Physiotherapy Statistics Exams

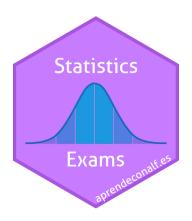




Table of contents

Pr	eface	3
1	2018-05-31 Descriptive Statistics and Regression	4
2	2022-05-06 Probability and Random Variables	8
3	2022-06-06 Descriptive Statistics and Regression	11
4	2022-06-06 Probability and Random Variables	14
5	2023-03-23 Descriptive Statistics and Regression	16
6	2023-04-27 Probability and Random Variables	20
7	2023-05-30 Descriptive Statistics and Regression	23
8	2023-05-30 Probability and Random Variables	27

Preface

Statistics exam collection of the Physiotherapy degree.

1 2018-05-31 Descriptive Statistics and Regression

Exercise 1.1.

The ages of a sample of patients of a physical therapy clinic are:

- a. Compute the quartiles.
- b. Draw the box plot and identify outliers (do not group data into intervals).
- c. Split the sample into two groups, patients younger and older than 65. In which group is the mean more representative. Justify the answer.
- d. Which distribution is less symmetric, the one of patients younger than 65 or the one of patients older?
- e. Which age is relatively smaller with respect to its group, 50 years in the group of patients younger than 65 or 72 years in the group of patients older than 65?

Use the following sums for the computations.

Younger than 65: $\sum x_i = 953$ years, $\sum x_i^2 = 52475$ years², $\sum (x_i - \bar{x})^3 = -30846.51$ years³ and $\sum (x_i - \bar{x})^4 = 939658.83$ years⁴. Older than 65: $\sum x_i = 588$ years, $\sum x_i^2 = 43530$ years², $\sum (x_i - \bar{x})^3 = 1485$ years³ and $\sum (x_i - \bar{x})^4 = 26983.5$ years⁴.

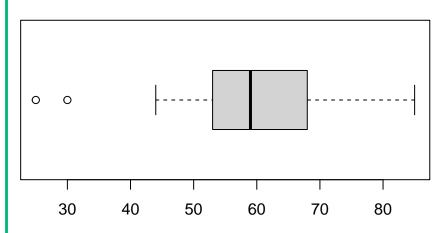
Solution

- a. $Q_1=53$ years, $Q_2=59$ years and $Q_3=68$ years. b. There are 2 outliers: 25, 30.

	30											1
59	61	63	63	63	66	68	70	71	72	74	82	85

Warm-up time	15	35	22	28	21	18	25	30	23	20
Injuries	42	2	16	6	17	29	10	3	12	20

Boxplot of patients ages



a. Let x be the age in patients younger than 65 and y the age in patients older than 65.

$$\begin{split} \bar{x} &= 52.9444 \text{ years}, \, s_x^2 = 112.1636 \text{ years}^2, \, s_x = 10.5907 \text{ years and } cv_x = 0.2. \\ \bar{y} &= 73.5 \text{ years}, \, s_y^2 = 39 \text{ years}^2, \, s_y = 6.245 \text{ years and } cv_y = 0.085. \end{split}$$

The mean is more representative in patients older than 65 since the coefficient of variation is smaller.

- b. $g_{1x}=-1.4426$ and $g_{1y}=0.7621$, thus the distribution of ages of people younger than 65 is less symmetric.
- c. The standard scores are $z_x(50) = -0.278$ and $z_y(72) = -0.2402$, thus 50 years is relative smaller in the group of people younger than 65.

Exercise 1.2.

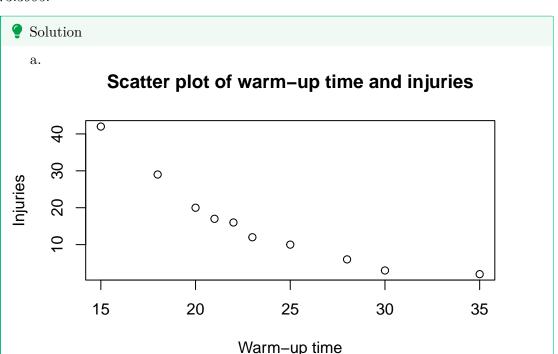
The table below shows the number of injuries of several teams during a league and the average warm-up time of its players.

- a. Draw the scatter plot.
- b. Which regression model is more suitable to predict the number of injuries as a function of the warm-up time, the logarithmic or the exponential? Use that regression model to predict the expected number of injuries for a team whose players warm-up 20 minutes a day.
- c. Which regression model is more suitable to predict the warm-up time as a function of the number of injuries, the logarithmic or the exponential? Use that regression

model to predict the warm-up time required to have no more than 10 injuries in a league.

d. Are these predictions reliable? Which one is more reliable?

Use the following sums for the computations (X warm-up time and Y number of injuries): $\begin{array}{l} \sum x_i = 237, \ \sum \log(x_i) = 31.3728, \ \sum y_j = 157, \ \sum \log(y_j) = 24.0775, \\ \sum x_i^2 = 5937, \ \sum \log(x_i)^2 = 98.9906, \ \sum y_j^2 = 3843, \ \sum \log(y_j)^2 = 66.3721, \\ \sum x_i y_j = 3115, \ \sum x_i \log(y_j) = 519.1907, \ \sum \log(x_i) y_j = 465.8093, \ \sum \log(x_i) \log(y_j) = 31.3728, \ \sum \log(x_i) \log(y_j) = 31.3728, \ \sum \log(x_i) \log(x_i) \log(y_j) = 31.3728, \ \sum \log(x_i) \log(x$ 73.3995.



a. $\bar{x} = 23.7 \text{ min}, s_x^2 = 32.01 \text{ min}^2.$ $\frac{1}{\log(x)} = 3.1373 \log(\min), \ s_{\log(x)}^2 = 0.0565 \log(\min)^2.$

 $\underline{\bar{y} = 15.7}$ injuries, $s_y^2 = 137.81$ injuries².

 $\begin{array}{l} \overline{\log(y)} = 2.4078 \, \log(\text{injuries}), \, s_{\log(y)}^2 = 0.8399 \, \log(\text{injuries})^2. \\ s_{x \log(y)} = -5.1446, \, s_{\log(x)y} = -2.6744 \end{array}$

Exponential determination coefficient: $r^2 = 0.9844$

Logarithmic determination coefficient: $r^2 = 0.9185$

So the exponential regression model es better to predict the number of injuries as a function of the warm-up time.

Exponential regression model: $y = e^{6.2168 + -0.1607x}$.

Prediction: y(20) = 20.1341 injuries.

b. The logarithmic model is better to predict the warm-up time as a function of the number or injuries.

Logarithmic regression model: $x=164.1851+-47.3292\log(y)$. Prediction: x(10)=55.2056112 min.

c. Both predictions are very reliable since de deternation coefficient is very high but the last one is a little less reliable as it is for a value further from the data range.

2 2022-05-06 Probability and Random Variables

Exercise 2.1.

8% of people in a population consume cocaine. It is also known that 4%: of people who consume cocaine have a heart attack and 10%: of people who have a heart attack consume cocaine.

- a. Construct the probability tree for the random experiment of drawing a random person from the population and measuring if he or she consumes cocaine and if he or she has a heart attack.
- b. Compute the probability that a random person of the population does not consume cocaine and does not have a heart attack.
- c. Are the events of consuming cocaine and having a heart attack dependent?
- d. Compute the relative risk and the odds ratio of suffering a heart attack consuming cocaine. Which association measure is more suitable for this study? Interpret it.

Solution

Let C the event of consuming cocaine and H the event of having a heart attack.

a.

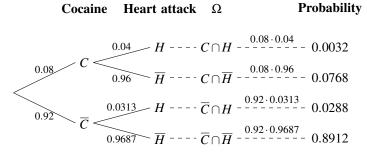


Figure 2.1: Probability tree

a.
$$P(\overline{C} \cap \overline{H}) = 0.8912$$
.

- b. The events are dependent as $P(C) = 0.08 \neq P(C|H) = 0.1$.
- c. RR(H) = 1.2778 and OR(H) = 1.2894. The odds ratio is more suitable as the study is retrospective. That means that the odds of having a heart attack is 1.2894 times greater if a person consumes cocaine.

Exercise 2.2.

A basketball player scores 12 points per game on average.

- a. What is the probability that the player scores more than 4 points in a quarter?
- b. If the player plays 10 games in a league, what is the probability of scoring less than 6 points in some game?

Solution

- a. Let X be the points scored in a quarter by the player. Then $X \sim P(3)$, and P(X > 4) = 0.1847.
- b. Let Y be the number of points scored in a game by the player. Then $Y \sim P(12)$ and P(Y < 6) = 0.0203.
 - Let Z be the number of games with less than 6 points scored by the player. Then $Z \sim B(10, 0.0203)$, and P(Z > 0) = 0.1858.

Exercise 2.3.

The creatine phosphokinase (CPK3) is an enzyme in the body that causes the phosphorylation of creatine. This enzyme is found in the skeletal muscle and can be measured in a blood analysis. The concentration of CPK3 in blood is normally distributed, and the interval centered at the mean with the reference values, that accumulates 99%: of the population, ranges from 40 to 308 IU/L in healthy adult males.

- a. Compute the mean and the standard deviation of the concentration of CPK3 in healthy males.
 - Note: If you are not able to compute the standard deviation, use $\sigma = 50$ UI/L for the next parts.
- b. A diagnostic test to detect muscular dystrophy gives a negative outcome when the concentration of CPK3 is below 300 UI/L. Compute the specificity of the test.
- c. If the concentration of CPK3 in people with muscular dystrophy also follows a normal distribution with mean $350~{\rm IU/L}$ and the same standard deviation, what is the sensitivity of the test?

d. Compute the predictive values of the test and interpret them assuming that the muscular dystrophy prevalence is 8%.

Solution

- a. $\mu = 174 \; \mathrm{IU/L}$ and $\sigma = 51.938 \; \mathrm{IU/L}$.
- b. Specificity = 0.9924.
- c. Sensitivity = 0.8321. The test is better to confirm the disease as the specificity is greater than the sensitivity.
- d. PPV = 0.9046. Thus, we can diagnose the disease with a positive outcome. NPV = 0.9855. Thus, we can rule out the disease with a negative outcome.

3 2022-06-06 Descriptive Statistics and Regression

Exercise 3.1.

The patients of a physiotherapy clinic were asked to assess their satisfaction in a scale from 0 to 10. The assessments are summarized in the table below.

Assessment	Patients
0 - 2	3
2 - 4	12
4 - 6	9
6 - 8	18
8 - 10	22

- a. Compute the interquartile range of the assessment and interpret it.
- b. If it is required an assessment greater than 5 in more than 50%: of patients for the clinic to remain open, will the clinic remain open?
- c. Is the assessment mean representative?
- d. Compute the coefficient of kurtosis of the assessment and interpret it. Is the kurtosis normal?
- e. If the assessment mean of another clinic is 6.8 and the standard deviation is 2.6, which assessment is relatively higher 6 in the first clinic or 6.2 in the second?

Use the following sums for the computations:

$$\sum x_i n_i = 408, \ \sum x_i^2 n_i = 3000, \ \sum (x_i - \bar{x})^3 n_i = -548.25 \text{ and } \sum (x_i - \bar{x})^4 n_i = 5140.45.$$

Solution

Let X be the patient assessment. a. $Q_1=4.4444,\,Q_3=9.0907$ and IQR=4.6463, so the central dispersion is moderate.

- a. F(5) = 0.2695, and the percentage of patients with an assessment greater than 5 is 73.05.
- b. $\bar{x}=6.375,\,s_x^2=6.2344,\,s_x=2.4969$ and cv=0.3917, thus the representative

Week	1	3	6	9	14	17	21	24
Grip strength	15	22	29	34	36	39	40	41

ity of the mean is moderate.

c. $g_2 = -0.9335$ and the distribution is flatter than a Gauss bell, but normal, as g_2 is between -2 and 2.

d. First clinic: z(6) = -0.1502

Second clinic: z(6.2) = -0.3077.

Thus, an assessment of 6 in the first clinic is relatively higher as its standard score is greater.

Exercise 3.2.

A study tries to determine the effectiveness a training program to increase the grip strength. The table below shows the grip strength in Kg in some weeks of the training program.

- a. Compute the regression coefficient of the grip strength on the weeks and interpret it.
- b. According to the logarithmic regression model, what is the expected grip strength after 5 and 25 weeks. Are these predictions reliable? Would these predictions be more reliable with the linear regression model?
- c. According to the exponential regression model, how many weeks are required to have a grip strength of 25 Kg?
- d. What percentage of the total variability of the weeks is explained by the exponential model?

Use the following sums (X=Weeks and Y=Grip strength):

$$\sum x_i = 95, \ \sum \log(x_i) = 16.7824, \ \sum y_j = 256, \ \sum \log(y_j) = 27.3423, \ \sum x_i^2 = 1629, \ \sum \log(x_i)^2 = 43.606, \ \sum y_j^2 = 8804, \ \sum \log(y_j)^2 = 94.3237, \ \sum x_i y_j = 3552, \ \sum x_i \log(y_j) = 342.9642, \ \sum \log(x_i) y_j = 608.4186, \ \sum \log(x_i) \log(y_j) = 60.047.$$

Solution

a. $\overline{x} = 11.875$ weeks, $s_x^2 = 62.6094$ weeks². $\overline{y} = 32$ Kg, $s_y^2 = 76.5$ Kg².

 $s_{xy} = 64 \text{ weeks-Kg.}$

Regression coefficient of Y on X: $b_{yx} = 1.0222 \text{ Kg/week}$. The grip strength increases 1.0222 Kg per week.

b. $\overline{\ln(x)}=2.0978$ ln (weeks), $s_{\ln(x)}^2=1.05$ ln (weeks)^2 and $s_{\ln(x)y}=8.9226$ ln (weeks)Kg.

Logarithmic regression model of Y on X: $y = 14.1729 + 8.498 \ln(x)$.

Predictions: y(5) = 27.8499 Kg and y(25) = 41.5268 Kg.

Logarithmic coefficient of determination: $r^2 = 0.9912$. The predictions are not reliable because the sample size is small.

Linear coefficient of determination: $r^2 = 0.8552$.

As the linear coefficient of determination is less than the logarithmic one, the predictions with the logarithmic model are more reliable.

- c. Exponential regression model of X on Y: $x=e^{-1.6345+0.1166y}$. Prediction: x(25)=3.6015 Weeks.
- d. As $r^2 = 0.9912$, the exponential models explains 99.12%: of the variability of the weeks.

4 2022-06-06 Probability and Random Variables

Exercise 4.1.

A diagnostic test for a cervical injury has a 99% of sensitivity and produces 80% of right diagnosis. Assuming that the prevalence of the injury is 10%

- a. Compute the specificity of the test.
- b. Can we rule out the injury with a negative outcome of the test?
- c. Can we diagnose the injury with a positive outcome of the test? What must the minimum prevalence of the injury be to diagnose the injury with a positive outcome of the test?

Solution

- a. Specificity = $P(-|\overline{D}) = 0.7789$.
- b. Negative predictive value = $P(\overline{D}|-) = 0.9986 > 0.5$, so we can rule out the injury with a negative outcome.
- c. Positive predictive value = P(D|+) = 0.3322 < 0.5, so we can not diagnose the injury with a positive outcome. The minimum prevalence required to be able to diagnose the injury with a positive outcome is P(D) = 0.1825.

Exercise 4.2.

A pharmacy sells two vaccines A and B against a virus. The A vaccine produces 5%: of side effects, while the B vaccine produces 2%: of side effects. The pharmacy has sold 10 units of the A vaccine and 100 units of the B vaccine.

- a. Compute the probability of having less than 2 side effects with the A vaccine.
- b. Compute the probability of having more than 3 side effects with the B vaccine.
- c. If we apply both vaccines to the same person at different moments, and assuming that the production of side effects of the vaccines are independent, what is the probability that this person will have any side effect?

Solution

- a. Let X be the number of side effects in 10 applications of A vaccine. Then, $X \sim B(10, 0.05)$ and P(X < 2) = 0.9139.
- b. Let Y be the number of side effects in 100 applications of B vaccine. Then, $Y \sim B(100, 0.02) \approx P(2)$ and P(Y > 3) = 0.1429.
- c. Let A and B the events of having side effects with vaccines A and B respectively. $P(A \cup B) = 0.069$.

Exercise 4.3.

The length of the femur bone is normally distributed in both men and women with a standard deviation of 4 cm. It is also known that the first quartile in women is 42.3 cm, while the third quartile in men is 50.7 cm.

- a. What is the difference between the means of the femur length of women and men? Remark: If you do not know how to compute the means, use a mean 44 cm for women and a mean 47 cm for men in the following parts.
- b. Compute the 60th percentile of the femur length in women. What percentage of men have a femur length less than the 60th percentile of women?
- c. If we pick a woman and man at random, what is the probability that neither of them has a femur length less than 45 cm?

Solution

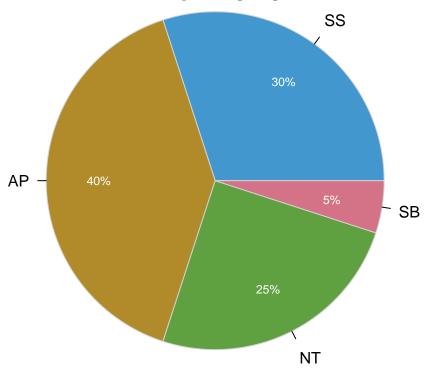
Let X and Y be the femur length of women and men respectively. Then $X \sim N(\mu_x, 4)$ and $Y \sim N(\mu_y, 4)$.

- a. $\mu_x = 44.91 \text{ cm} \text{ and } \mu_y = 48.02 \text{ cm}.$
- b. 60th percentile in women $P_{60}=45.9234$ cm, and P(Y<45.9234)=0.3001, that is, a 30.01 of men have a femur length less than the 60th percentile of women.
- c. $P(X \ge 45 \cap Y \ge 45) = 0.3805$.

5 2023-03-23 Descriptive Statistics and Regression

Exercise 5.1.

The chart below shows the percentage of grades in a Statistic course with 60 students.



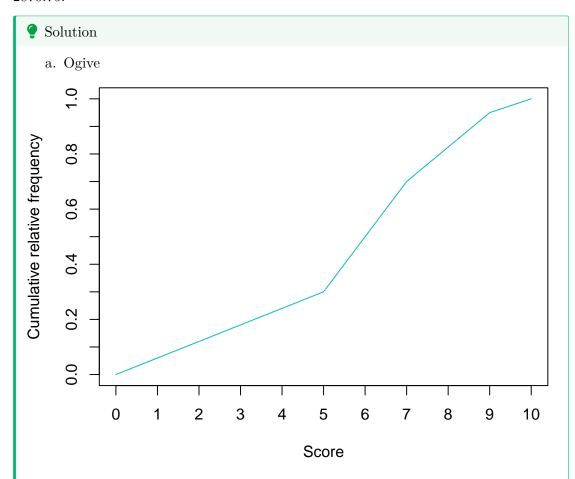
a. Plot the ogive of the score, assuming the following correspondence between grades and scores

Grade	Score
SS	[0, 5)
AP	[5, 7)
NT	[7, 9)
SB	[9, 10]

a. Compute the median and interpret it.

- b. How many students got a score greater than 8?
- c. Study the dispersion of the distribution.
- d. Study the skewness of the distribution. Is it normal?
- e. If we apply the transformation y = 10x + 5 to the scores, how changes the representativeness of the mean. And the skewness?

Use the following sums for the computations (X = Score): $\sum x_i n_i = 337.5$, $\sum x_i^2 n_i = 2207.25$, $\sum (x_i - \bar{x})^3 n_i = -172.55$ and $\sum (x_i - \bar{x})^4 n_i = 2870.75$.



- b. Me = 6 points.
- c. N(8) = 49.5 students.
- d. $\bar{x}=5.625$ points, $s_x^2=5.1469$ points², $s_x=2.2687$ points and $cv_x=0.4033$. Thus, there is a moderate dispersion with respect to the mean.
- e. $g_1 = -0.2463$ and therefore the distribution is a little bit left skewed.

f. $\bar{y} = 61.25$ points, $s_y^2 = 514.6875$ points², $s_y = 22.6867$ points and $cv_y = 0.3704$. As $cv_y < cv_x$ the representativeness of the mean increases. As the slope of the linear transformation is positive, the skewness does not change.

Exercise 5.2.

A study tries to determine if there is a relation between the gestation time (in weeks) and the age of the mother (in years). A sample of 40 mothers was taken and the sums below summarize the results (X=Age and Y=Gestation time):

```
\begin{array}{l} \sum x_i = 1262 \text{ years, } \sum \log(x_i) = 137.0078 \log(\text{years}), \ \sum y_j = 1583.6 \text{ weeks, } \sum \log(y_j) = 147.1305 \log(\text{weeks}), \\ \sum x_i^2 = 41862 \text{ years}^2, \ \sum \log(x_i)^2 = 471.4222 \log(\text{years})^2, \ \sum y_j^2 = 62734.685 \text{ weeks}^2, \\ \sum \log(y_j)^2 = 541.2096 \log(\text{weeks})^2, \\ \sum x_i y_j = 50116.7 \text{ years-weeks, } \sum x_i \log(y_j) = 4645.8 \text{ years-log(weeks), } \sum \log(x_i) y_j = 5428.9192 \log(\text{years}) \cdot \text{weeks, } \sum \log(x_i) \log(y_j) = 504.0696 \log(\text{years}) \cdot \log(\text{weeks}). \end{array}
```

- a. Which regression models, linear, exponential or logarithmic, explains better the relation between the age and the gestation time?
- b. Use the best model to predict the gestation time for a mother 45 years old. Is this prediction reliable?
- c. According to the linear model, how much increases or decreases the gestation time for every year of the mother?

Solution

a. Linear model: $\overline{x}=31.55$ years, $s_x^2=51.1475$ years². $\overline{y}=39.59$ weeks, $s_y^2=0.999$ weeks². $s_{xy}=3.853$ years-weeks. $r^2=0.2905$.

Exponential model: $\overline{\ln(y)}=3.6783$ ln(weeks), $s_{\ln(y)}^2=0.0006$ ln(weeks)² $s_{x\ln(y)}=0.0958$ years·ln(weeks). $r^2=0.2882$.

Logarithmic model: $\overline{\ln(x)}=3.4252$ ln(years), $s_{\ln(x)}^2=0.0536$ ln(years) $s_{\ln(x)y}=0.1195$ ln(years)weeks. $r^2=0.2668$.

As the linear coefficient of determination is greater, the linear model explains better the relation between de gestation time and the age of the mother.

b. Linear regression model of Y on X: y = 37.2133 + 0.0753x. Predictions: y(45) = 40.6032 weeks.

The predictions are not reliable because the coefficient of determination is pretty low.

c. Regression coefficient of Y on X: $b_{yx}=0.0753$ weeks/year. The gestation time increases 0.0753 weeks per year.

6 2023-04-27 Probability and Random Variables

Exercise 6.1.

A water source contaminated contains 0.1 amoebas per litre on average.

- a. What is the probability that 2 litres of water from this source contains more than one amoeba?
- b. If 5 persons drink 2 litres of water from this source, what is the probability of having some person infected with amoebas?
- c. If 100 persons drink half a litre of water from this source, what is the probability that less than 5 are infected with amoebas?

Solution

- a. Let X be the number of a moebas in 2 litres of contaminated water. Then $X \sim P(0.2)$ and P(X>1)=0.0175.
- b. The probability that a person who drank 2 litres of contaminated water is infected is $P(X \ge 1) = 0.1813$. Let Y be the number of persons infected with amoebas in a sample of 5 persons who drank 2 litres of contaminated water. Then $Y \sim B(5, 0.1813)$ and $P(Y \ge 1) = 0.6321$.
- c. Let U be the number of amoebas in half a litre of contaminated water. Then $U \sim P(0.05)$ and $P(U \geq 1) = 0.0488$. Let V be the number of persons infected with amoebas in a sample of 100 persons who drank half a litre of contaminated water. Then $V \sim B(100, 0.0488) \approx P(4.8771)$ and P(V < 5) = 0.4623.

Exercise 6.2.

Respiratory allergies affect 1 out of every 15 individuals in a population, while food intolerances affect 5% of individuals. Assuming that the two problems are independent,

- a. Compute the probability of having at least one of the problems.
- b. Compute the probability of having an allergy but not an intolerance.

- c. Compute the probability of having neither of the two problems.
- d. Compute the probability of having an allergy if you have an intolerance.

Solution

Let A the event of having respiratory allergies and B the event of having food intolerance.

- a. $P(A \cup B) = 0.1133$.
- b. P(A B) = 0.0633.
- c. $P(\overline{A} \cap \overline{B}) = 0.8867$.
- d. P(A|B) = 0.0667.

Exercise 6.3.

In a population of 20000 women, it is known that back width follows a normal distribution with mean 29 cm and standard deviation 2.4 cm.

- a. Compute the number of women with a back width greater than 32 cm.
- b. Compute the interquartile range of women's back width and interpret it.
- c. Compute the probability that a woman with a back width above the third quartile, has a back width above 32.

Solution

Let X be the back width, then $X \sim N(29, 2.4)$.

- a. P(X>32)=0.1056 and approximately 2113 persons have a back width greater than 32 cm.
- b. $Q_1 = 27.3812$ cm, $Q_3 = 30.6188$ cm, and IQR = 3.2376 cm.
- c. P(X > 32|X > 30.6188) = 0.4226.

Exercise 6.4.

A diagnostic test for prostate cancer has a specificity of 80% and produces 1.6% of false negatives. It is known that the prevalence of prostate cancer in a population is 2%.

a. Compute the sensitivity of the test. Does the outcome of the test depend on whether a man has prostate cancer?

- b. Is this a good test to diagnose the disease?
- c. What should be the minimum specificity of the test to diagnose the disease with a positive outcome?

Solution

- a. Sensitivity = P(+|D) = 0.2. The outcome of the test does not depend on the prostate cancer.
- b. Positive predictive value = P(D|+) = 0.02 < 0.5, so we can not confirm the prostate cancer with a positive outcome.
- c. Minimum specificity 0.9959.

7 2023-05-30 Descriptive Statistics and Regression

Exercise 7.1.

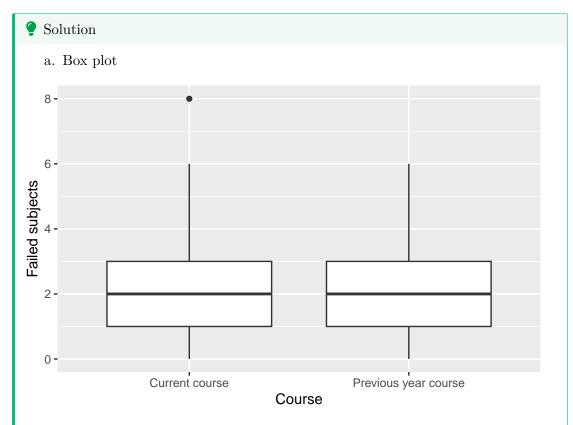
To see if the confinement due to COVID-19 influenced the performance of a course, the number of failed subjects of each student in the current course and in the previous year course has been counted, obtaining the table below.

- a. Draw the box plots of the failed subjects in the current and the previous year courses and compare them.
- b. Can we assume that both samples come from a normal population?
- c. In which sample is the mean more representative?
- d. Which number of failed subjects is relatively greater, 7 in the current course or 6 in the previous year course?

Use the following sums for the computations:

Previous year course: $\sum x_i = 84$ subjects, $\sum x_i^2 = 254$ subjects², $\sum (x_i - \bar{x})^3 = 122.99$ subjects³ y $\sum (x_i - \bar{x})^4 = 669.21$ subjects⁴. Current course: $\sum x_i = 91$ subjects, $\sum x_i^2 = 341$ subjects², $\sum (x_i - \bar{x})^3 = 301.16$ subjects³ y $\sum (x_i - \bar{x})^4 = 2012.88$ subjects⁴.

Failed subjects	Previous year course	Current course
0	7	8
1	15	12
2	11	8
3	5	7
4	4	3
5	2	2
6	1	2
8	0	1



a. Previous year course: $\bar{x}=1.8667$ subjects, $s_x^2=2.16$ subjects², $s_x=1.4697$ subjects, $g_1=0.8609$ and $g_2=0.1874$.

Current course: $\bar{x} = 2.1163$ subjects, $s_x^2 = 3.4516$ subjects², $s_x = 1.8578$ subjects, $g_1 = 1.0922$ and $g_2 = 0.9292$.

In both courses the coefficients of skewness and kurtosis are between -2 and 2, so we can assume that both samples come from a normal population.

b. Previous year course: cv = 0.7873.

Current year course: cv = 0.8779.

As the coefficient of variation of the previous year course is smaller, its mean is a little bit more representative.

c. Previous year course: z(6) = 2.8124.

Current course: z(7) = 2.6287.

Thus, 6 failed subjects is relative greater in the previous year course than 7 failed subjects in the current course.

Exercise 7.2.

The following table shows the reduction of inflammation in trauma (in percentage) for different doses of dexketoprofen given for 4 days (in mg).

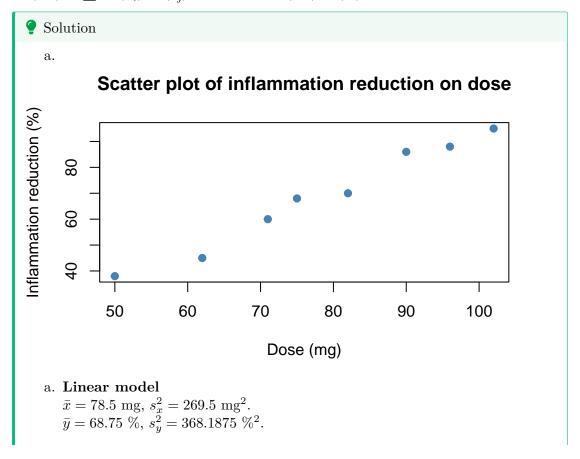
Dose (mg)	50	62	71	75	82	90	96	102
Inflammation reduction (%)	38	45	60	68	70	86	88	95

- a. Draw the scatter diagram of inflammation reduction versus dose of dexketoprofen.
- b. What percentage of inflammation reduction variability does the linear model explain? And the logarithmic model?
- c. According to the best of the two previous models, what is the expected percentage of inflammation reduction if we use 75 mg of dexketoprofen? Which dose should be administered to attain an inflammation reduction of 90%? Are these predictions reliable?

Use the following sums for the computation (X=Dose, Y=inflammation reduction): $\sum x_i = 628$ mg, $\sum \log(x_i) = 34.7152$ $\log(\text{mg})$, $\sum y_j = 550$ %, $\sum \log(y_j) = 33.4922$

 $\log(\%)$,

 $\sum x_i^2 = 51454 \text{ mg}^2, \; \sum \log(x_i)^2 = 151.0394 \; \log(\text{mg})^2, \; \sum y_j^2 = 40758 \; \%^2, \; \sum \log(y_j)^2 = 140.9659 \; \log(\%)^2,$



 $s_{xy}=311.625~{\rm mg}\cdot$ %. $r^2=0.9787,$ so the linear model explains 97.87% of the variability of the inflammation reduction.

Logarithmic model

 $\overline{\log(x)} = 4.3394 \, \log(\text{mg}), \, s_{\log(x)}^2 = 0.0496 \, \log(\text{mg})^2.$

 $s_{\log(x)y} = 4.1936 \log(\text{mg}) \cdot \%.$

 $r^2 = 0.9637$, so the logarithmic model explains 96.37% of the variability of the inflammation reduction.

b. The linear model fits better than the logarithmic one as its coefficient of determination is greater.

Regression line of Y on X: y = -22.0202 + 1.1563x.

Prediction: y(75) = 64.7029 %.

Regression line of X on Y: x = 20.3117 + 0.8464y.

Prediction: x(90) = 96.4855 mg.

Although the linear coefficient of determination is close to 1, the sample size is too small to consider these predictions reliable.

8 2023-05-30 Probability and Random Variables

Exercise 8.1.

In a city of 3.5 million inhabitants there are three urban transport systems: metro, bus and tram. In general, on a working day, the number of travelers is 1.500.000 for the metro, 750.000 for the bus and 450.000 for the tram. Moreover, it is known that 30% of metro travelers also use the bus, 10% of metro travelers also use the tram and 5% of metro travelers also use both bus and tram. Finally, 15% of bus travelers also use the tram. (Hint: An inhabitant can take or not the urban transport).

- a. Calculate the probability that, on a working day, an inhabitant uses only one of three transport systems.
- b. Calculate the probability that, in a working day, an inhabitant uses at least one transport system.
- c. When only a transport system is used, there is a 2% probability of having a delay of more than 5 minutes on a working day. However, the probability of having such a delay rises to 7% when combining more than one transport system in a working day. Calculate the probability that an inhabitant suffers a delay of more than 5 minutes on a working day.
- d. With the same information as in part (c) and knowing that a traveler suffered a delay of more than 5 minutes, calculate the probability that this traveler took more than one transport system.

Solution

Let M, B and T the events corresponding to picking the metro, bus and tram respectively.

- a. $P(\text{Only one transport}) = P(M \cap \overline{B} \cap \overline{T}) + P(\overline{M} \cap B \cap \overline{T}) + P(\overline{M} \cap \overline{B} \cap T) = 0.5185.$
- b. $P(\text{At least one transport}) = P(M \cup B \cup T) = 0.6793.$
- c. Let D the event of suffering a delay of more than 5 minutes. Then, P(D) = 0.0216.

d. P(More than one transport|D) = 0.5211.

Exercise 8.2.

To study the effectiveness of an exercise program in front of the computer to delay the onset of presbyopia, a sample of 1000 40-year-old people without presbyopia was taken and divided into two groups of equal size. One group followed the exercise program and the other did not. After 10 years there were 45 people who developed presbyopia in the group that followed the exercise program and 125 in the group that did not follow the program.

- a. Calculate the risk of developing presbyopia without following the exercise program.
- b. Calculate the relative risk of developing presbyopia if the exercise program is followed and interpret it.
- c. Calculate the odds ratio of developing presbyopia if the exercise program is followed and interpret it.
- d. What type of association measure is most appropriate in this study?
- e. If a sample of 20 40-year-old people without presbyopia is taken, what is the probability that at least 2 will develop presbyopia if they do not follow the exercise program?
- f. If a sample of 100 40-year-old people without presbyopia is taken, what is the probability that fewer than 2 will develop presbyopia if they follow the exercise program?

Solution

- a. Let D the event of suffering Presbyopia. Then $R_C(D) = 0.25$.
- b. RR(D) = 0.36. The risk of suffering Presbyopia following the exercise program is almost one third of the risk not following the exercise program.
- c. OR(D) = 0.2967. The odds of suffering Presbyopia following the exercise program is almost one fourth of the odds not following the exercise program.
- d. Let X be the number of persons not following the program of exercises that develops presbyopia in a sample of 20. Then $X \sim B(20, 0.25)$ and $P(X \ge 2) = 0.9757$.
- e. Let Y be the number of persons following the program of exercises that develops presbyopia in a sample of 100. Then $X \sim B(100,0.09) \approx P(9)$ and P(X < 2) = 0.0012.

Exercise 8.3.

For weightlifters, the recovery time after a lumbar vertebrae injury follows a normal distribution with a mean of 60 days and a standard deviation of 8 days. Knowing this, answer the following questions:

- a. If we take a sample consisting of 200 weightlifters, how many of them will recover from the injury in 52 days or less?
- b. How long should we expect to take a weightlifter in percentile 33 of our distribution to recover?
- c. If the Olympic Games take place in 84 days, and a competitor has just had his lumbar vertebrae injured, what is the probability for him not to take part in the competition?
- d. In 54 days, the qualifying tournament for the Weightlifting World Championship will take place, and 6 competitors have suffered lumbar vertebrae injury today. What is the probability that at least one of them will recover on time to compete in that tournament?

Solution

Let X be the recovery time after a lumbar vertebrae injury, then $X \sim N(60, 8)$.

- a. P(X < 52) = 0.1587 and approximately 32 weight lifters will recover in 52 days or less.
- b. $P_{33} = 56.4807$ days.
- c. P(X > 84) = 0.0013.
- d. Let Y the number of weight lifters that recover in less than 54 days in a sample of 6 weight lifters just injured. Then $Y \sim B(6, 0.2266)$ and $P(Y \ge 1) = 0.786$.