[X] = \sum x_i p_i = \frac{21}{6} = 3\frac{1}{2}$$

$$\text{and } \sigma^2 = E[X^2] - (E[X])^2 = \frac{91}{6} - \left(\frac{21}{6}\right)^2 = \frac{35}{12} = 2\frac{11}{12}$$

\Rightarrow The expected mean and variance for the score, X , are $\mu = 3\frac{1}{2}$ and $\sigma^2 = 2\frac{11}{12}$

(b) $Z = 3X - 6$

$$\Rightarrow E[Z] = E[3X - 6] = 3E[X] - 6 = 10\frac{1}{2} - 6 = 4\frac{1}{2}$$

$$\text{and } \text{Var}[Z] = \text{Var}[3X - 6] = 3^2 \text{Var}[X] = 9 \times \frac{35}{12} = 26\frac{1}{4}$$

\Rightarrow The expected mean and variance for the prize, $\$Z$, are $\mu = 4\frac{1}{2}$ and $\sigma^2 = 26\frac{1}{4}$

Interpretation of *expected value*

In the example above, assume that the die was rolled N times, where N is large, and the score recorded each time.

If the frequencies were *exactly* $\frac{1}{6}N$, then the mean of the recorded scores would equal $E[X] = 3.5$, and the variance of the recorded scores would equal $\text{Var}[X] = 2\frac{11}{12}$.

Also, if the value of the prize, Z , was calculated each time, then the mean of the values would equal $E[Z] = 4.5$, and the variance of the values would equal $\text{Var}[X] = 26\frac{1}{4}$.

Example: A die is rolled until a 6 is thrown, or until the die has been rolled 4 times. Find the expected mean and standard deviation of the number of times the die is rolled.

Solution: Expected mean $= \mu = E[X] = \sum x_i p_i$
 and expected variance $= \sigma^2 = E[X^2] - (E[X])^2 = \sum x_i^2 p_i - \mu^2$

number of rolls, x_i	outcome	probability p_i	$x_i p_i$	$x_i^2 p_i$
1	(six)	$\frac{1}{6}$	$\frac{1}{6}$	$\frac{1}{6}$
2	(not six)(six)	$\frac{5}{6} \times \frac{1}{6}$	$\frac{10}{36}$	$\frac{20}{36}$
3	(not six)(not six)(six)	$(\frac{5}{6})^2 \times \frac{1}{6}$	$\frac{75}{216}$	$\frac{225}{216}$
4	(not six)(not six)(not six) 4 th roll could be anything	$(\frac{5}{6})^3$	$\frac{500}{216}$	$\frac{2000}{216}$
			$\frac{671}{216}$	$\frac{2381}{216}$

$$\Rightarrow \mu = \sum x_i p_i = \frac{671}{216} = 3.11 \text{ to 3 S.F.}$$

$$\text{and } \sigma^2 = \sum x_i^2 p_i - \mu^2 = \frac{2381}{216} - \left(\frac{671}{216}\right)^2 = 1.372920953$$

$$\Rightarrow \sigma = \sqrt{1.372920953} = 1.17 \text{ to 3 S.F.}$$

Answer : Expected mean number of rolls is 3.11, with standard deviation 1.17

The discrete uniform distribution

Conditions for a discrete uniform distribution

- The discrete random variable X is defined over a set of n distinct values
- Each value is equally likely, with probability $1/n$.

Example: The random variable X is defined as the score on a single die. X is a discrete uniform distribution on the set $\{1, 2, 3, 4, 5, 6\}$

The probability distribution is

Score	1	2	3	4	5	6
Probability	$\frac{1}{6}$	$\frac{1}{6}$	$\frac{1}{6}$	$\frac{1}{6}$	$\frac{1}{6}$	$\frac{1}{6}$

Expected mean and variance

For a discrete uniform random variable, X defined on the set $\{1, 2, 3, 4, \dots, n\}$,

X	1	2	3	4	n
Probability	$\frac{1}{n}$	$\frac{1}{n}$	$\frac{1}{n}$	$\frac{1}{n}$			$\frac{1}{n}$

By symmetry we can see that the Expected mean $= \mu = E[X] = \frac{1}{2}(n+1)$,

$$\begin{aligned} \text{or } \mu &= E[X] = \sum x_i p_i = 1 \times \frac{1}{n} + 2 \times \frac{1}{n} + 3 \times \frac{1}{n} + \dots + n \times \frac{1}{n} \\ &= (1 + 2 + 3 + \dots + n) \times \frac{1}{n} = \frac{1}{2}n(n+1) \times \frac{1}{n} = \frac{1}{2}(n+1) \end{aligned} \quad \text{long winded method}$$

$$\begin{aligned} \text{Var}[X] &= \sigma^2 = E[X^2] - (E[X])^2 = \sum x_i^2 p_i - \mu^2 \\ &= \left(1^2 \times \frac{1}{n} + 2^2 \times \frac{1}{n} + 3^2 \times \frac{1}{n} + \dots + n^2 \times \frac{1}{n}\right) - \left(\frac{1}{2}(n+1)\right)^2 \\ &= (1^2 + 2^2 + 3^2 + \dots + n^2) \times \frac{1}{n} - \left(\frac{1}{2}(n+1)\right)^2 \\ &= \frac{1}{6}n(n+1)(2n+1) \times \frac{1}{n} - \frac{1}{4}(n+1)^2 \quad \text{since } \sum i^2 = \frac{1}{6}n(n+1)(2n+1) \\ &= \frac{1}{24}(n+1)[4(2n+1) - 6(n+1)] \\ &= \frac{1}{24}(n+1)(2n-2) = \frac{1}{12}(n+1)(n-1) \\ \Rightarrow \quad \text{Var}[X] &= \sigma^2 = \frac{1}{12}(n^2 - 1) \end{aligned}$$

These formulae can be quoted in an exam (if you learn them!).

Non-standard uniform distribution

The formulae can sometimes be used for non-standard uniform distributions.

Example: X is the score on a fair 10 sided spinner. Define $Y = 5X + 3$.

Find the mean and variance of Y .

Y is the distribution $\{8, 13, 18, \dots, 53\}$, all with the same probability $\frac{1}{10}$.

Solution: X is a (standard) discrete uniform distribution on the set $\{1, 2, 3, \dots, 10\}$

$$\Rightarrow E[X] = \frac{1}{2}(n+1) = 5\frac{1}{2}$$

$$\text{and } \text{Var}[X] = \frac{1}{12}(n^2 - 1) = \frac{99}{12} = 8\frac{1}{4}$$

$$\Rightarrow E[Y] = E[5X + 3] = 5E[X] + 3 = 30\frac{1}{2}$$

$$\text{and } \text{Var}[Y] = \text{Var}[5X + 3] = 5^2 \text{Var}[X] = 25 \times \frac{99}{12} = 206\frac{1}{4}$$

$$\Rightarrow \text{mean and variance of } Y \text{ are } 27\frac{1}{2} \text{ and } 206\frac{1}{4}.$$

10 The Normal Distribution $N(\mu, \sigma^2)$

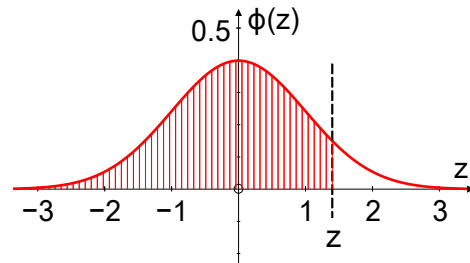
The standard normal distribution $N(0, 1^2)$

The diagram shows the standard normal distribution

Mean, $\mu, = 0$

Standard deviation, $\sigma, = 1$

The tables give the area, $\Phi(z)$, from $-\infty$ up to z ;
this is the probability $P(Z < z)$



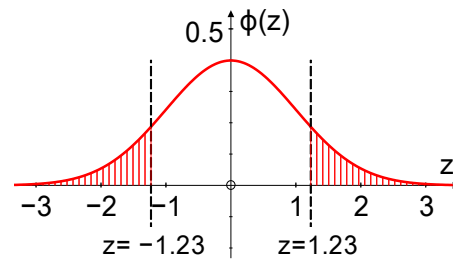
To find other probabilities, sketch the curve and use your head.

Example: $P(Z < -1.23)$

= area up to $-1.23 = \Phi(-1.23)$

= area beyond $+1.23 = 1 - \Phi(+1.23)$

= $1 - 0.8907 = 0.1093$ to 4 D.P. from tables

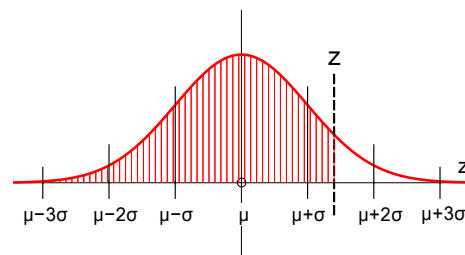


The general normal distribution $N(\mu, \sigma^2)$

Use of tables

To use the tables for a Normal distribution with mean μ and standard deviation σ

We use $Z = \frac{X - \mu}{\sigma}$ (see appendix) and look in the tables under this value of Z



Example: The length of life (in months) of Blowdri's hair driers is approximately Normally distributed with mean 90 months and standard deviation 15 months, $N(90, 15^2)$.

- (a) Each drier is sold with a 5 year guarantee. What proportion of driers fail before the guarantee expires?
- (b) The manufacturer decides to change the length of the guarantee so that no more than 1% of driers fail during the guarantee period. How long should he make the guarantee?

Solution:

(a) X is the length of life of drier $\Rightarrow X \sim N(90, 15^2)$

5 years = 60 months

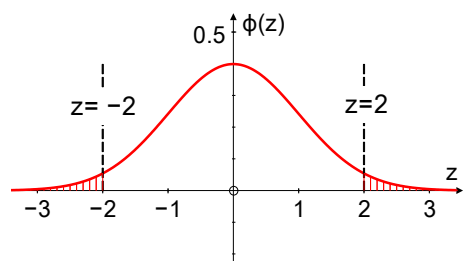
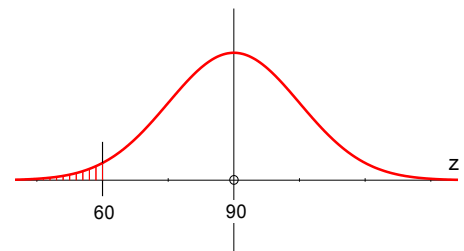
\Rightarrow we want $P(X < 60) = \text{area up to } 60$

$$\Rightarrow Z = \frac{X - \mu}{\sigma} = \frac{60 - 90}{15} = -2.0$$

so we want area to left of $Z = -2$

$$= \Phi(-2) = 1 - \Phi(2)$$

$$= 1 - 0.9772 = 0.0228 \text{ to 4 D.P. from tables.}$$



\Rightarrow the proportion of hair driers failing during the guarantee period is 0.0228 to 4 D.P.

(b) Let the length of the guarantee be t years

\Rightarrow we need $P(X < t) = 0.01$.

We need the value of Z such that $\Phi(Z) = 0.01$

From the tables $Z = -2.3263$ to 4 D.P. from tables

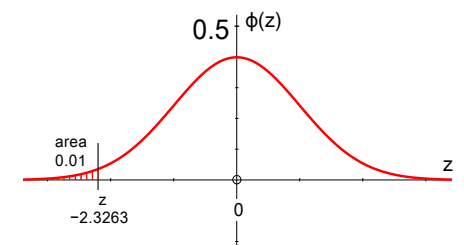
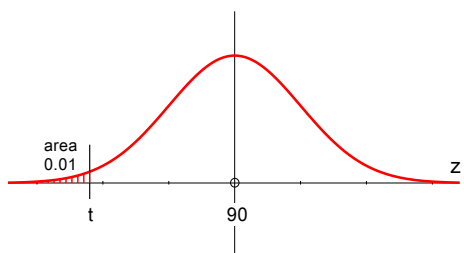
(remember to look in the small table after the Normal tables)

Standardising the variable

$$\Rightarrow Z = \frac{X - \mu}{\sigma} = \frac{t - 90}{15}$$

$$\Rightarrow \frac{t - 90}{15} = -2.3263 \text{ to 4 D.P. from tables}$$

$$\Rightarrow t = 55.1 \text{ to 3 S.F.}$$



so the manufacturer should give a guarantee period of 55 months (4 years 7 months)

Example: The results of an examination were Normally distributed. 10% of the candidates had more than 70 marks and 20% had fewer than 35 marks.

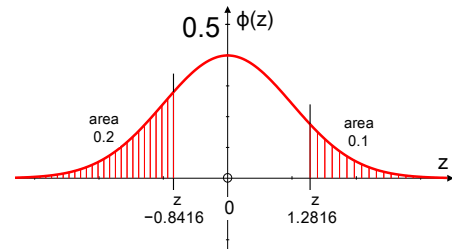
Find the mean and standard deviation of the marks.

Solution:

First we need the values from the tables

$$\Rightarrow \Phi(-0.8416) = 0.2,$$

$$\text{and } 1 - \Phi(1.2816) = 0.1$$



Using $Z = \frac{X - \mu}{\sigma}$ we have

$$-0.8416 = \frac{35 - \mu}{\sigma}$$

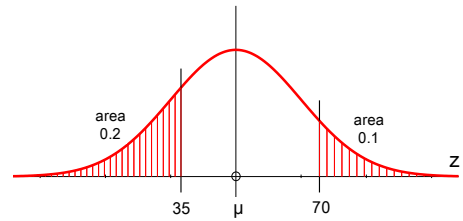
$$\Rightarrow \mu = 35 + 0.8416\sigma \quad \mathbf{I}$$

$$\text{and } 1.2816 = \frac{70 - \mu}{\sigma}$$

$$\Rightarrow \mu = 70 - 1.2816\sigma \quad \mathbf{II}$$

$$\mathbf{I} - \mathbf{II} \Rightarrow 0 = -35 + 2.1232\sigma$$

$$\Rightarrow \sigma = 16.5 \text{ and } \mu = 48.9 \text{ to 3 s.f.}$$



simultaneous equations

Example: The weights of chocolate bars are normally distributed with mean 205 g and standard deviation 2.6 g. The stated weight of each bar is 200 g.

- (a) Find the probability that a single bar is underweight.
- (b) Four bars are chosen at random. Find the probability that fewer than two bars are underweight.

Solution:

- (a) Let W be the weight of a chocolate bar, $W \sim N(205, 2.6^2)$.

$$Z = \frac{W - \mu}{\sigma} = \frac{200 - 205}{2.6} = -1.9230769\dots$$

$$P(W < 200) = P(Z < -1.92) = 1 - \Phi(1.92) = 1 - 0.9726$$

\Rightarrow probability of an underweight bar is 0.0274.

- (b) We want the probability that 0 or 1 bars chosen from 4 are underweight.

Let U be underweight and C be correct weight.

$$P(1 \text{ underweight}) = P(CCCU) + P(CCUC) + P(CUCC) + P(UCCC)$$

$$= 4 \times 0.0274 \times 0.9726^3 = 0.1008354753$$

$$P(0 \text{ underweight}) = 0.9726^4 = 0.7403600224$$

\Rightarrow the probability that fewer than two bars are underweight = 0.841 to 3 s.f.

11 Context questions and answers

Accuracy

You are required to *give your answers to an appropriate degree of accuracy*.

There is no hard and fast rule for this, but the following guidelines should never let you down.

1. If stated in the question give the required degree of accuracy.
2. When using a calculator, give 3 S.F.
unless finding S_{xx} , S_{xy} etc. in which case you can give more figures – you should use *all* figures when finding the PMCC or the regression line coefficients.
3. Sometimes it is appropriate to give a mean to 1 or 2 D.P. rather than 3 S.F.
4. When using the tables and doing simple calculations (which do not *need* a calculator), you should give 4 D.P.

Statistical models

Question 1

- (a) Explain briefly what you understand by
- (i) a statistical experiment,
 - (ii) an event.
- (b) State one advantage and one disadvantage of a statistical model.

Answer

- (a) (i) a test/investigation/process for collecting data to provide evidence to test a hypothesis.
(ii) A subset of possible outcomes of an experiment, e.g. a total of 6 on two dice.
- (b) Quick, cheap can vary the parameters and predict
Does not replicate real world situation in every detail.

Question 2

Statistical models can be used to describe real world problems. Explain the process involved in the formulation of a statistical model.

Answer

- Observe real world problem
- Devise a statistical model and collect data
- Compare observed against expected outcomes and test the model
- Refine model if necessary

Question 3

- (a) Write down two reasons for using statistical models.
- (b) Give an example of a random variable that could be modelled by
 - (i) a normal distribution,
 - (ii) a discrete uniform distribution.

Answer

- (a) To simplify a real world problem
To improve understanding / describe / analyse a real world problem
Quicker and cheaper than using real thing
To predict possible future outcomes
Refine model / change parameters possible *Any 2*
 - (b) (i) height, weight, etc. (ii) score on a face after rolling a fair die
-

Histograms

Question 1

Give a reason to justify the use of a histogram to represent these data.

Answer

The variable (minutes delayed) is **continuous**.

Averages

Question 1

Write down which of these averages, mean or median, you would recommend the company to use. Give a reason for your answer.

Answer

The median, because the data is **skewed**.

Question 2

State whether the newsagent should use the median and the inter-quartile range or the mean and the standard deviation to compare daily sales. Give a reason for your answer.

Answer

Median & IQR as the data is likely to be **skewed**

Question 3

Compare and contrast the attendance of these 2 groups of students.

Answer

Median 2nd group < Median 1st group;

Mode 1st group > Mode 2nd group;

2nd group had larger spread/IQR than 1st group

Only 1 student attends all classes in 2nd group

Question 4

Compare and contrast these two box plots.

Answer

Median of Northcliffe is greater than median of Seaview.

Upper quartiles are the same

IQR of Northcliffe is less than IQR of Seaview

Northcliffe positive skew, Seaview negative skew – using mean and median

or Northcliffe symmetrical, Seaview positive skew – using quartiles

Notice that it is possible to have different skewness depending on the way of measuring.

Range of Seaview greater than range of Northcliffe

any 3 acceptable comments

Skewness

Question 1

Comment on the skewness of the distribution of bags of crisps sold per day.
Justify your answer.

Answer

$Q_2 - Q_1 = 7$; $Q_3 - Q_2 = 11$; $Q_3 - Q_2 > Q_2 - Q_1$ so positive skew.

Question 2

Give two other reasons why these data are negatively skewed.

Answer

For negative skew; Mean < median < mode: $49.4 < 52 < 56$

$Q_3 - Q_2 < Q_2 - Q_1$: $8 < 17$

Question 3

Describe the skewness of the distribution. Give a reason for your answer.

Answer

No skew *or* slight negative skew.

$$0.22 = Q_3 - Q_2 \approx Q_2 - Q_1 = 0.23 \quad \text{or} \quad 0.22 = Q_3 - Q_2 < Q_2 - Q_1 = 0.23$$

or mean (3.23) \approx median (3.25), *or* mean (3.23) $<$ median (3.25)

When the skew is very slight, you can say “no skew” or “slight skew”.

Correlation

Question 1

Give an interpretation of your PMCC (−0.976)

Answer

Negative skew \Leftrightarrow as height increases, temperature decreases (**must be in context**).

Question 2

Give an interpretation of this value, PMCC = −0.862.

Answer

Negative skew \Leftrightarrow as sales at one petrol station increase, sales at the other decrease (**must be in context**).

Question 3

Give an interpretation of your correlation coefficient, 0.874.

Answer

Positive skew \Leftrightarrow taller people tend to be more confident (**must be in context**).

Question 4

Comment on the assumption that height and weight are independent.

Answer

Evidence (in question) suggests height and weight are positively correlated/linked, therefore the assumption of independence is not sensible (**must be in context**).

Regression

Question 1

Suggest why the authority might be cautious about making a prediction of the reconditioning cost of an incinerator which had been operating for 4500 hours since its last reconditioning.

Answer **4500** is well **outside** the **range of observed** values, and there is no evidence that the model will apply.

Question 2

Give an interpretation of the slope, 0.9368, and the intercept, 19, of your regression line.

Answer

The slope, b – for every **extra hour** of **practice 0.9368 fewer errors** will be made

The intercept, a – **without practice 19 errors** will be made.

Question 3

Interpret the value of b (coefficient of x in regression line).

Answer

3 extra ice-creams are sold for **every 1°C increase** in temperature

Question 4

At 1 p.m. on a particular day, the highest temperature for 50 years was recorded. Why you should not use the regression equation to predict ice cream sales on that day.

Answer

Temperature is likely to be **outside range** of **observed** values, so model not reliable.

Question 5

Interpret the value of a , (regression line)

Answer

$a = 562$ is the **number of chocolate bars sold per day** if **no money spent on advertising**

Question 6

Give a reason to support fitting a regression model of the form $y = a + bx$ to these data.

Answer

Points on the scatter graph lie **close to a straight line**.

Question 7 Give an interpretation of the value of b . *Answer* A flight costs **£2.03 (or about £2)** for every extra **100km** or about **2p** per **extra km**.

Discrete uniform distribution

Question 1

A discrete random variable is such that each of its values is assumed to be equally likely.

- (a) Write down the name of the distribution that could be used to model this random variable.
- (b) Give an example of such a distribution.
- (c) Comment on the assumption that each value is equally likely.
- (d) Suggest how you might refine the model in part (a).

Answer

- (a) Discrete uniform
 - (b) Tossing a fair die /coin and recording the score, drawing a card from a pack, and recording its value.
 - (c) Useful in theory – allows problems to be modelled, but the assumption might not be true in practice
 - (d) Carry out an experiment to find the probabilities – which might not fit the model.
-

Normal distribution

Question 1

The random variable X is normally distributed with mean 177.0, standard deviation 6.4.

It is suggested that X might be a suitable random variable to model the height, in cm, of adult males.

- (a) Give two reasons why this is a sensible suggestion.
- (b) Explain briefly why mathematical models can help to improve our understanding of real-world problems.

Answer

- (a) Male heights cluster round a central height of approx 177/178 cm
Height is a continuous random variable.
Most male heights lie within $177 \pm 3 \times 6.4$ (within 3 standard deviations of the mean).
- (b) Simplifies real world problems
Enable us to gain some understanding of real world problems more quickly/cheaply.

Question 2

Explain why the normal distribution may not be suitable to model the number of minutes that motorists are delayed by these roadworks.

Answer For this data skewness is 3.9 (given in the question), whereas a normal distribution is symmetrical and has no skew.

Question 3

Describe two features of the Normal distribution

Answer

Bell shaped curve; symmetrical about the mean; 95% of data lies within 2 s.d. of mean; etc. (any 2).

Question 4

Give a reason to support the use of a normal distribution in this case.

Answer

Since **mean and median are similar** (or **equal** or **very close**), the distribution is (nearly) **symmetrical** and a **normal distribution may be suitable**.

Allow mean or median close to mode/modal class \Rightarrow), the distribution is (nearly) **symmetrical** and a **normal distribution may be suitable**.

12 Appendix

$-1 \leq \text{P.M.C.C.} \leq 1$

Cauchy-Schwartz inequality

$$\begin{aligned} \text{Consider } & (a_1^2 + a_2^2)(b_1^2 + b_2^2) - (a_1b_1 + a_2b_2)^2 \\ &= a_1^2 b_1^2 + a_1^2 b_2^2 + a_2^2 b_1^2 + a_2^2 b_2^2 - a_1^2 b_1^2 - 2a_1b_1a_2b_2 - a_2^2 b_2^2 \\ &= a_1^2 b_2^2 - 2a_1b_1a_2b_2 + a_2^2 b_1^2 \\ &= (a_1b_2 - a_2b_1)^2 \geq 0 \\ \Rightarrow & (a_1^2 + a_2^2)(b_1^2 + b_2^2) - (a_1b_1 + a_2b_2)^2 \geq 0 \\ \Rightarrow & (a_1b_1 + a_2b_2)^2 \leq (a_1^2 + a_2^2)(b_1^2 + b_2^2) \end{aligned}$$

This proof can be generalised to show that

$$\begin{aligned} & (a_1b_1 + a_2b_2 + \dots + a_nb_n)^2 \leq (a_1^2 + a_2^2 + \dots + a_n^2)(b_1^2 + b_2^2 + \dots + b_n^2) \\ \text{or } & \left(\sum_{i=1}^n a_i b_i \right)^2 \leq \left(\sum_{i=1}^n a_i^2 \right) \left(\sum_{i=1}^n b_i^2 \right) \end{aligned}$$

P.M.C.C. between -1 and $+1$

In the above proof, take $a_i = (x_i - \bar{x})$, and $b_i = (y_i - \bar{y})$

$$\Rightarrow S_{xy} = \sum (x_i - \bar{x})(y_i - \bar{y}) = \sum a_i b_i$$

$$S_{xx} = \sum (x_i - \bar{x})^2 = \sum a_i^2 \quad \text{and} \quad S_{yy} = \sum (y_i - \bar{y})^2 = \sum b_i^2$$

$$\text{P.M.C.C.} = r = \frac{S_{xy}}{\sqrt{S_{xx}S_{yy}}}$$

$$\Rightarrow r^2 = \frac{S_{xy}^2}{S_{xx}S_{yy}} = \frac{(\sum a_i b_i)^2}{(\sum a_i^2)(\sum b_i^2)} \leq 1$$

using the Cauchy-Schwartz inequality

$$\Rightarrow -1 \leq r \leq +1$$

Regression line and coding

The regression line of y on x has equation $y = a + bx$, where $b = \frac{S_{xy}}{S_{xx}}$, and $a = \bar{y} - b\bar{x}$.

Using the coding $x = hX + m$, $y = gY + n$, the regression line for Y on X is found by writing $gY + n$ instead of y , and $hX + m$ instead of x in the equation of the regression line of y on x ,

$$\Rightarrow gY + n = a + b(hX + m)$$

$$\Leftrightarrow Y = \frac{a+bm-n}{g} + \frac{hb}{g} X \quad \dots \dots \dots \text{equation I.}$$

Proof

$$\bar{x} = h\bar{X} + m, \text{ and } \bar{y} = g\bar{Y} + n \quad \text{this is the effect of coding on the mean}$$

$$\Rightarrow (x - \bar{x}) = (hX + m) - (h\bar{X} + m) = (hX - h\bar{X}), \text{ and similarly } (y - \bar{y}) = (gY - g\bar{Y}).$$

Let the regression line of y on x be $y = a + bx$, and

let the regression line of Y on X be $Y = \alpha + \beta X$.

$$\text{Then } b = \frac{S_{xy}}{S_{xx}} \text{ and } a = \bar{y} - b\bar{x}$$

$$\text{also } \beta = \frac{S_{XY}}{S_{XX}} \text{ and } \alpha = \bar{Y} - \beta\bar{X}.$$

$$b = \frac{S_{xy}}{S_{xx}} = \frac{\sum(x-\bar{x})(y-\bar{y})}{\sum(x-\bar{x})^2} = \frac{\sum(hX-h\bar{X})(gY-g\bar{Y})}{\sum(hX-h\bar{X})^2} = \frac{hg \sum(X-\bar{X})(Y-\bar{Y})}{h^2 \sum(X-\bar{X})^2} = \frac{g \sum(X-\bar{X})(Y-\bar{Y})}{h \sum(X-\bar{X})^2}$$

$$\Rightarrow b = \frac{g}{h} \beta$$

$$\Rightarrow \beta = \frac{hb}{g}$$

$$\alpha = \bar{Y} - \beta\bar{X} = \frac{\bar{y}-n}{g} - \frac{hb}{g} \left(\frac{\bar{x}-m}{h} \right) = \frac{\bar{y}-b\bar{x}+bm-n}{g} = \frac{a+bm-n}{g} \quad \text{since } a = \bar{y} - b\bar{x}$$

and so

$$Y = \alpha + \beta X \Leftrightarrow Y = \frac{a+bm-n}{g} + \frac{hb}{g} X$$

which is the same as equation I.

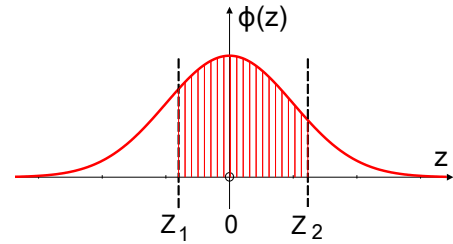
Normal Distribution, $Z = \frac{X - \mu}{\sigma}$

The standard normal distribution with mean 0 and standard deviation 1 has equation

$$\phi(z) = \frac{1}{\sqrt{2\pi}} e^{-\frac{1}{2}z^2}$$

The normal distribution tables allow us to find the area between Z_1 and Z_2 .

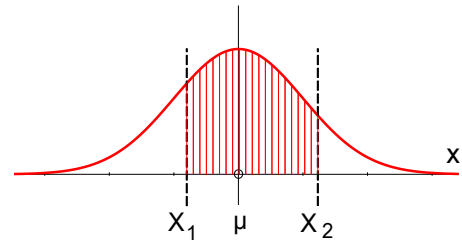
$$= \int_{Z_1}^{Z_2} \frac{1}{\sqrt{2\pi}} e^{-\frac{1}{2}z^2} dz$$



The normal distribution with mean μ and standard deviation σ has equation

$$f(x) = \frac{1}{\sigma\sqrt{2\pi}} e^{-\frac{1}{2}\left(\frac{x-\mu}{\sigma}\right)^2}$$

$$P(X_1 \leq X \leq X_2) = \int_{X_1}^{X_2} \frac{1}{\sigma\sqrt{2\pi}} e^{-\frac{1}{2}\left(\frac{x-\mu}{\sigma}\right)^2} dx$$



Using the substitution $z = \frac{x - \mu}{\sigma}$

$$dx = \sigma dz, \quad Z_1 = \frac{X_1 - \mu}{\sigma} \quad \text{and} \quad Z_2 = \frac{X_2 - \mu}{\sigma}$$

$$\Rightarrow P(X_1 \leq X \leq X_2) = \int_{Z_1}^{Z_2} \frac{1}{\sigma\sqrt{2\pi}} e^{-\frac{1}{2}z^2} \sigma dz$$

$$= \int_{Z_1}^{Z_2} \frac{1}{\sqrt{2\pi}} e^{-\frac{1}{2}z^2} dz$$

= the area under the standard normal curve, which we can find from the tables using

$$Z_1 = \frac{X_1 - \mu}{\sigma} \quad \text{and} \quad Z_2 = \frac{X_2 - \mu}{\sigma}.$$

$$\text{Thus } P(X_1 \leq X \leq X_2) = P(Z_1 \leq Z \leq Z_2) = \Phi(Z_2) - \Phi(Z_1)$$

Index

- accuracy, 48
- box plots, 19
- class boundaries, 7
- correlation, 30
 - context questions, 51
- cumulative frequency, 6
 - cumulative frequency curves, 8
- cumulative probability distribution, 40
- discrete uniform distribution, 42
 - context questions, 53
 - expected mean, 43
 - expected variance, 43
- exclusive events, 29
- expectation. *See* expected values
- expectation algebra, 40
- expected values, 40
 - expected mean, 40
 - expected variance, 20
 - interpretation, 41
- frequency distributions, 6
 - grouped frequency distributions, 7
- histograms, 9
 - context questions, 49
- independent events, 28
- interquartile range, 16, 22
- mean, 13
 - coding, 14
 - when to use, 15
- median
 - discrete lists and tables, 16
 - grouped frequency tables, 17
 - when to use, 15
- mode, 13
 - when to use, 15
- normal distribution, 44
 - context questions, 53
 - general normal distribution, 44
 - standard normal distribution, 44
 - standardising the variable, proof, 57
- outliers, 19
- percentiles, 19
- probability
 - diagrams for two dice etc, 26
 - number of arrangements, 29
 - rules, 25
 - Venn diagrams, 25
- probability distributions, 39
- product moment correlation coefficient, 30
 - between -1 and +1, 55
 - coding, 31
 - interpretation, 31
- quartiles
 - discrete lists and tables, 16
 - grouped frequency tables, 17
- random variables, 39
 - continuous, 39
 - discrete, 39
- range, 22
- regression, 33
 - context questions, 52
 - explanatory variable, 33
 - response variable, 33
- regression line, 33
 - interpretation, 34
- regression line and coding
 - proof, 56
- relative frequency, 25
- sample spaces, 25
- scatter diagrams, 30
 - line of best fit, 30
- skewness, 20
 - context questions, 50
- standard deviation, 22
- statistical modelling, 5
- statistical models
 - context questions, 48
- stem and leaf diagrams, 7
- tree diagrams, 27
- variables
 - continuous variables, 6
 - discrete variables, 6
 - qualitative variables, 6
 - quantitative variables, 6
- variance, 22
 - coding, 23