

Exam of Statistics for Data Science

It is forbidden to consult any material during the test, with the exception of a scientific calculator. Duration of written exam is 1.5h.

Exercise 1 (6 points). Let $X, Y \sim Ber(0.5)$ be i.i.d. random variables. Define the random variables:

$$U = X + Y \quad V = |X - Y|$$

- a Determine the joint and marginal distributions of U and V .
- b Find out whether U and V are dependent or independent.
- c Determine the covariance $Cov(U, V)$ and the correlation coefficient $Cor(U, V)$.

Solution. (a) see solution of Ex. 9.6 (a) at page 437 of [1].

(b) Dependent, because e.g.,

$$P(U = 0, V = 0) = \frac{1}{4} \neq P(U = 0)P(V = 0) = \frac{1}{4} \cdot \frac{1}{2} = \frac{1}{8}$$

(c) We have

$$E[U] = 0 \cdot \frac{1}{4} + 1 \cdot \frac{1}{2} + 2 \cdot \frac{1}{4} = 1 \quad E[V] = 0 \cdot \frac{1}{2} + 1 \cdot \frac{1}{2} = \frac{1}{2}$$

$$E[UV] = \sum_a a \cdot P(UV = a) = 0P(UV = 0) + 1P(U = 1, V = 1) + 2P(U = 1, V = 2) = \frac{1}{2}.$$

Therefore

$$Cov(U, V) = E[UV] - E[U]E[V] = \frac{1}{2} - \frac{1}{2} = 0$$

and then $Cor(U, V) = 0$ as well.

Exercise 2 (6 points). Suppose that x_1, \dots, x_n is a dataset, which is a realization of a random sample from a Rayleigh distribution, which is a continuous distribution with probability density function:

$$f_\theta = \frac{x}{\theta^2} e^{-\frac{x^2}{\theta^2}} \quad \text{for } x \geq 0.$$

In this case what is the maximum likelihood estimate of θ ?

Solution. The likelihood is

$$L(\theta) = \prod_{i=1}^n \frac{x_i}{\theta^2} e^{-\frac{x_i^2}{\theta^2}}$$

and then the log-likelihood is

$$\ell(\theta) = \sum_{i=1}^n (\log x_i - 2\log \theta - \frac{1}{2\theta^2} x_i^2)$$

The (log-)likelihood has maximum when the derivative (w.r.t. θ) is zero, i.e., when

$$\frac{d\ell(\theta)}{d\theta} = -\frac{2n}{\theta} + \frac{1}{\theta^3} \sum_{i=1}^n x_i^2 = 0$$

which occurs for

$$\theta = \sqrt{\frac{1}{2n} \sum_{i=1}^n x_i^2}$$

Exercise 3 (6 points). Consider a linear regression model *without intercept*:

$$Y_i = \beta x_i + U_i \quad \text{for } i = 1, \dots, n$$

where U_1, \dots, U_n are independent random variables with $E[U_i] = 0$ and $\text{Var}(U_i) = 2$. Consider the following three estimators for the parameter β :

$$\hat{\beta}_1 = \frac{1}{n} \sum_{i=1}^n Y_i / x_i$$

$$\hat{\beta}_2 = (\sum_{i=1}^n Y_i) / (\sum_{i=1}^n x_i)$$

$$\hat{\beta}_3 = (\sum_{i=1}^n x_i Y_i) / (\sum_{i=1}^n x_i^2)$$

Show that all three estimators are unbiased for β . Compute their variance and discuss their efficiency.

Solution. We observe $E[Y_i] = \beta x_i + E[U_i] = \beta x_i$ and $\text{Var}(Y_i) = \text{Var}(U_i) = 2$. We calculate:

- $E[\hat{\beta}_1] = \frac{1}{n} \sum_{i=1}^n E[Y_i] / x_i = \frac{1}{n} \sum_{i=1}^n \beta x_i / x_i = \beta$
- $E[\hat{\beta}_2] = (\sum_{i=1}^n E[Y_i]) / (\sum_{i=1}^n x_i) = (\sum_{i=1}^n \beta x_i) / (\sum_{i=1}^n x_i) = \beta$
- $E[\hat{\beta}_3] = (\sum_{i=1}^n x_i E[Y_i]) / (\sum_{i=1}^n x_i^2) = (\sum_{i=1}^n \beta x_i^2) / (\sum_{i=1}^n x_i^2) = \beta$

and:

- $\text{Var}(\hat{\beta}_1) = \frac{1}{n^2} \sum_{i=1}^n \text{Var}(Y_i) / x_i^2 = \frac{2}{n} \frac{1}{n} \sum_{i=1}^n \frac{1}{x_i^2}$
- $\text{Var}(\hat{\beta}_2) = (\sum_{i=1}^n \text{Var}(Y_i)) / (\sum_{i=1}^n x_i)^2 = 2n / (\sum_{i=1}^n x_i)^2 = \frac{2}{n} \frac{1}{(\frac{1}{n} \sum_{i=1}^n x_i)^2}$
- $\text{Var}(\hat{\beta}_3) = (\sum_{i=1}^n x_i^2 \text{Var}(Y_i)) / (\sum_{i=1}^n x_i^2)^2 = \frac{2}{\sum_{i=1}^n x_i^2} = \frac{2}{n} \frac{1}{\frac{1}{n} \sum_{i=1}^n x_i^2}$

Since $1/x^2$ is a convex function, by the Jensen's inequality

$$\frac{1}{n} \sum_{i=1}^n \frac{1}{x_i^2} \geq \frac{1}{(\frac{1}{n} \sum_{i=1}^n x_i)^2}$$

and then $\text{Var}(\hat{\beta}_1) \geq \text{Var}(\hat{\beta}_2)$. Since x^2 is also convex, by the Jensen's inequality:

$$\frac{1}{n} \sum_{i=1}^n x_i^2 \geq \left(\frac{1}{n} \sum_{i=1}^n x_i \right)^2$$

and then $\text{Var}(\hat{\beta}_2) \geq \text{Var}(\hat{\beta}_3)$.

Exercise 4 (6 points). Environmentalists have taken 16 samples from the wastewater of a chemical plant and measured the concentration of a certain carcinogenic substance. They found $\bar{x}_{16} = 2.24$ (ppm) and $s_{16}^2 = 1.12$, and want to use these data in a lawsuit against the plant. It may be assumed that the data are a realization of a normal random sample.

- a Construct the 97.5% one-sided confidence interval that the environmentalists made to convince the judge that the concentration exceeds legal limits.
- b The plant management uses the same data to construct a 97.5% one-sided confidence interval to show that concentrations are not too high. Construct this interval as well.

Hint. $t_{15,0.025} = 2.131$

Solution. See full solution of Ex. 24.6 (a,b) at page 470 of [1].

Exercise 5 (6 points). Write an R function to compute p-value in 2-sided t-test without using the pre-defined built-in `t.test()`.

Solution.

```
pvalue = function(data, m, n)  # data = vector of n values, m = actual mean
{
  xbar <- mean(data) # sample mean
  sbar <- sd(data) # sample variance
  t0 <- sqrt(n)*(xbar-m)/sbar # studentized mean
  v = pt(t0, n-1) # P(t <= t0)
  p <- min(v, 1-v) # lower tail
  return (2*p) # 2-sided
}
```

References

- [1] F.M. Dekking, C. Kraaikamp, H.P. Lopuhaä, and L.E. Meester. *A Modern Introduction to Probability and Statistics*. Springer, 2005.

