

**ADVANCED GCE UNIT
MATHEMATICS (MEI)**

Statistics 3

TUESDAY 5 JUNE 2007

4768/01

Afternoon

Time: 1 hour 30 minutes

Additional Materials:

Answer booklet (8 pages)

Graph paper

MEI Examination Formulae and Tables (MF2)

INSTRUCTIONS TO CANDIDATES

- Write your name, centre number and candidate number in the spaces provided on the answer booklet.
- Answer **all** the questions.
- You are permitted to use a graphical calculator in this paper.
- Final answers should be given to a degree of accuracy appropriate to the context.

INFORMATION FOR CANDIDATES

- The number of marks is given in brackets [] at the end of each question or part question.
- The total number of marks for this paper is 72.

ADVICE TO CANDIDATES

- Read each question carefully and make sure you know what you have to do before starting your answer.
- 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.

This document consists of **4** printed pages.

- 1 A manufacturer of fireworks is investigating the lengths of time for which the fireworks burn. For a particular type of firework this length of time, in minutes, is modelled by the random variable T with probability density function

$$f(t) = kt^3(2 - t) \quad \text{for } 0 < t \leq 2$$

where k is a constant.

- (i) Show that $k = \frac{5}{8}$. [2]
- (ii) Find the modal time. [2]
- (iii) Find $E(T)$ and show that $\text{Var}(T) = \frac{8}{63}$. [5]
- (iv) A large random sample of n fireworks of this type is tested. Write down in terms of n the approximate distribution of \bar{T} , the sample mean time. [3]
- (v) For a random sample of 100 such fireworks the times are summarised as follows.

$$\Sigma t = 145.2 \quad \Sigma t^2 = 223.41$$

Find a 95% confidence interval for the mean time for this type of firework and hence comment on the appropriateness of the model. [6]

- 2 The operator of a section of motorway toll road records its weekly takings according to the types of vehicles using the motorway. For purposes of charging, there are three types of vehicle: cars, coaches, lorries. The weekly takings (in thousands of pounds) for each type are assumed to be Normally distributed. These distributions are independent of each other and are summarised in the table.

Vehicle type	Mean	Standard deviation
Cars	60.2	5.2
Coaches	33.9	6.3
Lorries	52.4	4.9

- (i) Find the probability that the weekly takings for coaches are less than £40 000. [3]
- (ii) Find the probability that the weekly takings for lorries exceed the weekly takings for cars. [4]
- (iii) Find the probability that over a 4-week period the total takings for cars exceed £225 000. What assumption must be made about the four weeks? [5]
- (iv) Each week the operator allocates part of the takings for repairs. This is determined for each type of vehicle according to estimates of the long-term damage caused. It is calculated as follows: 5% of takings for cars, 10% for coaches and 20% for lorries. Find the probability that in any given week the total amount allocated for repairs will exceed £20 000. [6]

- 3 The management of a large chain of shops aims to reduce the level of absenteeism among its workforce by means of an incentive bonus scheme. In order to evaluate the effectiveness of the scheme, the management measures the percentage of working days lost before and after its introduction for each of a random sample of 11 shops. The results are shown below.

Shop	A	B	C	D	E	F	G	H	I	J	K
% days lost before	3.5	5.0	3.5	3.2	4.5	4.9	4.1	6.0	6.8	8.1	6.0
% days lost after	1.8	4.3	2.9	4.5	4.4	5.8	3.5	6.7	6.4	5.4	5.1

- (a) The management decides to carry out a t test to investigate whether there has been a reduction in absenteeism.
- (i) State clearly the hypotheses that should be used together with any necessary assumptions. [4]
- (ii) Carry out the test using a 5% significance level. [7]
- (b) Find a 95% confidence interval for the true mean percentage of days lost after the introduction of the incentive scheme and state any assumption needed. The management has set a target that the mean percentage should be 3.5. Do you think this has been achieved? Explain your answer. [7]
- 4 A machine produces plastic strip in a continuous process. Occasionally there is a flaw at some point along the strip. The length of strip (in hundreds of metres) between successive flaws is modelled by a continuous random variable X with probability density function $f(x) = \frac{18}{(3+x)^3}$ for $x > 0$. The table below gives the frequencies for 100 randomly chosen observations of X . It also gives the probabilities for the class intervals using the model.

Length x (hundreds of metres)	Observed frequency	Probability
$0 < x \leq 0.5$	21	0.2653
$0.5 < x \leq 1$	24	0.1722
$1 < x \leq 2$	12	0.2025
$2 < x \leq 3$	15	0.1100
$3 < x \leq 5$	13	0.1094
$5 < x \leq 10$	9	0.0874
$x > 10$	6	0.0532

- (i) Examine the fit of this model to the data at the 5% level of significance. [9]

You are given that the median length between successive flaws is 124 metres. At a later date the following random sample of ten lengths (in metres) between flaws is obtained.

239 77 179 221 100 312 52 129 236 42

- (ii) Test at the 10% level of significance whether the median length may still be assumed to be 124 metres. [9]

Permission to reproduce items where third-party owned material protected by copyright is included has been sought and cleared where possible. Every reasonable effort has been made by the publisher (OCR) to trace copyright holders, but if any items requiring clearance have unwittingly been included, the publisher will be pleased to make amends at the earliest possible opportunity.

OCR is part of the Cambridge Assessment Group. Cambridge Assessment is the brand name of University of Cambridge Local Examinations Syndicate (UCLES), which is itself a department of the University of Cambridge.

Mark Scheme 4768
June 2007

Q1	$f(t) = kt^3(2-t) \quad 0 < t \leq 2$			
(i)	$\int_0^2 kt^3(2-t)dt = 1$ $\therefore \left[k \left(\frac{2t^4}{4} - \frac{t^5}{5} \right) \right]_0^2 = 1$ $\therefore k \left(8 - \frac{32}{5} \right) - 0 = 1$ $\therefore k \times \frac{8}{5} = 1 \quad \therefore k = \frac{5}{8}$	M1 E1	Integral of $f(t)$, including limits (possibly implied later), equated to 1. Convincingly shown. Beware printed answer.	2
(ii)	$\frac{df}{dt} = \frac{5}{8}(6t^2 - 4t^3) = 0$ $\therefore 6t^2 - 4t^3 = 0$ $\therefore 2t^2(3 - 2t) = 0$ $\therefore t = (0 \text{ or } \frac{3}{2})$	M1 A1	Differentiate and set equal to zero. c.a.o.	2
(iii)	$E(T) = \int_0^{\frac{2}{8}} t^4(2-t)dt$ $= \left[\frac{5}{8} \left(\frac{2t^5}{5} - \frac{t^6}{6} \right) \right]_0^{\frac{2}{8}} = \frac{5}{8} \times \left(\frac{64}{5} - \frac{64}{6} \right) = \frac{4}{3}$ $E(T^2) = \int_0^{\frac{2}{8}} t^5(2-t)dt$ $= \left[\frac{5}{8} \left(\frac{2t^6}{6} - \frac{t^7}{7} \right) \right]_0^{\frac{2}{8}} = \frac{5}{8} \times \left(\frac{128}{6} - \frac{128}{7} \right) = \frac{40}{21}$ $\text{Var}(T) = \frac{40}{21} - \left(\frac{4}{3} \right)^2 = \frac{120 - 112}{63} = \frac{8}{63}$	M1 A1 M1 M1 A1	Integral for $E(T)$ including limits (which may appear later). Integral for $E(T^2)$ including limits (which may appear later). Convincingly shown. Beware printed answer.	5
(iv)	$\bar{T} \sim N\left(\frac{4}{3}, \frac{8}{63n}\right)$	B1 B1 B1	Normal distribution. Mean. ft c's $E(T)$. Correct variance.	3

(v)	$n = 100, \quad \bar{t} = \frac{145 \cdot 2}{100} = 1 \cdot 452,$ $s_{n-1}^2 = \frac{223 \cdot 41 - 100 \times 1 \cdot 452^2}{99} = 0 \cdot 12707$ <p>CI is given by $1 \cdot 452 \pm$</p> $1 \cdot 96 \times \frac{0 \cdot 3565}{\sqrt{100}}$ $= 1 \cdot 452 \pm 0 \cdot 0698 = (1 \cdot 382, 1 \cdot 522)$ <p>Since $E(T)$ ($= 4/3$) lies outside this interval it seems the model may not be appropriate.</p>	<p>B1</p> <p>M1</p> <p>B1</p> <p>M1</p> <p>A1</p> <p>E1</p>	<p>Both mean and variance. Accept sd = 0·3565</p> <p>ft c's $\bar{t} \pm$.</p> <p>ft c's s_{n-1}.</p> <p>c.a.o. Must be expressed as an interval.</p>	<p>6</p>
				18

Q2	$Ca \sim N(60.2, 5.2^2)$ $Co \sim N(33.9, 6.3^2)$ $L \sim N(52.4, 4.9^2)$		When a candidate's answers suggest that (s)he appears to have neglected to use the difference columns of the Normal distribution tables, penalise the first occurrence only.	
(i)	$P(Co < 40) = P\left(Z < \frac{40 - 33.9}{6.3} = 0.9683\right)$ $= 0.8336$	M1 A1 A1	For standardising. Award once, here or elsewhere. c.a.o.	3
(ii)	Want $P(L > Ca)$ i.e. $P(L - Ca > 0)$ $L - Ca \sim N(52.4 - 60.2 = -7.8,$ $4.9^2 + 5.2^2 = 51.05)$ $P(\text{this} > 0) = P\left(Z > \frac{0 - (-7.8)}{\sqrt{51.05}} = 1.0917\right)$ $= 1 - 0.8625 = 0.1375$	M1 B1 B1 A1	Allow $Ca - L$ provided subsequent work is consistent. Mean. Variance. Accept $sd = \sqrt{51.05} = 7.1449\dots$ c.a.o.	4
(iii)	Want $P(Ca_1 + Ca_2 + Ca_3 + Ca_4 > 225)$ $Ca_1 + \dots \sim N(60.2 + 60.2 + 60.2 + 60.2 = 240.8,$ $5.2^2 + 5.2^2 + 5.2^2 + 5.2^2 = 108.16)$ $P(\text{this} > 225) = P\left(Z > \frac{225 - 240.8}{\sqrt{108.16}} = -1.519\right)$ $= 0.9356$ Must assume that the weeks are independent of each other.	M1 B1 B1 A1 B1	Mean. Variance. Accept $sd = \sqrt{108.16} = 10.4$. c.a.o.	5
(iv)	$R \sim N(0.05 \times 60.2 + 0.1 \times 33.9 + 0.2 \times 52.4 = 16.88,$ $0.05^2 \times 5.2^2 + 0.1^2 \times 6.3^2 + 0.2^2 \times 4.9^2 = 1.4249)$ $P(R > 20) = P\left(Z > \frac{20 - 16.88}{\sqrt{1.4249}} = 2.613\right)$ $= 1 - 0.9955 = 0.0045$	M1 A1 M1 M1 A1 A1	Mean. For 0.05^2 etc. For $\times 5.2^2$ etc. Accept $sd = \sqrt{1.4249} = 1.1937$. c.a.o.	6
				18

<p>(b)</p>	<p>For “days lost after” $\bar{x}=4.6182, s_{n1} = 1.4851 (s_{n1}^2 = 2.2056)$</p> <p>CI is given by $4.6182 \pm$ 2.228 $\times \frac{1.4851}{\sqrt{11}}$ $= 4.6182 \pm 0.9976 = (3.620(6), 5.615(8))$</p> <p>Assume Normality of population of “days lost after”.</p> <p>Since 3.5 lies outside the interval it seems that the target has not been achieved.</p>	<p>B1</p> <p>M1</p> <p>B1</p> <p>M1</p> <p>A1</p> <p>E1</p> <p>E1</p>	<p>Do not allow $s_n = 1.4160 (s_n^2 = 2.0051)$.</p> <p>ft c’s $\bar{x} \pm$.</p> <p>ft c’s s_{n1}.</p> <p>c.a.o. Must be expressed as an interval. ZERO if not same distribution as test. Same wrong distribution scores maximum M1B0M1A0. Recovery to t_{10} is OK.</p>	<p>7</p>
				<p>18</p>

Q4																																									
(i)	<table border="1" style="width: 100%; border-collapse: collapse;"> <tr> <td style="text-align: center;">Obs</td> <td style="text-align: center;">21</td> <td style="text-align: center;">24</td> <td style="text-align: center;">12</td> <td style="text-align: center;">15</td> <td style="text-align: center;">13</td> <td style="text-align: center;">9</td> <td style="text-align: center;">6</td> </tr> <tr> <td style="text-align: center;">Exp</td> <td style="text-align: center;">26.53</td> <td style="text-align: center;">17.22</td> <td style="text-align: center;">20.25</td> <td style="text-align: center;">11.00</td> <td style="text-align: center;">10.94</td> <td style="text-align: center;">8.74</td> <td style="text-align: center;">5.32</td> </tr> </table>							Obs	21	24	12	15	13	9	6	Exp	26.53	17.22	20.25	11.00	10.94	8.74	5.32																		
Obs	21	24	12	15	13	9	6																																		
Exp	26.53	17.22	20.25	11.00	10.94	8.74	5.32																																		
	<p> $\therefore X^2 = \frac{(21 - 26.53)^2}{26.53} + \text{etc}$ $= 1.1527 + 2.6695 + 3.3611 + 1.4545 + 0.3879$ $+ 0.0077 + 0.0869$ $= 9.1203$ </p> <p>d.o.f. = 7 - 1 = 6 Refer to χ^2_6. Upper 5% point is 12.59 $9.1203 < 12.59 \therefore$ Result is not significant. Evidence suggests the model fits the data at the 5% level.</p>							<p>M1 A1 M1 A1 A1 M1 A1 E1 E1</p>	<p>Probabilities $\times 100$. All Expected frequencies correct.</p> <p>At least 4 values correct.</p> <p>No ft from here if wrong. No ft from here if wrong. ft only c's test statistic. ft only c's test statistic.</p>	9																															
(ii)	<table border="1" style="width: 100%; border-collapse: collapse;"> <thead> <tr> <th style="text-align: center;">Data</th> <th style="text-align: center;">Diff = data - 124</th> <th style="text-align: center;">Rank of diff </th> </tr> </thead> <tbody> <tr><td style="text-align: center;">239</td><td style="text-align: center;">115</td><td style="text-align: center;">9</td></tr> <tr><td style="text-align: center;">77</td><td style="text-align: center;">-47</td><td style="text-align: center;">3</td></tr> <tr><td style="text-align: center;">179</td><td style="text-align: center;">55</td><td style="text-align: center;">4</td></tr> <tr><td style="text-align: center;">221</td><td style="text-align: center;">97</td><td style="text-align: center;">7</td></tr> <tr><td style="text-align: center;">100</td><td style="text-align: center;">-24</td><td style="text-align: center;">2</td></tr> <tr><td style="text-align: center;">312</td><td style="text-align: center;">188</td><td style="text-align: center;">10</td></tr> <tr><td style="text-align: center;">52</td><td style="text-align: center;">-72</td><td style="text-align: center;">5</td></tr> <tr><td style="text-align: center;">129</td><td style="text-align: center;">5</td><td style="text-align: center;">1</td></tr> <tr><td style="text-align: center;">236</td><td style="text-align: center;">112</td><td style="text-align: center;">8</td></tr> <tr><td style="text-align: center;">42</td><td style="text-align: center;">-82</td><td style="text-align: center;">6</td></tr> </tbody> </table>							Data	Diff = data - 124	Rank of diff	239	115	9	77	-47	3	179	55	4	221	97	7	100	-24	2	312	188	10	52	-72	5	129	5	1	236	112	8	42	-82	6	
Data	Diff = data - 124	Rank of diff																																							
239	115	9																																							
77	-47	3																																							
179	55	4																																							
221	97	7																																							
100	-24	2																																							
312	188	10																																							
52	-72	5																																							
129	5	1																																							
236	112	8																																							
42	-82	6																																							
	<p> $W_- = 3 + 2 + 5 + 6 = 16$ </p> <p>Refer to Wilcoxon single sample (/paired) tables for $n = 10$. Lower two-tail 10% point is 10. $16 > 10 \therefore$ Result is not significant.</p> <p>Seems there is no evidence against the median length being 124.</p>							<p>M1 M1 A1 B1 M1 M1A1 E1 E1</p>	<p>For differences. For ranks of difference . All correct. ft from here if ranks wrong.</p> <p>Or $W_+ = 9 + 4 + 7 + 10 + 1 + 8 = 39$</p> <p>No ft from here if wrong.</p> <p>Or, if 39 used, upper point is 45. No ft from here if wrong. Or $39 < 45$. ft only c's test statistic. ft only c's test statistic.</p>	9																															
								18																																	

4768: Statistics 3

General Comments

Once again the overall standard of the scripts seen was pleasing: many candidates appeared well prepared for the paper.

As reported previously, it was noticeable that candidates' empathy with the use of correct mathematical notation was often poor. For example: integrals were often written without the terminator "dx" and the symbols "=" and " \Rightarrow " were treated as synonymous. Also, despite a comment in last June's report, many candidates continue to show a lack of appreciation of the level of detail of arithmetic required to convince the examiner that an answer printed in the question has been obtained genuinely.

Invariably all four questions were attempted, and attempted well, on the whole. Questions 2 and 4 were found to be particularly high scoring. There was no evidence to suggest that candidates found themselves short of time at the end.

Comments on Individual Questions

1 Continuous random variables; Central Limit Theorem; duration of fireworks.

- (i) Almost all candidates got off to a good start here, experiencing no difficulty with the fairly straightforward integral that was involved.
- (ii) The mode was found correctly, though most candidates were seen to disregard the root $t = 0$ without comment.
- (iii) The value of $E(T)$ was found easily. For $\text{Var}(T)$ the layout and organisation of work was untidy at times, and all too often candidates were insufficiently careful about showing the printed answer convincingly.
- (iv) Most candidates were able to write down the correct distribution here, based on the Central Limit Theorem.
- (v) This was generally well answered. When candidates got it wrong it was usually because they constructed the interval either using a t value instead of a Normal value or using an incorrect alternative to the sample standard deviation. Most spotted that the mean of the model lay outside the interval, thus calling the model into question.

2 Combinations of Normal distributions; motorway toll charges.

- (i) This part was found to be very straightforward.
- (ii) This part, too, was well answered. It was pleasing to note that fewer candidates slipped up with the inequality of the requirement than in the past.

- (iii) Usually the mean of the total takings was correct, but the variance was correct less often. Typically the error came about through a lack of proper understanding of the difference between $\text{Var}(4X)$ ($= 4^2\text{Var}(X)$) and $\text{Var}(X_1 + \dots + X_4)$ ($=\text{Var}(X_1) + \dots + \text{Var}(X_4)$). In several cases the former was used when it should have been the latter. Furthermore the notation of the former was often seen when the subsequent working seemed to indicate that the latter was intended. That the weeks should be independent of each other was not as well known as would have been liked.
- (iv) There were many correct solutions to this part. As one might expect the errors that were seen all related to the calculation of the required variance.

3 The t distribution: paired test for the population mean difference; confidence interval for a population mean; absenteeism in the workplace.

- (a)(i) For candidates at this level too many seemed unable to set down clearly and concisely the hypotheses for this paired t test. It is reasonable to expect them to be familiar with the conventional notation " μ " for the population mean (difference, in this case) and to define it as such. For the necessary assumption, "Normality" on its own was not enough; candidates were expected to be explicit in naming the population of differences.
- (ii) The t test itself was usually carried out successfully. Candidates seemed well versed in what they had to do. There is still the issue of encouraging candidates to express their final conclusion in suitable language.
- (b) Most candidates answered this part well, apart from the required assumption. This time it was the Normality of the "days lost after" that was needed, and again candidates were expected to be explicit in identifying that population. The majority of candidates were able to provide the required interpretation of their interval in relation to the target. However, some carried over the mean and standard deviation from part (a), a consequence of which was that their confidence interval included a negative part which would be difficult to interpret in context. A few candidates thought that, since their interval was higher than the target value, then the target had been surpassed.

4 Chi-squared test of goodness of fit; Wilcoxon single sample test for a population median; distance between flaws in lengths of plastic strip.

- (i) This was well answered. Many candidates seemed to be making good use of their calculators, obtaining a correct value of the test statistic with little fuss. There were only occasional errors over the number of degrees of freedom and hence the critical value. As in Question 3 the language of the conclusion sometimes left room for improvement, and there were some who thought they were fitting data to a model rather than the other way round.
- (ii) Apart from a handful of candidates who ill advisedly attempted a t test, this part of the question was well answered. The work submitted was well organised and easy to follow.