

Apply logistic regression Model to analyze the number of children someone has

Junming Zhang, Hairong Sun

Saturday, October 17, 2020

Abstract

Having children is one of the key factors of a family, and children play an important role in family relationship. Also, birth rate influences the future of a community because they will be the labour force of the community, which has an impact on the economical and social development, and even existence of the community. Therefore, it is necessary to have an indicator to suggest what factors influence if someone may have many children, and so the government can make policies to control birth rate based on the indicator. We tried to investigate such a factor by the dataset **Canadian General Social Survey (GSS)** (citation 5) of year 2017, a dataset generated in Canada, in which had some attributes like number of childrens one person has and some other features of a person which may help (total_children, marital_status, education, partner_education, income_family, selfRated_health and selfRated_mental_health). To generate such an indicator, we built a logistic regression model with the dataset, which predicts the possibility of many children (more than 2 children) one may have with some attributes I searched before which may help.

Introduction

Our goal is to find a model to predict the possibility one may have more than 2 children based on attributes may have strong affects, and test how strongly are these attributes correlated to the number of children one may have. To create such a model, we use the dataset **Canadian General Social Survey (GSS)** from year 2017, which has attributes that may contribute to the our research goal, like total_children, marital_status, education, partner_education, income_family, selfRated_health and selfRated_mental_health, and we built our model with these attributes. Total children is the attribute we are interested in. Education and partner_education are important because people who are enrolled in the tertiary education tend to postpone their marriage and have fewer children (citation 1). We consider income_family because income_family is a mirror of the economy development, which is related to birth rate, for instance, economic depression may mean low fertility (Pobric & Robinson, 2015). We look into marital_status, because the type of partnership may contribute to birth rate, for example, those who get married may have more children than those who cohabit (Martinez, Daniels, & Chandra, 2012). Finally but still important, the health status impacts the number of children one may have, for example, countries with higher HDI (Human Development Index, which involves life expectancy, education, and per capita income) may have lower fertility rate, which is reflected by the citation 4, countries with lower HDI have higher fertility rate and vice versa. Since health is related to HDI for life expectancy, we put the selfRated_health and selfRated_mental_health into the model. Therefore, we create a logistic regression model with these attributes to predict the possibility one may have many (more than 2) children, and analyze how is our prediction related to these attributes, and how strong is our interest and these attributes correlated.

Data

The dataset is obtained from **Canadian General Social Survey (GSS)** of year 2017, it contains all the attributes I listed in the Introduction section that are used to build the model. To make the dataset, they use a questionnaire and interview the respondents on phone call (Beaupré, 2020). A brief outline of the questionnaire is following (Beaupré, 2020):

- Entry component (respondent's date of birth)
- Family origins
- Leaving the parental home
- Conjugal history
- Intentions and reasons to form a union
- Respondent's children
- Fertility intentions
- Maternity/parental leave
- Organization and decision making within the household
- Arrangements and financial support after a separation/divorce
- Labour market new and education
- Health and subjective well-being
- Characteristics of respondent's dwelling
- Characteristics of respondent of spouse/partner

The questionnaire to build the dataset was delivered by telephone (Beaupré, 2020), and this questionnaire was helpful because it covers many details on the personal conditions of the respondent, like health and education, for which there is significant proof that influences the fertility rate. However, the previous nationalities of respondents are ignored, which may also introduce errors in the result, because some countries may have special cultures and religions that affect the fertility rate. There are also pros and cons for collecting data by phone. The benefit is that, since the most people have their own telephones today, so it is easy to connect and the data can be collected with lower costs. However, some people may not respond to the phone call, which leads to the non-participation error.

There are 81 variables/attributes and 20602 observations in the dataset. The variables generally cover many aspects about the living conditions and the personal conditions of the interviewee, which may suggest our interest, and we tested some of them which are possibly helpful according to the documents and references we found, and investigate the correlation between the variables in our scope and our interest. And also, the dataset has a large number of observations with respect to the place where the data were collected, and thus this makes the results (in Canada) can be found from the dataset more representative. However, since the dataset is only limited to one country (Canada), the variables do not reflect other factors may also have impacts but not suitable for just one country, like policy, war or peace, natural conditions, and so hard to reflect worldwide facts.

The data are collected with the stratified random sampling (simple random sampling without replacement in the stratum) method (Beaupré, 2020), a probability sampling approach. The target population for the dataset included all persons 15 years of age and older in Canada, excluding: 1. Residents of the Yukon, Northwest Territories, and Nunavut; and 2. Full-time residents of institutions (Beaupré, 2020). The frame of the survey is 1. Lists of telephone numbers in use (both landline and cellular) available to Statistics Canada from various sources (telephone companies, Census of population, etc.); and 2. The Address Register (AR): List of all dwellings within the ten provinces. The probability sampling method (collection approach for this dataset) decreases errors like generalization and more representative for the whole population. However,

there are some drawbacks of the dataset from both non-sampling error. The non-sampling error is mainly from (partial or total) non-participation. This is handled by adjusting the weights to less for non-participation cases (Beaupré, 2020), and in my implementation, I removed all rows with NA in the columns we needed to build the model.

Model

Here is the number of observations for each stratum (since the data are collected by stratified sampling without replacement), the stratum was divided based on the province the interviewee lived (Beaupré, 2020), we built our logistic regression model based on the stratification below:

```
## # A tibble: 10 x 2
## # Groups:   province [10]
##   province      n
##   <chr>      <int>
## 1 Alberta    1064
## 2 British Columbia 1490
## 3 Manitoba    708
## 4 New Brunswick   772
## 5 Newfoundland and Labrador 702
## 6 Nova Scotia    837
## 7 Ontario     3313
## 8 Prince Edward Island 421
## 9 Quebec      2191
## 10 Saskatchewan   675

##
## Call:
## svyglm(formula = if_many ~ as.factor(marital_status) + as.factor(education) +
##   as.factor(partner_education) + as.factor(income_family) +
##   as.factor(self_rated_health) + as.factor(self_rated_mental_health),
##   design = gss.design, family = "binomial")
##
## Weighted Residuals:
##   Min      1Q  Median      3Q      Max
## -1.5352 -0.6442 -0.5368  1.0985  6.0155
##
## Coefficients:
##
## (Intercept)                                Estimate
## as.factor(marital_status)Living common-law -0.43755
## as.factor(marital_status)Married           0.11963
## as.factor(marital_status)Separated         0.08655
## as.factor(marital_status)Single, never married -1.78663
## as.factor(marital_status)Widowed           0.15065
## as.factor(education)College, CEGEP or other non-university certificate or di... 0.10964
## as.factor(education)High school diploma or a high school equivalency certificate 0.32493
## as.factor(education)Less than high school diploma or its equivalent           0.63463
## as.factor(education)Trade certificate or diploma                             0.26572
## as.factor(education)University certificate or diploma below the bachelor's level 0.15623
## as.factor(education)University certificate, diploma or degree above the bach... 0.01973
## as.factor(partner_education)College, CEGEP or other non-university certificate or d... 0.08019
## as.factor(partner_education)High school diploma or a high school equivalency certi... 0.27359
## as.factor(partner_education)Less than high school diploma or its equivalent      0.49913
```

## as.factor(partner_education)Trade certificate or diploma	0.17458
## as.factor(partner_education)University certificate or diploma below the bachelor's level	0.15677
## as.factor(partner_education)University certificate, diploma or degree above the ba...	-0.18488
## as.factor(income_family)\$125,000 and more	0.09295
## as.factor(income_family)\$25,000 to \$49,999	0.04431
## as.factor(income_family)\$50,000 to \$74,999	-0.08042
## as.factor(income_family)\$75,000 to \$99,999	-0.10151
## as.factor(income_family)Less than \$25,000	-0.35354
## as.factor(selfRatedHealth)Excellent	-1.00447
## as.factor(selfRatedHealth)Fair	-0.74889
## as.factor(selfRatedHealth)Good	-0.91246
## as.factor(selfRatedHealth)Poor	-0.57634
## as.factor(selfRatedHealth)Very good	-1.01712
## as.factor(selfRatedMentalHealth)Excellent	1.11938
## as.factor(selfRatedMentalHealth)Fair	1.01850
## as.factor(selfRatedMentalHealth)Good	1.15656
## as.factor(selfRatedMentalHealth)Poor	0.72011
## as.factor(selfRatedMentalHealth)Very good	1.04924
##	Std. Error
## (Intercept)	0.80190
## as.factor(marital_status)Living common-law	0.15624
## as.factor(marital_status)Married	0.14709
## as.factor(marital_status)Separated	0.27844
## as.factor(marital_status)Single, never married	0.21058
## as.factor(marital_status)Widowed	0.28506
## as.factor(education)College, CEGEP or other non-university certificate or di...	0.06728
## as.factor(education)High school diploma or a high school equivalency certificate	0.06961
## as.factor(education)Less than high school diploma or its equivalent	0.08758
## as.factor(education)Trade certificate or diploma	0.08999
## as.factor(education)University certificate or diploma below the bachelor's level	0.11604
## as.factor(education)University certificate, diploma or degree above the bach...	0.08362
## as.factor(partner_education)College, CEGEP or other non-university certificate or d...	0.06902
## as.factor(partner_education)High school diploma or a high school equivalency certi...	0.06823
## as.factor(partner_education)Less than high school diploma or its equivalent	0.08684
## as.factor(partner_education)Trade certificate or diploma	0.09231
## as.factor(partner_education)University certificate or diploma below the bachelor's level	0.11428
## as.factor(partner_education)University certificate, diploma or degree above the ba...	0.08616
## as.factor(income_family)\$125,000 and more	0.06704
## as.factor(income_family)\$25,000 to \$49,999	0.07913
## as.factor(income_family)\$50,000 to \$74,999	0.07434
## as.factor(income_family)\$75,000 to \$99,999	0.07511
## as.factor(income_family)Less than \$25,000	0.13440
## as.factor(selfRatedHealth)Excellent	0.47467
## as.factor(selfRatedHealth)Fair	0.47697
## as.factor(selfRatedHealth)Good	0.47342
## as.factor(selfRatedHealth)Poor	0.48759
## as.factor(selfRatedHealth)Very good	0.47333
## as.factor(selfRatedMentalHealth)Excellent	0.64966
## as.factor(selfRatedMentalHealth)Fair	0.65492
## as.factor(selfRatedMentalHealth)Good	0.64903
## as.factor(selfRatedMentalHealth)Poor	0.69056
## as.factor(selfRatedMentalHealth)Very good	0.64929
##	t value
## (Intercept)	-1.771

## as.factor(marital_status)Living common-law	-2.801
## as.factor(marital_status)Married	0.813
## as.factor(marital_status)Separated	0.311
## as.factor(marital_status)Single, never married	-8.484
## as.factor(marital_status)Widowed	0.528
## as.factor(education)College, CEGEP or other non-university certificate or di...	1.630
## as.factor(education)High school diploma or a high school equivalency certificate	4.668
## as.factor(education)Less than high school diploma or its equivalent	7.246
## as.factor(education)Trade certificate or diploma	2.953
## as.factor(education)University certificate or diploma below the bachelor's level	1.346
## as.factor(education)University certificate, diploma or degree above the bach...	0.236
## as.factor(partner_education)College, CEGEP or other non-university certificate or d...	1.162
## as.factor(partner_education)High school diploma or a high school equivalency certi...	4.010
## as.factor(partner_education)Less than high school diploma or its equivalent	5.748
## as.factor(partner_education)Trade certificate or diploma	1.891
## as.factor(partner_education)University certificate or diploma below the bachelor's level	1.372
## as.factor(partner_education)University certificate, diploma or degree above the ba...	-2.146
## as.factor(income_family)\$125,000 and more	1.387
## as.factor(income_family)\$25,000 to \$49,999	0.560
## as.factor(income_family)\$50,000 to \$74,999	-1.082
## as.factor(income_family)\$75,000 to \$99,999	-1.351
## as.factor(income_family)Less than \$25,000	-2.630
## as.factor(selfRatedHealth)Excellent	-2.116
## as.factor(selfRatedHealth)Fair	-1.570
## as.factor(selfRatedHealth)Good	-1.927
## as.factor(selfRatedHealth)Poor	-1.182
## as.factor(selfRatedHealth)Very good	-2.149
## as.factor(selfRatedMentalHealth)Excellent	1.723
## as.factor(selfRatedMentalHealth)Fair	1.555
## as.factor(selfRatedMentalHealth)Good	1.782
## as.factor(selfRatedMentalHealth)Poor	1.043
## as.factor(selfRatedMentalHealth)Very good	1.616
##	Pr(> t)
## (Intercept)	0.07660
## as.factor(marital_status)Living common-law	0.00511
## as.factor(marital_status)Married	0.41604
## as.factor(marital_status)Separated	0.75592
## as.factor(marital_status)Single, never married	< 2e-16
## as.factor(marital_status)Widowed	0.59718
## as.factor(education)College, CEGEP or other non-university certificate or di...	0.10319
## as.factor(education)High school diploma or a high school equivalency certificate	3.08e-06
## as.factor(education)Less than high school diploma or its equivalent	4.55e-13
## as.factor(education)Trade certificate or diploma	0.00316
## as.factor(education)University certificate or diploma below the bachelor's level	0.17820
## as.factor(education)University certificate, diploma or degree above the bach...	0.81348
## as.factor(partner_education)College, CEGEP or other non-university certificate or d...	0.24537
## as.factor(partner_education)High school diploma or a high school equivalency certi...	6.11e-05
## as.factor(partner_education)Less than high school diploma or its equivalent	9.27e-09
## as.factor(partner_education)Trade certificate or diploma	0.05863
## as.factor(partner_education)University certificate or diploma below the bachelor's level	0.17012
## as.factor(partner_education)University certificate, diploma or degree above the ba...	0.03192
## as.factor(income_family)\$125,000 and more	0.16560
## as.factor(income_family)\$25,000 to \$49,999	0.57555
## as.factor(income_family)\$50,000 to \$74,999	0.27938

```

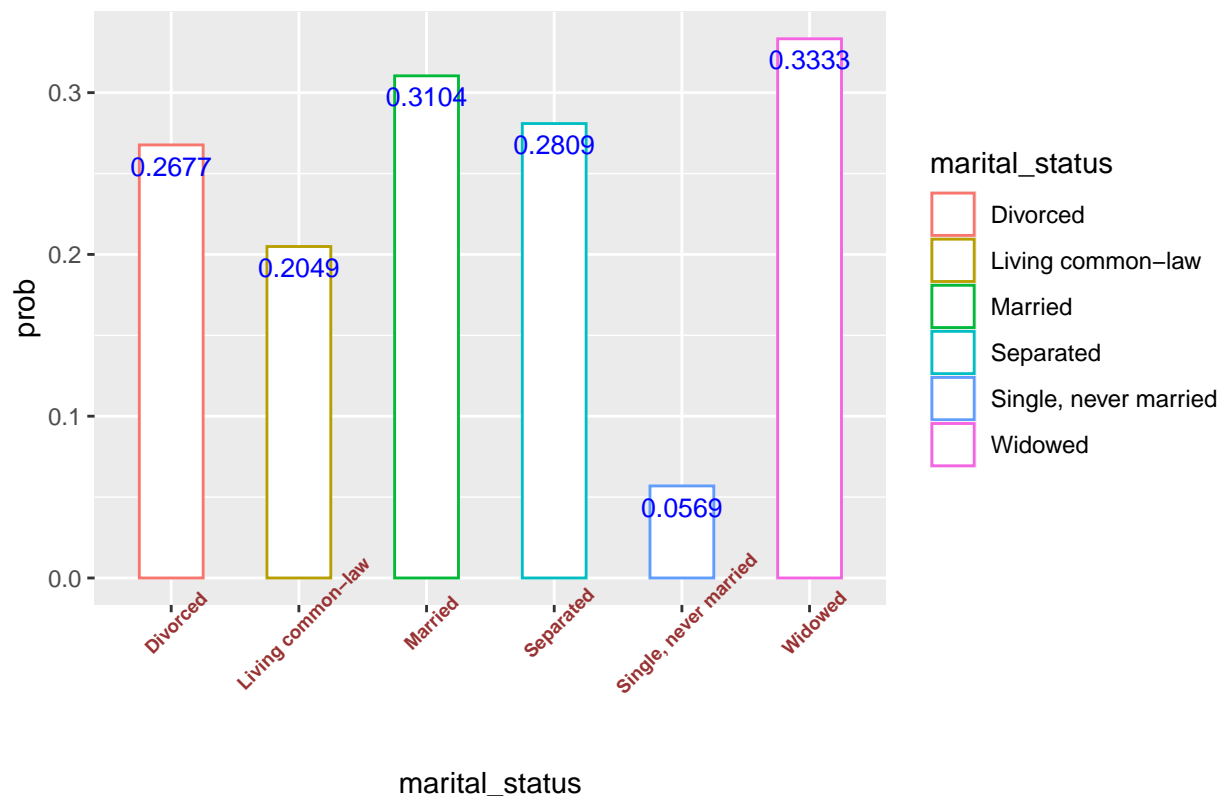
## as.factor(income_family)$75,000 to $99,999 0.17660
## as.factor(income_family)Less than $25,000 0.00854
## as.factor(selfRatedHealth)Excellent 0.03435
## as.factor(selfRatedHealth)Fair 0.11642
## as.factor(selfRatedHealth)Good 0.05396
## as.factor(selfRatedHealth)Poor 0.23722
## as.factor(selfRatedHealth)Very good 0.03167
## as.factor(selfRatedMentalHealth)Excellent 0.08491
## as.factor(selfRatedMentalHealth)Fair 0.11993
## as.factor(selfRatedMentalHealth)Good 0.07477
## as.factor(selfRatedMentalHealth)Poor 0.29706
## as.factor(selfRatedMentalHealth)Very good 0.10612
##
## (Intercept) .
## as.factor(marital_status)Living common-law **
## as.factor(marital_status)Married
## as.factor(marital_status)Separated
## as.factor(marital_status)Single, never married ***
## as.factor(marital_status)Widowed
## as.factor(education)College, CEGEP or other non-university certificate or di...
## as.factor(education)High school diploma or a high school equivalency certificate ***
## as.factor(education)Less than high school diploma or its equivalent ***
## as.factor(education)Trade certificate or diploma **
## as.factor(education)University certificate or diploma below the bachelor's level
## as.factor(education)University certificate, diploma or degree above the bach...
## as.factor(partner_education)College, CEGEP or other non-university certificate or d...
## as.factor(partner_education)High school diploma or a high school equivalency certi... ***
## as.factor(partner_education)Less than high school diploma or its equivalent ***
## as.factor(partner_education)Trade certificate or diploma .
## as.factor(partner_education)University certificate or diploma below the bachelor's level
## as.factor(partner_education)University certificate, diploma or degree above the ba... *
## as.factor(income_family)$125,000 and more
## as.factor(income_family)$25,000 to $49,999
## as.factor(income_family)$50,000 to $74,999
## as.factor(income_family)$75,000 to $99,999
## as.factor(income_family)Less than $25,000 **
## as.factor(selfRatedHealth)Excellent *
## as.factor(selfRatedHealth)Fair
## as.factor(selfRatedHealth)Good .
## as.factor(selfRatedHealth)Poor
## as.factor(selfRatedHealth)Very good *
## as.factor(selfRatedMentalHealth)Excellent .
## as.factor(selfRatedMentalHealth)Fair
## as.factor(selfRatedMentalHealth)Good .
## as.factor(selfRatedMentalHealth)Poor
## as.factor(selfRatedMentalHealth)Very good
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 0.998 on 12131 degrees of freedom
## Multiple R-squared:  0.001652, Adjusted R-squared:  -0.001722
## F-statistic: 0.6274 on 32 and 12131 DF, p-value: 0.9498

```

Results

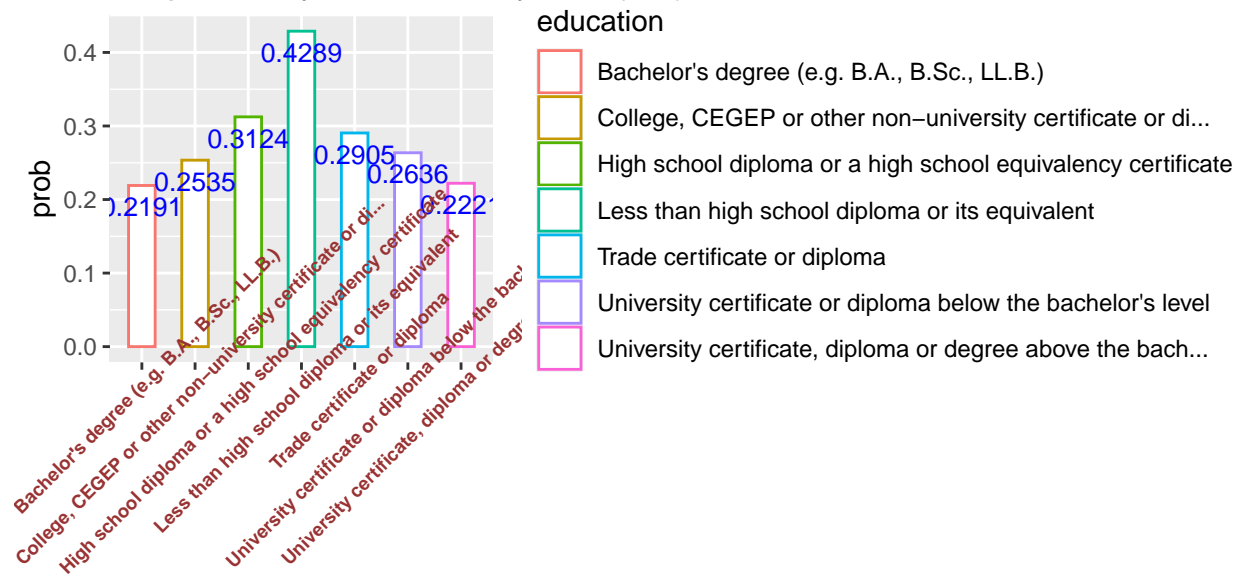
```
##          marital_status have_many total  prob
## 1          Divorced         68    254 0.2677
## 2   Living common-law       382   1864 0.2049
## 3           Married      2824   9097 0.3104
## 4           Separated        25     89 0.2809
## 5 Single, never married        45    791 0.0569
## 6           Widowed         26     78 0.3333
```

the possibility to have many kids (> 2) for each marital status



```
##          education have_many total
## 1   Bachelor's degree (e.g. B.A., B.Sc., LL.B.)       583   2661
## 2 College, CEGEP or other non-university certificate or di...   720   2840
## 3 High school diploma or a high school equivalency certificate   831   2660
## 4   Less than high school diploma or its equivalent       537   1252
## 5           Trade certificate or diploma                 278    957
## 6 University certificate or diploma below the bachelor's level   131    497
## 7 University certificate, diploma or degree above the bach...   290   1306
##      prob
## 1 0.2191
## 2 0.2535
## 3 0.3124
## 4 0.4289
## 5 0.2905
## 6 0.2636
## 7 0.2221
```

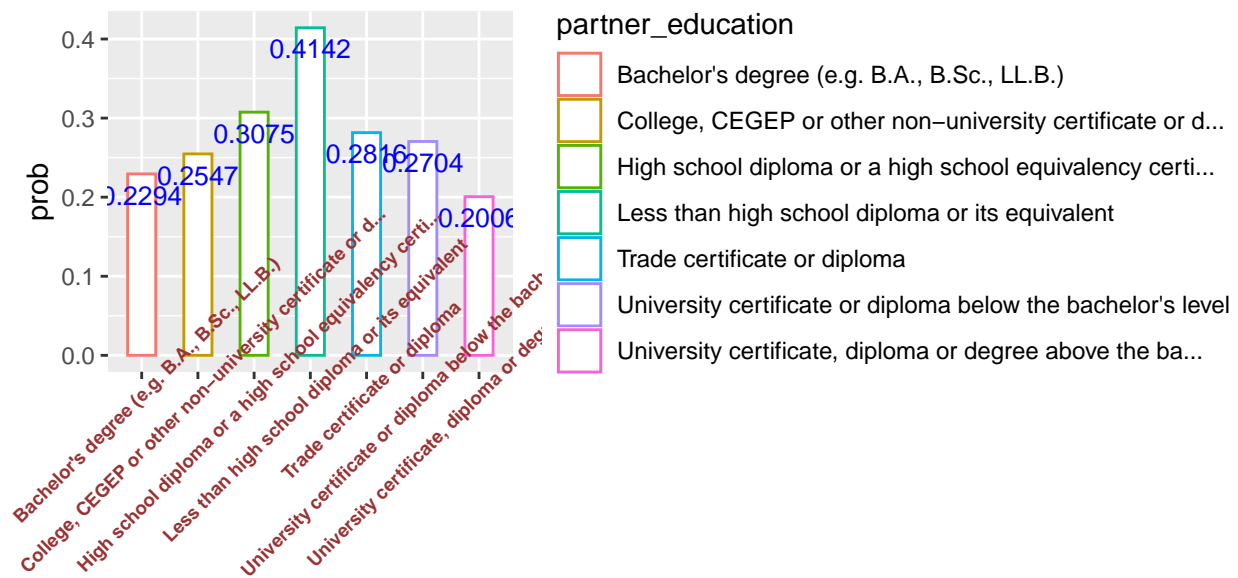
the possibility to have many kids (> 2) for each education status



education

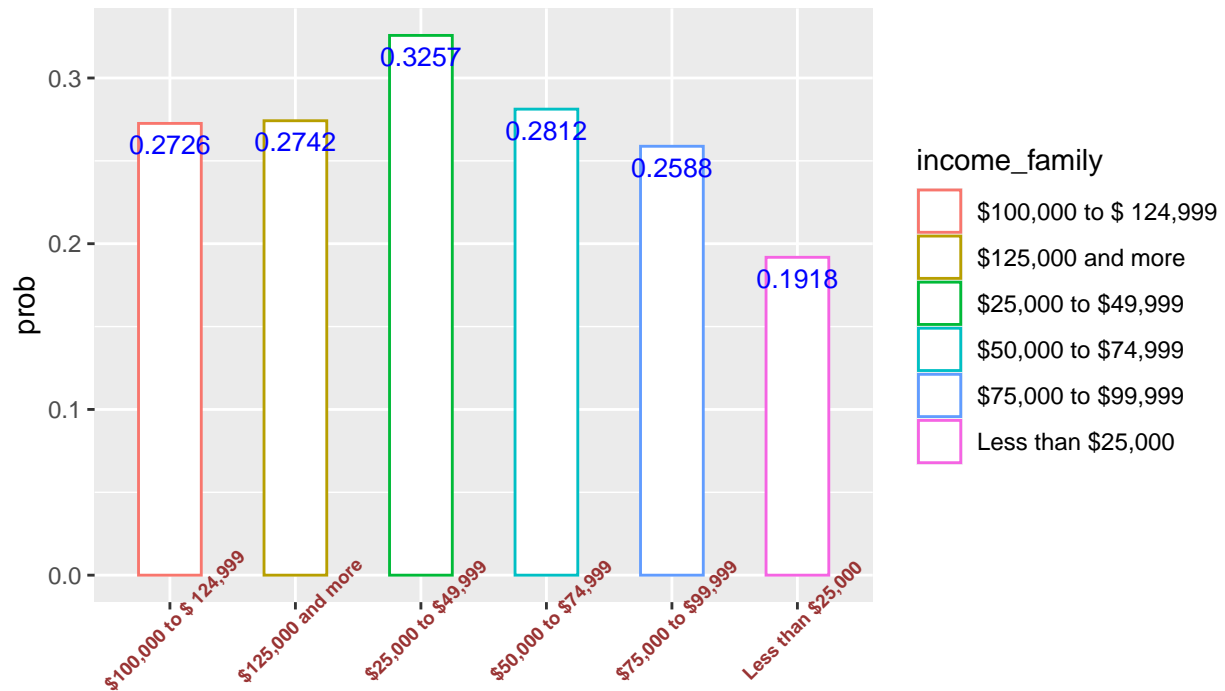
##	partner_education	have_many	total
## 1	Bachelor's degree (e.g. B.A., B.Sc., LL.B.)	576	2511
## 2	College, CEGEP or other non-university certificate or d...	657	2579
## 3	High school diploma or a high school equivalency certi...	944	3070
## 4	Less than high school diploma or its equivalent	541	1306
## 5	Trade certificate or diploma	254	902
## 6	University certificate or diploma below the bachelor's level	139	514
## 7	University certificate, diploma or degree above the ba...	259	1291
##	prob		
## 1	0.2294		
## 2	0.2547		
## 3	0.3075		
## 4	0.4142		
## 5	0.2816		
## 6	0.2704		
## 7	0.2006		

the possibility to have many kids (> 2) for each partner_education status



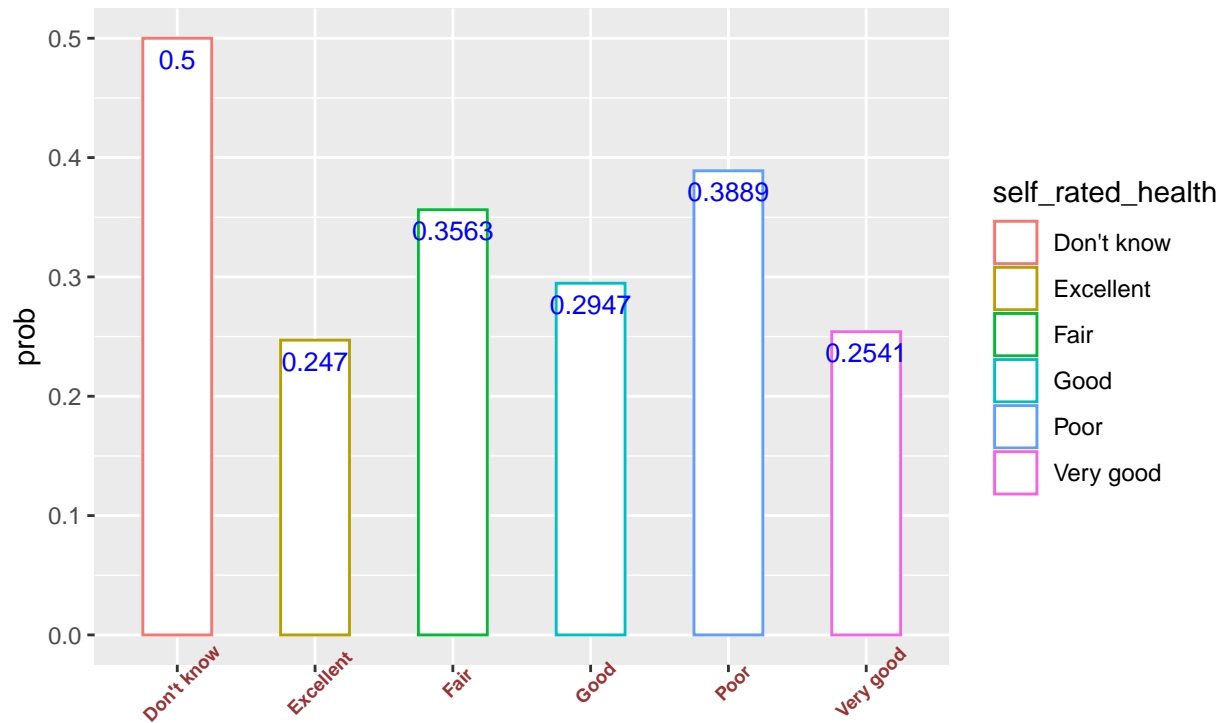
partner_education					
##	income_family	have_many	total	prob	
## 1	\$100,000 to \$ 124,999	467	1713	0.2726	
## 2	\$125,000 and more	1037	3782	0.2742	
## 3	\$25,000 to \$49,999	595	1827	0.3257	
## 4	\$50,000 to \$74,999	624	2219	0.2812	
## 5	\$75,000 to \$99,999	549	2121	0.2588	
## 6	Less than \$25,000	98	511	0.1918	

the possibility to have many kids (> 2) for each interval of family income



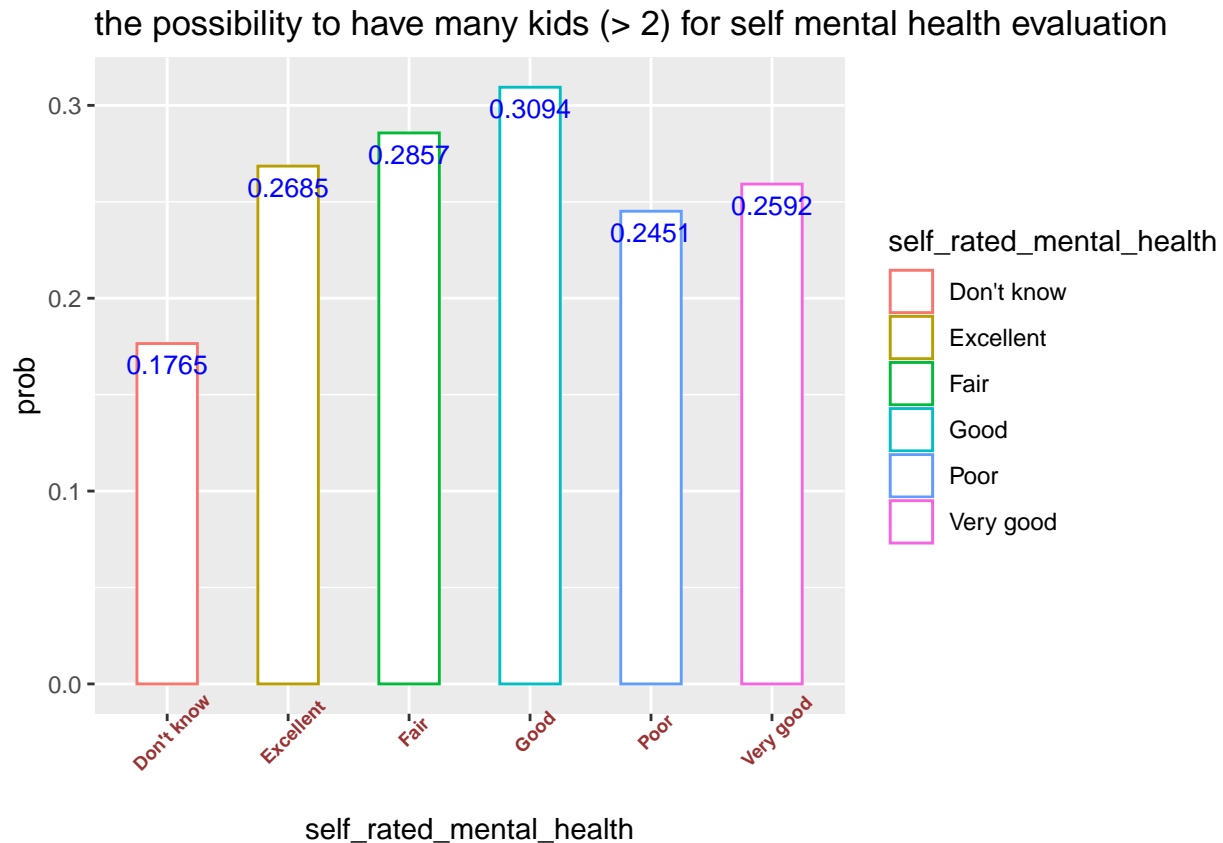
income_family					
##	selfRatedHealth	haveMany	total	prob	
## 1	Don't know	10	20	0.5000	
## 2	Excellent	692	2802	0.2470	
## 3	Fair	362	1016	0.3563	
## 4	Good	1065	3614	0.2947	
## 5	Poor	119	306	0.3889	
## 6	Very good	1122	4415	0.2541	

the possibility to have many kids (> 2) for self health evaluation



self Rated health

##	self Rated mental health	have many	total	prob
## 1	Don't know	3	17	0.1765
## 2	Excellent	1035	3855	0.2685
## 3	Fair	170	595	0.2857
## 4	Good	1022	3303	0.3094
## 5	Poor	25	102	0.2451
## 6	Very good	1115	4301	0.2592



Discussion

Weaknesses

Next Steps

In the next step, we may collect more data in other countries with different economic conditions and cultural backgrounds, like the country they immigrate from, the religious background to show the results more generally. Also, we can use the principal component analysis to narrow the variables of the model, which are strongly correlated with the interest. It is also worth considering to build a neural network model to make predictions based on our data, since the NN model is more robust to random cases.

References

1. NCHS Pressroom - 1997 Fact Sheet - Mothers Education and Birth Rate. (2009, November 17). Retrieved October 17, 2020, from <https://www.cdc.gov/nchs/pressroom/97facts/edu2birt.htm>
2. Pobric, A., & Robinson, G. M. (2015). Population ageing and low fertility: Recent demographic changes in Bosnia and Herzegovina. *Journal of Population Research*, 32(1), 23-43. doi:10.1007/s12546-014-9141-5
3. Martinez, G., Ph.D, Daniels, K., Ph.D, & Chandra, A., Ph.D. (2012, April 12). Fertility of Men and Women Aged 15–44 Years in the United States: National Survey of Family Growth, 2006–2010. Retrieved October 17, 2020, from <https://www.cdc.gov/nchs/data/nhsr/nhsr051.pdf>

4. Fertility rate, total (births per woman). (n.d.). Retrieved October 17, 2020, from <https://data.worldbank.org/indicator/SP.DYN.TFRT.IN>
5. 2017 General Social Survey (GSS): Families Cycle 31. (2017). Retrieved October 17, 2020, from Statistics Canada.
6. Beaupré, P. (2020). General Social Survey Cycle 31 : Families Public Use Microdata File Documentation and User's Guide. Ottawa, Canada: Authority of the Minister responsible for Statistics Canada. Retrieved October 17, 2020, from https://sda-artsci-utoronto-ca.myaccess.library.utoronto.ca/sdaweb/dli2/gss/gss31/gss31/more_doc/GSS31_User_Guide.pdf
7. Wickham et al., (2019). Welcome to the tidyverse. *Journal of Open Source Software*, 4(43), 1686, <https://doi.org/10.21105/joss.01686>
8. T. Lumley (2020) "survey: analysis of complex survey samples". R package version 4.0.
9. T. Lumley (2004) Analysis of complex survey samples. *Journal of Statistical Software* 9(1): 1-19
10. T. Lumley (2010) Complex Surveys: A Guide to Analysis Using R. John Wiley and Sons.
11. User1489975user1489975 1, BenBarnesBenBarnes 17.3k66 gold badges5151 silver badges7070 bronze badges, Mnelmnel 103k2525 gold badges241241 silver badges240240 bronze badges, Amrrsamrrs 5, RnoobRnoob 9031010 silver badges1212 bronze badges, DroneyDroney 11111 silver badge44 bronze badges, . . . Luchao QiLuchao Qi 3122 bronze badges. (1961, October 01). Omit rows containing specific column of NA. Retrieved October 18, 2020, from <https://stackoverflow.com/questions/11254524/omit-rows-containing-specific-column-of-na>
12. Ggplot2 barplots : Quick start guide - R software and data visualization. (n.d.). Retrieved October 18, 2020, from <http://www.sthda.com/english/wiki/ggplot2-barplots-quick-start-guide-r-software-and-data-visualization>