Paper Reference(s)

# 6684/01
# Edexcel GCE
## Statistics S2
## Gold Level G2

## Time: 1 hour 30 minutes

| **Materials required for examination papers** | **Items included with question** |
|---|---|
| Mathematical Formulae (Green) | Nil |

**Candidates may use any calculator allowed by the regulations of the Joint Council for Qualifications. Calculators must not have the facility for symbolic algebra manipulation, differentiation and integration, or have retrievable mathematical formulas stored in them.**

## Instructions to Candidates

Write the name of the examining body (Edexcel), your centre number, candidate number, the unit title (Statistics S2), the paper reference (6684), your surname, initials and signature.

## Information for Candidates

A booklet 'Mathematical Formulae and Statistical Tables' is provided.
Full marks may be obtained for answers to ALL questions.
There are 8 questions in this question paper. The total mark for this paper is 75.

## Advice to Candidates

You must ensure that your answers to parts of questions are clearly labelled.
You must show sufficient working to make your methods clear to the Examiner. Answers without working may gain no credit.

**Suggested grade boundaries for this paper:**

| A* | A | B | C | D | E |
|---|---|---|---|---|---|
| 65 | 55 | 47 | 36 | 26 | 17 |

1.   A bag contains a large number of counters. A third of the counters have a number 5 on them and the remainder have a number 1.

A random sample of 3 counters is selected.

(*a*)  List all possible samples.

**(2)**

(*b*)  Find the sampling distribution for the range.

**(3)**

_____

2.   A test statistic has a distribution B(25, *p*).

Given that

$$H_0 : p = 0.5, \quad H_1 : p \neq 0.5,$$

(*a*)  find the critical region for the test statistic such that the probability in each tail is as close as possible to 2.5%.

**(3)**

(*b*)  State the probability of incorrectly rejecting $H_0$ using this critical region.

**(2)**

_____

3.   Bacteria are randomly distributed in a river at a rate of 5 per litre of water. A new factory opens and a scientist claims it is polluting the river with bacteria. He takes a sample of 0.5 litres of water from the river near the factory and finds that it contains 7 bacteria. Stating your hypotheses clearly test, at the 5% level of significance, the claim of the scientist.

**(7)**

_____

**4.** Past records suggest that 30% of customers who buy baked beans from a large supermarket buy them in single tins. A new manager questions whether or not there has been a change in the proportion of customers who buy baked beans in single tins. A random sample of 20 customers who had bought baked beans was taken.

(*a*) Using a 10% level of significance, find the critical region for a two-tailed test to answer the manager's question. You should state the probability of rejection in each tail which should be less than 0.05.

**(5)**

(*b*) Write down the actual significance level of a test based on your critical region from part (*a*).

**(1)**

The manager found that 11 customers from the sample of 20 had bought baked beans in single tins.

(*c*) Comment on this finding in the light of your critical region found in part (*a*).

**(2)**

_____

**5.** In a game, players select sticks at random from a box containing a large number of sticks of different lengths. The length, in cm, of a randomly chosen stick has a continuous uniform distribution over the interval [7, 10].

A stick is selected at random from the box.

(*a*) Find the probability that the stick is shorter than 9.5 cm.

**(2)**

To win a bag of sweets, a player must select 3 sticks and wins if the length of the longest stick is more than 9.5 cm.

(*b*) Find the probability of winning a bag of sweets.

**(2)**

To win a soft toy, a player must select 6 sticks and wins the toy if more than four of the sticks are shorter than 7.6 cm.

(*c*) Find the probability of winning a soft toy.

**(4)**

_____

**6.** A company has a large number of regular users logging onto its website. On average 4 users every hour fail to connect to the company's website at their first attempt.

(*a*) Explain why the Poisson distribution may be a suitable model in this case.

**(1)**

Find the probability that, in a randomly chosen **2 hour** period,

(*b*) (i) all users connect at their first attempt,

(ii) at least 4 users fail to connect at their first attempt.

**(5)**

The company suffered from a virus infecting its computer system. During this infection it was found that the number of users failing to connect at their first attempt, over a 12 hour period, was 60.

(*c*) Using a suitable approximation, test whether or not the mean number of users per hour who failed to connect at their first attempt had increased. Use a 5% level of significance and state your hypotheses clearly.

**(9)**

---

**7.** (*a*) Define the critical region of a test statistic.

**(2)**

A discrete random variable $X$ has a Binomial distribution $B(30, p)$. A single observation is used to test $H_0 : p = 0.3$ against $H_1 : p \neq 0.3$

(*b*) Using a 1% level of significance find the critical region of this test. You should state the probability of rejection in each tail which should be as close as possible to 0.005.

**(5)**

(*c*) Write down the actual significance level of the test.

**(1)**

The value of the observation was found to be 15.

(*d*) Comment on this finding in light of your critical region.

**(2)**

---

**8.** The continuous random variable $X$ has probability density function given by

$$f(x) = \begin{cases} \dfrac{3}{32}(x-1)(5-x) & 1 \le x \le 5, \\ \\ 0 & \text{otherwise.} \end{cases}$$

(a) Sketch $f(x)$ showing clearly the points where it meets the $x$-axis.

**(2)**

(b) Write down the value of the mean, $\mu$, of $X$.

**(1)**

(c) Show that $E(X^2) = 9.8$.

**(4)**

(d) Find the standard deviation, $\sigma$, of $X$.

**(2)**

The cumulative distribution function of $X$ is given by

$$F(x) = \begin{cases} 0 & x < 1 \\ \\ \dfrac{1}{32}(a - 15x + 9x^2 - x^3) & 1 \le x \le 5 \\ \\ 1 & x > 5 \end{cases}$$

where $a$ is a constant.

(e) Find the value of $a$.

**(2)**

(f) Show that the lower quartile of $X$, $q_1$, lies between 2.29 and 2.31.

**(3)**

(g) Hence find the upper quartile of $X$, giving your answer to 1 decimal place.

**(1)**

(h) Find, to 2 decimal places, the value of $k$ so that

$$P(\mu - k\sigma < X < \mu + k\sigma) = 0.5.$$

**(2)**

---

**TOTAL FOR PAPER: 75 MARKS**

**END**

| Question Number | Scheme | Marks |
|---|---|---|
| **1.** (a) | (1, 1, 1), (5, 5, 5), (1, 5, 5), (1, 5, 1) | B1 |
| | (1,1,1); (5,5,5); (1, 5, 5); (5, 1, 5); (5, 5, 1) (5, 1, 1); (1, 5, 1); (1, 1, 5) | B1 |
| | | **(2)** |
| (b) | $r$: 0 and 4 | B1 |
| | $P(R = 0) = \dfrac{9}{27}$ or $\dfrac{1}{3}$    $P(R = 4) = \dfrac{18}{27}$ or $\dfrac{2}{3}$ | M1d A1 |
| | | **(3)** |
| | | **[5]** |

| | | | |
|---|---|---|---|
| **2.** (a) | $X \sim B(25,0.5)$    may be implied by calculations in part a or b | | M1 |
| | $P(X \le 7) = 0.0216$ | | |
| | $P(X \ge 18) = 0.0216$ | | |
| | CR $X \le 7$; $\cup$ $X \ge 18$ | | A1,A1 |
| | | | (3) |
| (b) | $P(\text{rejecting } H_0) = 0.0216 + 0.0216$ | | M1 |
| | $= 0.0432$ | awrt 0.0432/0.0433 | A1 |
| | | | (2) |
| | | | **Total 5** |

| Question Number | Scheme | Marks |
|---|---|---|
| **3** | <u>One tail test</u><br>Method 1<br>$\quad$ $H_o : \lambda = 5$ ($\lambda = 2.5$) $\qquad\qquad\qquad\qquad\qquad$ may use $\lambda$ or $\mu$<br>$\quad$ $H_1 : \lambda > 5$ ($\lambda > 2.5$)<br><br>$X \sim Po\,(2.5)$ $\qquad\qquad\qquad\qquad\qquad\qquad\qquad$ may be implied<br><br>$P(X \geq 7) = 1 - P(X \leq 6)$ $\quad$ [ $P(X \geq 5) = 1 - 0.8912 = 0.1088$ ] $\quad$ att $P(X \geq 7)$ \| $P(X \geq 6)$<br>$\qquad\qquad = 1 - 0.9858$ $\qquad$ $P(X \geq 6) = 1 - 0.9580 = 0.0420$<br><br>$\qquad\qquad = 0.0142$ $\qquad\qquad\qquad$ CR $X \geq 6$ $\qquad\qquad$ awrt 0.0142<br><br>$0.0142 < 0.05$ $\qquad$ $7 \geq 6$ or 7 is in critical region or 7 is significant<br><br>(Reject $H_0$.) There is significant evidence at the 5% significance level that the factory <u>is polluting the river</u> with bacteria.<br>**or**<br>The scientists claim is justified | B1<br>B1<br><br>M1<br><br>M1<br><br><br>A1<br><br>M1<br><br>B1<br><br><br>(7)<br>Total 7 |
| | Method 2<br>$H_o : \lambda = 5$ ($\lambda = 2.5$) $\qquad\qquad\qquad\qquad$ may use $\lambda$ or $\mu$<br>$H_1 : \lambda > 5$ ($\lambda > 2.5$)<br><br>$X \sim Po\,(2.5)$ $\qquad\qquad\qquad\qquad\qquad\qquad\qquad$ may be implied<br><br>$P(X < 7)$ $\qquad\qquad$ [$P(X < 5) = 0.8912$] $\qquad$ att $P(X < 7)$ \| $P(X < 6)$<br>$\qquad\qquad\qquad$ $P(X < 6) = 0.9580$<br><br>$\qquad\qquad = 0.9858$ $\qquad\qquad\qquad$ CR $X \geq 6$ $\qquad$ wrt 0.986<br><br>$0.9858 > 0.95$ $\qquad$ $7 \geq 6$ or 7 is in critical region or 7 is significant<br><br>(Reject $H_0$.) There is significant evidence at the 5% significance level that the factory <u>is polluting the river</u> with bacteria<u>.</u><br>**or**<br>The scientists claim is justified | B1<br>B1<br><br>M1<br><br><br><br>M1 A1<br><br>M1<br><br>B1<br><br>(7) |

| | | |
|---|---|---|
| Two tail test<br>Method 1<br><br>H$_o$ : λ = 5 (λ = 2.5)          may use λ or μ<br>H$_1$ : λ ≠ 5 (λ ≠ 2.5)<br><br>$X \sim$ Po (2.5) | | B1<br>B0<br><br>M1 |
| P($X \geq 7$) = 1 − P($X \leq 6$) | [P($X \geq 6$) = 1− 0.9580 = 0.0420]     att P($X \geq 7$) │P($X \geq 7$) | M1 |
| = 1 − 0.9858 | P($X \geq 7$) = 1− 0.9858 = 0.0142 | |
| = 0.0142 | CR $X \geq 7$          awrt 0.0142 | A1 |
| 0.0142 < 0.025 | 7 ≥ 7 or 7 is in critical region or 7 is significant│ | M1 |
| (Reject H$_0$.) There is significant evidence at the 5% significance level that the factory <u>is polluting the river </u>with bacteria.<br>**or**<br>The scientists claim is justified | | B1<br><br>(7) |
| Method 2<br>H$_o$ : λ = 5 (λ = 2.5)          may use λ or μ<br>H$_1$ : λ ≠ 5 (λ ≠ 2.5)<br><br>$X \sim$ Po (2.5) | | B1<br>B0<br><br>M1 |
| P($X < 7$) | [P($X < 6$) = 0.9580]          att P($X < 7$) │P($X < 7$)<br>P($X < 7$) =  0.9858 | |
| = 0.9858 | CR $X \geq 7$          awrt 0.986 | M1A1 |
| 0.9858 > 0.975 | 7 ≥ 7 or 7 is in critical region or 7 is significant│ | M1 |
| (Reject H$_0$.) There is significant evidence at the 5% significance level that the factory <u>is polluting the river </u>with bacteria.<br>**or**<br>The scientists claim is justified | | B1<br><br>(7) |

| Question Number | | Scheme | Marks |
|---|---|---|---|
| **4.** | **(a)** | $X \sim$ B(20, 0.3) | M1 |
| | | $\qquad\qquad\qquad\qquad\qquad\qquad$ P($X \le 2$) = 0.0355 | A1 |
| | | P($X \le 9$) = 0.9520 $\qquad$ so $\qquad\qquad\qquad$ P($X \ge 10$) = 0.0480 | A1 |
| | | Therefore the critical region is $\{X \le 2\} \cup \{X \ge 10\}$ | A1 A1 (5) |
| | **(b)** | 0.0355 + 0.0480 = 0.0835 $\qquad\qquad\qquad$ awrt (0.083 or 0.084) | B1 (1) |
| | **(c)** | 11 is in the critical region | B1ft |
| | | there is evidence of a <u>change/ increase</u> in the <u>proportion/number</u> of <u>customers buying single tins</u> | B1ft (2) |
| | | | **(8 marks)** |

| | | | |
|---|---|---|---|
| **5.** | **(a)** | $\dfrac{9.5 - 7}{10 - 7}$ | M1 |
| | | $= \dfrac{5}{6}$ $\qquad\qquad\qquad\qquad\qquad$ awrt 0.833 | A1 |
| | | | (2) |
| | **(b)** | P(Longest > 9.5) = 1 - P(all < 9.5) = $1 - \left(\dfrac{5}{6}\right)^3$ | M1 |
| | | $\qquad\qquad\qquad\qquad\qquad = \dfrac{91}{216}$ or 0.421 | A1 |
| | | | (2) |
| | **(c)** | P(a stick < 7.6) = $\dfrac{0.6}{3} = 0.2$ | B1 |
| | | Let $Y$ = number of sticks (out of 6) <7.6 $\quad$ then $Y \sim$ B(6, 0.2) | M1 |
| | | P($Y$ > 4) = 1 - P($Y \le 4$) | M1 |
| | | $\qquad$ = 1 - 0.9984 | |
| | | $\qquad\qquad$ = 0.0016 or $\dfrac{1}{625}$ | A1 |
| | | | (4) |
| | | | **8** |

| Question Number | | Scheme | Marks |
|---|---|---|---|
| 6. | (a) | Connecting occurs at random/independently, singly or at a constant rate | B1 (1) |
| | (b) | Po (8) | B1 |
| | (i) | $P(X = 0) = 0.0003$ | M1A1 |
| | (ii) | $P(X \geq 4) = 1 - P(X \leq 3)$ | M1 |
| | | $= 1 - 0.0424$ | A1 (5) |
| | | $= 0.9576$ | |
| | (c) | H0 : $\lambda = 4$ (48)  H1 : $\lambda > 4$ (48) | B1 |
| | | N(48,48) | M1 A1 |

Method 1

$P(X \geq 59.5) = P\left( Z \geq \dfrac{59.5 - 48}{\sqrt{48}} \right)$

$= P ( Z \geq 1.66)$
$= 1 - 0.9515$
$= 0.0485$

Method 2

$\dfrac{x - 0.5 - 48}{\sqrt{48}} = 1.6449$

M1 M1 A1

$x = 59.9$

A1

$0.0485 < 0.05$

Reject H$_0$. Significant. 60 lies in the Critical region — M1

The number of failed connections at the first attempt has increased. — A1 ft (9)

[15]

| Question Number | | Scheme | Marks |
|---|---|---|---|
| 7. | | | |
| | (a) | The set of values of the test statistic for which | B1 |
| | | the null hypothesis is rejected in a hypothesis test. | B1 (2) |
| | (b) | $X \sim B(30,0.3)$ | M1 |
| | | $P(X \leq 3) = 0.0093$ | |
| | | $P(X \leq 2) = 0.0021$ | A1 |
| | | $P(X \geq 16) = 1 - 0.9936 = 0.0064$ | |
| | | $P(X \geq 17) = 1 - 0.9979 = 0.0021$ | A1 |
| | | Critical region is $(0 \leq)x \leq 2$ or $16 \leq x(\leq 30)$ | A1A1 (5) |
| | (c) | Actual significance level 0.0021+0.0064=0.0085 or 0.85% | B1 (1) |
| | (d) | 15 (it) is not in the critical region | Bft 2, 1, 0 |
| | | not significant | |
| | | No significant evidence of a change in $p = 0.3$ | |
| | | accept H$_0$, (reject H$_1$) | |
| | | $P(x \geq 15) = 0.0169$ | (2) |
| | | | Total [10] |

| Question Number | Scheme | Marks |
|---|---|---|
| **8.** (a) | ∩ shape which does not go below the *x*-axis [condone missing patios] <br> Graph must end at the points (1,0) and (5,0) and the points labelled at 1 and 5 | B1 <br> B1 <br> (2) |
| (b) | $E(X) = 3$  (by symmetry) | B1 <br> (1) |
| (c) | $\left[E(X^2)\right] = \int x^2 f(x)dx = \dfrac{3}{32}\int \left(6x^3 - x^4 - 5x^2\right)dx$ | M1 |
| | $= \tfrac{3}{32}\left[\dfrac{6x^4}{4} - \dfrac{x^5}{5} - \dfrac{5x^3}{3}\right]_1^5$ | A1 |
| | $= \tfrac{3}{32}\left(\left[\dfrac{6\times 625}{4} - 625 - \dfrac{625}{3}\right] - \left[\dfrac{6}{4} - \dfrac{1}{5} - \dfrac{5}{3}\right]\right) = 9.8$  (*) | M1 <br> A1 cso <br> (4) |
| (d) | s.d. = $\sqrt{9.8 - E(X)^2}$ , <br>    $= 0.8944\ldots$ **awrt 0.894** | M1 <br> A1 <br> (2) |
| (e) | $F(1) = 0 \Rightarrow \tfrac{1}{32}(a - 15 + 9 - 1) = 0$,   leading to <u>$a = 7$</u> | M1 A1 <br> (2) |
| (f) | $F(2.29) = 0.2449\ldots$, $F(2.31) = 0.2515\ldots$ <br> Since $F(q_1) = 0.25$ and these values are either side of 0.25 then $2.29 < q_1 < 2.31$ | M1 A1 <br> A1 <br> (3) |
| (g) | Since the distribution is symmetric $q_3 = 5 - 1.3 = \underline{3.7}$         cao | B1 <br> (1) |
| (h) | We know $P(q_1 = 2.3 < X < 3.7 = q_3) = 0.5$ <br> so $k\sigma = 0.7$ <br> so $k = \dfrac{0.7}{0.894\ldots} = 0.7826\ldots = $ **awrt 0.78** | M1 <br> A1 <br> (2) |
| | | **17** |

**Examiner reports**

**Question 1**

This question was answered well by many candidates. In part (a) the vast majority of candidates were able to score one mark and many scored both marks as they could list all 8 combinations.

Part (b) caused some candidates problems as they failed to recognise that the range was the difference between the largest and smallest values. Some candidates found the sampling distribution of the mean or the sum rather than the range. Others thought that all samples were equally likely and hence obtained the wrong probabilities.

Those candidates who could identify the range usually went on to calculate the correct corresponding probabilities. A common error was to give the range as 0 and 1.

**Question 2**

Part (a) was a routine question with many fully correct answers being seen. However, a minority of candidates were unable to write down the critical regions correctly. The two most common errors were to write down the regions in terms of probabilities, e.g. $P(X \leq 7)$ or to simply write $X = 7$ and $X = 18$.

In part (b) many candidates were distracted by the word "incorrectly". We reject the null hypothesis 'incorrectly' when it is in fact true. In this case, this means that $p$ really is 0.5, so that we can use the distribution B(25, 0.5) to work out probabilities. This is the distribution that had already been used in part (a), so that all that needed to be done was to simply add the two probabilities that were used to identify the two parts of the critical region in part (a): $0.0216 + 0.0216 = 0.0432$. It was all too common to see correct working followed by incorrect and irrelevant work that invalidated the whole response. The most common incorrect postscripts were $0.05 - 0.0432 = 0.0068$ and $1 - 0.0432 = 0.9568$.

**Question 3**

The majority of candidates found this question straightforward. They were most successful if they used the probability method and compared it with 0.05. Those who attempted to use 95% were less successful and this is not a recommended route for these tests.

Most candidates knew how to specify the hypotheses with most candidates using 2.5 rather than 5. Some candidates used $p$, or did not use a letter at all, in stating their hypotheses, but most of the time they used $\lambda$. A minority found $P(X=7)$ and some worked with Po(5).

If using the critical region method, not all candidates showed clearly, either their working, or a comparison with the value of 7 and the CR $X > 7$. A sizeable minority of candidates failed to put their conclusion back into the given context. Reject H0 is not sufficient.

**Question 4**

This was a very well answered question. Candidates were able to use binomial tables and gave the answer to the required number of decimal places. As in previous years there were some candidates who confused the critical region with the probability of the test statistic being in that region but this error has decreased. Candidates were able to describe the acceptance of the hypothesis in context although sometimes it would be better if they just repeated the wording from the question which would help them avoid some of the mistakes seen. There were still a few candidates who did not give a reason in context at all.

In part (a) many candidates failed to read this question carefully assuming it was identical to similar ones set previously. Most candidates correctly identified B(20,0.3) to earn the method mark and many had 0.0355 written down to earn the first A mark, although in light of their subsequent work, this may often have been accidental. A majority of candidates did not gain the second A mark as they failed to respond to the instruction "state the probability of rejection in each case". In the more serious cases, candidates had shown no probabilities from the tables, doing all their work mentally, only writing their general strategy: "$P(X \leq c) < 0.05$". Whilst many candidates were able to write down the critical region using the correct notation there are still some candidates who are losing marks they should have earned, by writing $P(X \leq 2)$ for the critical region $X \leq 2$

Part (b) was usually correct.

Part (c) provided yet more evidence of candidates who had failed to read the question: "in the light of your critical region". Some candidates chose not to mention the critical region and a number of those candidates who identified that 11 was in the critical region did not refer to the manager's question.

**Question 5**

Part (a) was routine and the vast majority of candidates demonstrated familiarity with the probability density function of a rectangular distribution.

Part (b) required some careful initial thought that eluded a large majority of the candidates. The successful solutions fell into two camps. On the one hand were those candidates who identified the Binomial distribution $X \sim B\left(3, \dfrac{5}{6}\right)$ and then evaluated the probability $P(X \geq 1)$ by conventional means. On the other hand there were candidates who worked 'from first principles' and produced either an elaborate tree diagram or a list of all possible outcomes. This latter approach ('from first principles') is of course valid, but takes longer and is more susceptible to error.

There were many complete and correct solutions to part (c). Finding the length of a stick shorter than 7.6 cm was straightforward for most but many attempts reflected candidates' lack of understanding of what the question required of them in this part. Successful candidates realised it was a binomial situation using B(6, $p$) and used tables to find $1 - P(X \leq 4)$ or calculated $P(X = 5) + P(X = 6)$. Common errors included misinterpretation of $P(X > 4)$ as $1 - P(X \leq 3)$ or even $P(X = 4)$ and a small minority of candidates gave the answer as $\left(\dfrac{1}{5}\right)^4 = \dfrac{1}{625}$ which gained no marks as it is an incorrect method.

**Question 6**

The majority of candidates were familiar with the technical terms in part (a), but failed to establish any context.

Part (b) was a useful source of marks for a large proportion of the candidates. The only problems were occasional errors in detail. In part (i) a few did not spot the change in time scale and used Po(4) rather than Po(8). Some were confused by the wording and calculated $P(X = 8)$ rather than $P(X = 0)$. The main source of error for (ii) was to find $1 - P(X \leq 4)$ instead of $1 - P(X \leq 3)$.

In part (c) the Normal distribution was a well-rehearsed routine for many candidates with many candidates concluding the question with a clear statement in context.

The main errors were
- Some other letter (or none) in place of $\lambda$ or $\mu$
- Incorrect Normal distribution: e.g. N(60, 60)
- Omission of (or an incorrect) continuity correction
- Using 48 instead of 60
- Calculation errors

A minority of candidates who used the wrong distribution (usually Poisson) were still able to earn the final two marks in the many cases when clear working was shown. This question was generally well done with many candidates scoring full marks.

**Question 7**

Part (a) tested candidates' understanding of the critical region of a test statistic and responses were very varied, with many giving answers in terms of a 'region' or 'area' and making no reference to the null hypothesis or the test being significant. Many candidates lost at least one mark in part (b), either through not showing the working to get the probability for the upper critical value, i.e. $1 - P(X \leq 15) = P(X \geq 16) = 0.0064$, or by not showing any results that indicated that they had used B(30, 0.3) and just writing down the critical regions, often incorrectly. A minority of candidates still write their critical regions in terms of probabilities and lose the final two marks. Responses in part (c) were generally good with the majority of candidates making a comment about the observed value and their critical region. A small percentage of responses contained contradictory statements.

**Question 8**

The overall response to part (a) was disappointing. The shape of the graph is assumed knowledge from GCSE mathematics. If f($x$) is a quadratic expression, then the graph of $y$ = f($x$) is a parabola, whose shape is 'known'. Many candidates made a 'table of values' which they then used to plot the graph which then consisted of straight lines. Occasionally the correct parabola shape was seen with 1 and 5 marked on the $x$-axis but the graph continued below the $x$-axis indicating that the probability could be negative for values outside of the given range.

Although part (b) was intended to be straightforward, many candidates failed to see the symmetry of the probability density function (even with a correct diagram), and pursued integration to find the mean. Fortunately, most attempts were successful. It should be noted that in general if there is 1 mark for a question the answer should be able to be written down without any long calculations.

Part (c) was well answered by the majority of candidates, with responses ranging from the very concise to lengthy. Nearly all candidates were able to multiply the brackets and integrate successfully, although there were occasional sign errors. Sufficient working for the substitution of the limits was nearly always shown. Some candidates preferred to find Var($X$) first, only to add [ E($X$) ]$^2$ at the end.
Some candidates obtained an answer of almost 9.8, this is despite the fact that most modern calculators will work using exact fractions. It was in fact possible to work with exact decimals, the penultimate stage of working being 9.765625 − (−0.034375). However, some candidates chose to approximate the earlier decimals, resulting in a final answer that differed from 9.8.

Part (d) was usually answered correctly, with a significant minority giving their answer in surd form.
In part (e) the obvious methods, in this case solving either F(1) = 0 or F(5) = 1, were not universally adopted. For those who adopted the latter approach, a surprising proportion struggled to solve the linear equation: $\frac{1}{32}(a - 75 + 225 - 125) = 1$. Many candidates integrated f($x$), or occasionally F($x$), with mixed success.

The solution to part (f) had been well rehearsed by candidates, although some failed to give an adequate explanation for the final mark. An ideal correct response such as "0.25 is between 0.245 and 0.252, so the median must lie between 2.29 and 2.31" was rarely seen. A sizeable minority pursued the solution of F($x$) = $\frac{1}{4}$. Those who used a sign change method were sometimes successful, and a few candidates used their graphical calculators to state all three roots, before selecting the required one.

In part (g), many candidates did not appreciate the symmetric nature of the probability density function. However, those who did usually produced a correct answer of 3.7. Other candidates resorted to a repetition of their method for finding the lower quartile.

In part (h) those candidates who saw that the given probability referred to the IQR of the distribution often produced a neat solution. More often, however, candidates abandoned their solution, usually after much algebraic effort.

## Statistics for S2 Practice Paper Gold 2

| Qu | Max Score | Modal score | Mean % | Mean average scored by candidates achieving grade: | | | | | | | |
|----|-----------|-------------|--------|------|------|------|------|------|------|------|------|
| | | | | ALL | A* | A | B | C | D | E | U |
| 1 | 5 | | 62.4 | 3.12 | 4.01 | 3.62 | 2.55 | 2.03 | 1.66 | 1.31 | 0.85 |
| 2 | 5 | | 67.2 | 3.36 | 4.38 | 4.07 | 3.54 | 2.98 | 2.35 | 1.65 | 0.82 |
| 3 | 7 | | 64.9 | 4.54 | | 5.90 | 5.08 | 3.90 | 2.77 | 1.80 | 0.76 |
| 4 | 8 | | 62.4 | 4.99 | | 6.45 | 5.18 | 4.25 | 3.27 | 2.18 | 1.14 |
| 5 | 8 | | 59.5 | 4.76 | 6.80 | 6.07 | 4.88 | 3.93 | 2.92 | 1.84 | 0.99 |
| 6 | 15 | | 70.2 | 10.53 | 13.77 | 12.86 | 11.06 | 9.10 | 6.63 | 4.69 | 2.36 |
| 7 | 10 | | 64.5 | 6.45 | | 6.64 | 5.50 | 4.26 | 2.76 | 1.58 | 0.79 |
| 8 | 17 | | 56.2 | 9.56 | 14.37 | 12.56 | 9.54 | 7.32 | 5.53 | 3.61 | 2.38 |
| | 75 | | 63.1 | 47.31 | | 58.17 | 47.33 | 37.77 | 27.89 | 18.66 | 10.09 |