

# Homework 3

Chapter 6.6.

p. 95-96: ex. 1.

1) Let  $X_1, \dots, X_n \sim \text{Poisson}(\lambda)$  and let  $\hat{\lambda} = \frac{1}{n} \cdot \sum_{i=1}^n X_i$ . Find the bias, SE, MSE.

$\text{bias}(\hat{\lambda}) = E_n(\hat{\lambda}) - \lambda$ , if bias = 0, then we say unbiased, otherwise biased.

$$\begin{aligned}\text{bias}(\hat{\lambda}) &= E_n(\hat{\lambda}) - \lambda = E_n\left(\frac{1}{n} \cdot \sum_{i=1}^n X_i\right) - \lambda = \\ &= E_n\left(\frac{1}{n}\right) \cdot E\left(\sum_{i=1}^n X_i\right) - \lambda = \frac{1}{n} \cdot \sum_{i=1}^n E(X_i) - \lambda = \\ &= \frac{1}{n} \cdot \sum_{i=1}^n \lambda - \lambda = \frac{n \cdot \lambda}{n} - \lambda = \lambda - \lambda = 0 - \\ &\text{unbiased estimator}\end{aligned}$$

$$\begin{aligned}\widehat{SE}(\hat{\lambda}) &= SE\left(\frac{1}{n} \cdot \sum_{i=1}^n X_i\right) = \frac{1}{n} SE\left(\sum_{i=1}^n X_i\right) = \\ &= \frac{1}{n} \cdot \sqrt{\sum_{i=1}^n SE(X_i)^2} = \frac{1}{n} \cdot \sqrt{\sum_{i=1}^n \text{Var}(X_i)} = \\ &= \frac{1}{n} \cdot \sqrt{\sum_{i=1}^n \lambda} = \frac{1}{n} \cdot \sqrt{\lambda \cdot n} = \frac{\sqrt{n} \cdot \sqrt{\lambda}}{n} = \\ &= \sqrt{\frac{\lambda}{n}}\end{aligned}$$



MSE - Mean Square Error

$$\begin{aligned} \text{MSE} &= \text{bias}^2(\hat{\tau}) + \text{Var}(\hat{\tau}) = \\ &= \text{bias}^2(\hat{\tau}) + [\text{SE}(\hat{\tau})]^2 = \\ &= 0^2 + \left[\sqrt{\frac{\tau}{n}}\right]^2 = 0 + \frac{\tau}{n} = \frac{\tau}{n} \end{aligned}$$

Chapter 7.4.

p. 104-105: ex. 2, 3, 9.

2. Let  $X_1, \dots, X_n \sim \text{Bernoulli}(p)$   
Let  $Y_1, \dots, Y_m \sim \text{Bernoulli}(q)$

1) Find the plug-in estimator and  $\text{SE}_p$   
In Bernoulli distribution:

mean,  $\mu = p$

variance,  $\text{Var} = p \cdot (1-p)$

$$\begin{aligned} \text{Let } \hat{p} &= T(F) = \mu = \frac{1}{n} \sum_{i=1}^n X_i = \bar{X}_n \\ \text{The standard error is } \text{SE} &= \sqrt{\text{Var}(\bar{X}_n)} = \\ &= \sqrt{\text{Var}\left(\frac{1}{n} \sum_{i=1}^n X_i\right)} = \sqrt{\frac{1}{n^2} \cdot (\text{Var}(X_1) + \dots + \text{Var}(X_n))} = \\ &= \sqrt{\frac{n \cdot \text{Var}(X_i)}{n^2}} = \sqrt{\frac{\hat{p} \cdot (1-\hat{p})}{n}} \quad \text{or} \quad \sqrt{\frac{\bar{X}_n \cdot (1-\bar{X}_n)}{n}} \end{aligned}$$



2) Find an approximate 90% conf. interval for  $p$

Normal-based interval. For 90% conf. interval,  $Z_{\frac{1-0.9}{2}} = Z_{\frac{0.1}{2}} = Z_{0.05} = 1.64$

A normal-based 90% conf. interval for  $p$  is:

$$\hat{p} \pm Z_{0.05} \cdot \hat{SE}(\hat{p}) = \hat{p} \pm 1.64 \cdot \sqrt{\frac{\hat{p} \cdot (1-\hat{p})}{n}}$$

$$\text{or} \quad \bar{X}_n \pm 1.64 \cdot \sqrt{\frac{\bar{X}_n \cdot (1-\bar{X}_n)}{n}}$$

3) Find the plug-in estimator and  $\hat{SE}$  for  $(p-q)$

The plug-in estimator  $(p-q) = \int x dF_1(x) - \int x dF_2(x) = \theta$

$$\begin{aligned} \hat{\theta} &= \int x d\hat{F}_1(x) - \int x d\hat{F}_2(x) = \frac{1}{n} \sum_{i=1}^n X_i - \frac{1}{m} \sum_{i=1}^m Y_i = \\ &= \hat{p} - \hat{q} = \bar{X}_n - \bar{Y}_m \end{aligned}$$

$$\begin{aligned} \text{Var}(\hat{\theta}) &= \text{Var}(\hat{p} - \hat{q}) = \text{Var}(\hat{p}) + \text{Var}(\hat{q}) = \\ &= \frac{\hat{p}(1-\hat{p})}{n} + \frac{\hat{q}(1-\hat{q})}{m} \end{aligned}$$



$$SE(\hat{\theta}) = \sqrt{\frac{\hat{p}(1-\hat{p})}{n} + \frac{\hat{q}(1-\hat{q})}{m}}$$

4) Find an approximate 90% conf. interval for  $(p-q)$ .

$Z_{\frac{1-\alpha}{2}} = Z_{\frac{1-0.9}{2}} = Z_{\frac{0.1}{2}} = Z_{0.05} = 1.64$ , an approximate interval is:

$$(\hat{p}-\hat{q}) \pm 1.64 \cdot \sqrt{\frac{\hat{p}(1-\hat{p})}{n} + \frac{\hat{q}(1-\hat{q})}{m}}$$

or

$$(\bar{X}_n - \bar{Y}_m) \pm 1.64 \cdot \sqrt{\frac{\bar{X}_n \cdot (1-\bar{X}_n)}{n} + \frac{\bar{Y}_m \cdot (1-\bar{Y}_m)}{m}}$$

### 3. Simulation.

First, I defined the variable grid, which represents all the  $x$ -values from  $-3$  to  $+3$ . More length - more precise and less step will be between values. Next I define  $F_n$  -



cumulative distribution function, which will fit for all values between  $-3$  to  $+3$ . The graph will be smooth, it is true CDF.

Then I take 100 observation of standard Normal distribution. For plotting - I sort them. Then I create an estimate of CDF  $\hat{F}_n(x)$ :

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

Then I construct confidence bound

$$\alpha = 1 - 0.95 = 0.05$$

$$\epsilon_n = \sqrt{\frac{1}{2n} \log\left(\frac{2}{\alpha}\right)}$$

$$L(X) = \max\{\hat{F}_n(x) - \epsilon_n, 0\}$$

$$U(X) = \min\{\hat{F}_n(x) + \epsilon_n, 1\}$$

No matter what the true  $F$ , then

$$P(L(x) \leq F(x) \leq U(x) \text{ for all } x) \geq 1 - \alpha$$



```
1 #-----
2 #Generate 100 observations from a N(0,1) distribution.
3 #-----
4 par(mfrow=c(2,2))
5 grid <- seq(-3,3, length=1000)
6 Fn <- pnorm(grid)
7 plot(grid, Fn, type="l", xlab="x value",ylab="CDF", sub="Theoretical CDF")
8
9
10 n <- 100
11 x <- rnorm(n, mean = 0, sd = 1)
12 x <- sort(x)
13
14 Fn_hat <- (1:n)/n #ecdf(x)
15 plot(x, Fn_hat, type="s", xlab="x value",ylab="CDF", sub="Empirical CDF", xlim=c(-3,3))
16
17 plot(grid, Fn, type="l",xlab="x value",ylab="CDF")
18 lines(x,Fn_hat,type="s", lty=3,col=4,lwd=3)
19
```

```

20 ▾ #-----
21 #Compute a 95 percent confidence band for the CDF F
22 ▾ #-----
23 alph <- 1 - 0.95
24 epsilon <- sqrt(1 / (2 * n) * log(2 / alph)) # I did by the formulas on the page 99 in wasserman's Book
25
26 L <- pmax(Fn_hat - epsilon, 0)
27 U <- pmin(Fn_hat + epsilon, 1)
28
29 plot(grid, Fn, type="l", xlab="x", ylab="cdf", xlim=c(-3,3))
30 lines(x,L,lty=3,col=2,type="s")
31 lines(x,U,lty=3,col=2,type="s")
32

```

```

33 #-----
34 #Repeat this 1000 times and see how often
35 #the confidence band contains the true distribution function.
36 #-----
37 ans <- c()
38
39 for(i in 1:1000){
40   grid <- seq(-3,3, length=1000)
41   Fn <- pnorm(grid)
42
43   n <- 100
44   x <- rnorm(n, mean = 0, sd = 1)
45   x <- sort(x)
46
47   Fn_hat <- (1:n)/n #ecdf(x)
48
49
50   alph <- 1 - 0.95
51   epsilon <- sqrt(1 / (2 * n) * log(2 / alph))
52
53   L <- pmax(Fn_hat - epsilon, 0)
54   U <- pmin(Fn_hat + epsilon, 1)
55
56   fraction = c()
57   for (i in 1:100)
58   {
59     logic <- U[i]>=pnorm(x[i]) && L[i]<=pnorm(x[i])
60     fraction <- append(fraction, logic)
61   }
62
63   ans <- append(ans,all(fraction))
64 }
65
66 mean(ans) # About 0.955 which is >=0.95 (95% confidence interval)
67

```

33:90 [Untitled] ↕

Console Terminal × Jobs ×

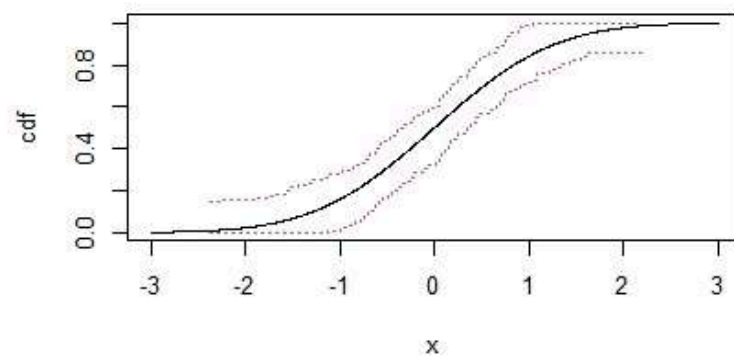
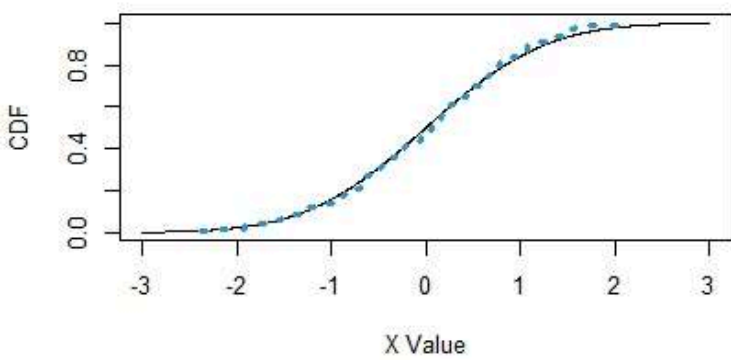
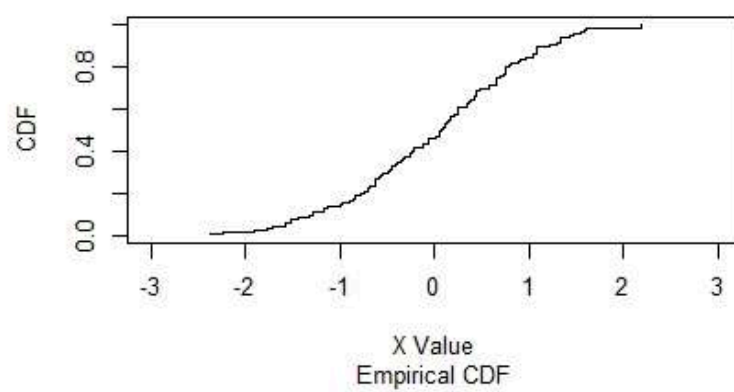
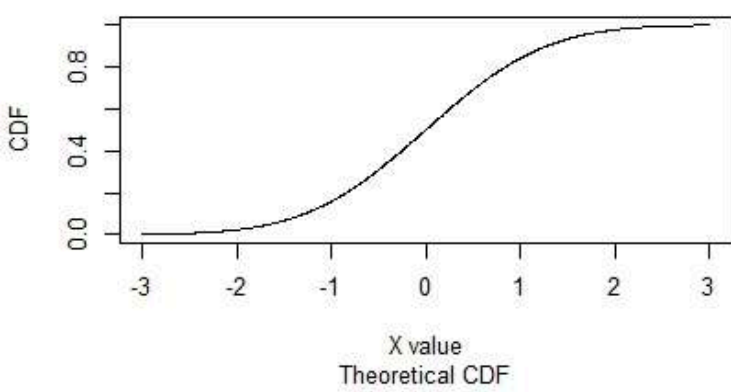
R 4.1.2 · ~/

```

+ ans <- append(ans,all(fraction))
+ }
>
> mean(ans)
[1] 0.955
> |

```







```

69 ▾ #-----
70 #Repeat using data from a Cauchy distribution.
71 ▾ #-----
72 ans <- c()
73
74 ▾ for(i in 1:1000){
75   grid <- seq(-3,3, length=1000)
76   Fn <- pnorm(grid)
77
78   n <- 100
79
80   x <- rcauchy(n) #I only changed the function for random sample
81   x <- sort(x)
82
83   Fn_hat <- (1:n)/n #ecdf(x)
84
85
86   alph <- 1 - 0.95
87   epsilon <- sqrt(1 / (2 * n) * log(2 / alph))
88
89   L <- pmax(Fn_hat - epsilon, 0)
90   U <- pmin(Fn_hat + epsilon, 1)
91
92   fraction = c()
93   for (i in 1:100)
94 ▾   {
95     logic <- U[i]>=pnorm(x[i]) && L[i]<=pnorm(x[i])
96     fraction <- append(fraction, logic)
97 ▾   }
98
99   ans <- append(ans,all(fraction))
100 ▾ }
101
102 mean(ans) # About 0.19 - 0.20 which is (< 0.95) lower than 95% confidence interval
103 |
104 #Zhetessov Nur M.
105

```



In Cauchy distribution, calculating the mean will provide no useful information, because mean is undefined, as the  $\text{Var}(X)$ . Therefore, I can't calculate confidence interval properly for Cauchy distribution. In experiment it shows only 20% match.

9.

$n = 100$  people

$m = 100$  people

$k_1 = 90$  recovered people

$k_2 = 85$  recovered people

$p_1$  - standard treatment

$p_2$  - new treatment

---

Estimate  $\theta = p_1 - p_2$ ,  $SE(\theta)$

80% conf interval.

95% conf interval.



$$p_1 = \frac{\# \text{ recovered by standard treat.}}{\# \text{ total}}$$

$$p_2 = \frac{\# \text{ recovered by new treat.}}{\# \text{ total}}$$

$$p_1 = \frac{90}{100} = 0.9 \text{ - probability of recovery by standard treat.}$$

$$p_2 = \frac{85}{100} = 0.85 \text{ - prob. recover. by new treat.}$$

An estimate is

$$\theta = p_1 - p_2 = 0.9 - 0.85 = 0.05$$

Standard error is

I showed it in exercise 2.

$$[SE(\theta)]^2 = \text{Var}(\theta) = \text{Var}(p_1 - p_2) =$$

$$= \text{Var}(p_1) + \text{Var}(p_2) = \frac{p_1(1-p_1)}{n} + \frac{p_2(1-p_2)}{m}$$

$$SE(\theta) = \sqrt{\frac{p_1(1-p_1)}{n} + \frac{p_2(1-p_2)}{m}}$$

$$SE(\theta) = \sqrt{\frac{0.9 \cdot 0.1}{100} + \frac{0.85 \cdot 0.15}{100}} = 0.0466$$



Z-score for 80% conf. interval is

$$Z_{\frac{1-0.8}{2}} = Z_{\frac{0.2}{2}} = Z_{0.1} = 1.282$$

Upper bound:

$$(p_1 - p_2) + Z_{0.1} * SE(\theta)$$

$$0.05 + 1.282 * 0.0466$$

$$0.110$$

Lower bound:

$$(p_1 - p_2) - Z_{0.1} * SE(\theta)$$

$$0.05 - 1.282 * 0.0466$$

$$-0.010$$

Conf. interval 80% is  $(-0.010; 0.110)$

Z-score for 95% is  $Z_{0.025} = 1.96$

Lower bound:  $(p_1 - p_2) - 1.96 * SE(\theta)$

$$0.05 - 1.96 * 0.0466$$

$$-0.041$$

Upper bound:  $0.05 + 1.96 * 0.0466$

$$0.1413 \quad \text{Interval is } (-0.041; 0.1413)$$



## Chapter 8.6

p. 116-118: ex. 6, 7d

6. In this our statistic is a function

$$g(\mu) = e^{\mu}$$

I have create a random normal sample  $\sim N(5, 1)$  of 100 observations.

Then I found  $\hat{\theta}$  by  $g(\bar{X}_{\hat{\theta}})$   
Using bootstrap method, I sample randomly our original sample and found its statistic by  $g(x)$ . All  $g(\theta^*)$  were contained in a vector.

To find conf. interval, I need SE. We find SE from vector with all  $\theta^*$  taking their mean and  $\sqrt{\text{Var}}$ . With it I can construct conf. interval:

$$\hat{\theta} \pm z * \sqrt{\frac{1}{n} \sum_{i=1}^n (\theta_i^* - \frac{1}{n} \sum_{j=1}^n \theta_j^*)^2}$$

Then I plot all  $\theta^*$  versus its frequency and



```
6Ex_Zhetessov.R x
Source on Save
Run
Source

1 #-----
2 #Create a data set mu=5 consisting of n=100 observations.
3 #-----
4
5 n <- 100
6 mu <- 5
7 x <- rnorm(n, mean=mu, sd=1)
8
9 #-----
10 # Use the bootstrap method
11 #-----
12 est_hat_theta <- exp(mean(x))
13
14 bar_theta <- c()
15 for(i in 1:100000){
16   rand_sampling <- sample(x, size = n, replace = TRUE)
17   star_theta = exp(mean(rand_sampling))
18
19   bar_theta <- append(bar_theta, star_theta)
20 }
21
22
23 #-----
24 #(a) Get the SE and 95 percent confidence interval for Theta
25 #-----
26 SE = sqrt(var(bar_theta))
27
28 alpha = 1 - (95/100)
29 z_score = abs(qnorm(alpha/2))
30
31 upp_bound = est_hat_theta + z_score * SE
32 Low_bound = est_hat_theta - z_score * SE
33
34 c(Low_bound,upp_bound)
35
36 #-----
37 #(b) Plot a histogram of the bootstrap replications. Compare this to the true sampling distribution
38 #-----
39 hist(bar_theta, breaks = 1000, col = "darkmagenta",freq = FALSE)
40 lines(x = density(x = bar_theta), col = "red", lwd = 3)
41
42 true_theta = exp(mu)
43 abline(v=true_theta,col="blue",lwd=2)
44
45 #Zhetessov Nur M.
```

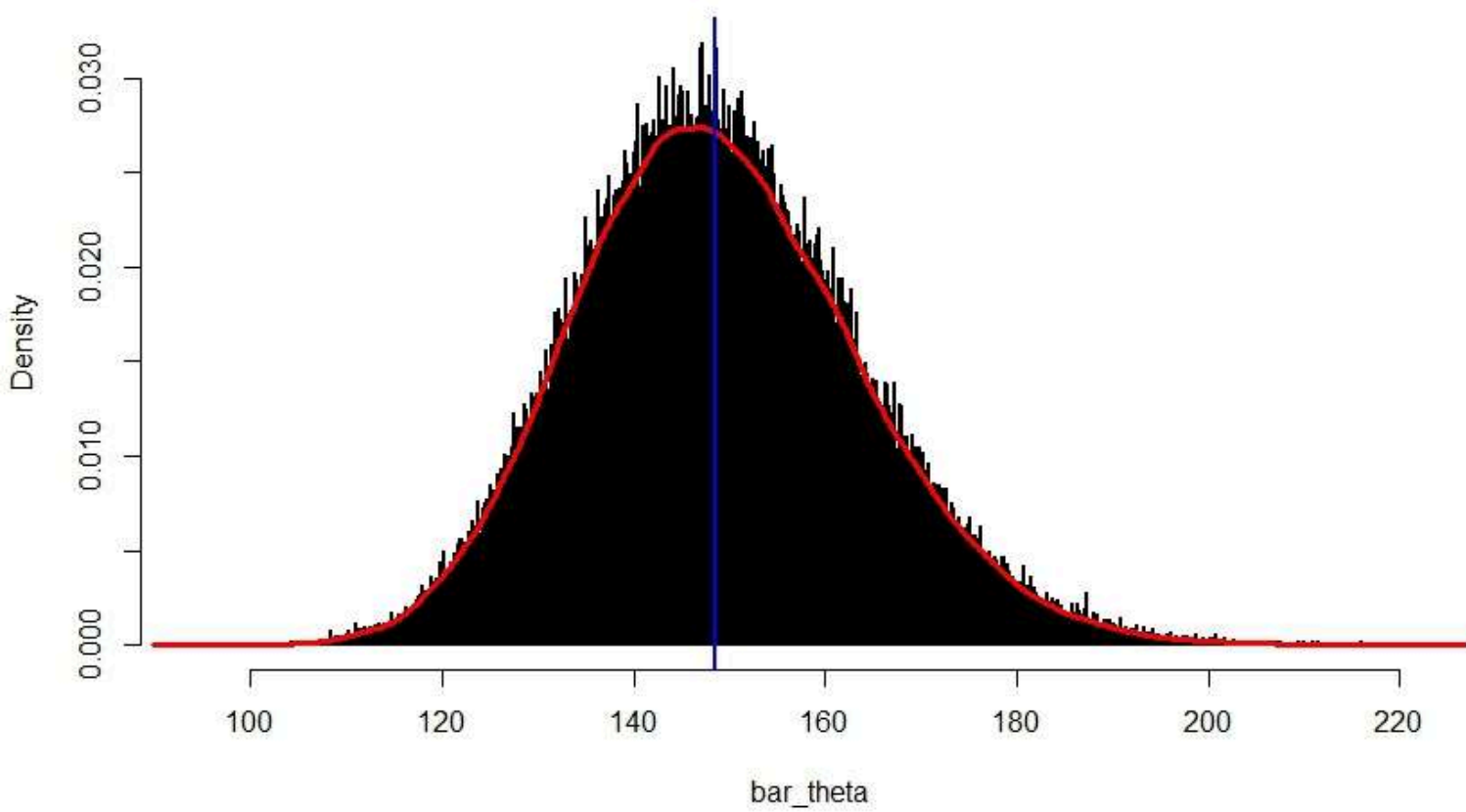
45:18 (Untitled) R Script

Console Terminal Jobs

```
R 4.1.2 ~/>
> Low_bound = est_hat_theta - z_score * SE
>
> c(Low_bound,upp_bound)
[1] 119.2606 176.9549
>
> #-----
> #(b) Plot a histogram of the bootstrap replications
```



Histogram of bar\_theta





draw a line - blue line - the true  $\theta$  which by statistic is  $e^{\mu} = e^5$ . I see that  $\geq 0.95$  of times our confidence interval contain the true parameter.

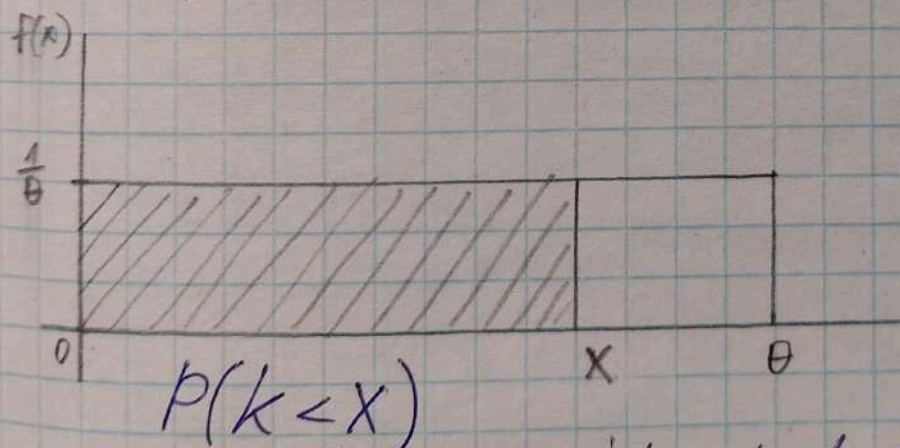
7. a)  $X_1, \dots, X_n \sim \text{Uniform}(0, \theta)$ .

$$\hat{\theta} = X_{\max} = \max\{X_1, \dots, X_n\}.$$

Generate data set of 50 with  $\theta=1$

In Uniform distribution the probability density function (PDF) is  $f(x) = \frac{1}{b-a}$  ( $a \leq x \leq b$ )

$$f(x) = \frac{1}{\theta - 0} = \frac{1}{\theta}$$



Assume all  $X_i$  are independent. The maximum

$X_i < \theta$ , if and only if all  $X_i < \theta$

$$P(\hat{\theta} < x) = P(\max\{X_1, \dots, X_n\} < x) = P(X_1 < x, X_2 < x, \dots, X_n < x) =$$



$$\begin{aligned}
 &= P(X_1 < x) \cdot P(X_2 < x) \cdots P(X_n < x) = \\
 &= (x-0) \cdot \left(\frac{1}{\theta-0}\right) \cdot (x-0) \cdot \left(\frac{1}{\theta-0}\right) \cdots (x-0) \cdot \left(\frac{1}{\theta-0}\right) = \\
 &= \left(x \cdot \frac{1}{\theta}\right)^n = \left(\frac{x}{\theta}\right)^n \text{ since } \theta=1, \text{ then } (x)^n \\
 &\text{where the distribution is:}
 \end{aligned}$$

$$F(x) = \begin{cases} 0, & x < 0 \\ \left(\frac{x}{\theta}\right)^n, & 0 < x < \theta \\ 1, & x > \theta \end{cases} \text{ for } 0 \leq \theta \leq 1$$

The confidence interval in computer simulation showed good coverage 0.95-lower bound to 1, 0.2 upper bound. which confirms the true  $\theta = 1$



```

7a_Ex.R x 6Ex_Zhetessov.R x
Source on Save
1 n<-50
2 x <- runif(n, 0, 1)
3
4 hat_theta = max(x)
5
6 bar_theta <- c()
7 for(i in 1:100000){
8   rand_sampling <- sample(x, size = n, replace = TRUE)
9   star_theta = max(rand_sampling)
10
11   bar_theta <- append(bar_theta, star_theta)
12 }
13
14 SE = sqrt(var(bar_theta))
15
16 alpha = 1 - 0.95
17 z = abs(qnorm(alpha/2))
18
19 Lower_bound = hat_theta - z * SE
20 Upper_bound = hat_theta + z * SE
21
22 Lower_bound
23 Upper_bound
24
25 hist(bar_theta, breaks = 100, col = "darkmagenta",freq = FALSE, xlim=c(0.8, 1))
26 #lines(x = density(x = bar_theta), col = "red", lwd = 3)
27
28 true_theta = max(1)
29 abline(v=true_theta,col="blue",lwd=2)
30
31 #Zhetessov Nur M.
32

```

31:18 (Top Level) ↕

```

Console Terminal x Jobs x
R 4.1.2 · C:/Users/Nur/Desktop/дз/Стат/HW3/
> z = abs(qnorm(alpha/2))
>
> Lower_bound = hat_theta - z * SE
> Upper_bound = hat_theta + z * SE
>
> Lower_bound
[1] 0.958237
> Upper_bound
[1] 1.021634
>
> hist(bar_theta, breaks = 100, col = "darkmagenta",freq = FALSE, xlim=c(0.8, 1))
> #lines(x = density(x = bar_theta), col = "red", lwd = 3)
>
> true_theta = max(1)
> abline(v=true_theta,col="blue",lwd=2)
>
> #Zhetessov Nur M.
>

```



Histogram of bar\_theta

