

# How does age or self rated life feelings affect a person's mental health

Group 49: Yingchen Tan, Xueqi Wang, XiWen Ran, Zezhou Han

2020, Oct, 19th

## Abstract

In this context, the group of four people discussed the Canadian General Social Survey data list which collects Canada lived respondent's personal information during 2017. The group is defined as the extent to describe one of the relationships in the dataset.

To test the hypothesis — the relationship between respondent's age and self rated feelings rate that affect their mental health. Respondents were randomly picked and answered with their personal information. They were asked to answer their age in number, self rated their life feelings rate from 1-10 and rate their mental health. Excluding N/A error variables, we also clean the data for variable “self\_rated\_mental\_health”, such that people rated “good, very good, excellent” means positive and “fair, poor” means negative. Based on these data, we first used histogram, then used a logit regression model to make the discussion.

The results show the effect in the variables which is the same as hypothesized: older people and people who rated higher on feelings life will get higher mental health scores.

These results suggest that people should pay more attention to today's young people, maybe because they were feeling hard on study or working. On this basis, we would consider analyzing more variables that could potentially affect people's mental health positively or negatively.

## Introduction

The project is based on the data of Canadian General Social Survey(GSS) and making the model to analyse our topic. Canadian General Social Survey is the data that collected people's personal and family information during 2017, which included family members, address location, income, etc.

We will discuss the topic of “How does age or self rated life feelings affect a person's mental health”.

There is a variable called “self\_rated\_mental\_health” which collects from each respondent, it represents the respondent's own self rated mental health score with positive side of “good, very good, excellent” and also negative side of “fair, poor”. And all group members were interested in this variable. This variable is based on their self rated, which is different from other “real” data, you can directly see the value of height, the revenue or the family members, etc. We would like to know the reason how respondents rank their score. To discuss this variable, we choose variables “age” and “feelings\_life”. We will compute the correlation between “age” and “feelings\_life” to confirm that it is independent or dependent on the ratio of “self\_rated\_mental\_health”.

Then we are using logistic regression to model our dependent variables by using R program, discussing their linear relationship by looking at the graph. We were hoping that we can find linear relationships between any variables and well reported them by observing the data, so making the logistic regression model, analyzing the result and discussion with conclusion will be our main works in the next steps.

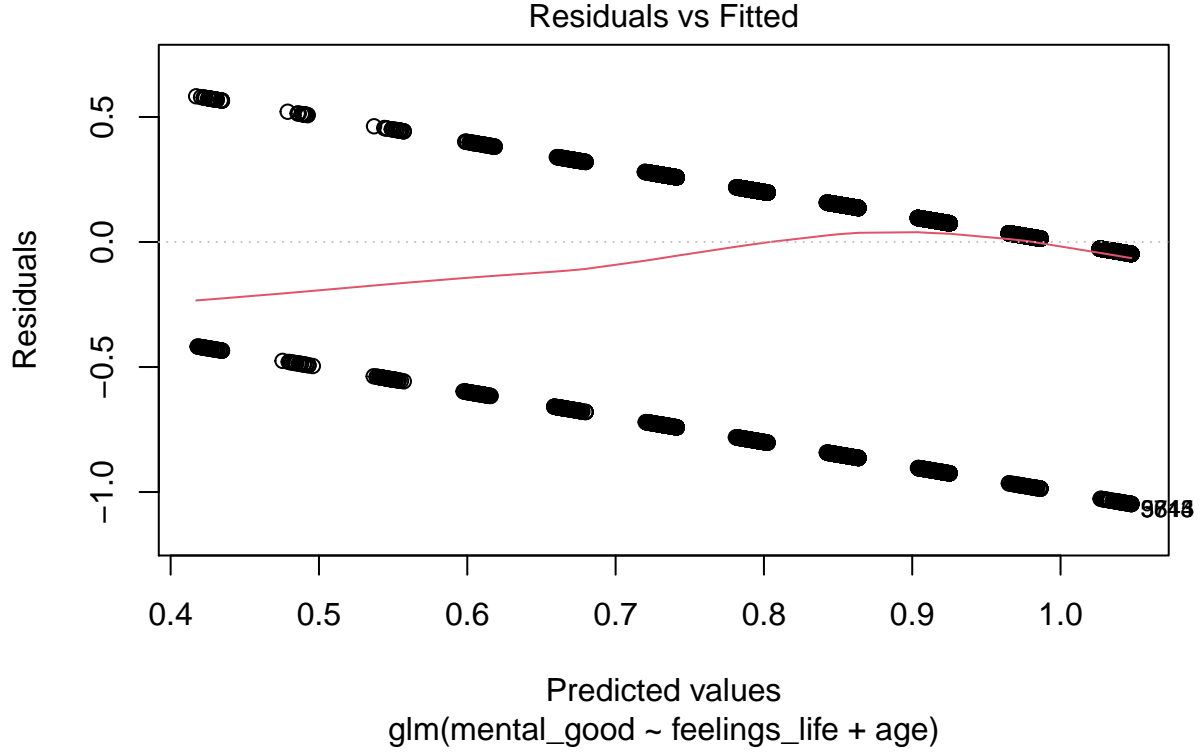
## Data

We could use our variables to see the correlation between the predictors and estimators. For example, we could see the relationship between the age and the feeling life index to see which age period of people are satisfied with their life. But the weakness of our pick on data is we still have a lot of NA response in the database so we will have some bias on our model.

By checking the dataset, we can say that the survey is from the Census of Canada. The target population for the 2017 GSS included all persons 15 years of age and older in Canada excluding the residents of Yukon, Northwest Territories, and Nunavut; and Full-time residents of institutions. In order to carry out sampling, each of the ten provinces were divided into strata(i.e geographic areas). Also, the non-CMA areas of each of the ten provinces were also grouped to form ten more strata, for a total of 27 strata. The payoff is this method takes the longest time cost compared to others. Since that's the Census of the country, I think there are no non-response situations during the experiment. But we still can make the hypothesis for if there are some people that make fake information or they just are not submitting their response Since we have some responses that are actually not available. In this case, the cost might be time costs and that also would cause bias with the final results. That would also cause Cognitive bias when the bias is large.

## Model

```
##
## Call:
## glm(formula = mental_good ~ feelings_life + age, data = gssd)
##
## Deviance Residuals:
##      Min       1Q   Median       3Q      Max
## -1.04770  -0.03415   0.07493   0.09227   0.58278
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept)  4.082e-01  9.963e-03  40.976 < 2e-16 ***
## feelings_life 6.131e-02  1.064e-03  57.605 < 2e-16 ***
## age          3.290e-04  9.868e-05   3.334 0.000857 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## (Dispersion parameter for gaussian family taken to be 0.06191006)
##
##      Null deviance: 1462.8  on 20272  degrees of freedom
## Residual deviance: 1254.9  on 20270  degrees of freedom
## AIC: 1136.3
##
## Number of Fisher Scoring iterations: 2
## [1] 0 10
## [1] 15 80
```



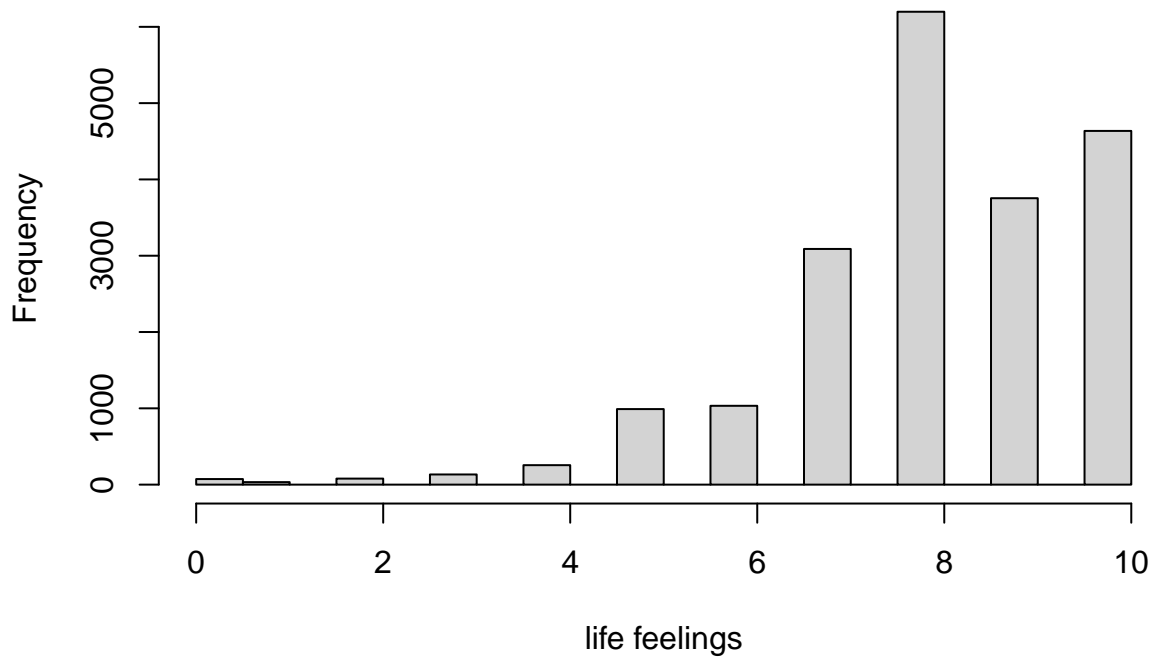
In this part, we use R to create a logistic regression model to find the probability of people's mental health in good status. We use `feelings_life` and `age` to be the explanatory variables. Since the age span of the sample is large, then age may be an important factor to affect people's mental health. Also, the variable `feelings_life` records their self rate to their life, which may directly show people's life satisfaction, so `feelings_life` should also be an important factor affecting people's mental health. Let  $p$  be the probability of people's mental health in good status, and `feelings_life` to be  $X_1$  and `age` to be  $X_2$ , then the coefficient of  $X_1$  is  $\hat{\beta}_1$  and the coefficient of  $X_2$  is  $\hat{\beta}_2$ .  $\hat{\beta}_0$  is the intercept. From the plot we can see the red line does not go over the data point range and it is not a straight line, which means using logistic regression assumption is appropriate.

From the summary table, we can see the  $\hat{\beta}_1 = 6.131 \times 10^{-2}$ ,  $\hat{\beta}_2 = 3.29 \times 10^{-4}$  and  $\hat{\beta}_0 = 4.082 \times 10^{-1}$ , so the formula of the linear regression model is  $\log(\hat{p}/(1 - \hat{p})) = 6.131 \times 10^{-2}X_1 + 3.29 \times 10^{-4}X_2 + 4.082 \times 10^{-1}$  for  $X_1 \in [0, 10]$  and  $X_2 \in [15, 80]$

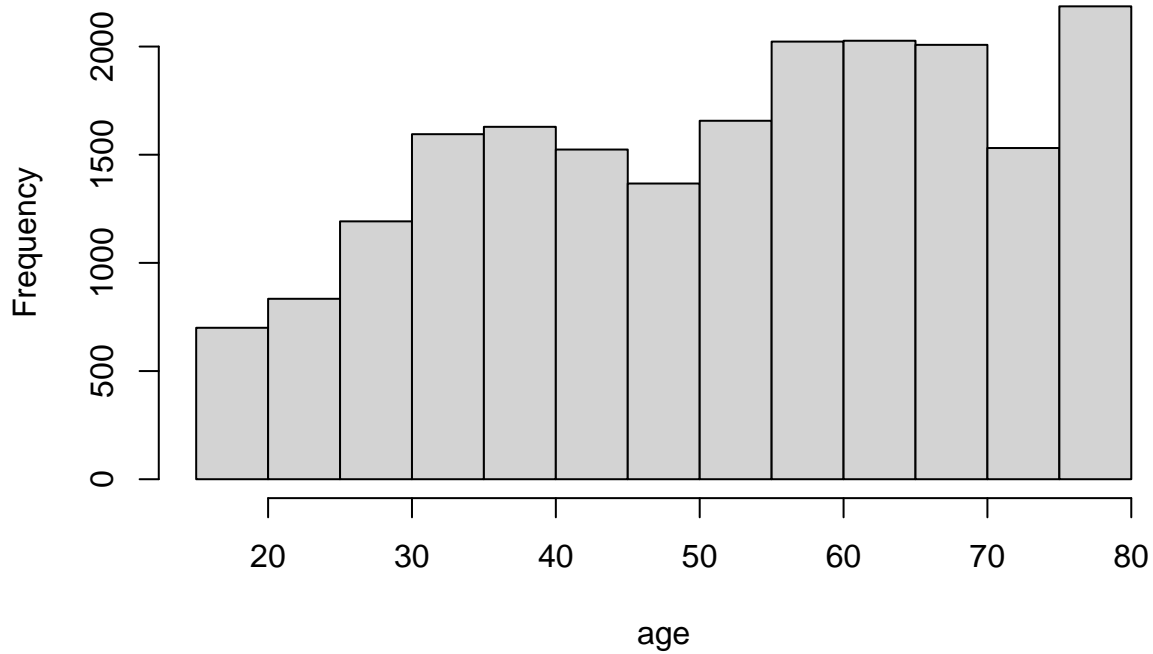
Also, we can see the p-value for the null hypothesis that  $\hat{\beta}_1 = 0$  and  $\hat{\beta}_0 = 0$  are both smaller than  $2 \times 10^{-16}$ , and the p-value for the null hypothesis that  $\hat{\beta}_2 = 0$  is 0.000857. Three numbers are very small numbers, so we can reject these three null hypotheses. It means that age and `feelings_life` are relating to the probability of people's mental health in good status and the model has an intercept with y-axis not in  $\log(\hat{p}/(1 - \hat{p}))=0$ .

## Results

### Histogram of life feelings



### Histogram of age



The histogram of the age of respondents shows that most of the respondents belong to the age brackets of twenty to eighty years. The histogram of individuals' feelings about life as a whole explores a scale of 0 to 10, where 0 means "Very dissatisfied" and 10 means "Very satisfied." It can be seen that most of the data values in the histogram are in the range of 5 to 10. This leads to the conclusion that a large number of individuals

are satisfied with their life.

The Residuals vs. Fitted plot displays the linear relationship between variables, which demonstrates equally spread residuals around a horizontal line without distinct patterns. The data points were simulated in a way that met the model assumptions precisely, indicating that the model was quite accurate. The variables “feelings\_life” and “age” are both statistically significant because the corresponding p-values are less than the significance level of 0.05. In short, “feelings\_life” and “age” are two variables that relate to the probability of individuals’ mental health in good status. When the scale of individuals’ feelings about life becomes larger, individuals’ mental health status would more likely be good. On the other hand, age is also a significant contributor to the fairly good status of mental health.

## Discussion

The result shows that maintaining a good status of mental health can be affected by individuals’ feelings about life as a whole. We find similar trends in age, and we can access the relationship between age and mental health status, which older people are more likely to achieve and sustain a state of good mental health. To avoid caveats, we clean up the inaccurate data, which is not appropriate for the purpose of getting the correct model. We keep the data values which have a strong influence on our observation results. The small world and the large world always provide the opportunity for a conversation between model and reality. The small world represents the model itself, while the large world is what we hope to deploy the model in, which is behind the presentation of the model’s data. There exist challenges in our statistical modeling. Although we attempt to minimize any uncertainty related to our data, there are still weaknesses that we must improve in the future.

## Weaknesses

One of the weaknesses of our work is that the model is limited by the quality of the data. Besides, the data for those individuals who answer “Don’t know” to particular questions cannot contribute to the analysis. As a result, when the data is not fulfilled, the result will be biased as well. For future improvement, we would consider analyzing more variables that could potentially affect people’s mental health positively or negatively. We should also make sure that the most significant variable in determining the mental health status of individuals would be provided. Thus, it is possible for us to analyze the difference between various influence factors.

## Next Steps

Since the data collected people’s personal and family information during 2017, there should be some different questions to explore in 2020. Based on our current results, it is essential to conduct a follow-up survey, which helps us to estimate the accuracy and precision of the results obtained. We would include the questions, such as asking individuals’ feelings about life or their feelings in general. In addition, there would be questions that relate to the effect of COVID-19. The important variable we are going to measure is to what extent, COVID-19 influences individuals’ day to day life, given a scale of 0 to 10.

## References

General social survey on Family (cycle 31), 2017 <https://sda-artsci-utoronto-ca.myaccess.library.utoronto.ca/sdaweb/html/gss.htm>

2017 General Social Survey: Families Cycle 31, PUMF-Data Dictionary, February 2020. [https://sda-artsci-utoronto-ca.myaccess.library.utoronto.ca/sdaweb/dli2/gss/gss31/gss31/more\\_doc/index.htm](https://sda-artsci-utoronto-ca.myaccess.library.utoronto.ca/sdaweb/dli2/gss/gss31/gss31/more_doc/index.htm)