## SM-4331 Advanced Statistics Class Test 1

2019/20 Semester 2 18 February 2020 Time allowed: 90 minutes

## Instructions:

- There are three (3) questions totalling 20 marks and one (1) bonus question for 5 marks. The total attainable marks is 20 only.
- Answer ALL questions on a separate answer sheet.
- Ensure that you have written your name and student number on your answer sheets that you are submitting.
- The use of calculators is allowed.

| Question: | 1 | 2 | 3 | 4 | Total |
|-----------|---|---|---|---|-------|
| Marks:    | 8 | 4 | 7 | 5 | 20    |

1. Let  $X_1, \ldots, X_n$  be a random sample from an exponential distribution with the density function

$$f(x|\lambda) = \begin{cases} \lambda^{-1}e^{-x/\lambda} & x \ge 0\\ 0 & x < 0 \end{cases}$$

where  $\lambda > 0$  is an unknown parameter.

(a) (3 marks) Find the MLE  $\hat{\mu}$  for  $\mu$ .

# Solution:

(½ mark) The log-likelihood function from a single observation is given by

$$l(\lambda|X_i) = \log\left(\lambda^{-1}e^{-X_i/\lambda}\right)$$
$$= -\log\lambda - X_i/\lambda$$

(½ mark) Now,  $l(\lambda|X_1,\ldots,X_n) = \sum_{i=1}^n l(\lambda|X_i) = -n\log\lambda - \sum X_i/\lambda$ . (1 mark) Differentiating this with respect to  $\lambda$  yields  $l'(\lambda) = -n/\lambda + \sum X_i/\lambda^2$ .

(1 mark) Equating this to zero and solving for lambda, we obtain  $\hat{\lambda} = \bar{X}_n$ .

(b) (3 marks) Find the Cramér-Rao lower bound for the variance of the unbiased estimators for  $\lambda$ .

## Solution:

(ecf) Score function (from a single observation) is  $S(\lambda) = -1/\lambda + X_i/\lambda^2$ . (½ mark) Obtain the derivative of the score function  $S'(\lambda) = 1/\lambda^2 - 2X_i/\lambda^3$ .

(1 mark) The Fisher information is therefore

$$\mathcal{I}_{X_i}(\lambda) = \mathrm{E}\left(-S'(\lambda)\right) = -1/\lambda^2 + 2\,\mathrm{E}(X_i)/\lambda^3 = 1/\lambda^2.$$

(½ mark) The Full Fisher information is  $\mathcal{I}_{\mathbf{X}}(\lambda) = n/\lambda^2$ .

(1 mark) The Cramér-Rao bound is

$$\operatorname{Var}(\hat{\mu}) \ge \mathcal{I}_{\mathbf{X}}(\lambda)^{-1} = \lambda^2/n.$$

Alternatively, full marks if reason that the sample mean estimator attains its Cramér-Rao bound and has variance  $\text{Var}(X_i)/n = \lambda^2/n$ .

(c) (2 marks) Find the MLE for  $\theta = \lambda^2$ , and show that it is biased.

## Solution:

(1 mark) Using the invariance property of the MLE, the MLE for  $\theta = \lambda^2$  is  $\hat{\theta} = \bar{X}_n^2$ .

(1 mark) Using the fact that  $\operatorname{Var}(\bar{X}_n) = \operatorname{E}(\bar{X}_n^2) - \operatorname{E}^2(\bar{X}_n) = \lambda^2/n$ , the expectation of  $\hat{\theta}$  is

$$E(\hat{\theta}) = E(\bar{X}_n^2) = Var(\bar{X}_n) + E^2(\bar{X}_n)$$
$$= \lambda^2/n + \lambda^2$$
$$= \frac{n+1}{n} \cdot \lambda^2 \neq \lambda^2$$

Therefore, it is biased.

2. (a) (1 mark) Let  $X_n$ ,  $Y_n$ , X and Y be random variables, g a continuous function, and c a real-valued constant. All of the following are results of Slutzky's theorem as  $n \to \infty$ , **except one**. Which one of the below statements is not necessarily true?

A. If 
$$X_n \xrightarrow{P} X$$
 and  $Y_n \xrightarrow{P} Y$ ,  $X_n + Y_n \xrightarrow{P} X + Y$ .

B. If 
$$X_n \xrightarrow{P} X$$
 and  $Y_n \xrightarrow{P} Y$ ,  $X_n Y_n \xrightarrow{P} XY$ .

C. If 
$$X_n \xrightarrow{D} X$$
,  $g(X_n) \xrightarrow{D} g(X)$ .

**D.** If 
$$X_n \xrightarrow{\mathbf{D}} X$$
 and  $Y_n \xrightarrow{\mathbf{D}} Y$ ,  $X_n + Y_n \xrightarrow{\mathbf{D}} X + Y$ .

(b) (3 marks) Let  $\{X_n\}$  be a sequence of random variables. Suppose that  $\mathrm{E}(X_n) \to c$  and  $\mathrm{Var}(X_n) \to 0$  as  $n \to \infty$ , where c is a fixed constant. Using Markov's inequality, show that  $X_n$  converges to c in probability.

#### **Solution:**

(1½ mark) We note that  $\operatorname{Var}(X_n) = \operatorname{E}(X_n^2) - \operatorname{E}^2(X_n) \to 0$  as  $n \to \infty$ , and since  $\operatorname{E}(X_n) \to c$ ,  $\operatorname{E}^2(X_n) \to c^2$  using Slutzky's Theorem. Therefore,  $\operatorname{E}(X_n^2) \to c^2$  also.

(1½ mark) Markov's inequality states that

$$P(|X_n - c| > \epsilon) \le \frac{E(|X_n - c|^2)}{\epsilon^2}$$

$$= \frac{E(X_n^2) + c^2 - 2c E(X_n)}{\epsilon^2}$$

$$\to \frac{c^2 + c^2 - 2c \cdot c}{\epsilon^2} = 0$$

as  $n \to \infty$ . Therefore  $X_n \xrightarrow{P} c$ . Note that  $E(|X_n - c|^2) \neq Var(X_n)$  because  $E(X_n) \neq c$ ; equality is attained only as the sequence progresses and n gets larger and larger.

3. A forester wants to estimate the total number of farm acres planted in trees for a state. Because the number of acres of trees varies considerably with the size of the farm, he decides to stratify on farm sizes. The N=240 farms in the state are placed in one of four categories according to the size. A stratified random sample of n=40farms, selected by using proportional allocation, yields the results shown in Table 1.

| $n_h$      | Strata I | Strata II | Strata III | Strata IV |
|------------|----------|-----------|------------|-----------|
| 1          | 97       | 125       | 142        | 167       |
| 2          | 42       | 67        | 310        | 220       |
| 3          | 25       | 256       | 495        | 780       |
| 4          | 105      | 310       | 320        | 655       |
| 5          | 27       | 220       | 196        | 540       |
| 6          | 45       | 142       | 256        |           |
| 7          | 53       | 155       | 440        |           |
| 8          | 67       | 96        | 510        |           |
| 9          | 125      | 47        | 396        |           |
| 10         | 92       | 236       |            |           |
| 11         | 86       | 352       |            |           |
| 12         | 43       | 190       |            |           |
| 13         | 59       |           |            |           |
| 14         | 21       |           |            |           |
| $N_h$      | 86       | 72        | 52         | 30        |
| $\sum x$   | 887      | 2196      | 3065       | 2362      |
| $\sum x^2$ | 70131    | 501464    | 1178157    | 1405314   |

Table 1: Total farm acres of trees for the stratified sample of n = 40 farms.

(a) (2 marks) Calculate the sample mean and (unbiased) sample variance of the number of farm acres planted in trees within each strata.

(½ marks) 
$$\bar{X}_1 = 63.35$$
,  $\bar{X}_2 = 183$ ,  $\bar{X}_3 = 340.55$ ,  $\bar{X}_4 = 472.4$ .  
(1 mark)  $s^2 = \frac{1}{n-1} \sum_{i=1}^n (x_i - \bar{x})^2 = \frac{1}{n-1} \left( \sum_{i=1}^n x_i^2 - n\bar{x}^2 \right)$ .  
(½ marks)  $s_1^2 = 1071.8$ ,  $s_2^2 = 9054.2$ ,  $s_3^2 = 16794.3$ ,  $s_4^2 = 72376.3$ .

(b) (2 marks) Using your answer to (a), calculate the estimate of the mean number of farm acres planted in trees from the stratified sample. Give an estimate of the variance of this estimator.

# Solution:

(1 mark) The mean is

$$\bar{X}_{st} = \sum_{h=1}^{4} \omega_h \bar{X}_h$$

where  $\omega_h = N_h/N$ . Thus, plugging in the numbers, we get  $\bar{X}_{st} = 210.44$ . (1 mark) The formula for the variance of this estimator is

$$\operatorname{Var}(\bar{X}_{st}) = \sum_{h=1}^{4} \frac{N_h - n_h}{N_h n_h} S_h^2.$$

We can estimate this value by substituting the sample variances  $s_h^2$  for the population variances  $S_h^2$ . Plugging in the numbers, we get  $\widehat{\text{Var}}(\bar{X}_{st}) = 14298.6$ .

(c) (3 marks) Estimate the total number of acres of trees on farms in the state, and construct a suitable 95% confidence interval for this estimate.

# Solution:

(1 mark) We are now interested in  $\tau = N\mu$ , which is estimated using  $\hat{\tau} = N\bar{X}_{st}$ . Using the figures obtained earlier, we get  $240 \times 210.44 = 50,506$  acres. (1 mark) Note that assuming asymptotic normality,

$$\bar{X}_{st} \approx N(\mu, \widehat{Var}(\bar{X}_{st}))$$

and therefore,

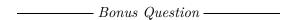
$$\hat{\tau} = N\bar{X}_{st} \approx N(N\mu, N^2 \widehat{\text{Var}}(\bar{X}_{st})).$$

(1 mark) A suitable 95% confidence interval for the total number of acress planted in trees is therefore  $\hat{\tau} \pm 1.96 \times N \sqrt{\widehat{\mathrm{Var}}(\bar{X}_{st})}$ , which is

$$50506 \pm 1.96 \times 240\sqrt{14298.63677} = (-5743, 106755)$$

But since there cannot be negative farm acres, the lower bound should be zero.

An alternative solution would be using the stratified totals. This would lead to the same estimate, in which the distribution of  $\hat{\tau}$  would be the sum of the normal distributions for each of the stratified totals.



4. (a) (1 mark) State, without proof, the pmf f(x) of a random variable X distributed according to  $X \sim \text{Bin}(n, p)$ .

#### Solution:

(one mark) 
$$f(x|p) = \binom{n}{x} p^x (1-p)^{n-x}$$
.

(b) (1 mark) What is the mean and variance of  $X \sim \text{Bin}(n, p)$ ?

#### Solution:

(½ mark) 
$$E(X) = np$$
  
(½ mark)  $Var(X) = np(1-p)$ 

(c) (2 marks) Let  $Y_1, \ldots, Y_n \stackrel{\text{iid}}{\sim} \operatorname{Bern}(p)$  be a random sample, and  $\hat{p} = X/n$  be an estimator for p, where  $X = \sum_{i=1}^n Y_i$ . Using the CLT, find the distribution of X.

# **Solution:**

(1 mark) By the CLT,  $\hat{p} \xrightarrow{\mathcal{D}} \mathcal{N}(p, p(1-p)/n)$  since  $\hat{p}$  is an estimator for the mean.

(1 mark) Therefore.

$$X = n\hat{p} \approx N(np, np(1-p)).$$

(d) (1 mark) What can you say about the distribution of  $X \sim \text{Bin}(n, p)$  as  $n \to \infty$ ?

## Solution:

(½ mark) We know that since  $X = \sum_{i=1}^{n} Y_i$  where each  $Y_i \sim \text{Bern}(p)$ , X must be distributed according to  $X \sim \text{Bin}(n, p)$ .

(½ mark) We notice that

$$X \approx N \left( \overbrace{np}, \overbrace{np(1-p)}^{\operatorname{Var}(X)} \right).$$

for large n. Therefore, for  $X \sim \text{Bin}(n,p)$  converges to N(np, np(1-p)) as n grows larger and larger. This is the normal approximation to the binomial distribution, typically used when n is larger than 30.