

# Assignment 4: Collaborating Together

## Introduction to Applied Data Science

### 2022-2023

Giovanna Tullume Carrion  
[g.e.tullumecarrion@student.uu.nl](mailto:g.e.tullumecarrion@student.uu.nl)  
<http://www.github.com/GioviEli>

June 2023

## Assignment 4: Collaborating Together

### Part 1: Contributing to another student's Github repository

In this assignment, you will create a Github repository, containing this document and the .pdf output, which analyzes a dataset individually using some of the tools we have developed.

This time, make sure to not only put your name and student e-mail in your Rmarkdown header, but also your Github account, as I have done myself.

However, you will also pair up with a class mate and contribute to each others' Github repository. Each student is supposed to contribute to another student's work by writing a short interpretation of 1 or 2 sentences at the designated place (this place is marked with **designated place**) in the other student's assignment.

This interpretation will not be graded, but a Github shows the contributors to a certain repository. This way, we can see whether you have contributed to a repository of a class mate.

**Question 1.1:** Fill in the **github username** of the class mate to whose repository you have contributed.

[TBD]

### Part 2: Analyzing various linear models

In this part, we will summarize a dataset and create a couple of customized tables. Then, we will compare a couple of linear models to each other, and see which linear model fits the data the best, and yields the most interesting results.

We will use a dataset called **GrowthSW** from the **AER** package. This is a dataset containing 65 observations on 6 variables and investigates the determinants of economic growth. First, we will try to summarize the data using the **modelsummary** package.

```
library(AER)
data(GrowthSW)
```

One of the variables in the dataset is **revolutions**, the number of revolutions, insurrections and coup d'états in country  $i$  from 1965 to 1995.

**Question 2.1:** Using the function `datasummary`, summarize the mean, median, sd, min, and max of the variables `growth`, and `rgdp60` between two groups: countries with `revolutions` equal to 0, and countries with more than 0 revolutions. Call this variable `treat`. Make sure to also write the resulting data set to memory. Hint: you can check some examples [here](#).

```
library(modelsummary); library(tidyverse)

# write your code here
GrowthSW$treat <- if_else(GrowthSW$revolutions > 0, "revolution", "no revolution")
datasummary(treat * (growth + rgdp60) ~ mean + median + sd + min + max, data=GrowthSW)
```

treat		mean	median	sd	min	max
no revolution	growth	2.46	2.29	1.28	0.42	6.65
	rgdp60	5283.32	5393.00	2439.39	1374.00	9895.00
revolution	growth	1.68	1.92	2.11	-2.81	7.16
	rgdp60	1988.67	1259.00	1698.18	367.00	6823.00

**Designated place:** type one or two sentences describing this table of a fellow student below. For example, comment on the mean and median growth of both groups. Then stage, commit and push it to their github repository.

### Part 3: Make a table summarizing reressions using `modelsummary` and `kable`

In question 2, we have seen that growth rates differ markedly between countries that experienced at least one revolution/episode of political stability and countries that did not.

**Question 3.1:** Try to make this more precise this by performing a t-test on the variable `growth` according to the group variable you have created in the previous question.

```
# write t test here
# Independent Sample T-Test

# t = (mean_1 - mean_2) / sqrt( sd_1**2 / count_1 + sd_2**2 / count_2 )

# I cannot for the life of me figure out how to read the data from the datasummary.
# So I am forced to calculate them here again.

# Sample 1
sample1_mean <- mean(GrowthSW$growth[GrowthSW$treat == "revolution"])
sample1_sd <- sd(GrowthSW$growth[GrowthSW$treat == "revolution"])
sample1_size <- sum(GrowthSW$treat == "revolution")

# Sample 2
sample2_mean <- mean(GrowthSW$growth[GrowthSW$treat == "no revolution"])
sample2_sd <- sd(GrowthSW$growth[GrowthSW$treat == "no revolution"])
sample2_size <- sum(GrowthSW$treat == "no revolution")

ttest_value <- (sample1_mean - sample2_mean) / sqrt( (sample1_sd**2 / sample1_size) + (sample2_sd**2 / sample2_size) )
ttest_value
```

```
## [1] -1.853087
```

**Question 3.2:** What is the  $p$ -value of the test, and what does that mean? Write down your answer below.

Degrees of freedom from two samples:  $\text{DoF} = \text{size\_sample1} + \text{size\_sample2} - 2$ . In our case DoF is thus 63. P-value in this case then is

We can also control for other factors by including them in a linear model, for example:

$$\text{growth}_i = \beta_0 + \beta_1 \cdot \text{treat}_i + \beta_2 \cdot \text{rgdp60}_i + \beta_3 \cdot \text{tradeshare}_i + \beta_4 \cdot \text{education}_i + \epsilon_i$$

**Question 3.3:** What do you think the purpose of including the variable `rgdp60` is? Look at `?GrowthSW` to find out what the variables mean.

`rgdp60` shows the GDP of the country in 1960. `Growth` shows the growth in GDP. `rgdp60` is added to know what is the initial value of the GDP, because only a growth does not say to much if you do not know where they originally came from.

We now want to estimate a stepwise model. Stepwise means that we first estimate a univariate regression  $\text{growth}_i = \beta_0 + \beta_1 \cdot \text{treat}_i + \epsilon_i$ , and in each subsequent model, we add one control variable.

**Question 3.4:** Write four models, titled `model1`, `model2`, `model3`, `model4` (using the `lm` function) to memory. Hint: you can also use the `update` function to add variables to an already existing specification.

Now, we put the models in a list, and see what `modelsummary` gives us:

```
list(model1, model2, model3, model4) |>
  modelsummary(stars=T,
# edit this to remove the statistics other than R-squared
# and N
)
```

```
## Error in eval(expr, envir, enclos): object 'model1' not found
```

**Question 3.5:** Edit the code chunk above to remove many statistics from the table, but keep only the number of observations  $N$ , and the  $R^2$  statistic.

**Question 3.6:** According to this analysis, what is the main driver of economic growth? Why?

**Question 3.7:** In the code chunk below, edit the table such that the cells (including standard errors) corresponding to the variable `treat` have a red background and white text. Make sure to load the `kableExtra` library beforehand.

```
library(kableExtra)
```

```
## Error in library(kableExtra): there is no package called 'kableExtra'
```

```
list(model1, model2, model3, model4) |>
  modelsummary(stars=T, gof_map = c("nobs", "r.squared"))
```

```
## Error in eval(expr, envir, enclos): object 'model1' not found
```

```
# use functions from modelsummary to edit this table
```

**Question 3.8:** Write a piece of code that exports this table (without the formatting) to a Word document.

**The End**