

Predicting Self-Rated Mental Health Based on Demographic and Family Traits

James Bao, Alan Chen, Xinyi Zhang, Zidong Yang

10/19/2020

The code used to generate this RMarkdown file can be found at https://github.com/cindy Zhang99/sta304_ps3/blob/main/sta304_ps3.Rmd.

Abstract

- An abstract is included and appropriately pitched to a general audience.
- The abstract answers: what was done, what was found, and why this matters (all at a high level).
- If your abstract is longer than four sentences then you need to think a lot about whether it is too long. It may be fine (there are always exceptions) but you should probably have a good reason.

Introduction

- The introduction is self-contained and tells a reader everything they need to know, including putting it into a broader context.
- Your introduction should provide a bit of broader context to motivate the reader, as well as providing a bit more detail about what you're interested in, what you did, what you found, why it's important, etc.
- A reader should be able to read only your introduction and have a good idea about the research that you carried out.
- It would be rare that you would have tables or figures in your introduction (again there are always exceptions but think deeply about whether yours is one).
- It must outline the structure of the report.

Data

The dataset we used in our modeling is the 2017 General Social Survey (Family cycle). The following sections will discuss how the data was collected, what the key features of the dataset are, and what the data looks like.

Data Collection

From February 1, 2017 to November 30, 2017, Statistics Canada gathered data on the Canadian family unit by conducting voluntary telephone interviews. Their target population was all non-institutionalized individuals living in Canada, aged 15 or older. Cross-sectional sampling was conducted in a two-stage design.

The stratified simple random sampling method was used in the first stage. Here, the sampling frame consisted of telephone numbers from the Census grouped as households using data from Statistic Canada’s dwelling frame. Strata were formed at the census metropolitan area (CMA) level and at the province level (i.e., large CMAs formed their own stratum, smaller CMAs were grouped together, and the non-CMA regions of each province were grouped together), forming a total of 27 non-overlapping strata. Finally, households were sampled randomly from each stratum such that the number sampled units from each stratum corresponded to the population sizes of each stratum. To reiterate, the sampled population for this first stage was the chosen households from each stratum.

The stratified simple random sampling method was also used in the second stage. Here, the sampling frame was a list of household members, aged 15 and older, from the households selected in the first stage. Then, one individual was randomly selected from each household, forming the sampled population. Approximately 43,000 individuals were contacted to participate in the survey.

Overall, the surveying method using two-stage simple random stratified sampling is effective in generating a sample that is geographically representative of the Canadian population. In addition to estimates about the Canadian population at large, the stratified sampling method also allows estimates to be made about subpopulations (at the province level).

Statistics Canada reported that the non-response rate was 47.6%. To reduce the effects of non-response bias, survey responses were adjusted based on the demographic characteristics of households that were non-responsive (by pulling their information from the 2016 Census). This ensures that the discrepancy between the target population and survey responses resulting from non-response is minimized. Furthermore, for the Family cycle of the GSS, responses were also adjusted for income and household size to make more accurate survey estimates for the variables of interest.

Statistics Canada did not disclose the true cost of conducting the survey but we can make some speculations based on the available information about their field work methodology. Surveying was conducted using Computer Assisted Telephone Interviewing (CATI) wherein interviewers read aloud the computerized questionnaire and immediately record the respondent’s answers. Although this allows for a reduction in costs compared to traditional in-person surveying, labor costs still include time spent computerizing the survey, training interviewers, and having interviewers administer the questionnaire. Other labor costs include designing the questionnaire and surveying methodology as well as conducting quality control (data consistency was checked by the CATI system during surveying and unresolved inconsistencies were handled afterwards by support staff). Non-labor costs likely included paying for equipment, phone service, offices, and so forth. Again, although we don’t have exact costs, we can conclude that the time and costs associated with conducting the GSS is a clear reason why it is only administered once a year.

Per the report on the 2017 GSS from Statistics Canada, extensive research and testing was conducted when designing the questionnaire. Consequently, a major strength of the questionnaire is that it contains focused questions that comprehensively and extensively capture the subject of interest (the Canadian family). Upon reading through the questionnaire made available by Statistics Canada, the wording of each question is precise and clear, leaving little room for ambiguity. Additionally, another strength of the survey is that a vast majority of questions were objective (dates, events, counts) removing potential response biases that occur with subjective questions. (Not all questions were objective however, in fact the variable of interest we will model in subsequent sections consists of subjective responses.) On the other hand, because of the specificity of the questions, the survey is very long with several dozens sections and several questions per section. Furthermore, as a result of the large scope of the target population, many questions in the survey did not apply to a large majority of respondents (e.g., number of grandchildren, questions about additional marriages, etc.). The data collected is also incomplete because participants were given the option to refuse to answer or answer “I don’t know” to each question since participation was voluntary.

Overall, the surveying methodology and distributed questionnaire were carefully designed in the interest of collecting accurate, representative data wherever possible.

Data Characteristics

The full dataset of responses to the 2017 General Social Survey (Family cycle) contains 20,602 observations for over 400 variables relating to the Canadian family. A large reason for our choice to use this dataset is because it is the most recent GSS cycle available for modeling. Other benefits of this dataset have been previously touched upon in the previous section. Namely, the data was checked for consistency in real time by the CATI system (as well as by survey support staff) so there is a certain measure of accuracy that other survey results lack. Additionally, the stratified simple random sampling method used to distribute the survey suggests that the results are representative of the Canadian population to some degree (in the geographical sense at the very least). A major weakness of the data is that it is not complete because of the voluntary nature of the surveying.

In the interest of space, we will only discuss the variables in the dataset that are relevant to our model. The variable we aim to predict is `self_rated_mental_health` while the factors that we chose to inform this prediction are age, sex, marital_status, and self_rated_health. We chose these factors based on the demographic information mentioned in mental health statistics (age, sex, and health) and based on what we suspected might contribute to mental health in the context of family composition (marital status). More explicitly, here is what information the chosen variables in the dataset represent:

- age (agedc in the original dataset): the exact age of the respondent (in decimals) at the time of the survey
- sex (sex): sex of the respondent, the options being “Male” or “Female”
- marital_status (marstat): marital status of the respondent, the options being “Single, never married”, “Married”, “Living common-law”, “Separated” (but still legally married), “Divorced”, or “Widowed”
- self_rated_health (srh_110): self-rated health, the options being “Excellent”, “Very good”, “Good”, “Fair”, and “Poor”
- self_rated_mental_health (srh_115): self-rated mental health, the options being “Excellent”, “Very good”, “Good”, “Fair”, and “Poor”

For age, there is a similar variable in the original dataset that uses only natural numbers (agec), however, we chose to use agedc (and renamed it to “age”) in the interest of accuracy. There are many variables related to marriage in the original dataset (totunc: total number of marriage and common-law unions, nmarevrc: number of marriages the respondent has had, etc.) but they don’t capture the same scope of information as marital_status does (for example, being divorced or widowed is not reflected in those variables). Consequently, we chose marital status as opposed to the other available variables related to marriage. For the other three variables we use (sex, self_rated_health or srh_110 in the original dataset, and self_rated_mental_health or srh_115 in the original dataset), there are no similar equivalents.

Data Visualization

```
library(tidyverse) # for data manipulation and plots
```

```
## -- Attaching packages -----  
  
## v ggplot2 3.3.2      v purrr  0.3.4  
## v tibble  3.0.3      v dplyr  1.0.2  
## v tidyr   1.1.2      v stringr 1.4.0  
## v readr   1.3.1      v forcats 0.5.0  
  
## -- Conflicts -----  
## x dplyr::filter() masks stats::filter()  
## x dplyr::lag()    masks stats::lag()
```

```
# load the csv, can be downloaded via utoronto
poll <- as_tibble(data.frame(read_csv("gss_cleaned.csv")))
```

```
## Parsed with column specification:
## cols(
##   .default = col_character(),
##   caseid = col_double(),
##   age = col_double(),
##   age_first_child = col_double(),
##   age_youngest_child_under_6 = col_double(),
##   total_children = col_double(),
##   age_start_relationship = col_double(),
##   age_at_first_marriage = col_double(),
##   age_at_first_birth = col_double(),
##   distance_between_houses = col_double(),
##   age_youngest_child_returned_work = col_double(),
##   feelings_life = col_double(),
##   hh_size = col_double(),
##   number_total_children_intention = col_double(),
##   number_marriages = col_double(),
##   fin_supp_child_supp = col_double(),
##   fin_supp_child_exp = col_double(),
##   fin_supp_lump = col_double(),
##   fin_supp_other = col_double(),
##   is_male = col_double(),
##   main_activity = col_logical()
##   # ... with 1 more columns
## )
```

```
## See spec(...) for full column specifications.
```

```
# choose pertinent variables
poll <- poll %>% select(age, sex, marital_status, self Rated_health,
                      total_children, self Rated_mental_health)

# clean up the data
poll <- poll[!grepl("Don't know", poll$self Rated_mental_health),]

# poll <- head(poll, 1000)
```

```
nrow(poll)
```

```
## [1] 20545
```

```
# table(poll$age)
table(poll$sex)
```

```
##
## Female    Male
## 11177     9368
```

```
table(poll$marital_status)
```

```
##
##           Divorced      Living common-law           Married
##           1759           2073           9481
##           Separated Single, never married           Widowed
##           641           4698           1887
```

```
table(poll$self_rated_mental_health)
```

```
##
## Excellent      Fair      Good      Poor Very good
##      6080      1296      5813      326      6924
```

Model

```
# load the csv, can be downloaded via utoronto
poll <- as_tibble(data.frame(read_csv("gss_cleaned.csv")))
```

```
## Parsed with column specification:
## cols(
##   .default = col_character(),
##   caseid = col_double(),
##   age = col_double(),
##   age_first_child = col_double(),
##   age_youngest_child_under_6 = col_double(),
##   total_children = col_double(),
##   age_start_relationship = col_double(),
##   age_at_first_marriage = col_double(),
##   age_at_first_birth = col_double(),
##   distance_between_houses = col_double(),
##   age_youngest_child_returned_work = col_double(),
##   feelings_life = col_double(),
##   hh_size = col_double(),
##   number_total_children_intention = col_double(),
##   number_marriages = col_double(),
##   fin_supp_child_supp = col_double(),
##   fin_supp_child_exp = col_double(),
##   fin_supp_lump = col_double(),
##   fin_supp_other = col_double(),
##   is_male = col_double(),
##   main_activity = col_logical()
##   # ... with 1 more columns
## )
```

```
## See spec(...) for full column specifications.
```

```
# choose pertinent variables
```

```
poll <- poll %>% select(age, sex, marital_status, self_rated_health,
                        self_rated_mental_health)
```

```
# clean up the data
```

```
poll$self_rated_mental_health %>% table()
```

```
## .
```

```
## Don't know    Excellent      Fair      Good      Poor    Very good
##           57      6080      1296      5813      326      6924
```

```
poll<-poll[!grepl("Don't know", poll$self_rated_mental_health),]
```

```
poll$self_rated_mental_health %>% table()
```

```
## .
```

```
## Excellent      Fair      Good      Poor    Very good
##           6080      1296      5813      326      6924
```

```
# poll <- head(poll, 1000)
```

```
model <- nnet::multinom(self_rated_mental_health ~ age + sex + marital_status + self_rated_health,
                        data = poll)
```

```
## # weights: 70 (52 variable)
```

```
## initial value 32875.988237
```

```
## iter 10 value 23611.931535
```

```
## iter 20 value 23441.438461
```

```
## iter 30 value 23049.583322
```

```
## iter 40 value 22914.075187
```

```
## iter 50 value 22822.824806
```

```
## iter 60 value 22775.879430
```

```
## final value 22775.875150
```

```
## converged
```

```
summary(model)
```

```
## Call:
```

```
## nnet::multinom(formula = self_rated_mental_health ~ age + sex +
## marital_status + self_rated_health, data = poll)
```

```
##
```

```
## Coefficients:
```

```
## (Intercept)      age      sexMale marital_statusLiving common-law
## Fair      1.0217425 -0.032798399 -0.3224026      -0.16949105
## Good      1.5874823 -0.011040378 -0.2400994      0.09958731
## Poor      1.8561011 -0.044664300 -0.2137392     -0.62177910
## Very good  0.9413021 -0.005288015 -0.2162900      0.04000497
## marital_statusMarried marital_statusSeparated
## Fair      -0.29251444      0.7299606
## Good      -0.01207268      0.4443427
## Poor      -0.92293929      0.5025271
```

```

## Very good          -0.03471334          0.1393385
##          marital_statusSingle, never married marital_statusWidowed
## Fair              0.224221668          0.17670603
## Good              0.130575216          0.22994484
## Poor              0.033859427          0.01458883
## Very good         -0.002477955          0.17862358
##          selfRatedHealthExcellent selfRatedHealthFair
## Fair              -3.360131          1.5800364
## Good              -3.050227          0.1066989
## Poor              -5.162473          -0.2307758
## Very good         -1.813127          -0.3809237
##          selfRatedHealthGood selfRatedHealthPoor
## Fair              -0.1643287          1.9008484
## Good              0.1545703          0.1368938
## Poor              -1.8749204          1.6135608
## Very good         -0.1348050          -0.1311144
##          selfRatedHealthVery good
## Fair              -1.6225589
## Good              -1.2432679
## Poor              -3.5250590
## Very good         0.3213029
##
## Std. Errors:
##          (Intercept)          age          sexMale marital_statusLiving common-law
## Fair          0.4695617 0.002429496 0.06897229          0.16272558
## Good          0.3101351 0.001459802 0.04227380          0.09946604
## Poor          0.5122282 0.004525478 0.12350622          0.28565061
## Very good     0.3325459 0.001338815 0.03893398          0.09150798
##          marital_statusMarried marital_statusSeparated
## Fair          0.12358331          0.1975624
## Good          0.07773531          0.1410396
## Poor          0.20958816          0.3020509
## Very good     0.07207928          0.1364599
##          marital_statusSingle, never married marital_statusWidowed
## Fair          0.13684546          0.15750158
## Good          0.08986249          0.10021966
## Poor          0.21528134          0.25746406
## Very good     0.08379589          0.09385189
##          selfRatedHealthExcellent selfRatedHealthFair
## Fair          0.4574703          0.4452440
## Good          0.2951180          0.2967677
## Poor          0.5107378          0.4420195
## Very good     0.3189361          0.3248795
##          selfRatedHealthGood selfRatedHealthPoor
## Fair          0.4437937          0.4540376
## Good          0.2921365          0.3112635
## Poor          0.4407436          0.4439970
## Very good     0.3188215          0.3405639
##          selfRatedHealthVery good
## Fair          0.4472054
## Good          0.2925541
## Poor          0.4643552
## Very good     0.3178609
##

```

```
## Residual Deviance: 45551.75
## AIC: 45655.75
```

```
head(fitted(model))
```

```
##      Excellent      Fair      Good      Poor Very good
## 1 0.6719425 0.014404347 0.09910454 0.0024202042 0.2121284
## 2 0.1951443 0.046536541 0.49248128 0.0062707489 0.2595671
## 3 0.2389554 0.012145612 0.16505036 0.0010447006 0.5828039
## 4 0.2601140 0.007720795 0.14990947 0.0005466639 0.5817091
## 5 0.1437139 0.082680958 0.52334153 0.0175120007 0.2327516
## 6 0.7074104 0.006452120 0.08074126 0.0006178245 0.2047784
```

```
input <- data.frame(selfRatedHealth = c("Excellent"), age = c(21.5), sex = c("Male"), marital_status = c("Married"))
predict(model, newdata = input, "probs")
```

```
##      Excellent      Fair      Good      Poor      Very good
## 0.658524299 0.028453356 0.107809344 0.007717323 0.197495678
```

The purpose of the model is to predict a person's self-rated mental health based on other characteristics presented in the dataset such as demographics, lifestyle, socioeconomic status, etc. Since the self-rated mental health data is represented in categories, the task at its core is a classification problem.

Some models that were considered were linear regression, naive bayes, binary logistic regression, and multinomial logistic regression. To begin, linear regression is not suitable because it is often difficult to find an accurate linear relationship between predictors and categories, not to mention the fact that linear regression is more suitable when the dependent variable is continuous. Naive Bayes would be a viable option since it is able to handle classification of more than two categories using joint probability and Bayes' Theorem. However, Naive Bayes only works well under the assumption that the explanatory variables are independent, but this is often not the case. Given the context of the data, it is highly anticipated that the characteristics of a person are correlated in some way or another.

In addition, generative models (eg. Naive Bayes) has a higher asymptotic error than discriminative models (eg. logistic regression), but it approaches the asymptotic error faster. In other words, discriminative models tend to perform better given large enough data and generative models will perform better on less data as it learns faster. Since the dataset is large, choosing a discriminative model would be more appropriate for this task. As a result, logistic regression is the model.

It is also important to note that the dependent variable has more than two categories. One way to handle this is to group multiple categories together until there are only two categories left which can be used for binary logistic regression. While this simplifies the complexity of implementing the model itself, there will be a loss in information from merging classes and it alters the original research task. The option to handle multi-classification is to use multinomial logistic regression which is extended from binary logistic regression.

Results

Discussion

Predicting mental health in the context of the nuclear family. Does being married have a positive or negative impact? ## Weaknesses ## Next Steps

Appendix

References

Interview method/survey size: <https://www.statcan.gc.ca/eng/survey/household/4501> Detailed information about GSS 2017: <https://www23.statcan.gc.ca/imdb/p2SV.pl?Function=getSurvey&Id=335816> Questionnaire: https://www23.statcan.gc.ca/imdb/p3Instr.pl?Function=assembleInstr&lang=en&Item_Id=335815#qb345205 Mental health statistics: <https://www.camh.ca/en/Driving-Change/The-Crisis-is-Real/Mental-Health-Statistics> Discriminative vs Generative Classifiers (Naive Bayes vs logistic regression): <https://ai.stanford.edu/~ang/papers/nips01-discriminativegenerative.pdf>