

1.  $[7 \times 1 = 7 \text{ marks}]$  **Circle the letter corresponding to the correct answer.**

(a) Consider an experiment in which data  $y_1, y_2, \dots, y_n$  are collected and the sample mean  $\bar{y} = \frac{1}{n} \sum_{i=1}^n y_i$  is used to estimate the mean of the population. Which of the following statements is FALSE?

A: The variability of the sampling distribution of the sample mean is affected by the sample size  $n$ .

B: The variability of the sampling distribution of the sample mean is affected by the standard deviation of the population.

☒ C: The location of the sampling distribution of the sample mean is affected by the sample size  $n$ .

D: How often the sample mean is within one unit of the population mean is affected by the standard deviation of the population.

E: The skewness of the sampling distribution of the sample mean is affected by the skewness of the population.

(b) Which of the following statements is TRUE?

A: An estimator  $\tilde{\theta}$  is a known quantity which can be used to estimate  $\theta$ .

☒ B: An estimator  $\tilde{\theta}$  is a random variable and its distribution is called the sampling distribution.

C: An estimate  $\hat{\theta}$  is a random variable and its distribution is called the sampling distribution.

D: The true value of the parameter  $\theta$  is known as soon as we have collected the data.

(c) The width of a 20% likelihood interval for  $\theta$  for data from a *Binomial*  $(n, \theta)$  distribution

☒ A: decreases as  $n$  increases.

B: increases as  $n$  increases.

C: does not change as  $n$  increases.

D: None of the above.

(d) A pivotal quantity for the unknown parameter  $\theta$  is

A: a function of the observed data  $y_1, y_2, \dots, y_n$ .

B: a function of  $Y_1, Y_2, \dots, Y_n$  whose distribution depends only on  $\theta$ .

☒ C: a function of  $Y_1, Y_2, \dots, Y_n$  and  $\theta$  whose distribution does not depend on  $\theta$ .

D: None of the above.

- (e) The width of a 95% confidence interval for  $\mu$  for data from a  $G(\mu, \sigma)$  distribution with known value of  $\sigma$
- A: does not change as  $n$  increases.
  - ☒ B: decreases as  $n$  increases.
  - C: increases as  $n$  increases.
  - ☐ D: increases as  $n$  decreases.
- (f) Suppose it is reasonable to model observed data  $y_1, y_2, \dots, y_n$  using a *Exponential* ( $\theta$ ) distribution. Let  $\text{thetahat} = \bar{y} = \frac{1}{n} \sum_{i=1}^n y_i$ . Which of the following functions in R should be used to calculate the relative likelihood function for  $\theta$ ?
- A: `ExpRLF <- function(x) {log((thetahat/x)^n*exp(n*(1-thetahat/x)))}`
  - B: `ExpRLF <- function(x) {log((x/thetahat)^n*exp(n*(1-x/thetahat)))}`
  - ☒ C: `ExpRLF<- function(x) {(thetahat/x)^n*exp(n*(1-thetahat/x))}`
  - D: `ExpRLF <- function(x) {(x/thetahat)^n*exp(n*(1-x/thetahat))}`
  - E: None of the above.
- (g) Suppose it is reasonable to model observed data  $y_1, y_2, \dots, y_{25}$  using a *Exponential* ( $\theta$ ) distribution. Suppose also that the maximum likelihood estimate of  $\theta$  is  $\hat{\theta} = 2$  is observed. Which of the following statements in R provides the maximum likelihood estimate of  $P(Y \leq 3)$  if  $Y \sim \text{Exponential}(\theta)$ ?
- A: `dexp(3,1/2)`
  - ☒ B: `pexp(3,1/2)`
  - C: `qexp(3,1/2)`
  - D: `rexp(3,1/2)`
  - E: None of the above.

2. [8 marks] Suppose  $y_1, y_2, \dots, y_n$  is an observed random sample from the distribution with probability density function

$$f(y; \theta) = \theta y^{\theta-1}, \quad 0 < y < 1, \quad \theta > 0.$$

(a) [4 marks] Derive the maximum likelihood estimate of  $\theta$  based on the data  $y_1, y_2, \dots, y_n$ . **Clearly show all your steps.**

The likelihood function is

$$L(\theta) = \prod_{i=1}^n f(y_i; \theta) = \prod_{i=1}^n \theta y_i^{\theta-1} = \theta^n \left( \prod_{i=1}^n y_i \right)^{\theta-1} \quad \text{for } \theta > 0$$

or more simply

$$L(\theta) = \theta^n \left( \prod_{i=1}^n y_i \right)^{\theta} \quad \theta > 0$$

The log likelihood is

$$l(\theta) = n \log \theta + \theta \log \left( \prod_{i=1}^n y_i \right) \quad \theta > 0$$

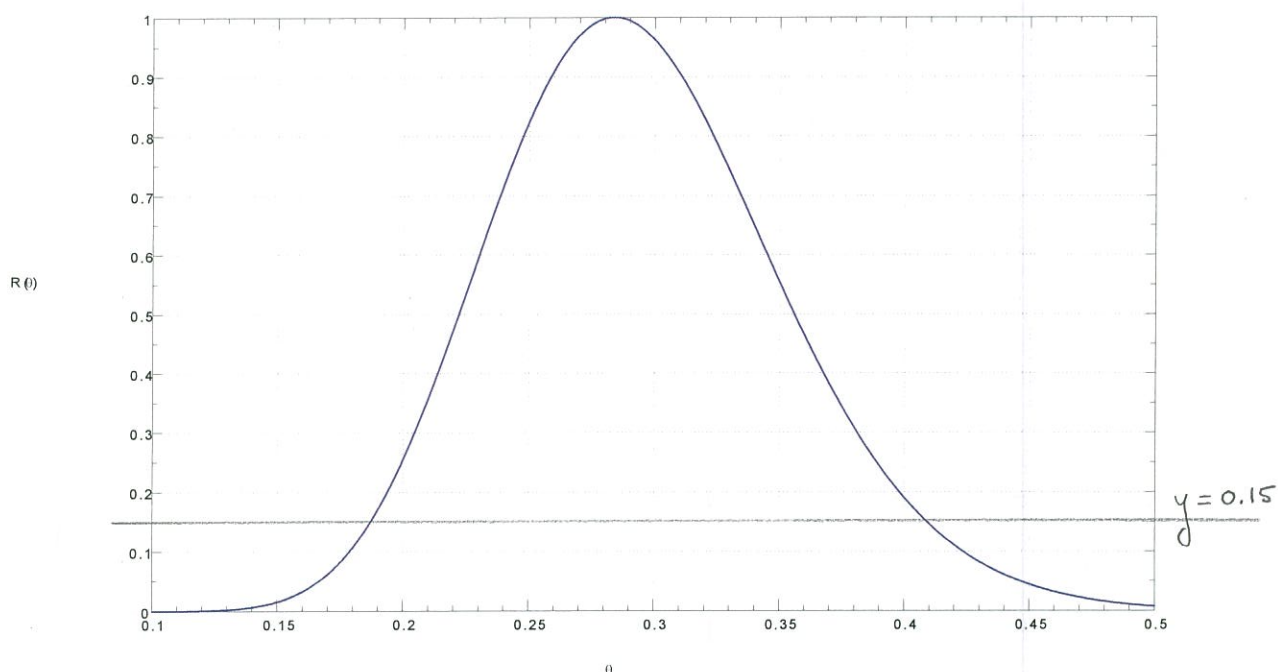
Solving

$$l'(\theta) = \frac{n}{\theta} + \log \left( \prod_{i=1}^n y_i \right) = \frac{n}{\theta} + \sum_{i=1}^n \log(y_i) = 0$$

gives the maximum likelihood estimate

$$\hat{\theta} = \frac{-n}{\log \left( \prod_{i=1}^n y_i \right)} = \frac{-n}{\sum_{i=1}^n \log(y_i)}$$

(b) [2 marks] The following is a plot of the relative likelihood function for  $\theta$  for a given set of data:



From the plot, graphically determine a 15% likelihood interval for  $\theta$ . Use two decimal places

$$[\underline{0.19}, \underline{0.41}]$$

(c) [2 marks] Complete the following R code to obtain the upper limit of the 15% likelihood interval for  $\theta$ :

```
RLF <- function(x) {exp(n*log(x/thetahat)+n*(1-x/thetahat))}
uniroot(function(x) RLF(x)-0.15, lower= 0.35, upper= 0.45)
```

Note: Lower limit must be a number  $\geq 0.3$  but  $< 0.41$ . The upper limit must be a number  $> 0.41$  but  $< 0.5$ .

3. [3 marks] In a study on the ability of rats to navigate a maze, a researcher records the length of time it took 21 different rats to find food at the end of the maze. Times in seconds are assumed to follow a  $G(\mu, \sigma)$  distribution. A previous study has shown that  $\sigma = 10.2$  seconds. Let  $y_i$  = time to complete the maze for the  $i$ th rat,  $i = 1, 2, \dots, 21$ . The researcher's data gave  $\sum_{i=1}^{21} y_i = 1068.4$  and  $\sum_{i=1}^{21} y_i^2 = 56453.7$ .

**Write your final answer only in the space provided.**

(a) [2 marks] A 90% confidence interval for  $\mu$  is (use 3 decimal places)

[ 47.215 , 54.538 ]

$$\bar{y} \pm 1.645 \times \frac{\sigma}{\sqrt{n}} = \frac{1068.4}{21} \pm 1.645 \left( \frac{10.2}{\sqrt{21}} \right) = 50.87619 \pm 3.661478 = [47.21471, 54.53767]$$

(b) [1 mark] Circle the letter corresponding to the best interpretation of the interval found in (a).

A. The probability that the true value of  $\mu$  is contained in the interval is 0.9.

☒ B. If the experiment is repeated many times, then we expect 90% of the intervals to contain the true value of  $\mu$ .

C. We are 90% confident that  $\mu = \hat{\mu}$ .

D. All of the above.

4. [ $7 \times 1 = 7$  marks] **Write your final answer only in the space provided.**

(a) **Without using Chi-squared tables** determine the following: (use 3 decimal places)

(i) If  $X \sim \chi^2(1)$  then  $P(X > 2.25) = \underline{0.134}$ .

$$P(X \geq 2.25) = 2P(Z \geq \sqrt{2.25}) = 2[1 - P(Z \leq 1.5)] = 2(1 - 0.93319) = 0.13362$$

(ii) If  $X \sim \chi^2(2)$  then  $P(X > 3.1) = \underline{0.212}$ .

$$P(X > 3.1) = e^{-3.1/2} = 0.212248$$

(b) **Using Chi-squared tables** determine the following: (use all decimal places available from the table)

(i) If  $X \sim \chi^2(18)$  then  $P(X > 27)$  lies between 0.05 and 0.1. (You must use values from the Chi-squared tables.)

$$P(X > 25.989) = 1 - 0.9 = 0.1$$

$$P(X > 28.869) = 1 - 0.95 = 0.05$$

(ii) If  $X \sim \chi^2(9)$  then the value of  $a$  such that  $P(X \leq a) = 0.05$  is  $a = \underline{3.325}$ .

(iii) If  $X \sim \chi^2(25)$  then the value of  $b$  such that  $P(X > b) = 0.025$  is  $b = \underline{40.646}$ .

(c) For the following questions specify **the distribution and its parameter(s)**:

(i) If  $X \sim G(-3, 2)$ ,  $Y \sim N(4, 16)$  and  $V \sim \text{Exponential}(2)$  independently then the distribution of

$$W = \left(\frac{X+3}{2}\right)^2 + \left(\frac{Y-4}{4}\right)^2 + V \text{ is } \underline{\chi^2(4)}.$$

$$\chi^2(1 + 1 + 2)$$

(ii) If  $X_i \sim \chi^2(i^2)$ ,  $i = 1, 2, 3$  independently then the distribution of  $\sum_{i=1}^3 X_i$  is  $\chi^2(14)$ .

$$\sum_{i=1}^3 i^2 = 1 + 4 + 9 = 14$$