

## Problem Set 4

### Submission:

Thursday, 03/03/2022, until 1 PM, to be uploaded on the NYU Brightspace course homepage.

### 1. A brief reminder on conditional distributions

[4 Points]

In this problem we recall the notion of **conditional distributions** for two random variables  $X$  and  $Y$ . This is in particular relevant for the notion of sufficient statistics.

- If  $X$  and  $Y$  are discrete random variables with values in  $\Omega_X$  and  $\Omega_Y$  respectively, then for any  $y \in \Omega_Y$  with  $\mathbf{P}[Y = y] > 0$ , the conditional distribution of  $X$  given  $Y = y$  is characterized by the conditional probability mass function

$$p_{X|Y=y}(x) = \frac{\mathbf{P}[X = x, Y = y]}{\mathbf{P}[Y = y]}, \quad x \in \Omega_X.$$

This is the distribution of  $X$  under the probability measure  $\mathbf{P}[\cdot | Y = y]$ .

- If  $X$  and  $Y$  are continuous real random variables, then for any  $y \in \mathbb{R}$  with  $f_Y(y) > 0$ , the conditional distribution of  $X$  given  $Y = y$  is characterized by the conditional probability density function

$$f_{X|Y=y}(x) = \frac{f_{X,Y}(x, y)}{f_Y(y)}, \quad x \in \mathbb{R},$$

where  $f_{X,Y}$  is the joint probability density function of  $(X, Y)$  and  $f_Y$  is the probability density function of  $Y$ .

- (a) A fair coin is tossed 4 times. Let  $X$  denote the number of times that heads comes up and  $Y = 1$  if heads comes up on the first toss and  $Y = 0$  otherwise. Determine the conditional distribution of  $X$  given  $Y = 1$  and the conditional distribution of  $Y$  given  $X = 3$ .
- (b) Consider jointly continuous random variables  $X, Y$  with density

$$f_{X,Y}(x, y) = 4ye^{-2y(x+1)} \mathbb{1}_{\{x, y > 0\}}.$$

Determine the conditional probability density function  $f_{X|Y=y}$ . What is the conditional distribution of  $X$  given  $Y = y$ ?

### 2. Sufficient statistics

[4 Points]

- (a) Let  $X_1, \dots, X_n \stackrel{\text{i.i.d.}}{\sim} \text{Pois}(\lambda)$  with unknown parameter  $\lambda > 0$ . Show that  $T(\mathbf{X}) = \sum_{j=1}^n X_j$  is a sufficient statistic for  $\lambda$

(i) ...directly, using the definition of sufficiency.

(ii) ...using the Neyman-characterization of sufficiency.

*Hint:* For (i), use the result of Problem 2 (a) on Problem set 1. You need to calculate  $\mathbf{P}[X_1 = x_1, \dots, X_n = x_n | T(\mathbf{X}) = t]$  for all possible values of  $x_1, \dots, x_n, t \in \mathbb{N}_0$ .

- (b) Let  $X_1, \dots, X_n \stackrel{\text{i.i.d.}}{\sim} \text{Pareto}(\lambda, a)$  with known  $a > 0$  and unknown  $\lambda > 0$ , where we say that  $X \sim \text{Pareto}(\lambda, a)$  if

$$f_X(x) = \frac{\lambda a^\lambda}{x^{\lambda+1}} \mathbb{1}_{(a, \infty)}(x).$$

Find a sufficient statistic  $T(\mathbf{X}) \in \mathbb{R}$  for  $\lambda$ , using the Neyman-characterization.

### 3. Method of moments

[4 Points]

- Consider  $X_1, \dots, X_n \stackrel{\text{i.i.d.}}{\sim} \mathcal{U}([0, \theta])$  with unknown  $\theta > 0$ . Calculate an estimator  $\hat{\theta}_n$  for  $\theta$  based on the method of moments. Check this estimator for consistency and unbiasedness.
- Consider  $X_1, \dots, X_n \stackrel{\text{i.i.d.}}{\sim} \text{Bin}(N, p)$ , where  $\theta = (N, p) \in \mathbb{N} \times (0, 1)$  is unknown (meaning both  $N$  and  $p$  are unknown). Determine an estimator  $(\hat{N}_n, \hat{p}_n)$  for  $(N, p)$  based on the method of moments.  
*Hint:* You need both  $\mathbf{E}_{(N,p)}[X_1]$  and  $\mathbf{E}_{(N,p)}[X_1^2]$ .

### 4. (R exercise) The empirical distribution

[4 Points]

*Note: Please provide your source code and images obtained with your solution.*

Suppose that  $X_1, \dots, X_n$  are i.i.d. real random variables such that  $X_1$  has cumulative distribution function  $F$ . Given a realization of these random variables, we can consider the *empirical cumulative distribution function*

$$\hat{F}_n(x) = \frac{1}{n} \sum_{i=1}^n \mathbb{1}_{\{X_i \leq x\}}.$$

- Explain why  $\hat{F}_n(x)$  fulfills the following:
  - For every  $x \in \mathbb{R}$ , one has  $\mathbf{E}[\hat{F}_n(x)] = F(x)$  and  $\text{Var}[\hat{F}_n(x)] = \frac{1}{n} F(x)(1 - F(x))$ .
  - For every  $x \in \mathbb{R}$ , one has  $\hat{F}_n(x) \xrightarrow[n \rightarrow \infty]{\mathbf{P}} F(x)$ .

*Hint:* For any event  $A$  in some probability space, one has  $\mathbf{E}[\mathbb{1}_A] = \mathbf{P}[A]$ .

- Use R to generate  $n \in \{5, 20, 500\}$  samples of i.i.d. random variables  $X_1, \dots, X_n$  following a  $\mathcal{U}([0, 1])$ -distribution or a  $\mathcal{E}(3)$ -distribution. Plot the empirical cumulative distribution function together with the graph of the cumulative distribution function  $F_{X_1}$ . What do you observe?

*Hint:* The R-command `ecdf()` calculates the empirical distribution function of a vector and `plot(ecdf())` plots the respective graph.

- Data on the magnitudes of earthquakes near Fiji are available from R, using the command `quakes`.<sup>1</sup> For help on this dataset type `?quakes`. Plot a histogram and the empirical cumulative distribution function for the *magnitudes*. Calculate the average  $\bar{X}_n$  and sample variance  $S_n^2$  for the magnitude.

*Hint:* The data set `quakes` is a data frame containing information on 5 observations (i.e. a table with 5 columns). To obtain a vector *only* containing the data in column 1, use `quakes[, 1]`.

- Suppose it is suggested that the data for the magnitudes  $X_1, \dots, X_n$  can be modelled by a  $\Gamma(\alpha, \beta)$  distribution. Find consistent estimators  $\hat{\alpha}_n$  and  $\hat{\beta}_n$  for  $\alpha$  and  $\beta$ , and calculate the estimates using the data from `quakes`. Plot the cumulative distribution function of  $\Gamma(\hat{\alpha}_n, \hat{\beta}_n)$  together with the empirical distribution function of the data. What do you observe?

*Hint:* Recall that for  $X \sim \Gamma(\alpha, \beta)$ , we have  $\mathbf{E}[X] = \frac{\alpha}{\beta}$  and  $\text{Var}[X] = \frac{\alpha}{\beta^2}$ .

<sup>1</sup>alternatively, you can find this data on <https://www.stat.cmu.edu/larry/all-of-statistics/data/fijiquakes.dat>