

CS-6313 – Statistical Methods for Data Science

Mini Project #3

Group No - 5

Manan Dalal (MUD200000)

Lipi Patel (LDP210000)

Contribution of Team Members:

Worked together with each other to solve both the problems. Manan worked on the scripts and made them efficient and readable. Lipi worked on analytic solutions. Both worked together on documenting and giving the final touch to the project.

Question 1:

Suppose we would like to estimate the parameter $\theta (> 0)$ of a Uniform $(0, \theta)$ population based on a random sample X_1, \dots, X_n from the population. In the class, we have discussed two estimators for θ — the maximum likelihood estimator, $\hat{\theta}_1 = X(n)$, where $X(n)$ is the maximum of the sample, and the method of moments estimator, $\hat{\theta}_2 = 2\bar{X}$, where \bar{X} is the sample mean. The goal of this exercise is to compare the mean squared errors of the two estimators to determine which estimator is better. Recall that the mean squared error of an estimator $\hat{\theta}$ of a parameter θ is defined as $E\{(\hat{\theta} - \theta)^2\}$. For the comparison, we will focus on $n = 1, 2, 3, 5, 10, 30$ and $\theta = 1, 5, 50, 100$.

- a) Explain how you will compute the mean squared error of an estimator using Monte Carlo simulation.

Solution:

- ⇒ For computing the Mean squared error, we must first set the population parameter.
- ⇒ This will enable us to simulate sample values which would allow us to calculate the estimated value.
- ⇒ The estimated value of the difference between estimator and the parameter squared will give us the Mean Squared Error.

- b) For a given combination of (n, θ) , compute the mean squared errors of both $\hat{\theta}_1$ and $\hat{\theta}_2$ using Monte Carlo simulation with $N = 1000$ replications. Be sure to compute both estimates from the same data.

Solution

- ⇒ To compute MSE's for both estimators for 1000 replications, we will first create a function `getMLEandMOM` which will give the the values of MLE and MOM for a single sample.
- ⇒ Then, we would create a function called `getMSE` which will compute the MSE for both the estimators for a given value of n and θ .
- ⇒ For understanding purposes, we then get the MSE for $n = 1$ and $\theta = 1$.

R Code

```
1 # Returns MLE and MOM for a sample
2 getMLEandMOM <- function(n, theta) {
3   sample = runif(n, min = 0, max = theta)
4   mom = 2 * mean(sample)
5   mle = max(sample)
6   return(c(mle, mom))
7 }
8
9 # Returns Mean squared Error of MLE and MOM for 1000 replications
10 getMSE <- function(n, theta) {
11   values = replicate(1000, getMLEandMOM(n, theta))
12   values = (values - theta)^2
13   values.mom = values[c(TRUE, FALSE)]
14   values.mle = values[c(FALSE, TRUE)]
15
16   return(c(mean(values.mle), mean(values.mom)))
17 }
18
19 mse_1_1 = getMSE(1,1)
20
21 > mse_1_1
[1] 0.3266865 0.3393726
```

- c) Repeat (b) for the remaining combinations of (n, θ) . Summarize your results graphically.

Solution

- ⇒ In order to fit all plots in one graph, we use the function `par()`.
- ⇒ Generating graphs for MSE for both MLE and MOM, theta with fixed n.
- ⇒ Generating graphs for MSE for both MLE and MOM, n with fixed theta.

R Code

```
22 # Question 1 - c
23
24 n = c(1,2,3,5,10,30)
25 theta = c(1,5,50,100)
26
27 mse <- list()
28 for(x in n) {
29   for(y in theta){
30     mse = append(mse, list(getMSE(x,y)))
31   }
32 }

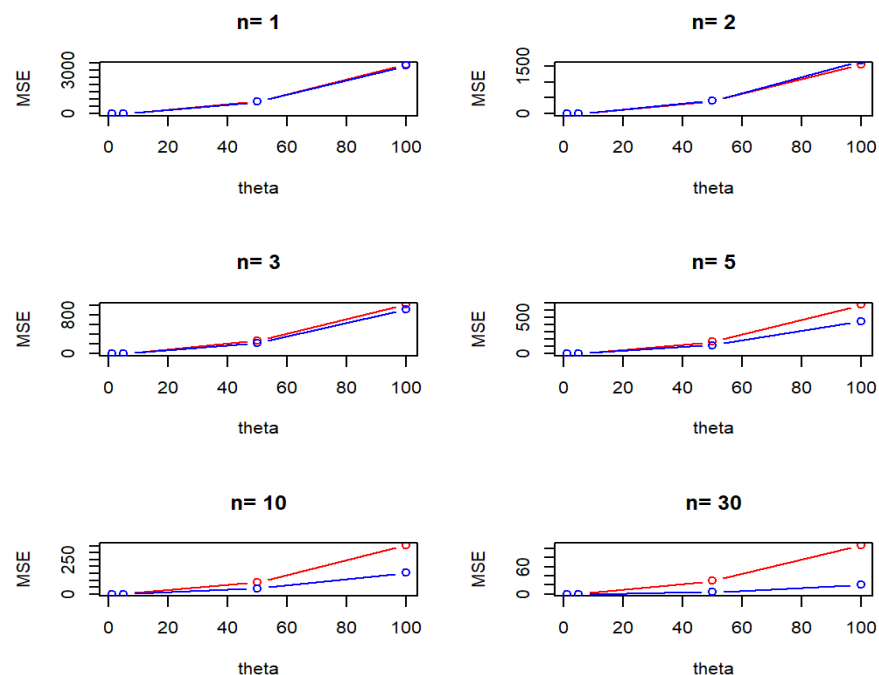
34 # Drawing graphs for fixed 'n' and varying 'theta'
35 par(mfrow = c(3,2))
36 t = 1
37
38 for(x in n) {
39   msel1 = list()
40   msel2 = list()
41   for(y in 0:3){
42     msel1 = append(msel1, unlist(mse[t + y])[1])
43     msel2 = append(msel2, unlist(mse[t + y])[2])
44   }
45
46   plot(theta, c(unlist(msel1)), type="b", xlab="theta", ylab="MSE",
47         col="red", main= paste("n=", x))
48   lines(theta, c(unlist(msel2)), type="b", col="blue")
49   legend("topleft", legend = c("MLE", "MOM"), col = c("red", "blue"),
50         text.col = c("black", "black"), lty = 1, pch = 1,
51         inset = 0.01, ncol = 1, cex = 0.6, bty = 'n')
52   t = t + 4
53 }
```

```

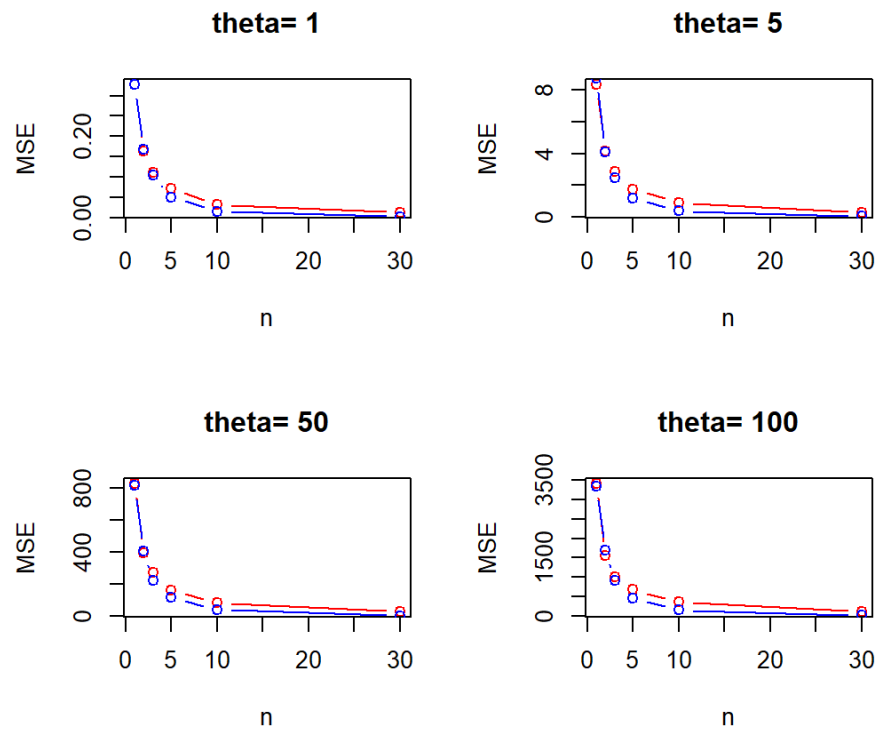
55 # Drawing graphs for fixed 'theta' and varying 'n'
56 par(mfrow = c(2,2))
57 t = 1
58
59 for(x in theta) {
60   msel1 = list()
61   msel2 = list()
62   for(y in 0:5){
63     msel1 = append(msel1, unlist(mse[t+4*y])[1])
64     msel2 = append(msel2, unlist(mse[t+4*y])[2])
65   }
66
67   plot(n, c(unlist(msel1)), type="b", xlab="n", ylab="MSE",
68        col="red", main= paste("theta=", x))
69   lines(n, c(unlist(msel2)), type="b", col="blue")
70   legend("topleft", legend = c("MLE", "MOM"), col = c("red", "blue"),
71         text.col = c("black", "black"), lty = 1, pch = 1,
72         inset = 0.01, ncol = 1, cex = 0.6, bty = 'n')
73   t = t + 1
74 }

```

Outputs:



Graph 1 : Graph with fixed value of n



Graph 2 : Graph with fixed value of theta

d) Based on (c), which estimator is better? Does the answer depend on n or θ ? Explain. Provide justification for all your conclusions.

Solution

- ⇒ The second graph shows us that the value of theta does not affect the graph as they are quite similar for all values of theta.
- ⇒ Thus, we can infer that the estimator would not depend on the value of theta.
- ⇒ By looking at Graph-1, we can evidently say that the Method of Moments(MOM) gives a very good estimate for smaller values of n (i.e. $n = 1, 2, 3$).
- ⇒ But, as the value of n increases (i.e. $n = 5, 10, 30$), the Maximum Likelihood Estimator performs better.
- ⇒ Thus, for larger values of n , MLE is better and would be the preferred choice of estimation when choosing among the two estimators.

Question-2:

Suppose the lifetime, in years, of an electronic component can be modeled by a continuous random variable with probability density function

$$f(x) = (\theta / x^{\theta+1}), x \geq 1, \\ 0, x < 1)$$

where $\theta > 0$ is an unknown parameter. Let X_1, \dots, X_n be a random sample of size n from this population.

a) Derive an expression for maximum likelihood estimator of θ .

Solution:

Let the likelihood function be

$$L(\theta) = \prod_{i=1}^n (\theta / x_i^{\theta+1})$$

Taking log on both sides:

$$\log(L(\theta)) = \log\left(\prod_{i=1}^n (\theta / x_i^{\theta+1})\right)$$

$$= \log\left(\theta^n \times \prod_{i=1}^n \frac{1}{x_i^{\theta+1}}\right)$$

$$= n \log \theta + \sum_{i=1}^n \log(x_i^{-\theta-1})$$

$$= n \log \theta - (\theta+1) \sum_{i=1}^n \log x_i$$

$$= n \log \theta - \theta \sum_{i=1}^n \log x_i - \sum_{i=1}^n \log x_i$$

$$\frac{d(L(\theta))}{d\theta} = \frac{n}{\theta} - \sum_{i=1}^n \log x_i$$

$$\therefore \frac{n}{\theta} - \sum_{i=1}^n \log x_i = 0$$

$$\therefore \frac{n}{\theta} = \sum_{i=1}^n \log x_i$$

$$\therefore \hat{\theta}_{MLE} = \frac{n}{\sum_{i=1}^n \log x_i}$$

- b) Suppose $n = 5$ and the sample values are $x_1 = 21.72$, $x_2 = 14.65$, $x_3 = 50.42$, $x_4 = 28.78$, $x_5 = 11.23$. Use the expression in (a) to provide the maximum likelihood estimate for θ based on these data.

Solution

Inputting the values given in the expression derived in (a), we get.

$$\begin{aligned}\hat{\theta}_{MLE} &= \frac{5}{\log(21.72) + \log(14.65) + \log(50.42) + \log(28.78) + \log(11.23)} \\ &= \frac{5}{\log(21.72 \times 14.65 \times 50.42 \times 28.78 \times 11.23)} \\ &= \frac{5}{\log(5137517.08)} \\ &= \frac{5}{15.45} \\ \hat{\theta}_{MLE} &= 0.3236\end{aligned}$$

- c) Even though we know the maximum likelihood estimate from (b), use the data in (b) to obtain the estimate by numerically maximizing the log-likelihood function using `optim` function in R. Do your answers match?

Solution

- ⇒ In R, in order to minimize a function, we can maximize the negative of that function.
- ⇒ We use `optim()` to perform this task.

- ⇒ Optim function takes various arguments such as:
- par – initial value of the parameters that need to be minimized.
 - fn – minimization function
 - method – optimization method to use
 - dat – data available.
- ⇒ When, we run the code below, we get the value of theta as .3064 which is close to the analytically derived value of 0.3236.

R Code

```
76 # Question-2 c
77
78 # Returns neg value of the derived function
79 negativeMLE <- function(par, dat) {
80   print(par,dat)
81   result = length(dat) * log(par) - (par+1) * sum(log(dat))
82   return(-result)
83 }
84
85 x <- c(21.42, 14,65, 50.42, 28.78, 11.23)
86
87 mle <- optim(par = 2, fn = negativeMLE, method = "L-BFGS-B",
88             hessian = TRUE, lower = 0.01, dat = x)
```

Output

```
> mle$par
[1] 0.3064923
```


- d) Use the output of numerical maximization in (c) to provide an approximate standard error of the maximum likelihood estimate and an approximate 95% confidence interval for θ . Are these approximations going to be good? Justify your answer.

Solution

- ⇒ The formula for Standard Error is given by:
 - $SE(\theta)^2 = 1/h$ $\Rightarrow h$ is the hessian function
- ⇒ From the following code we get the value for SE as 0.1251236.
- ⇒ A 95% confidence interval, means
 - $(1 - \alpha) = 0.95$
 - $\alpha = 0.05$
 - $\alpha/2 = 0.025$
 - $1 - \alpha/2 = 0.975$
- ⇒ The confidence interval formula is given by:
 - $\theta \pm Z_{\alpha/2} * SE(\theta)$
- ⇒ We can get $Z_{\alpha/2}$ by using the qnorm function.
- ⇒ The confidence interval we get from our calculations is:
 - (0.06125447, 0.55172985)
- ⇒ This lets us infer that out of 100 trials for population estimation, the true estimate value would lie within the interval 95 times.

R Code

```
90 # d
91
92 # Finding Standard Error
93 se <- (1/mle$hessian)^0.5
94
95 # Finding Confidence Interval
96 conf = mle$par + c(-1,1)*se*qnorm(0.975)
```

Output

```
> se
      [,1]
[1,] 0.1251236
> conf
[1] 0.06125447 0.55172985
```
