

# Informative Title Name

STA304 - Assignment 3

GROUP NUMBER: ADD YOUR NAMES HERE

November 5, 2021

## Introduction

<Here you should have a few paragraphs of text introducing the problem, getting the reader interested/ready for the rest of the report.>

<Introduce terminology.>

<Highlight hypotheses.>

<Optional: You can also include a description of each section of this report as a last paragraph.>

## Data

### Data Collection Process

The Data is a collection of two datasets. The first dataset contains data from the General Social Survey on Family (cycle 31) on 2017. Canada's General Social Survey (GSS) program conducts annual survey covering one topic in depth (citation). As such, this dataset contains mostly contains information pertaining to families. However, we will later investigate some common variables between this dataset and the Canada Election Study (CES) dataset. This will give us a list of factors in the general population that could be potentially associated to political affiliations

the Canada Election Study (CES) data which was collected in 2019. This data was collected from a questionnaire delivered to the people living in Canada through the Computer Assisted Telephone Interviewing(CATI). Phone calls were made potential interviewees during both the day and evening for every weekday. These questions asked some personal information such as their age and also political opinions such as "what is your opinion on Justin Trudeau?" (citation)

### Data cleaning

The regression model will be constructed, will be applied to both datasets, therefore only variables that appeared in both the datasets could be used. One way was examine if there were common terms present in the columns for both the dataset. Unfortunately most the columns the CES dataset were just question numbers. However, each column also had a label that stated the question itself. For example column "q2" had the label "IN what year were you born?". Hence we collected common words that appeared in the columns of the GSS data and the labels of the CES data to find possible topics that were common in both datasets.

Then for each topic we search the columns for both the dataset to see if any two column were describing identical or similar variables. These are the variables we found:

- Age: The age of the person was recorded in both data with the same column names.
- Sex: Gender was recorded in CES data and sex was recorded in the GSS data. Here we assumed the people who reported their gender to be male or female also would have the same sex as their gender.

Hence we removed all genders that were not male or female and renamed the q3 column in the CES data to sex.

- Province: Both the datasets recorded the province a person lived in.
- religious importance: Both datasets recorded how the level of importance for religion. The categorical values in both the datasets were the same with the exception of an additional value “Refused” being present in the CES data. We removed rows containing this value from the CES data.
- Aboriginal: Both datasets recorded if a person was aboriginal or not. This was an option for question 66a in CES data which asked which ethnicity the person belonged to. q66a\_15 recorded if the person belonged to an Aboriginal group. the value “(1) Selected” indicated the person identified as original and “(0) Not Selected”. We filtered out other values as the indicated that the question was skipped or if the person was not sure. Then we mutated “(1) Selected” to “Yes” and “(0) Not Selected” to “No” to match the values of the variable in both dataset. Both datasets had values “Don’t know” to indicate that a person was not sure. However, there very few rows with this value so we removed these rows.
- Education Level: Both datasets recorded the highest level of education completed by a person. However, these categories had slightly different names. For example, “University certificate, diploma or degree above the bach..” indicated that a person has a qualification above a bachelor’s degree. In the CES data there are Master’s degree as well as other degree above bachelor’s degree. Hence we put the education level in both datasets into three groups: “Above Bachelor”, “Below Bachelor” and “Bachelor”.
- Household size: Both the datasets recorded household size.

After finding the common variables, we removed rows with missing values in these columns for both the datasets. Then we renamed the columns of the CES data to match the names of the columns of the GSS data. This was done so that models constructed on the CES data could also be used on the GSS data for post-stratification. Then we took the common variables and constructed two datasets. One dataset containing the rows from the CES data and the other dataset containing the rows of the GSS data. Furthermore, in the first dataset, voting data was added to the CES variables. This is the outcome of interest

<Include a description of the important variables.>

Variable Name	Description
sex	Sex of the person
province	Province that a person is currently residing in
education	Level of education divided into three categories: ‘Above Bachelor’, ‘Below Bachelor’ and ‘Bachelor’
religion_importance	How important is their religion to a person: ‘Not important at all’, ‘Somewhat important’, ‘Very important’, ‘Not very important’ ‘Don’t know’
aboriginal	Is the person aboriginal
hh_size	How many people does a person share their living space with including themselves
age	Age of the person in years
vote	Party that the person will/might/has vote(d) for

<Include a description of the numerical summaries. Remember you can use `r` to use inline R code.>

*# Use this to create some plots. Should probably describe both the sample and population.*

<Include a clear description of the plot(s). I would recommend one paragraph for each plot.>

## Methods

<Include some text introducing the methodology, maybe restating the problem/goal of this analysis.>

## Model Specifics

<I will (incorrectly) be using a linear regression model to model the proportion of voters who will vote for Donald Trump. This is a naive model. I will only be using age, which is recorded as a numeric variable, to model the probability of voting for Donald Trump. The simple linear regression model I am using is:>

$$y = \beta_0 + \beta_1 x_{age} + \epsilon$$

<Where  $y$  represents the ....  $\beta_0$  represents....>

## Post-Stratification

<In order to estimate the proportion of voters....>

<To put math/LaTeX inline just use one set of dollar signs. Example:  $\hat{y}^{PS}$  >

*include.your.mathematical.model.here.if.you.have.some.math.to.show*

All analysis for this report was programmed using **R version 4.0.2**.

## Results

<Here you present your results. You may want to put them into a well formatted table. Be sure that there is some text describing the results.>

<Note: Alternatively you can use the `knitr::kable` function to create a well formatted table from your code. See here: <https://rmarkdown.rstudio.com/lesson-7.html>.>

<Remember you can use `r` to use inline R code.>

<Include an explanation/interpretation of the visualizations. Make sure to comment on the appropriateness of the assumptions/results.>

## Conclusions

<Here you should give a summary of the Hypotheses, Methods and Results>

<Highlight Key Results.>

<Talk about big picture.>

<Comment on any Weaknesses.>

<End with a concluding paragraph to wrap up the report.>

## Bibliography

1. Golemund, G. (2014, July 16) *Introduction to R Markdown*. RStudio. [https://rmarkdown.rstudio.com/articles\\_intro.html](https://rmarkdown.rstudio.com/articles_intro.html). (Last Accessed: January 15, 2021)
2. Dekking, F. M., et al. (2005) *A Modern Introduction to Probability and Statistics: Understanding why and how*. Springer Science & Business Media.
3. Allaire, J.J., et. el. *References: Introduction to R Markdown*. RStudio. <https://rmarkdown.rstudio.com/docs/>. (Last Accessed: January 15, 2021)