

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

James Bao, Alan Chen, Xinyi Zhang, Rose

10/19/2020

Abstract

Predicting mental health in the context of the nuclear family. Does having children have a positive (fulfilled) or negative impact (stressed)? Does being married have a positive or negative impact?

Introduction

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 strata, 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 strata such that the number sampled units from each strata corresponded to the population sizes of each strata. To reiterate, the sampled population for this first stage was the chosen households from each strata.

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 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 are what the variables in the dataset represent:

- `age`: Age of respondent with decimal at time of the survey interview

- sex: Sex of respondent
- marital_status: Marital status of the respondent
- self Rated health: Self rated health
- self Rated mental health: Self rated mental health

Data Visualization

-what does the data look like -plot the raw data

```
library(tidyverse)
```

```
## -- Attaching packages ----- tidyverse 1.3.0 --
```

```
## 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 ----- tidyverse_conflicts() --
## x dplyr::filter() masks stats::filter()
## x dplyr::lag()     masks stats::lag()
```

```
data <- 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.
```

```
data <- data[, c("age", "sex", "education", "religion_has_affiliation", "marital_status", "total_children")]
```

Model

Bayes' Theorem for Naive Bayes Classifier: $P(c|x) = \frac{P(x|c)P(c)}{P(x)}$

```
install.packages('tidyverse')
install.packages("sjstats")
install.packages("ROCR")
install.packages("brms")
install.packages("modelr")
install.packages("nnet")
install.packages("tidybayes")
require(nnet)
library(tidyverse) # for data manipulation and plots
```

```
# using 12 rows of dummy data until we get the real data
```

```
data <- tibble(
  sex = c("M", "F", "M", "M", "F", "F", "M", "F", "M", "F", "M"),
  martial_status = c("Married", "Divorced", "Single, never married", "Living common-law", "Widowed", "S",
  age = c("under 18", "18-34", "35-50", "50-70", "over 70", "under 18", "18-34", "35-50", "under 18", "1",
  education = c("Bachelor's", "High school", "College", "Trade Cert", "University degree above Bachelor",
  has_regilious_afflication = c("Y", "Y", "N", "Y", "Y", "N", "Y", "Y", "N", "Y", "Y", "N"),
  has_children = c("Y", "N", "Y", "N", "Y", "N", "Y", "N", "N", "N", "N"),
  selfrated_mental_health = c("Fair", "Good", "Very Good", "Excellent", "Poor", "Fair", "Good", "Very G
)
data
```

```
## # A tibble: 12 x 7
##   sex   martial_status age   education has_regilious_a~ has_children
##   <chr> <chr>          <chr> <chr>      <chr>          <chr>
## 1 M     Married        unde~ Bachelor~ Y             Y
## 2 F     Divorced        18-34 High sch~ Y             N
## 3 M     Single, never~ 35-50 College  N             Y
## 4 M     Living common~ 50-70 Trade Ce~ Y             N
## 5 F     Widowed          over~ Universi~ Y             Y
## 6 F     Separated        unde~ Universi~ N             N
## 7 F     Married          18-34 Bachelor~ Y             Y
## 8 M     Divorced        35-50 High sch~ Y             N
## 9 F     Single, never~ unde~ College  N             Y
## 10 M    Married          18-34 Trade Ce~ Y             N
## 11 F    Divorced        35-50 Bachelor~ Y             N
## 12 M    Single, never~ 18-34 High sch~ N             N
## # ... with 1 more variable: selfrated_mental_health <chr>
```

```
# add 1 after all x-y combinations to avoid zero frequency problem
```

```
# count y
poor <- data %>% filter(selfrated_mental_health=="Poor") %>% tally() + 1
fair <- data %>% filter(selfrated_mental_health=="Fair") %>% tally() + 1
```

```

good <- data %>% filter(selfrated_mental_health=="Good") %>% tally() + 1
vgood <- data %>% filter(selfrated_mental_health=="Very Good") %>% tally() + 1
excellent <- data %>% filter(selfrated_mental_health=="Excellent") %>% tally() + 1

total_mental <- poor + fair + good + vgood + excellent

# count sex given y
male_poor <- data %>% filter(sex=="M" & selfrated_mental_health=="Poor") %>% tally() + 1
male_fair <- data %>% filter(sex=="M" & selfrated_mental_health=="Fair") %>% tally() + 1
male_good <- data %>% filter(sex=="M" & selfrated_mental_health=="Good") %>% tally() + 1
male_vgood <- data %>% filter(sex=="M" & selfrated_mental_health=="Very Good") %>% tally() + 1
male_excellent <- data %>% filter(sex=="M" & selfrated_mental_health=="Excellent") %>% tally() + 1
female_poor <- data %>% filter(sex=="F" & selfrated_mental_health=="Poor") %>% tally() + 1
female_fair <- data %>% filter(sex=="F" & selfrated_mental_health=="Fair") %>% tally() + 1
female_good <- data %>% filter(sex=="F" & selfrated_mental_health=="Good") %>% tally() + 1
female_vgood <- data %>% filter(sex=="F" & selfrated_mental_health=="Very Good") %>% tally() + 1
female_excellent <- data %>% filter(sex=="F" & selfrated_mental_health=="Excellent") %>% tally() + 1

total_male <- male_poor + male_fair + male_good + male_vgood + male_excellent
total_female <- female_poor + female_fair + female_good + female_vgood + female_excellent

# count has_children given y
has_children_poor <- data %>% filter(has_children=="Y" & selfrated_mental_health=="Poor") %>% tally() + 1
has_children_fair <- data %>% filter(has_children=="Y" & selfrated_mental_health=="Fair") %>% tally() + 1
has_children_good <- data %>% filter(has_children=="Y" & selfrated_mental_health=="Good") %>% tally() + 1
has_children_vgood <- data %>% filter(has_children=="Y" & selfrated_mental_health=="Very Good") %>% tally() + 1
has_children_excellent <- data %>% filter(has_children=="Y" & selfrated_mental_health=="Excellent") %>% tally() + 1
no_children_poor <- data %>% filter(has_children=="N" & selfrated_mental_health=="Poor") %>% tally() + 1
no_children_fair <- data %>% filter(has_children=="N" & selfrated_mental_health=="Fair") %>% tally() + 1
no_children_good <- data %>% filter(has_children=="N" & selfrated_mental_health=="Good") %>% tally() + 1
no_children_vgood <- data %>% filter(has_children=="N" & selfrated_mental_health=="Very Good") %>% tally() + 1
no_children_excellent <- data %>% filter(has_children=="N" & selfrated_mental_health=="Excellent") %>% tally() + 1

total_has_children <- has_children_poor + has_children_fair + has_children_good + has_children_vgood + has_children_excellent
total_no_children <- no_children_poor + no_children_fair + no_children_good + no_children_vgood + no_children_excellent

# repeat for the rest of the independent var

## example: predict mental health state given male, has children
# calculating prob(poor)
class_prior_prob_poor <- poor / total_mental
likelihood_poor <- (male_poor/poor) * (has_children_poor/poor)
predictor_prior_prob_poor <- (total_male/total_mental) * (total_has_children/total_mental)
prob_poor <- likelihood_poor * class_prior_prob_poor / predictor_prior_prob_poor

prob_poor

##          n
## 1 0.1545455

```

```

# repeat this process to get probability for fair, good, vgood, excellent
# the classification will be the class with the largest probability

# note that the probability value themselves should not be taken seriously since naive bayes is a bad e

# 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$self Rated mental health %>% table()
# poll<-poll[!grepl("Don't know", poll$self Rated mental health),]
# poll$self Rated mental health %>% table()

poll <- head(poll, 1000)

model <- nnet::multinom(self Rated mental health ~ age + sex + marital_status + self Rated health,
                        data = poll)

## # weights: 84 (65 variable)
## initial value 1788.175950

```

```
## iter 10 value 1217.602465
## iter 20 value 1148.942962
## iter 30 value 1132.529809
## iter 40 value 1127.538341
## iter 50 value 1125.837525
## iter 60 value 1124.707510
## iter 70 value 1124.598186
## final value 1124.596386
## 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
## Excellent 11.53664 -0.03491665 -0.3001279 6.468537
## Fair 17.86858 -0.06991689 -0.8579523 5.360938
## Good 30.82132 -0.04660078 -0.7155401 6.414833
## Poor 19.37958 -0.08049740 -0.8638462 4.139046
## Very good 29.52018 -0.04744081 -0.4051695 6.480312
## marital_statusMarried marital_statusSeparated
## Excellent -10.41129 8.120901
## Fair -11.47128 8.029377
## Good -10.19707 8.709045
## Poor -12.60247 9.358431
## Very good -10.23275 8.639051
## marital_statusSingle, never married marital_statusWidowed
## Excellent -11.80296 -10.83217
## Fair -12.41079 -32.68344
## Good -11.79471 -10.55576
## Poor -12.95577 -11.37188
## Very good -11.92809 -10.12267
## self_rated_healthExcellent self_rated_healthFair
## Excellent 19.285429 4.568500
## Fair 11.904930 1.823893
## Good -1.277130 -13.048693
## Poor 9.787835 -1.207769
## Very good 0.876981 -12.550617
## self_rated_healthGood self_rated_healthPoor
## Excellent 3.824870 20.442358
## Fair -1.003075 17.799820
## Good -13.683148 2.302144
## Poor -2.742509 16.790272
## Very good -13.257852 3.238129
## self_rated_healthVery good
## Excellent 19.6606504
## Fair 13.8608590
## Good 0.6669313
## Poor 11.0024561
## Very good 3.1513423
##
```

```
## Std. Errors:
##      (Intercept)      age    sexMale marital_statusLiving common-law
## Excellent    1.341674 0.03416602 0.9647511                0.3878738
## Fair         1.389493 0.03503881 0.9951848                0.4966857
## Good         1.457356 0.03398857 0.9589977                0.3838304
## Poor         1.547312 0.03800455 1.0812613                0.9497979
## Very good    1.462708 0.03411494 0.9626099                0.3772835
##      marital_statusMarried marital_statusSeparated
## Excellent          0.8837379                0.4439544
## Fair              0.9283008                0.5300380
## Good              0.8802635                0.3810928
## Poor              1.1571840                0.6383493
## Very good         0.8801744                0.3917895
##      marital_statusSingle, never married marital_statusWidowed
## Excellent          0.6973325                1.227993e+00
## Fair              0.7537025                5.160951e-08
## Good              0.6899751                1.215301e+00
## Poor              0.9234452                1.464508e+00
## Very good         0.6928557                1.217357e+00
##      selfRated_healthExcellent selfRated_healthFair
## Excellent          0.2249380                1.091411
## Fair              0.3954234                1.104556
## Good              0.6398070                1.225476
## Poor              0.7110687                1.228300
## Very good         0.6286038                1.232652
##      selfRated_healthGood selfRated_healthPoor
## Excellent          0.6793876                0.2920380
## Fair              0.7110747                0.2814524
## Good              0.8871095                0.6459217
## Poor              0.8093229                0.4182572
## Very good         0.8902829                0.6550287
##      selfRated_healthVery good
## Excellent          0.2185480
## Fair              0.3073876
## Good              0.6227630
## Poor              0.7042581
## Very good         0.6168416
##
## Residual Deviance: 2249.193
## AIC: 2379.193
```

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

```
##      Don't know Excellent      Fair      Good      Poor Very good
## 1 2.330000e-08 0.6731361 0.020308530 0.1021620 0.0036780279 0.2007154
## 2 1.077170e-02 0.1891944 0.028231660 0.5032988 0.0041965547 0.2643069
## 3 2.055971e-09 0.2375888 0.015130081 0.1872810 0.0006473774 0.5593528
## 4 4.255623e-09 0.2773835 0.009949682 0.1805215 0.0003579040 0.5317875
## 5 2.111469e-10 0.1779567 0.056832472 0.4743416 0.0098516022 0.2810176
## 6 8.034049e-09 0.6514559 0.008720240 0.1077107 0.0007880384 0.2313251
```

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


##	Don't know	Excellent	Fair	Good	Poor	Very good
##	9.792516e-09	6.228990e-01	3.206166e-02	8.984887e-02	8.030243e-03	2.471602e-01

Results

Discussion

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>