

## Introduction to statistical inference – Project work – Generation 2022/23

Name and Surname: \_\_\_\_\_ Student ID number: \_\_\_\_\_

Professor: \_\_\_\_\_ Grade: \_\_\_\_\_

Using the dataset *GenderEqualityIndex.sav* conduct the analysis, answer the questions and tasks. The analysis can be done using SPSS or R. If the analysis is done in R, it is desirable to print the code. If the analysis is done in SPSS, there is no need to print the output. **Note: Depending on the software used and software version, the results among students might differ.**

The *GenderEqualityIndex.sav* is a dataset that provides the data on the European Gender Equality Index – a composite indicator which aims to measure the level of (in)equality among males and females in the European Union. Besides the name of the country, the overall value of the EGDI, you are provided with indicator data. Based on the pre-calculated values per pillar or dimension (Work, Money, Knowledge, Time, Power, and Health), the countries have been classified in the following manner:

Pillar or dimension	Code	Value
Work, Money, Knowledge, Time, Power, and Health	1	High
	2	Medium
	3	Low

- Obtain the parameters of descriptive statistics for the indicators *Participation* and *Segregation and quality of work*:

Indicator	Mean	SD	Median	Minimum	Maximum
<i>Participation</i>					
<i>Segregation and quality of work</i>					

- Indicator with the higher mean is \_\_\_\_\_
  - Indicator with the smaller variance is \_\_\_\_\_
  - Is the median of *Participation* greater than the median of *Segregation and quality of work*? Yes No
  - Range of *Segregation and quality of work* is \_\_\_\_\_
  - Interquartil range of *Participation* is \_\_\_\_\_
- How many countries belong to each group based on the calculated pillar *Time*?

Low	Medium	High

3. Is there statistically significant difference between the country groups based on *Power* for the indicator *Financial resources*?

The first step in the analysis is to conduct the \_\_\_\_\_ test.

$H_0$  \_\_\_\_\_

$H_1$  \_\_\_\_\_

Value of the statistics \_\_\_\_\_  $p$  value \_\_\_\_\_

We conclude \_\_\_\_\_

To test the differences we use \_\_\_\_\_

$H_0$  \_\_\_\_\_

$H_1$  \_\_\_\_\_

Value of the statistics \_\_\_\_\_  $p$  value \_\_\_\_\_

We REJECT ACCEPT (we could not reject) the  $H_0$

Conclusion

---



---

What would be your next steps if the  $H_0$  is rejected? Just explain the procedure without conducting it.

---



---

4. Is there difference in the values of the indicator *Not at-risk-of-poverty (%)* between countries which have medium and high *Work*?

The first step in the analysis is to conduct the \_\_\_\_\_ test.

$H_0$  \_\_\_\_\_

$H_1$  \_\_\_\_\_

Value of the statistics \_\_\_\_\_  $p$  value \_\_\_\_\_

We conclude \_\_\_\_\_

To inspect whether there are differences we use \_\_\_\_\_

$H_0$  \_\_\_\_\_

$H_1$  \_\_\_\_\_

Value of the statistics \_\_\_\_\_  $p$  value \_\_\_\_\_

We REJECT ACCEPT (we could not reject) the  $H_0$

Conclusion

\_\_\_\_\_

\_\_\_\_\_

5. Is there statistically significant correlation between the indicators *Economic situation* and *Segregation*.

The first step in the analysis is to conduct the \_\_\_\_\_ test.

$H_0$  \_\_\_\_\_

$H_1$  \_\_\_\_\_

Value of the statistics (*Economic situation*) \_\_\_\_\_  $p$  value \_\_\_\_\_

Value of the statistics (*Segregation*) \_\_\_\_\_  $p$  value \_\_\_\_\_

We conclude \_\_\_\_\_

To inspect whether there is significant correlation we use \_\_\_\_\_

$H_0$  \_\_\_\_\_

$H_1$  \_\_\_\_\_

Value of the statistics \_\_\_\_\_ The obtained correlation coefficient is \_\_\_\_\_  $p$  value \_\_\_\_\_

We REJECT ACCEPT (we could not reject) the  $H_0$

Conclusion \_\_\_\_\_

\_\_\_\_\_

6. Does the *Power* level impact the *Money* level? How many countries are there with low *Power* and medium *Money* level?

$H_0$  \_\_\_\_\_

$H_1$  \_\_\_\_\_

We use the \_\_\_\_\_ test

Value of the statistics \_\_\_\_\_  $p$  value \_\_\_\_\_

We REJECT ACCEPT (we could not reject) the  $H_0$

Conclusion \_\_\_\_\_

\_\_\_\_\_

There are \_\_\_\_\_ countries with low *Power* and medium *Money* level.

7. Is there statistically significant difference between countries with low and medium *Knowledge* regarding the values of the indicator *Attainment and participation*?

The first step in the analysis is to conduct the \_\_\_\_\_ test.

$H_0$  \_\_\_\_\_

$H_1$  \_\_\_\_\_

Value of the statistics \_\_\_\_\_  $p$  value \_\_\_\_\_

We conclude \_\_\_\_\_

Therefore, we will use the \_\_\_\_\_

As a pre-test, we conduct the \_\_\_\_\_ test.

$H_0$  \_\_\_\_\_

$H_1$  \_\_\_\_\_

Value of the statistics \_\_\_\_\_  $p$  value \_\_\_\_\_

We REJECT ACCEPT (we could not reject) the  $H_0$  of the pre-test

We conduct the test

$H_0$  \_\_\_\_\_

$H_1$  \_\_\_\_\_

Value of the statistics \_\_\_\_\_  $p$  value \_\_\_\_\_

We REJECT ACCEPT (we could not reject) the  $H_0$

Conclusion

\_\_\_\_\_

\_\_\_\_\_

Mean values of the indicator *Attainment and participation* per group are:

Low *Knowledge* \_\_\_\_\_ Medium *Knowledge* \_\_\_\_\_

8. Is the mean value of the indicator *Care activities* on the population 75?

The first step in the analysis is to conduct the \_\_\_\_\_ test.

$H_0$  \_\_\_\_\_

$H_1$  \_\_\_\_\_

Value of the statistics \_\_\_\_\_  $p$  value \_\_\_\_\_

We conclude \_\_\_\_\_

In the next step we will use \_\_\_\_\_ test

$H_0$  \_\_\_\_\_

$H_1$  \_\_\_\_\_

Value of the statistics \_\_\_\_\_  $p$  value \_\_\_\_\_

We REJECT ACCEPT (we could not reject) the  $H_0$

Conclusion \_\_\_\_\_

\_\_\_\_\_

\_\_\_\_\_

9. Create a linear regression model and model how *Graduates of tertiary education (%)* impact *Attainment and participation*. Answer the following questions:

A) What is the obtained model equation?

\_\_\_\_\_

B) Is the coefficient related to the impact of *Graduates of tertiary education (%)* statistically significant?

$H_0$  \_\_\_\_\_

$H_1$  \_\_\_\_\_

Value of the statistics \_\_\_\_\_  $p$  value \_\_\_\_\_

We conclude \_\_\_\_\_

\_\_\_\_\_

C) Is the overall model statistically significant?

$H_0$  \_\_\_\_\_

$H_1$  \_\_\_\_\_

Value of the statistics \_\_\_\_\_  $p$  value \_\_\_\_\_

We conclude \_\_\_\_\_

\_\_\_\_\_

D) Comment on the model quality

$R^2 =$  \_\_\_\_\_ The  $R^2$  indicates that \_\_\_\_\_

E) Can the created model be used for predictions?

Answer \_\_\_\_\_

10. Is there statistically significant difference between the *Health* groups if we observe the values of the indicator *Social activities*.

The first step in the analysis is to conduct the \_\_\_\_\_ test.

$H_0$  \_\_\_\_\_

$H_1$  \_\_\_\_\_

Value of the statistics \_\_\_\_\_  $p$  value \_\_\_\_\_

We conclude \_\_\_\_\_

To test the differences we use \_\_\_\_\_

$H_0$  \_\_\_\_\_

$H_1$  \_\_\_\_\_

Value of the statistics \_\_\_\_\_  $p$  value \_\_\_\_\_

We ☐ REJECT ☐ ACCEPT (we could not reject) the  $H_0$

Conclusion \_\_\_\_\_

\_\_\_\_\_

If the  $H_0$  is rejected, how would you proceed with the analysis? You do not have to conduct the analysis.

\_\_\_\_\_

\_\_\_\_\_

\_\_\_\_\_

11. Create a linear regression model and model how *Social*, *Access*, and *Political* impact *Economic*. Answer the following questions:

A) What is the obtained model equation?

\_\_\_\_\_

B) Is the coefficient related to the impact of *Social* statistically significant?

$H_0$  \_\_\_\_\_

$H_1$  \_\_\_\_\_

Value of the statistics \_\_\_\_\_  $p$  value \_\_\_\_\_

We conclude \_\_\_\_\_

C) Is the coefficient related to the impact of *Access* statistically significant?

$H_0$  \_\_\_\_\_

$H_1$  \_\_\_\_\_

Value of the statistics \_\_\_\_\_  $p$  value \_\_\_\_\_

We conclude \_\_\_\_\_

D) Is the coefficient related to the impact of *Political* statistically significant?

$H_0$  \_\_\_\_\_

$H_1$  \_\_\_\_\_

Value of the statistics \_\_\_\_\_  $p$  value \_\_\_\_\_

We conclude \_\_\_\_\_

E) Is the overall model statistically significant?

$H_0$  \_\_\_\_\_

$H_1$  \_\_\_\_\_

Value of the statistics \_\_\_\_\_  $p$  value \_\_\_\_\_

We conclude \_\_\_\_\_

\_\_\_\_\_



F) Comment on the model quality

Adjusted  $R^2$  = \_\_\_\_\_ The  $R^2$  indicates that \_\_\_\_\_

\_\_\_\_\_

G) Is there a problem of multicollinearity in the model?

Answer \_\_\_\_\_

H) Is there a problem of autocorrelation in the model?

Value of the DW test \_\_\_\_\_

Knowing that  $dl=0.969$   $du=1.414$  draw the Durbin Watson critical range and make the conclusion:

\_\_\_\_\_

Durbin Watson critical range

