
OCR SINGLE STATISTICS REVISION SHEET

The OCR single maths A level is examined with three compulsory 2 hour exams. Each paper carries equal weight ($33\frac{1}{3}\%$) and each paper is marked out of 100 marks.

Paper 1 contains only pure mathematics.

Paper 2 contains pure mathematics and statistics (roughly 50 marks for each section).

Paper 3 contains pure mathematics and mechanics (roughly 50 marks for each section).

Therefore your A level is roughly $\frac{2}{3}$ pure maths, $\frac{1}{6}$ statistics, and $\frac{1}{6}$ mechanics. This revision sheet *should* cover all of the statistics you need. *Please* get in contact if you spot anything missing.

I hope you find this revision sheet useful and wish you the very best of luck with your studies.

J.M.S.

“Without data, all you are is just another person with an opinion.”

Representation Of Data

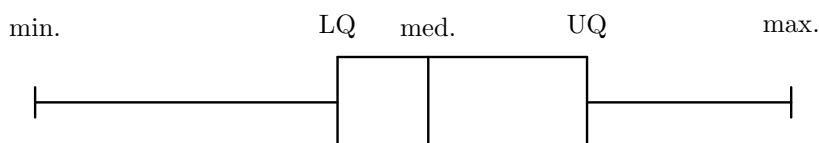
- You must be happy constructing unordered, back-to-back and ordered stem and leaf diagrams. They show the overall distribution of the data and back-to-back diagrams allow you to compare two sets of data.
- Cumulative frequency graphs. The cumulative frequency is a “running total” of the frequencies as you go up the values. For example

x	f	⇒	Create	⇒
$0 \leq x < 5$	8	⇒	Cumulative	⇒
$5 \leq x < 10$	13	⇒	Frequency	⇒
$10 \leq x < 15$	17	⇒		
$15 \leq x < 20$	10			

x (upper limit of)	cum. freq
5	8
10	21
15	38
20	48

Plot the second of these tables and join it with a smooth curve to form the *cumulative frequency curve*. From this the median and the two quartiles can be found.

- Once these values are found we can draw a *box and whisker diagram*. The box and whisker diagram uses five values: the minimum, the maximum, the lower quartile, the upper quartile and the median. It is good for showing spread and comparing two quantities.



- Histograms are usually drawn for continuous data in classes. If the classes have equal widths, then you merely plot amount against frequency.
- If the classes do *not* have equal widths then we need to create a new column for *frequency density*. Frequency density is defined by $f.d. = \frac{\text{frequency}}{\text{class width}}$. The *area* of the bars are what represents the frequency, *not* the height.
- Frequency polygons are made by joining together the mid-points of the bars of a histogram with a ruler.

Measures Of Location

- The *mean* (arithmetic mean) of a set of data $\{x_1, x_2, x_3 \dots x_n\}$ is given by

$$\bar{x} = \frac{\text{sum of all values}}{\text{the number of values}} = \frac{\sum x}{n}.$$

When finding the mean¹ of a frequency distribution the mean is given by

$$\frac{\sum(xf)}{\sum f} = \frac{\sum(xf)}{n}.$$

- If a set of numbers is arranged in ascending (or descending) order the *median* is the number which lies half way along the series. It is the number that lies at the $\left(\frac{n+1}{2}\right)^{\text{th}}$ position. Thus the median of $\{13, 14, 15, 15\}$ lies at the $2\frac{1}{2}$ position \Rightarrow average of 14 and 15 \Rightarrow median = 14.5.
- The *mode* of a set of numbers is the number which occurs the most frequently. Sometimes no mode exists; for example with the set $\{2, 4, 7, 8, 9, 11\}$. The set $\{2, 3, 3, 3, 4, 5, 6, 6, 6, 7\}$ has two modes 3 and 6 because each occurs three times. One mode \Rightarrow “unimodal”. Two modes \Rightarrow “bimodal”. More than two modes \Rightarrow “multimodal”.

	ADVANTAGES	DISADVANTAGES
MEAN	<ul style="list-style-type: none"> ★ The best known average. ★ Can be calculated exactly. ★ Makes use of all the data. ★ Can be used in further statistical work. 	<ul style="list-style-type: none"> ★ Greatly affected by extreme values. ★ Can't be obtained graphically. ★ When the data are discrete can give an impossible figure (2.34 children).
MEDIAN	<ul style="list-style-type: none"> ★ Can represent an actual value in the data. ★ Can be obtained even if some of the values in a distribution are unknown. ★ Unaffected by irregular class widths and unaffected by open-ended classes. ★ Not influenced by extreme values. 	<ul style="list-style-type: none"> ★ For grouped distributions its value can only be estimated from an ogive. ★ When only a few items available or when distribution is irregular the median may not be characteristic of the group. ★ Can't be used in further statistical calculations.
MODE	<ul style="list-style-type: none"> ★ Unaffected by extreme values. ★ Easy to calculate. ★ Easy to obtain from a histogram. 	<ul style="list-style-type: none"> ★ May exist more than one mode. ★ Can't be used for further statistical work. ★ When the data are grouped its value cannot be determined exactly.

Measures Of Spread

- The simplest measure of spread is the *range*. Range = $x_{\max} - x_{\min}$.
- The interquartile range is simply the upper quartile take away the lower quartile. Both of these values are usually found from a cumulative frequency graph (above).
- The *sum of squares from the mean* is called the *sum of squares* and is denoted

$$S_{xx} = \sum (x - \bar{x})^2 = \sum x^2 - n\bar{x}^2.$$

For example given the data set $\{3, 6, 7, 8\}$ the mean is 6; $\sum x^2 = 9 + 36 + 49 + 64 = 158$; so $S_{xx} = \sum x^2 - n\bar{x}^2 = 158 - 4 \times 6^2 = 14$.²

¹Statistics argues that the average person has one testicle and that 99.999% of people have more than the average number of arms...

²Or we could have done $S_{xx} = \sum (x - \bar{x})^2 = (3 - 6)^2 + (6 - 6)^2 + (7 - 6)^2 + (8 - 6)^2 = 14$.

- The *standard deviation* (σ) is defined: $\sigma = \sqrt{\text{variance}} = \sqrt{\frac{S_{xx}}{n}} = \sqrt{\frac{\sum x^2}{n} - \bar{x}^2}$.
- *Example:* Given the set of data {5, 7, 8, 9, 10, 10, 14} calculate the standard deviation. Firstly we note that $\bar{x} = 9$.

$$\begin{aligned}\sigma &= \sqrt{\frac{\sum x^2}{n} - \bar{x}^2} = \sqrt{\frac{(5^2 + \dots + 14^2)}{7} - 9^2} \\ &= \sqrt{\frac{615}{7} - 81} = 2.6186\dots\end{aligned}$$

- When dealing with frequency distributions such as

x	1	2	3	4	5
f	4	5	7	5	4

, we *could* calculate σ by writing out the data³ and carrying out the calculations as above, but this is clearly slow and inefficient. To our rescue comes a formula for σ that allows direct calculation from the table. This is

$$\sigma = \sqrt{\frac{\sum(x^2f)}{n} - \bar{x}^2}.$$

- *Example:* Calculate mean and sd for the above frequency distribution. For easy calculation we need to add certain columns to the usual x and f columns thus;

x	f	xf	x^2f
1	4	4	4
2	5	10	20
3	7	21	63
4	5	20	80
5	4	20	100
$n = \sum f = 25$		$\sum(xf) = 75$	$\sum(x^2f) = 267$.

$$\text{So } \bar{x} = \frac{\sum(xf)}{n} = \frac{75}{25} = 3 \text{ and } \sigma = \sqrt{\frac{\sum(x^2f)}{n} - \bar{x}^2} = \sqrt{\frac{267}{25} - 3^2} = 1.2961\dots$$

- *Linear Coding.* Given the set of data {2, 3, 4, 5, 6} we can see that $\bar{x} = 4$ and it can be calculated that $\sigma = 1.414$ (3dp). If we add 20 to all the data points we can see that the mean becomes 24 and the standard deviation will be unchanged. If the data set is multiplied by 3 we can see that the mean becomes 12 and the standard deviation would become three times as large (4.743 (3dp)).
- If, instead of being given $\sum x$ and $\sum x^2$, you were given $\sum(x - a)$ and $\sum(x - a)^2$ for some constant a , you just use the substitution $u = x - a$ and use $\sum u$ and $\sum u^2$ to work out the mean of u and the standard deviation of u . Then, using the above paragraph, we know $\bar{x} = \bar{u} + a$ and $\sigma_x = \sigma_u$.

Probability

- An *independent event* is one which has no effect on subsequent events. The events of spinning a coin and then cutting a pack of cards are independent because the way in which the coin lands has no effect on the cut. For two *independent events* A & B

$$\mathbb{P}(A \text{ and } B) = \mathbb{P}(A) \times \mathbb{P}(B).$$

³{1, 1, 1, 1, 2, 2, 2, 2, 2, 3, 3, 3, 3, 3, 3, 3, 3, 3, 4, 4, 4, 4, 4, 4, 5, 5, 5, 5}!!!

For example a fair coin is tossed and a card is then drawn from a pack of 52 playing cards. Find the probability that a head and an ace will result.

$$\mathbb{P}(\text{head}) = \frac{1}{2}, \quad \mathbb{P}(\text{ace}) = \frac{4}{52} = \frac{1}{13}, \quad \text{so } \mathbb{P}(\text{head and ace}) = \frac{1}{2} \times \frac{1}{13} = \frac{1}{26}.$$

- *Mutually Exclusive Events.* Two events which cannot occur at the same time are called mutually exclusive. The events of throwing a 3 or a 4 in a single roll of a fair die are mutually exclusive. For any two mutually exclusive events

$$\mathbb{P}(A \text{ or } B) = \mathbb{P}(A) + \mathbb{P}(B).$$

For example a fair die with faces of 1 to 6 is rolled once. What is the probability of obtaining either a 5 or a 6?

$$\mathbb{P}(5) = \frac{1}{6}, \quad \mathbb{P}(6) = \frac{1}{6}, \quad \text{so } \mathbb{P}(5 \text{ or } 6) = \frac{1}{6} + \frac{1}{6} = \frac{1}{3}.$$

- *Non-Mutually Exclusive Events.* When two events can both happen they are called non-mutually exclusive events. For example studying English and studying Maths at A Level are non-mutually exclusive. By considering a Venn diagram of two events A & B we find

$$\mathbb{P}(A \text{ or } B) = \mathbb{P}(A) + \mathbb{P}(B) - \mathbb{P}(A \text{ and } B),$$

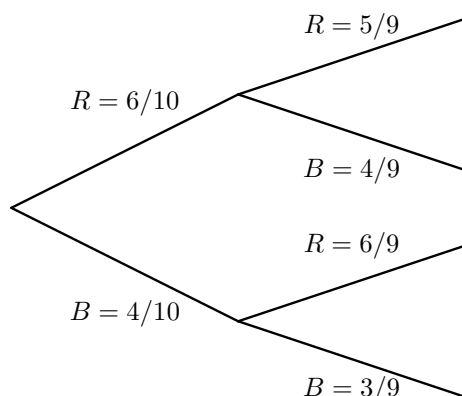
$$\mathbb{P}(A \cup B) = \mathbb{P}(A) + \mathbb{P}(B) - \mathbb{P}(A \cap B).$$

- *Tree Diagrams.* These may be used to help solve probability problems when more than one event is being considered. The probabilities on any branch section must sum to one. You multiply along the branches to discover the probability of that branch occurring.

For example a box contains 4 black and 6 red pens. A pen is drawn from the box and it is not replaced. A second pen is then drawn. Find the probability of

- two red pens being obtained.
- two black pens being obtained.
- one pen of each colour being obtained.
- two red pens *given* that they are the same colour.

Draw tree diagram to discover:



$$(i) \mathbb{P}(\text{two red pens}) = \frac{6}{10} \times \frac{5}{9} = \frac{30}{90} = \frac{1}{3}.$$

$$(ii) \mathbb{P}(\text{two black pens}) = \frac{4}{10} \times \frac{3}{9} = \frac{12}{90} = \frac{2}{15}.$$

$$(iii) \mathbb{P}(\text{one of each colour}) = 1 - \frac{30}{90} - \frac{12}{90} = \frac{8}{15}.$$

$$(iv) \mathbb{P}(\text{two reds} \mid \text{same colour}) = \frac{1/3}{1/3 + 2/15} = \frac{5}{7}.$$

- *Conditional Probability.* In the above example we see that the probability of two red pens is $\frac{1}{3}$, but the probability of two red pens *given that both pens are the same colour* is $\frac{5}{7}$. This is known as conditional probability. $\mathbb{P}(A \mid B)$ mean the probability of A *given* that B has happened. It is governed by

$$\mathbb{P}(A \mid B) = \frac{\mathbb{P}(A \text{ and } B)}{\mathbb{P}(B)} = \frac{\mathbb{P}(A \cap B)}{\mathbb{P}(B)}.$$

For example if there are 120 students in a year and 60 study Maths, 40 study English and 10 study both then

$$\mathbb{P}(\text{study English} \mid \text{study Maths}) = \frac{\mathbb{P}(\text{study Maths \& English})}{\mathbb{P}(\text{study Maths})} = \frac{10/120}{60/120} = \frac{1}{6}.$$

- A is independent of B if $\mathbb{P}(A) = \mathbb{P}(A \mid B) = \mathbb{P}(A \mid B')$. (i.e. whatever happens in B the probability of A remains unchanged.) For example flicking a coin and then cutting a deck of cards to try and find an ace are independent because

$$\mathbb{P}(\text{cutting ace}) = \mathbb{P}(\text{cutting ace} \mid \text{flick head}) = \mathbb{P}(\text{cutting ace} \mid \text{flick tail}) = \frac{1}{13}.$$

Permutations And Combinations

- Factorials are defined $n! \equiv n \times (n-1) \times (n-2) \times \dots \times 2 \times 1$. Many expressions involving factorials simplify with a bit of thought. For example $\frac{n!}{(n-2)!} = n(n-1)$. Also there is a convention that $0! = 1$.
- The number of ways of arranging n different objects in a line is $n!$ For example how many different arrangements are there if 4 different books are to be placed on a bookshelf? There are 4 ways in which to select the first book, 3 ways in which to choose the second book, 2 ways to pick the third book and 1 way left for the final book. The total number of different ways is $4 \times 3 \times 2 \times 1 = 4!$
- Permutations. The number of ways of selecting r objects from n when *the order of the selection matters* is ${}^n P_r$. It can be calculated by

$${}^n P_r = \frac{n!}{(n-r)!}.$$

For example in how many ways can the gold, silver and bronze medals be awarded in a race of ten people? The order in which the medals are awarded matters, so the number of ways is given by ${}^{10} P_3 = 720$.

In another example how many words of four letters can be made from the word CONSIDER? This is an arrangement of four out of eight different objects where the order matters so there are ${}^8 P_4 = \frac{8!}{4!} = 1680$ different words.

- Combinations. The number of ways of selecting r objects from n when *the order of the selection does not matter* is ${}^n C_r$. It can be calculated by

$${}^n C_r \equiv \binom{n}{r} = \frac{n!}{r!(n-r)!}.$$

For example in how many ways can a committee of 5 people be chosen from 8 applicants? Solution is given by ${}^8 C_5 = \frac{8!}{5! \times 3!} = 56$.

In another example how many ways are there of selecting your lottery numbers (where one selects 6 numbers from 49)? It does not matter which order you choose your numbers, so there are ${}^{49} C_6 = 13\,983\,816$ possible selections.

- If letters are repeated in a 'word', then you just divide through by the factorials of each repeat. Therefore there are $\frac{11!}{4! \times 4! \times 2!}$ arrangements of the word 'MISSISSIPPI'.⁴

⁴As a non-mathematical aside, find two fruits that are anagrams of each other.

- You must be good at ‘choosing committee’ questions [be on the lookout, they can be in disguise]. For example how many ways are there of choosing a committee of 3 women and 4 men from a group containing 10 women and 5 men? There are $\binom{10}{3}$ ways of choosing the women (the order doesn’t matter) and $\binom{5}{4}$ ways of choosing the men. Therefore overall there are $\binom{5}{4} \times \binom{10}{3}$ ways of choosing the committee.
- Example: If I deal six cards from a standard deck of cards, in how many ways can I get exactly four clubs? Well there are $\binom{13}{4}$ ways of getting the clubs, and $\binom{39}{2}$ ways of getting the non-clubs, so therefore the answer to the original question is $\binom{13}{4} \times \binom{39}{2}$.
- When considering arranging objects in a line we start from the principle that there are $n!$ ways of arranging n objects. In the harder examples you need to be cunning.

For example three siblings join a queue with 5 other people making 8 in total.

1. How many way are there of arranging the 8 in a queue? Easy; 8!
 2. How many ways are there of arranging the 8, such that the siblings are together? We, we imaging the three siblings tied together. There are therefore 6! ways of arranging the 5 and the bundle of siblings and then there are 3! ways of arranging the siblings in the bundle. Therefore the answer is $6! \times 3!$
 3. How many ways are there of arranging the siblings so they are not together? There are 5! ways of arranging the five without the siblings. There are then 6 places for the first sibling to go, 5 for the second, and 4 for the third. Therefore $5! \times 6 \times 5 \times 4$.
- To calculate probabilities we go back to first principles and remember that probability is calculated from the number of ways of getting what we want over the total number of possible outcomes. So in the above example, if the 8 are arranged randomly in a line, what is the probability of the siblings being together? $\mathbb{P}(\text{together}) = \frac{6! \times 3!}{8!}$.

Going back to the four club question if it asked for the probability of getting exactly four clubs if I dealt exactly six cards from the pack, the answer would be $\frac{\binom{13}{4} \times \binom{39}{2}}{\binom{52}{6}}$. The $\binom{52}{6}$ represents the total number of ways I can deal six cards from the 52.

Probability Distributions

- A random variable is a quantity whose value depends on chance. The outcome of a random variable is usually denoted by a capital letter (e.g. X). We read $\mathbb{P}(X = 2)$ as the probability that the random variable takes the value 2. For a fair die, $\mathbb{P}(X = 5) = \frac{1}{6}$.
- For discrete random variables they are usually presented in a table. For example for a fair die:

x	1	2	3	4	5	6
$\mathbb{P}(X = x)$	$\frac{1}{6}$	$\frac{1}{6}$	$\frac{1}{6}$	$\frac{1}{6}$	$\frac{1}{6}$	$\frac{1}{6}$

- In general, for any event, the probability distribution is of the form

x	x_1	x_2	x_3	x_4	x_5	x_6	\dots
$\mathbb{P}(X = x)$	p_1	p_2	p_3	p_4	p_5	p_6	\dots

- As before, it is crucial that we remember the probabilities sum to one. This can be useful at the start or problems where a constant must be evaluated. For example in:

x	1	2	3	4
$\mathbb{P}(X = x)$	k	k	$2k$	$4k$

We discover $k + k + 2k + 4k = 1$, so $k = \frac{1}{8}$.

Binomial And Geometric Distributions

- The **binomial distribution** is applicable when you have fixed number of *repeated, independent* ‘trials’ such that each trial can be viewed as ‘success’ (p) or ‘fail’ ($q = 1 - p$). In justifying a binomial distribution you must not just quote the previous sentence; you must apply it to the situation in the question. For example: “Binomial is applicable because the probability of each tulip flowering is independent of each other tulip and the probability of flowering is a constant”.
- For example if I throw darts at a dart board and my chance of hitting a double is 0.1 and I throw 12 darts at the board and my chance of hitting a double is independent of all the other throws then a binomial distribution will be applicable. We let X be the number of doubles I hit. X can therefore take the values $\{0, 1, 2, \dots, 11, 12\}$; i.e. there are 13 possible outcomes. $p = 0.1$; the probability of success and $q = 1 - p = 0.9$; the probability of failure. We write $X \sim B(n, p)$ which here is $X \sim B(12, 0.1)$.
- I would always advise you to visualise the tree diagram. From this we can ‘see’ that $\mathbb{P}(X = 12) = 0.1^{12}$ and $\mathbb{P}(X = 0) = 0.9^{12}$. In general

$$\mathbb{P}(X = x) = \binom{n}{x} \times p^x \times q^{n-x}.$$

So in the example, the probability I hit exactly 7 doubles is $\mathbb{P}(X = 7) = \binom{12}{7} \times 0.1^7 \times 0.9^5$.

- For questions such as $\mathbb{P}(X \leq 5)$ or $\mathbb{P}(X \geq 8)$ you must be able to use the tables in the formula book. The tables always give $\mathbb{P}(X \leq \text{something})$. You must be able to convert probabilities to this form and then read off from the table. For $X \sim B(10, 0.35)$.

$$\mathbb{P}(X \leq 7) = 0.9952,$$

$$\mathbb{P}(X < 5) = \mathbb{P}(X \leq 4) = 0.7515,$$

$$\mathbb{P}(X \geq 7) = 1 - \mathbb{P}(X \leq 6) = 1 - 0.9740 = 0.0260,$$

$$\mathbb{P}(X > 3) = 1 - \mathbb{P}(X \leq 3) = 1 - 0.5138 = 0.4862.$$

- The **geometric distribution** is applicable when you are looking for how long you wait until an event has occurred. The events must be *repeated, independent and success/fail*. Potentially you could wait forever until a success occurs; something to look for if you are unsure what distribution to apply. Similar to the binomial you must justify *in the context of the question*.
- Going back to the darts example, we could rephrase it as how long must I wait until I hit a double? Let X be the number of throws until I hit a double. We write $X \sim \text{Geo}(0.1)$. X can take the values $\{1, 2, 3, \dots\}$.
- Obviously $\mathbb{P}(X = 1) = 0.1$. Less obviously $\mathbb{P}(X = 4) = 0.9^3 \times 0.1$ (I must have three failures and *then* my success). In general

$$\mathbb{P}(X = x) = q^{x-1} \times p.$$

- There are no tables for the geometric distribution because there does not need to be. To calculate $\mathbb{P}(X \geq 5)$ we must have had 4 failures. Therefore $\mathbb{P}(X \geq 5) = q^4 = (1-p)^4$. Also to calculate $\mathbb{P}(X \leq 6)$ we use the fact that $\mathbb{P}(X \leq 6) = 1 - \mathbb{P}(X \geq 7) = 1 - q^6 = 1 - (1-p)^6$. In general

$$\mathbb{P}(X \geq x) = (1-p)^{x-1} \quad \text{and} \quad \mathbb{P}(X \leq x) = 1 - (1-p)^x.$$

Expectation And Variance Of A Random Variable

- The expected value of the event is denoted $\mathbb{E}(X)$ or μ . It is defined

$$\mathbb{E}(X) = \mu = \boxed{\sum x\mathbb{P}(X = x)}.$$

For example for a fair die with

x	1	2	3	4	5	6
$\mathbb{P}(X = x)$	$\frac{1}{6}$	$\frac{1}{6}$	$\frac{1}{6}$	$\frac{1}{6}$	$\frac{1}{6}$	$\frac{1}{6}$

we find:

$$\begin{aligned} \mathbb{E}(X) &= \left(1 \times \frac{1}{6}\right) + \left(2 \times \frac{1}{6}\right) + \left(3 \times \frac{1}{6}\right) + \left(4 \times \frac{1}{6}\right) + \left(5 \times \frac{1}{6}\right) + \left(6 \times \frac{1}{6}\right) \\ &= 3\frac{1}{2}. \end{aligned}$$

- The variance of an event is denoted $\text{Var}(X)$ or σ^2 and is defined

$$\text{Var}(X) = \sigma^2 = \mathbb{E}(X^2) - (\mathbb{E}(X))^2 = \mathbb{E}(X^2) - \mu^2 = \boxed{\sum x^2\mathbb{P}(X = x) - \mu^2}.$$

So for the *biased* die with distribution

x	1	2	3	4	5	6
$\mathbb{P}(X = x)$	$\frac{1}{3}$	$\frac{1}{6}$	0	0	$\frac{1}{6}$	$\frac{1}{3}$

we find that

$$\mathbb{E}(X) = \left(1 \times \frac{1}{3}\right) + \left(2 \times \frac{1}{6}\right) + (3 \times 0) + (4 \times 0) + \left(5 \times \frac{1}{6}\right) + \left(6 \times \frac{1}{3}\right) = 3\frac{1}{2}$$

and

$$\begin{aligned} \text{Var}(X) &= \sum x^2\mathbb{P}(X = x) - \mu^2 \\ &= \left(1^2 \times \frac{1}{3}\right) + \left(2^2 \times \frac{1}{6}\right) + (3^2 \times 0) + (4^2 \times 0) + \left(5^2 \times \frac{1}{6}\right) + \left(6^2 \times \frac{1}{3}\right) - 3\frac{1}{2}^2 \\ &= 17\frac{1}{6} - 3\frac{1}{2}^2 = 4\frac{11}{12}. \end{aligned}$$

- The expectation of a binomial distribution $B(n, p)$ is np . The variance of $B(n, p)$ is npq .
- The expectation of a geometric distribution $\text{Geo}(p)$ is $\frac{1}{p}$.

Correlation

- The Product Moment Correlation Coefficient is a number (r) calculated on a set of bivariate data that tells us how correlated two data sets are.
- The value of r is such that $-1 < r < 1$. If $r = 1$ you have perfect positive linear correlation. If $r = -1$ you have perfect negative linear correlation. If $r = 0$ then there exists no correlation between the data sets.
- It is defined

$$r = \frac{\sum(x_i - \bar{x})(y_i - \bar{y})}{\sqrt{\sum(x_i - \bar{x})^2 \sum(y_i - \bar{y})^2}} = \frac{S_{xy}}{\sqrt{S_{xx}S_{yy}}}$$

where we define the individual components as

$$\begin{aligned} S_{xx} &= \sum x^2 - \frac{1}{n} (\sum x)^2, \\ S_{yy} &= \sum y^2 - \frac{1}{n} (\sum y)^2, \\ S_{xy} &= \sum xy - \frac{1}{n} \sum x \sum y. \end{aligned}$$

- So to calculate r for the data set

x	14	12	16	18	21	13	15	17
y	1	2	4	5	2	8	5	6

we write the data in columns and add extra ones. We then sum the columns and calculate from these sums. Note that in the above example $n = 8$ (i.e. the number of pairs, not the number of individual data pieces).

x	y	x^2	y^2	xy
14	1	196	1	14
12	2	144	4	24
16	4	256	16	64
18	5	324	25	90
21	2	441	4	42
13	8	169	64	104
15	5	225	25	75
17	6	289	36	102
126	33	2044	175	515

Therefore

$$\begin{aligned} S_{xx} &= \sum x^2 - \frac{1}{n} (\sum x)^2 = 2044 - \frac{126^2}{8} = 59.5, \\ S_{yy} &= \sum y^2 - \frac{1}{n} (\sum y)^2 = 175 - \frac{33^2}{8} = 38.875, \\ S_{xy} &= \sum xy - \frac{1}{n} \sum x \sum y = 515 - \frac{126 \times 33}{8} = -4.75. \end{aligned}$$

Therefore

$$r = \frac{S_{xy}}{\sqrt{S_{xx}S_{yy}}} = \frac{-4.75}{\sqrt{59.5 \times 38.875}} = -0.09876\dots$$

Therefore the data has very, very weak negative correlation. Basically it has no *meaningful* correlation.

- It can be shown that if one (or both) of the variables are transformed in a linear fashion i.e. if we replace the x values by, say, $\frac{x-4}{3}$ (or any transformation formed by $+$, $-$, \div or \times with constants) then the value of r will be unchanged.

- You need to be able to calculate Spearman's rank correlation coefficient (r_s). You will be given a table and you will need to (in the next 2 columns) rank the data. If two data points are tied then you (e.g. the 2nd and 3rd are tied) then you rank them both 2.5.

%	IQ	Rank %	Rank IQ	d	d^2
89	143	2.5	1	1.5	2.25
55	89	7	8	-1	1
72	102	5	6	-1	1
91	136	1	2	-1	1
89	126	2.5	3	-0.5	0.25
30	60	9	9	0	0
71	115	6	4	2	4
53	100	8	7	1	1
78	103	4	5	-1	1

Now $r_s = 1 - \frac{6 \sum d^2}{n(n^2 - 1)}$. $\sum d^2$ is just the sum of the d^2 column in the table and n is the number of pairs of data; here $n = 9$. We therefore find $r_s = 1 - \frac{6 \times 11.5}{9(81-1)} = 0.9041\dot{6}$. Therefore we see a strong degree of positive association.

- If r_s is close to -1 then strong negative association. If close to zero then no meaningful association/agreement.

Regression

- For any set of bivariate data (x_i, y_i) there exist two possible regression lines; 'y on x' and 'x on y'.
- If neither is controlled (see below) then if you want to predict y from a given value of x , you use the 'y on x' line. If you want to predict x from a given value of y , you use the 'x on y' line.
- The 'y on x' line is defined

$$y = a + bx \quad \text{where} \quad b = \frac{S_{xy}}{S_{xx}} \quad \text{and} \quad a = \bar{y} - b\bar{x}.$$

- The 'x on y' line is defined

$$x = a' + b'y \quad \text{where} \quad b' = \frac{S_{xy}}{S_{yy}} \quad \text{and} \quad a' = \bar{x} - b'\bar{y}.$$

- Both regression lines pass through the average point (\bar{x}, \bar{y}) .
- In the example in the book (P180) the height of the tree is the dependent variable and the circumference of the tree is the independent variable. This is because the experiment has been constructed to see how the height of the tree depends on its circumference.
- If one variable is being controlled by the experimenter (e.g. x), it is called a controlled variable. If x is controlled you would never use the 'x on y' regression line. Only use the 'y on x' line. You would use this to predict y from x (expected) and x from y (not expected)

Continuous Random Variables

- A continuous random variable (crv) is usually described by means of a probability density function (pdf) which is defined for all real x . It must satisfy

$$\int_{-\infty}^{\infty} f(x) dx = 1 \quad \text{and} \quad f(x) \geq 0 \text{ for all } x.$$

- Probabilities are represented by areas under the pdf. For example the probability that X lies between a and b is

$$\mathbb{P}(a < X < b) = \int_a^b f(x) dx.$$

It is worth noting that for any specific value of X , $\mathbb{P}(X = \text{value}) = 0$ because the area of a single value is zero.

- The median is the value m such that

$$\int_{-\infty}^m f(x) dx = \frac{1}{2}.$$

That is; the area under the curve is cut in half at the value of the median. Similarly the lower quartile (Q_1) and upper quartile (Q_3) are defined

$$\int_{-\infty}^{Q_1} f(x) dx = \frac{1}{4} \quad \text{and} \quad \int_{-\infty}^{Q_3} f(x) dx = \frac{3}{4}.$$

- The expectation of X is defined

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

Compare this to the discrete definition of $\sum x \mathbb{P}(X = x)$. Always be on the lookout for symmetry in the distribution before carrying out a long integral; it could save you a lot of time. You should therefore always sketch the distribution if you can.

- The variance of X is defined

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

Again, compare this to the discrete definition of $\sum x^2 \mathbb{P}(X = x) - \mu^2$. Don't forget to subtract μ^2 at the end; someone always does!

- The main use for this chapter is to give you the basics you may need for the normal distribution. The normal distribution is by far the most common crv.

The Normal Distribution

- The normal distribution (also known as the Gaussian distribution⁵) is the most common crv. It is found often in nature; for example daffodil heights, human IQs and pig weights can all be modelled by the normal curve. A normal distribution can be summed up by two parameters; its mean (μ) and its variance (σ^2). For a random variable X we say $X \sim N(\mu, \sigma^2)$.

⁵I do wish we would call it the Gaussian distribution. Carl Friedrich Gauss. Arguably the greatest mathematician ever. German...

- As with all crvs probabilities are given by areas; i.e. $\mathbb{P}(a < X < b) = \int_a^b f(x) dx$. However the $f(x)$ for a normal distribution is complicated and impossible to integrate exactly. We therefore need to use tables to help us. Since there are an infinite number of $N(\mu, \sigma^2)$ distributions we use a special one called the standard normal distribution. This is $Z \sim N(0, 1^2)$.
- The tables given to you work out the areas to the left of a value. The notation used is $\Phi(z) = \int_{-\infty}^z f(z) dz$. So $\Phi(0.2)$ is the area to the left of 0.2 in the standard normal distribution. The tables do not give $\Phi(\text{negative value})$ so there are some tricks of the trade you must be comfortable with. These and they are always helped by a sketch and remembering that the area under the whole curve is one. For example

$$\begin{aligned}\Phi(z) &= 1 - \Phi(-z) \\ \mathbb{P}(Z > z) &= 1 - \Phi(z)\end{aligned}$$

- Real normal distributions are related to the standard distribution by

$$Z = \frac{X - \mu}{\sigma} \quad (\dagger).$$

So if $X \sim N(30, 16)$ and we want to answer $\mathbb{P}(X > 24)$ we convert $X = 24$ to $Z = (24 - 30)/4 = -1.5$ and answer $\mathbb{P}(Z > -1.5) = \mathbb{P}(Z < 1.5) = 0.9332$.

- Another example; If $Y \sim N(100, 5^2)$ and we wish to calculate $\mathbb{P}(90 < Y < 105)$. Converting to $\mathbb{P}(-2 < Z < 1)$ using \dagger . Then finish off with

$$\mathbb{P}(-2 < Z < 1) = \Phi(1) - \Phi(-2) = \Phi(1) - (1 - \Phi(2)) = 0.8413 - (1 - 0.9772) = 0.8185.$$

- You must also be able to do a ‘reverse’ lookup from the table. Here you don’t look up an area from a z value, but look up a z value from an area.

For example find a such that $\mathbb{P}(Z < a) = 0.65$. Draw a sketch as to what this means; to the left of some value a the area is 0.65. Therefore, reverse looking up we discover $a = 0.385$.

- Harder example; Find b such that $\mathbb{P}(Z > b) = 0.9$. Again a sketch shows us that the area to the right of b must be 0.9, so b must be negative. Considering the sketch carefully, we discover $\mathbb{P}(Z < -b) = 0.9$, so reverse look up tells us $-b = 1.282$, so $b = -1.282$.

- Reverse look up is then combined with \dagger in questions like this. For $X \sim N(\mu, 5^2)$ it is known $\mathbb{P}(X < 20) = 0.8$; find μ . Here you will find it easier if you draw both a sketch for the X and also for Z and marking on the important points. The z value by reverse look up is found to be 0.842. Therefore by \dagger we obtain, $0.842 = (20 - \mu)/5$, so $\mu = 15.79$.

- Harder example; $Y \sim (\mu, \sigma^2)$ you know $\mathbb{P}(Y < 20) = 0.25$ and $\mathbb{P}(Y > 30) = 0.4$. You should obtain two \dagger equations;

$$-0.674 = \frac{20 - \mu}{\sigma} \quad \text{and} \quad 0.253 = \frac{30 - \mu}{\sigma} \quad \Rightarrow \quad \mu = 27.27 \text{ and } \sigma = 10.79.$$

- The binomial distribution can sometimes be approximated by the normal distribution. If $X \sim B(n, p)$ and $np > 5$ and $nq > 5$ then we can use $V \sim N(np, npq)$ as an approximation. Because we are going from a discrete distribution to a continuous, a continuity correction must be used.

- For example if $X \sim B(90, \frac{1}{3})$ we can see $np = 30 > 5$ and $nq = 60 > 5$ so we can use $V \sim N(30, 20)$. Some examples of the conversions:

$$\begin{aligned}\mathbb{P}(X = 29) &\approx \mathbb{P}(28.5 < V < 29.5), \\ \mathbb{P}(X > 25) &\approx \mathbb{P}(V > 25.5), \\ \mathbb{P}(5 \leq X < 40) &\approx \mathbb{P}(4\frac{1}{2} < V < 39\frac{1}{2}).\end{aligned}$$

The Poisson Distribution

- The Poisson distribution is a discrete random variable (like the binomial or geometric distribution). It is defined

$$\mathbb{P}(X = x) = e^{-\lambda} \frac{\lambda^x}{x!}.$$

X can take the values $0, 1, 2, \dots$ and the probabilities depend on only one parameter, λ . Therefore we find

x	0	1	2	3	\dots
$\mathbb{P}(X = x)$	$e^{-\lambda} \frac{\lambda^0}{0!}$	$e^{-\lambda} \frac{\lambda^1}{1!}$	$e^{-\lambda} \frac{\lambda^2}{2!}$	$e^{-\lambda} \frac{\lambda^3}{3!}$	\dots

- For a Poisson distribution $\mathbb{E}(X) = \text{Var}(X) = \lambda$. We write $X \sim \text{Po}(\lambda)$.
- As for the binomial we use tables to help us and they are given (for various different λ s) in the form $\mathbb{P}(X \leq x)$. So if $\lambda = 5$ and we wish to discover $\mathbb{P}(X < 8)$ we do $\mathbb{P}(X < 8) = \mathbb{P}(X \leq 7) = 0.8666$. Also note that if we want $\mathbb{P}(X \geq 4)$ we would use the fact that probabilities sum to one, so $\mathbb{P}(X \geq 4) = 1 - \mathbb{P}(X \leq 3) = 1 - 0.2650 = 0.7350$.
- The Poisson distribution can be used as an approximation to the binomial distribution provided $n > 50$ and $np < 5$. If these conditions are met and $X \sim B(n, p)$ we use $W \sim \text{Po}(np)$. [No continuity correction required since we are approximating a discrete by a discrete.]
- For example with $X \sim B(60, \frac{1}{30})$ both conditions are met and we use $W \sim \text{Po}(2)$. Therefore some example of some calculations:

$$\mathbb{P}(X \leq 3) \approx \mathbb{P}(W \leq 3) = 0.8571 \text{ (from tables)}$$

$$\mathbb{P}(3 < X \leq 7) \approx \mathbb{P}(3 < W \leq 7) = \mathbb{P}(W \leq 7) - \mathbb{P}(W \leq 3) = 0.9989 - 0.8571 = 0.1418.$$

- The normal distribution can be used as an approximation to the to the Poisson distribution if $\lambda > 15$. So if $X \sim \text{Po}(\lambda)$ we use $Y \sim N(\lambda, \lambda)$. However, here we *are* approximating a discrete by a continuous, so a continuity correction must be applied.
- For example if $X \sim \text{Po}(50)$ we can use $Y \sim N(50, 50)$ since $\lambda > 15$. To calculate $\mathbb{P}(X = 49)$ we would calculate (using $Z = (X - \mu)/\sigma$)

$$\begin{aligned} \mathbb{P}(X = 49) &\approx \mathbb{P}(48.5 < Y < 49.5) = \mathbb{P}(-0.212 < Z < -0.071) \\ &= \mathbb{P}(0.071 < Z < 0.212) \\ &= \Phi(0.212) - \Phi(0.071) \\ &= 0.5840 - 0.5283 = 0.0557. \end{aligned}$$

Similarly

$$\begin{aligned} \mathbb{P}(X < 55) &\approx \mathbb{P}(Y < 54.5) \\ &= \mathbb{P}\left(Z < \frac{54.5 - 50}{\sqrt{50}}\right) \\ &= \mathbb{P}(Z < 0.6364) \\ &= 0.738. \end{aligned}$$

Sampling

- If a sample is taken from an underlying population you can view the mean of this sample as a random variable in its own right. This is a subtle point and you should dwell on it! If you can't get to sleep sometime, you should lie awake thinking about it. (I had to.)
- If the underlying population has $\mathbb{E}(X) = \mu$ and $\text{Var}(X) = \sigma^2$, then the distribution of the mean of the sample, \bar{X} , is

$$\mathbb{E}(\bar{X}) = \mu \text{ (the same as the underlying)} \text{ and } \text{Var}(\bar{X}) = \frac{\sigma^2}{n}.$$

This means that the larger your sample, the less likely it is that the mean of this sample is a long way from the population mean. So if you are taking a sample, make it as big as you can!

- If your sample is sufficiently large (roughly > 30) the central limit theorem (CLT) states that the distribution of the sample mean is approximated by

$$\bar{X} \sim N\left(\mu, \frac{\sigma^2}{n}\right)$$

no matter what the underlying distribution is.

- If the underlying population is discrete you need to include a $\frac{1}{2n}$ correction factor when using the CLT. For example $\mathbb{P}(\bar{X} > 3.4)$ for a discrete underlying with a sample size of 45 would mean you calculate $\mathbb{P}(\bar{X} > 3.4 + \frac{1}{90})$.
- If the underlying population is a normal distribution then no matter how large the sample is (e.g. just 4) we can say

$$\bar{X} \sim N\left(\mu, \frac{\sigma^2}{n}\right).$$

- If you have the whole population data available to you then to calculate the mean you use $\mu = \frac{\sum x}{n}$ and to calculate the variance you use

$$\sigma^2 = \frac{\sum x^2}{n} - \bar{x}^2 = \frac{\sum x^2 - n\bar{x}^2}{n}.$$

However you do not usually have all the data. It is more likely that you merely have a sample from the population. From this sample you may want to estimate the population mean and variance. As you would expect your best estimate of the population mean is the mean of the sample $\frac{\sum x}{n}$. However the best estimate of the population variance is not the variance of the sample. You must calculate s^2 where

$$s^2 = \frac{\sum x^2 - n\bar{x}^2}{n-1} = \frac{n}{n-1} \left(\frac{\sum x^2 - n\bar{x}^2}{n} \right) = \frac{n}{n-1} \left(\frac{\sum x^2}{n} - \bar{x}^2 \right).$$

Some textbooks use $\hat{\sigma}$ to mean s ; they both mean 'the unbiased estimator of the population σ '. So

$$\text{(Estimate of population variance)} = \frac{n}{n-1} \times \text{(Sample variance)}.$$

- You could be given raw data ($\{x_1, x_2, \dots, x_n\}$) in which you just do a direct calculation. Or summary data ($\sum x^2, \sum x$ and n). Or you could be given the sample variance and n . From all of these you should be able to calculate s^2 . It should be clear from the above section how to do this.

Continuous Hypothesis Testing

- In *any* hypothesis test you will be testing a ‘null’ hypothesis H_0 against an ‘alternative’ hypothesis H_1 . In S2, your H_0 will only *ever* be one of these three:

$$H_0 : p = \text{something}$$

$$H_0 : \lambda = \text{something}$$

$$H_0 : \mu = \text{something}$$

Don’t deviate from this and you can’t go wrong. Notice that it does *not* say $H_0 = p = \text{something}$.

- The book gives three approaches to continuous hypothesis testing, but they are all essentially the same. You always compare the probability of what you have seen (under H_0) and anything more extreme, and compare this probability to the significance level. If it is less than the significance level, then you reject H_0 and if it is greater, then you accept H_0 .
- Remember we connect the real (X) world to the standard (Z) world using $Z = \frac{X-\mu}{\sigma}$.
- You can do this by:
 - Calculating the probability of the observed value and anything more extreme and comparing to the significance level.
 - Finding the critical Z -values for the test and finding the Z -value for the observed event and comparing. (e.g. critical Z -values of 1.96 and -1.96 ; if observed Z is 1.90 we accept H_0 ; if observed is -2.11 the reject H_0 .)
 - Finding the critical values for \bar{X} . For example critical values might be 17 and 20. If X lies between them then accept H_0 ; else reject H_0 .
- Example: P111 Que 8. Using method 3 from above.

Let X be the amount of magnesium in a bottle. We are told $X \sim N(\mu, 0.18^2)$. We are taking a sample of size 10, so $\bar{X} \sim N(\mu, \frac{0.18^2}{10})$. Clearly

$$H_0 : \mu = 6.8$$

$$H_1 : \mu \neq 6.8.$$

We proceed assuming H_0 is correct. Under H_0 , $\bar{X} \sim N(6.8, \frac{0.18^2}{10})$. This is a 5% two-tailed test, so we need $2\frac{1}{2}\%$ at each end of our normal distribution. The critical Z values are (by reverse lookup) $Z_{\text{crit}} = \pm 1.960$. To find how these relate to \bar{X}_{crit} we convert thus

$$Z_{\text{crit}} = \frac{\bar{X}_{\text{crit}} - \mu}{\sqrt{\frac{\sigma^2}{n}}}$$

$$1.960 = \frac{\bar{X}_{\text{crit}} - 6.8}{\sqrt{\frac{0.18^2}{10}}}$$

and $-1.960 = \frac{\bar{X}_{\text{crit}} - 6.8}{\sqrt{\frac{0.18^2}{10}}}$

These solve to $\bar{X}_{\text{crit}} = 6.912$ and $\bar{X}_{\text{crit}} = 6.688$. The observed \bar{X} is 6.92 which lies just outside the acceptance region. We therefore reject H_0 and conclude that the amount of magnesium per bottle is probably *different* to 6.8. [The book is in error in claiming that we conclude it is bigger than 6.8.]

Discrete Hypothesis Testing

- For any test with discrete variables, it is usually best to find the critical value(s) for the test you have set and hence the critical region. The critical value is the first value *at which you would reject* the null hypothesis.
- For example if testing $X \sim B(16, p)$ we may test (at the 5% level)

$$\begin{aligned}H_0 &: p = \frac{5}{6} \\H_1 &: p < \frac{5}{6}.\end{aligned}$$

We are looking for the value at the lower end of the distribution (remember the “<” acts as an arrow telling us where to look in the distribution). We find $\mathbb{P}(X \leq 11) = 0.1134$ and $\mathbb{P}(X \leq 10) = 0.0378$. Therefore the critical value is 10. Thus the critical region is $\{0, 1, 2, \dots, 9, 10\}$. So when the result for the experiment is announced, if it lies in the critical region, we reject H_0 , else accept H_0 .

- Another example: If testing $X \sim B(20, p)$ at the 10% level with

$$\begin{aligned}H_0 &: p = \frac{1}{6} \\H_1 &: p \neq \frac{1}{6}.\end{aligned}$$

Here we have a two tailed test with 5% at either end of the distribution. At the lower end we find $\mathbb{P}(X = 0) = 0.0261$ and $\mathbb{P}(X \leq 1) = 0.1304$ so the critical value is 0 at the lower end. At the upper end we find $\mathbb{P}(X \leq 5) = 0.8982$ and $\mathbb{P}(X \leq 6) = 0.9629$. Therefore

$$\begin{aligned}\mathbb{P}(X \geq 6) &= 1 - \mathbb{P}(X \leq 5) = 1 - 0.8982 = 0.1018 \\ \mathbb{P}(X \geq 7) &= 1 - \mathbb{P}(X \leq 6) = 1 - 0.9629 = 0.0371\end{aligned}$$

So at the upper end we find $X = 7$ to be the critical value. [Remember that at the upper end, the critical value is always one more than the upper of the two values where the gap occurs; here the gap was between 5 and 6 in the tables, so 7 is the critical value.] The critical region is therefore $\{0, 7, 8, \dots, 20\}$.

- There is a Poisson example in the ‘Errors in hypothesis testing’ section.

Errors In Hypothesis Testing

- A Type I error is made when a true null hypothesis is rejected.
- A Type II error is made when a false null hypothesis is accepted.
- For continuous hypothesis tests, the $\mathbb{P}(\text{Type I error})$ is just the significance level of the test. [This fact should be obvious; if not think about it harder!]
- For a Type II error, you must consider something like the example on page 140/1 which is superbly explained. From the original test, you will have discovered the acceptance and the rejection region(s). When you are told the real mean of the distribution and asked to calculate the $\mathbb{P}(\text{Type II error})$, you must use the new, real mean and the old standard deviation (with a new normal distribution; e.g. $N(\mu_{\text{new}}, \sigma_{\text{old}}^2/n)$) and work out the probability that the value lies within the old acceptance region. [Again, the book is *very* good on this and my explanation is poor.]

- For discrete hypothesis tests, the $\mathbb{P}(\text{Type I error})$ is not merely the stated significance level of the test. The stated value (e.g. 5%) is merely the ‘notional’ value of the test. The true significance level of the test (and, therefore, the $\mathbb{P}(\text{Type I error})$) is the probability of all the values in the rejection region, given the truth of the null hypothesis.

For example in a binomial hypothesis test we might have discovered the rejection region was $X \leq 3$ and $X \geq 16$. If the null hypothesis was “ $H_0: p = 0.3$ ”, then the true significance level of the test would be $\mathbb{P}(X \leq 3 \text{ or } X \geq 16 \mid p = 0.3)$.

- To calculate $\mathbb{P}(\text{Type II error})$ you would, given the true value for p (or λ for Poisson), calculate the probability of the *complementary* event. So in the above example, if the true value of p was shown to be 0.4, you would calculate $\mathbb{P}(3 < X < 16 \mid p = 0.4)$.
- Worked example for Poisson: A hypothesis is carried out to test the following:

$$H_0 : \lambda = 7$$

$$H_1 : \lambda \neq 7$$

$$\alpha = 10\%$$

Two tailed test.

Under H_0 , $X \sim \text{Po}(7)$. We discover the critical values are $X = 2$ and $X = 13$. The critical region is therefore $X \leq 2$ and $X \geq 13$.

Therefore $\mathbb{P}(\text{Type I error})$ and the true value of the test is therefore

$$\begin{aligned} \mathbb{P}(X \leq 2 \text{ or } X \geq 13 \mid \lambda = 7) &= \mathbb{P}(X \leq 2) + \mathbb{P}(X \geq 13) \\ &= \mathbb{P}(X \leq 2) + 1 - \mathbb{P}(X \leq 12) \\ &= 0.0296 + 1 - 0.9730 \\ &= 0.0566 = 5.66\%. \end{aligned}$$

Given that the true value of λ was shown to be 10, then $\mathbb{P}(\text{Type II error})$ would be

$$\begin{aligned} \mathbb{P}(2 < X < 13 \mid \lambda = 10) &= \mathbb{P}(X \leq 12) - \mathbb{P}(X \leq 2) \\ &= 0.7916 - 0.0028 \\ &= 0.7888 = 78.88\%. \end{aligned}$$