

STATS/CSE 780

Homework Assignment 1

Pratheepa Jeganathan

30 January, 2023

Instruction

- **Due before 10:00 PM on Tuesday, January 31, 2023.**
- **Submit a copy of PDF with your report (2-3 pages) and technical supplemental material (less than 10 pages) to Avenue to Learn using the link that was emailed to you.**
 - Technical supplemental material can only include R or Python codes for the results reported.
- **Late penalty for assignments: 15% will be deducted from assignments each day after the due date (rounding up).**
- **Assignments won't be accepted after 48 hours after the due date.**

Assignment Standards

Your assignment must conform to the Assignment Standards listed below.

- RMarkdown or L^AT_EX is strongly recommended to write the report, and RMarkdown must be used to write the supplemental material.
- Report is about the results by applying data science methods and how you interpret or discuss the results. Don't show in the report how you do the analysis using R/Python.

- Technical supplemental material is how you produce the report results using R/Python. Don't print chunk messages, warnings, or extended data frames in the PDF.
- Write your name and student number on the title page. We will not grade assignments without the title page.
- Eleven-point font (times or similar) must be used with 1.5 line spacing and margins of at least 1-inch all around (The RMarkdown for this assignment used these formats).
- Your report may not exceed **three pages**, inclusive of tables and figures. It would help if you chose the tables and figures accordingly for the report. You can also keep other tables and figures in the supplementary material (**less than 10 pages**) and refer to the report. In addition, you may use one page for the bibliography and one page for the title page.
- You may discuss homework problems with other students, but you have to prepare the written assignments yourself.
- No screenshots are accepted for any reason.
- The writing and referencing should be appropriate to the graduate level.
- Various tools, including publicly available internet tools, may be used by the instructor to check the originality of submitted work.
- If you use ChatGPT for R/Python coding, write in the supplementary material. There is no point deduction.

You can use either R (ggplot2 and R packages for data transformation), which is covered in class, or Python (matplotlib and Python modules for data transformation).

[Statistics Canada](#) is the national statistical office. In this assignment, you will use Statistics Canada's publicly available data to explore Canada's economy, society and environment.

Visit the [data site](#). On the left, you will see filtering options. On the right, you will see available data related to selected options. For this assignment, you can use the data in tabular format (Table).

On the left (filtering options)

1. Choose Province or Territory (you will need at least one quantitative and categorical variable by provinces and territory).
2. Choose one subject.
3. Choose frequency (you will need at least ten-time points.)

On the right (table)

4. Choose the data description. The link will take you to a page displaying a dashboard.

On the data dashboard page,

5. Choose the reference time (you will need at least ten-time points).
 6. At the bottom of the page, you can find "how to cite". Keep the doi in your reference manager to cite in your report.
 7. On the top right of the dashboard, you have a button for the "Download option."
- Choose the "Download entire table" option from the pop-up window.
 - Unzip the downloaded zip file.
 - You can access the data and description in the unzipped folder.

If there are too many variables and samples, you can choose a subset of data after downloading to make the following plots.

- (i) Briefly describe your chosen dataset and clearly explain where it was sourced.
- (ii) Clearly explain data transformation and the preprocessing methods you used to tidy the data.
- (iii) Choose one (quantitative) variable for the following analysis. Then, use an appropriate visualization method to describe the trend of the variable in the selected time frequency across provinces. Finally, clearly describe any statistical transformation used for visualization and interpret the results.
- (iv) Aggregate (aggregate over the provinces) the selected variable (from iii) for Canada and inspect the trend over the selected time frequency using an appropriate visualization method. Interpret the results.
- (v) You can use either Shiny in R, which is covered in class, or Streamlit for the following analysis.
 - Choose a categorical variable with more than two categories—product type or health status, income status, etc.
 - Use an appropriate plot to show the change of quantitative variable (from iii or any other quantitative variable) in the selected time frequency across provinces when the user chooses the category.
 - The supplementary material must include the R or Python code you used to create the app.
 - You must submit a link to your Shiny App or Streamlit community cloud. We (Instructor or TA) must have access to the app when we grade it; otherwise, no points for the app (only for the code if provided) are given. Describe your app in the report.

For all the questions, write a clear and concise interpretation of the plots and clearly state what conclusions can be drawn from the plots or graphs — these conclusions should be cast in the context of the chosen dataset.

- Plots must be readable.
- Choose an appropriate font size for plots.
- Label all aesthetics and axes in the plot.
- Use appropriate statistical transformation for plots.

Grading scheme

(i)	Data descrip- tion	Describe the chosen dataset (background of the dataset and the variables) [3]
(ii)	Data transfor- mation	Did you choose all the downloaded variables and observations or a subset? Describe the reasons for using all the data or the subset. [2]
	Pre-processing	How did you identify missing values? How did you represent the missing values in tidy data? How did you identify outliers? If there were any outliers, how did you handle them? [4]
(iii)	Plot	Appropriate plot, the plot is readable, appropriate font size, label all aesthetics and axes, use appropriate statistical transformation, interpretation (how to read the plot), conclusion (any interesting patterns) [4]
(iv)	Plot	Appropriate plot, the plot is readable, appropriate font size, label all aesthetics and axes, use appropriate statistical transformation, interpretation (how to read the plot), conclusion (any interesting patterns) [4]
(v)	Shiny app	Link to shiny app works, the output is an appropriate plot, shiny app reacts to the user inputs (categories), description of the app is written in the report or the app [4]
	Plot	Plot is readable, appropriate font size, label all aesthetics and axes, use appropriate statistical transformation, interpretation (how to read one of the plots), conclusion (any interesting pattern in one of the plots) [4]
References		Reference list starts on a new page, references are appropriate and list out in the report [2]
Supplementary material		Supplementary material starts on a new page, code readability, all codes are within the margins, the R codes and the outputs for the questions are presented [3]
	Shiny app or Streamlit	Shiny app codes (don't execute the codes when you create the PDF) or Streamlit workflow [2]

The maximum points for this assignment is 32. We will convert this to 100%.