

Project Instructions

Intro to Statistics and Data Analysis with R (0560.1823)

Adi Sarid

January 2020

Background

The following document contains instructions to the project in the Introduction to Statistics and Data Analysis with R course.

The project has a weight of 30% of your final grade.

Goal

The goal of the project is to demonstrate and practice the different elements we have been talking about, which are a part of most data analysis/data science projects.

Methods

In this project you will handle the different phases of a data analysis project:

1. Data **Import** (reading the data into R).
2. Data **Tidying** (arranging the data into something you can work with)
3. Understanding the data:
 - a. **Transforming** variables.
 - b. **Visualizing** (using ggplot2 to show distribution of variables, relationships between variables, and to hypothesize).
 - c. **Modelling**: using a few of the tools we have learned during the course (like hypothesis testing, regression, analysis of variance, etc.) to examine your hypothesis.
4. **Communicate** your findings via an RMarkdown report.

Instructions and Schedule

The project should be performed individually.

Choosing a dataset

First, you should select a dataset on which you will perform the project. I recommend using a data set from either Kaggle or from tidyTuesday. You can select something else.

In any case, please do not choose something “too popular” (e.g., no built-in R datasets). **Check with me before you start**, so that I can confirm your dataset: email adi@sarid-ins.co.il with:

- The dataset name
- Source (a url with the data and documentation of the dataset)
- A direct link to download the raw data you will be using

Please choose a dataset and email me the details by **January 12th 2020** (you can do it earlier and start the work on the project once I email you back with the green light).

Draft submission (initial submission)

I'm giving you the opportunity to submit a draft of your work for a review. If needed, I will suggest revisions, and allow you to correct your work in order to improve your grade. The draft submissions will be graded and you can decide based on the grade you received if you want to improve them in order to improve your grade.

You should submit your draft by **February 1st 2020**.

Please submit your file to moodle as `statintro_draft_studentname_studentID.zip` which budles an Rmd version, data files, and a knitted html version of your report. The Rmd should compile standalone in every computer.

Final submission

Final submissions should be made by **February 20th 2020**.

Please submit your file to moodle as `statintro_final_studentname_studentID.zip` which budles an Rmd version, data files, and a knitted html version of your report. The Rmd should compile standalone in every computer.

Grading

You will be graded along the following lines:

- Data import and tidying (10%): Your ability to use the proper methods to import the data, and tidy it towards the next stages.
- Visualizations (25%): Your ability to utilize visualizations to articulate your hypothesis and to illustrate different patterns and relationships in the data. You should be able to match the proper types of charts to what ever it is you are trying to show.
- Modelling (25%): Your ability to match the appropriate statistical tests/models to the problem, verifying (or highlighting) certain assumptions which are valid or invalid in this case. Please provide at least two relevant models/hypothesis tests that we learned.
- Communication, documentation, explanations (25%): You should be able to explain the different steps you are doing, lead the reader in a logical and appealing manner, explain your results, and highlight the research or business implications of your findings.
- Code (15%): Readability, proper use and proper documentation of code.

Good luck!