

MATHS

Boost your performance and confidence with these topic-based exam questions

Practice questions created by actual examiners and assessment experts

Detailed mark scheme

Suitable for all boards

Designed to test your ability and thoroughly prepare you

4. Statistics & Probability 4.4 Probability Distributions







IB Maths DP

4. Statistics & Probability

CONTENTS
4.1 Correlation & Regression
4.1.1 Bivariate Data
4.1.2 Correlation & Regression
4.2 Statistics Toolkit
4.2.1 Sampling & Data Collection
4.2.2 Statistical Measures
4.2.3 Frequency Tables
4.2.4 Linear Transformations of Data
4.2.5 Outliers
4.2.6Univariate Data
4.2.7 Interpreting Data
4.3 Probability
4.3.1 Probability & Types of Events
4.3.2 Conditional Probability
4.3.3 Bayes' Theorem
4.3.4 Sample Space Diagrams APERS PRACTICE
4.4 Probability Distributions
4.4.1 Discrete Probability Distributions
4.4.2 Mean & Variance
4.5 Binomial Distribution
4.5.1 The Binomial Distribution
4.5.2 Calculating Binomial Probabilities
4.6 Normal Distribution
4.6.1 The Normal Distribution
4.6.2 Calculations with Normal Distribution
4.6.3 Standardisation of Normal Variables
4.7 Further Probability Distributions
4.7.1 Probability Density Function



4.1 Correlation & Regression

4.1.1 Bivariate Data

Scatter Diagrams

What does bivariate data mean?

- **Bivariate data** is data which is collected on **two variables** and looks at how one of the factors affects the other
 - Each data value from one variable will be **paired** with a data value from the other variable
 - The two variables are often related, but do not have to be

What is a scatter diagram?

- A scatter diagram is a way of graphing bivariate data
 - One variable will be on the x-axis and the other will be on the y-axis
 - The variable that can be **controlled** in the data collection is known as the **independent** or **explanatory variable** and is plotted on the *x*-axis
 - The variable that is **measured** or discovered in the data collection is known as the **dependent** or **response variable** and is plotted on the y-axis
- Scatter diagrams can contain outliers that do not follow the trend of the data

Exam Tip

- If you use scatter diagrams in your Internal Assessment then be aware that finding outliers for bivariate data is different to finding outliers for univariate data
 - (x, y) could be an outlier for the bivariate data even if x and y are not outliers for their separate univariate data



Correlation

What is correlation?

- Correlation is how the two variables change in relation to each other
 - Correlation could be the result of a causal relationship but this is not always the case
- Linear correlation is when the changes are proportional to each other
- Perfect linear correlation means that the bivariate data will all lie on a straight line on a scatter diagram
- When describing correlation mention
 - The type of the correlation
 - **Positive correlation** is when an **increase** in one variable results in the other variable **increasing**
 - **Negative correlation** is when an **increase** in one variable results in the other variable **decreasing**
 - No linear correlation is when the data points don't appear to follow a trend
 - The strength of the correlation
 - Strong linear correlation is when the data points lie close to a straight line
 - Weak linear correlation is when the data points are not close to a straight line
- If there is strong linear correlation you can draw a line of best fit (by eye)
 - The line of best fit will pass through the mean point $(\overline{x}, \overline{y})$
 - If you are asked to draw a line of best fit
 - Plot the mean point
 - Draw a line going through it that follows the trend of the data



What is the difference between correlation and causation?



- It is important to be aware that just because correlation exists, it does not mean that the change in one of the variables is **causing** the change in the other variable
 - Correlation does not imply causation!
- If a change in one variable **causes** a change in the other then the two variables are said to have a **causal relationship**
 - Observing correlation between two variables does **not always** mean that there is a causal relationship
 - There could be **underlying factors** which is causing the correlation
 - Look at the two variables in question and consider the context of the question to decide if there could be a causal relationship
 - If the two variables are temperature and number of ice creams sold at a park then it is likely to be a causal relationship
 - Correlation may exist between global temperatures and the number of monkeys kept as pets in the UK but they are unlikely to have a causal relationship





Worked Example

A teacher is interested in the relationship between the number of hours her students spend on a phone per day and the number of hours they spend on a computer. She takes a sample of nine students and records the results in the table below.

Hours spent on a phone per day	7.6	7.0	8.9	3.0	3.0	7.5	2.1	1.3	5.8
Hours spent on a computer per day	1.7	1.1	0.7	5.8	5.2	1.7	6.9	7.1	3.3

a)

7

Draw a scatter diagram for the data.



b)

Describe the correlation.

Strong negative linear correlation

c) Draw a line of best fit.









Linear Regression

What is linear regression?

- If strong linear correlation exists on a scatter diagram then the data can be modelled by a linear model
 - Drawing lines of best fit by eye is not the best method as it can be difficult to judge the best position for the line
- The least squares regression line is the line of best fit that minimises the sum of the squares of the gap between the line and each data value
- It can be calculated by either looking at:
 - vertical distances between the line and the data values
 - This is the **regression line of y on x**
 - horizontal distances between the line and the data values
 - This is the **regression line of** x **on** y

How do I find the regression line of y on x?

- The **regression line of y on x** is written in the form y = ax + b
- a is the gradient of the line
 - It represents the change in y for each individual unit change in x
 - If a is **positive** this means y **increases** by a for a unit increase in x
 - If a is **negative** this means y **decreases** by |a| for a unit increase in x
- bis the y intercept
 - It shows the value of y when x is zero
- You are expected to use your GDC to find the equation of the regression line
 - Enter the bivariate data and choose the model "ax + b"
 - Remember the **mean point** $(\overline{x}, \overline{y})$ will lie on the regression line

How do I find the regression line of x on y?

- The **regression line of x on y** is written in the form x = cy + d
- c is the gradient of the line
 - \circ It represents the change in x for each individual unit change in y
 - If c is **positive** this means x **increases** by c for a unit increase in y
 - If c is **negative** this means x **decreases** by |c| for a unit increase in y
- disthex-intercept
 - \circ It shows the value of x when y is zero
- You are expected to use your GDC to find the equation of the regression line
 - It is found the same way as the regression line of y on x but with the two data sets **switched around**
 - Remember the **mean point** $(\overline{x}, \overline{y})$ will lie on the regression line

How do I use a regression line?

- The regression line can be used to decide what type of correlation there is if there is no scatter diagram
 - If the gradient is **positive** then the data set has **positive correlation**
 - If the gradient is **negative** then the data set has **negative correlation**



- The regression line can also be used to **predict** the value of a **dependent variable** from an **independent variable**
 - The equation for the y on x line should only be used to make predictions for y
 - Using a y on x line to predict x is not always reliable
 - \circ The equation for the x on y line should only be used to make predictions for x
 - Using an x on y line to predict y is not always reliable
 - Making a prediction within the range of the given data is called interpolation
 - This is usually reliable
 - The stronger the correlation the more reliable the prediction
 - Making a prediction outside of the range of the given data is called **extrapolation**
 - This is much less reliable
 - The prediction will be more reliable if the number of data values in the original sample set is bigger
- The y on x and x on y regression lines intersect at the mean point $(\overline{x}, \overline{y})$

Exam Tip

- Once you calculate the values of a and b store then in your GDC
 - This means you can use the full display values rather than the rounded values when using the linear regression equation to predict values

EXAM PAPERS PRACTICE

• This avoids rounding errors



Worked Example

The table below shows the scores of eight students for a maths test and an English test.

Maths (x)	7	18	37	52	61	68	75	82
English (y)	5	3	9	12	17	41	49	97

a)

Write down the value of Pearson's product-moment correlation coefficient, r.



b)

Write down the equation of the regression line of y on x, giving your answer in the form y = ax + b where a and b are constants to be found.



c)

Write down the equation of the regression line of x on y, giving your answer in the form x = cy + d where c and d are constants to be found.



d)

Use the appropriate regression line to predict the score on the maths test of a student who got a score of 63 on the English test.







PMCC

What is Pearson's product-moment correlation coefficient?

- Pearson's product-moment correlation coefficient (PMCC) is a way of giving a numerical value to a **linear relationship** of bivariate data
- The PMCC of a sample is denoted by the letter r
 - $r \operatorname{can} take any value such that <math>-1 \le r \le 1$
 - A positive value of r describes positive correlation
 - A negative value of r describes negative correlation
 - r = 0 means there is **no linear correlation**
 - r=1 means perfect positive linear correlation
 - r = -1 means **perfect negative linear** correlation
 - The closer to 1 or -1 the stronger the correlation



How do I calculate Pearson's product-moment correlation coefficient (PMCC)?

- You will be expected to use the statistics mode on your GDC to calculate the PMCC
- The formula can be useful to deepen your understanding

$$r = \frac{S_{xy}}{S_x S_y}$$

11 of 100



•
$$S_{xy} = \sum_{i=1}^{n} x_i y_i - \frac{1}{n} \left(\sum_{i=1}^{n} x_i \right) \left(\sum_{i=1}^{n} y_i \right)$$
 is linked to the **covariance**
• $S_x = \sqrt{\sum_{i=1}^{n} x_i^2 - \frac{1}{n} \left(\sum_{i=1}^{n} x_i \right)^2}$ and $S_y = \sqrt{\sum_{i=1}^{n} y_i^2 - \frac{1}{n} \left(\sum_{i=1}^{n} y_i \right)^2}$ are linked to the **variances**

• You **do not need to learn this** as using your GDC will be expected

When does the PMCC suggest there is a linear relationship?

- Critical values of r indicate when the PMCC would suggest there is a linear relationship
 - In your exam you will be given critical values where appropriate
 - Critical values will depend on the size of the sample
- If the **absolute value** of the **PMCC** is **bigger** than the **critical value** then this suggests a linear model is appropriate





Types of Data

What are the different types of data?

- Qualitative data is data that is usually given in words not numbers to describe something
 For example: the colour of a teacher's car
- Quantitative data is data that is given using numbers which counts or measures something
 - For example: the number of pets that a student has
- Discrete data is quantitative data that needs to be counted
 - Discrete data can only take **specific values** from a set of (usually finite) values
 - For example: the number of times a coin is flipped until a 'tails' is obtained
- Continuous data is quantitative data that needs to be measured
 - Continuous data can take **any value** within a range of infinite values
 - For example: the height of a student
- Age can be discrete or continuous depending on the context or how it is defined
 - If you mean how many years old a person is then this is discrete
 - If you mean how long a person has been alive then this is continuous

What is the difference between a population and a sample?

- The population refers to the whole set of things which you are interested in
 - For example: if a vet wanted to know how long a typical French bulldog slept for in a day then the population would be all the French bulldogs in the world
- A sample refers to a subset of the population which is used to collect data from
 - For example: the vet might take a sample of French bulldogs from different cities and record how long they sleep in a day
- A sampling frame is a list of all members of the population
 - For example: a list of employees' names within a company
- Using a sample instead of a population:
 - Is quicker and cheaper
 - Leads to less data needing to be analysed
 - Might not fully represent the population
 - Might introduce bias



Sampling Techniques

What is a random sample and a biased sample?

- A **random sample** is where every member of the population has an equal chance of being included in the sample
- A **biased sample** is one from which misleading conclusions could be drawn about the population
 - Random sampling is an attempt to minimise bias

What sampling techniques do I need to know?

Simple random sampling

- **Simple random sampling** is where every group of members from the population has an **equal probability** of being selected for the sample
- To carry this out you would...
 - uniquely number every member of a population
 - randomly select *n* different numbers using a random number generator or a form of lottery (where numbers are selected randomly)
- Effectiveness:
 - Useful when you have a small population or want a small sample (such as children in a class)
 - It can be time-consuming if the sample or population is large
 - This can not be used if it is not possible to number or list all the members of the population (such as fish in a lake)

Systematic sampling XAM PAPERS PRACTICE

- **Systematic sampling** is where a sample is formed by choosing members of a population at regular intervals using a list
- To carry this out you would...

size of sample (n)

- choose a random starting point between 1 and k
- select every *k*th member after the first one
- Effectiveness:
 - Useful when there is a natural order (such as a list of names or a conveyor belt of items)
 - Quick and easy to use
 - This can not be used if it is not possible to number or list all the members of the population (such as penguins in Antarctica)

Stratified sampling

• **Stratified sampling** is where the population is divided into disjoint groups and then a random sample is taken from each group



- The proportion of a group that is sampled is equal to the proportion of the population that belong to that group
- To carry this out you would...
 - Calculate the number of members sampled from each stratum

 - $\frac{\text{size of sample }(n)}{\text{size of population }(N)} \times \text{number of members in the group}$
 - Take a random sample from each group
- Effectiveness:
 - Useful when there are very different groups of members within a population
 - The sample will be representative of the population structure
 - The members selected from each stratum are chosen randomly
 - This can not be used if the population can not be split into groups or if the groups overlap

Quota sampling

- Quota sampling is where the population is split into groups (like stratified sampling) and members of the population are selected until each quota is filled
- To carry this out you would...
 - Calculate how many people you need from each group
 - Select members from each group until that guota is filled
 - The members do not have to be selected randomly
- Effectiveness:
 - Useful when collecting data by asking people who walk past you in a public place or when a sampling frame is not available 🧹 DDACTI
 - This can introduce bias as some members of the population might choose not to be included in the sample

Convenience sampling

- Convenience sampling is where a sample is formed using available members of the population who fit the criteria
- To carry this out you would...
 - Select members that are easiest to reach
- Effectiveness:
 - Useful when a list of the population is not possible
 - This is unlikely to be representative of the population structure
 - This is likely to produce biased results

What are the main criticisms of sampling techniques?

- Most sampling techniques can be improved by taking a larger sample
- Sampling can introduce bias so you want to minimise the bias within a sample • To minimise bias the sample should be as close to random as possible
- A sample only gives information about those members



• Different samples may lead to different conclusions about the population



Worked Example

Mike is a biologist studying mice in an open enclosure. He has access to approximately 540 field mice and 260 harvest mice. Mike wants to sample 10 mice and he wants the proportions of the two types of mice in his sample to reflect their respective proportions of the population.

a)

Calculate the number of field mice and harvest mice that Mike should include in his sample.





Reliability of Data

How can I decide if data is reliable?

- Data from a sample is reliable if similar results would be obtained from a different sample from the same population
- The sample should be **representative** of the population
- The sample should be **big enough**
 - Sampling a small proportion of a population is unlikely to be reliable

What can cause data to be unreliable?

- If the sample is **biased**
 - It is **not random**
- If errors are made when collecting data
 - Numbers could be recorded incorrectly, duplicated or missed out
- If the person collecting the data favours some members over others
 - They might seek out members who will lead to a desired outcome
 - They might exclude members if they would cause the sample to oppose the desired outcome
- If a significant proportion of **data is missing**
 - Some data may be unavailable
 - Some members might decide not to be part of the sample
 - This will mean the results are not necessarily representative of the population



4.2.2 Statistical Measures

Mean, Mode, Median

What are the mean, mode and median?

- Mean, median and mode are measures of central tendency
 They describe where the centre of the data is
 - They describe where the centre of the
- They are all types of **averages**
- In statistics it is important to be specific about which average you are referring to
- The **units** for the mean, mode and median are the **same** as the units for the data

How are the mean, mode, and median calculated for ungrouped data?

- The mode is the value that occurs most often in a data set
 - It is possible for there to be **more than one mode**
 - It is possible for there to be **no mode**
 - In this case **do not** say the mode is zero
- The median is the middle value when the data is in order of size
 - If there are two values in the middle then the median is the **midpoint** of the two values
- The mean is the sum of all the values divided by the number of values

$$\overline{x} = \prod_{i=1}^{n} \sum_{i=1}^{n} x_i$$

- Where $\sum_{i=1}^{n} x_i = x_1 + x_2 + \dots + x_n$ is the sum of the *n* pieces of data
- \circ The mean can be represented by the symbol μ
- Your **GDC** can calculate these statistical measures if you input the data using the statistics mode







Quartiles & Range

What are quartiles?

- Quartiles are measures of location
- Quartiles divide a population or data set into four equal sections
 - The lower quartile, Q_1 splits the lowest 25% from the highest 75%
 - \circ The **median, Q**₂ splits the lowest 50% from the highest 50%
 - \circ The **upper quartile**, **Q**₃ splits the lowest 75% from the highest 25%
- There are different methods for finding quartiles
 - $\circ~$ Values obtained by hand and using technology may differ
- You will be expected to use your GDC to calculate the quartiles

What are the range and interquartile range?

- The range and interquartile range are both measures of dispersion
 - They describe how spread out the data is
- The range is the largest value of the data minus the smallest value of the data
- The interquartile range is the range of the central 50% of data
 - It is the upper quartile minus the lower quartile

$$IQR = Q_3 - Q_1$$

- This is given in the formula booklet
- The units for the range and interquartile range are the same as the units for the data







Standard Deviation & Variance

What are the standard deviation and variance?

- The standard deviation and variance are both measures of dispersion
 They describe how spread out the data is in relation to the mean
- The variance is the mean of the squares of the differences between the values and the mean
 - $\circ~$ Variance is denoted σ^2
- The standard deviation is the square-root of the variance
 - $\circ~$ Standard deviation is denoted σ
- The units for the standard deviation are the same as the units for the data
- The units for the variance are the square of the units for the data

How are the standard deviation and variance calculated for ungrouped data?

• In the exam you will be expected to use the statistics function on your **GDC** to calculate the standard deviation and the variance

ΞE

• Calculating the standard deviation and the variance by hand may deepen your understanding

The formula for variance is
$$\sigma^2 = \sum_{i=1}^{k} f_i(x_i - \mu)^2$$

• This can be rewritten as
EXAM $P_{\sigma^2} = \sum_{i=1}^{k} f_i x_i^2$
 $n = \frac{1}{n} - \frac{1}{n} - \frac{1}{n}$

- The formula for **standard deviation** is $\sigma = \sqrt{-1}$
 - This can be rewritten as

$$\sigma = \sqrt{\frac{\sum_{i=1}^{k} f_i x_i^2}{n} - \mu^2}$$

n

• You **do not** need to learn these formulae as you will use your GDC to calculate these







Ungrouped Data

How are frequency tables used for ungrouped data?

- Frequency tables can be used for ungrouped data when you have lots of the same values within a data set
 - They can be used to collect and present data easily
- If the value 4 has a frequency of 3 this means that there are three 4's in the data set

How are measures of central tendency calculated from frequency tables with ungrouped data?

- The mode is the value that has the highest frequency
- The median is the middle value
 - Use cumulative frequencies (running totals) to find the median
- The mean can be calculated by
 - Multiplying each value x_i by its frequency f_i
 - Summing to get $\Sigma f_i x_i$
 - Dividing by the total frequency $n = \Sigma f_i$
 - This is given in the formula booklet



• Your **GDC** can calculate these statistical measures if you input the values and their frequencies using the statistics mode

How are measures of dispersion calculated from frequency tables with ungrouped data?

- The range is the largest value of the data minus the smallest value of the data
- The interquartile range is calculated by

$$IQR = Q_3 - Q_1$$

- The **quartiles** can be found by using your GDC and inputting the values and their frequencies
- The standard deviation and variance can be calculated by hand using the formulae
 - Variance

$$\sigma^2 = \frac{\sum_{i=1}^{k} f_i x_i^2}{n} - \mu^2$$

• Standard deviation



$$\sigma = \sqrt{\frac{\sum_{i=1}^{k} f_i x_i^2}{n} - \mu^2}$$

- You **do not need to learn** these formulae as you will be expected to use your GDC to find the standard deviation and variance
 - You may want to see these formulae to deepen your understanding

Exam Tip

- Always check whether your answers make sense when using your GDC
 - The value for a measure of central tendency should be within the range of data















Grouped Data

How are frequency tables used for grouped data?

- Frequency tables can be used for grouped data when you have lots of the same values within the same interval
 - Class intervals will be written using inequalities and without gaps
 - $10 \le x < 20$ and $20 \le x < 30$
 - If the class interval $10 \le x \le 20$ has a frequency of 3 this means there are three values in that interval
 - You do not know the **exact data values** when you are given grouped data

How are measures of central tendency calculated from frequency tables with grouped data?

- The modal class is the class that has the highest frequency
 - This is for equal class intervals only
- The median is the middle value
 - The exact value can not be calculated but it can be estimated by using a **cumulative** frequency graph
- The exact mean can not be calculated as you do not have the raw data
- The mean can be estimated by
 - Identifying the mid-interval value (midpoint) x_i for each class
 - Multiplying each value by the class frequency f_i
 - Summing to get $\Sigma f_i x_i$
 - Dividing by the total frequency $n = \Sigma f_i$
 - This is given in the formula booklet

EXAM PAPE
$$\sum_{i=1}^{k} f_{i} f_{i}^{X}$$
 PRACTICE
 $\overline{x} = \frac{1}{n}$

• Your **GDC** can estimate the mean if you input the mid-interval values and the class frequencies using the statistics mode

How are measures of dispersion calculated from frequency tables with grouped data?

- The exact range can not be calculated as the largest and smallest values are unknown
- The interquartile range can be estimated by

$$IQR = Q_3 - Q_1$$

- Estimates of the quartiles can be found by using a cumulative frequency graph
- The **standard deviation** and **variance** can be estimated using the mid-interval values x_i in the formulae
 - Variance



$$\sigma^2 = \frac{\sum_{i=1}^{k} f_i x_i^2}{n} - \mu^2$$

• Standard deviation

$$\sigma = \sqrt{\frac{\sum_{i=1}^{k} f_i x_i^2}{n} - \mu^2}$$

- You **do not need to learn** these formulae as you will be expected to use your GDC to estimate the standard deviation and variance using the mid-interval values
 - You may want to use these formulae to deepen your understanding

Exam Tip

- As you can only estimate statistical measures from a grouped frequency table it is good practice to indicate that the values are not exact
 - You can do this by rounding values rather than leaving as surds and fractions
 - $\overline{x} = 0.333$ (3sf) rather than \overline{x}





Worked Example

The table below shows the heights in cm of a group of 25 students.



Linear Transformations of Data

Why are linear transformations of data used?

- Sometimes data might be very large or very small
- You can apply a linear transformation to the data to make the values more manageable
 - You may have heard this referred to as:
 - Effects of constant changes
 - Linear coding
- Linear transformations of data can affect the statistical measures

How is the mean affected by a linear transformation of data?

- Let \overline{X} be the **mean** of some data
- If you multiply each value by a constant k then you will need to multiply the mean by k • Mean is $k\overline{x}$
- If you add or subtract a constant a from all the values then you will need to add or subtract the constant a to the mean
 - Mean is $\overline{x} \pm a$

How is the variance and standard deviation affected by a linear transformation of data?

- Let σ² be the variance of some data
 σ is the standard deviation
- If you multiply each value by a constant k then you will need to multiply the variance by k²
 - Variance is $k^2 \sigma^2$
 - You will need to multiply the standard deviation by the absolute value of k
 - Standard deviation is $|k|\sigma$
 - If you add or subtract a constant a from all the values then the variance and the standard deviation stay the same
 - Variance is σ^2
 - Standard deviation is σ

Exam Tip

- If you forget these results in an exam then you can look in the HL section of the formula booklet to see them written in a more algebraic way
 - Linear transformation of a single variable



• where E(...) means the mean and Var(...) means the variance







Outliers

What are outliers?

- Outliers are extreme data values that do not fit with the rest of the data
 - They are either a lot bigger or a lot smaller than the rest of the data
- Outliers are defined as values that are more than 1.5 x IQR from the nearest quartile
 - *x* is an outlier if *x* < *Q*₁ 1.5 × IQR or *x* > *Q*₃ + 1.5 × IQR
- Outliers can have a big effect on some statistical measures

Should I remove outliers?

- The decision to remove outliers will depend on the context
- Outliers **should be removed** if they are found to be **errors**
 - The data may have been recorded incorrectly
 - For example: The number 17 may have been recorded as 71 by mistake
- Outliers should not be removed if they are a valid part of the sample
 - The data may need to be checked to verify that it is not an error
 - For example: The annual salaries of employees of a business might appear to have an outlier but this could be the director's salary







Box Plots

Univariate data is data that is in **one variable**.

What is a box plot (box and whisker diagram)?

- A box plot is a graph that clearly shows key statistics from a data set
 - It shows the median, quartiles, minimum and maximum values and outliers
 - It does not show any other individual data items
- The middle 50% of the data will be represented by the box section of the graph and the lower and upper 25% of the data will be represented by each of the whiskers
- Any outliers are represented with a cross on the outside of the whiskers
 - If there is an outlier then the whisker will end at the value before the outlier
- Only one axis is used when graphing a box plot
- It is still important to make sure the axis has a clear, even scale and is labelled with units



What are box plots useful for?

- Box plots can clearly show the shape of the distribution
 - If a box plot is symmetrical about the median then the data could be **normally distributed**
- Box plots are often used for **comparing two sets of data**
 - Two box plots will be drawn next to each other using the same axis
 - They are useful for **comparing data** because it is easy to see the main shape of the distribution of the data from a box plot
 - You can easily compare the medians and interquartile ranges

Exam Tip

- In an exam you can use your GDC to draw a box plot if you have the raw data
 - You calculator's box plot can also include outliers so this is a good way to check






Cumulative Frequency Graphs

What is cumulative frequency?

- The cumulative frequency of x is the running total of the frequencies for the values that are less than or equal to x
- For grouped data you use the upper boundary of a class interval to find the cumulative frequency of that class

What is a cumulative frequency graph?

- A cumulative frequency graph is used with data that has been organised into a **grouped frequency** table
- Some coordinates are plotted
 - The x-coordinates are the upper boundaries of the class intervals
 - The y-coordinates are the **cumulative frequencies** of that class interval
- The coordinates are then joined together by hand using a **smooth increasing curve**

What are cumulative frequency graphs useful for?

- They can be used to **estimate** statistical measures
 - Draw a horizontal line from the y-axis to the curve
 - For the median: draw the line at 50% of the total frequency
 - For the lower quartile: draw the line at 25% of the total frequency
 - For the upper quartile: draw the line at 75% of the total frequency
 - For the pth percentile: draw the line at p% of the total frequency
 - Draw a **vertical line** down from the curve to the x-axis
 - This **x-value** is the relevant statistical measure
- They can used to estimate the number of values that are bigger/small than a given value
 - Draw a vertical line from the given value on the x-axis to the curve
 - Draw a **horizontal line** from the curve to the *y*-axis
 - This value is an estimate for how many values are less than or equal to the given value
 - To estimate the number that is greater than the value subtract this number from the total frequency
 - They can be used to estimate the interquartile range IQR = $Q_3 Q_1$
 - They can be used to construct a **box plot** for grouped data













Histograms

What is a (frequency) histogram?

- A frequency histogram clearly shows the frequency of class intervals
 - The classes will have **equal class intervals**
 - The **frequency** will be on the *y*-axis
 - The bar for a class interval will begin at the lower boundary and end at the upper boundary
- A frequency histogram is **similar to a bar chart**
 - A bar chart is used for qualitative or discrete data and has gaps between the bars
 - A frequency histogram is used for continuous data and has no gaps between bars

What are (frequency) histograms useful for?

- They show the **modal class** clearly
- They show the shape of the distribution
 - It is important the class intervals are of equal width
- They can show whether the variable can be modelled by a normal distribution
 - If the shape is symmetrical and bell-shaped





Worked Example

The table below and its corresponding histogram show the mass, in kg, of some new born bottlenose dolphins.

Mass, <i>m</i> kg	Frequency
$4 \le m < 8$	4
$8 \le m < 12$	15
$12 \le m < 16$	19
$16 \le m < 20$	10
$20 \le m < 24$	6

a)

Draw a frequency histogram to represent the data.





4.2.7 Interpreting Data

Interpreting Data

How do l interpret statistical measures?

- The mode is useful for qualitative data
 - It is not as useful for quantitative data as there is not always a unique mode
- The mean includes all values
 - $\circ \ \ \mathsf{It} \, \mathsf{is} \, \mathsf{affected} \, \mathsf{by} \, \mathsf{outliers}$
 - A smaller/larger mean is preferable depending on the scenario
 - A smaller mean time for completing a puzzle is better
 - A bigger mean score on a test is better
- The median is not affected by outliers
 - It does not use all the values
- The range gives the full spread of the all of the data
 - It is affected by outliers
- The **interquartile range gives the spread of the middle 50%** about the median and is not affected by outliers
 - It does not use all the values
 - A bigger IQR means the data is more spread out about the median
 - A smaller IQR means the data is more centred about the median
- The standard deviation and variance use all the values to give a measure of the average spread of the data about the mean
 - They are affected by outliers
 - A bigger standard deviation means the data is more spread out about the mean
 - $\circ~$ A smaller standard deviation means the data is more centred about the mean

How do I choose which diagram to use to represent data?

- Box plots
 - Can be used with ungrouped univariate data
 - Shows the range, interquartile range and quartiles clearly
 - Very useful for comparing data patterns quickly
- Cumulative frequency graphs
 - Can be used with continuous grouped univariate data
 - Shows the running total of the frequencies that fall below the upper bound of each class
- Histograms
 - $\circ~$ Can be used with continuous grouped univariate data
 - Used with equal class intervals
 - Shows the frequencies of the group
- Scatter diagrams
 - Can be used with ungrouped **bivariate** data
 - Shows the graphical relationship between the variables

How do I compare two or more data sets?

- Compare a measure of central tendency
 - If the data contains outliers use the median
 - If the data is roughly symmetrical use the mean



- Compare a measure of dispersion
 - $\circ~$ If the data contains outliers use the interquartile range
 - If the data is roughly symmetrical use the standard deviation
- Consider whether it is better to have a smaller or bigger average
 - This will depend on the context
 - A smaller average time for completing a puzzle is better
 - A bigger average score on a test is better
- Consider whether it is better to have a smaller or bigger spread
 Usually a smaller spread means it is more consistent
- Always relate the **comparisons to the context** and consider reasons
 - Consider the sampling technique and the data collection method

Worked Example The box plots below show the waiting times for the two doctor surgeries, HealthHut and FitFirst. HealthHut FitFirst ō 10 20 40 50 60 Waiting time (minutes) Compare the two distributions of waiting times in context. EMPAREM PAPERS PRACTICE · a measure of central tendency · a measure of dispersion Health Hut's median waiting time is smaller than Fit First's (20 < 24). On average patients get seen quicker at HealthHut. Fit First's interquartile range is smaller than Health Hut's (13 < 19). There is less variability of waiting times at FitFirst



4.3 Probability

4.3.1 Probability & Types of Events

Probability Basics

What key words and terminology are used with probability?

- An experiment is a repeatable activity that has a result that can be observed or recorded
 - Trials are what we call the repeats of the experiment
- An outcome is a possible result of a trial
- An **event** is an outcome or a collection of outcomes
 - Events are usually denoted with capital letters: A, B, etc
 - \circ n(A) is the number of outcomes that are included in event A
 - An event can have one or more than one outcome
- A sample space is the set of all possible outcomes of an experiment
 - \circ This is denoted by U
 - n(U) is the total number of outcomes
 - It can be represented as a **list** or a **table**

How do I calculate basic probabilities?

- If all outcomes are equally likely then probability for each outcome is the same
 - Probability for each outcome is $\frac{1}{n(U)}$
- **Theoretical probability** of an event can be calculated without using an experiment by dividing the number of outcomes of that event by the total number of outcomes

EXAM PAP(A) =
$$n(A)$$
 RACTICE

- This is given in the **formula booklet**
- Identifying all possible outcomes either as a list or a table can help
- Experimental probability (also known as relative frequency) of an outcome can be calculated using results from an experiment by dividing its frequency by the number of trials
 - **Relative frequency** of an outcome is

Total number of trials (*n*)

How do I calculate the expected number of occurrences of an outcome?

- Theoretical probability can be used to calculate the expected number of occurrences of an outcome from *n* trials
- If the probability of an outcome is p and there are n trials then:
 - The expected number of occurrences is **np**
 - This does not mean that there will exactly np occurrences
 - If the experiment is repeated multiple times then we expect the number of occurrences to average out to be *np*

What is the complement of an event?

- The probabilities of all the outcomes add up to 1
- Complementary events are when there are **two events** and **exactly one** of them will occur One event has to occur but both events can not occur at the same time
- The complement of event A is the event where event A does not happen



- This can be thought of as **not A**
- This is denoted A'

$$P(A) + P(A') = 1$$

This is in the **formula booklet**

It is commonly written as P(A') = 1 - P(A)

What are different types of combined events?

- The intersection of two events (A and B) is the event where both A and B occur
 - This can be thought of as **A and B**
 - \circ This is denoted as $A \cap B$
- The union of two events (A and B) is the event where A or B or both occur
 - This can be thought of as **A or B**
 - \circ This is denoted $A \cup B$
- The event where A occurs given that event B has occurred is called conditional probability
 - This can be thought as **A given B**
 - \circ This is denoted $A \mid B$

How do I find the probability of combined events?

• The probability of A or B (or both) occurring can be found using the formula

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

This is given in the formula booklet

You subtract the probability of A and B both occurring because it has been

included twice (once in P(A) and once in P(B))

• The probability of A and B occurring can be found using the formula

$$P(A \cap B) = P(A)P(B|A)$$

A rearranged version is given in the **formula booklet** Basically you multiply the probability of *A* by the probability of *B* then happening

Exam Tip

• In an exam drawing a Venn diagram or tree diagram can help even if the question does not ask you to



Worked Example

Dave has two fair spinners, A and B. Spinner A has three sides numbered 1, 4, 9 and spinner B has four sides numbered 2, 3, 5, 7. Dave spins both spinners and forms a two-digit number by using the spinner A for the first digit and spinner B for the second digit.

T is the event that the two-digit number is a multiple of 3.

a)

 $\label{eq:listallthepossible} two-digit\,numbers.$





Independent & Mutually Exclusive Events

What are mutually exclusive events?

- Two events are mutually exclusive if they cannot both occur
 - For example: when rolling a dice the events "getting a prime number" and "getting a 6" are mutually exclusive
- If A and B are mutually exclusive events then:
 - $\circ P(A \cap B) = 0$

What are independent events?

- Two events are independent if one occurring does not affect the probability of the other occurring
 - For example: when flipping a coin twice the events "getting a tails on the first flip" and "getting a tails on the second flip" are independent
- If A and B are independent events then:

• P(A|B) = P(A) and P(B|A) = P(B)

- If A and B are independent events then:
 - $\circ P(A \cap B) = P(A)P(B)$
 - This is given in the **formula booklet**

This is a useful formula to test whether two events are statistically independent

How do I find the probability of combined mutually exclusive events?

If A and B are mutually exclusive events then

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

This is given in the formula booklet PRACTICE

This occurs because $P(A \cap B) = 0$

• For any two events A and B the events $A \cap B$ and $A \cap B'$ are **mutually exclusive** and A is the **union** of these two events

• $P(A) = P(A \cap B) + P(A \cap B')$

This works for any two events A and B





a)

Worked Example

A student is chosen at random from a class. The probability that they have a dog is 0.8, the probability they have a cat is 0.6 and the probability that they have a cat or a dog is 0.9.

Find the probability that the student has both a dog and a cat.

Let D be event "has a dog" and C be "has a cat" $P(D \cup C) = P(D) + P(C) - P(D \cap C)$ $0.9 = 0.8 + 0.6 - P(D \cap C)$ $P(D \cap C) = 0.5$

b)

Two events, Q and R, are such that P(Q) = 0.8 and $P(Q \cap R) = 0.1$. Given that Q and R are independent, find P(R).



c)

Two events, S and T, are such that P(S) = 2P(T). ACTICE Given that S and T are mutually exclusive and that $P(S \cup T) = 0.6$ find P(S) and P(T).

S and T mutually exclusive
$$\Rightarrow P(S \cup T) = P(S) + P(T)$$

 $0.6 = P(S) + P(T)$
 $0.6 = 2P(T)^{+} + P(T)$
 $0.6 = 3P(T)$
 $P(S) = 2P(T)$
 $P(T) = 0.2$ and $P(S) = 0.4$



4.3.2 Conditional Probability

Conditional Probability

What is conditional probability?

- **Conditional probability** is where the probability of an **event** happening can vary depending on the outcome of a prior event
- The event A happening given that event B has happened is denoted A|B
- A common example of conditional probability involves selecting multiple objects from a bag without replacement
 - The probability of selecting a certain item changes depending on what was selected before

This is because the total number of items will change as they are not replaced once they have been selected

How do I calculate conditional probabilities?

- Some conditional probabilities can be calculated by using counting outcomes
 - Probabilities without replacement can be calculated like this
 - For example: There are 10 balls in a bag, 6 of them are red, two of them are selected without replacement

To find the probability that the second ball selected is red given that the first one is red count how many balls are left:

A red one has already been selected so there are 9 balls left and 5 are red so the

- probability is g
- You can use sample space diagrams to find the probability of A given B:
 - reduce your sample space to just include outcomes for event B
 - find the proportion that also contains outcomes for event A
- There is a formula for conditional probability that you should use

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

- This is given in the formula booklet
- This can be rearranged to give $P(A \cap B) = P(B)P(A \mid B)$
- By symmetry you can also write $P(A \cap B) = P(A)P(B|A)$

How do I tell if two events are independent using conditional probabilities?

- If A and B are two events then they are independent if:
 - $\circ P(A | B) = P(A) = P(A | B')$
- Equally you can still use $P(A \cap B) = P(A)P(B)$ to test for independence
 - This is given in the **formula booklet**







Bayes' Theorem

What is Bayes' Theorem

- Bayes' Theorem allows you switch the order of conditional probabilities
 - $\circ~$ If you know $\mathrm{P}(B)$, $\mathrm{P}(B')$ and $\mathrm{P}(A \,|\, B)$ then Bayes' Theorem allows you to find $\mathrm{P}(B | A)$
- Essentially if you have a **tree diagram** you will already know the conditional probabilities of the **second branches**
 - Bayes' Theorem allows you to find the **conditional probabilities** if you **switch the order** of the events
- For any two events A and B Bayes' Theorem states:

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

- This is given in the **formula booklet**
- This formula is derived using the formulae:

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

$$P(A) = P(B \cap A) + P(B' \cap A)$$

$$P(B \cap A) = P(B)P(A | B) \text{ and } P(B' \cap A) = P(B')P(A | B')$$

- Bayes' Theorem can be **extended** to **mutually exclusive events** $B_1, B_2, ..., B_n$ and any other event A
 - In your exam you will have a **maximum of three** mutually exclusive events

$$P(B_{i}|A) = \frac{P(B_{i})P(A|B_{i})}{P(B_{1})P(A|B_{1}) + P(B_{2})P(A|B_{2}) + P(B_{3})P(A|B_{3})}$$

This is given in the **formula booklet**

How do I calculate conditional probabilities using Bayes' Theorem?

- Start by drawing a tree diagram
 - Label **B**₁ & **B**₂ (& **B**₃ if necessary) on the **first set** of branches
 - Label **A & A'** on the **second set** of branches
- The questions will give you enough information to label the probabilities on this tree
- Identify the probabilities needed to use Bayes' Theorem
 - The probabilities will come in pairs: $P(B_i)$ and $P(A|B_i)$







SWAP B_i for $\mathsf{B}_1\,,\mathsf{B}_2$ or B_3 DEPENDING ON THE PROBABILITY NEEDED

THE DENOMINATOR DOES NOT CHANGE



• In an exam you are less likely to make a mistake when using the formula if you draw a tree diagram first



Worked Example

Lucy is doing a quiz. For each question there's a 45% chance that it is about music, 30% chance that it is about TV and 25% chance that it is about literature. The probability that Lucy answers a question correctly is 0.6 for music, 0.95 for TV and 0.4 for literature.



Draw a tree diagram to represent this information.





Venn Diagrams

What is a Venn diagram?

- A Venn diagram is a way to illustrate **events** from an **experiment** and are particularly useful when there is an overlap between possible **outcomes**
- A Venn diagram consists of
 - a **rectangle** representing the **sample space (U)**
 - The rectangle is labelled U
 - Some mathematicians instead use S or ξ
 - a circle for each event
 - Circles may or may not overlap depending on which **outcomes** are shared between **events**
- The numbers in the circles represent either the **frequency** of that event or the **probability** of that event
 - If the **frequencies** are used then they should **add up to the total frequency**
 - If the probabilities are used then they should add up to 1

What do the different regions mean on a Venn diagram?

- A' is represented by the regions that are **not in** the A circle
- $A \cap B$ is represented by the region where the A and B circles **overlap**
- $A \cup B$ is represented by the regions that are in A or B or both
- Venn diagrams show 'AND' and 'OR' statements easily
- Venn diagrams also instantly show **mutually exclusive** events as these circles will **not overlap**
- Independent events can not be instantly seen
 - You need to use probabilities to deduce if two events are independent







How do I solve probability problems involving Venn diagrams?

- Draw, or add to a given Venn diagram, filling in as many values as possible from the information provided in the question
- It is usually helpful to work from the centre outwards
 - Fill in intersections (overlaps) first
- If two events are independent you can use the formula

$$\circ P(A \cap B) = P(A)P(B)$$

- To find the conditional probability P(A | B)
 - Add together the frequencies/probabilities in the *B* circle This is your denominator
 - Out of those frequencies/probabilities add together the ones that are also in the A circle
 - This is your numerator
 - Evaluate the fraction





Exam Tip

 \bigcirc

- If you struggle to fill in a Venn diagram in an exam:
 - Label the missing parts using algebra
 - Form equations using known facts such as: the sum of the probabilities should be 1
 P(A∩B)=P(A)P(B) if A and B are independent events





Worked Example

40 people are asked if they have sugar and/or milk in their coffee. 21 people have sugar, 25 people have milk and 7 people have neither.

a)

Draw a Venn diagram to represent the information.



b)

One of the 40 people are randomly selected, find the probability that they have sugar but not milk with their coffee.



c)

Given that a person who has sugar is selected at random, find the probability that they have milk with their coffee.

Given that sugar has been selected we only want the S circle as our total. Out of the S circle 13 also have milk $P(M|S) = \frac{13}{21}$



Tree Diagrams

What is a tree diagram?

- A tree diagram is another way to show the outcomes of combined events
 They are very useful for intersections of events
- The events on the branches must be **mutually exclusive**
 - $\circ~$ Usually they are an event and its complement
- The probabilities on the second sets of branches **can depend** on the outcome of the first event
 - These are conditional probabilities
- When selecting the items from a bag:
 - The second set of branches will be the **same** as the first if the items **are replaced**
 - The second set of branches will be the **different** to the first if the items **are not replaced**

How are probabilities calculated using a tree diagram?

- To find the probability that two events happen together you **multiply** the corresponding probabilities on their branches
 - It is helpful to find the probability of all combined outcomes once you have drawn the tree
- To find the probability of an event you can:
 - **add together** the probabilities of the **combined outcomes** that are part of that event For example: $P(A \cup B) = P(A \cap B) + P(A \cap B') + P(A' \cap B)$
 - **subtract** the probabilities of the combined outcomes that are not part of that event from 1

For example: $P(A \cup B) = 1 - P(A' \cap B')$



Do I have to use a tree diagram?

- If there are **multiple events** or trials then a tree diagram can get big
- You can break down the problem by using the words **AND/OR/NOT** to help you find probabilities without a tree



• You can speed up the process by only drawing parts of the tree that you are interested in

Which events do I put on the first branch?

- If the events A and B are independent then the order does not matter
- If the events A and B are **not independent** then the **order does matter**
 - If you have the probability of **A given B** then put **B on the first set** of branches
 - $\circ~$ If you have the probability of ${\pmb B}$ given ${\pmb A}$ then put ${\pmb A}$ on the first set of branches

Exam Tip

- In an exam do not waste time drawing a full tree diagram for scenarios with lots of events unless the question asks you to
 - Only draw the parts that you are interested in







Worked Example

20% of people in a company wear glasses. 40% of people in the company who wear glasses are right-handed. 50% of people in the company who don't wear glasses are right-handed.

a)

Draw a tree diagram to represent the information.



b)

One of the people in the company are randomly selected, find the probability that they are right-handed.

c)

Given that a person who is right-handed is selected at random, find the probability that they wear glasses.

$$P(\alpha|R) = \frac{P(\alpha R)}{P(R)} = \frac{0.08}{0.48}$$
$$P(\alpha|R) = \frac{1}{6}$$



4.4 Probability Distributions

4.4.1 Discrete Probability Distributions

Discrete Probability Distributions

What is a discrete random variable?

- A random variable is a variable whose value depends on the outcome of a random event
 - The value of the random variable is not known until the event is carried out (this is what is meant by 'random' in this case)
- Random variables are denoted using upper case letters (X, Y, etc)
- Particular outcomes of the event are denoted using lower case letters (x, y, etc)
- P(X=x) means "the probability of the random variable X taking the value x"
- A **discrete** random variable (often abbreviated to DRV) can only take **certain values** within a set
 - Discrete random variables usually count something
 - Discrete random variables usually can only take a finite number of values but it is possible that it can take an infinite number of values (see the examples below)
- Examples of discrete random variables include:
 - The number of times a coin lands on heads when flipped 20 times this has a finite number of outcomes: {0,1,2,...,20}
 - The number of emails a manager receives within an hour this has an infinite number of outcomes: {1,2,3,...}
 - The number of times a dice is rolled until it lands on a 6 CTT C E this has an infinite number of outcomes: {1,2,3,...}
 - The number that a dice lands on when rolled once this has a finite number of outcomes: {1,2,3,4,5,6}

What is a probability distribution of a discrete random variable?

- A discrete probability distribution fully describes all the values that a discrete random variable can take along with their associated probabilities
 - This can be given in a **table**
 - Or it can be given as a **function** (called a discrete probability distribution function or "pdf")
 - They can be represented by **vertical line graphs** (the possible values for along the horizontal axis and the probability on the vertical axis)
- The sum of the probabilities of all the values of a discrete random variable is 1
 - This is usually written $\sum P(X=x) = 1$
- A **discrete uniform distribution** is one where the random variable takes a finite number of values each with an **equal probability**
 - If there are n values then the probability of each one is





How do I calculate probabilities using a discrete probability distribution?

- First draw a table to represent the probability distribution
 - If it is given as a function then find each probability
 - If any probabilities are unknown then use algebra to represent them
- Form an equation using $\sum P(X=x) = 1$
 - Add together all the probabilities and make the sum equal to 1
- To find P(X=k) EXAM PAPERS PRACTICE
 - If k is a possible value of the random variable X then P(X = k) will be given in the table
 If k is not a possible value then P(X = k) = 0
- To find $P(X \le k)$
 - Identify all possible values, x_i , that X can take which satisfy $x_i \le k$
 - Add together all their corresponding probabilities
 - $\circ P(X \le k) = \sum_{x_i \le k} P(X = x_i)$
 - Some mathematicians use the notation F(x) to represent the cumulative distribution $F(x) = P(X \le x)$
- Using a similar method you can find P(X < k), P(X > k) and $P(X \ge k)$
- As all the probabilities add up to 1 you can form the following equivalent equations:

• P(X < k) + P(X = k) + P(X > k) = 1

$$\circ P(X > k) = 1 - P(X \le k)$$

$$\circ P(X \ge k) = 1 - P(X < k)$$

How do I know which inequality to use?

- $P(X \le k)$ would be used for phrases such as:
 - At most, no greater than, etc



- P(X < k) would be used for phrases such as:
 Fewer than
- P(X≥k) would be used for phrases such as:
 At least, no fewer than, etc
- P(X > k) would be used for phrases such as:
 - Greater than, etc





4.4.2 Mean & Variance

Expected Values E(X)

What does E(X) mean and how do I calculate E(X)?

- E(X) means the expected value or the mean of a random variable X
 - $\circ~$ The expected value does not need to be an obtainable value of X
 - For example: the expected value number of times a coin will land on tails when flipped 5 times is 2.5
- For a **discrete** random variable, it is calculated by:
 - Multiplying each value of X with its corresponding probability
 - Adding all these terms together

$$E(X) = \sum x P(X = x)$$

This is given in the **formula booklet**

- Look out for **symmetrical** distributions (where the values of X are symmetrical and their probabilities are symmetrical) as the mean of these is the same as the median
 - For example: if X can take the values 1, 5, 9 with probabilities 0.3, 0.4, 0.3 respectively then by symmetry the mean would be 5

How can I decide if a game is fair?

- Let X be the random variable that represents the gain/loss of a player in a game
 X will be negative if there is a loss
- Normally the expected gain or loss is calculated by **subtracting** the **cost to play** the game from the **expected value** of the **prize**
- If E(X) is **positive** then it means the player can **expect to make a gain**
- If E(X) is **negative** then it means the player can **expect to make a loss**
- The game is called fair if the expected gain is 0
 - $\circ E(X) = 0$





?

Daphne pays \$5 to play a game where she wins a prize of \$1, \$5, \$10 or \$100. The random variable W represents the amount she wins and has the probability distribution shown in the following table:

	W	1	5	10	100	
	P(W = w)	0.35	0.5	0.05	0.01	
a)						
Calculate the	expected value c	of Daphr	ne's priz	e.		
	Formula booklet	Expected valu discrete rando variable X	e of a m E	$(X) = \sum x \mathbf{P}(X =$	x)	
	$E(W) = \sum w P(W)$	= w)				
	= × 0.35 +	5 × 0 5	+ 10 × 0.0	05 + 100×	0.01	
[Expected value	= \$4.35				
b)	athortho gorooid	foir		1		
Determine wi				h .	٥	
	A game is tair	is expec	ted gair	n/loss is	U	
	Prize - cost					
	4.35 - 5 = -	0.65				
E	XAM PA	Pko	25	RA	CTICE	
	Expected loss is	5 4 0.0	5 50	game is	not tair	



Variance Var(X)

What does Var(X) mean and how do I calculate Var(X)?

- Var(X) means the variance of a random variable X
 - The **standard deviation** is the **square root** of the variance
 - This provides a **measure of the spread** of the outcomes of *X*
 - The variance and standard deviation can **never be negative**
- The variance of X is the **mean of the squared difference** between X and the mean

$$\operatorname{Var}(X) = \operatorname{E}(X - \mu)^2$$

- This is given in the formula booklet
- This formula can be rearranged into the more useful form:

$$Var(X) = E(X^2) - [E(X)]^2$$

- This is given in the **formula booklet**
 - Compare this formula to the formula for the variance of a set of data
- This formula works for both **discrete** and **continuous** X

How do I calculate E(X²) for discrete X?

- E(X²) means the expected value or the mean of the random variable defined as X²
- For a **discrete** random variable, it is calculated by:
 - Squaring each value of X to get the values of X^2
 - Multiplying each value of X² which its corresponding probability
 - Adding all these terms together

 $E(X^{2}) = \sum x^{2}P(X = x)$ This is given in the formula booklet as part of the formula for Var(X) $Var(X) = \sum x^{2}P(X = x) - \mu^{2}$

• **E(f(X))** can be found in a similar way

Is E(X²) equal to E(X)²?

- Definitely not!
 - They are only equal if X can only take one value
- E(X²) is the mean of the values of X²
- E(X)² is the square of the mean of the values of X
- To see the difference
 - $\circ~$ Imagine a random variable X that can only take 1 and -1 with equal chance
 - E(X) = 0 so **E(X)² = 0**
 - The square values are 1 and 1 so **E(X²) = 1**



Exam Tip

 \bigcirc

2

- In an exam you can enter the probability distribution into your GDC using the statistics mode
 - Enter the possible values as the data
 - Enter the probabilities as the frequencies
- You can then calculate the mean and variance just like you would with data

Worked Example

The score on a game is represented by the random variable S defined below.

		S	0	1	2	10				
		P(S=s)	0.4	0.3	0.25	0.05				
Calcula	te Var(S).									
	Calculate	E(5)								
	Formula booklet Expected value of a discrete random variable X $E(X) = \sum_{x} P(X = x)$									
	$E(s) = \sum_{s} P(s=s) = 0 \times 0.4 + 1 \times 0.3 + 2 \times 0.25 + 10 \times 0.05 = 1.3$									
	Calculate E(S ²)									
	$E(S^{*}) = \sum S^{*}P(S = s) = 0^{2} \times 0.4 + 1^{2} \times 0.3 + 2^{2} \times 0.25 + 10^{2} \times 0.05 = 6.3$									
	Calculate	Var(S)								
	Formula b	ooklet Variance	e	Va	r(X) = E(X -	$(-\mu)^{2} = \mathbb{E}(X^{2}) - [\mathbb{E}(X)]^{2}$				
	$V_{ar}(5) = E(5^2) - [E(5)]^2 = 6.3 - 1.3^2$									
	Var(5) =	4.61								



Transformation of a Single Variable

How do I calculate the expected value and variance of a transformation of X?

- Suppose X is **transformed** by the function f to form a new variable T = f(X)
 - This means the function f is applied to all possible values of X
- Create a **new probability distribution table**
 - The top row contains the values $t_i = f(x_i)$
 - The bottom row still contains the values $P(X = x_i)$ which are unchanged as:

$$P(X = x_i) = P(f(X) = f(x_i)) = P(T = t_i)$$

Some values of T may be equal so you can add their probabilities together

- The **mean** is calculated in the same way
 - $\circ E(T) = \sum t P(X = x)$
- The **variance** is calculated using the same formula
 - $Var(T) = E(T^2) [E(T)]^2$

Are there any shortcuts?

- There are formulae which can be used if the transformation is linear
 - T = aX + b where a and b are constants
- If the transformation is **not linear** then there are **no shortcuts**
 - \circ You will have to first find the probability distribution of T

What are the formulae for $E(aX \pm b)$ and $Var(aX \pm b)$?

- If a and b are constants then the following formulae are true:
 - E(aX±b) = aE(X)±b PAPERS PRACTICE
 - $Var(aX \pm b) = a^2 Var(X)$

These are given in the **formula booklet**

- This is the same as linear transformations of data
 - The mean is affected by multiplication and addition/subtraction
 - The variance is affected by multiplication but not addition/subtraction
- Remember division can be written as a multiplication

$$\circ \quad \frac{X}{a} = \frac{1}{a}X$$







4.5 Binomial Distribution

4.5.1 The Binomial Distribution

Properties of Binomial Distribution

What is a binomial distribution?

- A binomial distribution is a discrete probability distribution
- + A discrete random variable X follows a binomial distribution if it counts the number of
 - ${\it successes}$ when an experiment satisfies the following conditions:
 - There are a fixed finite number of trials (n)
 - The outcome of each trial is **independent** of the outcomes of the other trials
 - There are exactly two outcomes of each trial (success or failure)
 - The probability of success is constant (p)
- If X follows a binomial distribution then it is denoted $X \sim B(n, p)$
 - *n* is the **number of trials**
 - p is the probability of success
- The probability of failure is 1 p which is sometimes denoted as q
- The formula for the probability of *r* successful trials is given by:

•
$$P(X = r) = {}^{n}C_{r} \times p^{r}(1 - p)^{n - r}$$
 for $r = 0, 1, 2, ..., n$
 ${}^{n}C_{r} = \frac{n!}{r!(n - r)!}$ where $n! = n \times (n - 1) \times (n - 2) \times ... \times 3 \times 2 \times 1$

• You will be expected to use the distribution function on your **GDC to calculate** probabilities with the binomial distribution **DRACTICE**

What are the important properties of a binomial distribution?

• The expected number (mean) of successful trials is

$$E(X) = np$$

- You are given this in the **formula booklet**
- The variance of the number of successful trials is

$$Var(X) = np(1-p)$$

- You are given this in the **formula booklet**
- Square root to get the **standard deviation**
- The distribution can be represented visually using a vertical line graph
 - If p is close to 0 then the graph has a tail to the right
 - $\circ~$ If p is close to 1 then the graph has a tail to the left
 - $\circ~$ If p is close to 0.5 then the graph is roughly symmetrical
 - If p = 0.5 then the graph is symmetrical






Modelling with Binomial Distribution

How do I set up a binomial model?

- Identify what a trial is in the scenario
 - For example: rolling a dice, flipping a coin, checking hair colour
- Identify what the successful outcome is in the scenario
 For example: rolling a 6, landing on tails, having black hair
- Identify the parameters
 - *n* is the number of trials and *p* is the probability of success in each trial
- Make sure you clearly state what your random variable is
 - $\circ~$ For example, let X be the number of students in a class of 30 with black hair

What can be modelled using a binomial distribution?

- Anything that satisfies the four conditions
- For example: let T be the number of times a fair coin lands on tails when flipped 20 times:
 - A trial is flipping a coin: There are 20 trials so **n = 20**
 - We can assume each coin flip does not affect subsequent coin flips: they are **independent**
 - A success is when the coin lands on tails: **Two outcomes** tails or not tails (heads)
- The coin is fair: The probability of tails is constant with p = 0.5
- Sometimes it might seem like there are more than two outcomes
 - For example: let Y be the number of yellow cars that are in a car park full of 100 cars Although there are more than two possible colours of cars, here the trial is whether a car is yellow so there are two outcomes (yellow or not yellow) Y would still need to fulfil the other conditions in order to follow a binomial distribution
- Sometimes a sample may be taken from a population
 - For example: 30% of people in a city have blue eyes, a sample of 30 people from the city is taken and X is the number of them with blue eyes

As long as the population is large and the sample is random then it can be assumed that each person has a 30% chance of having blue eyes

What can not be modelled using a binomial distribution?

- Anything where the number of trials is **not fixed** or is **infinite**
 - $\circ~$ The number of emails received in an hour
 - The number of times a coin is flipped until it lands on heads
- Anything where the outcome of **one trial affects** the outcome of the **other trials**
 - The number of caramels that a person eats when they eat 5 sweets from a bag containing 6 caramels and 4 marshmallows
 - If you eat a caramel for your first sweet then there are less caramels left in the bag when you choose your second sweet
 - Anything where there are **more than two outcomes** of a trial A person's shoe size
 - The number a dice lands on when rolled
 - Anything where the **probability of success changes**



The number of times that a person can swim a length of a swimming pool in under a minute when swimming 50 lengths

The probability of swimming a lap in under a minute will decrease as the person gets tired

The probability is **not constant**

Exam Tip

• An exam question might involve different types of distributions so make it clear which distribution is being used for each variable





Worked Example

It is known that 8% of a large population are immune to a particular virus. Mark takes a sample of 50 people from this population. Mark uses a binomial model for the number of people in his sample that are immune to the virus.

~	١
a)

 $State \, the \, distribution \, that \, Mark \, uses.$





State two assumptions that Mark must make in order to use a binomial model.



C)

Calculated the expected number of people in the sample that are immune to the virus.

Formula booklet	Binomial distribution $X \sim B(n, p)$	
E(X) = 50 × 0.08	Mean	E(X) = np
4 people		



4.5.2 Calculating Binomial Probabilities

Calculating Binomial Probabilities

Throughout this section we will use the random variable $X \sim B(n, p)$. For binomial, the probability of X taking a non-integer or negative value is always zero. Therefore any values of X mentioned in this section will be assumed to be non-negative integers.

How do I calculate P(X = x): the probability of a single value for a binomial distribution?

- You should have a GDC that can calculate binomial probabilities
- You want to use the "Binomial Probability Distribution" function
 - This is sometimes shortened to BPD, Binomial PD or Binomial Pdf
- You will need to enter:
 - The 'x' value the value of x for which you want to find P(X = x)
 - The 'n' value the **number of trials**
 - The 'p' value the **probability of success**
- Some calculators will give you the option of **listing the probabilities** for **multiple values** of x at once
- There is a formula that you can use but you are expected to be able to use the distribution function on your GDC

•
$$P(X=x) = {}^{n}C_{n} \times p^{x}(1-p)^{n-x}$$

$${}^{n}C_{x} = \frac{n!}{r!(n-r)!}$$

How do I calculate $P(a \le X \le b)$: the cumulative probabilities for a binomial distribution?

- You should have a GDC that can calculate cumulative binomial probabilities
 - Most calculators will find $P(a \le X \le b)$
 - Some calculators can only find $P(X \le b)$ The identities below will help in this case
- You should use the "Binomial Cumulative Distribution" function
 - This is sometimes shortened to BCD, Binomial CD or Binomial Cdf
- You will need to enter:
 - The lower value this is the **value a**

This can be zero in the case $P(X \le b)$

- The upper value this is the **value b** This can be *n* in the case $P(X \ge a)$
- The 'n' value the **number of trials**
- The 'p' value the **probability of success**

How do I find probabilities if my GDC only calculates $P(X \le x)$?

- To calculate $P(X \le x)$ just enter x into the cumulative distribution function
- To calculate P(X < x) use:
 - $P(X < x) = P(X \le x 1)$ which works when X is a binomial random variable $P(X < 5) = P(X \le 4)$



- To calculate P(X > x) use:
 - $P(X > x) = 1 P(X \le x)$ which works for any random variable X $P(X > 5) = 1 - P(X \le 5)$
- To calculate $P(X \ge x)$ use:
 - $P(X \ge x) = 1 P(X \le x 1)$ which works when X is a binomial random variable $P(X \ge 5) = 1 - P(X \le 4)$
- To calculate $P(a \le X \le b)$ use:
 - $P(a \le X \le b) = P(X \le b) P(X \le a 1)$ which works when X is a binomial random variable

 $\mathsf{P}(5 \leq X \leq 9) = \mathsf{P}(X \leq 9) - \mathsf{P}(X \leq 4)$

What if an inequality does not have the equals sign (strict inequality)?

- For a binomial distribution (as it is discrete) you could **rewrite all strict inequalities** (< and >) as **weak inequalities** (≤ and ≥) by using the identities for a binomial distribution
 - $P(X < x) = P(X \le x 1)$ and $P(X > x) = P(X \ge x + 1)$
 - For example: $P(X < 5) = P(X \le 4)$ and $P(X > 5) = P(X \ge 6)$
- It helps to think about the **range of integers** you want
 - Identify the smallest and biggest integers in the range
- If your range has no minimum or maximum then use 0 or n
 - $\circ P(X \le b) = P(0 \le X \le b)$
 - $\circ P(X \ge a) = P(a \le X \le n)$
- P(a < X ≤ b) = P(a+1 ≤ X ≤ b) • P(5 < X ≤ 9) = P(6 ≤ X ≤ 9)
 P(b ≤ X ≤ 9) = P(b ≤ X ≤ 9)
- $P(a \le X < b) = P(a \le X \le b 1)$ • $P(5 \le X < 9) = P(5 \le X \le 8)$
- $P(a < X < b) = P(a+1 \le X \le b-1)$ • $P(5 < X < 9) = P(6 \le X \le 8)$

Exam Tip

- If the question is in context then write down the inequality as well as the final answer
 - This means you still might gain a mark even if you accidentally type the wrong numbers into your GDC







4.6 Normal Distribution

4.6.1 The Normal Distribution

Properties of Normal Distribution

The binomial distribution is an example of a discrete probability distribution. The normal distribution is an example of a **continuous** probability distribution.

What is a continuous random variable?

- A continuous random variable (often abbreviated to CRV) is a random variable that can take **any value** within a range of infinite values
 - Continuous random variables usually measure something
 - For example, height, weight, time, etc

What is a continuous probability distribution?

- A continuous probability distribution is a probability distribution in which the random variable *X* is continuous
- The probability of X being a **particular value is always zero**
 - P(X=k) = 0 for any value k
 - Instead we define the probability density function f(x) for a specific value
 This is a function that describes the relative likelihood that the random variable would be close to that value
 - We talk about the **probability** of X being within a **certain range**
- A continuous probability distribution can be represented by a continuous graph (the values for X along the horizontal axis and probability **density** on the vertical axis)
- The area under the graph between the points x = a and x = b is equal to P(a ≤ X ≤ b)
 The total area under the graph equals 1
- As P(X = k) = 0 for any value k, it does not matter if we use strict or weak inequalities
 P(X ≤ k) = P(X < k) for any value k when X is a continuous random variable

What is a normal distribution?

- A normal distribution is a continuous probability distribution
- The **continuous random variable** *X* can follow a normal distribution if:
 - The distribution is **symmetrical**
 - The distribution is **bell-shaped**
- If X follows a normal distribution then it is denoted $X \sim N(\mu, \sigma^2)$
 - μ is the **mean**
 - σ^2 is the **variance**
 - $\circ \sigma$ is the standard deviation
- If the mean changes then the graph is translated horizontally
- If the **variance** increases then the graph is **widened horizontally** and **made taller vertically** to maintain the same area
 - A small variance leads to a tall curve with a narrow centre
 - A large variance leads to a short curve with a wide centre





What are the important properties of a normal distribution?

- The **mean** is μ
- The **variance** is σ^2
 - If you need the standard deviation remember to square root this
- The normal distribution is symmetrical about
 - Mean = Median = Mode = μ
- There are the results:
 - Approximately **two-thirds (68%)** of the data lies within **one standard deviation** of the mean $(\mu \pm \sigma)$
 - Approximately **95%** of the data lies within **two standard deviations** of the mean ($\mu \pm 2\sigma$)
 - Nearly all of the data (99.7%) lies within three standard deviations of the mean ($\mu \pm 3\sigma$)





Modelling with Normal Distribution

What can be modelled using a normal distribution?

- A lot of real-life continuous variables can be modelled by a normal distribution provided that the population is large enough and that the variable is **symmetrical** with **one mode**
- For a normal distribution X can take any real value, however values far from the mean (more than 4 standard deviations away from the mean) have a probability density of **practically zero**
 - This fact allows us to model variables that are not defined for all real values such as height and weight

What can not be modelled using a normal distribution?

- Variables which have **more than one mode** or **no mode**
 - \circ For example: the number given by a random number generator
- Variables which are not symmetrical
 - $\circ~$ For example: how long a human lives for

Exam Tip

• An exam question might involve different types of distributions so make it clear which distribution is being used for each variable





4.6.2 Calculations with Normal Distribution

Calculating Normal Probabilities

Throughout this section we will use the random variable $X \sim N(\mu, \sigma^2)$. For X distributed normally, X can take any real number. Therefore any values mentioned in this section will be assumed to be real numbers.

How do I find probabilities using a normal distribution?

- The **area under a normal curve** between the points x = a and x = b is equal to the **probability** P(a < X < b)
 - Remember for a normal distribution you do not need to worry about whether the inequality is strict (< or >) or weak (≤ or ≥)

 $P(a < X < b) = P(a \le X \le b)$

• You will be **expected to use** distribution functions on your **GDC** to find the probabilities when working with a normal distribution

How do I calculate P(X = x): the probability of a single value for a normal distribution?

- The probability of a single value is always zero for a normal distribution
 - You can picture this as the area of a single line is zero
- P(X=x)=0
- Your GDC is likely to have a "Normal Probability Density" function
 - This is sometimes shortened to NPD, Normal PD or Normal Pdf
 - IGNORE THIS FUNCTION for this course!
 - This calculates the **probability density function** at a point **NOT the probability**

How do I calculate P(a < X < b): the probability of a range of values for a normal distribution?

- You need a GDC that can calculate cumulative normal probabilities
- You want to use the "Normal Cumulative Distribution" function
- \circ This is sometimes shortened to NCD, Normal CD or Normal Cdf
- You will need to enter:
 - \circ The 'lower bound' this is the value a
 - $\circ~$ The 'upper bound' this is the value b
 - $\circ~$ The ' μ' value this is the mean
 - \circ The ' σ ' value this is the standard deviation
- Check the order carefully as some calculators ask for standard deviation before mean
 - Remember it is the standard deviation
 - so if you have the **variance** then **square root it**
- Always sketch a quick diagram to visualise which area you are looking for

How do I calculate P(X > a) or P(X < b) for a normal distribution?

- You will still use the "Normal Cumulative Distribution" function
- P(X > a) can be estimated using an **upper bound that is sufficiently bigger** than the **mean**
 - Using a value that is more than 4 standard deviations **bigger than the mean** is quite accurate
 - Or an easier option is just to input lots of 9's for the upper bound (99999999... or 10⁹⁹)
- P(X < b) can be estimated using a lower bound that is sufficiently smaller than the mean



- Using a value that is more than 4 standard deviations **smaller than the mean** is quite accurate
- Or an easier option is just to input lots of 9's for the lower bound with a negative sign (-99999999... or -10⁹⁹)

Are there any useful identities?

- $P(X < \mu) = P(X > \mu) = 0.5$
- As P(X=a) = 0 you can use:
 - $\circ P(X < a) + P(X > a) = 1$
 - $\circ P(X > a) = 1 P(X < a)$
 - P(a < X < b) = P(X < b) P(X < a)
- These are useful when:
 - The mean and/or standard deviation are unknown
 - You only have a diagram
 - You are working with the **inverse distribution**

Exam Tip

Check carefully whether you have entered the standard deviation or variance
 into your GDC









Inverse Normal Distribution

Given the value of P(X < a) how do I find the value of a?

- Your GDC will have a function called "Inverse Normal Distribution"
 - Some calculators call this InvN
- Given that P(X < a) = p you will need to enter:
 - The 'area' this is the value p Some calculators might ask for the 'tail' - this is the left tail as you know the area to the left of a
 - The ' μ ' value this is the mean
 - The ' σ ' value this is the standard deviation

Given the value of P(X > a) how do I find the value of a?

- If your calculator **does** have the **tail option** (left, right or centre) then you can use the "Inverse Normal Distribution" function straightaway by:
 - Selecting 'right' for the tail
 - Entering the area as 'p'
- If your calculator **does not** have the **tail option** (left, right or centre) then:
 - Given P(X > a) = p
 - Use P(X < a) = 1 P(X > a) to rewrite this as

$$P(X < a) = 1 - p$$

Then use the method for P(X < a) to find a

) Exam Tip

- Always check your **answer makes sense**
 - If P(X < a) is less than 0.5 then a should be smaller than the mean
 - If P(X < a) is more than 0.5 then a should be bigger than the mean
 - A sketch will help you see this







4.6.3 Standardisation of Normal Variables

Standard Normal Distribution

What is the standard normal distribution?

- The standard normal distribution is a normal distribution where the mean is 0 and the standard deviation is 1
 - \circ It is denoted by Z
 - $Z \sim N(0, 1^2)$

Why is the standard normal distribution important?

- Any **normal distribution curve** can be transformed to the standard normal distribution curve by a **horizontal translation** and a **horizontal stretch**
- Therefore we have the relationship:

$$\circ \ Z = \frac{X - \mu}{\sigma}$$

• Where
$$X \sim \mathrm{N}(\mu, \sigma^2)$$
 and $Z \sim \mathrm{N}(0, 1^2)$

• Probabilities are related by:

•
$$P(a < X < b) = P\left(\frac{a-\mu}{\sigma} < Z < \frac{b-\mu}{\sigma}\right)$$

• This will be useful when the mean or variance is unknown

• Some mathematicians use the function $\Phi(z)$ to represent P(Z < z)

z-values

What are z-values (standardised values)? **RACTICE**

- For a normal distribution $X \sim N(\mu, \sigma^2)$ the z-value (standardised value) of an x-value tells you how many standard deviations it is away from the mean
 - If z = 1 then that means the x-value is 1 standard deviation bigger than the mean
 - If z = -1 then that means the x-value is 1 standard deviation smaller than the mean
- If the x-value is more than the mean then its corresponding z-value will be positive
- If the x-value is less than the mean then its corresponding z-value will be negative
- The z-value can be calculated using the formula:

$$\circ z = \frac{x-\mu}{\sigma}$$

• This is given in the formula booklet

• z-values can be used to compare values from different distributions



Finding Sigma and Mu

How do I find the mean (μ) or the standard deviation (σ) if one of them is unknown?

- If the mean or standard deviation of $X \sim N(\mu, \sigma^2)$ is unknown then you will need to use the standard normal distribution
- You will need to use the formula

•
$$z = \frac{x - \mu}{\sigma}$$
 or its rearranged form $x = \mu + \sigma z$

• You will be given a **probability for a specific value** of

•
$$P(X < x) = p \text{ or } P(X > x) = p$$

- To find the unknown parameter:
- STEP 1: Sketch the normal curve
 - Label the known value and the mean
- STEP 2: Find the *z*-value for the given value of *x*
 - Use the **Inverse Normal Distribution** to find the value of z such that P(Z < z) = p or P(Z > z) = p
 - $\circ~$ Make sure the direction of the inequality for Z is consistent with the inequality for X
 - Try to use lots of decimal places for the z-value or store your answer to avoid rounding errors

You should use at least one extra decimal place within your working than your intended degree of accuracy for your answer

- STEP 3: Substitute the known values into $z = \frac{x \mu}{\sigma}$ or $x = \mu + \sigma z$
 - You will be given and one of the parameters (μ or σ) in the question
 - You will have calculated zin STEP 2 ERS PRACTICE
- STEP 4: Solve the equation

How do I find the mean (μ) and the standard deviation (σ) if both of them are unknown?

- If **both** of them are **unknown** then you will be given two probabilities for two specific values of **x**
- The process is the same as above
 - You will now be able to **calculate two z -values**
 - You can form **two equations** (rearranging to the form $x = \mu + \sigma z$ is helpful)
 - You now have to **solve the two equations simultaneously** (you can use your calculator to do this)
 - Be careful not to mix up which z-value goes with which value of x





It is known that the times, in minutes, taken by students at a school to eat their lunch can be modelled using a normal distribution with mean μ minutes and standard deviation σ minutes.

Given that 10% of students at the school take less than 12 minutes to eat their lunch and 5% of the students take more than 40 minutes to eat their lunch, find the mean and standard deviation of the time taken by the students at the school.





4.7 Further Probability Distributions

4.7.1 Probability Density Function

Calculating Probabilities using PDF

A **continuous random variable** can take *any* value in an interval so is typically used when continuous quantities are involved (time, distance, weight, etc)

What is a probability density function (p.d.f.)?

- For a continuous random variable, a function can be used to model probabilities
 This function is called a **probability density function** (p.d.f.), denoted by f(x)
- For f(x) to represent a p.d.f. the following conditions must apply
 - $f(x) \ge 0$ for **all** values of x
 - The area under the graph of y = f(x) must total 1
- In most problems, the **domain** of x is restricted to an interval, a ≤ X ≤ b say, with all values of x outside of the interval having f(x)=0

How do I find probabilities using a probability density function (p.d.f.)?

• The probability that the continuous random variable X lies in the interval $a \le X \le b$, where X has the probability density function f(x), is given by

$$P(a \le X \le b) = \int_{a}^{b} f(x) \, dx$$

- P(a ≤ X ≤ b) = P(a < X < b)
 For any continuous random variable (including the normal distribution) P(X = n) = 0
 One way to think of this is that a = b in the integral above
- For **linear** functions it can be easier to find the probability using the area of geometric shapes
 - Rectangles: A = bh
 - Triangles: $A = \frac{1}{2}(bh)$
 - Trapezoids: $A = \frac{1}{2}(a+b)h$

How do I determine whether a function is a pdf?

- Some questions may ask for justification of the use of a given function for a probability density function
 - In such cases check that the function meets the two conditions
 - $f(x) \ge 0$ for **all** values of x
 - total area under the graph is 1

How do l use a pdf to find probabilities?

STEP 1

Identify the **probability density function**, f(x) - this may be given as a **graph**, an **equation** or as a **piecewise function**



e.g.
$$f(x) = \begin{cases} 0.02x & 0 \le x \le 10\\ 0 & \text{otherwise} \end{cases}$$

Identify the **limits** of X for a particular problem Remember that $P(a \le X \le b) = P(a < X < b)$

STEP 2

Sketch, or use your GDC to draw, the graph of y = f(x) Look for basic shapes (rectangles, triangles and trapezoids) as finding these areas is easier without using integration Look for symmetry in the graph that may make the problem easier Break the area required into two or more parts if it makes the problem easier

STEP 3

Find the area(s) required using basic shapes or integration and answer the question

- Trickier problems may involve finding a limit of the integral given its value
 - ∘ i.e. Find one of the boundaries in the domain of X, given the probability e.g. Find the value of a given that $P(0 \le X \le a) = 0.09$







Worked Example

The continuous random variable, X, has probability density function.

$$f(x) = \begin{cases} 0.08x & 0 \le x \le 5\\ 0 & \text{otherwise} \end{cases}$$

a)

Show that f(x) can represent a probability density function.



b)

Find, both geometrically and using integration, $P(0 \le X \le 2)$.







Median & Mode of a CRV

What is meant by the median of a continuous random variable?

• The **median**, *m*, of a **continuous random variable**, *X*, with **probability density function** *f*(*x*) is defined as the value of *X* such that

$$P(X < m) = P(X > m) = 0.5$$

- Since P(X = m) = 0 this can also be written as $P(X \le m) = P(X \ge m) = 0.5$
- IF the p.d.f. is symmetrical (i.e. the graph of *y* = *f*(*x*) is symmetrical) then the **median** will be halfway between the lower and upper limits of *x*
 - In such cases the graph of y = f(x) has **axis** of **symmetry** in the line x = m

How do I find the median of a continuous random variable?

• The **median**, *m*, of a continuous random variable, *X*, with probability density function *f*(*x*) is defined as the value of *X* such that

$$\int_{-\infty}^{m} f(x) \, \mathrm{d}x = \frac{1}{2}$$

or

$$\int_{m}^{\infty} f(x) \, \mathrm{d}x = \frac{1}{2}$$

- The equation that should be used will depend on the information in the question
 - If the graph of y = f(x) is symmetrical, symmetry may be used to deduce the median
 - This may often be the case if f(x) is **linear** and the **area under the graph** is a basic **shape** such as a **rectangle**

How do I find the median of a continuous random variable with a piecewise p.d.f.?

- For **piecewise functions**, the **location** of the **median** will determine **which equation** to use in order to find it
 - For example

if
$$f(x) = \begin{cases} \frac{1}{5}x & 0 \le x \le 2\\ \frac{2}{15}(5-x) & 2 \le x \le 5\\ 0 & \text{otherwise} \end{cases}$$

then $\int_{0}^{2} \int_{0}^{1} x \, dx = 0.4$ so the median must lie in the interval $2 \le x \le 5$

so to find the median, m, solve $\int_{2}^{m} \frac{2}{15} (5-x) dx = 0.1$

('0.4 of the area' already used for $0 \le x \le 2$)

Use a GDC to plot the function and evalutae integral(s)

What is meant by the mode of a continuous random variable?



• The mode of a continuous random variable, X, with probability density function f(x) is the value of x that produces the greatest value of f(x)

How do I find the mode of a continuous random variable?

- This will depend on the **type** of **function** *f*(*x*); the easiest way to find the **mode** is by considering the **shape** of the **graph** of *y* = *f*(*x*)
- If the graph is a curve with a maximum point, the mode can be found by differentiating and solving f'(x) = 0
 - If there is more than one solution to f'(x) = 0 then further work may be needed in deducing the mode
 - There could be **more than one** mode
 - Look for **valid values** of *x* from the **domain** of the p.d.f.
 - Use the **second derivative** (f''(x)) to **deduce** the **nature** of each **stationary point Check** the **values** of f(x) at the **lower** and **upper limits** of x, one of these could be the **maximum value** f(x) reaches
- If the graph of y = f(x) is **symmetrical**, symmetry may be used to deduce the mode
 - For a symmetrical p.d.f., median = mode = mean







Worked Example

The continuous random variable X has probability function f(x) defined as

$$f(x) = \frac{1}{64} (16x - x^3) \quad 0 \le x \le 4$$

a)

Find the median of X, giving your answer to three significant figures.



b)

Find the exact value of the mode of X.



Differentiate, solving
$$f'(x)=0$$
 to find the mode
 $f'(x) = \frac{1}{64}(16-3x^2)$
 $16-3x^2=0$
Using a GOC (ensure you get exact answers)
 $x = \pm \frac{14\sqrt{3}}{3}$
Clearly from sketch of graph, $x = \frac{1}{3}\sqrt{3}$ is a (local) maximum
Also, $x = -\frac{4\sqrt{3}}{3}$ does not lie in the interval $0 \le x \le 4$
 \therefore Mode = $\frac{14}{3}\sqrt{3}$





Mean & Variance of a CRV

What are the mean and variance of a continuous random variable?

- E(X) is the **expected value**, or **mean**, of the **continuous random variable** X
 - \circ E(X) can also be denoted by μ
- Var(X) is the **variance** of the continuous random variable X
 - Var(X) can also be denoted by σ^2
 - $\circ~$ The standard deviation, $\sigma,$ is the square root of the variance

How do I find the mean and variance of a continuous random variable?

• The mean is given by

$$\mu = \mathrm{E}(X) = \int_{-\infty}^{\infty} x f(x) \, \mathrm{d}x$$

- This is given in the **formula booklet**
- If the graph of y = f(x) has **axis** of **symmetry**, x = a, then E(X) = a
- The variance is given by

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

where $E(X^2) = \int_{-\infty}^{\infty} x^2 f(x) dx$

- This is given in the formula booklet
- Another version of the variance is given in the **formula booklet**

$$Var(x) = \int_{-\infty}^{\infty} (x - \mu)^2 f(x) \, dx = \int_{-\infty}^{\infty} x^2 f(x) \, dx - \mu^2$$

- but the first version above is usually more practical for solving problems
- Be careful about confusing $E(X^2)$ and $[E(X)]^2$

•
$$E(X^2) = \int_{-\infty}^{\infty} x^2 f(x) dx$$
 "mean of the squares"
• $[E(X)]^2 = \left[\int_{-\infty}^{\infty} x f(x) dx\right]^2$ "mean of the squares"

How do I find the mean and variance of a linear transformation of a continuous random variable?

• For the **continuous random variable**, *X*, with **mean** E(*X*) and **variance** Var(*X*) then

$$E(aX+b) = aE(X) + b$$

and

$$\operatorname{Var}(aX+b) = a^2 \operatorname{Var}(X)$$



Exam Tip

 \bigcirc

• Using your **GDC** to draw the graph of *y* = *f*(*x*) can **highlight** any **symmetrical** properties which **reduce** the **work** involved in finding the **mean** and **variance**











$$\sigma = \sqrt{Var(x)} \quad Var(x) = E(x^2) - [E(x)]^2$$

$$E(x^2) = \int_0^2 \infty^2 [1 \cdot 5\infty^2 (1 - 0 \cdot 5\infty)] dx$$

$$Using GDC, \quad E(x^2) = 1 \cdot 6$$

$$\therefore \sigma = \sqrt{1 \cdot 6} - (1 \cdot 2)^2 = \sqrt{0 \cdot 16} = 0 \cdot 4$$

$$E(x^2) \int C E(x)$$

$$\sigma = 0 \cdot 4$$

