

**BEng and MEng EXAMINATIONS 2014**

**PART II : STATISTICS (COMPUTING C245)**

**Date : Friday 16th May 2014 10 - 12**

*DO NOT OPEN THIS PAPER UNTIL THE INVIGILATOR TELLS YOU TO*

**Answer ALL questions**

*Statistical data sheets are provided*

*[Before starting, please make sure that the paper is complete; there should be a total of FOUR questions. Ask the invigilator for a replacement if your copy is faulty.]*

1. Choose one answer for each part. Partial credit may be awarded for working if an incorrect answer is selected. There is no negative marking.

- (i) For the data below, calculate the sample median,  $x_{(\{n+1\}/2)}$ ; mean,  $\bar{x} = \sum_{i=1}^n x_i/n$ ; interquartile range (IQR),  $x_{(3\{n+1\}/4)} - x_{(\{n+1\}/4)}$ ; and standard deviation,  $s_n = \sqrt{\sum_{i=1}^n (x_i - \bar{x})^2/n}$ .

0 2 3 3 4 7 9.

Which one of the following statements about these sample statistics is true?

- (a) The median is less than the mean and the IQR is less than the standard deviation;
  - (b) The mean is less than the median and the IQR is less than the standard deviation;
  - (c) The median is less than the mean and the standard deviation is less than the IQR;
  - (d) The mean is less than the median and the standard deviation is less than the IQR;
  - (e) The median is equal to the mean and the IQR is equal to the standard deviation.
- (ii) For two events  $A$  and  $B$ , we have  $P(A) = 0.2$ ,  $P(B) = 0.3$ , and  $P(\bar{A}|\bar{B}) = 0.8$ . What is  $P(\bar{B}|A)$  equal to?
- (a) 0.3;      (b) 0.4;      (c) 0.5;      (d) 0.6;      (e) 0.7;      (f) 0.8.
- (iii) A die is rolled three times. What is the probability that we observe at least one 3?
- (a) 0.5;      (b) 0.58;      (c) 0.42;      (d) 0.33;      (e) 0.67.
- (iv) A bag contains 5 red balls and 1 black ball. Balls are randomly drawn from the bag until the black ball is found, and the number of draws which were required is noted. How much greater is the expected number of draws required if any red ball drawn from the bag is put back in the bag rather than set aside?
- (a) 0;      (b) 1;      (c) 2;      (d) 2.5;      (e) 3.
- (v) Suppose  $X_1, X_2, \dots, X_n$  are  $n$  independent random variables which all follow the same distribution  $P_X$  which has mean  $\mu$  and variance  $\sigma^2$ . Let  $S_n = \sum_{i=1}^n X_i$  be the sum of the  $n$  samples, and define  $T_n = S_n/\sqrt{n}$ . By the Central Limit Theorem, which of the following is an approximate distribution for  $T_n$ ?
- (a)  $N(\sqrt{n}\mu, \sigma^2)$ ;   (b)  $N(\mu, n\sigma^2)$ ;   (c)  $N(\frac{\mu}{\sqrt{n}}, \sigma^2)$ ;   (d)  $N(\sqrt{n}\mu, \frac{\sigma^2}{\sqrt{n}})$ ;   (e)  $N(\mu, \frac{\sigma^2}{n})$ .

2. If  $X$  is a discrete random variable on the positive integers  $\{1, 2, 3, \dots\}$ , it can be shown that

$$E(X) = \sum_{x=1}^{\infty} P(X \geq x).$$

The Geometric distribution describes such a random variable. Let  $X \sim \text{Geometric}(p)$ , so that  $X$  has probability mass function

$$p(x) = p(1-p)^{x-1}, \quad x = 1, 2, 3, \dots$$

- (i) Find  $P(X \geq x)$  and hence show that  $E(X) = \frac{1}{p}$ . [*Hint: Recall that the infinite sum of a geometric progression with first term  $a$  and common ratio  $r$  is  $a/(1-r)$ .*]
  - (ii) Find the number of rolls of a die we would expect to require until we finally observe a 6.
  - (iii) Suppose that after five rolls of the die we have not yet observed a 6. If we decide to roll the die another five times, what is the probability that we will have observed a 6?
  - (iv) If  $X$  is a discrete random variable on the positive integers, prove the formula  $E(X) = \sum_{x=1}^{\infty} P(X \geq x)$ . [*Hint: Express  $P(X \geq x)$  as a summation, and then change the order of the two summations.*]
3. The number of emails sent to an IT helpdesk in a small department between the hours of 10 a.m. and 11 a.m. was recorded each day for a period of 100 working days. The data are summarised in the following table.

Number of emails	0	1	2	3	$\geq 4$
Frequency	14	39	29	11	7

Before collecting the data, it had been believed that the number of emails in an hour was well described by a Poisson distribution, and that two emails per hour were to be expected on average.

Recall that a Poisson random variable with mean parameter  $\lambda$  has probability mass function

$$p(x) = \frac{\lambda^x e^{-\lambda}}{x!}, \quad x = 0, 1, 2, \dots$$

- (i) For a sample  $X_1, \dots, X_n$  from the Poisson distribution, show that the sample mean is the maximum likelihood estimator (MLE) for  $\lambda$ .
- (ii) Calculate the value of the MLE for  $\lambda$  for the data above. When calculating the sample mean of the data, treat the category “ $\geq 4$ ” as taking fixed value 4.5.
- (iii) Using the MLE for  $\lambda$ , perform a goodness-of-fit test at the 95% significance level to assess whether the number of emails per hour does indeed follow a Poisson distribution.
- (iv) Suppose that instead of estimating  $\lambda$  we assume 2 emails per hour are received on average. Repeat the Poisson goodness-of-fit test under this assumption, and comment on the comparison.

[C245 2014]

4. Let  $X_1, \dots, X_n$  be a random sample of independent observations from the continuous uniform distribution  $U(0, b)$ , where the parameter  $b$  is unknown. The probability density function for each  $X_i$  is, thus,

$$f(x) = \begin{cases} \frac{1}{b}, & 0 \leq x \leq b \\ 0, & \text{otherwise.} \end{cases}$$

- (i) Find the cumulative distribution function,  $F(x) = P(X_i < x)$  for  $x \in [0, b]$ .
- (ii) Let  $X_{(n)} = \max_i X_i$  be the largest value in the sample. Find the cumulative distribution function  $P(X_{(n)} < x)$  for  $X_{(n)}$ . Hence derive the probability density function of  $X_{(n)}$ .  
[Hint: Notice that  $X_{(n)} < x \iff \forall i, X_i < x$ ].
- (iii) Calculate the expected value  $E(X_{(n)})$  as a function of  $b$  and  $n$ .
- (iv) Find the bias of the maximum sample value as an estimator for the upper limit of the population  $b$ . From this result, suggest a revised estimator of  $b$  using  $X_{(n)}$  which is unbiased. Remember in constructing your estimator that  $b$  is unknown.

END OF PAPER

	EXAMINATION QUESTIONS/SOLUTIONS 2013-2014	Course <b>Comp245</b>
Question 1.		Marks & seen/unseen
Parts	<p>(i) <math>x_{((n+1)/2)}=3, \bar{x}=4, \text{IQR}=x_{(6)} - x_{(2)}=5, s_n = 2\sqrt{2}</math>. Answer is <u>(c)</u>.</p> <p>(ii) Using Bayes Theorem,</p> $P(\bar{B} A) = \frac{P(A \cap \bar{B})}{P(A)} = \frac{P(A \bar{B})P(\bar{B})}{P(A)} = \frac{(1 - P(\bar{A} \bar{B}))(1 - P(B))}{P(A)}$ $= \frac{(1 - 0.8)(1 - 0.3)}{0.2} = 0.7.$ <p>Answer is <u>(e)</u>.</p> <p>(iii) The probability of not observing a 3 is <math>(5/6)^3 \approx 0.58</math>, so the probability of at least one 3 is <math>1 - 0.58 = 0.42</math>. Answer is <u>(c)</u>.</p> <p>(iv) Let <math>X</math> be the number of draws until the black ball is found. Without replacement, <math>X</math> is equally likely to be any number between 1 and 6, so <math>E(X)=3.5</math>. With replacement, <math>X \sim \text{Geometric}(1/6)</math>, so <math>E(X)=6</math>. Answer is <u>(d)</u>.</p> <p>(v) By the Central Limit Theorem, an approximate distribution for <math>S_n</math> is <math>N(n\mu, n\sigma^2)</math>. The distribution of <math>T_n = S_n / \sqrt{n}</math> is then, approximately, <math>N(\sqrt{n}\mu, \sigma^2)</math>. Answer is <u>(a)</u>.</p>	<div>seen sim.</div> <div>seen sim.</div> <div>seen sim.</div> <div>unseen</div> <div>unseen</div> <div>4 marks each</div>
	Setter's initials MM <div>Checker's initials</div>	Page number 1 of 4

	EXAMINATION QUESTIONS/SOLUTIONS 2013-2014	Course <b>Comp245</b>
Question 2.		Marks & seen/unseen
Parts	<p>(i) <math>P(X \geq x) = \sum_{k=x}^{\infty} p(k) = p \sum_{k=x}^{\infty} (1-p)^{k-1} = p \times \frac{(1-p)^{x-1}}{1-(1-p)} = (1-p)^{x-1}.</math></p> <p>We then have <math>E(X) = \sum_{x=1}^{\infty} P(X \geq x) = \sum_{x=1}^{\infty} (1-p)^{x-1} = \frac{1}{1-(1-p)} = \frac{1}{p}.</math></p> <p>(ii) If <math>X</math> is the number of rolls required, then <math>X \sim \text{Geom}\left(\frac{1}{6}\right)</math>, so <math>E(X) = \frac{1}{\frac{1}{6}} = 6.</math></p> <p>(iii) This is the same as the probability of failing to roll a 6 in five attempts; the first five unsuccessful rolls have no bearing on subsequent attempts (the Geometric is memoryless).</p> <p>The probability is <math>P(X \leq 10   X &gt; 5) = 1 - P(X \geq 11   X &gt; 5) = 1 - \left(1 - \frac{1}{6}\right)^5 = 0.598.</math></p> <p>(iv) <math>\sum_{x=1}^{\infty} P(X \geq x) = \sum_{x=1}^{\infty} \sum_{k=x}^{\infty} p(k) = \sum_{k=1}^{\infty} \sum_{x=1}^k p(k) = \sum_{k=1}^{\infty} kp(k) = E(X).</math></p>	<p>seen 4 marks</p> <p>unseen 4 marks</p> <p>seen sim. 3 marks</p> <p>seen sim. 4 marks</p> <p>unseen 5 marks</p>
	Setter's initials MM <span style="margin-left: 100px;">Checker's initials</span>	Page number 2 of 4

	EXAMINATION QUESTIONS/SOLUTIONS 2013-2014	Course <b>Comp245</b>																														
Question 3.		Marks & seen/unseen																														
Parts	<p>(i)</p> $L(\lambda) = \prod_{i=1}^n \frac{\lambda^{x_i} e^{-\lambda}}{x_i!} = \frac{\lambda^{\sum_{i=1}^n x_i} e^{-n\lambda}}{\prod_{i=1}^n x_i!} \implies \ell(\lambda) = \left(\sum_{i=1}^n x_i\right) \log(\lambda) - n\lambda - \log\left(\prod_{i=1}^n x_i!\right)$ $\implies \frac{d}{d\lambda} \ell(\lambda) = \frac{\sum_{i=1}^n x_i}{\lambda} - n \implies \hat{\lambda} = \frac{\sum_{i=1}^n x_i}{n} = \bar{x}.$ <p>This is a maximum since</p> $\implies \frac{d^2}{d\lambda^2} \ell(\lambda) = -\frac{\sum_{i=1}^n x_i}{\lambda^2}$ <p>is negative for all <math>\lambda &gt; 0</math>.</p> <p>(ii) <math>\hat{\lambda} = \bar{x} = 1.615</math>.</p> <p>(iii) Assuming a Poisson(1.61) distribution, we would have the following expected counts.</p> <table><tr><td>Number of emails (<math>i</math>)</td><td>0</td><td>1</td><td>2</td><td>3</td><td><math>\geq 4</math></td></tr><tr><td>Observed Frequency (<math>O_i</math>)</td><td>14</td><td>39</td><td>29</td><td>11</td><td>7</td></tr><tr><td>Expected Frequency (<math>E_i</math>)</td><td>19.89</td><td>32.12</td><td>25.94</td><td>13.96</td><td>8.09</td></tr></table> <p>These give a chi-square statistic of</p> $X^2 = \sum_{i=0}^4 \frac{(O_i - E_i)^2}{E_i} = 4.35$ <p>which is less than <math>\chi^2_{3,0.95} = 7.81</math>. So insufficient evidence to reject the null hypothesis that the email counts follow a Poisson distribution.</p> <p>(iv) Now assuming a Poisson(2) distribution, we would have the following expected counts.</p> <table><tr><td>Number of emails (<math>i</math>)</td><td>0</td><td>1</td><td>2</td><td>3</td><td><math>\geq 4</math></td></tr><tr><td>Expected Frequency (<math>E_i</math>)</td><td>13.53</td><td>27.07</td><td>27.07</td><td>18.04</td><td>14.29</td></tr></table> <p>Here,</p> $X^2 = \sum_{i=0}^4 \frac{(O_i - E_i)^2}{E_i} = 11.88$ <p>which is greater than <math>\chi^2_{4,0.95} = 9.49</math>, so we reject the null hypothesis that the email counts follow a Poisson(2) distribution. This suggests the mean number of emails per hour is not close to 2.</p>	Number of emails ( $i$ )	0	1	2	3	$\geq 4$	Observed Frequency ( $O_i$ )	14	39	29	11	7	Expected Frequency ( $E_i$ )	19.89	32.12	25.94	13.96	8.09	Number of emails ( $i$ )	0	1	2	3	$\geq 4$	Expected Frequency ( $E_i$ )	13.53	27.07	27.07	18.04	14.29	<div>seen</div> <div>7 marks</div> <div>3 marks</div> <div>seen sim.</div> <div>5 marks</div> <div>seen sim.</div> <div>5 marks</div>
Number of emails ( $i$ )	0	1	2	3	$\geq 4$																											
Observed Frequency ( $O_i$ )	14	39	29	11	7																											
Expected Frequency ( $E_i$ )	19.89	32.12	25.94	13.96	8.09																											
Number of emails ( $i$ )	0	1	2	3	$\geq 4$																											
Expected Frequency ( $E_i$ )	13.53	27.07	27.07	18.04	14.29																											
	Setter's initials MM	Page number 3 of 4																														
	Checker's initials																															

	EXAMINATION QUESTIONS/SOLUTIONS 2013-2014	Course <b>Comp245</b>
Question 4.		Marks & seen/unseen
Parts	<p>(i) <math>X_i \sim U(0, b)</math> has cumulative distribution function</p> $F_{X_i}(x) = \begin{cases} 0, & x \leq 0 \\ \frac{x}{b}, & 0 < x < b \\ 1, & x \geq b. \end{cases}$ <p>(ii) By independence of the <math>\{X_i\}</math>,</p> $F_{X_{(n)}}(x) = P(X_{(n)} < x) = \prod_{i=1}^n P(X_i < x) = \left(\frac{x}{b}\right)^n, \quad 0 \leq x \leq b.$ <p>The density is the first derivative of this distribution function, giving</p> $f_{X_{(n)}}(x) = \frac{d}{dx} F_{X_{(n)}}(x) = \frac{d}{dx} \left(\frac{x}{b}\right)^n = \frac{n}{b} \left(\frac{x}{b}\right)^{n-1}, \quad 0 \leq x \leq b.$ <p>(iii)</p> $E(X_{(n)} b) = \int_{-\infty}^{\infty} x f_{X_{(n)}}(x) dx = n \int_0^b \left(\frac{x}{b}\right)^n dx = \frac{n}{n+1} \frac{x^{n+1}}{b^n} \Big _0^b = \left(\frac{n}{n+1}\right)b.$ <p>(iv) As an estimator of <math>b</math>,</p> $\text{Bias}(X_{(n)}) = E(X_{(n)} b) - b = -\frac{b}{n+1}.$ <p>Since <math>E(X_{(n)} b) = \left(\frac{n}{n+1}\right)b</math>, consider a revised estimator <math>T = \left(\frac{n+1}{n}\right)X_{(n)}</math>. Then clearly <math>E(T b) = b</math>, so <math>T</math> is now an unbiased estimator for <math>b</math>.</p>	<p>seen 2 marks</p> <p>unseen 3 marks</p> <p>seen 4 marks</p> <p>seen 5 marks</p> <p>unseen 3 marks</p> <p>unseen 3 marks</p>
	Setter's initials MM	Page number 4 of 4