

# CSE 544, Spring 2017, Probability and Statistics for Data Science

## Assignment 3: Statistical Inference

(8 questions, 70 points total)

Due: 3/27, in class

I/We understand and agree to the following:

- (a) Academic dishonesty will result in an 'F' grade and referral to the Academic Judiciary.
- (b) Late submission, beyond the 'due' date/time, will result in a score of 0 on this assignment.

(write down the name of all collaborating students on the line below)

### 1. MME for Gamma distribution

(Total 5 points)

The Gamma( $x, y$ ) distribution has mean  $x + y$  and variance  $x + y^2$ . Find the MME for  $x$  and  $y$ .

$$\text{Gamma}(x, y) \rightarrow \begin{array}{l} \text{mean} = xy = E[X] \\ \text{Variance} = xy^2 = \text{Var}[X] \end{array}$$

Equating  $E[X]$  &  $\text{Var}[X]$  to corresponding sample mo-  
exp. var

we get,  $E[X] = \frac{\sum X_i}{n} = xy$

$$\text{Var}[X] = \frac{\sum (X_i - \bar{X})^2}{n} = xy^2$$

$$x = \frac{\bar{X}}{y}$$

$$xy^2 = \frac{\sum (X_i - \bar{X})^2}{n}$$

$$ny = \bar{X}$$

MME of  $y$

$$\hat{y} = \frac{\sum (X_i - \bar{X})^2}{n \bar{X}}$$

MME of  $x$

$$\hat{x} = \frac{\bar{X}}{\hat{y}} = \frac{n \bar{X}^2}{\sum (X_i - \bar{X})^2}$$

## 2. Properties of estimators

(Total 5 points)

Recall the three properties of estimators introduced in class: (i)  $\text{bias}(\hat{\theta}) = E[\hat{\theta}] - \theta$ , (ii)  $\text{se}(\hat{\theta}) = \sqrt{\text{Var}[\hat{\theta}]}$ ,

and (iii)  $\text{MSE}(\hat{\theta}) = E[(\hat{\theta} - \theta)^2]$ . Find these quantities for  $\hat{\theta} = \frac{1}{n} \sum_{i=1}^n X_i$ , where  $X_i \sim \text{Poisson}(\theta)$ .

$$(i) \text{bias}(\hat{\theta}) = E[\hat{\theta}] - \theta$$

$$= E\left(\frac{1}{n} \sum X_i\right) - \theta$$

$$= \frac{1}{n} \sum_{i=1}^n E[X_i] - \theta = \frac{\sum E[X_i]}{n} - \theta$$

$$= \frac{n \cdot 1}{n} - 1 = 1 - 1 = 0$$

( $E[X_i]$  of Poisson = 1)

$$\hat{\theta} = \frac{1}{n} \sum X_i$$

where  $X_i \sim \text{Poisson}$

known case  
 $\theta = 1$

$$= \theta - \theta$$

$$= 0$$

$$(ii) \text{se}(\hat{\theta}) = \sqrt{\text{Var}(\hat{\theta})}$$

$$= \sqrt{\text{Var}\left(\frac{1}{n} \sum X_i\right)}$$

$$= \sqrt{\frac{1}{n^2} \sum \text{Var}(X_i)}$$

$$= \sqrt{\frac{\theta}{n}}$$

( $\text{Var of Poisson} = 1$ )

$$(iii) \text{MSE}(\hat{\theta}) = E[(\hat{\theta} - \theta)^2]$$

$$= E[\hat{\theta}^2] - 2\theta E[\hat{\theta}] + E[\theta^2]$$

$$= E[\hat{\theta}^2] - E[\hat{\theta}]^2 + E[\hat{\theta}]^2 - 2\theta E[\hat{\theta}] + \theta^2$$

$$= \text{Var}(\hat{\theta}) + (E[\hat{\theta}] - \theta)^2$$

$$= \frac{\theta}{n} + (\text{bias}(\hat{\theta}))^2$$

$$= \frac{\theta}{n}$$

$$= \frac{\theta}{n} + 0$$

$$= \frac{\theta}{n}$$

$$E[(\hat{\theta} - \theta)^2]$$

### 3. Plug-in estimates

(a) The kurtosis for a random variable  $X$  is defined as  $Kurt[X] = E[(X - \mu)^4] / \sigma^4$ . Derive the plug-in estimate of the kurtosis. (Total 10 points) (3 points)

(b) Show that the plug-in estimator of the variance of  $X$  is  $\hat{\sigma}^2 = \frac{1}{n} \sum_{i=1}^n (X_i - \bar{X}_n)^2$ , where  $\bar{X}_n$  is the sample mean,  $\bar{X}_n = \frac{1}{n} \sum_{i=1}^n X_i$ . (3 points)

(c) Find the bias of  $\hat{\sigma}^2$  in terms of the true variance,  $\sigma^2$ . (4 points)

$$(a) Kurt[X] = \frac{E[(X - \mu)^4]}{\sigma^4}$$

$$E[X] = \int x f(x) dx$$

$$\hat{\mu} = \frac{\sum X_i}{n} = \bar{X}$$

$$\hat{\sigma}^2 = \frac{1}{n} \sum (X_i - \hat{\mu})^2$$

$$= \frac{\int (x - \mu)^4 dF(x)}{\left( \int (x - \mu)^2 dF(x) \right)^2}$$

$$= \frac{\sum (X_i - \hat{\mu})^4}{n^3 \hat{\sigma}^4}$$

Remove  $\frac{1}{n} \rightarrow \frac{\sum}{n}$

$$= \frac{\sum (X_i - \hat{\mu})^4}{n^3 \hat{\sigma}^4} = \frac{\sum (X_i - \bar{X})^4}{\left( \sum (X_i - \bar{X})^2 \right)^2}$$

$$(b) \hat{\sigma}^2 = Var(X) = \int x^2 d\hat{F}_n(x) - \left( \int x d\hat{F}_n(x) \right)^2$$

$$= E[X^2] - E[X]^2 = \frac{1}{n} \sum X_i^2 - \left( \frac{\sum X_i}{n} \right)^2$$

$$= \frac{1}{n} \sum X_i^2 - \bar{X}^2$$

$$= \frac{1}{n} \sum X_i^2 - 2\bar{X}^2 + \bar{X}^2$$

$$= \frac{1}{n} \sum X_i^2 - \frac{2}{n} \left( \sum X_i \right) \bar{X} + \frac{\bar{X}^2}{n}$$

$$= \frac{1}{n} \sum (X_i - \bar{X})^2$$



$$\begin{aligned}
(c) E[\hat{\sigma}^2] &= E\left[\frac{1}{n} \sum (X_i - \bar{X})^2\right] \\
&= E\left[\frac{1}{n} \sum (X_i - \mu) - (\bar{X} - \mu)\right]^2 \\
&= E\left[\frac{1}{n} \sum (X_i - \mu)^2 - \frac{2}{n^2} \sum (\bar{X} - \mu) \sum (X_i - \mu) + \frac{1}{n} \sum (\bar{X} - \mu)^2\right] \\
&= E\left[\frac{1}{n} \sum (X_i - \mu)^2 - \frac{2}{n} (\bar{X} - \mu) (\sum X_i - n\mu) + \frac{1}{n} n (\bar{X} - \mu)^2\right] \\
&= E\left[\frac{1}{n} \sum (X_i - \mu)^2 - 2(\bar{X} - \mu)^2 + (\bar{X} - \mu)^2\right] \\
&= E\left[\frac{1}{n} \sum (X_i - \mu)^2 - (\bar{X} - \mu)^2\right] \\
&= E\left[\frac{\sum (X_i - \mu)^2}{n}\right] - E[(\bar{X} - \mu)^2] \\
&= \sigma^2 - \text{Var}(\bar{X}) \\
&= \sigma^2 - \text{Var}\left(\frac{\sum X_i}{n}\right) \\
&= \sigma^2 - \frac{1}{n^2} \sum \text{Var} X_i \\
&= \sigma^2 - \frac{\sigma^2}{n}
\end{aligned}$$

$$E[\hat{\sigma}^2] - \sigma^2 = -\frac{\sigma^2}{n}$$

#### 4. Functionals of the empirical distribution function

(Total 10 points)

Let  $X_1, X_2, \dots, X_n$  be i.i.d. RVs with true CDF  $F$ . Let  $\hat{F}$  be their empirical distribution function. Let  $a$  and  $b$  be some numbers with  $a < b$ . Define  $\theta = F(b) - F(a)$ . Let  $\hat{\theta} = \hat{F}_n(b) - \hat{F}_n(a)$ .

(a) Find  $se(\hat{\theta})$

(5 points)

(b) Find an approximate  $(1-\alpha)$  CI for  $\hat{\theta}$ .

(5 points)

$$\theta = F(b) - F(a)$$

$$\hat{\theta} = \hat{F}_n(b) - \hat{F}_n(a) = \frac{1}{n} \sum 1(X_i \in (a, b])$$

$$(a) \text{Var}(\hat{\theta}) = \text{Var}\left(\frac{1}{n} \sum 1(X_i \in (a, b])\right)$$

$$= \frac{1}{n^2} \text{Var}\left(\sum 1(X_i \in (a, b])\right)$$

$$= \frac{1}{n} \text{Var}(1(X_i \in (a, b]))$$

$$= \frac{1}{n} (F(b) - F(a))(1 - F(b) + F(a))$$

$$= \frac{\theta(1-\theta)}{n}$$

$$se = \sqrt{\frac{\theta(1-\theta)}{n}}$$

$$\text{Var}(\hat{F}_n) = \frac{F(n)(1-F(n))}{n}$$

$$(b) P\left(-z_{\alpha} \leq \frac{\hat{\theta} - E[\hat{\theta}]}{se(\hat{\theta})} \leq z_{\alpha}\right)$$

$$= P\left(-z_{\alpha} \leq \frac{\hat{\theta} - \theta}{\sqrt{\frac{\theta(1-\theta)}{n}}} \leq z_{\alpha}\right)$$

$$= P\left(-z_{\alpha} \sqrt{\frac{\theta(1-\theta)}{n}} + \theta \leq \hat{\theta} \leq z_{\alpha} \sqrt{\frac{\theta(1-\theta)}{n}} + \theta\right)$$

CR for  $\hat{\theta}$  is

$$L(N) = \theta - z_{\alpha} \sqrt{\frac{\theta(1-\theta)}{n}}$$

$$R(N) = \theta + z_{\alpha} \sqrt{\frac{\theta(1-\theta)}{n}}$$

### 5. Histogram, meet the empirical distribution function

(Total 5 points)

Let  $X_1, X_2, \dots, X_n$  be i.i.d. RVs with true CDF  $F$  with sample space  $[0, 1]$ . Let  $\hat{F}_n$  be the associated empirical distribution function. Let  $\hat{f}_n$  be the empirical pdf based on a histogram on the range  $[0, 1]$  with some bin size  $h$ , as in class. For some  $x$  in the range  $(0, 1)$ , show that  $\hat{f}_n(x) \approx d\hat{F}_n(x)$ . Use the fact that the derivative of a function,  $g()$ , at  $x$ , is  $\lim_{\Delta x \rightarrow 0} \frac{g(x+\Delta x) - g(x)}{\Delta x}$ .

Histogram ps

$$\hat{f}_n(x) = \frac{\sum 1(X_i \in (r_k, r_{k+1}])}{n(r_{k+1} - r_k)}$$

for  $x \in (r_k, r_{k+1}]$

$$\begin{aligned} \hat{f}_n(x) &= \frac{1}{nh} \sum 1(X_i \in (r_k, r_{k+1}]) \\ &= \frac{1}{nh} \sum (1(X_i \leq r_{k+1}) - 1(X_i \leq r_k)) \\ &= \frac{1}{nh} \sum (1(X_i \leq x+h) - 1(X_i \leq x)) \\ &= \frac{1}{n} \sum 1(X_i \leq x+h) - \frac{1}{n} \sum 1(X_i \leq x) \\ &= \frac{\hat{F}_n(x+h) - \hat{F}_n(x)}{h} \\ &\approx d\hat{F}_n(x) \end{aligned}$$



6. Normal-based CI for  $\hat{F}_n$

(Total 10 points)

$X_1, X_2, \dots, X_n$  be i.i.d. RVs with true CDF  $F$ . Let  $\hat{F}_n$  be their empirical distribution function.

(a) Derive  $\text{Var}[\hat{F}_n]$ .

(4 points)

(b) Derive the Normal-based  $(1-\alpha)$  CI for  $\hat{F}_n$ .

(6 points)

(a)  $n$  i.i.d. Bernoulli RV,  $\text{sum} = \sum I(X_k \leq x)$   
 $\rightarrow$  success prob =  $F(x)$

$\Rightarrow n\hat{F}_n(x) = \text{Binomial RV}$

$\Rightarrow$  Variance of Binomial Distribution =  $np(1-p)$

$$\Rightarrow \text{Var}(\hat{F}_n) = \frac{F(x)(1-F(x))}{n} = \text{Var}(\hat{F}_n)$$

(b)  $(1-\alpha)$  CI for  $\hat{F}_n$

$$\hat{F}_n(x) \pm z_{\alpha/2} \sqrt{\frac{\hat{F}_n(x)(1-\hat{F}_n(x))}{n}}$$