

SOC252 Winter 2023

Assignment 1

Due date: Friday 3 February, 11:59pm

Details

There are **100 points** in total.

You will need to submit both your answers to the questions and accompanying R code. You should submit:

- your Quarto file; and
- the knitted PDF resulting from your Quarto file.

Please submit both files via Quercus.

Remember to:

- Label the answers to each question
- Label any graphs clearly with suitable axis labels and titles
- Comment your code so that it is easy to understand

Question 1 (30 points)

This question relates to the GSS dataset. We will be looking at how age at the time of first birth varies by education and current age. Note there are a few different education variables in the GSS dataset but for this question, we will be focusing on the binary `has_bachelor_or_higher` variable.

a)

Report the following descriptive statistics:

- i) What proportion of respondents have a non-missing observation for their age at the time of the birth of their first child?
- ii) What proportion of respondents have a non-missing observation for their highest level of education (`has_bachelor_or_higher`)?
- iii) For those respondents who have a non-missing education value:
 - What is the number of respondents by education group (at least a Bachelor's degree, less than a Bachelor's degree) that have a non-missing observation for age at first birth?
 - What is the proportion of respondents by education group (at least a Bachelor's degree, less than a Bachelor's degree) that have a non-missing observation for age at first birth?

Comment briefly on your calculations.

b)

For parts b) and c), we will be looking at the subset of respondents who have an education level reported, so you can filter out those respondents who have missing values of education.

Plot histograms of age at first birth by education level (at least a Bachelor's degree, less than a Bachelor's degree), with both histograms shown on the same chart but colored in different colors. Use `geom_histogram(position = "dodge")` so that the histograms are plotted next to each other. Interpret your chart.

c)

- i) Calculate the correlation between age and age at first birth. Interpret your finding.
- ii) Calculate the mean age of first birth by age group (`age_group`) and education level (`has_bachelor_or_higher`). Note that this question uses the categorical variable `age_group` (not the continuous variable `age`)
- iii) Create a line chart of the results from part iii), plotting mean age of first birth (y axis) versus age group (x axis), with a separate line (and different color) for education level. Comment on your chart. Does the pattern over age agree with your findings from part i)? Why or why not?

Question 2 (40 points)

This question relates to the country indicators dataset. You will be looking at two country's fertility rate (TFR) and child mortality rate. Specifically, the two variables are

- `tfr` = total fertility rate, which is the average number births per woman in that particular country and year
- `child_mort` = under-five child mortality rate, which is the number of deaths to children aged 5 or less per 1,000 live births.

a)

Choose two countries. For each country, find

- i) the mean TFR and child mortality rate
- ii) the year which has the highest TFR
- iii) the correlation between TFR and child mortality

b)

Focusing on child mortality, report the percent increase or decrease in the child mortality rate observed from 2009 to 2017.

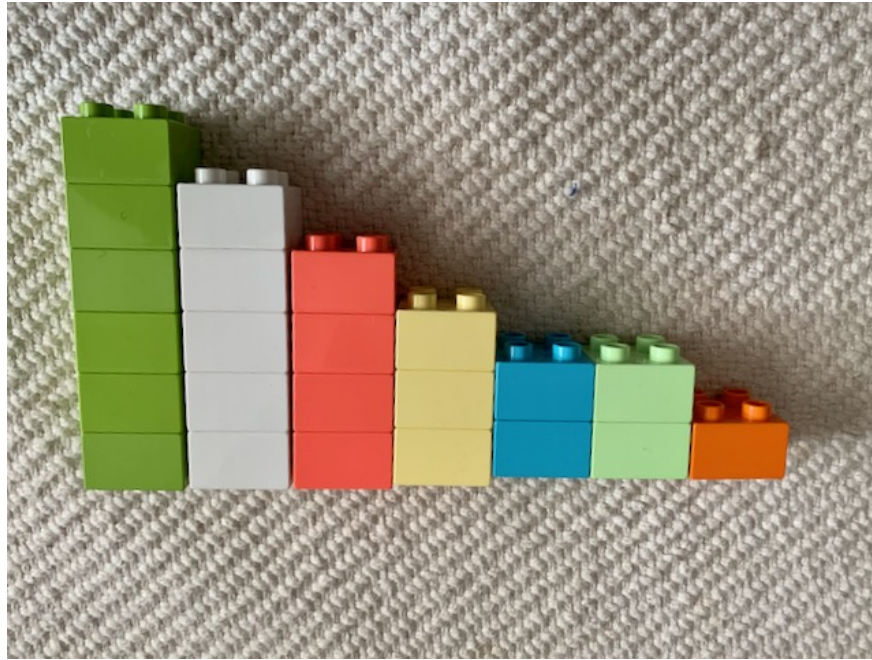
c)

Make one graph (only one!) that illustrates the trend in child mortality over time in the two countries of interest.

In all parts, comment briefly on what you observe.

Question 3 (30 points)

The photo below shows counts of lego blocks of various colors, organized from smallest to largest.



Recreate this photo as a bar chart in ggplot. Notes:

- You will first need to create your own dataset which summarizes the information contained in this photo
- The x axis should show the color names, and the y axis should show the count of bricks.
- You are not required to change the fill color of the bars in the chart, but if you would like to, the following vector of color codes represent colors that are similar to the block colors (in order from left to right): `c("#70961c", "white", "#ee5e4f", "#d5c47c", "#008db3", "#a5d395", "#d35800")`