



Stockholm  
University

**BASIC STATISTICS FOR ECONOMISTS, STE101. REEXAM SOLUTIONS**

Department of statistics

Edgar Bueno

2023–08–15

**Part one. Multiple choice**

- Which of the following is **correct** as an interpretation of the mean,  $\bar{x}$ :
  - It is the most likely value;
  - It is the most frequently occurring value;
  - It is the center of gravity of the observations;
  - It is the middle observation;
  - None of the above.
- A researcher is studying the relation between the continent in which a country is located (Africa, Asia, Europe and Latin America) and its Human Development Index —HDI— (Very high, High, Medium, Low). Using a statistical software, the researcher carries out a test of independence between both variables and obtains a  $p$ -value smaller than 0.00001. Considering a significance level  $\alpha = 0.01$ , which of the following is **correct** regarding the test that has been implemented:
  - the continent and the HDI are independent;
  - the continent and the HDI are dependent;
  - the continent and the HDI are negatively correlated;
  - the continent and the HDI are positively correlated;
  - the information provided is not enough for making a conclusion.
- Let  $X_1, \dots, X_n$  be a large random sample from a distribution  $f_X(x)$  with expectation  $\mu_X$  and variance  $\sigma_X^2$ . Let also  $\bar{X} = \sum_{i=1}^n X_i/n$  be the sample mean and  $S_X^2$  be the sample variance. Which of the following sentences is **not correct**:
  - if  $f_X(x)$  is not the normal distribution then  $\sqrt{n}(\bar{X} - \mu_X)/\sigma_X$  follows exactly a  $t$  distribution;
  - if  $f_X(x)$  is the normal distribution then  $\sqrt{n}(\bar{X} - \mu_X)/\sigma_X$  follows exactly a normal distribution;
  - if  $f_X(x)$  is the normal distribution then  $\sqrt{n}(\bar{X} - \mu_X)/S_X$  follows exactly a  $t$  distribution;
  - if  $f_X(x)$  is not the normal distribution then  $\sqrt{n}(\bar{X} - \mu_X)/S_X$  follows approximately a  $t$  distribution;
  - if  $f_X(x)$  is not the normal distribution then  $\sqrt{n}(\bar{X} - \mu_X)/S_X$  follows approximately a normal distribution.

4. Let us consider a random experiment with sample space given by the seasons of the year, i.e.  $S = \{Spring, Summer, Fall, Winter\}$ . Which of the following is a probability on  $S$ ?
- (a)  $P(Spring) = 1/5$ ;  $P(Summer) = 1/5$ ;  $P(Fall) = 1/5$ ;  $P(Winter) = 1/5$ ;  
 (b)  $P(Spring) = 1$ ;  $P(Summer) = 3/4$ ;  $P(Fall) = 2/4$ ;  $P(Winter) = 1/4$ ;  
 (c)  $P(Spring) = 1/4$ ;  $P(Summer) = 1/4$ ;  $P(Fall) = 1/4$ ;  $P(Winter) = 0$ ;  
 (d)  $P(Spring) = 1/4$ ;  $P(Summer) = 2/4$ ;  $P(Fall) = 3/4$ ;  $P(Winter) = 1$ ;  
 (e)
5. A researcher has asked the thirteen married men in a small community about the brideprice they had to pay to the bride's family when they got married. The brideprice values (in USD) are
- 20000 3000 10000 20000 13000 0 4000 20000 63000 8000 3000 12000 31000
- What is the **range** of the brideprice?
- (a) -9500;  
 (b) 9500;  
 (c) 11000;  
 (d) 16500;  
 (e)
6. An ice-cream shop offers 10 different flavors. How many combinations of 2 scoops can be made if the order is important and the flavors can be used more than once?
- (a) 20;  
 (b) 45;  
 (c) 55;  
 (d) 90;  
 (e)
7. In a card game, the player has three possible outcomes: win, tie or lose. If the player wins (which happens with probability 0.19), he gets two dollars; if the player loses (which happens with probability 0.47), he loses one dollar; in the case of a tie, the player neither wins nor loses any money. What is the expected amount of money of the player at the end of one game?
- (a) -0.39;  
 (b)    
 (c) 0.00;  
 (d) 0.33;  
 (e) 1.00.

8. Coffee Inc. is a company that imports two types of coffee to Sweden. The number of sacks of the type *Arabica* imported every month can be described by a normally distributed random variable with expectation 50 and variance 4. The number of sacks of the type *Robusta* imported every month can be described by a normally distributed random variable with expectation 45 and variance 9. The covariance between the number of sacks of both types is 5. What is the probability that the total imports during one month exceed 100 sacks?
- (a) 0.00;
  - (b) 0.08;
  - (c) ;
  - (d) 0.35
  - (e) 0.41.
9. A teacher knows by experience that the number of points students get in the final exam can be adequately modeled by a normal distribution. A sample of nine students has been selected and their score in the exam has been measured. The sample mean is  $\bar{x}_s = 53.6$  and the sample variance is  $S_{x,s}^2 = 492.5$ . A 90% confidence interval for the expected number of points of the students in the exam is:
- (a) (12.3, 94.9);
  - (b) ;
  - (c) (41.4, 65.8);
  - (d) (43.3, 63.9);
  - (e) (44.1, 63.1).
10. One week before the local elections of a city, a candidate, Mrs. A, believes that more than 30% of the voters support her. In order to verify her claim, the campaign has selected a sample of 100 voters. 37 out of the 100 voters in the sample claim that they will vote for Mrs. A. The value of the statistic for testing the alternative that the proportion of voters for Mrs. A is larger than 30% is:
- (a) ;
  - (b) 1.66;
  - (c) 1.98;
  - (d) 14.50;
  - (e) 30.03;

11. The teacher of a course in statistics wants to explain the score of students in the final exam (variable *exam*) in terms of the score in a previous home assignment (variable *assignment*) through a linear regression of the form:

$$exam = \beta_0 + \beta_1 assignment + \epsilon$$

The following table shows the scores of the eight students in the course:

<i>Assignment</i>	42	48	50	50	51	55	59	67
<i>Exam</i>	38	43	57	33	81	50	48	84

The estimated intercept of the regression line of interest is:

- (a) -548.7;  
 (b)   
 (c) 0.6;  
 (d) 1.5;  
 (e) 11.4.
12. The following table summarizes the scores of 170 students in an exam of statistics:

Points	[0, 40)	[40, 50)	[50, 60)	[60, 70)	[70, 80)	[80, 90)	[90, 100)
Frequency	51	17	22	34	21	17	8

The teacher of the course wants to test whether the scores in the previous table can be considered as a random sample from a (truncated) normal distribution. If it was, the probability in each class would be as given in the following table. (**Hint:** This is a goodness of fit test.)

Points	[0, 40)	[40, 50)	[50, 60)	[60, 70)	[70, 80)	[80, 90)	[90, 100)
Probability	0.33	0.17	0.17	0.14	0.10	0.06	0.03

Having a significance level  $\alpha = 0.05$ , what is the critical value:

- (a) 1.97;  
 (b)   
 (c) 14.07;  
 (d) 18.51;  
 (e) 389.92.

## Part one. Multiple choice

1. See Section 1.1 in the lecture notes.
2. The null hypothesis in a test of independence is that both variables are independent. The alternative is that they are dependent. As  $p\text{-value} < \alpha$ , the null hypothesis is rejected.
3. See Section 6.2 in Newbold et al. or Chapter 7 in the lecture notes.
4. In all cases (a) to (d) we have  $P(S) \neq 1$ .
5.  $\text{range}_x = x_{(13)} - x_{(1)} = 63000 - 0 = 63000$ .
6. We are selecting  $x = 2$  flavors out of  $n = 10$ , which gives  $n^x = 100$ .
7. Let  $X =$  "Amount of money of the player at the end of one game". We have  $P_X(2) = 0.19$ ,  $P_X(-1) = 0.47$  and  $P_X(0) = 0.34$ . Then

$$\mu_X = \sum_x xP_X(x) = 2 \cdot 0.19 + (-1) \cdot 0.47 + 0 \cdot 0.34 = -0.09$$

8. Let  $X =$  "Number of sacks of the type *Arabica* imported during one month" and  $Y =$  "Number of sacks of the type *Robusta* imported during one month". We have  $X \sim N(\mu_X, \sigma_X^2)$  with  $\mu_X = 50$  and  $\sigma_X^2 = 4$  and  $Y \sim N(\mu_Y, \sigma_Y^2)$  with  $\mu_Y = 45$  and  $\sigma_Y^2 = 9$ . We have also that  $\sigma_{XY} = 5$ .

Let  $W = X + Y =$  "Total sacks imported during one month". We have  $W \sim N(\mu_W, \sigma_W^2)$  with  $\mu_W = \mu_X + \mu_Y = 95$  and  $\sigma_W^2 = \sigma_X^2 + \sigma_Y^2 + 2\sigma_{XY} = 23$ . Then  $\sigma_W = 4.8$ . This yields

$$P(W > 100) = 1 - P(Z < 1.04) = 0.15.$$

9. We have

$$\hat{\sigma}_{\bar{X}}^2 = \frac{S_{x,s}^2}{n} = \frac{492.5}{9} = 54.72.$$

The confidence interval is then

$$\bar{x}_s \pm t_{n-1, \alpha/2} \hat{\sigma}_{\bar{X}} = 53.6 \pm 1.860 \cdot 7.40 = (39.8, 67.4).$$

10. We have  $\bar{x} = \hat{p} = 37/100 = 0.37$  and  $\hat{\sigma}_{\bar{X}}^2 = \hat{p}(1 - \hat{p})/n = 0.37(1 - 0.37)/100 = 0.002331$ . Therefore the test statistic is  $t_{obs} = (\bar{x} - \mu_0)/\hat{\sigma}_{\bar{X}} = (0.37 - 0.3)/0.04828 = 1.45$ .
11. The estimated slope is  $b_1 = S_{xy,s}/S_{x,s}^2 = 86.79/57.64 = 1.5$ .  
The estimated intercept is  $b_0 = \bar{y} - b_1\bar{x} = 54.25 - 1.5 \cdot 52.75 = -25.2$ .
12.  $\chi_{K-1, \alpha}^2 = \chi_{6, 0.05}^2 = 12.59$ .

## Part two. Complete solution

13. (a) On the first semester of a year, 232 students took the final exam in a course of statistics. 141 students passed the exam. Considering the set of students taking the exam as a random sample, find a 95% confidence interval for the passing rate (i.e. the proportion of students passing the exam). (5p.)
- (b) On the second semester of the year, 261 students took the final exam of the course. 139 students passed the exam. Considering the set of students taking the exam as a random sample, find a 95% confidence interval for the passing rate. (5p.)
- (c) Assuming that the sets of students taking the exam each semester are independent of each other, find a 95% confidence interval for the difference of the passing rate between the first and the second semester of the year. (5p.)
- (d) Assuming that all 493 students taking the exam during the year are independent of each other and that the passing rate is equal to 0.5, find the (approximated) probability that the number of students passing the exam is larger than 280. (**Hint:** What is the distribution of the number of students passing the exam?) (5p.)

### Solution

- (a). We have  $n = 232$  and  $\alpha = 0.05$ , which yields  $\hat{p} = 141/232 = 0.6078$  and  $t_{n-1, \alpha/2} = 1.97$ . The confidence interval is

$$\hat{p} \pm t_{n-1, \alpha/2} \sqrt{\frac{\hat{p}(1-\hat{p})}{n}} = 0.6078 \pm 1.97 \sqrt{\frac{0.6078(1-0.6078)}{232}} = (54.5\%, 67.1\%).$$

- (b). We have  $n = 261$  and  $\alpha = 0.05$ , which yields  $\hat{p} = 139/261 = 0.5326$  and  $t_{n-1, \alpha/2} = 1.97$ . The confidence interval is

$$\hat{p} \pm t_{n-1, \alpha/2} \sqrt{\frac{\hat{p}(1-\hat{p})}{n}} = 0.5326 \pm 1.97 \sqrt{\frac{0.5326(1-0.5326)}{261}} = (47.2\%, 59.3\%).$$

- (c). Let  $X$  and  $Y$  be, respectively, the random variables indicating whether a student selected at random from the first and the second semester pass the exam or not. We get  $\bar{d}_s = \hat{p}_x - \hat{p}_y = 0.6078 - 0.5326 = 0.0752$  and

$$\hat{\sigma}_{\bar{D}}^2 = \frac{\hat{p}_x(1-\hat{p}_x)}{n_x} + \frac{\hat{p}_y(1-\hat{p}_y)}{n_y} = \frac{0.6078(1-0.6078)}{232} + \frac{0.5326(1-0.5326)}{261} = 0.001981.$$

Also,  $t_{n_x+n_y-2, \alpha/2} = 1.96$ . The confidence interval is

$$\bar{d}_s \pm t_{n_x+n_y-2, \alpha/2} \hat{\sigma}_{\bar{D}} = 0.0752 \pm 1.96 \cdot 0.0445 = (-1.2\%, 16.3\%).$$

- (d). Let  $Y =$  “number of students passing the exam during the year”, we have that  $Y \sim \text{Bin}(n, P)$  with  $n = 493$  and  $P = 0.5$ . As  $n$  is large, we have that  $Y \stackrel{\text{approx}}{\sim} N(nP, nP(1-P)) = N(246.5, 123.25)$ . And, approximately,  $P(Y > 280) = 1 - P(Z < 3.02) = 0.0013$ .

14. In order to fit the linear regression that explains “sleeping time” (variable *sleep*, in hours per week) in terms of “working time” (variable *work*, in hours per week) and a dummy variable indicating whether the individual has kids younger than three

$$young\_kid = \begin{cases} 1 & \text{if the individual has kids younger than 3} \\ 0 & \text{otherwise} \end{cases}$$

data on  $n = 10$  individuals was collected and it is shown in the following table.

work	young_kid	sleep
34	0	35
29	0	48
43	0	49
41	1	52
32	0	54
29	1	56
44	0	57
16	0	58
59	0	60
6	0	65

The estimated regression is

$$\widehat{sleep} = 57 - 0.12 work + 1.00 young\_kid.$$

- (a) Interpret the three estimated coefficients. (5p.)
- (b) Predict the number of hours sleeping for a person working forty hours per week who has one kid younger than three. (5p.)
- (c) Find the ten fitted values and the ten residuals. (5p.)
- (d) Calculate the coefficient of determination and the adjusted coefficient of determination. (**Note:** The adjusted coefficient of determination may take an unexpected value.) (5p.)

## Solution

- (a). We expect a person who does not work ( $work = 0$ ) and does not have young kids ( $young\_kid = 0$ ) to sleep around 57 hours per week.

Having the same status regarding having young kids, we expect a person who works one hour more than another one, to sleep 0.12 hours (seven minutes) less.

Among two persons who work the same amount of time, we expect one person having young kids to sleep one hour more than one person who does not have young kids.

- (b). We have  $\widehat{sleep} = 57 - 0.12 work + 1.00 young\_kid = 57 - 0.12 \times 40 + 1.00 \times 1 = 53.2$ .

- (c). The fitted value associated to the first individual is

$$\widehat{sleep}_1 = 57 - 0.12 work_1 + 1.00 young\_kid_1 = 57 - 0.12 \times 34 + 1.00 \times 0 = 52.92.$$

The associated residual is  $e_1 = sleep_1 - \widehat{sleep}_1 = 35 - 52.92 = -17.92$ . The remaining fitted values and residuals are found in an analogous way. They are presented in the following table.

work	young_kid	sleep	$\widehat{sleep}$	residual
34	0	35	52.92	-17.92
29	0	48	53.52	-5.52
43	0	49	51.84	-2.84
41	1	52	53.08	-1.08
32	0	54	53.16	0.84
29	1	56	54.52	1.48
44	0	57	51.72	5.28
16	0	58	55.08	2.92
59	0	60	49.92	10.08
6	0	65	56.28	8.72

- (d). We have  $n = 10$ ,  $K = 2$ ,

$$SSE = \sum_s (sleep_i - \widehat{sleep}_i)^2 = (-17.92)^2 + \dots + 8.72^2 = 577.8$$

and

$$SST = \sum_s (sleep_i - \overline{sleep}_i)^2 = (35 - 53.4)^2 + \dots + (65 - 53.4)^2 = 608.4.$$

This yields  $R^2 = 1 - SSE/SST = 5.03\%$  and  $\bar{R}^2 = 1 - (SSE/(n - K - 1))/(SST/(n - 1)) = -22.1\%$ .