



Teacher Support Materials 2009

Statistics GCE

Paper Reference SS04

Copyright © 2009 AQA and its licensors. All rights reserved.

Permission to reproduce all copyrighted material has been applied for. In some cases, efforts to contact copyright holders have been unsuccessful and AQA will be happy to rectify any omissions if notified.

The Assessment and Qualifications Alliance (AQA) is a company limited by guarantee registered in England and Wales (company number 3644723) and a registered charity (registered charity number 1073334). Registered address: AQA, Devas Street, Manchester M15 6EX.
Dr Michael Cresswell, Director General.

Question 1

An air ambulance is based at a hospital. There is a probability of 0.005 that, on any one day, there is a delay in the air ambulance leaving to deal with an emergency.

- (a) Assuming that the probability of a delay is independent from day to day, specify a distribution that might be used to model the number of days on which a delay occurs during a period of 200 days. (1 mark)
- (b) Use a distributional approximation to find the probability that, during a period of 200 days, there are three or more days on which a delay occurs. (3 marks)
- (c) Explain why the distributional approximation that you used in part (b) is appropriate. (2 marks)

Student Response

1		Leave blank
	a) $X \sim B(200, 0.005)$ ✓	1
	b) $X \sim B(200, 0.005) \rightarrow X \sim P(1)$ ✓	
	$200 \times 0.005 = 1$	
	PROBABILITY $X \geq 3$	
	$= 1 - \sqrt{2}$ ✓	3
	$= 1 - 0.91969$	
	$= 0.0803$ ✓	
	c) n is large ✓ and p is very small ✓	
	$n = 200$ $p = 0.005$	
	$np \leq 10$	2
	\therefore do binomial approximation to poisson distribution. ✓	(6)

Commentary

This turned out to be a "love it or loathe it" starter question. Many candidates failed to identify the correct binomial model so made little progress. A common error for those who started off correctly was to use a normal approximation which was quite inappropriate in this case. The example shows an ideal solution. This candidate provided all the necessary method and explanation without excessive detail and gave the final answer in part (b) to the three significant figures required by the rubric.

Mark scheme

Q	Solution	Marks	Comments
1	$X =$ number of days out of 200 when delay occurs.		
(a)	$X \sim B(200, 0.005)$	B1	
(b)	$B(200, 0.005) \approx \text{Po}(1)$	B1	
	$P(X \geq 3) = 1 - P(X \leq 2)$	M1	
	$= 1 - 0.9197 = 0.0803$	A1	awrt
(c)	Poisson approximation to binomial with large n and very small p .	E1 E1	
		6	

Question 2

Students on an environmental science course are investigating nitrate pollution in a river in an agricultural region. The level of pollution becomes a concern when the mean concentration of nitrate exceeds 30 milligrams per litre of water.

The river is divided into a large number of sections of equal length.

- (a) One student takes samples of water at 8 randomly chosen locations along one of these sections and analyses the samples for nitrate concentration. Her results, in milligrams of nitrate per litre of water, are:

30 34 34 37 28 30 34 35

Carry out a test to investigate whether the nitrate pollution in this section of the river is a cause for concern. Assume that the data are drawn from a normal population and use the 1% significance level. (8 marks)

- (b) The students carry out similar investigations to that in part (a) on 42 sections. Their tests indicate that the mean concentration of nitrate exceeds 30 milligrams per litre of water in 16 sections.

- (i) Carry out a test, at the 1% significance level, to determine whether the level of nitrate concentration is a cause for concern in less than 60 per cent of sections of this river. (7 marks)
- (ii) State **one** assumption that must be made for your conclusion in part (b)(i) to be valid. (1 mark)

Student Response

Leave blank

2) $M = 30 \text{ mg per litre}$ ✓

a) $H_0: M = 30 \text{ mg}$ ✓ $\sigma = 7$ ✓ $B1$ ✓

$H_1: M > 30 \text{ mg}$ ✓ $n = 8$ $S = 3.0589$ ✓

$\bar{x} = 32.75$ ✓ $B1$ 1% $s.l.$ ✓

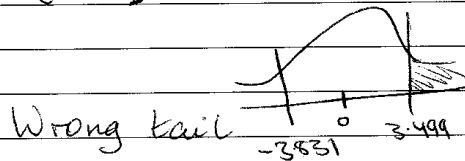
T.S.

$$z = \frac{30 - 32.75}{\frac{3.0589}{\sqrt{8}}} = -3.831 \times$$

$M1 \text{ ml}$

C.V.

$$t = 3.499 \times$$



accept H_0 and conclude that the nitrate pollution is not a cause for concern

as it ~~does not exceed 30mg~~, there is not significant evidence at the 1% level to suggest it exceeds 30mg. $A0$

5

b) $H_0: p = 0.6$ ✓

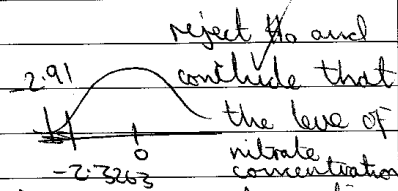
$H_1: p < 0.6$ ✓ $1 \text{ tail @ } 1\%$ $p = \frac{16}{42} = 0.38$ ✓

T.S.

$$z = \frac{0.38 - 0.6}{\sqrt{\frac{0.6 \times 0.4}{42}}} = -2.91$$

C.V.

$$z = -2.3263$$



reject H_0 and conclude that the level of nitrate concentration

is a cause for concern in less than 60% of the sections.

7

Commentary

There were some excellent solutions to this question with many candidates scoring full or nearly full marks on both hypothesis tests. Many errors were simply due to carelessness, but there were some more fundamental mistakes at the interpretation stage. Part (a) of this example illustrates both types of error. The candidate recognised the need for a t-statistic with 7 degrees of freedom, but obtained the wrong critical value. A sign error in calculating the test statistic meant that critical value and test statistic were at opposite ends of the distribution so that, although the candidate reached the correct conclusion, it was based on an incorrect argument and the final mark was not earned. A critical value and a test statistic with opposite signs in a one-tailed test should warn the candidate to check back for errors.

The candidate produced a flawless solution to part (b) with the diagram being particularly helpful in clarifying the final explanation. This candidate used the distribution of a proportion, but equally successful solutions were produced using the normal approximation to a binomial distribution or an exact binomial distribution. As no method was specified in the question, any appropriate choice of method could gain full marks.

Mark Scheme

Q	Solution	Marks	Comments
2(a)	$\bar{x} = 32.75$ $s = 3.059$	B1	Allow 3.06 for sd.
	$H_0 : \mu = 30$ $H_1 : \mu > 30$	B1	Both
	Test statistic = $\frac{32.75 - 30}{\frac{3.059}{\sqrt{8}}}$	M1 m1	
	$= 2.543$	A1	2.54 to 2.544
	$\nu = 8 - 1 = 7$	B1	
	$t = 2.998$	B1	
	2.543 < 2.998 so cannot reject H_0 . There is not enough evidence at the 1% significance level to say that the level of nitrate pollution is a cause for environmental concern.	A1√	ft on ts and critical value.
	(b)(i) $H_0 : p = 0.6$ $H_1 : p < 0.6$	B1	Both
	Under H_0 , $X \sim B(42, 0.6)$ $\approx N(25.2, 10.08)$	B2	B1 mean; B1 variance.
	Test statistic = $\frac{16.5 - 25.2}{\sqrt{10.08}} = -2.74$ or $\frac{16 - 25.2}{\sqrt{10.08}} = -2.90$	M1 A1	-2.75 to -2.73 -2.91 to -2.89
Using proportions, $Y \sim N(0.6, 0.005714)$	(B2)		
Test statistic = $\frac{\frac{16}{42} - 0.6}{\sqrt{0.005714}} = -2.90$	(M1) (A1)		
$z = -2.3263$ $ts < \text{critical value}$ so H_0 can be rejected. There is evidence at the 1% level that nitrate pollution is a cause for concern in less than 60% of the length of rivers in the region.	B1 A1√	Accept 2.33; ignore sign ft on ts and z.	
(ii) Assume that the sections tested are chosen at random.	B1		
		16	

Question 3

There used to be a cash machine outside a village store. During the time the store was open, the number of people per hour who used the machine was modelled by a Poisson distribution with mean 13.5.

For security reasons the machine was then moved inside the store. On a particular day after the machine was moved, 47 people used the machine during a 4-hour period.

- (a) Assuming that a Poisson distribution continues to be a suitable model, construct an approximate 90% confidence interval for the mean number of people using the machine during a 4-hour period when the store is open. (4 marks)
- (b) Comment on the suggestion that fewer people use the machine in its new location. (3 marks)

Student Response

③	X/4		Leave blank
a	mean per hour = 13.5		
	mean per 4 hours = $13.5 \times 4 = 54$ ✓		
	$X \sim \text{Po}(54)$ approx $\sim N(54, 54)$		
	90% C.I for μ : B1		
	$54 \pm 1.6449 \times \sqrt{54}$		1
	54 ± 12.09		
	= 41.91 to 66.09	B1	
b	It is possible On that day less people used the machine, but 47 is within the confidence intervals so it is impossible to say that the machine is used less in its new location. B1		2
			③

Commentary

The example above illustrates a very common error in solutions to this question where the wrong figure was used as a basis for the confidence interval in part (a). The value of 13.5 users per hour when the machine was outside the store was stated to be a population mean, so it was not appropriate to treat it as a sample value.

Most candidates realised that some scaling was necessary to match the two time periods involved, usually, as in the example, by finding the mean number of users for a four-hour period when the machine was outside the store. This candidate recognised that the comparison in part (b) required use of the number 47, which referred to the number of users after the move. It was, however, quite common to see the confidence interval calculated as above, and then an argument based on the fact that 54 lies within the confidence interval.

Mark Scheme

Q	Solution	Marks	Comments
3(a)	<p>Y = Number of people using the machine in a 4-hour period.</p> <p>$Y \sim \text{Po}(\lambda) \approx N(\lambda, \lambda)$</p> <p>$z = 1.6449$</p> <p>standard error = $\sqrt{47}$</p> <p>90% confidence limits for λ are:</p> <p>$47 \pm 1.6449 \times \sqrt{47}$</p> <p>giving (35.7, 58.3)</p>	<p>B1</p> <p>B1</p> <p>M1</p> <p>A1</p>	<p>Accept 1.645, 1.64</p> <p>(35.7 to 35.8, 58.2 to 58.3)</p>
(b)	<p>Mean uses per 4 hours when outside = 54 or 90% CI for mean uses per hour when inside is $\left(\frac{35.7}{4}, \frac{58.3}{4}\right)$</p> <p>Mean when outside lies within CI for mean when inside.</p> <p>Not enough evidence to say that fewer people use the machine in its new location while the store is open.</p> <p>There may still be a reduction in use because it is not available outside opening hours.</p> <p>Conclusion from confidence interval may be suspect if model used is poor.</p>	<p>B1</p> <p>B2</p>	<p>Comparison of CI with previous mean.</p> <p>Any two valid points.</p>
		7	

Question 4

Yvonne was one of two candidates in a local election. On the day of the election, voters were interviewed as they left their polling stations and asked which candidate they had voted for.

From a random sample of 320 females, 188 said that they had voted for Yvonne.

From a random sample of 260 males, 117 said that they had voted for Yvonne.

- (a) Assuming that the answers given were truthful, construct an approximate 95% confidence interval for:
- the proportion of females who voted for Yvonne;
 - the proportion of males who voted for Yvonne. (7 marks)
- (b) State, with a reason, whether females were more likely than males to vote for Yvonne. (2 marks)
- (c) At the close of voting, Yvonne believed that she had won the election. Comment, with justifications, on her belief. (4 marks)

Student Response

<p><u>b</u> males are less likely to vote for yvonne as my upper level of the interval is lower than that of the lower limit of the female confidence interval.</p>	2
<p><u>c</u></p> $320 + 260 = 580$ $188 + 117 = 305$	
$\frac{305}{580} = 0.525862069$ <p>Sample values only</p> <p><u>0.526</u> Yvonne had over half the votes in her favour. F!</p>	1
(10)	

Commentary

Part (a) of this question was very well done, often with efficient use of statistical calculators, so that most candidates were arguing from correct confidence intervals in the remaining parts of the question.

The example shows an excellent solution to part (b), making clear use of the confidence intervals. Many candidates attempted to compare the two intervals, but stated that the upper limit for females was higher than that for males, and similarly with the lower limits. This could still have been true if the intervals had a large amount of overlap, so did not prove the point convincingly.

This candidate's attempt at part (c) exemplifies a very common approach which based the argument on the sample values only. The conclusion reached could only be justified if the 580 people in the two samples were the only voters in the election. Some candidates went a stage further and constructed another confidence interval based on the combined sample. This led to a better argument, but was still not reliable. A crucial point in trying to forecast the result of the election was the ratio of males to females who actually voted and this could be quite different from the ratio of sample sizes.

Mark Scheme

Q	Solution	Marks	Comments
4(a)(i)	$\hat{p} = \frac{188}{320} = 0.5875$ $z = 1.96$ 95% confidence limits for p are: $0.5875 \pm 1.96 \times \sqrt{\frac{0.5875 \times 0.4125}{320}}$ giving (0.534, 0.641)	B1 B1 M1 m1 A1	Allow 0.587, 0.588 Here or in (ii) Here or in (ii) m1 for standard error (0.533 to 0.534, 0.64 to 0.642)
(ii)	$\hat{q} = \frac{117}{260} = 0.45$ 95% confidence limits for q are: $0.45 \pm 1.96 \times \sqrt{\frac{0.45 \times 0.55}{260}}$ giving (0.390, 0.510)	M1 A1	(0.389 to 0.39, 0.51 to 0.511)
(b)	The lower bound of the CI for p is greater than the upper bound of the CI for q . It seems likely that females were more likely than males to vote for Yvonne.	E1 B1	
(c)	Two candidates so Yvonne needs > 50% of votes. Lower bound for $q > 0.5$ so it is likely that a majority of females voted for her, but the proportion of males could be well under a half. Result may depend on numbers of males and females who vote. Yvonne's belief seems over-optimistic.	E1 E1 B1	
		13	

Question 5

A company produces low calorie sweetener tablets. Each tablet produced is sealed in a packet. The weight, X milligrams, of a tablet is normally distributed with mean 110 and standard deviation 5. The weight, Y milligrams, of an empty packet is normally distributed with mean 370 and standard deviation 12.

You may assume that X and Y are independent random variables.

- (a) Find the probability that the weight of a packet containing a tablet is less than 500 milligrams. (4 marks)
- (b) By considering the variable $Y - 3X$, find the probability that the weight of a packet is more than three times the weight of the tablet it contains. (7 marks)

Student Response

		Leave blank
5	$X \sim N(110, 5^2)$ $Y \sim N(370, 12^2)$	
a)	$P(X+Y < 500)$ $= 0.938$	4
	$110 + 370 = 480$ mean $5^2 + 12^2 = \sqrt{169} = 13$ stand deviation $X+Y \sim N(480, 13^2)$	
b)	$\text{mean} = 370 - 3 \times 110 = 40$ $\text{s.d} = 12^2 + (3) \times 5^2 = 8.660$	3
	$P(Y - 3X > 0) = 0.999$	(7)

Commentary

Many candidates are confident about working with the sum of two independent, normally distributed variables and there were plenty of completely correct solutions to part (a) as shown in the example.

Calculating the variance for $Y - 3X$ was a cause of many errors. Most recognised that the separate variances must be added, but many forgot to square the coefficient of X , as illustrated in this example. This candidate started off correctly to find the required probability and scored the first method mark; another method mark was available if the candidate had written down a formula to show how the probability was obtained.

Mark Scheme

Q	Solution	Marks	Comments
5	$X \sim N(110, 5^2)$ $Y \sim N(370, 12^2)$		
(a)	$X + Y \sim N(110 + 370, 5^2 + 12^2)$ $= N(480, 169)$ $P(X + Y < 500) = \Phi\left(\frac{500 - 480}{13}\right)$ $= \Phi(1.538)$ $= 0.938$	M1 A1 M1 A1	Means and variances added. cao 1.538 to 1.54 awrt 0.937 to 0.938
(b)	$3X \sim N(3 \times 110, 9 \times 5^2)$ $Y - 3X \sim N(370 - 330, 12^2 + 225)$ $= N(40, 369)$ $P(Y > 3X) = P(Y - 3X > 0)$ $= 1 - \Phi\left(\frac{0 - 40}{\sqrt{369}}\right)$ $= \Phi(2.082)$ $= 0.981$	M1 M2 A1 M1 m1 A1	M1 mean; M1 variance. cao or equivalent. awrt 2.08 awrt
		11	

Question 6

A short-stay car park in a shopping area has spaces marked out for 90 cars. A local councillor notices that there are always some vacant spaces. He puts forward a plan to create a garden and seating area using part of the car park. This would reduce the number of parking spaces to 78.

- (a) From a random sample of 14 users of the car park, 11 say that the car park will be too small if this plan is carried out. Carry out a test to determine whether more than half of the users of the car park think it will be too small. Use an exact distribution and the 5% significance level. (6 marks)

- (b) The number of occupied spaces, x , in the car park is recorded on each of 16 randomly chosen occasions during shopping hours. The results may be summarised as follows:

$$\bar{x} = 59.9 \quad s = 7.83$$

- (i) Construct a 95% confidence interval for the mean, μ , of the number of spaces occupied in the car park during shopping hours. Assume that the sample is drawn from a normal population. (4 marks)
- (ii) The councillor claims that the value of μ is no more than 65. State, with a reason, whether this claim is plausible. (2 marks)
- (c) It is found that the number of occupied spaces during shopping hours is best modelled by a Poisson distribution with mean μ .
- (i) Comment on the validity of your confidence interval found in part (b)(i). (2 marks)
- (ii) Taking μ to be 65, use a distributional approximation to find the probability that more than 78 spaces are occupied in the car park at any one time. (4 marks)
- (d) Use your results to explain whether or not the car park will be too small if the plan to create a garden and seating area is carried out. (4 marks)

Student Response

bii	As the whole of the confidence interval for μ is less than 65 then the claim that $\mu \leq 65$ is plausible as it is not over 65. It is between (55.7, 64.1) (56, 64)	2
cii	$X \sim Po(65) \rightarrow X \sim N(65, 65)$ ✓ 78.5 ✓ $P(X \geq 78) \quad Z = \frac{78 - 65}{\sqrt{65}} = 1.67$ ✓ $1 - 0.95254 = 0.04746$ ✓	4
d.	as cii shows there is a low probability of 78 spaces being taken up at any one time ✓ Part bi also shows only ✓ E1 between 55.7, 64.1 spaces are taken up at any one time. It is only the users that think it will be too small. Also three distributions have been used to test this so overall the car park will not be too small with only 78 carparking spaces. B1 ✓	2

Commentary

Candidates often find difficulty with hypothesis tests where they are required to use an exact probability distribution and this was the case in part (a) of this question where use of the binomial distribution $B(14, 0.5)$ was required. Many used a normal approximation which did not answer the question as set. Part (b) was generally well done although some candidates used a z -value rather than the appropriate t -value.

A good deal of the question involved the selection of appropriate information and results and their interpretation. The above example shows a good answer to part (b) (ii), making use of the confidence interval just calculated.

There were very few clear answers to part (c) (i). Some explained that a Poisson distribution may be approximated by a normal distribution, but only the best candidates extracted the fact that we were dealing with a distribution with a large mean which made the approximation appropriate. This candidate's answer illustrates one of various attempts to provide supporting evidence to say that either the confidence interval was still valid or that it was not. Other common suggestions were to invoke the central limit theorem in favour of validity, or to say that the interval was not valid because the sample was too small.

The example shows a perfect solution to part (c) (ii) – the most common reasons for loss of marks here were to forget the continuity correction or to go in the wrong direction and use 77.5. This candidate's use of a diagram is a good way of clarifying the point.

The final part of the question required careful selection and use of previous results. Many quoted the result of part (a), which could be a factor to consider in making the final decision as to whether to reduce the size of the car park, but was not relevant in determining whether or not it would be too small. The example shows a common misinterpretation of the confidence interval found in part (b) (i). The candidate treats it as being the range for the number of spaces occupied at any one time, rather than for the mean number of occupied spaces. The contribution of this confidence interval to the argument was that it had earlier indicated that μ was likely to be no more than 65, so that we were looking at a worst case scenario in part (c) (ii). Very few candidates realised that the probability found in (c) (ii) would be even smaller if a lower value of μ was used. Most candidates quoted their result from (c) (ii) and used it as evidence to support their conclusion.

Mark Scheme

Q	Solution	Marks	Comments
6(a)	$H_0 : p = 0.5$ $H_1 : p > 0.5$ Under H_0 , $X \sim B(14, 0.5)$ $P(X \geq 11) = 1 - P(X \leq 10)$ $= 1 - 0.9713$ $= 0.0287$ $0.0287 < 5\%$ so result is significant at the 5% level. Evidence suggests that more than half of car park users think it will be too small.	B1 B1 M1 A1 E1 A1√	Both. Accept 0.029 ft on probability.
(b)(i)	$v = 15$; $t = 2.131$ 95% confidence limits for μ are: $59.9 \pm 2.131 \times \frac{7.83}{4}$ giving (55.7, 64.1)	B1 M1 m1 A1	sd divided by 4. (awrt, 64 to 64.1)
(ii)	65 is above upper confidence limit. Seems likely that the claim is true.	E1 B1	
(c)(i)	Parent population is (discrete) Poisson but with large mean so closely approximated by normal distribution. Reasonable to accept that CI is valid.	E1 B1	
(ii)	$Y \sim \text{Po}(65) \approx N(65, 65)$ $P(Y > 78) = 1 - \Phi\left(\frac{78.5 - 65}{\sqrt{65}}\right)$ $= 1 - \Phi(1.674)$ $= 1 - 0.9529 = 0.0471$	B1 M1 m1 A1	Continuity correction attempted. 0.047 to 0.0475; cao
(d)	Probability that reduced car park is too small = 0.0471 (with $\mu = 65$). μ taken at higher value than expected so true probability could be smaller. It seems likely that the car park will be large enough most of the time. May be some occasions when it is full / no allowance for increasing business	E1 E1 B1 E1	Significance of μ used. Extra point based on size of probability from (c)(ii).
		22	