

# Problem Sheet 04

Alexander Pieper (815402)  
alexander.pieper@uni-potsdam.de

21. November 2021

## Exercise 1

Let  $(x_1, \dots, x_n) \in (0, \infty)^n$  be a realisation of independent on  $[0, \theta]$  uniformly distributed random variables  $X_1, \dots, X_n$ . What is Maximum Spacing Estimator in this case? Using the data set provide on Moodle computer the unknown parameter  $\theta$  via the Maximum spacing estimator for the three different sets of samples (note that they are of different sizes).

**Answer:**

$$X_1, \dots, X_n \sim U(0, \theta) \Rightarrow \Theta = \mathbb{R}^+ \quad (1)$$

The Maximum Spacing Estimator is given by (Lecture 8, p.24):

$$\widehat{\theta}_{MS} = \operatorname{argmax}_{\theta \in \Theta} \prod_{i=1}^{n+1} D_i(\Theta), \quad (2)$$

where  $D_i(\Theta) = F_\theta(X_{(i)}) - F_\theta(X_{(i-1)})$ . Also we define  $X_{(0)} = -\infty \rightarrow F_\theta(X_{(0)}) = 0$  and also  $X_{(n+1)} = -\infty \rightarrow F_\theta(X_{(n+1)}) = 1$

The Distribution function of our uniform distribution is given as

$$F(t) = \begin{cases} 0 & , \text{ for } t < 0 \\ \frac{t}{\theta} & , \text{ for } t \in [0, \theta] \\ 1 & , \text{ for } t > \theta. \end{cases} \quad (3)$$

Since all  $X_{(i)} \in [0, \theta]$  we have  $F_\theta(X_{(i)}) = \frac{X_{(i)}}{\theta}$  and thus

$$\widehat{\theta}_{MS} = \operatorname{argmax}_{\theta \in \Theta} \prod_{i=1}^{n+1} F_\theta(X_{(i)}) - F_\theta(X_{(i-1)}) \quad (4)$$

$$= \operatorname{argmax}_{\theta \in \Theta} \log \left( \prod_{i=1}^{n+1} F_\theta(X_{(i)}) - F_\theta(X_{(i-1)}) \right) \quad \text{same argmax} \quad (5)$$

$$= \operatorname{argmax}_{\theta \in \Theta} \sum_{i=1}^{n+1} \log (F_\theta(X_{(i)}) - F_\theta(X_{(i-1)})) \quad \text{property of log()} \quad (6)$$

$$= \operatorname{argmax}_{\theta \in \Theta} \sum_{i=1}^n \log (F_\theta(X_{(i)}) - F_\theta(X_{(i-1)})) + \log (F_\theta(X_{(n+1)}) - F_\theta(X_{(n)})) \quad \text{drag out last element} \quad (7)$$

$$= \operatorname{argmax}_{\theta \in \Theta} \sum_{i=1}^n \log \left( \frac{X_{(i)}}{\theta} - \frac{X_{(i-1)}}{\theta} \right) + \log \left( 1 - \frac{X_{(n)}}{\theta} \right) \quad F_\theta(X_{(n+1)}) = 1 \quad (8)$$

$$= \operatorname{argmax}_{\theta \in \Theta} \sum_{i=1}^n \log \left( \frac{X_{(i)} - X_{(i-1)}}{\theta} \right) + \log \left( 1 - \frac{X_{(n)}}{\theta} \right) \quad (9)$$

To get the argmax, we calculate the first derivative w.r.t  $\theta$ .

$$\frac{\partial}{\partial \theta} \sum_{i=1}^n \log \left( \frac{X_{(i)} - X_{(i-1)}}{\theta} \right) + \log \left( 1 - \frac{X_{(n)}}{\theta} \right) = \sum_{i=1}^n \left[ -\frac{1}{\theta} \right] + \frac{X_{(n)}}{(\theta - X_{(n)})\theta} \quad (10)$$

$$= -\frac{n}{\theta} + \frac{X_{(n)}}{(\theta - X_{(n)})\theta} \quad (11)$$

Setting it equal to 0 and solving it for  $\theta$ .

$$-\frac{n}{\theta} + \frac{X_{(n)}}{(\theta - X_{(n)})\theta} \stackrel{!}{=} 0 \quad (12)$$

$$\Leftrightarrow \frac{n}{\theta} = \frac{X_{(n)}}{(\theta - X_{(n)})\theta} \quad (13)$$

$$\Leftrightarrow n = \frac{X_{(n)}}{(\theta - X_{(n)})} \quad \theta > 0 \quad (14)$$

$$\Leftrightarrow n(\theta - X_{(n)}) = X_{(n)} \quad (15)$$

$$\Leftrightarrow \theta = \frac{X_{(n)}}{n} + X_{(n)} \quad (16)$$

$$\Rightarrow \widehat{\theta}_{MS} = \frac{n+1}{n} X_{(n)} \quad (17)$$

Now we have to verify that this is indeed a maximum. To do so, we compute the 2nd derivative and show that it is  $< 0$ , at least at  $\widehat{\theta}_{MS}$ .

$$\frac{\partial^2}{\partial \theta^2} \sum_{i=1}^n \log \left( \frac{X_{(i)} - X_{(i-1)}}{\theta} \right) + \log \left( 1 - \frac{X_{(n)}}{\theta} \right) = -\frac{n}{\theta^2} + \frac{X_{(n)}(X_{(n)} - 2\theta)}{\theta^2 (\theta - X_{(n)})^2} \quad (18)$$

$$= -\frac{n}{\frac{(n+1)^2}{n^2} X_{(n)}^2} + \frac{X_{(n)}(X_{(n)} - 2\frac{n+1}{n} X_{(n)})}{\frac{(n+1)^2}{n^2} X_{(n)}^2 \left( \frac{n+1}{n} X_{(n)} - X_{(n)} \right)^2} \quad \text{using } \widehat{\theta}_{MS} \text{ as } \theta \quad (19)$$

$$= -\frac{n}{\frac{(n+1)^2}{n^2} X_{(n)}^2} + \frac{-\frac{n+1}{n} X_{(n)}^2}{\frac{(n+1)^2}{n^2} X_{(n)}^2 \left( \frac{n+1}{n} X_{(n)} - X_{(n)} \right)^2} \quad (20)$$

$$= -\frac{n^3}{(n+1)^2 X_{(n)}^2} + \frac{-X_{(n)}^2}{\frac{n+1}{n} X_{(n)}^2 \left( \frac{1}{n} X_{(n)} \right)^2} \quad (21)$$

$$= -\frac{n^3}{(n+1)^2 X_{(n)}^2} + \frac{-X_{(n)}^2}{\frac{n+1}{n^3} X_{(n)}^4} \quad (22)$$

$$= -\frac{n^3}{(n+1)^2 X_{(n)}^2} - \frac{1}{\frac{n+1}{n^3} X_{(n)}^2} \quad (23)$$

$$= -\frac{n^3}{(n+1) X_{(n)}^2} - \frac{n^3}{(n+1) X_{(n)}^2} \quad (24)$$

$$= -\frac{2n^3}{(n+1) X_{(n)}^2} < 0, \text{ since } n > 0 \& X_{(n)} > 0 \quad (25)$$

Hence,  $\widehat{\theta}_{MS} = \frac{n+1}{n} X_{(n)}$  is a global maximum

**Interpretation:** The maximum spacing estimator for  $\theta$  is a bit larger than the largest value in the sample set. The Maximum Likelihood estimator of the same Problem is  $\widehat{\theta}_{ML} = X_{(n)}$  [Ghosha & Jammalamadaka, 2000, p.74]. Therefore the difference between these two estimators is only the correction term  $\frac{n+1}{n}$  in the MS Estimator.

Personally, I think in this case, this correction factor makes a lot of sense, because when you only have very few samples (say 3), it is very likely that the true parameter  $\theta$  is not the largest value, but a bit larger than that. Please see the Jupyter notebook in the next page for a visual interpretation on this estimator.

```
In [ ]: import matplotlib.pyplot as plt
```

Firstly, we read the sampleset as lists of observations

```
In [ ]: with open('sampleset_1_problemsheet4_ex1.txt', 'r') as handle:
    sampleset_1 = handle.readlines()
    sampleset_1 = [float(i.split('\n')[0]) for i in sampleset_1]
with open('sampleset_2_problemsheet4_ex1.txt', 'r') as handle:
    sampleset_2 = handle.readlines()
    sampleset_2 = [float(i.split('\n')[0]) for i in sampleset_2]
with open('sampleset_3_problemsheet4_ex1.txt', 'r') as handle:
    sampleset_3 = handle.readlines()
    sampleset_3 = [float(i.split('\n')[0]) for i in sampleset_3]
```

Now we need to calculate the  $\widehat{\theta}_{MS} = \frac{n+1}{n}X_{(n)}$  for each of the datasets.

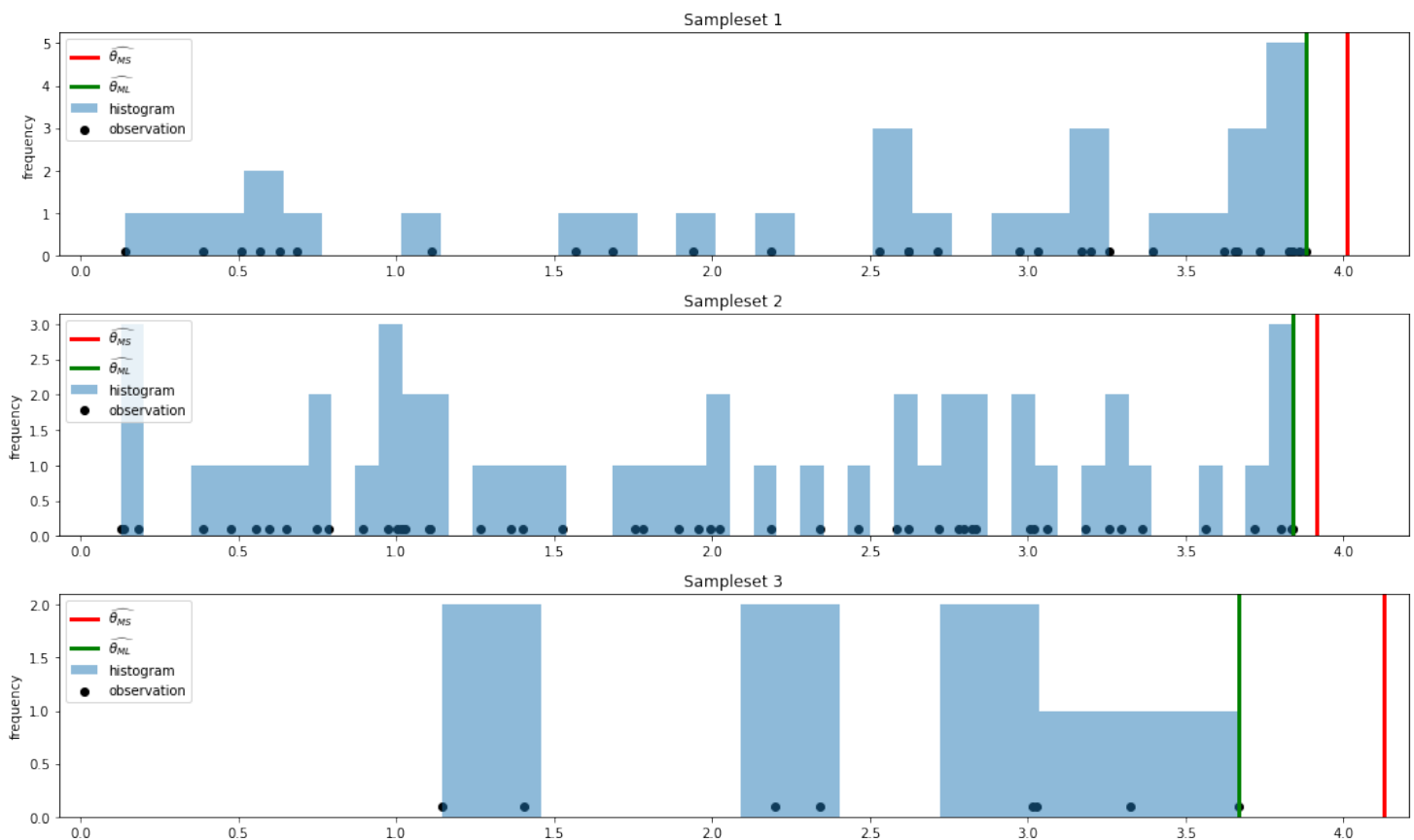
```
In [ ]: theta_hats_ms = {}
theta_hats_ml = {}
i = 1
for set in [sampleset_1, sampleset_2, sampleset_3]:
    theta_hats_ms['sampleset_' + str(i)] = (len(set) + 1)/(len(set)) * sorted(set)[-1]
    theta_hats_ml['sampleset_' + str(i)] = sorted(set)[-1]
    i += 1
```

```
In [ ]: print(f'The estimated Theta via Maximum Spacing method of the first sampleset is: {theta_hats_ms["sampleset_1"]}')
print(f'The estimated Theta via Maximum Spacing method of the second sampleset is: {theta_hats_ms["sampleset_2"]}')
print(f'The estimated Theta via Maximum Spacing method of the third sampleset is: {theta_hats_ms["sampleset_3"]}')
```

The estimated Theta via Maximum Spacing method of the first sampleset is: 4.0118133333333335  
The estimated Theta via Maximum Spacing method of the second sampleset is: 3.91578  
The estimated Theta via Maximum Spacing method of the third sampleset is: 4.1274

## Visualisation

```
In [ ]: fig, ax = plt.subplots(3,1,figsize = (15,9), sharex=True)
i = 0
for set in [sampleset_1, sampleset_2, sampleset_3]:
    ax[i].hist(set, bins = len(set), label = 'histogram', alpha = 0.5)
    ax[i].scatter(set, [0.1 for i in set], color = 'black', label = 'observation')
    ax[i].axvline(theta_hats_ms['sampleset_' + str(i+1)], color = 'red', label = r'$\widehat{\theta}_{MS}$', linewidth = 3)
    ax[i].axvline(theta_hats_ml['sampleset_' + str(i+1)], color = 'green', label = r'$\widehat{\theta}_{ML}$', linewidth = 3)
    ax[i].set(title = 'Sampleset ' + str(i+1), ylabel = 'frequency')
    ax[i].tick_params(labelbottom=True)
    ax[i].legend(loc = 'upper left')
    i += 1
fig.tight_layout()
plt.show()
```



## Analysis

The estimated parameters  $\theta$  are relatively close to 4, for all three sample sets. Looking at the Histograms and the observations of each Sampleset, we can also see this property of the Maximum spacing Estimator, that the smaller the size  $n$ , the further away is  $\theta$  from the largest value of the set  $X_{(n)}$ .

The MLE, which is  $X_{(n)}$  for this parameter, would not take the sample size into consideration at all. (Ghosha & Jammalamadaka, p.74)

## References

- [Ghosha & Jammalamadakab, 2000] Kaushik Ghosha, S. Rao Jammalamadakab, A general estimation method using spacings, (2000), Elsevier - Journal of statistical planning and inference 93 (2001) 71–82