

Quiz

Find a match:

- | | |
|-------------------------------------|---------------------|
| 1. Cumulative distribution function | A. Waiting time |
| 2. Quantile function | B. Separability |
| 3. Pooled variance | C. Approximation |
| 4. Exponential distribution | D. Antiderivative |
| 5. Degrees of freedom | E. Weighted average |
| 6. Central Limit Theorem | F. Goodness-of-fit |
| 7. Independence | G. Constraints |
| 8. Correlation | H. Inverse |

Solution:

1. D (C.d.f. is an antiderivative of p.d.f.)
2. H (Quantile function is an inverse of c.d.f.)
3. E (Pooled variance is a weighted average of several sample variances with weights being degrees of freedom)
4. A (Exponential distribution shows a waiting time in a Poisson process of the next event)
5. G (Degrees of freedom decrease by a number of constraints in a sample)
6. C (The CLT allows approximation of distributions in a large sample with the normal one)
7. B (Independence of random variables is defined as a separability of joint distribution into marginal ones)
8. F (Correlation is a goodness-of-fit metric for a linear regression problem)

Problem 1

A random sample of 400 married couples was selected from a large population of married couples.

- Heights of married men are approximately normally distributed with mean 70 inches and standard deviation 3 inches.
 - Heights of married women are approximately normally distributed with mean 65 inches and standard deviation 2.5 inches.
 - There were 20 couples in which wife was taller than her husband, and there were 380 couples in which wife was shorter than her husband.
1. Find a 95% confidence interval for the proportion of married couples in the population for which the wife is taller than her husband.
 2. Suppose that a married man is selected at random and a married woman is selected at random. Find the approximate probability that the woman will be taller than the man.
 3. Based on your answers to 1 and 2, are the heights of wives and their husbands independent? Explain your reasoning.

Solution:

1. $(1 - \alpha) \cdot 100\%$ confidence interval for p :

$$CI_{1-\alpha}(p) = \hat{P} \pm z_{\alpha/2} \cdot \sqrt{\frac{\hat{P}(1 - \hat{P})}{n}},$$

when \hat{P} can be approximated with normal distribution. The condition of approximation:

$$\begin{cases} n\hat{P} > 5, \\ n(1 - \hat{P}) > 5. \end{cases}$$

A value of sample proportion:

$$\hat{p} = \frac{20}{400} = 0.05.$$

Conditions of approximation are satisfied:

$$\begin{cases} 400 \cdot 0.05 = 20 > 5, \\ 400 \cdot 0.95 = 380 > 5, \end{cases}$$

so the formula for CI can be applied.

From problem statement:

$$\begin{aligned}\text{CI}_{95\%}(p) &= 0.05 \pm z_{0.025} \cdot \sqrt{\frac{0.05 \cdot 0.95}{400}} = \\ &= 0.05 \pm 1.96 \cdot \sqrt{\frac{0.05 \cdot 0.95}{400}} = \\ &= 0.05 \pm 0.021 = \boxed{(0.029; 0.071)}.\end{aligned}$$

2. Let X be a height of women, and let Y be a height of men. Those random variables are independent. From problem statement:

$$X \sim \mathcal{N}(65, 2.5^2), \quad Y \sim \mathcal{N}(70, 3^2).$$

The difference of heights is distributed as follows:

$$X - Y \sim \mathcal{N}(65 - 70, 2.5^2 + (-1)^2 \cdot 3^2) = \mathcal{N}(-5, 3.905^2).$$

The required probability:

$$P(X > Y) = P(X - Y > 0) = P\left(Z > \frac{0 - (-5)}{3.905}\right) \approx 1 - \Phi(1.28) \approx 1 - 0.9 = \boxed{0.1}.$$

3. Under the assumption of independent choice of both women and men, the proportion of couples, where wives are taller than husbands, should be equal to 0.1, as it was found in task 2.

However, the confidence interval for the real proportion of couples with such distinction is (0.029; 0.071), as it was found in task 1.

The actual confidence interval does not include 0.1, so there is a significant evidence, that couples are constituted out of men and women, which were chosen NOT independently.

Problem 2

Suppose 2000 points are selected independently at a random from the unit square

$S = \{(x, y) : 0 \leq x \leq 1, 0 \leq y \leq 1\}$. Let W be the number of points that fall into the set $A = \{(x, y) : x^2 + y^2 < 1\}$.

1. How is W distributed?
2. Find the mean, variance and standard deviation of W .
3. Estimate probability that W is greater than 1600.

Solution:

1. Some chosen point from the set S can be either in the set A , or out of the set A . So an event of a point being in the set A is a Bernoulli trial with a probability of success:

$$p = \frac{\mu(A \cap S)}{\mu(S)} = \frac{\pi}{4}.$$

A number of successes in the series of Bernoulli trials is a binomial random variable. Since there are 2000 points:

$$W \sim \text{Bin}\left(2000, \frac{\pi}{4}\right).$$

2. Moments of a random variable with binomial distribution are well-known:

$$E(W) = 2000 \cdot \frac{\pi}{4} = \boxed{500\pi \approx 1570.8},$$

$$V(W) = 2000 \cdot \frac{\pi}{4} \cdot \left(1 - \frac{\pi}{4}\right) = \boxed{500\pi - 125\pi^2 \approx 337.1},$$

$$\sigma(W) = \sqrt{V(W)} = \boxed{\sqrt{500\pi - 125\pi^2} \approx 18.36}.$$

3. According to the De Moivre–Laplace theorem:

$$W \sim \text{Bin}\left(2000, \frac{\pi}{4}\right) \stackrel{d}{\approx} W_{\text{CLT}} \sim \mathcal{N}(500\pi, 500\pi - 125\pi^2),$$

and since we approximate discrete distribution with continuous, we should conduct continuity correction. Point $W = 1600$ is excluded:

$$\begin{aligned} P(W > 1600) &\approx P(W_{\text{CLT}} > 1600.5) = P\left(Z > \frac{1600.5 - 1570.8}{18.36}\right) \approx 1 - \Phi(1.62) \approx \\ &\approx 1 - 0.947 = \boxed{0.053}. \end{aligned}$$

Problem 3

Distribution of X is uniform $\mathcal{U}(-a, a)$. Sample of size $n = 2$ is available.

Consider $\hat{a} = c \cdot (|X_1| + |X_2|)$ as a class of estimators for the parameter a .

Find c such that

1. Estimator \hat{a} is unbiased.
2. Estimator \hat{a} is the most efficient in the class. (In terms of mean square error.)

Solution:

1. An estimator is unbiased if:

$$\mathbb{E}(\hat{a}) = a.$$

Let's find $\mathbb{E}(\hat{a})$. Since X_1 and X_2 are i.i.d., and using linearity of expected value:

$$\mathbb{E}(\hat{a}) = \mathbb{E}(c \cdot (|X_1| + |X_2|)) = c \cdot 2 \mathbb{E}(|X|).$$

A p.d.f. of X is:

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

According to LOTUS:

$$\mathbb{E}(|X|) = \int_{-\infty}^{+\infty} |x| f(x) dx = \frac{1}{2a} \int_{-a}^a |x| dx = \frac{1}{2a} \cdot 2 \int_0^a x dx = \frac{1}{a} \cdot \frac{x^2}{2} \Big|_0^a = \frac{a}{2}.$$

Thus:

$$\mathbb{E}(\hat{a}) = c \cdot 2 \cdot \frac{a}{2} = ac.$$

Applying this result to the definition of unbiasedness:

$$ac = a, \quad \implies \quad \boxed{c = 1}.$$

2. Here we need to minimize the MSE of an estimator:

$$c^* = \arg \min_c \text{MSE}(\hat{a}).$$

The MSE:

$$\text{MSE}(\hat{a}) = \text{V}(\hat{a}) + \text{Bias}^2(\hat{a}).$$

The bias is found with $E(\hat{a})$ from task 1:

$$\text{Bias}(\hat{a}) = E(\hat{a}) - a = ac - a = a(c - 1).$$

Since X_1 and X_2 are i.i.d., and since variance is a quadratic function:

$$V(\hat{a}) = V(c \cdot (|X_1| + |X_2|)) = c^2 \cdot 2V(|X|).$$

Using common variance identity:

$$V(|X|) = E(|X|^2) - E(|X|)^2 = E(X^2) - \frac{a^2}{4}.$$

According to LOTUS:

$$E(X^2) = \int_{-\infty}^{+\infty} x^2 f(x) dx = \frac{1}{2a} \int_{-a}^a x^2 dx = \frac{1}{2a} \cdot \frac{x^3}{3} \Big|_{-a}^a = \frac{a^2}{3}.$$

Thus:

$$\begin{aligned} V(|X|) &= \frac{a^2}{3} - \frac{a^2}{4} = \frac{a^2}{12}, \\ V(\hat{a}) &= c^2 \cdot 2 \cdot \frac{a^2}{12} = \frac{c^2 a^2}{6}, \\ \text{MSE}(\hat{a}) &= \frac{c^2 a^2}{6} + a^2(c - 1)^2. \end{aligned}$$

Let's represent $\text{MSE}(\hat{a})$ as a quadratic polynomial with variable c :

$$\text{MSE}(\hat{a}) = \frac{7a^2}{6}c^2 - 2a^2c + a^2.$$

A graph of this function is a parabola with a vertex being a minimum point. A vertex point of the polynomial $ax^2 + bx + c$ is:

$$x^* = -\frac{b}{2a},$$

so the minimal c is:

$$c^* = -\frac{-2a^2}{2 \cdot 7a^2/6} = \boxed{\frac{6}{7}}.$$

Problem 4

Consider random variables X and Y with joint density function

$$f(x, y) = \begin{cases} \frac{1}{2} + cx, & x + y \leq 1, \ x \geq 0, \ y \geq 0, \\ 0, & \text{otherwise.} \end{cases}$$

1. Find c .
2. Find $f_X(x)$. Evaluate $E(X)$.
3. Write down an expression for $f_{Y|X}(x, y)$. Find $E(Y | X = x)$.

Solution:

1. Using normalization condition:

$$\iint_{-\infty}^{+\infty} f(x, y) dx dy = 1,$$

$$\begin{aligned} \iint_{-\infty}^{+\infty} f(x, y) dx dy &= \int_0^1 \int_0^{1-x} \left(\frac{1}{2} + cx \right) dy dx = \int_0^1 \left(\frac{1}{2} + cx \right) (1-x) dx = \\ &= \int_0^1 \left(\frac{1}{2} + \left(c - \frac{1}{2} \right) x - cx^2 \right) dx = \left. \frac{x}{2} \right|_0^1 + \left(c - \frac{1}{2} \right) \left. \frac{x^2}{2} \right|_0^1 - c \cdot \left. \frac{x^3}{3} \right|_0^1 = \\ &= \frac{1}{2} + \frac{c}{2} - \frac{1}{4} - \frac{c}{3} = \frac{c}{6} + \frac{1}{4}. \end{aligned}$$

Applying this result to the condition:

$$\frac{c}{6} + \frac{1}{4} = 1, \quad \implies \quad \boxed{c = \frac{9}{2}}.$$

2. Marginal distribution by definition:

$$f_X(x) = \int_{-\infty}^{+\infty} f(x, y) dy = \int_0^{1-x} \left(\frac{1}{2} + \frac{9x}{2} \right) dy \cdot I_{\{0 \leq x \leq 1\}} = \boxed{\frac{(1+9x)(1-x)}{2} \cdot I_{\{0 \leq x \leq 1\}}}.$$

Expected value by definition:

$$\begin{aligned} \mathbb{E}(X) &= \int_{-\infty}^{+\infty} x f_X(x) dx = \int_0^1 x \cdot \frac{(1+9x)(1-x)}{2} dx = \\ &= \frac{1}{2} \left[\int_0^1 x dx + \int_0^1 8x^2 dx - \int_0^1 9x^3 dx \right] = \frac{1}{2} \left[\frac{x^2}{2} \Big|_0^1 + \frac{8x^3}{3} \Big|_0^1 - \frac{9x^4}{4} \Big|_0^1 \right] = \\ &= \frac{1}{2} \cdot \left(\frac{1}{2} + \frac{8}{3} - \frac{9}{4} \right) = \boxed{\frac{11}{24}}. \end{aligned}$$

3. Conditional p.d.f. by definition:

$$f_{Y|X}(y | x) = \frac{f(x, y)}{f_X(x)} = \frac{\frac{1+9x}{2} \cdot I_{\{x+y \leq 1, x \geq 0, y \geq 0\}}}{\frac{(1+9x)(1-x)}{2} \cdot I_{\{0 \leq x \leq 1\}}} = \boxed{\frac{1}{1-x} \cdot I_{\{x+y \leq 1, x > 0, y \geq 0\}}}.$$

Conditional expectation by definition:

$$\begin{aligned} \mathbb{E}(Y | X = x) &= \int_{-\infty}^{+\infty} y \cdot f_{Y|X}(y | x) dy = \int_0^{1-x} y \cdot \frac{1}{1-x} dy \cdot I_{\{0 \leq x \leq 1\}} = \\ &= \frac{1}{1-x} \cdot I_{\{0 \leq x \leq 1\}} \cdot \int_0^{1-x} y dy = \frac{1}{1-x} \cdot I_{\{0 \leq x \leq 1\}} \cdot \frac{y^2}{2} \Big|_0^{1-x} = \\ &= \frac{1}{1-x} \cdot \frac{(1-x)^2}{2} \cdot I_{\{0 \leq x \leq 1\}} = \boxed{\frac{1-x}{2} \cdot I_{\{0 \leq x \leq 1\}}}. \end{aligned}$$

Problem 5

Internal angles $\theta_1, \theta_2, \theta_3, \theta_4$ of a certain quadrilateral, located on the ground, were measured by the aerial system. It is assumed that those observations x_1, x_2, x_3, x_4 were taken with minor and independent errors, which have zero mean and identical variance σ^2 .

1. Find the LSE of $\theta_1, \theta_2, \theta_3, \theta_4$.
2. Find an unbiased estimate of σ^2 in the case, described in part 1.
3. Let's assume now that the considered quadrilateral is a parallelogram with $\theta_1 = \theta_3$ and $\theta_2 = \theta_4$. How values of internal angles LSE would change? Find an unbiased estimate of σ^2 in this particular case.

Solution:

1. A model of observation is based upon a small error ε with $E(\varepsilon) = 0$ and $V(\varepsilon) = \sigma^2$:

$$X_i = \theta_i + \varepsilon_i \quad \forall i = \overline{1, 4},$$

so we know that $E(X_i) = \theta_i$ and $V(X_i) = \sigma^2$.

The LSE is determined by minimizing RSS function:

$$\text{RSS} = \sum_{i=1}^4 \varepsilon_i^2 = \sum_{i=1}^4 (X_i - \theta_i)^2.$$

But here we haven't take into account the fact that we are operating with a quadrilateral. The sum of all angles is 2π :

$$\sum_{i=1}^4 \theta_i = 2\pi.$$

Thus, we need to solve a constrained minimization problem:

$$\begin{aligned} \hat{\theta}_1, \hat{\theta}_2, \hat{\theta}_3, \hat{\theta}_4 &= \arg \min_{\theta_1, \theta_2, \theta_3, \theta_4} \sum_{i=1}^4 (X_i - \theta_i)^2, \\ \text{s.t. } \sum_{i=1}^4 \theta_i &= 2\pi. \end{aligned}$$

Let's build a Lagrange function \mathcal{L} :

$$\mathcal{L} = \sum_{i=1}^4 (X_i - \theta_i)^2 - \lambda \left(\sum_{i=1}^4 \theta_i - 2\pi \right).$$

By necessary condition of extremum:

$$\begin{cases} \left. \frac{\partial \mathcal{L}}{\partial \theta_i} \right|_{\theta_i = \hat{\theta}_i, \lambda = \hat{\lambda}} = 2(\hat{\theta}_i - X_i) - \hat{\lambda} = 0, \\ \left. \frac{\partial \mathcal{L}}{\partial \lambda} \right|_{\theta_i = \hat{\theta}_i, \lambda = \hat{\lambda}} = 2\pi - \sum_{i=1}^4 \hat{\theta}_i = 0. \end{cases}$$

Summing up the first equation 4 times for all i we get:

$$2 \sum_{i=1}^4 \hat{\theta}_i - 2 \sum_{i=1}^4 X_i = \hat{\lambda}.$$

Substituting it into the second equation we get:

$$\begin{cases} \hat{\theta}_i = X_i + \frac{1}{4} \left(2\pi - \sum_{i=1}^4 X_i \right), \\ \hat{\lambda} = \pi - \frac{1}{2} \sum_{i=1}^4 X_i. \end{cases}$$

We also need to prove that values of $\hat{\theta}_i$ are indeed arguments of minimum. Second derivatives:

$$\left. \frac{\partial^2 \mathcal{L}}{\partial \theta_i^2} \right|_{\theta_i = \hat{\theta}_i} = 2, \quad \left. \frac{\partial^2 \mathcal{L}}{\partial \theta_i \partial \theta_j} \right|_{\theta_i = \hat{\theta}_i, \theta_j = \hat{\theta}_j, i \neq j} = 0.$$

So the second differential of \mathcal{L} then:

$$d^2 \mathcal{L} = \sum_{i=1}^4 2d\theta_i^2 > 0.$$

By sufficient condition it is indeed a minimum. Q.E.D.

From problem statement:

$$\boxed{\hat{\theta}_i = x_i + \frac{1}{4} \left(2\pi - \sum_{i=1}^4 x_i \right)}.$$

Let's note that if there haven't been a constraint on a sum of angles, the minimization would have given a simple result of $\hat{\theta}_i = x_i$.

2. The descriptive statistic of σ^2 is an RSS function. Let's explore an expected value its estimate for a one i :

$$\begin{aligned} \mathbb{E} \left[\left(X_i - \hat{\theta}_i \right)^2 \right] &= \mathbb{E} \left[\frac{1}{4} \left(\sum_{i=1}^4 X_i - 2\pi \right) \right]^2 = \frac{1}{16} \mathbb{E} \left(\sum_{i=1}^4 X_i - 2\pi \right)^2 = \\ &= \frac{1}{16} \mathbb{E} \left(\sum_{i=1}^4 X_i - \mathbb{E} \left(\sum_{i=1}^4 X_i \right) \right)^2 = \frac{1}{16} \mathbb{V} \left(\sum_{i=1}^4 X_i \right) = \frac{1}{16} \cdot 4\sigma^2 = \frac{\sigma^2}{4}. \end{aligned}$$

It means that the sum of 4 observations will give unbiased result:

$$\mathbb{E} \left[\sum_{i=1}^4 \left(X_i - \hat{\theta}_i \right)^2 \right] = \sigma^2.$$

So the required estimator is:

$$\boxed{\hat{\sigma}^2 = \sum_{i=1}^4 \left(X_i - \hat{\theta}_i \right)^2}.$$

3. A new Lagrange function with $\theta_1 = \theta_3$ and $\theta_2 = \theta_4$:

$$\mathcal{L} = (X_1 - \theta_1)^2 + (X_2 - \theta_2)^2 + (X_3 - \theta_1)^2 + (X_4 - \theta_2)^2 - 2\lambda (\theta_1 + \theta_2 - \pi).$$

By necessary condition of extremum:

$$\begin{cases} \left. \frac{\partial \mathcal{L}}{\partial \theta_1} \right|_{\theta_1=\hat{\theta}_1, \theta_2=\hat{\theta}_2, \lambda=\hat{\lambda}} = 2 \left(\hat{\theta}_1 - X_1 \right) + 2 \left(\hat{\theta}_1 - X_3 \right) - \hat{\lambda} = 0, \\ \left. \frac{\partial \mathcal{L}}{\partial \theta_2} \right|_{\theta_1=\hat{\theta}_1, \theta_2=\hat{\theta}_2, \lambda=\hat{\lambda}} = 2 \left(\hat{\theta}_2 - X_2 \right) + 2 \left(\hat{\theta}_2 - X_4 \right) - \hat{\lambda} = 0, \\ \left. \frac{\partial \mathcal{L}}{\partial \lambda} \right|_{\theta_1=\hat{\theta}_1, \theta_2=\hat{\theta}_2, \lambda=\hat{\lambda}} = 2 \left(\pi - \hat{\theta}_1 - \hat{\theta}_2 \right) = 0. \end{cases}$$

Summing up the first 2 equations we get:

$$2 \left(\hat{\theta}_1 + \hat{\theta}_2 \right) - \sum_{i=1}^4 X_i = \hat{\lambda}.$$

Substituting it into the third equation we get:

$$\begin{cases} \hat{\theta}_1 = \frac{X_1 + X_3}{2} + \frac{1}{4} \left(2\pi - \sum_{i=1}^4 X_i \right), \\ \hat{\theta}_2 = \frac{X_2 + X_4}{2} + \frac{1}{4} \left(2\pi - \sum_{i=1}^4 X_i \right), \\ \hat{\lambda} = 2\pi - \sum_{i=1}^4 X_i. \end{cases}$$

We also need to prove that values of $\hat{\theta}_i$ are indeed arguments of minimum. Second derivatives:

$$\frac{\partial^2 \mathcal{L}}{\partial \theta_i^2} \Big|_{\theta_i = \hat{\theta}_i} = 4, \quad \frac{\partial^2 \mathcal{L}}{\partial \theta_i \partial \theta_j} \Big|_{\theta_i = \hat{\theta}_i, \theta_j = \hat{\theta}_j, i \neq j} = 0.$$

So the second differential of \mathcal{L} then:

$$d^2 \mathcal{L} = \sum_{i=1}^2 4d\theta_i^2 > 0.$$

By sufficient condition it is indeed a minimum. Q.E.D.

From problem statement:

$$\begin{cases} \hat{\theta}_1 = \frac{x_1 + x_3}{2} + \frac{1}{4} \left(2\pi - \sum_{i=1}^4 x_i \right), \\ \hat{\theta}_2 = \frac{x_2 + x_4}{2} + \frac{1}{4} \left(2\pi - \sum_{i=1}^4 x_i \right). \end{cases}$$

Let's explore one term of estimated RSS. Estimator $\hat{\theta}_1$ is unbiased:

$$\mathbb{E}(\hat{\theta}_1) = \mathbb{E} \left(\frac{X_1 + X_3}{2} + \frac{1}{4} \left(2\pi - \sum_{i=1}^4 X_i \right) \right) = \frac{\theta_1 + \theta_1}{2} + \frac{1}{4} (2\pi - 2(\theta_1 + \theta_2)) = \theta_1.$$

It means that $\mathbb{E}(X_1 - \hat{\theta}_1) = 0$ and inherently:

$$\mathbb{E} \left[(X_1 - \hat{\theta}_1)^2 \right] = \mathbb{V} (X_1 - \hat{\theta}_1).$$

Exploring further:

$$\begin{aligned} \mathbb{E} \left[(X_1 - \hat{\theta}_1)^2 \right] &= \mathbb{E} \left(\frac{3X_1}{4} - \frac{X_3}{4} + \frac{X_2}{4} + \frac{X_4}{4} - \frac{\pi}{2} \right)^2 = \\ &= \mathbb{V} \left(\frac{3X_1}{4} - \frac{X_3}{4} + \frac{X_2}{4} + \frac{X_4}{4} \right) = \left[\left(\frac{3}{4} \right)^2 + 3 \cdot \left(\frac{1}{4} \right)^2 \right] \sigma^2 = \frac{3\sigma^2}{4}. \end{aligned}$$

Adding up the rest of terms will give:

$$\mathbb{E} \left[(X_1 - \hat{\theta}_1)^2 + (X_2 - \hat{\theta}_2)^2 + (X_3 - \hat{\theta}_1)^2 + (X_4 - \hat{\theta}_2)^2 \right] = 3\sigma^2.$$

So the required estimator is:

$$\hat{\sigma}^2 = \frac{1}{3} \left[(X_1 - \hat{\theta}_1)^2 + (X_2 - \hat{\theta}_2)^2 + (X_3 - \hat{\theta}_1)^2 + (X_4 - \hat{\theta}_2)^2 \right].$$

Problem 6

Suppose that student's grade for a statistics exam, X , has continuous uniform distribution at the interval $[0, 100]$. But less than 25 points means "failed", and more than 80 points is "excellent", hence the final grade Y is calculated as follows:

$$Y = \begin{cases} 0, & X < 25 \\ X, & 25 \leq X < 80 \\ 100, & X \geq 80 \end{cases}$$

1. Find c.d.f. of Y . Sketch the plot.
2. Find p.d.f. of Y . Sketch the plot.
3. Find mean and variance of X and Y .
4. Find $E(Y \mid Y > 0)$.
5. Find $\text{Corr}(X, Y)$.

Solution:

1. A random variable X has following p.d.f. and c.d.f.:

$$f_X(x) = \frac{1}{100} \cdot I_{\{0 \leq x \leq 100\}},$$

$$F_X(x) = \begin{cases} 1, & x \geq 100, \\ \frac{x}{100}, & 0 \leq x < 100, \\ 0, & x < 0. \end{cases}$$

The c.d.f. of Y coincides with that of X for $y \in [25, 80)$:

$$F_Y(y) = \frac{y}{100}, \text{ for } 25 \leq y < 80.$$

$F_Y(y)$ has a sharp increase in $y = 0$ by a value of $P(X < 25) = \frac{1}{4}$, and is maintained on the same level up to $y = 25$:

$$F_Y(y) = \frac{1}{4}, \text{ for } 0 \leq y < 25.$$

Also, $F_Y(y)$ does not change to the right of $y = 80$ until it hits $y = 100$, where it has a sharp increase by a value of $P(X \geq 80) = \frac{1}{5}$:

$$F_Y(y) = \frac{4}{5}, \text{ for } 80 \leq y < 100.$$

Thus, the c.d.f. of Y has a following view:

$$F_Y(y) = \begin{cases} 1, & y \geq 100, \\ \frac{4}{5}, & 80 \leq y < 100, \\ \frac{y}{100}, & 25 \leq y < 80, \\ \frac{1}{4}, & 0 \leq y < 25, \\ 0, & y < 0. \end{cases}$$

A graph of $F_Y(y)$ is shown in the fig. 1.

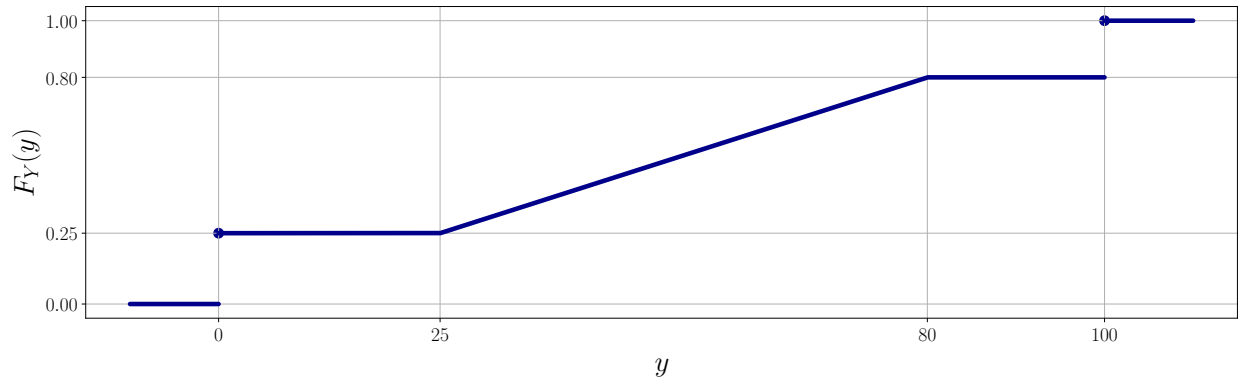


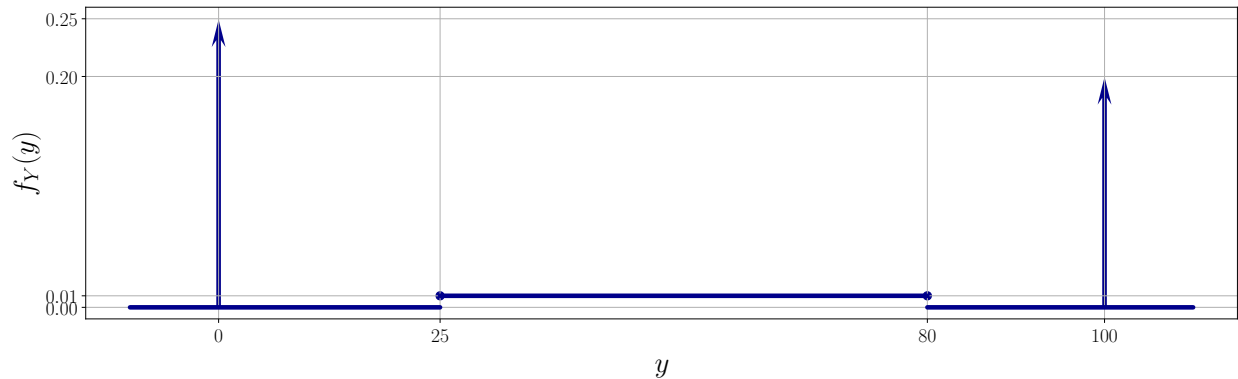
Figure 1: C.d.f. of the random variable Y .

2. Since $F_Y(y)$ has points of discontinuity $y = 0$ and $y = 100$, it only has generalized p.d.f. Taking derivative of c.d.f.:

$$f_Y(y) = \frac{dF_Y(y)}{dy} = \left[\frac{1}{100} \cdot I_{\{25 \leq y \leq 80\}} + \frac{1}{4} \delta(y) + \frac{1}{5} \delta(y - 100) \right],$$

where $\delta(y)$ is a Dirac delta function, defined as:

$$\delta(y) \simeq \begin{cases} +\infty, & y = 0, \\ 0, & y \neq 0, \end{cases} \quad \text{constrained by} \quad \int_{-\infty}^{+\infty} \delta(y) dy = 1.$$

Figure 2: Generalized p.d.f. of the random variable Y .

A graph of $f_Y(y)$ is shown in the fig. 2.

The $f_Y(y)$ possesses the most important property of probability densities – normalization by 1:

$$\begin{aligned} \int_{-\infty}^{+\infty} f_Y(y) dy &= \int_{25}^{80} \frac{1}{100} dy + \int_{-\infty}^{+\infty} \frac{1}{4} \delta(y) dy + \int_{-\infty}^{+\infty} \frac{1}{5} \delta(y - 100) dy = \\ &= \frac{y}{100} \Big|_{25}^{80} + \frac{1}{4} \int_{-\infty}^{+\infty} \delta(y) dy + \frac{1}{5} \underbrace{\int_{-\infty}^{+\infty} \delta(\eta) d\eta}_{\eta=y-100} = \frac{55}{100} + \frac{1}{4} + \frac{1}{5} = 1. \end{aligned}$$

3. The mean and variance of a random variable $X \sim \mathcal{U}(a, b)$:

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

In our case:

$$\mathbb{E}(X) = \frac{0+100}{2} = \boxed{50}, \quad \mathbb{V}(X) = \frac{(b-a)^2}{12} = \frac{(100-0)^2}{12} = \boxed{\frac{2500}{3} \approx 833.333}.$$

There are 2 ways to calculate mean and variance of Y .

(I) Direct calculation

By definition:

$$\mathbb{E}(Y) = \int_{-\infty}^{+\infty} y f_Y(y) dy = \int_{25}^{80} \frac{y}{100} dy + \frac{1}{4} \int_{-\infty}^{+\infty} y \delta(y) dy + \frac{1}{5} \int_{-\infty}^{+\infty} y \delta(y - 100) dy.$$

Using a sifting property of delta-function:

$$\int_{-\infty}^{+\infty} \varphi(y) \delta(y - b) dy = \varphi(b),$$

we can proceed with calculation expected value of Y :

$$E(Y) = \frac{y^2}{200} \Big|_{25}^{80} + \frac{1}{4} y \Big|_{y=0} + \frac{1}{5} y \Big|_{y=100} = \frac{231}{8} + 0 + 20 = \boxed{\frac{391}{8} = 48.875}.$$

In order to calculate variance $V(Y)$, we need to find $E(Y^2)$:

$$\begin{aligned} E(Y^2) &= \int_{-\infty}^{+\infty} y^2 f_Y(y) dy = \int_{25}^{80} \frac{y^2}{100} dy + \frac{1}{4} \int_{-\infty}^{+\infty} y^2 \delta(y) dy + \frac{1}{5} \int_{-\infty}^{+\infty} y^2 \delta(y - 100) dy = \\ &= \frac{y^3}{300} \Big|_{25}^{80} + \frac{1}{4} y^2 \Big|_{y=0} + \frac{1}{5} y^2 \Big|_{y=100} = \frac{19855}{12} + 0 + 2000 = \\ &= \frac{43855}{12} \approx 3654.583. \end{aligned}$$

Variance then:

$$V(Y) = E(Y^2) - E(Y)^2 = \frac{43855}{12} - \left(\frac{391}{8}\right)^2 = \boxed{\frac{243037}{192} \approx 1265.818}.$$

(II) Total expectation

A random variable Y consists of 3 parts exhaustive, which is reflected in the total expectation equation:

$$\begin{aligned} E(Y) &= E(Y | X < 25) \cdot P(X < 25) + E(Y | 25 \leq X < 80) \cdot P(25 \leq X < 80) + \\ &+ E(Y | X \geq 80) \cdot P(X \geq 80) = 0 \cdot \frac{1}{4} + \frac{25+80}{2} \cdot \frac{55}{100} + 100 \cdot \frac{1}{5} = \\ &= 0 + \frac{231}{8} + 20 = \boxed{\frac{391}{8} = 48.875}. \end{aligned}$$

The same for $E(Y^2)$:

$$\begin{aligned} E(Y^2) &= E(Y^2 | X < 25) \cdot P(X < 25) + \\ &+ E(Y^2 | 25 \leq X < 80) \cdot P(25 \leq X < 80) + E(Y^2 | X \geq 80) \cdot P(X \geq 80). \end{aligned}$$

While the first and the third terms are clear, the second one in the equation above requires more calculations:

$$\begin{aligned} E(Y^2 | 25 \leq X < 80) &= E(X^2 | 25 \leq X < 80) = \\ &= V(X | 25 \leq X < 80) + E(X | 25 \leq X < 80)^2 = \\ &= \frac{(80 - 25)^2}{12} + \left(\frac{25 + 80}{2}\right)^2 = \frac{9025}{3} \approx 3008.333. \end{aligned}$$

$E(Y^2)$ then:

$$E(Y^2) = 0 \cdot \frac{1}{4} + \frac{9025}{3} \cdot \frac{55}{100} + 100^2 \cdot \frac{1}{5} = \frac{43855}{12} \approx 3654.583.$$

Variance then:

$$V(Y) = E(Y^2) - E(Y)^2 = \frac{43855}{12} - \left(\frac{391}{8}\right)^2 = \boxed{\frac{243037}{192} \approx 1265.818}.$$

4. Total expectation:

$$E(Y) = E(Y | Y > 0) \cdot P(Y > 0) + E(Y | Y \leq 0) \cdot P(Y \leq 0).$$

Since $E(Y | Y \leq 0) = 0$:

$$E(Y | Y > 0) = \frac{E(Y)}{P(Y > 0)} = \frac{391}{8} : \left(1 - \frac{1}{4}\right) = \boxed{\frac{391}{6} \approx 65.167}.$$

5. Using the definition of correlation coefficient:

$$\text{Corr}(X, Y) = \frac{\text{Cov}(X, Y)}{\sqrt{V(X) \cdot V(Y)}} = \frac{E(XY) - E(X) \cdot E(Y)}{\sqrt{V(X) \cdot V(Y)}}.$$

The only term we don't know is $E(XY)$. Since the joint distribution is unknown, let's use total expectation:

$$\begin{aligned} E(XY) &= E(XY | X < 25) \cdot P(X < 25) + E(XY | 25 \leq X < 80) \cdot P(25 \leq X < 80) + \\ &+ E(XY | X \geq 80) \cdot P(X \geq 80) = E(0 | X < 25) \cdot P(X < 25) + \\ &+ E(X^2 | 25 \leq X < 80) \cdot P(25 \leq X < 80) + E(100X | X \geq 80) \cdot P(X \geq 80). \end{aligned}$$

From paragraph 3II we know that

$$E(X^2 | 25 \leq X < 80) = \frac{9025}{3} \approx 3008.333,$$

and by linearity of expected value:

$$E(100X \mid X \geq 80) = 100E(X \mid X \geq 80) = 100 \cdot \frac{80 + 100}{2} = 100 \cdot 90 = 9000.$$

Overall

$$E(XY) = 0 \cdot \frac{1}{4} + \frac{9025}{3} \cdot \frac{55}{100} + 9000 \cdot \frac{1}{5} = \frac{41455}{12} \approx 3454.583.$$

Correlation coefficient then:

$$\text{Corr}(X, Y) = \frac{\frac{41455}{12} - 50 \cdot \frac{391}{8}}{\sqrt{\frac{2500}{3} \cdot \frac{243037}{192}}} \approx \boxed{0.984}.$$

Problem 7

Let X and Y be two independent standard normal random variables. Find

1. $P(|X + Y| > |X - Y|)$.
2. $P(|X + Y| > 2|X - Y|)$.

Solution:

1. Using independence:

$$\begin{aligned}
 P(|X + Y| > |X - Y|) &= P((X + Y)^2 > (X - Y)^2) = \\
 &= P(X^2 + 2XY + Y^2 > X^2 - 2XY + Y^2) = P(4XY > 0) = \\
 &= P(X > 0) \cdot P(Y > 0) + P(X < 0) \cdot P(Y < 0) = \\
 &= \frac{1}{2} \cdot \frac{1}{2} + \frac{1}{2} \cdot \frac{1}{2} = \boxed{\frac{1}{2}}.
 \end{aligned}$$

2. There are 2 ways to calculate $P(|X + Y| > 2|X - Y|)$.

(I) Strict calculation

$$\begin{aligned}
 P(|X + Y| > 2|X - Y|) &= P((X + Y)^2 > 4(X - Y)^2) = \\
 &= P((X + Y + 2(X - Y)) \cdot (X + Y - 2(X - Y)) > 0) = \\
 &= P((3X - Y) \cdot (3Y - X) > 0) = \\
 &= P(3X - Y > 0 \cap 3Y - X > 0) + \\
 &\quad + P(3X - Y < 0 \cap 3Y - X < 0).
 \end{aligned}$$

Random variables $3X - Y$ and $3Y - X$ are not independent, so the probabilities of intersections can not be separated. Let's consider a region $(3X - Y) \cdot (3Y - X) > 0$. Its probability is a volume under a joint p.d.f. inside this region. This is illustrated in the fig. 3 and in the fig. 4.

Since X and Y are independent components of the vector $(X \ Y)^T \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$, their joint p.d.f. is radially symmetrical. The volume of a considered region is determined by the angle those lines cover.

The angle of one sector is $\arctan 3 - \arctan \frac{1}{3}$, so the required probability:

$$P(|X + Y| > 2|X - Y|) = \frac{\arctan 3 - \arctan \frac{1}{3}}{\pi} = \boxed{\frac{\arctan \frac{4}{3}}{\pi} \approx 0.295}.$$

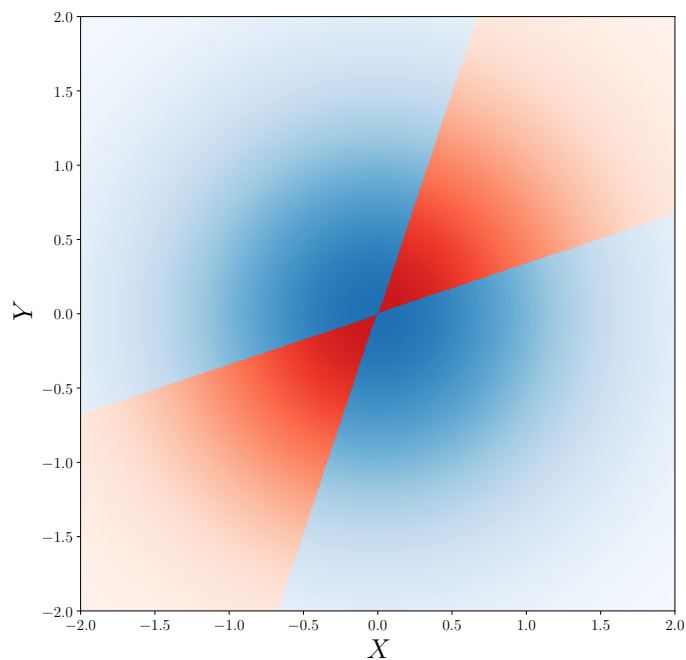


Figure 3: Probability density function of $(X \ Y)^\top \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$ (top view).

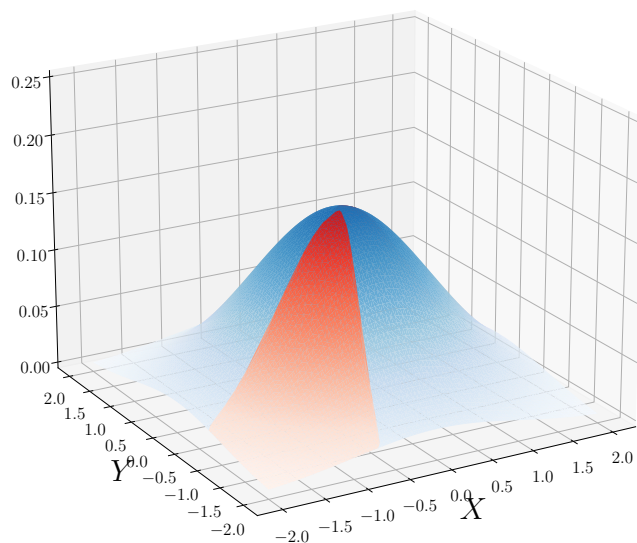


Figure 4: Probability density function of $(X \ Y)^\top \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$ (side view).

(II) Approximation

Let's consider variables $X + Y$ and $X - Y$. They have identical distribution:

$$X \pm Y \sim \mathcal{N}(0, 1^2 + (\pm 1)^2) = \mathcal{N}(0, 2).$$

Inherently, they are uncorrelated:

$$\begin{aligned}\text{Cov}(X + Y, X - Y) &= \text{Cov}(X, X) - \text{Cov}(X, Y) + \text{Cov}(Y, X) - \text{Cov}(Y, Y) = \\ &= \text{V}(X) - \text{V}(Y) = 2 - 2 = 0.\end{aligned}$$

Since X and Y are components of the bivariate normal vector $(X \ Y)^\top \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$, $X + Y$ and $X - Y$ are also components of a vector with bivariate normal distribution. In such case uncorrelatedness means independence – this is the unique property of multivariate normal distributions.

Thus, $X + Y$ and $X - Y$ are independent. The squares of standardized variables are also independent, and moreover, are χ^2 -distributed:

$$Q_1 = \left(\frac{X + Y}{\sqrt{2}} \right)^2 \sim \chi_1^2, \quad Q_2 = \left(\frac{X - Y}{\sqrt{2}} \right)^2 \sim \chi_1^2.$$

The initial probability can be rewritten with F -distribution:

$$\begin{aligned}\text{P}(|X + Y| > 2 |X - Y|) &= \text{P}((X + Y)^2 > 4(X - Y)^2) = \text{P}(Q_1 > 4Q_2) = \\ &= \text{P}\left(\frac{Q_1/1}{Q_2/1} > 4\right) = \text{P}(F_{1,1} > 4) \in \boxed{(0.25, 0.5)}.\end{aligned}$$

The exact result was found in 2I.

Problem 8

Two random variables are given: $X \sim \mathcal{N}(0, 9)$ and $Y \sim \mathcal{N}(0, 4)$. $\text{Corr}(X, Y) = -1$.
Evaluate $P(2X + Y > 3)$.

Solution:

Fig. 5 shows how a bivariate distribution changes with the increase of correlation coefficient between components.

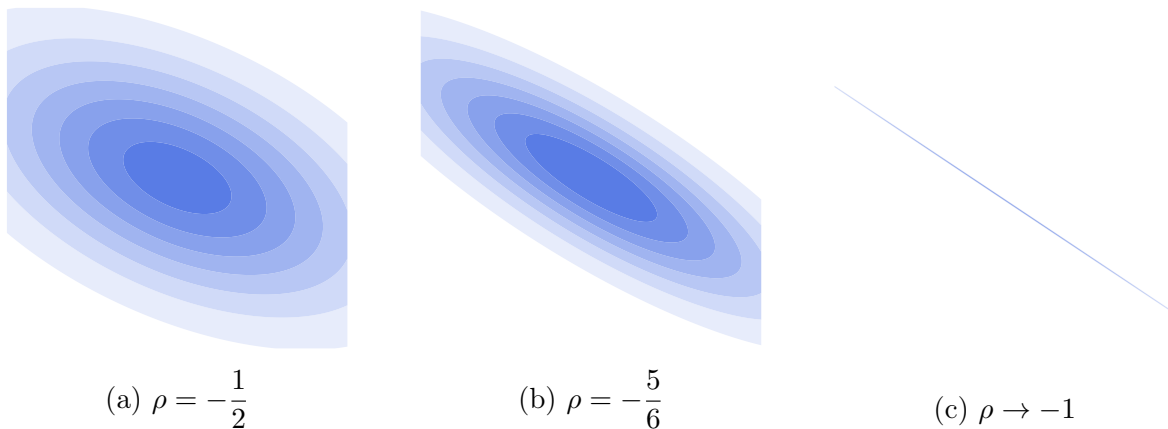


Figure 5: P.d.f. of a bivariate normal distribution with $\sigma_X = 3$ and $\sigma_Y = 2$.

The correlation of -1 means that there is a linear dependency between X and Y , but the change happens in opposite directions.

It means that standardized variables are opposite in sign:

$$\frac{X - \mu_X}{\sigma_X} = -\frac{Y - \mu_Y}{\sigma_Y},$$

$$Y = \mu_Y - \frac{\sigma_Y}{\sigma_X}(X - \mu_X).$$

From problem statement:

$$Y = -\frac{2}{3}X.$$

Thus, the probability:

$$\begin{aligned} P(2X + Y > 3) &= P\left(2X - \frac{2}{3}X > 3\right) = P\left(\frac{4}{3}X > 3\right) = P\left(X > \frac{9}{4}\right) = \\ &= P\left(Z > \frac{9/4 - 0}{3}\right) = 1 - \Phi(0.75) \approx 1 - 0.773 = \boxed{0.227}. \end{aligned}$$

Problem 9

The sample from bivariate normal distribution with random variables X and Y is following:

X	1.59	-2.20	-0.06	-1.45	-1.02	-2.59	-1.14	-3.25
Y	3.24	0.44	-1.14	5.40	2.09	5.33	1.25	8.72

Find 90% confidence interval for a population correlation coefficient ρ .

Solution:

$(1 - \alpha) \cdot 100\%$ confidence interval for ρ :

$$CI_{1-\alpha}(\rho) = \left(\tanh \left(\operatorname{artanh}(\hat{\rho}) - z_{\alpha/2} \cdot \frac{1}{\sqrt{n-3}} \right); \tanh \left(\operatorname{artanh}(\hat{\rho}) + z_{\alpha/2} \cdot \frac{1}{\sqrt{n-3}} \right) \right),$$

Values of sample mean:

$$\bar{x} = \frac{1}{8} \sum_{i=1}^8 x_i \approx -1.265, \quad \bar{y} = \frac{1}{8} \sum_{i=1}^8 y_i \approx 3.166.$$

Values of corrected sums:

$$\begin{aligned} SS_{xx} &= \sum_{i=1}^8 x_i^2 - 8 \cdot (-1.265)^2 \approx 16.283, \\ SS_{yy} &= \sum_{i=1}^8 y_i^2 - 8 \cdot 3.166^2 \approx 71.340, \\ SS_{xy} &= \sum_{i=1}^8 x_i y_i - 8 \cdot (-1.265) \cdot 3.166 \approx -17.240. \end{aligned}$$

Value of sample correlation coefficient:

$$\hat{\rho} = \frac{SS_{xy}}{\sqrt{SS_{xx} \cdot SS_{yy}}} = \frac{-17.240}{\sqrt{16.283 \cdot 71.340}} \approx -0.506.$$

Value of Fisher-transformed correlation coefficient:

$$\operatorname{artanh}(\hat{\rho}) = \frac{1}{2} \ln \left(\frac{1 + \hat{\rho}}{1 - \hat{\rho}} \right) = \frac{1}{2} \ln \left(\frac{1 - 0.506}{1 + 0.506} \right) \approx -0.557.$$

Confidence interval for Fisher-transformed ρ :

$$\begin{aligned} \text{CI}_{90\%}(\text{artanh}(\rho)) &= -0.557 \pm z_{0.05} \cdot \frac{1}{\sqrt{8-3}} = \\ &= -0.557 \pm 1.645 \cdot \frac{1}{\sqrt{8-3}} = \\ &= -0.557 \pm 0.736 = (-1.293; 0.179). \end{aligned}$$

Applying inverse Fisher-transform to the confidence interval above gives required interval:

$$\begin{aligned} \text{CI}_{90\%}(\rho) &= \tanh(\text{CI}_{90\%}(\text{artanh}(\rho))) = \frac{e^{2\text{CI}_{90\%}(\text{artanh}(\rho))} - 1}{e^{2\text{CI}_{90\%}(\text{artanh}(\rho))} + 1} = \\ &= \left(\frac{e^{2 \cdot (-1.293)} - 1}{e^{2 \cdot (-1.293)} + 1}, \frac{e^{2 \cdot 0.179} - 1}{e^{2 \cdot 0.179} + 1} \right) = \boxed{(-0.860; 0.177)}. \end{aligned}$$

Problem 10

Consider observations in the table below:

x	0	2	6	-3	4	1	-2	5	-1
y	8	2	0	6	1	5	7	3	4

1. Find Spearman's rank correlation coefficient r_s .
2. Find sample correlation coefficient r and compare it with r_s .

Solution:

1. Let's rank our sample and calculate differences d :

x	0	2	6	-3	4	1	-2	5	-1
y	8	2	0	6	1	5	7	3	4
rank(x)	4	6	9	1	7	5	2	8	3
rank(y)	9	3	1	7	2	6	8	4	5
d	-5	3	8	-6	5	-1	-6	4	-2
d^2	25	9	64	36	25	1	36	16	4

Spearman's rank correlation coefficient:

$$r_s = 1 - \frac{6 \sum_{i=1}^n d_i^2}{n(n^2 - 1)} = 1 - \frac{6 \cdot 216}{9 \cdot (81 - 1)} = \boxed{-0.8}.$$

2. Values of sample mean:

$$\bar{x} = \frac{1}{9} \sum_{i=1}^9 x_i = \frac{4}{3}, \quad \bar{y} = \frac{1}{9} \sum_{i=1}^9 y_i = 4.$$

Values of corrected sums:

$$SS_{xx} = \sum_{i=1}^9 x_i^2 - 9 \cdot \left(\frac{4}{3}\right)^2 = 80,$$

$$SS_{yy} = \sum_{i=1}^9 y_i^2 - 9 \cdot 4^2 = 60,$$

$$SS_{xy} = \sum_{i=1}^9 x_i y_i - 9 \cdot \frac{4}{3} \cdot 4 = -56.$$

Value of sample correlation coefficient:

$$r = \frac{SS_{xy}}{\sqrt{SS_{xx} \cdot SS_{yy}}} = \frac{-56}{\sqrt{80 \cdot 60}} \approx \boxed{-0.808}.$$

Clearly, correlation coefficients are close due to the absence of prominent outliers:

$$\boxed{r \approx r_s}.$$