

18.443 Problem Set 3 Spring 2015
Statistics for Applications
Due Date: 2/27/2015
prior to 3:00pm

Problems from John A. Rice, Third Edition. [*Chapter.Section.Problem*]

1. Problem 8.10.21.

Suppose that X_1, X_2, \dots, X_n are i.i.d. with density function

$$f(x | \theta) = \begin{cases} e^{-(x-\theta)}, & \text{if } x \geq \theta \\ 0, & \text{otherwise} \end{cases}$$

(a). Find the method of moments estimate of θ .

The first moment of X is

$$\begin{aligned} \mu_1 &= E[X] = \int_{\theta}^{\infty} x e^{-(x-\theta)} dx \\ &= \theta + \int_0^{\infty} y e^{-y} dy \\ &= \theta + [(y)(-e^{-y})]_{y=0}^{y=\infty} + \int_0^{\infty} e^{-y} dy \\ &= \theta + 1 \end{aligned}$$

(The second line follows by transforming to $y = x - \theta$; the third line follows from integration-by-parts.)

Equating the sample first moment to the population first moment:

$$\begin{aligned} \mu_1 &= \hat{\mu}_1 \\ \theta + 1 &= \frac{1}{n} \sum_{i=1}^n X_i = \bar{X} \\ \implies \hat{\theta} &= \bar{X} - 1 \end{aligned}$$

(b). Find the mle of θ . The likelihood of the data is

$$\begin{aligned} \text{lik}(\theta) &= f(X_1, \dots, X_n | \theta) \\ &= \prod_{i=1}^n f(X_i | \theta) \\ &= \prod_{i=1}^n [e^{-(X_i-\theta)} \mathbf{1}_{[\theta, \infty)}(X_i)] \\ &= [e^{-\sum_{i=1}^n (X_i-\theta)}] \prod_{i=1}^n [\mathbf{1}_{[0, X_i]}(\theta)] \\ &= [e^{-\sum_{i=1}^n X_i} e^{n\theta}] [\mathbf{1}_{[0, \min(X_1, \dots, X_n)]}(\theta)] \end{aligned}$$

$\text{lik}(\theta)$ is maximized by maximizing θ subject to $\theta \leq X_i$,

for all $i = 1, \dots, n$

i.e., $\hat{\theta}_{MLE} = \min(X_1, \dots, X_n)$

(c). Find a sufficient statistic for θ . Consider

$$T(X_1, \dots, X_n) = \min(X_1, \dots, X_n)$$

The distribution function of T , $F_T(t)$ satisfies

$$\begin{aligned} [1 - F_T(t)] &= P(T > t) \\ &= P(X_1 > t, X_2 > t, \dots, X_n > t) \\ &= \prod_{i=1}^n P(X_i > t) \\ &= \prod_{i=1}^n [e^{-(t-\theta)}] \\ &= [e^{-n(t-\theta)}] \end{aligned}$$

for values $t \geq \theta$.

The density of T is simply the derivative:

$$f_T(t | \theta) = ne^{-n(t-\theta)}, \quad t \geq \theta.$$

The conditional density of the sample given $T = t$ is

$$\begin{aligned} f(X_1, \dots, X_n | T, \theta) &= \frac{f(X_1, \dots, X_n | \theta)}{f_T(t | \theta)} \\ &= \frac{[e^{-\sum_{i=1}^n X_i} e^{n\theta}] [\mathbf{1}_{[0, \min(X_1, \dots, X_n)]}(\theta)]}{[e^{-n(t-\theta)}] \mathbf{1}_{[0, t]}(\theta)} \\ &= [e^{-\sum_{i=1}^n (X_i - t)}] \left[\prod_{i=1}^n \mathbf{1}_{[t, \infty)}(X_i) \right] \end{aligned}$$

The density function does not depend on θ , so

$$T = \min(X_1, \dots, X_n) \text{ is sufficient for } \theta.$$

2. Problem 8.10.45. **A Random walk Model for Chromatin**

The html in Rproject3.zip "Rproject3//Rproject3_rmd_rayleigh_theory.html" details estimation theory for a sample from a Rayleigh distribution.

(a). MLE of θ :

Data consisting of:

$$R_1, R_2, \dots, R_n$$

are i.i.d. *Rayleigh*(θ) random variables. The likelihood function is

$$\begin{aligned} \text{lik}(\theta) &= f(r_1, \dots, r_n | \theta) = \prod_{i=1}^n f(r_i | \theta) \\ &= \prod_{i=1}^n \left[\frac{r_i}{\theta^2} \exp\left(-\frac{r_i^2}{2\theta^2}\right) \right] \end{aligned}$$

The log-likelihood function is

$$\begin{aligned} \ell(\theta) &= \log[\text{lik}(\theta)] \\ &= \left[\sum_{i=1}^n \log(r_i) \right] - 2n \log(\theta) - \frac{1}{\theta^2} \sum_{i=1}^n [r_i^2 / 2] \end{aligned}$$

The mle solves $\frac{d}{d\theta}\ell(\theta) = 0$:

$$\begin{aligned} 0 &= \frac{d}{d\theta}(\ell(\theta)) \\ &= -2n\left(\frac{1}{\theta}\right) + 2\left(\frac{1}{\theta^3}\right) \sum_1^n [r_i^2/2] \\ \implies \hat{\theta}_{MLE} &= \left(\frac{1}{n} \sum_1^n [r_i^2/2]\right)^{1/2} \end{aligned}$$

(b). Method of moments estimate:

The first moment of the *Rayleigh*(θ) distribution is

$$\begin{aligned} \mu_1 &= E[R \mid \theta] = \int_0^\infty r f(r \mid \theta) dr \\ &= \int_0^\infty r \frac{r}{\theta^2} \exp\left(-\frac{r^2}{2\theta^2}\right) dr \\ &= \frac{1}{\theta^2} \int_0^\infty r^2 \exp\left(-\frac{r^2}{2\theta^2}\right) dr \\ &= \frac{1}{\theta^2} \int_0^\infty v \cdot \exp\left(-\frac{v}{2\theta^2}\right) \left[\frac{dv}{2\sqrt{v}}\right] \quad (\text{change of variables: } v = r^2) \\ &= \frac{1}{2\theta^2} \int_0^\infty v^{\frac{3}{2}-1} \cdot \exp\left(-\frac{v}{2\theta^2}\right) dv \\ &= \frac{1}{2\theta^2} \Gamma\left(\frac{3}{2}\right) (2\theta^2)^{\frac{3}{2}} \\ &= \sqrt{2}\theta \Gamma\left(\frac{3}{2}\right) = \sqrt{2}\theta \times \left(\frac{1}{2}\right) \Gamma\left(\frac{1}{2}\right) \\ &= \theta \times \frac{\sqrt{\pi}}{\sqrt{2}} \end{aligned}$$

(using the facts that $\Gamma(n+1) = n\Gamma(n)$ and $\Gamma(\frac{1}{2}) = \sqrt{\pi}$)

The MOM estimate solves:

$$\begin{aligned} \mu_1 &= \hat{\mu}_1 = \frac{1}{n} \sum R_i = \bar{R} \\ \theta \times \frac{\sqrt{\pi}}{\sqrt{2}} &= \bar{R} \\ \implies \hat{\theta}_{MOM} &= \bar{R} \times \frac{\sqrt{2}}{\sqrt{\pi}} \end{aligned}$$

(c). Approximate Variance of the MLE and method of moments estimate.

The approximate variance of the MLE is $Var(\hat{\theta}_{MLE}) \approx \frac{1}{nI(\theta)}$

where

$$\begin{aligned} I(\theta) &= E\left[-\frac{d^2}{d\theta^2}(\log(f(x \mid \theta)))\right] \\ &= E\left[-\frac{d^2}{d\theta^2}\left[\log\left(\frac{x}{\theta^2} \exp\left(-\frac{x^2}{2\theta^2}\right)\right)\right]\right] \\ &= E\left[-\frac{d}{d\theta}\left[-2\left(\frac{1}{\theta}\right) - \left(\frac{x^2}{2}\right)(-2)\theta^{-3}\right]\right] \\ &= E\left[-\left[\left(\frac{2}{\theta^2}\right) + (x^2)\right](-3)\theta^{-4}\right] \\ &= 3\theta^{-4}E[x^2] - \left(\frac{2}{\theta^2}\right) = 3\theta^{-4}(2\theta^2) - \left(\frac{2}{\theta^2}\right) \\ &= \frac{4}{\theta^2} \end{aligned}$$

So, $Var(\hat{\theta}_{MLE}) \approx \frac{\theta^2}{4n}$

Variance of the MOM estimate of Rayleigh Distribution Parameter:

The MOM estimate

$$\hat{\theta}_{MOM} = \bar{R} \times \frac{\sqrt{2}}{\sqrt{\pi}}$$

has variance:

$$Var(\hat{\theta}_{MOM}) = (\frac{\sqrt{2}}{\sqrt{\pi}})^2 Var(\bar{R}) = (\frac{2}{\pi}) \frac{Var(R)}{n}$$

$$\begin{aligned} Var(R) &= E[R^2] - (E[R])^2 \\ &= 2\theta^2 - (\sqrt{\frac{\pi}{2}}\theta)^2 \\ &= \theta^2(2 - \frac{\pi}{2}) \end{aligned}$$

So, $Var(\hat{\theta}_{MOM}) = \theta^2(2 - \frac{\pi}{2})(\frac{2}{\pi})(\frac{1}{n}) = \theta^2(\frac{4}{\pi} - 1)(\frac{1}{n}) \approx \frac{\theta^2}{n} \times 0.2732$

This exceeds the approximate $Var(\hat{\theta}_{MLE}) \approx \frac{\theta^2}{n} \times 0.25$

See the R script file:

Rproject3_script4_Chromatin_solution.r

3. Problem 8.10.51 **Double Exponential (Laplace) Distribution**

The double exponential distribution is

$$f(x | \theta) = \frac{1}{2}e^{|x-\theta|}, -\infty < x < \infty.$$

For an iid sample of size $n = 2m + 1$, show that the mle of θ is the median of the sample.

Let X_1, \dots, X_n denote the sample random variables with outcomes x_1, \dots, x_n . The likelihood function of the data is

$$\begin{aligned} lik(\theta) &= \prod_{i=1}^n f(x_i | \theta) = \prod_{i=1}^n [\frac{1}{2}e^{-|x_i-\theta|}] \\ &= (\frac{1}{2})^n e^{-\sum_{i=1}^n |x_i-\theta|} \end{aligned}$$

This is maximized by minimizing the sum in the exponent:

$$g(\theta) = \sum_{i=1}^n |x_i - \theta|$$

Note that $g(\theta)$ is a continuous function of θ and its derivative exists at all points θ that are not equal to any x_i

$$\begin{aligned} g'(\theta) &= \frac{d}{d\theta} g(\theta) = \sum_{i=1}^n [-1 \times \mathbf{1}(x_i > \theta) + (+1) \times \mathbf{1}(x_i < \theta)] \\ &= (-1) \times [\sum_{i=1}^n \mathbf{1}(x_i > \theta)] + (+1) \times \sum_{i=1}^n \mathbf{1}(x_i < \theta) \\ &= \begin{cases} \text{positive} & \text{if } \theta > \text{median}(x_i) \\ \text{negative} & \text{if } \theta < \text{median}(x_i) \end{cases} \end{aligned}$$

It follows that $g(\theta)$ is minimized at $\theta = \text{median}(x_i)$. A graph of $g(\theta)$ is piecewise linear with slope changes at each of the x_i values; the slope at any given θ (not equal to an x_j) is

$$\text{count}(x_i < \theta) - \text{count}(x_i > \theta).$$

4. **Problem 8.10.58 Gene Frequencies of Haptoglobin Type**

Gene frequencies are in equilibrium, the genotypes AA , Aa , and aa occur with probabilities $(1 - \theta)^2$, $2\theta(1 - \theta)$, and θ^2 . Plato et al. published the following data on Haptoglobin Type in a sample of 190 people

Haptoglobin Type		
Hp1-1	Hp1-2	Hp2-2
10	68	112

This is precisely the same problem as Example 8.5.1.A of the text and class notes which corresponds to count data: $(X_1, X_2, X_3) \sim \text{Multinomial}(n = 3, p = ((1 - \theta)^2, 2\theta(1 - \theta), \theta^2))$ distribution.

(a). Find the mle of θ

- $(X_1, X_2, X_3) \sim \text{Multinomial}(n, p = ((1 - \theta)^2, 2\theta(1 - \theta), \theta^2))$

- Log Likelihood for θ

$$\begin{aligned} \ell(\theta) &= \log(f(x_1, x_2, x_3 | p_1(\theta), p_2(\theta), p_3(\theta))) \\ &= \log\left(\frac{n!}{x_1!x_2!x_3!} p_1(\theta)^{x_1} p_2(\theta)^{x_2} p_3(\theta)^{x_3}\right) \\ &= x_1 \log((1 - \theta)^2) + x_2 \log(2\theta(1 - \theta)) \\ &\quad + x_3 \log(\theta^2) + (\text{non-}\theta \text{ terms}) \\ &= (2x_1 + x_2) \log(1 - \theta) + (2x_3 + x_2) \log(\theta) + (\text{non-}\theta \text{ terms}) \end{aligned}$$

- First Differential of log likelihood:

$$\ell'(\theta) = -\frac{(2x_1 + x_2)}{1 - \theta} + \frac{(2x_3 + x_2)}{\theta}$$

$$\implies \hat{\theta} = \frac{2x_3 + x_2}{2x_1 + 2x_2 + 2x_3} = \frac{2x_3 + x_2}{2n} = \frac{2(112) + 68}{2(190)} = 0.76842$$

(b). Find the asymptotic variance of the mle.

- $\text{Var}(\hat{\theta}) \longrightarrow \frac{1}{E[-\ell''(\theta)]}$

- Second Differential of log likelihood:

$$\begin{aligned}\ell''(\theta) &= \frac{d}{d\theta} \left[-\frac{(2x_1 + x_2)}{1 - \theta} + \frac{(2x_3 + x_2)}{\theta} \right] \\ &= -\frac{(2x_1 + x_2)}{(1 - \theta)^2} - \frac{(2x_3 + x_2)}{\theta^2}\end{aligned}$$

- Each of the X_i are *Binomial*($n, p_i(\theta)$) so

$$\begin{aligned}E[X_1] &= np_1(\theta) = n(1 - \theta)^2 \\ E[X_2] &= np_2(\theta) = n2\theta(1 - \theta) \\ E[X_3] &= np_3(\theta) = n\theta^2\end{aligned}$$

- $E[-\ell''(\theta)] = \frac{2n}{\theta(1 - \theta)}$
- $\hat{\sigma}_{\hat{\theta}}^2 = \frac{\hat{\theta}(1 - \hat{\theta})}{2n} = \frac{0.76842(1 - 0.76842)}{2 \times 190} = 0.0004682898 = (.02164)^2$

Parts (c), (d), and (e): see the R script

Rproject3_script1_multinomial_simulation_Problem_8_57.r

MIT OpenCourseWare
<http://ocw.mit.edu>

18.443 Statistics for Applications
Spring 2015

For information about citing these materials or our Terms of Use, visit: <http://ocw.mit.edu/terms>.