

# EE 5111: Estimation

Jan - May 2021

Mini Project 2

March 22, 2021

## 1 Performance of MLE

In the following exercises, we evaluate the performance of MLE for the estimation of certain parameters of three different distributions which are of practical interest.

1. Consider the following model:

$$x_n \sim \text{Weibull}(k, \theta) \quad n = 1, \dots, N, \quad (1)$$

where  $k$  and  $\theta$  are the shape and scale parameter of the Weibull RV with the probability density function (p.d.f.) given in (2).

$$f(x) = \begin{cases} \frac{k}{\theta} \left(\frac{x}{\theta}\right)^{k-1} e^{-(x/\theta)^k} & x \geq 0 \\ 0 & x < 0. \end{cases} \quad (2)$$

Compute the MLE of the parameter  $\theta$  for  $k = 1$  and  $N = 1, 10, 10^2, 10^3, 10^4$ . Generate data using  $\theta = 2$ .

2. The Gumbel distribution (Type I extreme value distribution) is widely used for characterising the distribution of the maximum of sequences of appropriately normalised independent and identically distributed random variables. Generate RVs  $Y_n, n = 1, \dots, N$  where  $Y_n = \max\{X_{n1}, \dots, X_{nK}\}$  and  $X_{nk} \sim \exp(\lambda)$ . Here,  $\exp(\lambda)$  represents the exponential distribution with mean  $\frac{1}{\lambda}$ . Then for  $K$  large, the RV  $Y_n$  follows the Gumbel distribution whose p.d.f. is given by

$$f_Y(y) = \frac{1}{\sigma} \exp\left(\frac{-(y - \mu)}{\sigma}\right) \exp\left(-\exp\left(\frac{-(y - \mu)}{\sigma}\right)\right), \quad (3)$$

where  $\mu = 14.9787$  and  $\sigma = 5$  are the location parameter and the scale parameter respectively. For  $\lambda = 0.2$ ,  $K = 20$  generate samples of the form  $Y_n$ . Assuming that you know the values of  $\mu$ , estimate the MLE of  $\sigma$  ( $\hat{\sigma}_{gev}$ ) using only the samples  $\{Y_n, n = 1, \dots, N\}$ .

3. The famous Pickands-Balkema-de Haan theorem proves that for any RV  $X_{nk}$  in the previous question, and a high threshold  $d$  ( $d \rightarrow F_{X_{nk}}^{-1}(1)$ ), we can approximately characterize the excess value  $\{X_{nk} - d | X_{nk} - d > 0\}$  by a generalized Pareto distribution (GPD) with

parameter  $\sigma$  which is same as the scale parameter in Question 2<sup>1</sup>. Let  $\sum_{n=1}^N \sum_{k=1}^K \mathbb{I}_{X_{nk} > d} = L$  be the number of RVs  $X_{nk}$  satisfying the required condition<sup>2</sup>. The p.d.f. of the GPD is given by

$$f_S(s) = \sigma^{-1} \exp\left(\frac{-s}{\sigma}\right). \quad (4)$$

Using the samples  $\{X_{nk}; 1 \leq n \leq N, 1 \leq k \leq K\}$  where  $X_{nk} \sim \exp(\lambda = 0.2)$ , generate RVs<sup>3</sup>  $\{S_\ell; 1 \leq \ell \leq L\}$  for  $d = 23$  and compute the MLE of the GPD parameter  $\sigma$  ( $\hat{\sigma}_{gpd}$ ). Compare this with the estimate of  $\sigma$  obtained from the previous question.

Present the following results for each of the above experiments<sup>4</sup>.

- Tabulate the values of  $\mathbb{E}[\hat{\theta}]$ ,  $\mathbb{E}[\hat{\sigma}_{gev}]$  and  $\mathbb{E}[\hat{\sigma}_{gpd}]$  against the number of samples  $N$ . What do you infer?
- Tabulate the values of  $\text{Var}(\hat{\theta})$ ,  $\text{Var}(\hat{\sigma}_{gev})$  and  $\text{Var}(\hat{\sigma}_{gpd})$  against the number of samples  $N$ . What do you infer?
- Plot the CDF of the estimates for  $N = 1, 10, 100, 1000, 10000$  samples. Ensure that you take enough realizations to get a smooth CDF. What can you say about the CDF? Justify. Do you observe the following relations?

$$\sqrt{N}(\hat{\theta} - \theta) \sim \mathcal{N}(0, I(\theta)^{-1})$$

$$\sqrt{N}(\hat{\sigma}_{gev} - \sigma) \sim \mathcal{N}(0, I(\sigma)^{-1})$$

$$\sqrt{L}(\hat{\sigma}_{gpd} - \sigma) \sim \mathcal{N}(0, I(\sigma)^{-1})$$

Here,  $I$  denotes the Fisher information<sup>5</sup>.

- Plot the PDF of the estimates for  $N = 1, 10, 100, 1000, 10000$  samples. Ensure that you take enough realizations to get a smooth PDF. What can you say about the PDF convergence?

## Submission

You are required to submit the solutions for this problem no later than March 23, 2021 (11:59 pm). Upload a compressed file containing the program/programs your four member team has written for the mini project along with a 1-2 page report in Moodle. The report should include the final expressions for the MLE and Fisher information you have used in all the three cases, the final plots for each of the questions (1)-(5) and a one-two line inference on the results observed. Note that you need not include any derivations in this report. The viva for each team will be conducted jointly on a date and time convenient for all the members.

<sup>1</sup>Here,  $A_k|b_k > 0$  represents the RV  $A_k$  given the condition  $b_k > 0$  is satisfied.

<sup>2</sup>Here  $\mathbb{I}_E$  is the indicator function for the event  $E$ .

<sup>3</sup>Note that here you are asked to use all the samples satisfying  $X_{nk} > d$  from the  $K \times N$  samples you had used for generating maxima samples in the previous question.

<sup>4</sup>Whenever the closed form solution for MLE is not available, use Newton Raphson or any other appropriate numerical method and mention the same in your report.

<sup>5</sup>Feel free to use any numerical integration package to evaluate the Fisher information if you could not derive the closed form expression for the same.