# Predicting 2019 Canadian Federal Election Results If All Eligible Individuals Voted

Maxwell Garrett

Decmeber 7th 2020

## Github URL

## Abstract

Data from the General Social Survey (2017) (Statistics Canada 2017) and Canadian Election Study (2019) Phone (Stephenson et al. 2020) datasets are used to build logistic regression models with post-stratification of the popular vote results for Canadian political parties. The popular vote for the 2019 federal Canadian election is predicted and compared with the actual recorded vote ("43rd General Election: Official Voting Results (Raw Data)" 2019).

## Keywords

Keywords: Post-stratification, logistic regression, Canadian election, voter turnout

## Introduction

It is important for politicians to know how they can best represent the interests of their constituents. Canadian election polling data provides important information on voting preferences of Canadians on important issues (Stephenson et al. 2020). It can be difficult to obtain large and representative samples of citizens' voting preferences. Often due to time and resource constraints, the sampling procedure results in certain subpopulations being under- or over-represented, with respect to the target population. Therefore, methods for using non-representative data to produce reliable voter models in an election are important.

In this report, we use the technique of post-stratification with logistic regression models for predicting the 2019 Canadian election outcome. This technique is used to make the sampled survey data representative of the population. We will use Statistics Canada's General Social Survey administered in 2017 as census data to understand the proportion of the population in subsets of the population. As well, we will use Canadian Election Study's 2019 phone dataset as survey data to understand the voting intentions of individuals in different subgroups of the population.

The two datasets mentioned above will be used to build counterfactual models predicting the potential Canadian election results if all eligible voters in Canada had voted in the 2019 election. We will be specifically modeling the popular vote. We are interested in investigating if the popular vote would significantly change in

the hypothetical scenario where all eligible voters participated in the election. In the methodology section, we describe the data, and the models of the election intentions. We then show our resulting election predictions based on the models' fits in the results section. Finally, we go over what our results indicate about the importance of increasing voter turnout.

## Methodology

In this section, we will first explain features of the datasets that we are using and then we will explain the popular vote proportion models we have developed.

### Data

The first dataset we discuss is the General Social Survey (2017) dataset. We will be using this dataset to calculate representative population proportions to use in post-stratification. We selected sex, place_birth_canada, province as variables of interest. This dataset's target population is all individuals above the age of 15 (inclusive) in Canada excluding individuals in the Yukon, Northwest Territories, and Nunavut (Beaupré 2020). As well, individuals in full-time institutions are excluded (Beaupré 2020).

Table 1: Characteristic table of selected variables from GSS data.

| characteristic | percent |
|---|---|
| Female | 54.39 |
| Born in Canada | 79.87 |
| Alberta | 8.36 |
| British Columbia | 12.25 |
| Manitoba | 5.77 |
| New Brunswick | 6.47 |
| Newfoundland and Labrador | 5.33 |
| Nova Scotia | 6.92 |
| Ontario | 27.28 |
| Prince Edward Island | 3.43 |
| Quebec | 18.59 |
| Saskatchewan | 5.61 |

The variable sex, an individual's biological sex, was selected for analysis as there is no missing data for this variable and this variable has similar categories in the CES dataset allowing for straightforward post-stratification. Table 1 tells us that 54% of the sample is female, the rest male. The variable place_birth_canada, indicating if an individual was born in Canada, was selected as we thought that an individual immigrating to Canada may have an impact on their voting preference. As well, this variable has very few missing values which means we keep almost all of the information from the dataset. Some missing values were replaced if the respondent had indicated the macroregion which they were born in. In this case, if the macroregion did not include Canada then they were marked as 'Born outside of Canada'. Approximately 80% of the respondents indicate that they are born in Canada according to Table 1. The last variable selected is province which is the current province that the respondent resides in. We selected this as we thought that location where an individual lives may play into their voting preference. As well, this variable is missing zero values and is similar to the province categories provided in the CES dataset. A drawback of this variable is that it unfortunately does not include territories while the CES dataset does, this means that we do not include territories in our analysis. In Table 1, we see the proportion of individuals reporting to reside in each province.

The second dataset we discuss is the Canadian Election Study (2019) phone survey (Stephenson et al. 2020).

We will be using this dataset to form our model on voting preference. This survey was performed over phone through stratified random sampling of Canadian phone numbers (Stephenson et al. 2020). The phone calls were performed the day after the election (Stephenson et al. 2020). The variables selected from this dataset were sex, place_birth_canada, province, and voting_pref. This dataset only includes individuals eighteen years old or older (Stephenson et al. 2020).

Table 2: Characteristic table of selected variables from CES data.

| characteristic | percent |
|---|---|
| Female | 42.52 |
| Born in Canada | 85.20 |
| Alberta | 7.21 |
| British Columbia | 20.54 |
| Manitoba | 6.77 |
| New Brunswick | 4.84 |
| Newfoundland and Labrador | 4.62 |
| Nova Scotia | 4.93 |
| Ontario | 20.70 |
| Prince Edward Island | 4.87 |
| Quebec | 18.73 |
| Saskatchewan | 6.77 |

The variable sex was created based on the variable q3, gender, in the dataset. This variable was reduced to sex as the method for determining gender was not sufficient, because it only involved the interviewer deducing the respondent's gender based on voice pitch (Stephenson et al. 2020). This variable was chosen as it is easily modified to fit with the values found for sex in the GSS dataset. In table 2 we see that approximately 43% of the respondents were identified to be female. The variable place_birth_canada was created based on the variable q64, the country of birth, by reducing the variable to be a binary indicating only if an individual was born in Canada or not. This was done to match the GSS dataset and allow easier computations later on. It was found that 85% of the respondents were born in Canada according to table 2. The province variable was created based on the variable q4, the province where they reside. Respondents who indicated living in Nunavut, Northwest Territories, and Yukon were removed as the census dataset does not include those regions. In table 2, we see the percent of the total respondents surveyed that reside in each province. The variable voting_pref indicates the voting preference of the respondent. It was created based on q11 and q12 which describe the voting preference of the respondent, q12 data was used if q11 data was missing. The variable has popular party names indicated and then all other responses indicated as other.

## Model

We are looking to predict the popular vote for the 2019 Canadian election had all potential voters in Canada voted. This will be modeled using multiple logistic regression models, one for each party in the data and one for spoiled ballots. For each model, we have a binary response variable created which indicates 1 if the respondent plans to vote for that option and 0 otherwise.

Our logistic regression models will model the log-odds of voting for a party. For example, one model will be modeling the log-odds of voting for the Liberal Party of Canada. There will be a model for voting Liberal, Conservative, NDP, BQ, Green Party, People's Party, other, and spoiling the ballot. The model formula is the following:

$$log(\frac{p}{1-p}) = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \beta_3 x_3 + \beta_4 x_4 + \beta_5 x_5 + ... + \beta_{14} x_{14} + \epsilon$$

where $x_1$ is sex equaling male, $x_2$ is sex equaling female, $x_3$ is place_birth_canada equaling 'Born in Canada', $x_4$ is place_birth_Canada equaling 'Born outside of Canada', and $x_5$ to $x_14$ are the provinces of Canada

excluding Nunavut, Yukon, and Northwest Territories. As well, $p$ represents the percent of votes given the sex, birth place, and province being for a party. Beta values indicate the change in log-odds expected if a given category is true except $\beta_0$ which is the intercept of the model.

The model was selected as the predictor variables used are available in both datasets. As well, it is plausible that the predictor variables are associated with voting preference as they subset the population into large groups that share some common characteristics.

**Post-stratification**

We use the census dataset (GSS) to find the number of respondents in each bin of the population. There are forty bins created total, each bin based on sex, birth place, and province. For each political party's model, we predict the proportion of a given bin that would vote for the party. We multiply this proportion by the number of respondents in the bin and divide by the total number of respondents to GSS. We then add up all the proportion estimates for bins, this is our estimate for the proportion of the population voting for a political party. As well, we use the standard error for each bin estimate from the model to calculate a 95% confidence interval of the proportion of voters for a given party.

# Results

Table 3: Predicted percent of popular vote using logistic regression models.

| Voting Preference | Percent of The Popular Vote | Lower Limit | Upper Limit |
|---|---|---|---|
| Liberal | 35.70 | 30.15 | 41.26 |
| Conservative | 32.00 | 26.74 | 37.26 |
| NDP | 15.34 | 11.14 | 19.53 |
| BQ | 3.58 | 2.66 | 4.51 |
| Green Party | 10.54 | 7.04 | 14.03 |
| People's Party | 1.56 | 0.23 | 2.90 |
| Other | 1.19 | -0.03 | 2.42 |
| Spoil | 3.58 | 2.66 | 4.51 |

In table 3, we have the predicted percent of the popular vote that each party would receive had everyone voted according to our models. As well, we have the lower and upper bounds of the 95% confidence intervals around these predicted popular votes. For the liberal party, we expect the percent of the popular vote to be 35.70%. For the Conservative party, we expect the percent of the popular vote to be 32.00%. For the NDP, we expect the percent of the popular vote to be 15.34%. For the BQ, we expect the percent of the popular vote to be 3.58%. For the Green Party, we expect the percent of the popular vote to be 10.54%. For the People's Party, we expect the percent of the popular vote to be 1.56%. We expect that 1.19% of people vote for any other party option not listed. Finally, we expect that 3.58% of voters spoil their ballot.
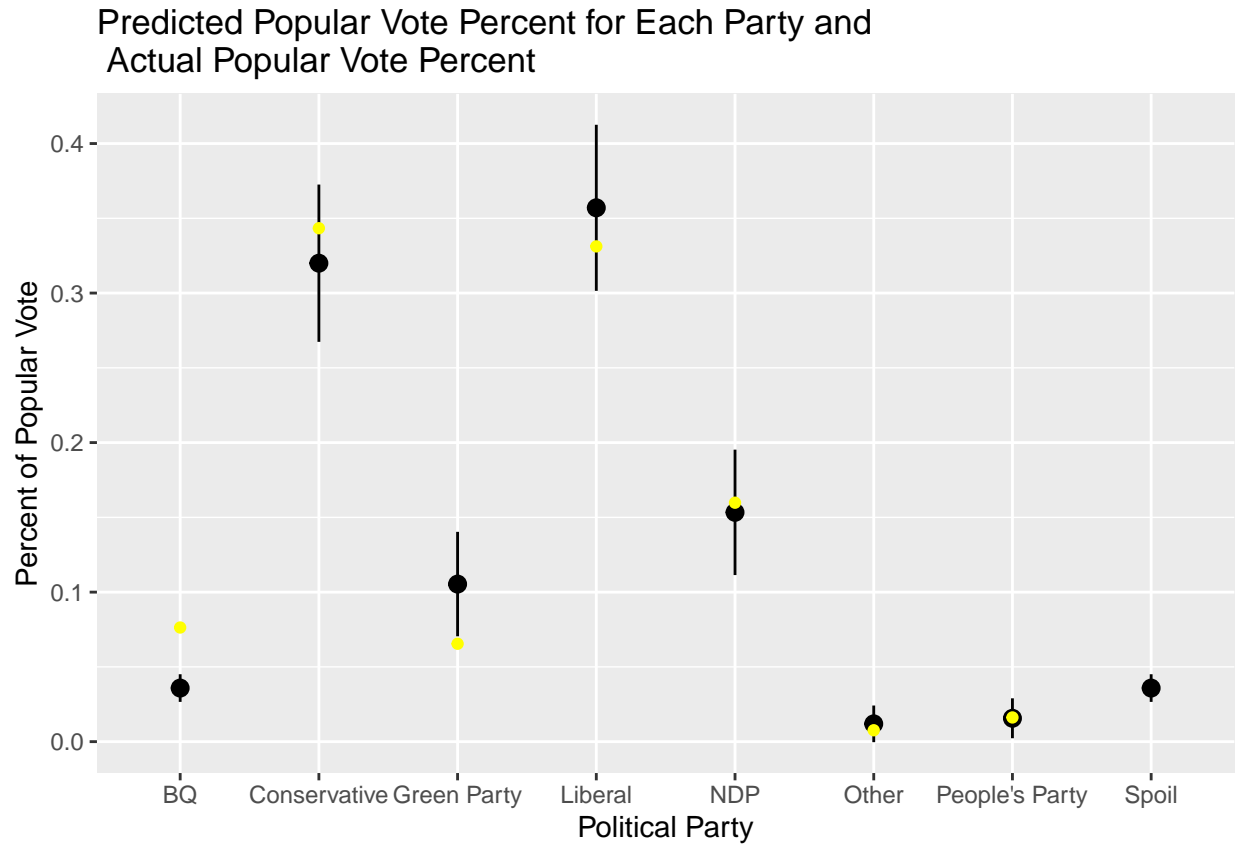
Table 4: Actual 2019 election results provided by Stats Canada.

| Party Voted For | Percent of The Popular Vote |
|---|---|
| Liberal | 33.12 |
| Conservative | 34.34 |
| NDP | 15.98 |
| BQ | 7.63 |
| Green Party | 6.55 |

| Party Voted For | Percent of The Popular Vote |
| --- | --- |
| People's Party | 1.62 |
| Other | 0.76 |
| Spoil | NA |

Table 4 provides the actual popular vote results from the 2019 Canadian Election ("43rd General Election: Official Voting Results (Raw Data)" 2019). We lack data for the number of spoiled ballots in this data.

Figure 1: Plot of predicted popular vote confidence intervals, and the actual popular vote collected in the 2019 election.



In figure 1, we have the predicted popular vote percentages plotted in black. As well, we have the 95% confidence interval around the predicted value marked with the black line. The actual popular vote percentages from the election results are marked in yellow points. We have no data for spoiled ballots therefore there is no yellow point for that category. We see that the proportion of votes for the Conservative, Liberal, NDP, Other, and People's Party are quite similar to our predictions. All these parties actual results fall within the 95% confidence interval. We notice that BQ received a higher proportion of votes than what was expected if all eligible voters voted. As well, the Green Party was expected to receive a higher proportion of overall votes than what was actually received had all eligible voted.

# Discussion

## Summary

In summary, we created eight logistic regression models predicting the proportion of voters for each party for each population group. We then used the technique of post-stratification to estimate the proportion of voters for each party across Canada. Finally, we compared the predicted election result confidence intervals with the actual election results received in 2019 in a plot.

## Conclusions

In our results, we see the scenario where every eligible citizen in Canada excluding Nunavut, Northwest Territories, and Yukon. In figure 1, we see the plotted predicted results and actual results from the 2019 federal election. There is a noticeable difference in the predicted results from the actual results for the Bloc Quebecois Party and the Green Party. For the BQ Party, we predict less proportional support had everyone voted while the Green Party we predict more proportional support had everyone voted. Other parties seem to have quite similar results. Based on the variance in results for the BQ Party and Green Party though we conclude that there is an impact of having everyone vote on the election results.

As well, based on these results we would predict a Liberal win of the popular vote, as the predicted popular vote is higher than all others. This differs from the actual popular vote winner the Conservative Party. This as well supports our conclusion that there is a noticeable impact of having evreyone vote.

In the 2019 election, there was a voter turnout of 67% of Canadians ("Appendix – Report on the 43rd General Election of October 21, 2019" 2019). In our models, we are predicting a turnout of 100%. This research area is of importantance as if turnout is low there may be sub-populations underrepresented by the parties elected. As we want representative democracies, we want voter turnout as close to 100% as possible.

## Weaknesses & Next Steps

The first important weakness to consider is that our census dataset (GSS) did not include provinces such as Nunavut, Northwest Territories and Yukon (Beaupré 2020). As well, the GSS dataset does not include institutionalized individuals (Beaupré 2020). Both these groups have many people who are eligible voters. By excluding these groups, we do not get models that fully represent the population of eligible voters. A future step would involve finding a dataset that includes these groups. This would greatly expand the ability of the models to predict the voting outcomes for the country.

The second important weakness to consider is that we are looking at the popular vote, but this does not necessarily provide information on the actual election outcome due to the voting system. The voting system involves each riding in Canada electing a party representative ("FAQs - General Questions" n.d.). This means that our popular vote predictions do not reflect who would be elected differently due to all eligible voters voting. A future step would involve finding a dataset that includes voting preferences with riding information on the respondent. This would allow us to accurately predict election outcomes for each riding and therefore more accurately predicting the election result.

The third important weakness is the handling of sex and gender in both datasets. The CES dataset has the phone operator record gender, but only ask for the gender of an individual when unsure (Stephenson et al. 2020). This is not an accurate way to record sex or gender as voice is not a reliable indicator of this. An improvement to this problem would be finding a dataset that collects gender through asking directly in the interview. This way an accurate answer by the respondent is given.

# References

"43rd General Election: Official Voting Results (Raw Data)." 2019. Elections Canada. https://www.elections.ca/res/rep/off/ovr2019app/51/data_donnees/table_tableau08.csv.

Alexander, Rohan, and Sam Caetano. 2020. "GSS Cleaning Code for 2017 Data."

"Appendix – Report on the 43rd General Election of October 21, 2019." 2019. Elections Canada. 2019. https://www.elections.ca/content.aspx?section=res&dir=rep/off/sta_ge43&document=app1&lang=e.

Beaupré, Pascale. 2020. *Cycle 31 : Families Public Use Microdata File Documentation and User's Guide.* Statistics Canada. http://www.chass.utoronto.ca/.

"FAQs - General Questions." n.d. Elections Canada. Accessed 2020. https://www.elections.ca/content.aspx?section=vot&dir=faq&document=faqgen&lang=e#gen1.

Henry, Lionel, and Hadley Wickham. 2020. *Rlang: Functions for Base Types and Core R and 'Tidyverse' Features.* https://CRAN.R-project.org/package=rlang.

Hodgetts, Paul A., and Rohan Alexander. 2020. *CesR: Access the Ces Datasets a Little Easier.*

Lumley, Thomas. 2004. "Analysis of Complex Survey Samples." *Journal of Statistical Software* 9 (1): 1–19.

———. 2010. *Complex Surveys: A Guide to Analysis Using R: A Guide to Analysis Using R.* John Wiley; Sons.

———. 2020. "Survey: Analysis of Complex Survey Samples."

Statistics Canada. 2017. "2017 General Social Survey: Families Cycle 31." http://www.chass.utoronto.ca/.

Stephenson, Laura B, Allison Harell, Daniel Rubenson, and Peter John Loewen. 2020. "2019 Canadian Election Study - Phone Survey." Harvard Dataverse. https://doi.org/10.7910/DVN/8RHLG1.

Wickham, Hadley, Mara Averick, Jennifer Bryan, Winston Chang, Lucy D'Agostino McGowan, Romain François, Garrett Grolemund, et al. 2019. "Welcome to the tidyverse." *Journal of Open Source Software* 4 (43): 1686. https://doi.org/10.21105/joss.01686.

Xie, Yihui. 2014. "Knitr: A Comprehensive Tool for Reproducible Research in R." In *Implementing Reproducible Computational Research*, edited by Victoria Stodden, Friedrich Leisch, and Roger D. Peng. Chapman; Hall/CRC. http://www.crcpress.com/product/isbn/9781466561595.

———. 2015. *Dynamic Documents with R and Knitr.* 2nd ed. Boca Raton, Florida: Chapman; Hall/CRC. https://yihui.org/knitr/.

———. 2020. *Knitr: A General-Purpose Package for Dynamic Report Generation in R.* https://yihui.org/knitr/.