

## S2 Cheat Sheet

Chapter	Usual types of questions	Tips	What can go ugly
1 – Binomial Distribution	<ul style="list-style-type: none"> <li>Finding the probability of a certain number of successes.</li> <li>Finding the probability of some range of successes, e.g. <math>P(1 &lt; X \leq 3)</math>, <math>P(X &gt; 5)</math>, <math>P(X \leq 2)</math></li> <li>Be able to list the assumptions made in order to model a scenario using a Binomial Distribution.</li> <li>Calculating mean = <math>np</math> and variance = <math>npq</math> (where <math>q = 1 - p</math>)</li> <li>Sneaky Geometric distribution questions (see right)</li> <li>Calculating an unknown value of <math>p</math> from context.</li> <li>Calculating an unknown value of <math>n</math> from context.</li> <li>Solve problems in which tables have to be used, but the probability of success is <math>&gt; 0.5</math> (i.e. not in table), by instead counting the number of failures.</li> <li>Solving “double inequalities”, e.g. “smallest value of <math>k</math> such that <math>P(X \geq k) \leq 0.1</math>”</li> </ul>	<ul style="list-style-type: none"> <li><math>P(X = x) = \binom{n}{x} p^x (1 - p)^{n-x}</math></li> <li>Remember the two edge cases where you don't need to use the full formula:               <ol style="list-style-type: none"> <li>0 successes: <math>(1 - p)^n</math></li> <li><math>n</math> successes: <math>p^n</math></li> <li>Thus “at least 1 success”: <math>1 - (1 - p)^n</math></li> </ol> </li> <li>For worded questions, always start your working by writing out your distribution, e.g. <math>X \sim B(20, 0.1)</math></li> <li>Sometimes you require the probability of a range when the cumulative table can't be used, because the value of <math>p</math> is not a nice number. This involves subtracting the opposite cases from 1: e.g. <math>P(X \geq 2) = 1 - P(X = 1) - P(X = 0)</math></li> <li>Remember your table requires <math>\leq</math>, so if you have <math>P(X &lt; 5)</math>, use <math>P(X \leq 4)</math>. See on the right regarding problems of ‘flipping’ your inequality.</li> <li>Assumptions of a Binomial Distribution:               <ol style="list-style-type: none"> <li>Fixed number of trials.</li> <li>Probability of success constant.</li> <li>Each trial is independent (ensure you put in context!)</li> <li>Each trial has two outcomes (‘success’ and ‘failure’)</li> </ol> </li> <li>When <math>p</math> is unknown: “An unfair coin with probability of heads <math>p</math> is tossed 20 times. The probability of seeing no heads is 0.1. Determine <math>p</math>”. Since this is an ‘edge case’: <math>P(X = 0) = (1 - p)^{20}</math>. Thus:               <math display="block">(1 - p)^{20} = 0.1</math> <math display="block">1 - p = \sqrt[20]{0.1}</math> <math display="block">p = 1 - \sqrt[20]{0.1} = 0.109</math> </li> <li>When <math>n</math> is unknown: “I play a game for which the probability of winning is 0.7. If I win every game, what is the smallest number of times I play such that the probability of winning every game is less than 0.01?” Again an edge case so: <math>P(X = n) = 0.7^n</math> <math display="block">0.7^n &lt; 0.01</math> <math display="block">n \log 0.7 &lt; \log 0.01</math> <math display="block">n &gt; \frac{\log 0.01}{\log 0.7} = 12.9</math> <p>Thus at least 13 games required. Notice that the direction of the inequality reversed because we divided by a negative number (<math>\log 0.07</math>).</p> </li> <li>Sometimes you'll get a part of a question which requires some non-Binomial probabilistic calculation, particularly involving some number of failures before a success is obtained, e.g. “Bob keeps firing arrows at a target until he gets a Bullseye. The probability he gets a Bullseye is 0.4. What's the probability he hits the Bullseye on the 4<sup>th</sup> shot”: <math>0.6^3 \times 0.4</math></li> </ul>	<ul style="list-style-type: none"> <li>Misreading terms like “at least” or “more than”. Make sure you get <math>&lt;</math> vs <math>\leq</math> right.</li> <li>Incorrectly ‘flipping’ a probability for <math>&gt;</math> and <math>\geq</math>, i.e. being one off when you like up a value in the cumulative Binomial table. <math>P(X &gt; 1) = 1 - P(X \leq 1)</math> <math>P(X \geq 1) = 1 - P(X = 0)</math> Just think logically about what the opposite of “more than 1” is and so on.</li> <li>Similarly, incorrect switching from the number of successes to the number of failures, usually by forgetting to replace the value with <math>n</math> minus it, or not preserving the strictness/non-strictness of the inequality.</li> </ul>

This is related to something called the Geometric Distribution which isn't formally covered in the syllabus.

- When switching from the number of successes  $X$  to the number of failures  $Y$  (so that the probability is less than 0.5 for the purposes of using tables), flip the inequality (but preserve  $<$  vs  $\leq$ ) and if the number of successes was  $k$ , use  $n - k$  for number of failures:

$$P(X < k) = P(Y > n - k)$$

$$P(X \geq k) = P(Y \leq n - k)$$

e.g. "In Joe's café 70% of customers buy a cup of tea. In a random sample of 20 customers find the probability that more than 15 buy a cup of tea."

$$X \sim B(20, 0.7)$$

$$\therefore Y \sim B(20, 0.3)$$

$$P(X > 15) = P(Y < 20 - 15)$$

$$= P(Y < 5)$$

$$= P(Y \leq 4)$$

- "Given  $X \sim B(50, 0.6)$ , find the smallest value of  $k$  such that  $P(X < k) \geq 0.9$ "

We can only use the table if the probability is less than 0.5. This question requires a great deal of care, particularly with the effect of switching from  $X$  to  $Y$  and getting  $<$  vs  $\leq$  right!

$$X \sim B(50, 0.6)$$

$$Y \sim B(50, 0.4)$$

$$P(X < k) = P(Y > 50 - k) \geq 0.9$$

$$1 - P(Y \leq 50 - k) \geq 0.9$$

$$P(Y \leq 50 - k) \leq 0.1$$

$$50 - k \leq 15$$

$$k \geq 35 \quad \therefore \quad k = 35$$

<p>2 – Poisson Distribution</p>	<ul style="list-style-type: none"> <li>• Be able to state the conditions under which a Poisson distribution may be used.</li> <li>• As with the Binomial Distribution, find the probability of a particular number of events happening within a given time frame, or within some range using tables.</li> <li>• State the mean and variance of a Poisson distribution.</li> <li>• Be able to scale a Poisson distribution to a different time period.</li> <li>• Feed the probability obtained from a Poisson distribution into a Binomial Distribution, and make subsequent calculations.</li> <li>• Approximate a Binomial Distribution using a Poisson Distributions, and be able to state the conditions under which we can do so.</li> <li>• Again, forming and solving 'double inequalities': e.g. "smallest value of <math>k</math> such that <math>P(X \geq k) \leq 0.1</math>"</li> </ul>	<ul style="list-style-type: none"> <li>• <math>P(X = x) = \frac{e^{-\lambda} \lambda^x}{x!}</math> While this is in the formula booklet, the easy way to remember it is that reading left to right and then down, the <math>\lambda</math> repeats consecutively, as does the <math>x</math>.</li> <li>• <math>E(X) = \lambda</math> and <math>Var(X) = \lambda</math>. The fact these are the same sometimes provides a justification for why a Poisson distribution would be suitable to model certain data. See 'Wordy Questions' page.</li> <li>• As with the Binomial Distribution, ensure you state the distribution for any wordy question, e.g. <math>X \sim Po(5)</math>.</li> <li>• Conditions required for Poisson: <ul style="list-style-type: none"> <li>a. Events occur independently (e.g. volcano eruptions might not be modelled using Poisson because volcano less likely to erupt immediately after previous eruption, thus eruptions not independent).</li> <li>b. Events occur singly in time.</li> <li>c. A fixed rate for which events occur.</li> </ul> </li> <li>• A very common occurrence is that you will need to scale to another time period. e.g. "A printer jams on average 0.3 times an hour. Find the probability over a 5 hour period the printer jams at least 4 times." <i>Just scaling the 3 times an hour to a 15 times every 5 hours:</i>  <math display="block">X \sim Po(0.3 \times 5) \rightarrow X \sim Po(1.5)</math> <math display="block">P(X \geq 4) = 1 - P(X \leq 3)</math> <math display="block">= 1 - 0.9344 = 0.0656</math> </li> <li>• Another common question is to feed the value calculated from a Poisson question, into a Binomial Distribution. e.g. "Defects occur in planks of wood with rate 0.5 per 100cm. If Bob buys 6 blanks each of length 100cm, find prob that fewer than 2 of planks contain at most 3 defects." <i>First find probability a plank of 100cm contains at most 3 defects:</i>  <math display="block">X \sim Po(5) \quad P(X \leq 3) = 0.2650</math> <i>Then feed into a Binomial Distribution:</i> <i>Let <math>Y</math> be the number planks with at most 3 defects. <math>Y \sim B(6, 0.2650)</math></i>  <math display="block">P(Y &lt; 2) = P(Y \leq 1) = P(Y = 0) + P(Y = 1)</math> <math display="block">= 0.735^6 + 6 \times 0.265 \times 0.735^5</math> <i>Notice that we couldn't use tables for the Binomial here because of the non-nice <math>p</math> value.</i> </li> <li>• As with the Binomial Distribution, we can get nasty 'double inequality' questions: "While a popcorn bag is in the microwave, an average of 5 pops can be heard per second. What's the minimum number of pops heard such that there's less than a 10% chance of hearing more than this number of pops?"  <math display="block">X \sim Po(5)</math> <math display="block">P(X &gt; k) &lt; 0.1</math> <math display="block">1 - P(X \leq k) &lt; 0.1</math> <math display="block">P(X \leq k) &gt; 0.9</math> <math display="block">k = 8</math> Note we had to take care with <math>X &gt; k</math> vs <math>X \geq k</math> and <math>&lt; 0.1</math> vs <math>\leq 0.1</math>. </li> <li>• For Binomial <math>\rightarrow</math> Poisson approximations, see notes on Normal Approximations.</li> </ul>	<ul style="list-style-type: none"> <li>• Since these questions are very similar to Binomial questions (except we use a different table and evaluate the probabilities differently), the same potential problems can arise, e.g. not correctly flipping an inequality from <math>P(X &gt; k)</math> to <math>1 - P(X \leq k)</math> and so on.</li> <li>• Not realising a Poisson question is a Poisson question! Remember that any mention of 'rate' implies Poisson rather than Binomial.</li> <li>• Also not realising when both a Poisson Distribution and a Binomial Distribution have to be used within the same question!</li> </ul>
---------------------------------	--	---	--

<p>3 – Continuous Random Variables</p>	<ul style="list-style-type: none"> <li>Find the probability of some range of values given a probability density function or cumulative distribution function. e.g. <math>P(X &gt; 2)</math>, <math>P(1 &lt; X &lt; 2)</math></li> <li>Appreciate that <math>P(X = k) = 0</math> if <math>X</math> is a continuous random variable.</li> <li>State the probability density function given a graph.</li> <li>Comment on the skew of a distribution.</li> <li>Calculate <math>E(X)</math>, <math>Var(X)</math>, the median/quartiles and the mode of a probability density function.</li> <li>Convert from <math>f(x)</math> to <math>F(x)</math> and vice versa, potentially involving multiple ranges.</li> <li>Be able to calculate <math>E(X^2)</math> (surprisingly common!)</li> <li>Be able to calculate the median and quartiles.</li> <li>Be able to calculate the mode, either by finding the turning point, or by inspection of the graph.</li> </ul>	<ul style="list-style-type: none"> <li>Remember that <math>f(x)</math> is the probability density function for continuous variables, and that this is not a ‘probability’ as such: we only get a probability when we integrate <math>f(x)</math> over some range.</li> <li>Relatedly, if <math>X</math> is a continuous variable, then <math>P(X = k) = 0</math> because the probability of a specific value is infinitely small (e.g. no one has an ‘exact’ height of 1.5m).</li> <li>Think of <math>F(x)</math> as “the running total of the probability up to <math>x</math>”.</li> </ul> $F(x) = P(X < x)$ <ul style="list-style-type: none"> <li>To find a probability over a range:</li> </ul> $P(a < X < b) = \int_a^b f(x) dx$ $P(X > a) = \int_a^{\infty} f(x) dx$ <p>For the latter, if you know the probability is 0 after some value <math>b</math> (this will always be the case in exams), we can use <math>b</math> instead of <math>\infty</math>.</p> <ul style="list-style-type: none"> <li>If <math>F(x)</math> is known, then you could calculate say <math>P(X &gt; 2)</math> using <math>P(X &gt; 2) = 1 - P(X \leq 2) = 1 - F(2)</math></li> <li>When asked to find the value of some constant <math>k</math> used in a p.d.f., use the fact that <math>\int_{-\infty}^{\infty} f(x) dx = 1</math>, i.e. the area under the whole probability function is 1.</li> <li>However, if the cumulative distribution function is given, then there is no need to integrate, just use <math>F(b) = 1</math> where <math>b</math> is the highest possible value (since <math>P(X \leq b) = 1</math>).</li> <li>Remember that it doesn’t matter for continuous variables if you use <math>X &gt; k</math> or <math>X \geq k</math> (but it does matter for discrete variables!).</li> <li><math>E(X) = \int_{-\infty}^{\infty} x f(x) dx</math></li> <li><math>Var(X) = E(X^2) - E(X)^2</math></li> </ul> $= \int_{-\infty}^{\infty} x^2 f(x) dx - \mu^2$ <ul style="list-style-type: none"> <li>To go from <math>f(x)</math> to <math>F(x)</math>, integrate. e.g.</li> </ul> $f(x) = \begin{cases} \frac{2}{3}x & 1 < x \leq 2 \\ 0 & \text{otherwise} \end{cases}$ <p>Then in the <math>1 &lt; x \leq 2</math> range:</p> $F(x) = \int_1^x \frac{2}{3}t dt = \left[ \frac{1}{3}t^2 \right]_1^x = \frac{1}{3}x^2 - \frac{1}{3}$ <p>Therefore the full cumulative distribution function is:</p> $F(x) = \begin{cases} 0 & x \leq 1 \\ \frac{1}{3}x^2 - \frac{1}{3} & 1 < x \leq 2 \\ 1 & x > 2 \end{cases}$ <p>Note the extra ‘1’ row required since the running total of the probability by the time you get to 2 is 1. Note that the use of <math>t</math> was to avoid a clash with the <math>x</math> used as the upper limit of the integral. But the mark scheme permits <math>\int_2^x \frac{2}{3}x dx</math>, so do this way if you find it less confusing.</p>	<ul style="list-style-type: none"> <li>A common error is doing <math>P(X \geq 10) = 1 - P(X \leq 9)</math> for continuous variables. This is true for discrete variables, but the opposite of “more than 10m tall” is not “under 9m tall”!</li> </ul> $P(X \geq 10) = 1 - P(X < 10)$ <p>Note that the <math>&lt;</math> vs <math>\leq</math> does not matter since <math>X</math> is continuous.</p> <ul style="list-style-type: none"> <li>Note paying attention to whether the question gives the cumulative distribution function <math>F(x)</math> or the probability distribution function <math>f(x)</math>. This will completely change the approach to use for answering a question! If you’re finding the median/quartiles and you’re already given <math>F(x)</math>, don’t integrate!</li> <li>When finding the cumulative distribution function, forgetting the rows for the two ‘ends’ of the ranges. You should have one more row in <math>F(x)</math> than <math>f(x)</math>. The same applies when going from <math>F(x)</math> to <math>f(x)</math>: you should have one less row.</li> <li>When finding the cumulative distribution function, then in the example on the left, it might have been tempted to go straight from <math>\frac{2}{3}x</math> to</li> </ul>
--	---	---	---

When  $f(x)$  has multiple rows, see the note on the right about ensuring you add the running total up to that range.

- To go from  $F(x)$  to  $f(x)$  just differentiate. Don't forget that the '1' row disappears.
- To find median or quartiles: use  $F(Q_1) = 0.25, F(Q_2) = 0.5, F(Q_3) = 0.75$ . Either  $F(x)$  will already be given, or you will have to determine it from  $f(x)$  first.

Sometimes you have to determine which range the quartile/median occurs in first by evaluating the borderline values (although this is not necessary if you only have one range).

e.g. If:

$$F(x) = \begin{cases} 0 & x < 0 \\ \frac{1}{4}x^2 & 0 \leq x \leq 1 \\ \frac{1}{20}x^4 + \frac{1}{5} & 1 < x \leq 2 \\ 1 & x > 2 \end{cases}$$

and we wished to find the median  $Q_2$ , then it might be in the  $0 \leq x < 1$  range or  $1 < x \leq 2$  range. However  $F(1) = \frac{1}{4}$ , i.e. the running total of the probability up to 1 is 0.25, thus the median wouldn't have yet occurred, and thus it's in the  $1 < x \leq 2$  range.

Then using  $F(Q_2) = 0.5$ :

$$\frac{1}{20}x^4 + \frac{1}{5} = \frac{1}{2}$$

and so on.

- The mode is the value of  $x$  such that  $f(x)$  is at its maximum. The mode can be calculated in two different ways: (and usually you can only use one of the two)
  - For curved graphs, finding the turning points using

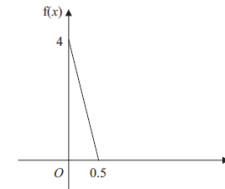
$$\frac{df(x)}{dx} = 0$$

- Using the graph. e.g. If  $f(x) = \begin{cases} 1 - \frac{1}{2}x & 0 < x < 2 \\ 0 & \text{otherwise} \end{cases}$ , then we can see from the sketch that the probability is greatest when  $x = 0$ , so the mode is 0.

- If asked to find the probability density function of a given graph, ensure you don't just give the equation of the line you see: you need to use the full curly brace construction covering all values:

Then  $f(x) = 4 - 8x$  is not enough, we need to write:

$$f(x) = \begin{cases} 4 - 8x & 0 \leq x \leq 4 \\ 0 & \text{otherwise} \end{cases}$$



- When asked to find skew of a probability density function, if two of the mode, median and mean have been found, compare in the usual S1 way (remembering if positive skew that  $mean > median > mode$ ). If you only have the graph look at the shape: a 'positive tail' mean positive skew.

$\frac{1}{3}x^2$  without properly evaluating  $\int_1^x \frac{2}{3}t dt$

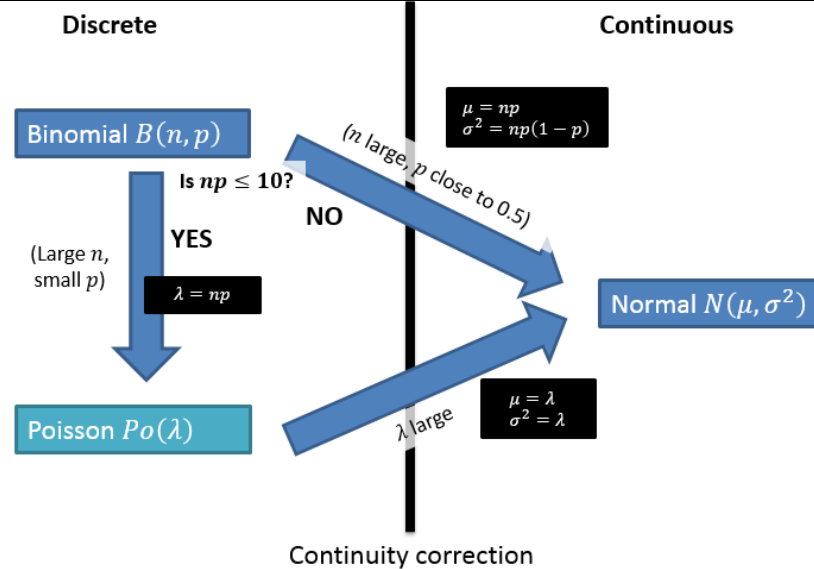
This will result in a missing constant.

- When finding the cumulative distribution function from  $f(x)$  where there are multiple ranges, forgetting to add on the running total up to the start of the range being considered. e.g. If you had ranges  $0 < x \leq 1$  and  $1 < x \leq 2$  in  $f(x)$ , then when finding  $F(x)$  in the latter range,  $F(x) = F(1) + \int_1^x f(t) dt$ . This is because you want the area up to 1 and then the area between 1 and  $x$ . i.e. Don't forget the  $F(1)$ !
- When finding the mode, accidentally giving the probability density of the mode as the answer rather than the mode itself (e.g. Jan 2011 Q5d: answer is 0 not 4).

<p>4 – Continuous Uniform Distribution</p>	<ul style="list-style-type: none"> <li>Find the probability of a range for a continuous uniform distribution, e.g. <math>P(2 &lt; X &lt; 3)</math> or <math>P(X &gt; 3)</math>.</li> <li>Sketch the probability function of a continuous uniform distribution.</li> <li>Find the mean or variance of a continuous uniform distribution:</li> <li>Find the <math>a</math> and <math>b</math> of <math>X \sim U(a, b)</math> when the mean and variance of the distribution is given.</li> <li>Be able to calculate <math>E(X^2)</math></li> <li>The probability calculated from a uniform distribution may be fed into a Binomial distribution.</li> </ul>	<ul style="list-style-type: none"> <li>If <math>X \sim U(a, b)</math> then:           <math display="block">E(X) = \frac{a + b}{2}</math> <math display="block">Var(X) = \frac{(b - a)^2}{12}</math> <p>I remember the variance as “a twelfth of the squared difference”.</p> </li> <li>Suppose <math>E(X) = 2</math> and <math>Var(X) = 27</math> and <math>X \sim U(a, b)</math> where <math>a</math> and <math>b</math> are unknown. Then:           <math display="block">\frac{a + b}{2} = 2 \rightarrow a + b = 4</math> <math display="block">\frac{(a - b)^2}{12} = 27 \rightarrow a - b = 18</math> <p>We can then solve these simultaneously.</p> </li> <li>The key is just remembering the area of the rectangle (when you sketch the p.d.f.) is 1. Therefore if <math>X \sim U(3, 5)</math> then since the width is 2, the height is clearly <math>\frac{1}{2}</math>. We would specify the probability density function as:           <math display="block">f(x) = \begin{cases} \frac{1}{2} &amp; 3 \leq x \leq 5 \\ 0 &amp; \text{otherwise} \end{cases}</math> </li> <li>For the example above, we could then find <math>P(X &gt; 4.2)</math> by just considering the rectangular area involved:           <math display="block">P(X &gt; 4.2) = 0.8 \times 0.5 = 0.4</math> </li> <li><math>E(X^2) = \int_{-\infty}^{\infty} x^2 f(x) dx</math>            For the above example, <math>E(X^2) = \int_3^5 \frac{1}{2} x^2 dx = \left[ \frac{1}{6} x^3 \right]_3^5 = \frac{49}{3}</math> </li> <li>“I pick 10 real numbers randomly from 12 to 17. Find the probability that at least 5 of these numbers are greater than 15.5.”            If <math>X</math> is the number picked each time, <math>P(X &gt; 15.5) = (17 - 15.5) \times \frac{1}{5} = 0.3</math>. Then if <math>Y</math> is the number of times a number greater than 15.5 was picked, <math>Y \sim B(10, 0.3)</math>, and calculate <math>P(Y \geq 5)</math>.         </li> <li>Don't forget your rules of coding from S1:           <math display="block">Var(7X) = 49Var(X)</math> <math display="block">Var(2 - 5X) = 25Var(X)</math> </li> </ul>	<ul style="list-style-type: none"> <li>Suppose that <math>X \sim U(3, 5)</math>. Then what is <math>P(4.2 &lt; X &lt; 6)</math>? You might be tempted to calculate <math>(6 - 4.2) \times \frac{1}{2}</math>, but the probability above a value of 5 is 0, thus:           <math display="block">P(4.2 &lt; X &lt; 6) = P(4.2 &lt; X &lt; 5) = (5 - 4.2) \times \frac{1}{2}</math>           i.e. we ‘truncate’ any part of the range which is outside the range of the uniform distribution.         </li> </ul>
--	---	---	---

5 – Normal Approximations

- Be able to approximate a Binomial Distribution using a Normal Distribution.
- Be able to approximate a Poisson Distribution using a Normal Distribution.
- Be able to approximate a Binomial Distribution using a Poisson Distribution.
- Be able to give the conditions under which such approximations can be made.
- Justify why we need continuity corrections.



- This diagram may seem like quite a lot to memorise, but all you need to memorise for carrying out the majority of approximations is this: If you have a Binomial Distribution, is  $np \leq 10$ ? If yes use Poisson, else use Normal. In terms of converting between the distributions, the mean and variance of the Normal/Poisson approximation is just the mean and variance of the original distribution.
- If asked why a continuity correction is needed (and suppose the original distribution is Poisson), say: "Poisson is discrete, but Normal is continuous".
- For continuity corrections, we want to go from a discrete  $X$  to a continuous version of it  $Y$ . You will never get a continuity correction wrong if you carry out these two simple steps:
  - Make sure your inequality uses  $\leq$  or  $\geq$  instead of  $<$  or  $>$ . i.e. Ensure inequality is non-strict.
  - 'Extend' your range by 0.5 at each end. i.e. If you visualise your inequality as a line on the number line, it should be 0.5 longer each end.  
 Examples:  $P(X < 4) = P(X \leq 5) = P(Y \leq 5.5)$   
 $P(3 < X \leq 5) = P(4 \leq X \leq 5) = P(3.5 \leq Y \leq 5.5)$   
 $P(X = 5) = P(4.5 \leq Y \leq 5.5)$
- I prefer to do the continuity correction immediately, i.e. before you either reverse the direction of the inequality or standardise. e.g.  
 $P(X > 5) = P(X \geq 6) = P(Y \geq 5.5) = 1 - P(Y < 5.5) = 1 - P(Z < \dots)$
- The number of marks effectively tells you what approximation you are using: If at least 6 marks, it's a Normal Approximation (because of the many steps of converting the distribution, standardising and continuity corrections), otherwise it's Binomial  $\rightarrow$  Poisson.

- If asked to give the conditions under which a Binomial Distribution can be approximated using a Poisson Distribution, do NOT say  $np \leq 10$ . This condition is a rule of thumb only: the actual condition is "n is large, p is small" (from which  $np \leq 10$  stems).
- All manner of things can go wrong with continuity corrections. This might be forgetting to convert to  $\geq$  or  $\leq$  first (i.e. incorrectly going from  $P(X < 4)$  to  $P(Y < 4.5)$ ), or making your range 0.5 smaller rather than 0.5 larger, e.g. incorrectly from  $P(X \leq 4)$  to  $P(Y \leq 3.5)$ .
- See note on the left about the perils of scaling the value instead of  $\lambda$  in the case of the Poisson Distribution.
- In the formula for  $Z$ , accidentally dividing by the variance rather than the standard deviation.

		<ul style="list-style-type: none"> <li>• Example Normal Approximation: “The number of houses sold by an estate agent follows a Poisson distribution, with a mean of 2 per week. The estate agent will receive a bonus if he sells more than 25 houses in the next 10 weeks. Use a suitable approximation to estimate the probability that the estate agent receives a bonus.” <ul style="list-style-type: none"> <li>a. Note first that you might be tempted to scale the 25 houses in 10 weeks to 5 houses in 2 weeks and stick with the original <math>\lambda = 2</math>. The catastrophic flaw in doing this is that the continuity correction affects the range differently depending on whether you’re using the original or scaled value. If not scaling: <math>P(X &gt; 25) = P(X \geq 26) = P(Y \geq 25.5)</math> If scaling: <math>P(X &gt; 5) = P(X \geq 6) = P(Y \geq 5.5)</math> In the latter incorrect case the 0.5 has a greater effect on the smaller value of 6 compared with the larger value of 26, so the probability will be too high.</li> <li>b. Step 1: Determine what approximation to use. In this example we have a Poisson Distribution, which always goes to Normal. If it were Binomial, you’d first determine if <math>np \leq 10</math>.</li> <li>c. Step 2: Identify original distribution. As discussed, we scale <math>\lambda</math> (rather than the 25), so: <math>X \sim Po(20)</math></li> <li>d. Step 3: Write the approximation, potentially with reference to a new continuous variable <math>Y</math> which is the continuous version of <math>X</math>, i.e.: <math>Y \sim N(20,20)</math> As discussed, use the mean and variance of the original distribution.</li> <li>e. Step 4: If necessary, carry out continuity correction to get a probability in terms of <math>Y</math>: <math>P(X &gt; 25) = P(X \geq 26) = P(Y \geq 25.5)</math></li> <li>f. Step 5: Use your S1 knowledge and find the probability by first standardising. Don’t forget that you’re dividing by the standard deviation, not the variance: <math display="block">P(Y \geq 25.5) = P\left(Z \geq \frac{25.5 - 20}{\sqrt{20}}\right) = P(Z \geq 1.23)</math> <math display="block">= 1 - P(Z \leq 1.23) = 1 - 0.8907 = 0.1093</math></li> </ul> </li> </ul>	
6 – Populations and Samples	<ul style="list-style-type: none"> <li>• Be able to define key terms such as ‘sampling distribution’, ‘sampling frame’, ‘population’, ‘sample’, ‘statistic’.</li> <li>• Be able to describe what the population is or sampling frame is given the context, and identify reasons for differences between the two.</li> <li>• Be able to list possible samples.</li> <li>• Be able to calculate the</li> </ul>	<ul style="list-style-type: none"> <li>• Key definitions: <ul style="list-style-type: none"> <li>a. <b>Statistic:</b> “A random variable (1) which is some function of the sample and not dependent on any population parameters (1)” - I think the ‘random variable’ bit is a bit pernicious (as does Wikipedia), but c’est la vie! If 1 mark, the second part is important.</li> <li>b. <b>‘Population’:</b> The collection of all items.</li> <li>c. <b>‘Sample’:</b> Some subset of the population which is intended to be representative of the population.</li> <li>d. <b>‘Census’:</b> When the entire population is sampled.</li> <li>e. <b>‘Sampling unit’:</b> Individual member or element of the population or sampling frame.</li> <li>f. <b>‘Sampling frame’:</b> A list of all sampling units or all the population.</li> <li>g. <b>Sampling distribution:</b> All possible samples are chosen from a population (1); the values of a statistic and the associated probabilities is a sampling distribution (1).</li> </ul> </li> <li>• It’s important you get your head around what the sampling distribution actually is: It gives the distribution over possible values of the statistic as we take different samples. So if for example</li> </ul>	<ul style="list-style-type: none"> <li>• When finding the sampling distribution, forgetting that different orderings are different possibilities, e.g. (1,1,2) and (1,2,1) should both be considered.</li> <li>• Even though you can subtract from 1 to find the last probability in your sampling distribution, if you have time, you might want to calculate it ‘the long way’ to check your answer, as probabilities</li> </ul>



	<p>sampling distribution for a variety of statistics, such as median (very common!), range, maximum and mode.</p> <ul style="list-style-type: none"> <li>Be able to identify when the sampling distribution is a Binomial Distribution or otherwise, and specify this distribution.</li> </ul>	<p>the statistic was the 'range' of the sample, then this range is likely to vary as we take different samples. As these ranges <u>vary across samples</u>, it forms a distribution.</p> <ul style="list-style-type: none"> <li>The sampling frame is the list of things in the population that are available for sampling, e.g. "The ID numbers", "The list of car registration numbers". The mark scheme seems to particularly like it when you refer to some identifying property of the things in the sampling frame. The sampling frame may be different from the population, because some things in the population may not be available for sampling. e.g. If sampling people who've visited a medical practice, "some people may have left the area but hadn't deregistered".</li> <li>When listing outcomes, it helps to be systematic in listing them so you don't miss any. Note that <u>different orderings count as distinct possibilities</u>. e.g. "You have a large collection of 1p, 2p and 5p coins, and take 3 coins. Find all samples in which the maximum is 5." We may want to first list the possibilities where 5 appears once, 5 appears twice, and so on..." (5,1,1), (1,5,1), (1,1,5), (5,1,2), (5,2,1), (1,5,2), (2,5,1), (1,2,5), (2,1,5), (5,2,2), (2,5,2), (2,2,5), (5,5,1), (5,1,5), (1,5,5), (5,5,2), (5,2,5), (2,5,5), (5,5,5)</li> <li>e.g. "A bag contains a large number of 1p and 2p coin, of which 40% are 1p and 60% are 2p. A sample of 2 coins. Find the sampling distribution of the sample maximum." When finding the sampling distribution, it may help to have a table as follows to organise your working, such that the outcomes for each possible value of the statistic are grouped:</li> </ul> <table border="1" data-bbox="741 683 1794 786"> <thead> <tr> <th>Possibilities</th> <th>Statistic (Maximum)</th> <th>Probability</th> </tr> </thead> <tbody> <tr> <td>(1,1)</td> <td>1</td> <td><math>0.4^2 = 0.16</math></td> </tr> <tr> <td>(1,2), (2,1), (2,2)</td> <td>2</td> <td><math>1 - 0.16 = 0.84</math></td> </tr> </tbody> </table> <p>Notice that we didn't need to do any complicated calculation for the last probability, because it was just 1 minus the others! Had we had to calculate it fully, then <math>2 \times 0.4 \times 0.6 + 0.6^2 = 0.84</math></p> <ul style="list-style-type: none"> <li>On the rare occasion you get a question asking for the sampling distribution, where you don't actually have to do any calculation, but just have to consider what well-known distribution you get as the sample varies: <i>"A factory produces components. Each component has a unique identity number and it is assumed that 2% of the components are faulty. On a particular day, a quality control manager wishes to take a random sample of 50 components. A statistic <math>F</math> represents the number of faulty components in the sample. Specifying the sampling distribution of <math>F</math>."</i> We know a sampling distribution is the possible values of the statistic as we take different samples of 50 light bulbs. If the statistic is the count of light bulbs, we can see this count varies Binomially between 0 and 50. Thus <math>F \sim B(50, 0.02)</math></li> </ul>	Possibilities	Statistic (Maximum)	Probability	(1,1)	1	$0.4^2 = 0.16$	(1,2), (2,1), (2,2)	2	$1 - 0.16 = 0.84$	<p>should obviously all add up to 1.</p>
Possibilities	Statistic (Maximum)	Probability										
(1,1)	1	$0.4^2 = 0.16$										
(1,2), (2,1), (2,2)	2	$1 - 0.16 = 0.84$										

<p>7 – Hypothesis Testing</p>	<ul style="list-style-type: none"> <li>• Be able to define key terms such as ‘critical region’, ‘hypothesis test’, ‘significance level of a test’.</li> <li>• Determine a critical region (one or two ranges depending on whether one or two tailed)</li> <li>• Understand when we need the probability to be strictly within the significance level, and when to be close to it as possible either side.</li> <li>• Be able to calculate the ‘actual level of significance’.</li> <li>• Carry out a hypothesis test involving Binomial and Poisson distributions.</li> <li>• Carry out a hypothesis test where a normal approximation is required.</li> </ul>	<ul style="list-style-type: none"> <li>• Key terms whose definitions you need to remember: <ul style="list-style-type: none"> <li>a. <b>Critical Region:</b> The <u>range of values</u> such that <u>the null hypothesis is rejected</u>.</li> <li>b. <b>Hypothesis Test:</b> a procedure to examine a value of a population parameter proposed by the null hypothesis</li> <li>c. <b>Significance level:</b> “the probability of rejecting <math>H_0</math> if <math>H_0</math> is true”, or “the probability of incorrectly rejecting <math>H_0</math>”.</li> </ul> </li> <li>• The “or more extreme” thing often confuses people: is “more extreme” below the value or above it? We can always tell this by what side of the mean we’re on. If <math>\lambda = 7</math> and we’re interested in seeing 10 hits to a website “or more extreme”, then since 10 is above the mean of 7, clearly we want “10 or above” to get the tail. Ensure “10 or above” is <math>X \geq 10</math> not <math>X &gt; 10</math>.</li> <li>• <b>Identifying the critical values from a table:</b>  Firstly, note if the test is one-tailed (involving <math>&lt;</math> or <math>&gt;</math>) or two-tailed (involving <math>\neq</math>), as you need to halve the significance level each side in the latter case.  At the left tail, find the closest value under the significance level (or half it) in the table, or if you are explicitly told to use the closest value to the significance level, do that. This is the lower end of the critical region.  However, at the right tail, first find the closest value above it (or again if explicitly told, the closest value either side), but then <u>go one above it</u>. This is the critical value.   Example: <i>Under the null hypothesis <math>X \sim B(10, 0.35)</math>. For a two-tailed test with significance level 5%, what are the critical values?</i>  Looking at the tables, at left end <math>X = 0</math>, and at right end, closest value with probability above 0.975 is 7, thus going one above <math>X = 8</math>.   If we had been asked for the ‘closest value’ to 0.025 and 0.975, then we’d then get <math>X = 0</math> and <math>X = 7</math> instead.  <u>It is vitally important you specify the probability of being in each tail</u> to evidence that you have used the table, e.g. “<math>P(X \geq 8) = 0.9952</math>”</li> <li>• For the critical region, don’t forget to provide a lower limit or upper limit in the case of the Binomial Distribution, as the outcomes are finite.  For the previous example: <math>X = 0</math> or <math>8 \leq X \leq 10</math>. Mark schemes usually condone the lack of <math>\leq 10</math>, but don’t take any chances.</li> <li>• The <u>actual level of significance</u> is the actual probability of being in the critical region. You should have already written out the probabilities of being in each part of the critical region, so it’s then just a case of adding the two probabilities.</li> <li>• The mark scheme for a hypothesis test without a normal approximation is as follows: <ul style="list-style-type: none"> <li>a. Specifying <math>H_0</math> and <math>H_1</math> (1 mark)</li> <li>b. Specifying the distribution for <math>X</math> under the null hypothesis, e.g. <math>X \sim B(10, 0.2)</math> (1 mark)  The <math>p</math> or <math>\lambda</math> will be your population parameter under the null hypothesis.</li> <li>c. 2 marks for either: Determining the probability of the observed value or more extreme (e.g. <math>P(X \geq 5) = 1 - P(X \leq 4) = \dots</math> <u>or</u> determining the critical region.</li> </ul> </li> </ul>	<ul style="list-style-type: none"> <li>• Forgetting to halve the significance level for two-tailed tests.</li> <li>• Being one off the critical value (particularly at the right tail). See notes on the left. Pay careful attention to whether it says “use the closest value”, or not.</li> <li>• Getting your continuity correction wrong (see notes for Chapter 5).</li> </ul>
-------------------------------	--	--	--

		<p>d. Using your probability to state whether <math>H_0</math> is rejected or not, ensuring you directly compare your probability with significance level. e.g. "0.0723 &gt; 0.05, so not significant (<math>H_0</math> is not rejected)" (1 mark)</p> <p>e. Put this conclusion in context. "Bob is not justified in his claim that the rate of flamingo attacks has increased."</p> <ul style="list-style-type: none"> <li>• The mark scheme for a hypothesis test with a normal approximation is usually broken down as the following: <ul style="list-style-type: none"> <li>a. Specifying <math>H_0</math> and <math>H_1</math> (1 mark)</li> <li>b. Possibly a mark for specifying your distribution of <math>X</math>, i.e. <math>X \sim B(\dots)</math> or <math>X \sim Po(\dots)</math></li> <li>c. Specifying the distribution of <math>Y</math> for your normal approximation, <math>Y \sim N(\dots, \dots)</math> (1 mark)</li> <li>d. Doing the continuity correction: <math>P(X \geq 40) = P(Y \geq 39.5)</math> (1 mark)</li> <li>e. Standardising to get <math>P(Z \dots)</math>. Note at this stage <math>&gt;</math> vs <math>\geq</math> does not matter as variable is continuous. (1 mark)</li> <li>f. As above, 2 marks for your two-part conclusion.</li> </ul> </li> <li>• Note that it is possible for the actual level of significance to be greater than the level of significance, if you were asked to find the closest value to 0.025/0.975, etc rather than those strictly below/above these values.</li> </ul>	
--	--	--	--

## Wordy/interpretation questions:

- Definitions (see the respective chapter notes above):  
*Statistic, population, census, sampling frame, sampling unit, sampling distribution, critical region, hypothesis test, significance level.*
- Modelling assumptions:
  - “List the assumptions made when modelling as a Binomial distribution”. See Chapter 1.
  - “State two conditions under which a Poisson distribution is a suitable model to use in statistical work”. See Chapter 2.
- (Surprisingly common!) “Describe the skewness of  $X$ , giving a reason for your answer.” *The mode and mean had previously been calculated in this question so “Positive skew as mean > mode”. 1 mark for skew type, 1 for reason.*
- Approximation related:
  - “Write down the conditions under which the Poisson distribution can be used as an approximation to the Binomial distribution.”  
 $\lambda$  is large
  - “Write down the two conditions needed to approximate the Binomial distribution by the Poisson distribution.”  
 $n$  is large,  $p$  is small (NOT  $np \leq 10$ )
  - “Write down which of the approximations used in part (a) is the most accurate estimate of the probability [One was Poisson, one Binomial]. You must give a reason for your answer.”  
*Normal approximation (1) because  $n$  is large and  $p$  close to half (1).  
OR Normal because mean  $\neq$  variance*
- “State the probability of incorrectly rejecting  $H_0$  using this critical region.”  
*Add the probabilities of your critical region(s).*
- [Given that  $X$  is a continuous variable] “Determine  $P(X = 3)$ ”. *Answer is 0.*
- “Identify a sampling frame.”  
*See notes on Chapter 6. Remember the mark scheme likes reference to an ‘identifier’, e.g. “list of unique identification numbers of the cookers.”*
- [In context of cookers being tested] “Give one reason, other than to save time and cost, why a sample is taken rather than a census.”  
*There would be no cookers left to sell (i.e. the idea that testing something in a sample destroys it/makes it unsellable).*

- “Identify the sampling units”. (*Using above example*) *A cooker*
- “A researcher took a sample of 100 voters from a certain town and asked them who they would vote for in an election. The proportion who said they would vote for Dr Smith was 35%. State the population and the statistic in this case. What do you understand by the sampling distribution of this statistic.”  
*Population is the residents of the town. Statistic is percentage/proportion who vote for Dr Smith. Sampling distribution here is number of people who voted for Dr Smith in all possible samples of 100.*
- [List of possible statistics given] “State, giving a reason which of the following is not a statistic based on this sample.”  
*The one involving the  $\mu$  and  $\sigma$ , because these are population parameters, whereas a statistic must be based on the sample only.*
- “Suggest a suitable model to describe the number of vehicles passing the fixed point in a 15 s interval.”  $X \sim Po(9)$  (*depending on what the average rate in the question is, which often you have to scale*)
- [Using distribution pictured] “state, giving your reason, whether  $E(X) < 3$ ,  $E(X) = 3$  or  $E(X) > 3$ .”  
 *$E(X) < 3$  because the graph tails to the left (i.e. negative skew).*
- “Find, to 2 decimal places, the value of  $k$  so that  $P(\mu - k\sigma < X < \mu + k\sigma) = 0.5$ .” *This means within  $k$  standard deviations either side of the mean, so by definition, we want  $P(-k < Z < k)$ . Since this is middle 50%, upper end is upper quartile so find when  $P(Z = k) = 0.75$ .*
- “Comment on this finding in the light of your critical region found in part (a).” *11 is in the critical region (1 mark) therefore there is evidence of a change/increase in the proportion/number of customers buying single tins (1 mark).*
- [Mean and variance of some data was calculated] “Explain how the answers from part (c) support the choice of a Poisson distribution as a model.”  
*For a Poisson model, Mean = Variance ; For these data  $3.69 \approx 3.73$*

