

Quiz

Two measurements of the side of the square were produced. Suppose the two measurements X_1 and X_2 are two independent random variables with mean a and variance σ^2 . The true length of the side of the square is a . Find MSE for the following estimator of the area of the square: X_1X_2 .

Solution:

Since measurements are independent:

$$E(X_1X_2) = E(X_1)E(X_2) = a^2,$$

which means that X_1X_2 is an unbiased estimator of a^2 :

$$\text{Bias}(X_1X_2) = E(X_1X_2) - a^2 = a^2 - a^2 = 0.$$

Also, we will need an expectation of X_i^2 . From variance and mean identity:

$$E(X_i^2) = V(X_i) + E(X_i)^2 = \sigma^2 + a^2.$$

MSE is a sum of bias squared and variance. Let's calculate a variance of X_1X_2 (using independence of X_1 and X_2):

$$\begin{aligned} V(X_1X_2) &= E(X_1^2X_2^2) - E(X_1X_2)^2 = E(X_1^2)E(X_2^2) - E(X_1)^2E(X_2)^2 = \\ &= (\sigma^2 + a^2)^2 - (a^2)^2 = \sigma^4 + 2a^2\sigma^2 + a^4 - a^4 = \sigma^4 + 2a^2\sigma^2. \end{aligned}$$

MSE of X_1X_2 then:

$$\text{MSE}(X_1X_2) = \text{Bias}^2(X_1X_2) + V(X_1X_2) = 0^2 + \sigma^4 + 2a^2\sigma^2 = \boxed{\sigma^4 + 2a^2\sigma^2}.$$

Problem 1

Come up with any linear p.d.f. (non-uniform)

- Can you conclude which parameter is greater without calculations – mean or median?
- Calculate them explicitly.
- Discuss the result. Why p.d.f. is unbalanced in terms of “mass” around its median?

Solution:

Let's consider random variable X with p.d.f.:

$$f(x) = \left(1 - \frac{x}{2}\right) \cdot I_{\{0 \leq x \leq 2\}}.$$

- The distribution $f(x)$ is continuous unimodal, and its right tail is longer than the left one (skewed to the right), which means that

$$E(X) > \text{median}(X).$$

See the illustration in the fig. 1.

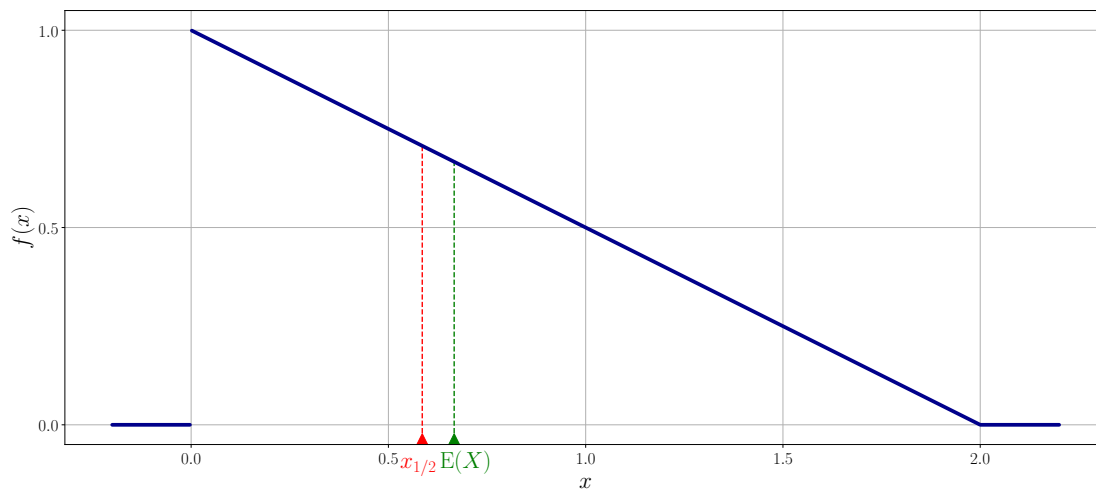


Figure 1: P.d.f. $f(x) = \left(1 - \frac{x}{2}\right) \cdot I_{\{0 \leq x \leq 2\}}$ with mean $E(X)$ and median $x_{1/2}$.

- Mean:

$$E(X) = \int_{-\infty}^{\infty} x f(x) dx = \int_0^2 x \left(1 - \frac{x}{2}\right) dx = \frac{x^2}{2} \Big|_0^2 - \frac{x^3}{6} \Big|_0^2 = \frac{2}{3}.$$

Let $x_{1/2} = \text{median}(X)$. Then:

$$\begin{aligned} \int_0^{x_{1/2}} \left(1 - \frac{x}{2}\right) dx &= \frac{1}{2}, \\ x \Big|_0^{x_{1/2}} - \frac{x^2}{4} \Big|_0^{x_{1/2}} &= \frac{1}{2}, \\ x_{1/2}^2 - 4x_{1/2} + 2 &= 0. \end{aligned}$$

$$x_{1/2} = \boxed{2 - \sqrt{2}}.$$

- (c) It's often unclear why mean and median of some random variable are different quantities. From a physical analogy we know that expected value is the center of mass of a distribution. Let's try to balance the distribution around its median.

We have a lever with an fulcrum at the point $x_{1/2} = 2 - \sqrt{2}$. Masses on both sides from fulcrum are identical and equal to $1/2$. Now we need to find their centers of mass.

The left one:

$$x_1 = \frac{\int_0^{x_{1/2}} x f(x) dx}{\int_0^{x_{1/2}} f(x) dx} = \frac{\int_0^{2-\sqrt{2}} \left(x - \frac{x^2}{2}\right) dx}{\frac{1}{2}} = \frac{2}{3} (\sqrt{2} - 1).$$

The right one:

$$x_2 = \frac{\int_{x_{1/2}}^2 x f(x) dx}{\int_{x_{1/2}}^2 f(x) dx} = \frac{\int_{2-\sqrt{2}}^2 \left(x - \frac{x^2}{2}\right) dx}{\frac{1}{2}} = \frac{2}{3} (3 - \sqrt{2}).$$

Distances from median to those points:

$$\begin{aligned} x_{1/2} - x_1 &= \frac{8 - 5\sqrt{2}}{3} \approx 0.31, \\ x_2 - x_{1/2} &= \frac{\sqrt{2}}{3} \approx 0.47, \end{aligned}$$

which means that the right part produces bigger moment and a lever will rotate clockwise. If the fulcrum had been in a point of expected value, lever would have been stable.

Mean and median are such values that minimize following expressions:

$$E(X) = \arg \min_{M_2} \int_X (x - M_2)^2 \cdot f(x) dx,$$
$$\text{median}(X) = \arg \min_{M_1} \int_X |x - M_1| \cdot f(x) dx,$$

and you can see that in contrast to median, mean considers distance to the point unevenly – with quadratic penalty. So the further point from expected value, the higher its “cost”. This is reflected in mechanical analogy – we consider not only a mass, but also a distance to that mass.

Problem 2

Let $\{X_1, \dots, X_n\}$ be a random sample from a $\text{Bin}(m, \pi)$ distribution, with both m and π unknown. Find the method of moments estimators for m – the number of trials, and π – the probability of success.

Solution:

Moment is a quantitative measure, related to the shape of p.d.f.'s graph.

Population k^{th} raw moments of a random variable X are calculated as follows:

$$\mu_k = \mathbf{E}(X^k).$$

Moments about mean $\mathbf{E}(X)$ are called central and calculated the following way:

$$\tilde{\mu}_k = \mathbf{E}\left((X - \mathbf{E}(X))^k\right).$$

They describe the shape of the function, independently of translation. The best known moments are: the 1st raw one – mean, and the 2nd central – variance σ^2 . If the function represents mass, then the 1st moment is the center of the mass, and the 2nd moment is the rotational inertia.

Let's discuss a few of higher-order moments. The 3rd standardized central moment is called skewness and shows the asymmetry of a function:

$$\text{Skew}(X) = \frac{\tilde{\mu}_3}{\sigma^3}.$$

Signs of skewness and their effect on p.d.f. are reflected in the fig. 2.

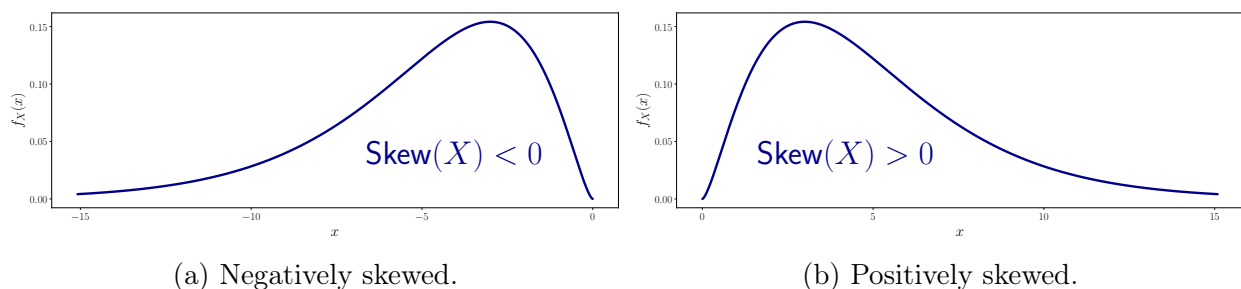


Figure 2: Probability density functions of random variables with opposite skewness.

The 4th standardized central moment is called excess kurtosis and shows the sharpness of distribution peak:

$$\text{Kurt}(X) = \frac{\tilde{\mu}_4}{\sigma^4} - 3,$$

where -3 shift is used to manipulate the excess kurtosis of standard normal distribution to be 0, since the sharpness of a peak is estimated with a reference to that of $\mathcal{N}(0, 1)$. In order to make correct comparisons of two distributions with excess kurtosis, their variance should be identical.

If $\text{Kurt}(X) > 0$, distribution peak is sharper than standard normal's one, if $\text{Kurt}(X) < 0$, distribution peak is smoother. The example is illustrated in the fig. 3.

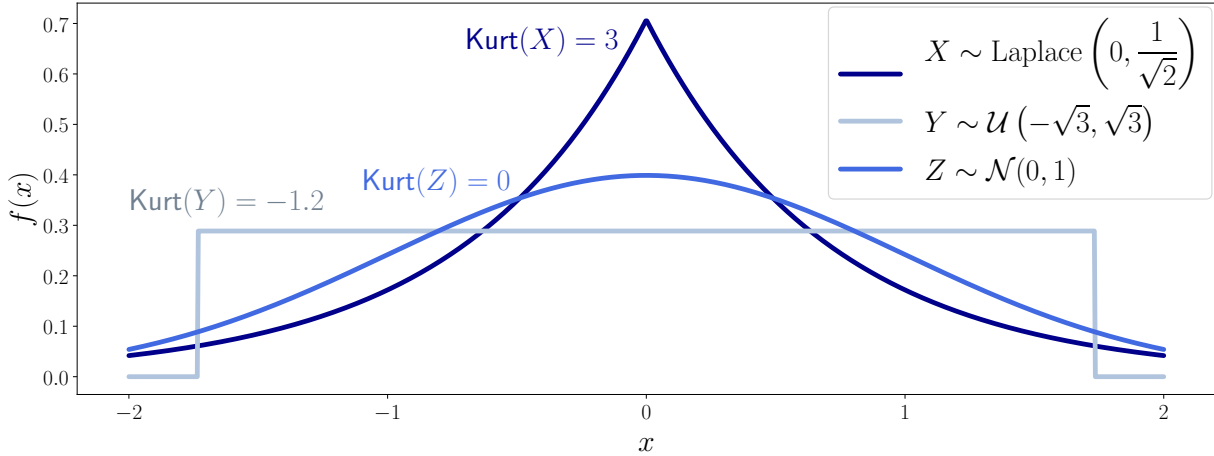


Figure 3: Comparison of excess kurtosis for Laplace X , uniform Y and normal Z distributions with zero mean and variance 1.

Population moments can be acquired from moment-generating function $M_X(t)$ of a random variable X , which is defined as

$$M_X(t) = E(e^{tX}),$$

which is basically two-sided Laplace transform with parameter $-t$ of a p.d.f. $f_X(x)$ in continuous case. We can get moments μ_k as follows:

$$E(X^k) = \left. \frac{d^k}{dt^k} M_X(t) \right|_{t=0}.$$

Let's note some important properties of moment-generation functions. If X_1, \dots, X_n are independent random variables, then:

$$M_{\sum_{i=1}^n X_i}(t) = \prod_{i=1}^n M_{X_i}(t).$$

For linear transformation $\alpha X + \beta$, where $\alpha, \beta \in \mathbb{R}$:

$$M_{\alpha X + \beta}(t) = e^{\beta t} M_X(\alpha t).$$

In addition to population moments μ_k , which we can get from population distribution of a random variable X , there are also sample moments M_k , which are calculated from a sample of n i.i.d. observations X_1, \dots, X_n :

$$M_k = \frac{1}{n} \sum_{i=1}^n X_i^k = \overline{X^k}.$$

The method of moments (MM) is to equate population and sample moments of the same degree k :

$$\mu_k = M_k,$$

and out of this system to get an estimation of a required parameter. Number of equations depends on the number of unknown parameters.

In our problem we have 2 unknown parameters – m and π . Thus, we need to make equations for the 1st and the 2nd moments:

$$\begin{cases} \mu_1 = M_1, \\ \mu_2 = M_2, \end{cases} \implies \begin{cases} E(X) \Big|_{m=\hat{m}, \pi=\hat{\pi}} = \overline{X}, \\ E(X^2) \Big|_{m=\hat{m}, \pi=\hat{\pi}} = \overline{X^2}, \end{cases} \implies \begin{cases} \hat{m}\hat{\pi} = \overline{X}, \\ \hat{m}\hat{\pi}(1 - \hat{\pi}) + (\hat{m}\hat{\pi})^2 = \overline{X^2}. \end{cases}$$

Here we used the fact that for $X \sim \text{Bin}(m, \pi)$: $E(X) = m\pi$ and $E(X^2) = V(X) + E(X)^2 = m\pi(1 - \pi) + (m\pi)^2$. Also, when we equated population and sample moments, we have made an estimation for required parameters m and π (true parameter shall not change, depending on the values from sample), thus putting hats over them.

Substituting \overline{X} into the second equation:

$$\hat{m}\hat{\pi}(1 - \hat{\pi}) + \overline{X}^2 = \overline{X^2},$$

$$\hat{m}\hat{\pi}(1 - \hat{\pi}) = \overline{X^2} - \overline{X}^2 = \frac{1}{n} \sum_{i=1}^n (X_i - \overline{X})^2 = \hat{\sigma}^2,$$

where $\hat{\sigma}^2$ is a biased sample variance. In the end we have two simple equations:

$$\begin{cases} \hat{m}\hat{\pi} = \overline{X}, \\ \hat{m}\hat{\pi}(1 - \hat{\pi}) = \hat{\sigma}^2, \end{cases} \implies \begin{cases} \hat{m} = \frac{\overline{X}}{\hat{\pi}}, \\ 1 - \hat{\pi} = \frac{\hat{\sigma}^2}{\overline{X}}, \end{cases} \implies \boxed{\begin{cases} \hat{m} = \frac{\overline{X}^2}{\overline{X} - \hat{\sigma}^2}, \\ \hat{\pi} = 1 - \frac{\hat{\sigma}^2}{\overline{X}}. \end{cases}}$$

Problem 3

Suppose that we have a random sample $\{X_1, \dots, X_n\}$ from a uniform distribution. Find the method of moments estimator of θ if

- (a) $X \sim \mathcal{U}(0, \theta)$,
- (b) $X \sim \mathcal{U}(-\theta, \theta)$.

Solution:

Let's remember that for a random variable $X \sim \mathcal{U}(a, b)$, population moments are:

$$\mathbb{E}(X) = \frac{a+b}{2} \quad \text{and} \quad \mathbb{V}(X) = \frac{(b-a)^2}{12}.$$

- (a) Equation on first moments gives:

$$\mathbb{E}(X) \Big|_{\theta=\hat{\theta}} = \bar{X}, \quad \implies \quad \frac{0+\hat{\theta}}{2} = \bar{X}, \quad \implies \quad \hat{\theta} = \boxed{2\bar{X}}.$$

- (b) Let's try to repeat previous calculations, but for symmetric distribution:

$$\mathbb{E}(X) \Big|_{\theta=\hat{\theta}} = \bar{X}, \quad \implies \quad \frac{-\hat{\theta}+\hat{\theta}}{2} = \bar{X}, \quad \implies \quad 0 = \bar{X}.$$

We end up with degenerate result, which can be “approximately correct” for large sample sizes.

Let's try to use an equation of another order, the 2nd for instance:

$$\begin{aligned} \mathbb{E}(X^2) \Big|_{\theta=\hat{\theta}} &= \mathbb{V}(X) \Big|_{\theta=\hat{\theta}} + \mathbb{E}(X)^2 \Big|_{\theta=\hat{\theta}} = \bar{X}^2, \quad \implies \quad \frac{(\hat{\theta} - (-\hat{\theta}))^2}{12} + 0^2 = \bar{X}^2, \\ &\implies \quad \hat{\theta}^2 = 3\bar{X}^2, \quad \implies \quad \hat{\theta} = \boxed{\sqrt{3\bar{X}^2}}. \end{aligned}$$

Let's find an estimator for θ^2 and calculate its bias:

$$\mathbb{E}(X^2) \Big|_{\theta^2=\hat{\theta}^2} = \bar{X}^2, \quad \implies \quad \frac{(\sqrt{\hat{\theta}^2} - (-\sqrt{\hat{\theta}^2}))^2}{12} = \bar{X}^2, \quad \implies \quad \hat{\theta}^2 = 3\bar{X}^2.$$

Estimator $\hat{\theta}^2$ is clearly unbiased:

$$\mathbb{E}(\hat{\theta}^2) = \mathbb{E}(3\bar{X}^2) = 3 \cdot \frac{1}{n} \cdot \mathbb{E}(n \cdot \bar{X}^2) = 3 \cdot \frac{\theta^2}{3} = \theta^2.$$

Does it mean that the estimator $\hat{\theta}$ is also unbiased? We can't calculate bias directly, but the answer to the question is NO, which can be seen from Jensen's inequality:

$$\mathbb{E}(\sqrt{3\bar{X}^2}) \leq \sqrt{\mathbb{E}(3\bar{X}^2)},$$

which holds since square root $f(x) = \sqrt{x}$ is a concave function. Moreover, equality is achieved when function f is not strictly convex or concave, in other words when it's linear. Since we have strictly concave function, in our case:

$$\mathbb{E}(\hat{\theta}) < \sqrt{\mathbb{E}(\hat{\theta}^2)} = \sqrt{\theta^2} = \theta, \quad \implies \quad \text{Bias}(\hat{\theta}) < 0.$$

Thus, the estimator $\hat{\theta}$ always underestimates true parameter θ .

The standard formulation of Jensen's inequality for 2 points x_1 and x_2 is following:

$$\forall t \in [0, 1] : \quad f(tx_1 + (1-t)x_2) \leq tf(x_1) + (1-t)f(x_2),$$

where f is a convex function. The inequality formalizes the statement that the secant line of a convex function lies above the graph of the function, illustrated in the fig. 4.

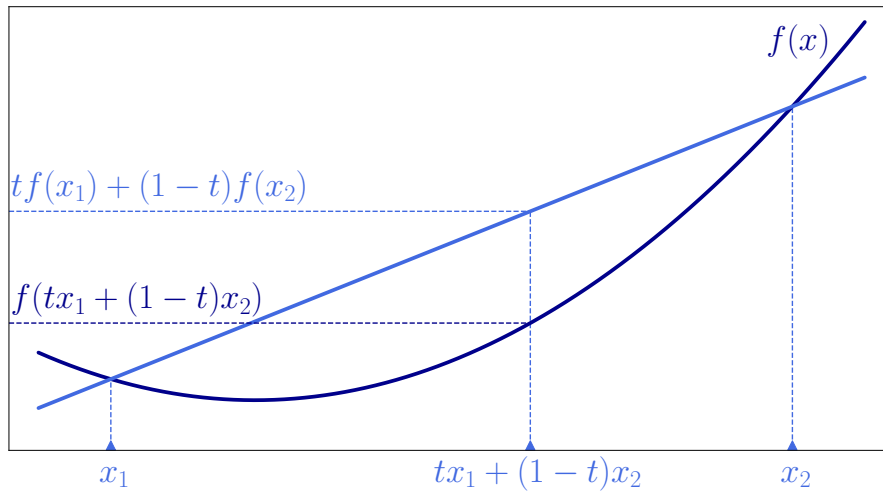


Figure 4: Jensen's inequality for 2 points.

If we would split the interval (x_1, x_2) with not one t , but with several t_i , such that $\sum_{i=1}^n t_i = 1$, number of which could tend to infinity, Jensen's inequality can be expressed in probabilistic statement, where t_i have physical meanings of probabilities:

$$f(\mathbb{E}(X)) \leq \mathbb{E}(f(X)),$$

where f is a convex function. Inequality inverts for a concave function.

Problem 4

Suppose that you are given observations y_1, y_2, y_3 and y_4 such that:

$$\begin{aligned}y_1 &= \alpha + \beta + \varepsilon_1, \\y_2 &= -\alpha + \beta + \varepsilon_2, \\y_3 &= \alpha - \beta + \varepsilon_3, \\y_4 &= -\alpha - \beta + \varepsilon_4.\end{aligned}$$

The variables $\varepsilon_i, i \in \{1, 2, 3, 4\}$, are independent and normally distributed with mean 0 and variance σ^2 .

- Find the least squares estimators of the parameters α and β .
- Verify that the least squares estimators in (a) are unbiased.
- Find the variance of the least squares estimator of the parameter α .

Solution:

- I. Matrices.** Let's use a known solution for OLS in matrix notation:

$$\hat{\boldsymbol{\beta}}_{\text{OLS}} = (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{y},$$

where \mathbf{X} is a matrix of regressors, \mathbf{y} is a vector of regressands, $\hat{\boldsymbol{\beta}}_{\text{OLS}}$ is vector of OLS parameters estimates:

$$\mathbf{X} = \begin{pmatrix} 1 & 1 \\ -1 & 1 \\ 1 & -1 \\ -1 & -1 \end{pmatrix}, \quad \mathbf{y} = \begin{pmatrix} y_1 \\ y_2 \\ y_3 \\ y_4 \end{pmatrix}, \quad \hat{\boldsymbol{\beta}}_{\text{OLS}} = \begin{pmatrix} \hat{\alpha} \\ \hat{\beta} \end{pmatrix}.$$

Pseudoinverse matrix:

$$\begin{aligned}(\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top &= \left(\begin{pmatrix} 1 & -1 & 1 & -1 \\ 1 & 1 & -1 & -1 \end{pmatrix} \begin{pmatrix} 1 & 1 \\ -1 & 1 \\ 1 & -1 \\ -1 & -1 \end{pmatrix} \right)^{-1} \begin{pmatrix} 1 & -1 & 1 & -1 \\ 1 & 1 & -1 & -1 \end{pmatrix} = \\ &= \begin{pmatrix} 4 & 0 \\ 0 & 4 \end{pmatrix}^{-1} \begin{pmatrix} 1 & -1 & 1 & -1 \\ 1 & 1 & -1 & -1 \end{pmatrix} = \begin{pmatrix} \frac{1}{4} & 0 \\ 0 & \frac{1}{4} \end{pmatrix} \begin{pmatrix} 1 & -1 & 1 & -1 \\ 1 & 1 & -1 & -1 \end{pmatrix} = \\ &= \begin{pmatrix} \frac{1}{4} & -\frac{1}{4} & \frac{1}{4} & -\frac{1}{4} \\ \frac{1}{4} & \frac{1}{4} & -\frac{1}{4} & -\frac{1}{4} \end{pmatrix}.\end{aligned}$$

Thus, OLS estimates:

$$\begin{cases} \hat{\alpha} = \frac{y_1 - y_2 + y_3 - y_4}{4}, \\ \hat{\beta} = \frac{y_1 + y_2 - y_3 - y_4}{4}. \end{cases}$$

II. Direct calculation. Let's minimize square of the second norm of error vector ϵ explicitly. Error vector is:

$$\epsilon = (\epsilon_1 \quad \epsilon_2 \quad \epsilon_3 \quad \epsilon_4)^\top,$$

and quantity to minimize:

$$\text{ESS} = \|\epsilon\|^2 = \epsilon_1^2 + \epsilon_2^2 + \epsilon_3^2 + \epsilon_4^2,$$

$$\hat{\alpha}, \hat{\beta} = \arg \min_{\alpha, \beta} \text{ESS}.$$

Expressing ESS via regressands and parameters:

$$\text{ESS} = (y_1 - \alpha - \beta)^2 + (y_2 + \alpha - \beta)^2 + (y_3 - \alpha + \beta)^2 + (y_4 + \alpha + \beta)^2.$$

By necessary condition of extremum:

$$\begin{cases} \left. \frac{\partial \text{ESS}}{\partial \alpha} \right|_{\alpha=\hat{\alpha}, \beta=\hat{\beta}} = -2(y_1 - \hat{\alpha} - \hat{\beta}) + 2(y_2 + \hat{\alpha} - \hat{\beta}) - 2(y_3 - \hat{\alpha} + \hat{\beta}) + 2(y_4 + \hat{\alpha} + \hat{\beta}) = 0, \\ \left. \frac{\partial \text{ESS}}{\partial \beta} \right|_{\alpha=\hat{\alpha}, \beta=\hat{\beta}} = -2(y_1 - \hat{\alpha} - \hat{\beta}) - 2(y_2 + \hat{\alpha} - \hat{\beta}) + 2(y_3 - \hat{\alpha} + \hat{\beta}) + 2(y_4 + \hat{\alpha} + \hat{\beta}) = 0, \end{cases}$$

$$\begin{cases} 4\hat{\alpha} = y_1 - y_2 + y_3 - y_4, \\ 4\hat{\beta} = y_1 + y_2 - y_3 - y_4, \end{cases} \implies \begin{cases} \hat{\alpha} = \frac{y_1 - y_2 + y_3 - y_4}{4}, \\ \hat{\beta} = \frac{y_1 + y_2 - y_3 - y_4}{4}. \end{cases}$$

Here we need to prove that values of $\hat{\alpha}$ and $\hat{\beta}$ are indeed arguments of minimum. Second derivatives:

$$\left. \frac{\partial^2 \text{ESS}}{\partial \alpha^2} \right|_{\alpha=\hat{\alpha}, \beta=\hat{\beta}} = 8, \quad \left. \frac{\partial^2 \text{ESS}}{\partial \alpha \partial \beta} \right|_{\alpha=\hat{\alpha}, \beta=\hat{\beta}} = \left. \frac{\partial^2 \text{ESS}}{\partial \beta \partial \alpha} \right|_{\alpha=\hat{\alpha}, \beta=\hat{\beta}} = 0, \quad \left. \frac{\partial^2 \text{ESS}}{\partial \beta^2} \right|_{\alpha=\hat{\alpha}, \beta=\hat{\beta}} = 8.$$

By sufficient condition of the minimum, let's see if the Hessian matrix \mathcal{H} is positive-definite (all minors are positive):

$$\mathcal{H} = \begin{pmatrix} \left. \frac{\partial^2 \text{ESS}}{\partial \alpha^2} \right|_{\alpha=\hat{\alpha}, \beta=\hat{\beta}} & \left. \frac{\partial^2 \text{ESS}}{\partial \alpha \partial \beta} \right|_{\alpha=\hat{\alpha}, \beta=\hat{\beta}} \\ \left. \frac{\partial^2 \text{ESS}}{\partial \beta \partial \alpha} \right|_{\alpha=\hat{\alpha}, \beta=\hat{\beta}} & \left. \frac{\partial^2 \text{ESS}}{\partial \beta^2} \right|_{\alpha=\hat{\alpha}, \beta=\hat{\beta}} \end{pmatrix} = \begin{pmatrix} 8 & 0 \\ 0 & 8 \end{pmatrix} \succ 0.$$

Q.E.D.

- (b) Estimator is unbiased, if its expected value is equal to estimated parameter.

$$\mathbb{E}(\hat{\alpha}) = \mathbb{E}\left(\frac{y_1 - y_2 + y_3 - y_4}{4}\right) = \frac{1}{4}(\mathbb{E}(y_1) - \mathbb{E}(y_2) + \mathbb{E}(y_3) - \mathbb{E}(y_4)).$$

Regressands y_i are random variables, since they are linear combination of parameters (constants) and errors ε_i , which are random variables themselves. Since ε_i are zero-mean, expected value of y_i is defined by terms of α and β :

$$\mathbb{E}(\hat{\alpha}) = \frac{1}{4}((\alpha + \beta) - (-\alpha + \beta) + (\alpha - \beta) - (-\alpha - \beta)) = \boxed{\alpha}.$$

$$\begin{aligned}\mathbb{E}(\hat{\beta}) &= \mathbb{E}\left(\frac{y_1 + y_2 - y_3 - y_4}{4}\right) = \frac{1}{4}(\mathbb{E}(y_1) + \mathbb{E}(y_2) - \mathbb{E}(y_3) - \mathbb{E}(y_4)) = \\ &= \frac{1}{4}((\alpha + \beta) + (-\alpha + \beta) - (\alpha - \beta) - (-\alpha - \beta)) = \boxed{\beta}.\end{aligned}$$

- (c) Variance of the regressand y_i is equal to the variance ε_i , since constant terms do not affect variance. Since ε_i are independent, y_i are also independent.

$$\begin{aligned}\mathbb{V}(\hat{\alpha}) &= \mathbb{V}\left(\frac{y_1 - y_2 + y_3 - y_4}{4}\right) = \frac{1}{4^2}(\mathbb{V}(y_1) + (-1)^2 \cdot \mathbb{V}(y_2) + \mathbb{V}(y_3) + (-1)^2 \cdot \mathbb{V}(y_4)) = \\ &= \frac{1}{16}(\sigma^2 + \sigma^2 + \sigma^2 + \sigma^2) = \boxed{\frac{\sigma^2}{4}}.\end{aligned}$$

Problem 5

A coin was tossed 10 times. Faces of a coin turned out as follows:

$$H, H, H, H, T, H, H, H, T, H.$$

What is the most likely probability of getting heads after one toss?

Solution:

Let the result of one throw of a coin be a Bernoulli random variable with probability of success (getting heads) p :

$$X \sim \text{Bernoulli}(p).$$

The result of an experiment with 10 throws has the following probability (considering outcomes in sample collectively independent):

$$\begin{aligned} & \mathbb{P}(X_1 = 1, X_2 = 1, \dots, X_9 = 0, X_{10} = 1) = \\ &= \mathbb{P}(X_1 = 1) \cdot \mathbb{P}(X_2 = 1) \cdot \dots \cdot \mathbb{P}(X_9 = 0) \cdot \mathbb{P}(X_{10} = 1) = p^8(1-p)^2. \end{aligned}$$

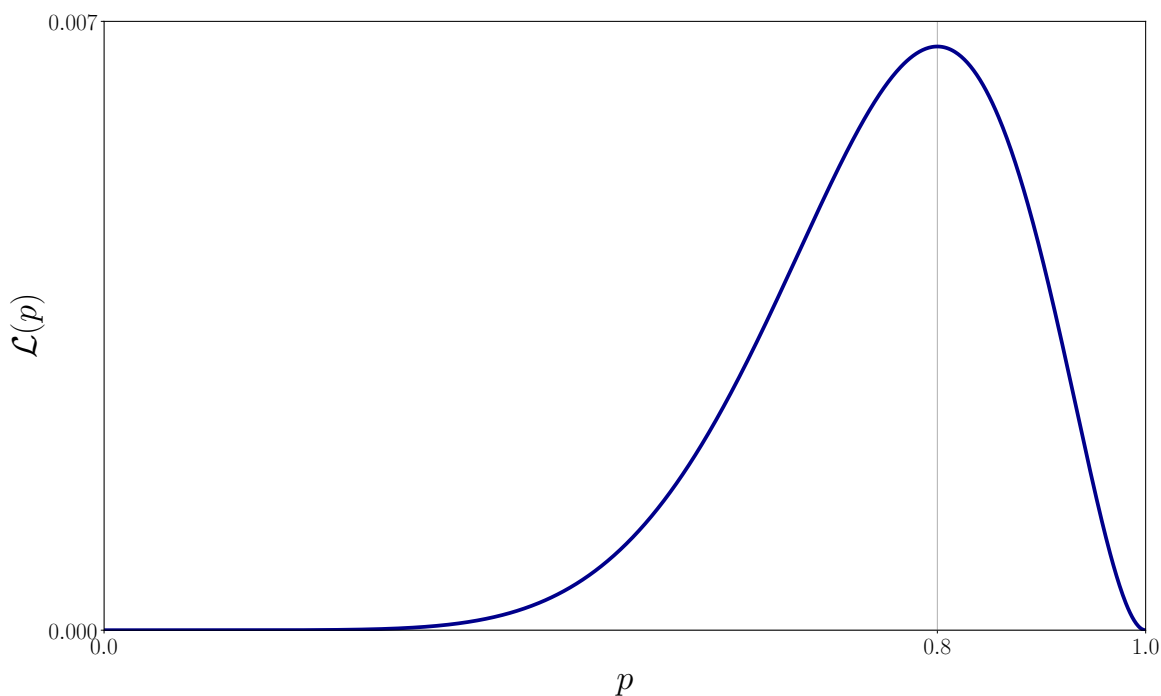


Figure 5: Likelihood function $\mathcal{L}(p)$.

Now the joint probability is not the function of random variable, but the function of parameter p , which shows the likelihood of getting that sample $\{X_1, \dots, X_{10}\}$ with a particular value of parameter p .

This function is called likelihood function $\mathcal{L}(p)$ and achieves its maximum in a value of \hat{p} , called maximum likelihood estimator of p :

$$\mathcal{L}(p) = p^8(1-p)^2.$$

So we need to find the argument of maximum of $\mathcal{L}(p)$:

$$\hat{p} = \arg \max_p \mathcal{L}(p).$$

The intuition suggests that the correct value is $\hat{p} = 0.8$, and it's validated by the graph of $\mathcal{L}(p)$, shown in the fig. 5, but let's derive it explicitly.

According to the FOC of extrema, we should find a derivative of $\mathcal{L}(p)$ and equalize it to zero. In those points the MLE \hat{p} should be found. Additionally, we also should verify that the extremum is indeed maximum with the SOC (second-order derivative is less than zero in a point \hat{p}).

But the likelihood $\mathcal{L}(p)$ is a product (by definition) and is kinda hard to take derivative from. That's why log-likelihood $l(p)$ is introduced:

$$l(p) = \ln \mathcal{L}(p).$$

In log-transformation the product converts into the sum, which is pretty easy to differentiate. Also, the argument of maximum doesn't change after the transform, since the logarithm is a monotonic function:

$$\hat{p} = \arg \max_p \mathcal{L}(p) = \arg \max_p \ln \mathcal{L}(p) = \arg \max_p l(p).$$

The log-likelihood function:

$$l(p) = 8 \ln p + 2 \ln(1-p).$$

According to the FOC of extrema:

$$\left. \frac{dl(p)}{dp} \right|_{p=\hat{p}} = \frac{8}{\hat{p}} - \frac{2}{1-\hat{p}} = 0,$$

$$\boxed{\hat{p} = 0.8}.$$

According to the SOC of extrema:

$$\left. \frac{d^2 l(p)}{dp^2} \right|_{p=\hat{p}} = -\frac{8}{\hat{p}^2} - \frac{2}{(1-\hat{p})^2} = -62.5 < 0,$$

which means that \hat{p} is indeed an argument of maximum.

Problem 6

Let $\{X_1, \dots, X_n\}$ be a random sample from $\text{Exp}(\lambda)$ distribution.

- (a) Derive the MLE of λ .
- (b) State the MLE of $\theta = \lambda^3$.

Solution:

- (a) The p.d.f. of a random variable with distribution $\text{Exp}(\lambda)$ is

$$f(x) = \lambda e^{-\lambda x} \cdot I_{\{x \geq 0\}}.$$

The likelihood function:

$$\mathcal{L}(\lambda) = \prod_{i=1}^n f(X_i; \lambda) = \prod_{i=1}^n \lambda e^{-\lambda X_i} \cdot I_{\{X_i \geq 0\}} = \lambda^n \exp\left(-\lambda \sum_{i=1}^n X_i\right) \cdot \prod_{i=1}^n I_{\{X_i \geq 0\}}.$$

If at least one X_i is less than 0, then we are in the wrong to use $\text{Exp}(\lambda)$ model. If the sample is fine, we can drop indicators from consideration.

The log-likelihood function:

$$l(\lambda) = \ln \mathcal{L}(\lambda) = n \ln \lambda - \lambda \sum_{i=1}^n X_i.$$

According to the FOC of extrema, possible values of the MLE $\hat{\lambda}$ are found as follows:

$$\left. \frac{dl(\lambda)}{d\lambda} \right|_{\lambda=\hat{\lambda}} = \frac{n}{\hat{\lambda}} - \sum_{i=1}^n X_i = 0.$$

$$\boxed{\hat{\lambda} = \frac{1}{\bar{X}}}.$$

According to the SOC of extrema:

$$\left. \frac{d^2 l(\lambda)}{d\lambda^2} \right|_{\lambda=\hat{\lambda}} = -\frac{n}{\hat{\lambda}^2} = -n\bar{X}^2 < 0,$$

which means that $\hat{\lambda}$ is indeed an argument of maximum.

- (b) According to the invariance principle of the MLE:

$$\hat{\theta} = \hat{\lambda}^3 = \left(\frac{1}{\bar{X}} \right)^3.$$

Problem 7

Suppose that X is a discrete random variable with the following probability mass function:

x	0	1	2	3
$P_X(x)$	$\frac{2\theta}{3}$	$\frac{\theta}{3}$	$\frac{2(1-\theta)}{3}$	$\frac{1-\theta}{3}$

where $0 \leq \theta \leq 1$ is a parameter. The following 10 independent observations were taken from such a distribution:

$$(3, 0, 2, 1, 3, 2, 1, 0, 2, 1).$$

What is the maximum likelihood estimate of θ .

Solution:

Since the sample is $(3, 0, 2, 1, 3, 2, 1, 0, 2, 1)$, the likelihood is

$$\begin{aligned} \mathcal{L}(\theta) &= P(X = 3) \cdot P(X = 0) \cdot P(X = 2) \cdot P(X = 1) \cdot P(X = 3) \cdot \\ &\quad \cdot P(X = 2) \cdot P(X = 1) \cdot P(X = 0) \cdot P(X = 2) \cdot P(X = 1). \end{aligned}$$

Substituting from the probability distribution given above, we have

$$\mathcal{L}(\theta) = \prod_{i=1}^n P_X(x_i; \theta) = \left(\frac{2\theta}{3}\right)^2 \left(\frac{\theta}{3}\right)^3 \left(\frac{2(1-\theta)}{3}\right)^3 \left(\frac{1-\theta}{3}\right)^2$$

The log-likelihood function is following:

$$\begin{aligned} l(\theta) &= \log \mathcal{L}(\theta) = \sum_{i=1}^n \log P(x_i; \theta) = \\ &= 2 \left(\log \frac{2}{3} + \log \theta \right) + 3 \left(\log \frac{1}{3} + \log \theta \right) + 3 \left(\log \frac{2}{3} + \log(1-\theta) \right) + \\ &\quad + 2 \left(\log \frac{1}{3} + \log(1-\theta) \right) = 5 \log \frac{1}{3} + 5 \log \frac{2}{3} + 5 \log \theta + 5 \log(1-\theta), \end{aligned}$$

For MLE $\hat{\theta}$ the derivative of $l(\theta)$ with respect to θ equals zero:

$$\left. \frac{dl(\theta)}{d\theta} \right|_{\theta=\hat{\theta}} = \frac{5}{\hat{\theta}} - \frac{5}{1-\hat{\theta}} = 0,$$

and the solution gives us the MLE, which is $\boxed{\hat{\theta} = 0.5}$. The method of moments estimation is $\hat{\theta}_{MM} = \frac{5}{12}$, which is different from MLE.

Problem 8

Let $\{X_1, \dots, X_n\}$ be a random sample from $\mathcal{U}(0, \theta)$ distribution. Find the MLE of θ .

Solution:

The p.d.f. of a random variable with distribution $\mathcal{U}(0, \theta)$ is

$$f(x) = \frac{1}{\theta} \cdot I_{\{0 \leq x \leq \theta\}}.$$

The likelihood function:

$$\mathcal{L}(\theta) = \prod_{i=1}^n f(X_i; \theta) = \prod_{i=1}^n \frac{1}{\theta} \cdot I_{\{0 \leq X_i \leq \theta\}} = \left(\frac{1}{\theta}\right)^n \cdot \prod_{i=1}^n I_{\{0 \leq X_i \leq \theta\}}.$$

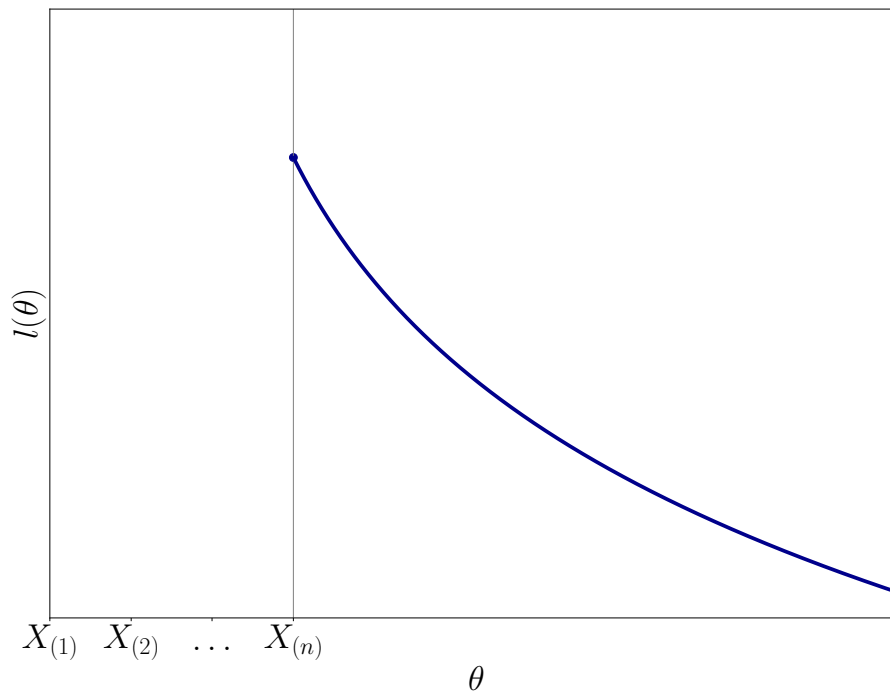


Figure 6: Log-likelihood function $l(\theta)$.

The log-likelihood function:

$$l(\theta) = \ln \mathcal{L}(\theta) = -n \ln \theta + \begin{cases} 0, & X_i \in [0, \theta] \quad \forall i \in \overline{1, n}, \\ -\infty, & \text{otherwise.} \end{cases}$$

The log-likelihood is a function of θ , but the condition in cases is built around X_i . Let's rearrange it.

Having in mind that $\forall i \in \overline{1, n} : X_i \geq 0$ (otherwise the presupposition to use $\mathcal{U}(0, \theta)$ model is wrong):

$$l(\theta) = -n \ln \theta + \begin{cases} 0, & \theta \geq X_i \quad \forall i \in \overline{1, n}, \\ -\infty, & \text{otherwise.} \end{cases}$$

The condition of θ being greater than all X_i means that it's greater than the maximal X_i . In notation of order statistics: $\theta \geq X_{(n)}$. This cutoff is shown in the fig. 6.

The log-likelihood function $l(\theta)$ reaches its maximum in $\theta = X_{(n)}$, which means that it's the maximum likelihood estimator of θ :

$$\boxed{\hat{\theta} = X_{(n)} \equiv \max(X_1, \dots, X_n)}.$$

Problem 9

Suppose that independent observations X and Y are taken from distributions $\text{Gamma}\left(a, \frac{1}{\eta}\right)$ and $\text{Gamma}\left(b, \frac{1}{\eta}\right)$ respectively, where both a and b are known and positive.

Note: The probability density function of a random variable that follows $\text{Gamma}(\alpha, \beta)$ distribution is

$$f(x) = \frac{\beta^\alpha}{\Gamma(\alpha)} x^{\alpha-1} e^{-\beta x}, \quad x > 0,$$

where $\Gamma(\alpha)$ is a function, which represents generalization of the factorial to non-integer numbers:

$$\Gamma(\alpha) = \int_0^\infty x^{\alpha-1} e^{-x} dx, \quad \Re(\alpha) > 0.$$

Its factorial-like behaviour is expressed via following property of recurrence:

$$\Gamma(\alpha + 1) = \alpha \Gamma(\alpha).$$

In particular, in case of α being a positive integer:

$$\Gamma(\alpha) = (\alpha - 1)!$$

- (a) Find the maximum likelihood estimator (MLE) of η .
- (b) Show that the MLE of η is unbiased and find its variance.
- (c) Compare the MLE with the alternative estimator

$$\hat{H} = \frac{1}{2} \left(\frac{X}{a} + \frac{Y}{b} \right)$$

Which one is better?

Solution:

- (a) Let's construct the likelihood function $\mathcal{L}(\eta)$, which is the product of p.d.f.s with known outcomes X and Y and population parameters $\left(a, \frac{1}{\eta}\right)$ and $\left(b, \frac{1}{\eta}\right)$ respectively:

$$\begin{aligned}\mathcal{L}(\eta) &= f\left(X; a, \frac{1}{\eta}\right) \cdot f\left(Y; b, \frac{1}{\eta}\right) = \\ &= \left(\frac{1}{\eta}\right)^a \frac{1}{\Gamma(a)} X^{a-1} e^{-\frac{X}{\eta}} \cdot \left(\frac{1}{\eta}\right)^b \frac{1}{\Gamma(b)} Y^{b-1} e^{-\frac{Y}{\eta}} = \\ &= \left(\frac{1}{\eta}\right)^{a+b} e^{-\frac{X+Y}{\eta}} \frac{X^{a-1} Y^{b-1}}{\Gamma(a)\Gamma(b)}.\end{aligned}$$

Building log-likelihood $l(\eta)$, and using the property of the logarithm of product:

$$l(\eta) = \ln \mathcal{L}(\eta) = -(a+b) \ln \eta - \frac{X+Y}{\eta} + \ln \left(\frac{X^{a-1} Y^{b-1}}{\Gamma(a)\Gamma(b)} \right).$$

Let's find the maximum of the function $l(\eta)$. Firstly, we have to find all possible extrema via Fermat's theorem of stationary points – in those points the derivative of the function $l(\eta)$ equals zero. The derivative of $l(\eta)$ is:

$$\frac{dl}{d\eta} = -\frac{a+b}{\eta} + \frac{X+Y}{\eta^2}.$$

Thus, stationary points are:

$$\frac{dl}{d\eta}(\hat{\eta}) = 0 \quad \Longleftrightarrow \quad -\frac{a+b}{\hat{\eta}} + \frac{X+Y}{\hat{\eta}^2} = 0 \quad \Longleftrightarrow \quad \boxed{\hat{\eta} = \frac{X+Y}{a+b}}.$$

We have to prove that $\hat{\eta}$ is the point of maximum. To do that, we can find the sign of the second derivative of $l(\eta)$ in this point. The second derivative of $l(\eta)$ is:

$$\frac{d^2l}{d\eta^2} = \frac{a+b}{\eta^2} - \frac{2(X+Y)}{\eta^3}.$$

Substituting $\hat{\eta}$ gives:

$$\frac{d^2l}{d\eta^2}(\hat{\eta}) = \frac{a+b}{\hat{\eta}^2} - \frac{2(X+Y)}{\hat{\eta}^3} = -\frac{(a+b)^3}{(X+Y)^2} < 0.$$

Since the second derivative of $l(\eta)$ in $\hat{\eta}$ is negative, $\hat{\eta}$ is the point of maximum. Thus, we have proven that $\hat{\eta}$ is the maximum likelihood estimator of η .

- (b) If the expected value of the estimator equals exactly the estimated parameter, this estimator is considered to be unbiased. In other words, we have to prove that $E(\hat{\eta}) = \eta$.

Using the linearity of expected value:

$$E(\hat{\eta}) = E\left(\frac{X + Y}{a + b}\right) = \frac{1}{a + b} (E(X) + E(Y)).$$

The expected value of outcome X by definition is the expected value of the variable from its population. X was taken from the population with p.d.f. $f\left(x; a, \frac{1}{\eta}\right)$:

$$E(X) = \int_{-\infty}^{+\infty} x f\left(x; a, \frac{1}{\eta}\right) dx = \int_0^{\infty} \left(\frac{1}{\eta}\right)^a \frac{1}{\Gamma(a)} x^a e^{-\frac{x}{\eta}} dx.$$

This integral looks very similar to the definition of gamma function from the note – $\Gamma(\alpha) = \int_0^{\infty} x^{\alpha-1} e^{-x} dx$. In order to get the exponent without scaling parameter $\frac{1}{\eta}$ let's switch to the new variable y :

$$y = \frac{x}{\eta}, \quad dy = \frac{dx}{\eta}, \quad (0, \infty)_x \rightarrow (0, \infty)_y.$$

Substituting y into the previous integral and putting constants out of integral:

$$E(X) = \left(\frac{1}{\eta}\right)^a \frac{1}{\Gamma(a)} \int_0^{\infty} \eta^a y^a e^{-y} dy \cdot \eta = \frac{\eta}{\Gamma(a)} \int_0^{\infty} y^a e^{-y} dy = \frac{\eta}{\Gamma(a)} \cdot \Gamma(a + 1).$$

Using the gamma function's property of recurrence – $\Gamma(\alpha + 1) = \alpha\Gamma(\alpha)$:

$$E(X) = \eta \frac{\Gamma(a + 1)}{\Gamma(a)} = \eta a.$$

Since Y is the outcome from identical population, but with parameter b – $f\left(x; b, \frac{1}{\eta}\right)$, the result of its expected value will be same, but with respective replacement $a \rightarrow b$:

$$E(Y) = \eta b.$$

Thus, the expected value of the estimator $\hat{\eta}$:

$$E(\hat{\eta}) = \frac{1}{a + b} (E(X) + E(Y)) = \frac{1}{a + b} (\eta a + \eta b) = \boxed{\eta}.$$

It means that the estimator $\hat{\eta}$ is unbiased.

Now let's find the variance of $\hat{\eta}$. Since variance is the quadratic function, and outcomes X and Y are independent:

$$\mathbf{V}(\hat{\eta}) = \mathbf{V}\left(\frac{X+Y}{a+b}\right) = \frac{1}{(a+b)^2} (\mathbf{V}(X) + \mathbf{V}(Y)).$$

By definition of variance – $\mathbf{V}(X) = \mathbf{E}(X^2) - \mathbf{E}(X)^2$. The expected value of X^2 is calculated identically to the $\mathbf{E}(X)$:

$$\mathbf{E}(X^2) = \int_{-\infty}^{+\infty} x^2 f\left(x; a, \frac{1}{\eta}\right) dx = \int_0^{\infty} \left(\frac{1}{\eta}\right)^a \frac{1}{\Gamma(a)} x^{a+1} e^{-\frac{x}{\eta}} dx.$$

Substituting $y = \frac{x}{\eta}$ gives:

$$\begin{aligned} \mathbf{E}(X^2) &= \left(\frac{1}{\eta}\right)^a \frac{1}{\Gamma(a)} \int_0^{\infty} \eta^{a+1} y^{a+1} e^{-y} dy \cdot \eta = \frac{\eta^2}{\Gamma(a)} \int_0^{\infty} y^{a+1} e^{-y} dy = \frac{\eta^2}{\Gamma(a)} \cdot \Gamma(a+2) = \\ &= \eta^2(a+1) \frac{\Gamma(a+1)}{\Gamma(a)} = \eta^2(a+1)a. \end{aligned}$$

Thus, the variance $\mathbf{V}(X)$:

$$\mathbf{V}(X) = \mathbf{E}(X^2) - \mathbf{E}(X)^2 = \eta^2(a+1)a - (\eta a)^2 = \eta^2 a.$$

Similarly for Y (same explanation as in $\mathbf{E}(Y)$):

$$\mathbf{V}(Y) = \eta^2 b.$$

Finally, the variance of the estimator $\hat{\eta}$:

$$\mathbf{V}(\hat{\eta}) = \frac{1}{(a+b)^2} (\mathbf{V}(X) + \mathbf{V}(Y)) = \frac{1}{(a+b)^2} (\eta^2 a + \eta^2 b) = \boxed{\frac{\eta^2}{a+b}}.$$

- (c) Let's find out whether the estimator \hat{H} is unbiased. Using linearity of the expected value:

$$\mathbf{E}(\hat{H}) = \mathbf{E}\left(\frac{1}{2} \left(\frac{X}{a} + \frac{Y}{b}\right)\right) = \frac{\mathbf{E}(X)}{2a} + \frac{\mathbf{E}(Y)}{2b} = \frac{\eta a}{2a} + \frac{\eta b}{2b} = \eta.$$

Since $\mathbf{E}(\hat{H}) = \eta$, the estimator is unbiased.

The efficiency of estimators is compared via their Mean-Squared Error (MSE), which is by definition: $\text{MSE}(\hat{\theta}) = \text{V}(\hat{\theta}) + \text{Bias}(\hat{\theta})^2$. Since both estimators ($\hat{\eta}$ and \hat{H}) have zero bias, we should compare their variances.

Let's find the variance of \hat{H} . Since variance is the quadratic function, and outcomes X and Y are independent:

$$\text{V}(\hat{H}) = \text{V}\left(\frac{1}{2}\left(\frac{X}{a} + \frac{Y}{b}\right)\right) = \frac{\text{V}(X)}{4a^2} + \frac{\text{V}(Y)}{4b^2} = \frac{\eta^2 a}{4a^2} + \frac{\eta^2 b}{4b^2} = \frac{\eta^2}{4} \left(\frac{1}{a} + \frac{1}{b}\right).$$

Building the ration of variances gives:

$$\frac{\text{V}(\hat{\eta})}{\text{V}(\hat{H})} = \frac{\eta^2}{a+b} \cdot \frac{4ab}{\eta^2(a+b)} = \frac{4ab}{(a+b)^2}.$$

Let's compare numerator and denominator of the last ratio:

$$\begin{array}{rcl} (a+b)^2 & \vee & 4ab \\ a^2 + 2ab + b^2 & \vee & 4ab \\ a^2 - 2ab + b^2 & \vee & 0 \\ (a-b)^2 & \geq & 0 \\ & \Downarrow & \\ (a+b)^2 & \geq & 4ab \end{array}$$

It means that

$$\frac{\text{V}(\hat{\eta})}{\text{V}(\hat{H})} = \frac{4ab}{(a+b)^2} \leq 1 \quad \implies \quad \text{V}(\hat{\eta}) \leq \text{V}(\hat{H}).$$

Since $4ab = (a+b)^2$ only when $a = b$, estimators $\hat{\eta}$ and \hat{H} are identical in the case of $a = b$. When $a \neq b$, the MLE $\hat{\eta}$ is better than \hat{H} (in terms of MSE).

Thus, the answer is following:

$$\boxed{\begin{cases} \hat{\eta} \text{ is more efficient than } \hat{H}, & a \neq b, \\ \hat{\eta} \equiv \hat{H}, & a = b. \end{cases}}$$

Problem 10

Find maximum likelihood estimator of parameter θ from sample $\{X_1, \dots, X_n\}$ with Laplace($\theta, 1$) distribution and p.d.f.:

$$f(x) = \frac{1}{2}e^{-|x-\theta|}.$$

Solution:

The likelihood function:

$$\mathcal{L}(\theta) = \prod_{i=1}^n f(X_i; \theta) = \prod_{i=1}^n \frac{1}{2}e^{-|X_i-\theta|} = \left(\frac{1}{2}\right)^n \exp\left(-\sum_{i=1}^n |X_i - \theta|\right).$$

The log-likelihood function:

$$l(\theta) = \ln \mathcal{L}(\theta) = -n \ln 2 - \sum_{i=1}^n |X_i - \theta|.$$

Let $g(\theta) = \sum_{i=1}^n |X_i - \theta|$. Since $-n \ln 2$ is a constant:

$$\hat{\theta} = \arg \max_{\theta} l(\theta) = \arg \min_{\theta} g(\theta).$$

Due to a nature of absolute value function, let's explore $g(\theta)$ without differentiation. To arrange the sample, let's use order statistics $X_{(i)}$ instead of common X_i .

- If $n = 1$, $g(\theta) = |X_{(1)} - \theta|$, see fig. 7. The minimum is clearly in $X_{(1)}$.
- If $n = 2$, $g(\theta) = |X_{(1)} - \theta| + |X_{(2)} - \theta|$, see fig. 8. The minimum is in the interval $[X_{(1)}, X_{(2)}]$.
- If $n = 3$, $g(\theta) = |X_{(1)} - \theta| + |X_{(2)} - \theta| + |X_{(3)} - \theta|$, see fig. 9. The minimum is in $X_{(2)}$.
- ...
- If n is arbitrary, $g(\theta) = \sum_{i=1}^n |X_i - \theta|$, see fig. 10. The minimum by induction will be in the median MED of the sample.

Thus, the result:

$$\boxed{\hat{\theta} = \text{MED}(X_1, \dots, X_n)}.$$

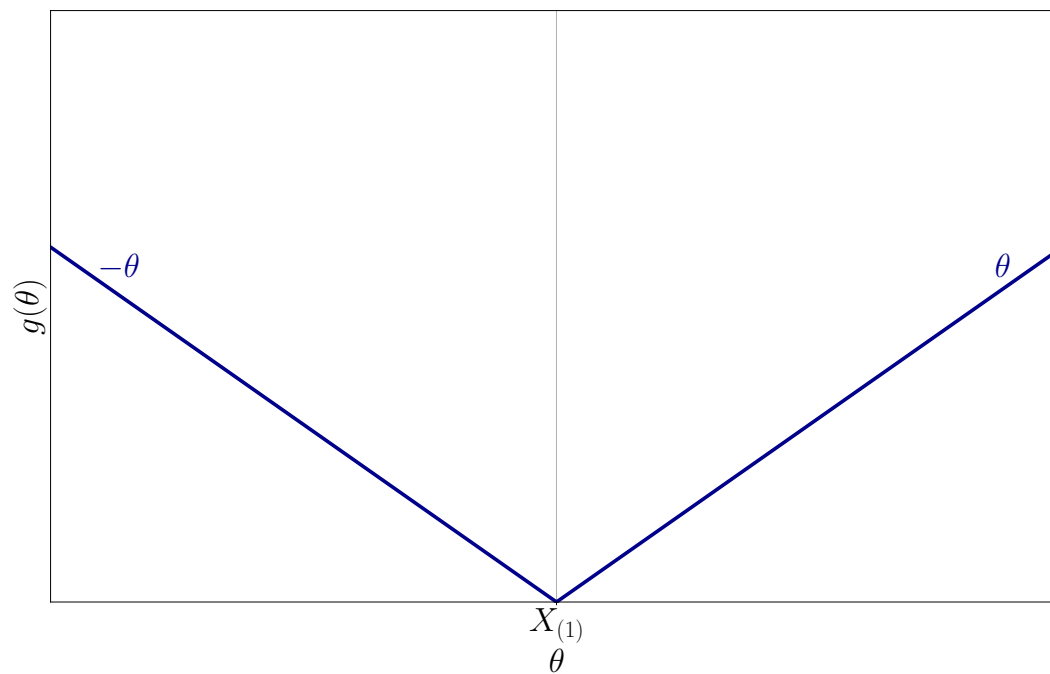


Figure 7: $g(\theta) = \sum_{i=1}^n |X_i - \theta|$ for $n = 1$.

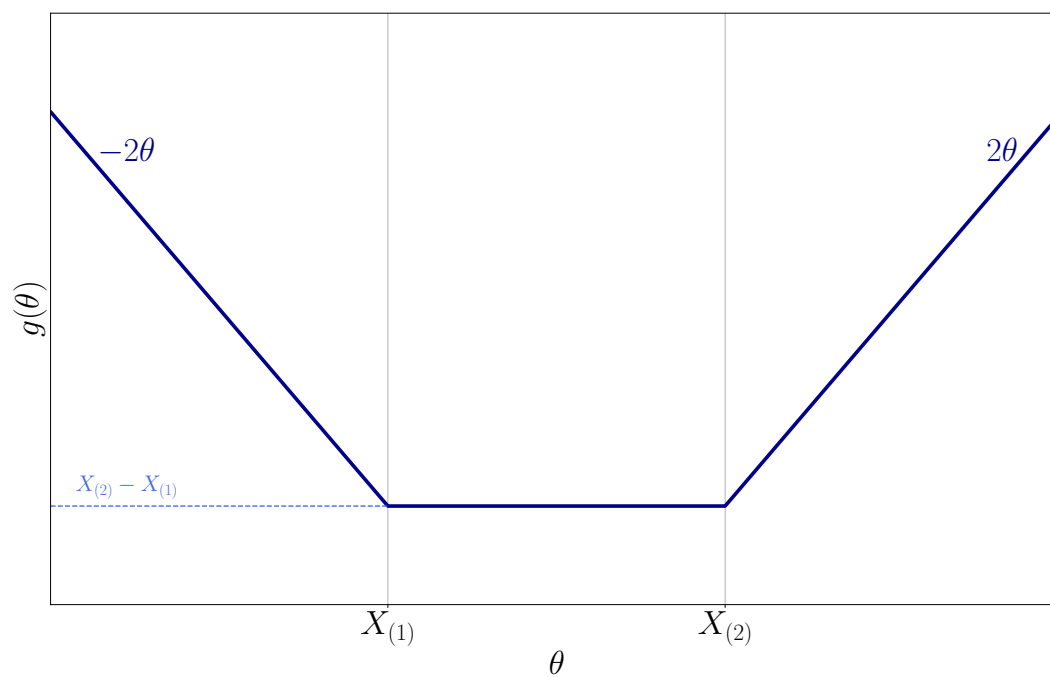


Figure 8: $g(\theta) = \sum_{i=1}^n |X_i - \theta|$ for $n = 2$.

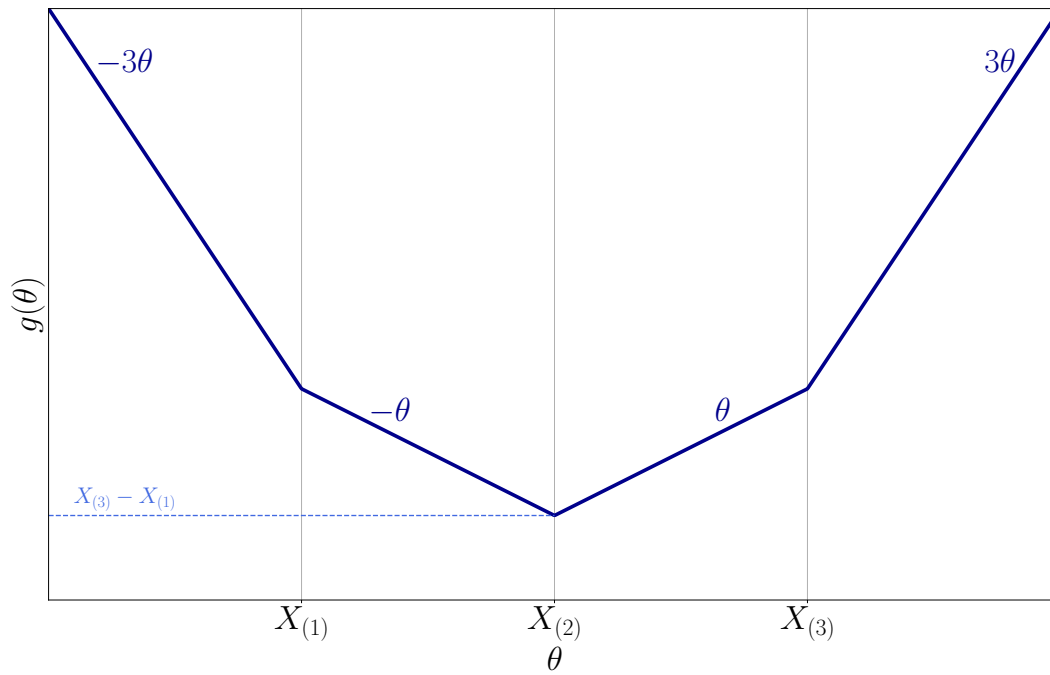


Figure 9: $g(\theta) = \sum_{i=1}^n |X_i - \theta|$ for $n = 3$.

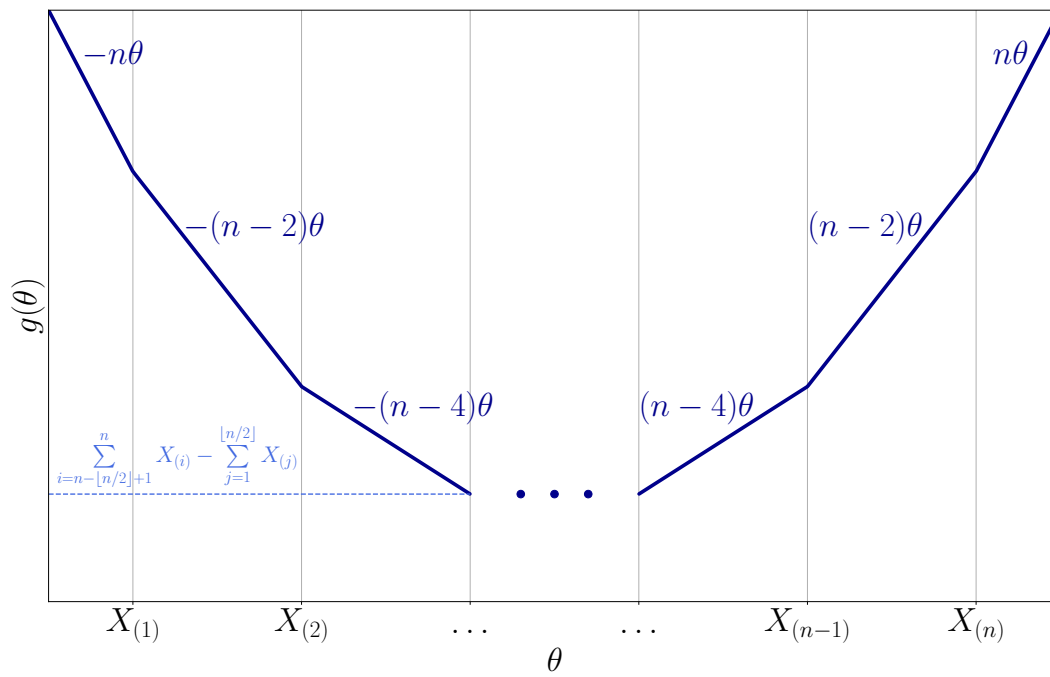


Figure 10: $g(\theta) = \sum_{i=1}^n |X_i - \theta|$ for arbitrary n .