# OCR
Oxford Cambridge and RSA

## Tuesday 9 June 2015 – Morning

## A2 GCE MATHEMATICS (MEI)

**4768/01** Statistics 3

### QUESTION PAPER

**Duration:** 1 hour 30 minutes

### INSTRUCTIONS TO CANDIDATES

These instructions are the same on the Printed Answer Book and the Question Paper.

*   The Question Paper will be found inside the Printed Answer Book.
*   Write your name, centre number and candidate number in the spaces provided on the Printed Answer Book. Please write clearly and in capital letters.
*   **Write your answer to each question in the space provided in the Printed Answer Book**. Additional paper may be used if necessary but you must clearly show your candidate number, centre number and question number(s).
*   Use black ink. HB pencil may be used for graphs and diagrams only.
*   Read each question carefully. Make sure you know what you have to do before starting your answer.
*   Answer **all** the questions.
*   Do **not** write in the bar codes.
*   You are permitted to use a scientific or graphical calculator in this paper.
*   Final answers should be given to a degree of accuracy appropriate to the context.

### INFORMATION FOR CANDIDATES

This information is the same on the Printed Answer Book and the Question Paper.

*   The number of marks is given in brackets **[ ]** at the end of each question or part question on the Question Paper.
*   You are advised that an answer may receive **no marks** unless you show sufficient detail of the working to indicate that a correct method is being used.
*   The total number of marks for this paper is **72**.
*   The Printed Answer Book consists of **12** pages. The Question Paper consists of **4** pages. Any blank pages are indicated.

### INSTRUCTION TO EXAMS OFFICER / INVIGILATOR

*   Do not send this Question Paper for marking; it should be retained in the centre or recycled. Please contact OCR Copyright should you wish to re-use this document.

OCR is an exempt Charity

**Turn over**

**1 (a)** A stratified sample of pupils at secondary schools in a particular local authority is to be chosen in order to collect information on absenteeism. In the local authority there are 4 secondary schools, A, B, C and D, with 1310, 1453, 843 and 1110 pupils respectively.

    **(i)** How many pupils should be chosen from each school in a stratified sample of 500 so that each school is represented proportionally? **[3]**

    **(ii)** Suggest two possible criteria for stratification other than by school. **[2]**

    **(iii)** State one advantage of choosing a stratified sample. **[1]**

**(b)** At a large secondary school, the median number of half days absent per pupil per year (based on several years' records) was known to be 23. Last year the school carried out a drive to lower the number of absences. A random sample of 12 pupils had been absent for the following numbers of half days during the year.

$$14 \quad 10 \quad 15 \quad 13 \quad 35 \quad 9 \quad 24 \quad 19 \quad 30 \quad 26 \quad 29 \quad 8$$

A Wilcoxon single sample test is to be carried out to see if the drive has been successful.

    **(i)** Why might a Wilcoxon test be appropriate? **[1]**

    **(ii)** What distributional assumption is needed for the test? **[1]**

    **(iii)** Carry out the test, using a 5% level of significance. **[10]**

**2** The distribution of the random variable $X$ is thought to be well modelled by the following probability density function:

$$f(x) = \begin{cases} k(1+x) & \text{for } 0 \leqslant x < 5, \\ 0 & \text{elsewhere,} \end{cases}$$

where $k$ is a positive constant.

    **(i)** Find the value of $k$. **[3]**

    **(ii)** Show that $P(a \leqslant X < a+1) = \frac{1}{35}(2a+3)$ for $0 \leqslant a \leqslant 4$. **[2]**

A random sample of 50 observations of $X$ is summarised as follows.

| $x$ | $0 \leqslant x < 1$ | $1 \leqslant x < 2$ | $2 \leqslant x < 3$ | $3 \leqslant x < 4$ | $4 \leqslant x < 5$ |
|---|---|---|---|---|---|
| Frequency | 1 | 5 | 7 | 20 | 17 |

    **(iii)** Test at the 10% level of significance whether the distribution of $X$ is well modelled by $f(x)$. **[10]**

    **(iv)** With reference to your calculations in part **(iii)** discuss briefly the outcome of the test. **[2]**

**3**  In agricultural research the oil content, as a percentage of the whole grain, of a cereal can be measured using near infra-red spectroscopy. An investigation into the effect of a particular treatment on the oil content of a certain cereal is being carried out. A sample of 10 plots of land is chosen and each plot is divided in half. In one half of each plot the cereal is grown with the treatment and in the other half the cereal is grown without the treatment. Subsequently the percentage oil content of the cereal for each half of each plot is measured and the results are as follows.

| Plot | A | B | C | D | E | F | G | H | I | J |
|---|---|---|---|---|---|---|---|---|---|---|
| With treatment | 41.1 | 44.3 | 42.4 | 48.2 | 52.4 | 54.6 | 35.9 | 33.6 | 51.1 | 47.0 |
| Without treatment | 42.5 | 37.7 | 42.1 | 32.4 | 42.7 | 41.5 | 36.9 | 31.7 | 52.6 | 41.2 |

A paired $t$ test with a 5% level of significance is to be used to see if the treatment appears to make any difference to the mean percentage oil content of the cereal.

  (i)  Explain what is meant by a 5% level of significance in a hypothesis test. **[2]**

  (ii)  State the conditions necessary for the test to be carried out. **[3]**

  (iii)  Assuming the conditions stated in part **(ii)** are met, carry out the test. **[10]**

  (iv)  Find a 90% confidence interval for the population mean difference in the percentage oil content with and without the treatment. **[4]**

**4**  Paul has been trying a new route to work in the mornings. He collects a large random sample of times, in minutes, and calculates a 95% confidence interval for the population mean time by this route. The confidence interval is $(45.369, 47.231)$ and the sample variance is 20.3.

  (i)  Explain what is meant by a 95% confidence interval for a population mean. **[1]**

  (ii)  Calculate the sample mean and the sample size. **[4]**

Paul reverts to his usual route and the time, in minutes, to travel to work each morning is modelled by a random variable which is Normally distributed with mean 41.3 and variance 11.7. The time, in minutes, for Paul to travel home each evening is modelled by a random variable which is Normally distributed with mean 44.8 and variance 14.2. In the rest of this question all journeys are by Paul's usual route and may be assumed to be independent of each other.

  (iii)  Calculate the probability that, on a randomly chosen day, Paul's total travelling time will be less than 90 minutes. **[3]**

  (iv)  Calculate the probability that, on a randomly chosen day, the time for Paul to travel home will be more than 5 minutes longer than the time to travel to work. **[4]**

  (v)  Calculate the probability that, in a randomly chosen five-day week, the mean time for Paul to travel to work on Monday and Tuesday will be more than 3 minutes longer than his mean time to travel to work on Wednesday, Thursday and Friday. **[6]**

**END OF QUESTION PAPER**

| Question | | | Answer | Marks | Guidance |
|---|---|---|---|---|---|
| 1 | (a) | (i) | $\dfrac{500}{4716} \times ...$ <br><br> = 138.88…, 154.05…, 89.37…, 117.68… <br> = 139,      154,      89,      118 | M1 <br><br><br> A1 <br> A1 <br> **[3]** | Correct factor used for at least one school. <br><br><br> All correct and given to at least 1 dp. <br> FT any errors in the previous line provided that sum = 500. |
| 1 | (a) | (ii) | e.g.    Sex (gender) <br>         Year group | B1 <br> B1 <br><br> **[2]** | Allow reasonable alternatives including ethnicity, birth date, <br> distance from school |
| 1 | (a) | (iii) | e.g. Provides information on each stratum (as well as the <br>       population). | B1 <br><br> **[1]** | Or representative |
| 1 | (b) | (i) | We have no information about the background population. | E1 <br> **[1]** | o.e. Must include "population" o.e. |
| 1 | (b) | (ii) | Symmetry. | B1 <br> **[1]** | |
| 1 | (b) | (iii) | $H_0$: $m = 23$      $H_1$: $m < 23$ <br><br><br><br> where $m$ is the population median number of days absent. | B1 <br><br><br><br> B1 | Both. Accept hypotheses in words, but must include "population". <br> Do NOT allow symbols other than $m$ unless clearly and explicitly <br> stated to be a <u>population</u> <u>median</u>. <br> Adequate definition of $m$ to include "population". |

| Absences | −23 | Rank of \|diff\| |
|---|---|---|
| 14 | −9 | 7 |
| 10 | −13 | 10 |
| 15 | −8 | 6 |
| 13 | −10 | 8 |
| 35 | 12 | 9 |
| 9 | −14 | 11 |
| 24 | 1 | 1 |
| 19 | −4 | 3 |
| 30 | 7 | 5 |
| 26 | 3 | 2 |
| 29 | 6 | 4 |
| 8 | −15 | 12 |

M1    for subtracting 23.

M1    for ranks. <br> A1    ft if ranks wrong.

| Question | | | Answer | Marks | Guidance | |
|---|---|---|---|---|---|---|
| | | | $W_+ = 1 + 2 + 4 + 5 + 9 = 21$ | B1 | $(W_- = 3 + 6 + 7 + 8 + 10 + 11 + 12 = 57)$ | |
| | | | Refer to Wilcoxon single sample tables for $n = 12$. | M1 | No ft from here if wrong. | |
| | | | Lower 5% point is 17 (or upper is 61 if 57 used). | A1 | i.e. a 1-tail test. No ft from here if wrong. | |
| | | | Result is not significant. | A1 | ft only c's test statistic. Dependent on all 3 M marks | |
| | | | Insufficient evidence to suggest that the median number of days absent has been reduced. | A1 | ft only c's test statistic. Dependent on all 3 M marks. Conclusion in context to include "on average" o.e. | |
| | | | | **[10]** | | |
| 2 | (i) | | Require $\int_0^5 k(1+x)\,dx = 1$ | | | |
| | | | $\int_0^5 k(1+x)\,dx = k\left(x + \dfrac{x^2}{2}\right)\Big|_0^5 = k\left(5 + \dfrac{25}{2}\right) - k \times 0 = \dfrac{35k}{2}$ | M1<br><br>A1 | Set up correct integral, including limits which may appear later. Allow method based on area, e.g., a trapezium.<br>Integral correctly evaluated, or correct area obtained, in terms of $k$. | |
| | | | $\therefore \dfrac{35k}{2} = 1 \qquad \therefore k = \dfrac{2}{35} \qquad$ Not 0.057…. | A1 | Set equal to 1 and rearranged for $k$. | |
| | | | | **[3]** | | |
| 2 | (ii) | | $P(a \le X < a+1) = \int_a^{a+1} \dfrac{2}{35}(1+x)\,dx$ | M1 | Set up correct integral, including limits which may appear later. Allow method based on area, e.g., a trapezium. Allow candidate's value of $k$. | |
| | | | $= \dfrac{2}{35}\left((a+1) + \dfrac{(a+1)^2}{2} - a - \dfrac{a^2}{2}\right)$ | | | |
| | | | $= \dfrac{2}{35}\left(1 + \dfrac{2a+1}{2}\right) = \dfrac{1}{35}(2a+3)$ | A1 | AG. Must be shown convincingly. | |
| | | | | **[2]** | | |
| 2 | (iii) | | $H_0$: The model is suitable / fits the data.<br>$H_1$: The model is not suitable / does not fit the data. | B1 | Both hypotheses. Must be the right way round.<br>Do not accept "data fit model" oe. | |
| | | | Expected frequencies are: $50 \times \left(\dfrac{3,\ 5,\ 7,\ 9,\ 11}{35}\right)$ | M1 | | |
| | | | $= \left(\dfrac{30,\ 50,\ 70,\ 90,\ 110}{7}\right)$ | A1 | Accept either fractions or decimals. | |
| | | | $= 4.2857,\ 7.1428/9,\ 10,\ 12.8571,\ 15.7142/3$ | | | |

| Question | | | Answer | Marks | Guidance | |
|---|---|---|---|---|---|---|
| | | | Merge first 2 cells:   Obs f = 6,   Exp f = 11.4285 | M1 | Merge first 2 cells | |
| | | | $X^2 = 2.5786 + 0.9 + 3.9683 + 0.1052$ | M1 | Calculation of $X^2$. Independent of previous mark. | |
| | | | $= 7.552$ | A1 | Awrt 7.55 | |
| | | | Refer to $\chi_3^2$. | M1 | Allow correct df (= cells – 1) from differently grouped table and ft. critical value only.  Otherwise, no ft if wrong. | |
| | | | Upper 10% point is 6.251. | A1 | No ft from here if wrong. $P(X^2 > 7.552) = 0.0562$. <br> If cells not merged $\chi_4^2$ 10% point is 7.779; $P(X^2 > 8.135) = 0.0867$. | |
| | | | Significant. | A1 | ft only c's test statistic. | |
| | | | Sufficient evidence to suggest that the pdf of $X$ is not well modelled by f($x$).. | A1 | ft only c's test statistic. <br> Do not accept "data do not fit model" oe. | |
| | | | | **[10]** | | |
| 2 | (iv) | | e.g. The model overestimates for $0 \le x < 2$. | E1 | Any 2 points relating to or explaining the outcome of the test. | |
| | | | The model underestimates for $3 \le x < 4$. | E1 | Other possibilities might include: <br> The test would not have been significant at 5%. | |
| | | | "Large discrepancy" but no direction E1 max | | The sample is a bit small making it difficult to assess. | |
| | | | | **[2]** | | |
| 3 | (i) | | 5% represents the probability of rejecting the null hypothesis … | E1 | | |
| | | | … when it is, in fact, true. | E1 | | |
| | | | | **[2]** | | |
| 3 | (ii) | | Must assume:   Normality of population … | B1 | Ignore references to unknown variance and/or sample size. | |
| | | | … of <u>differences</u>. | B1 | | |
| | | | Sample is random. | B1 | | |
| | | | | **[3]** | | |
| 3 | (iii) | | $H_0$: $\mu_D = 0$ <br> $H_1$: $\mu_D \ne 0$ | B1 | Both. Accept alternatives e.g. $\mu_A - \mu_B$ etc provided adequately defined. <br> Hypotheses in words only must include "population". Do NOT allow "$\bar{X} = ...$" or similar unless $\bar{X}$ is clearly and explicitly stated to be a <u>population</u> mean. | |
| | | | Where $\mu_D$ is the (population) mean difference in percentage oil content. | B1 | For adequate verbal definition. Allow absence of "population" if correct notation $\mu$ is used. | |

| Question | | | Answer | Marks | Guidance | |
|---|---|---|---|---|---|---|
| | | | <u>MUST</u> be PAIRED COMPARISON $t$ test.<br>Differences (with – without) are:<br>−1.4  6.6  0.3  15.8  9.7  13.1  −1.0  1.9  −1.5  5.8<br>$\bar{x} = 4.93$  $s_{n-1} = 6.310(4)$  $(s_{n-1}^2 = 39.822(3))$ | M1<br>A1 | Allow "without – with" if consistent with alternatives for<br>hypotheses above.<br>Do not allow $s_n = 5.9886$ ($s_n^2 = 35.8401$). | |
| | | | Test statistic is $\dfrac{4.93 - 0}{\dfrac{6.310}{\sqrt{10}}}$ | M1 | Allow c's $\bar{x}$ and/or $s_{n-1}$.<br>Allow alternative: $0 + $ (c's 2.262) $\times \dfrac{6.3104}{\sqrt{10}}$ (= 4.514) for<br>subsequent comparison with $\bar{x}$.<br>(Or $\bar{x} - $ (c's 2.262) $\times \dfrac{6.3104}{\sqrt{10}}$ (= 0.416) for comparison with 0.) | |
| | | | $= 2.470(4)$. | A1 | c.a.o. but ft from here in any case if wrong.<br>Use of $10 - \bar{x}$ scores M1A0, but ft. | |
| | | | Refer to $t_9$.<br>Double-tailed 5% point is ±2.262. | M1<br>A1 | No ft from here if wrong.<br>Must be minus 2.262 for "without – with" unless absolute values<br>are being compared. No ft from here if wrong. P($|t| > 2.4704$) =<br>0.03554. | |
| | | | Significant. | A1 | ft only c's test statistic as long as it includes their $s\big/\sqrt{n}$ | |
| | | | Sufficient evidence to suggest that the treatment appears to<br>make a difference to the mean percentage oil content of the<br>cereal. | A1 | ft only c's test statistic as above<br>Conclusion in context to include "on average" o.e. | |
| | | | | **[10]** | | |
| 3 | (iv) | | CI is given by  $4.93 \pm$ | M1 | ZERO/4 if not same distribution as test. Same wrong distribution<br>scores maximum M1B0M1A0. Recovery to $t_9$ is OK.<br>Allow c's $\bar{x}$. | |
| | | | 1.833<br>$\times \dfrac{6.3104}{\sqrt{10}}$ | B1<br>M1 | 1.833 seen.<br>Allow c's $s_{n-1}$. | |
| | | | $= 4.93 \pm 3.6577 = (1.271(9), 8.588(1))$ | A1 | c.a.o. Must be expressed as an interval. | |
| | | | | **[4]** | | |

| Question | | | Answer | Marks | Guidance | | |
|---|---|---|---|---|---|---|---|
| 4 | (i) | | In repeated sampling, 95% of all confidence intervals constructed in this way will contain the true mean. | E1 [1] | | | |
| 4 | (ii) | | Mean = (45.369 + 47.231)/2=46.3 $$47.231 = 46.3 + \sqrt{20.3} \times 1.96 / \sqrt{n}$$ $$n = \frac{1.96^2 \times 20.3}{0.931^2} \approx 89.97(2) \qquad \therefore n = 90$$ | B1 B1 M1  A1  [4] | cao Sight of 1.96. Or equivalent.  Must be an integer. FT candidate's mean. | | |
| 4 | (iii) | | Time to work $X \sim N(41.3, \ 11.7)$ Time to home $Y \sim N(44.8, \ 14.2)$ $X + Y \sim N(86.1, \ 25.9)$  $$P(X + Y < 90) = \Phi\left(\frac{90 - 86.1}{\sqrt{25.9}} = 0.7663\right) = 0.7783$$ | B1 B1  B1  [3] | Mean Variance  cao | | |
| 4 | (iv) | | Require $P(Y - X > 5)$ $Y - X \sim N(3.5, \ 25.9)$  $$P(Y - X > 5) = 1 - \Phi\left(\frac{5 - 3.5}{\sqrt{25.9}} = 0.2947\right) = 1 - 0.6159 = 0.3841$$ | M1 B1 B1  A1  [4] | Allow equivalent alternatives, e.g. $Y > X + 5$ or $X - Y < -5$ Mean Variance  cao | | |
| 4 | (v) | | Require $P\left(\frac{X_1 + X_2}{2} > \frac{X_3 + X_4 + X_5}{3} + 3\right)$    Mean $= 3(41.3 + 41.3) - 2(41.3 + 41.3 + 41.3) = 0$   Variance $= \frac{1}{4}(11.7 + 11.7) + \frac{1}{9}(11.7 + 11.7 + 11.7) = 9.75$ | M1  A1  B1  M1 A1 | For considering some $\overline{X_2}$ and $\overline{X_3}$  For $\overline{X_2} - \overline{X_3} > 3$  For 0  For 1/4, 1/9, and 11.7 seen o.e. For 9.75 | | |

| Question | | | Answer | Marks | Guidance | |
|---|---|---|---|---|---|---|
| | | | $P(Z > \dfrac{3-0}{\sqrt{9.75}}) = 1 - 0.8317 = 0.1683$ | A1 | cao | |
| | | | | **[6]** | | |