**2020 – 2021 Project Continuous Data Analysis/Statistical Modelling**

Group 3: Dereje Mengist Belete, Feihong Du, Echo Capwell Forbang, Serge Martin Nkoumnga, Lin Tang, Shengmin Zhang

# Introduction

United States Presidential election is widely concerned throughout the world. The republican candidate Trump Donald defeated the democratic candidate Hillary Clinton and won the election in 2016, which is far away from the result of the public opinion poll. In this study, we would like to investigate the factors, like socio-economic and demographic characteristics, which might impact the vote rate to Trump in 2016 election.

The dataset we used in this study is from the Kaggle website. In this dataset, we have 3141 observation indicating the election result and status of different areas of America. We implemented the regression approach which constituted a linear regression and a categorical regression of selected variables on the percentage of votes for the Republican party (GOP), to gain some insight in the association of certain variables with the vote rate of Trump.

# Protocol

### Research Question

For this project we will analyse part of the "2016 US Presidential Election Dataset"to examine the effect of the per capital money income obtained in the past 12 months on the amount of the percentage of GOP votes during the 2016 presidential elections, while accounting for other potential covariables, using a multiple linear regression model.

Did the economic characteristic (per capital money income obtained in the past 12 months) influence the percentage of GOP votes during the 2016 presidential elections?

What is the impact of socio-economic and demographic characteristics including race, gender, education, household, veterans and number of firms on the percentage of GOP votes during the 2016 presidential elections?

Moreover, we would like to find a way to predict the undecided counties (vote rate: 48%-52%) where our candidate may put more effort to on next election.

**Study Method**

A new dataset will be created based on the variables we choose to study. And then the distribution of each individual variable will be examined by univariate procedure, if it is necessary the missing value will be removed and the skewed data will be transformed. Bivariate relationships between the continuous variables will be examined with correlation and scatterplot matrix. Then we will split the dataset to training dataset and testing dataset (50:50) and normalize them respectively. The interactions between two variables will be created as new columns.

After the data preprocessing, the forward stepwise selection will be applied on the training dataset to find our "final model". After all the available variables access to the model, we will check whether the "final model" can fulfill all the assumptions.

In order to build logistic regression model, we will categorize the vote rate into 2 level: win (>50%) and lose(<50%). Then Compare your results/conclusions with those of the linear model to adapt the final model.

Fit the final linear model on the "testing" dataset to evaluate the model performance.

**Data Exploration**

We choose the income (per capital money income obtained in the past 12 months) as our key variable and other variables are white race, household, gender, education, firms and veterans. The distribution of all individual variables is examined by the univariate procedure. Since the firms and veterans are heavily right skewed (skewness: 16.15, 8.27 respectively), the log-transform is applied to those two variables. Missing value is checked by the mean procedure and there is no missing value in our dataset. Bivariate relationships between the continuous variables will be examined with correlation and scatterplot matrix.

We split the whole dataset data into training (1571 observations) and test (1570 observations) datasets with a proportion of 50:50. Both datasets were then standardized independently to have uniformity in interpreting our coefficients.

**Simple linear regression**

Simple linear regression (SLR) was used to assess the performance of key predictor (income) affecting the percentage of votes obtained by GOP candidate, the formula is as follows:

$$y_i = \beta_0 + \beta_1 x_i + \varepsilon_i \qquad \varepsilon_i \sim N(0, \delta^2), \ i = 1, \dots, n$$

$y_i = vote\ rate$ (Percentage of votes for Republican Party)

$x_i =$ income (Per capita money income in past 12 months (dollars), 2009-2013)

The results showed that income affect the votes obtained by GOP candidate marginally

significant (p-value <0.0001), our model then can be written as

$$y_i = 0.767 - 0.0000055x_i + \varepsilon_i \qquad \varepsilon_i \sim N(0, \delta^2) , i = 1, \dots, n$$

which means income increase one unit will result in the decrease of GOP support of 0.0000055%.

**Multiple linear regression**

We used forward stepwise regression method to build our model, absolute t-value, F-value and p-value were considered as the access criterions when we judge whether a new variable could be included or removed. Accordingly, newly added variable is of lowest residual sum of squares (SSR) and of highest R square, absolute t-values, F-values and p-value at 0.05 significance level in each step. Multicollinearity and confounding effect were checked by VIFs and p-value after each selection. With the addition of a new relevant variable included into our model, the already fitted variables did not change much in the significance and magnitude of the coefficients, hence no confounding effects. Variables which did not fulfill those criteria were omitted from the model.

According to the criterion we set, our model was built with main effects of income, white race, veterans, and firms. A similar approach was implemented to include influential two-way interactions terms into the model. The following interactions were included: household*education, white*education, income*firms, white*household, income*household, income*gender, education*veterans, white*firms and firms*veterans.

**Checking assumptions**

Multicollinearity: There is no multicollinearity amongst the predictors according to the VIF method, where we set VIF < 10. Normality: We used the QQ-plot of (studentized) residuals to check for the normality. There were slight deviations at the tails observed. However, normality can be assumed. Linearity: It was verified through the plot of (studentized) residuals versus predicted values. Points look randomly scattered around 0. No evidence of nonlinear pattern, hence linearity could be assumed. Furthermore, equality of variance (Homoscedasticity) was checked by examining the plot of (squared) residuals versus predicted values. Plot indicates that points equally randomly scattered around 0 point and therefore no evidence of heteroscedasticity.

**Checking for outliers**

Using the Cook's Distance threshold (4/n = 0.0025, n denotes the number of observations) to remove the outliers till convergence. The t-value, F-value and VIF value are the basis of our judgement. There were 93 outliers detected and removed in first time. Then the model was refitted by the new dataset without outliers. When checking the estimate coeffecients' siginificance, the p-value of white*household and firms*veterans are 0.4253 and 0.5833. These two variables were then removed from

our model. Checking the outlier again and we found only 2 outliers which are slight beyond the threshold. The assumptions for linear regression were assessed once more and arrived at the conclusion that, all assumptions seem to be fulfilled (Appendix II).

## Interpretation

A comparation across all statistics presented in Table 1 indicated that race (white) was the strongest direct predictor of gop across mutiple indices, followed by income, Intwhiteeducation, and Veterans, Intincomehousehold, Inteducationveterans, Intincomefirms, firms, Inthouseholdeducation, Intincomegender, Intwhitefirms. Race obtained the largest T-value (t-value = 18.73, p < .0001), demonstrating that it made the largest contribution to the regression equation, while

*Table 1 Statistical results of final model*

| Parameter Estimates | | | | | | |
|---|---|---|---|---|---|---|
| Variable | DF | Parameter Estimate | Standard Error | t Value | Pr > \|t\| | Variance Inflation |
| Intercept | 1 | 0.66781 | 0.00292 | 228.53 | <.0001 | 0 |
| ($x_1$) Income | 1 | -0.00000585 | 5.656001E-7 | -10.35 | <.0001 | 1.43514 |
| ($x_2$) White | 1 | 0.00365 | 0.00019506 | 18.73 | <.0001 | 1.48792 |
| ($x_3$) Veterans | 1 | -0.02913 | 0.00300 | -9.71 | <.0001 | 3.06874 |
| ($x_4$) Firms | 1 | -0.00793 | 0.00191 | -4.15 | <.0001 | 2.99856 |
| ($x_5$) Household*Education | 1 | 0.00798 | 0.00202 | 3.94 | <.0001 | 1.65466 |
| ($x_6$) White*Education | 1 | -0.00026287 | 0.00002623 | -10.02 | <.0001 | 1.43653 |
| ($x_7$) Income*Firms | 1 | -0.00000121 | 2.679091E-7 | -4.51 | <.0001 | 1.94733 |
| ($x_8$) Income*Household | 1 | 0.00001319 | 0.00000265 | 4.98 | <.0001 | 1.74923 |
| ($x_9$) Income*Gender | 1 | 6.895835E-7 | 2.030597E-7 | 3.40 | 0.0007 | 1.26568 |
| ($x_{10}$) Education*Veterans | 1 | -0.00152 | 0.00034197 | -4.46 | <.0001 | 1.58034 |
| ($x_{11}$) White*Firms | 1 | 0.00028353 | 0.00009679 | 2.93 | 0.0034 | 1.39653 |

holding all other predictor variables constant. The squared structure coefficient ($R^2$ = .2755) demonstrated that race (white) explained the largest amount (27.55%) of variance in y, the predicted values of gop.
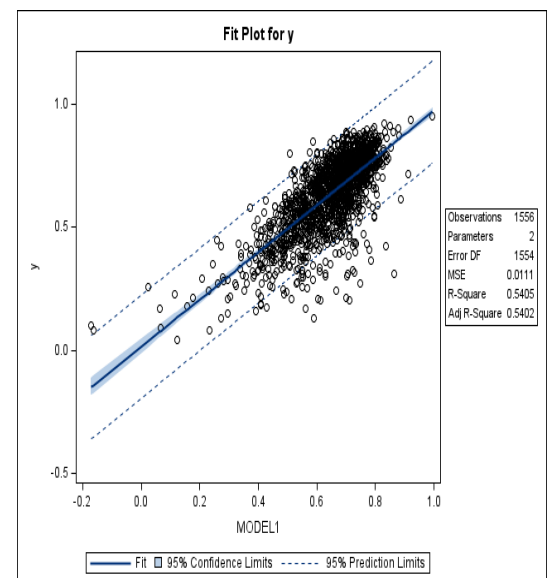
A good and standard statistical package has been used in this model, the R-squared value of the final multiple linear model is .5716 with all p < 0.01, which is quite an excellent, highly statistically significant result. About 58% ($R^2$%) of variations in y can be attributed to variations in x. Therefore y is reliably predictably with the multiple linear regression model:

$$Y_i = 0.67 - 5.85 * 10^{-6}x_1 + 3.65 * 10^{-3}x_2 - 0.03x_3 \\ - 0.01x_4 + 0.01x_5 - 2.63 * 10^{-4}x_6 \\ - 1.21 * 10^{-6} x_7 + 1.32 * 10^{-5}x_8 + 6.9 \\ * 10^{-7}x_9 - 1.52 * 10^{-3}x_{10} + 2.84 \\ * 10^{-4}x_{11}$$

By applying the final model we built on the testing dataset, the results ($R^2$=.5405, MSE = .0111) still shows its robust and reliance. The corresponding plot is shown as below.



Fit Plot for y

| | |
|---|---|
| Observations | 1556 |
| Parameters | 2 |
| Error DF | 1554 |
| MSE | 0.0111 |
| R-Square | 0.5405 |
| Adj R-Square | 0.5402 |

The overall findings supported how both race (white) was the most significant direct contributor and income was the second most important direct contributor to predicting variance in gop, as reflected across different T-values, F-values and VIF. This might be because the white people are currently the major communities of voters, and income has a negative contribution on y, the predicted values of gop, which is corresponding to the reasearch that indicates, in poor counties, income is associated

with Republican voting, while in many rich states, the relation between income and vote choice is nearly zero[1-5].
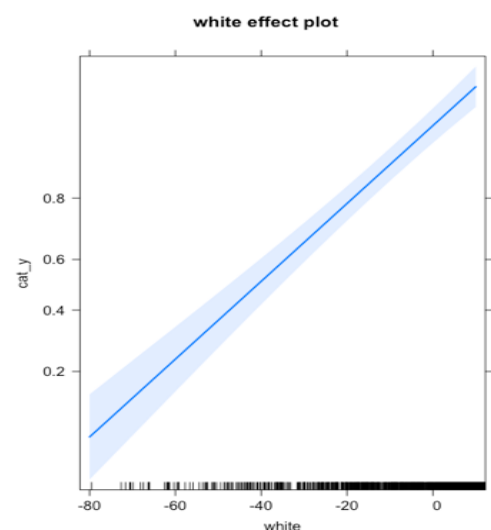
**General Linear Model (GLM)**

We dichotomized the response variable into two categories: win (vote rate greater than 50%) and lose (vote rate less than 50%), representing by 1 and 0 respectively. The GLM is based on the final model we have found in continuous part, fitting by the dichotomized training dataset. The formula of logistic regression can be written as following:

$$P(Y=1) = \frac{1}{1+e^{-(\beta_0+\beta_1 x_1+\beta_2 x_2+...+\beta_p x_p)}}$$

For easier understanding, it also can be transformed as following equation:

$$log(\frac{P(Y=1)}{1-P(Y=1)}) = log(\frac{P(Y=1)}{P(y=0)}) = \beta_0 + beta_1 x_1 + ... + \beta_p x_p$$

Obviously, the left part of the equation is the log odds. We choose white race as an example to interpret our GLM. The effect of white race has linear relationship with the cat_y, displayed as right figure, in which the slope is the estimate coefficient of white race (0.04949103). The conditional odds of the white race would be exp(0.04949103) = 1.050736, indicating that when other variables keep constants the odds would increase 5.0736% if the white race increase one unit. What should be noticed here is that the firms and veterans had been log transformed in previous data process, so when discussing about those two variables, the conditional odds ratios would be the estimate coefficients directly. Using the testing dataset for prediction, the accuracy of our model is 87.66% and the precision of win and lose are 88.8% and 79% respectively.



Based on the estimate coeffecients of our fitted model, the interaction of househould and education (estimate coeffecient 0.219) shows highest positive effect on the vote rate to our candidate in 2016, while the log-transformed firms (estimate coefficient -0.8772) shows the highest negative effect on it. Moreover, the income and log-transformed veterans do not show significant effects to the vote rate in our model.

**Discussion**

Our study performed on county level data on the US elections 2016. In our study, we fitted the final model with training dataset and then applied it to test dataset for prediction. Overlap, t-test and variance test were used to test the similarity and difference. Overlap ranges from 0 (no overlap) to 1 (complete overlap). Since our sample size is greater than 75, we use the Dhat4 = 0.87 of overlap coefficient which implies, the boundaries of both response variables from actual and predicted dataset greatly coincide with one another. Significant difference was found when conducting t- and var- test, which means the prediction is marginally different from observations. Nevertheless, we may improve our model by including more relevant variables in the future.

In order to get the "undecided" counties which we need to focus for the next election, the linear regression model was applied to train the whole dataset with the variables already introduced. We then performed function of 'predict' to get perdicted values corresponding to each observation. Afterwards, we categoried the predicted values and defined interval located between 48%-52% as "undecided" situations. According to those criterion, we got 145 counties among 1738 counties, where the candiate may put more effort in the next election (see Appendix III).

**Reference**

1. Why do we need to re-use training parameters to transform test data?
https://sebastianraschka.com/faq/docs/scale-training-test.html
2. Jim Frost (2019), Using Confidence Intervals to Compare Means.
https://statisticsbyjim.com/hypothesis-testing/confidence-intervals-compare-means/
3. Dr. Frank Wood (2010), Inference in Regression Analysis.
http://www.stat.columbia.edu/~fwood/Teaching/w4315/Spring2010/lecture_4.pdf
4. Nathans, Laura L., Oswald, Frederick L. and Nimon, Kim. "Interpreting Multiple Linear Regression: A Guidebook of Variable Importance." *Practical Assessment, Research & Evaluation,* 17, no. 9 (2012) Practical Assessment, Research & Evaluation:
https://hdl.handle.net/1911/71096

5. Andrew Gelman, Lane Kenworthy, Yu-Sung Su(2010). Income Inequality and Partisan Voting in the United States. 91(5):1204-1219.

**Schedule**

The distribution of tasks of the whole program is as follow:

- 2020.11.28 Get an overall understanding of the research topic and work for the protocol together.

- 2020.11.28 - 2020.12.6 Feihong Du, Lin Tang and Shengmin Zhang work for the code part.

- 2020.12.7 Discuss and solve problems for the code part and work together to form the final version.

- 2020.12.7 – 2020.12.10 Dereje Mengist Belete, Echo Capwell Forbang and Serge Martin Nkoumnga work for the manuscript.

- 2020.12.10 Feihong Du, Lin Tang and Shengmin Zhang work together to revise the manuscript.

- 2020.12.11 All members work together to form the final report.

**APPENDIX**

**Appendix I: Variables selected and their descriptions**

| Variable | | Codes | Description |
|---|---|---|---|
| Outcome Covariates (X) | GOP | Per vote gop | "Percentage of votes for Republican party" |
| | **Income** | INC910213 | "Per capita money income in past 12 months (2013 dollars), 2009-2013" |
| | **Household** | HSD310213 | "Household, 2009-2013" |
| | **Gender** | SEX255214 | "Female persons, percent, 2014" |
| | **Education** | EDU635213 | "High school graduate or higher, percent of persons age 25+, 2009-2013" |
| | **Firms** | SBO001207 | "Total number of firms, 2007" |
| | **Veterans** | VET605213 | "Veterans, 2009-2013" |
| | **White** | RHI125214 | "White alone, percent, 2014" |

**Appendix II: Plots of linear model assumptions**

**Appendix III: Counties need to be treated seriously for the next election**

| | | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|
| Adams | Brazor | Clevel | Greenv | Indian | Linn C | Nantuc | Richmo | St. He | Warren |
| Aiken | Calcas | Cumber | Grenad | Island | Lownde | Nash C | Roosev | St. Ja | Washin |
| Alachu | Calhou | Danvil | Hamilt | Jasper | Lucas | Nevada | Russel | St. Ma | Washoe |
| Alaska | Cape M | Darlin | Hampsh | Jeffer | Macomb | Ocean | San Ju | St. Ta | Wayne |
| Albema | Cass C | Delawa | Hardee | Kern C | Madera | Okaloo | Sangam | Talbot | Westmo |
| Anoka | Champa | Early | Harris | Kings | Madiso | Olmste | Santa | Talbot | Will C |
| Anson | Chicka | El Dor | Haywoo | Knox C | Manate | Orange | Scotla | Taliaf | Willia |
| Atkins | Chicot | Escamb | Hendry | La Sal | Mareng | Ouachi | Sedgwi | Terreb | Wilson |
| Atlant | Chitte | Fairfa | Hidalg | Lake C | Martin | Passai | Semino | Thurst | York C |
| Attala | Cibola | Fayett | Hoke C | Lancas | Meriwe | Peoria | Shelby | Tollan | |
| Bell C | Clacka | Frankl | Holmes | Larime | Minera | Pima C | Sonoma | Troup | |
| Berkel | Claibo | Galves | Hopewe | Lauder | Monroe | Pitt C | Spaldi | Tulare | |
| Bernal | Claren | Genese | Housto | Lee Co | Monter | Platte | Spokan | Tuscal | |
| Bexar | Clarke | Glouce | Iberia | Lenoir | Montgo | Provid | St. Ch | Ventur | |
| Bossie | Clear | Greene | Imperi | Lexing | Moreho | Rapide | St. Fr | Waltha | |

```sas
1  /*Statistic Molding Project 2020*/
2
3
4  %let dir='/folders/myfolders/StatmodProject2020/';
5  libname moldproj &dir;
6
7  proc import out=moldproj.county
8         datafile="/folders/myfolders/StatmodProject2020/county_facts_dictionary.csv" dbms=csv replace;
9     getnames=yes;
10     datarow=2;
11 run;
12
13 proc import out=moldproj.election
14         datafile="/folders/myfolders/StatmodProject2020/US_County_Level_Presidential_Results_12-16_.csv"
15         dbms=csv replace;
16     getnames=yes;
17     datarow=2;
18 run;
19 proc print data=moldproj.election (obs=5);
20 run;
21
22 /*Choose relevant variables and create a new dataset*/
23
24
25 %let x = white income household gender education firms veterans;
26
27 data election (keep=county y &x);
28     set moldproj.election;
29     rename county_name=county per_gop_2016=y RHI125214=white INC910213=income
30         HSD310213=household SEX255214=gender EDU635213=education SBO001207=firms
31         VET605213=veterans;
32 run;
33
34 proc print data=election (obs=10);
35 run;
36
37 *#####################################;
38 *### STEP 1: Descriptives Analysis ###;
39 *#####################################;
40
41 /*1. Univariate exploration of the data*/
42 proc univariate plot data=election;
43     var y &x;
44 run;
45
46 /* Log-transform the right skewed data:  */
47 /* firms (skewness: 16.1494055) */
48 /* veterans (skewness: 8.26233385) */
49 data election;
50     set election;
51
52     if firms ne 0 then
53         firms=log(firms);
54     veterans=log(veterans);
55 run;
56
57 proc means data=election nmiss;
58     var y &x;
59 run;
60
61 proc univariate plot data=election;
62     var &x;
63 run;
64
65 /*2. Bivariate relationships*/
66 title "Correlation matrix and scatter-plots";
67
68 proc sgscatter data=election;
69     matrix &x / diagonal=(histogram normal);
70 run;
71
72 title;
73
74 data election;
75     set election;
76     deel=firms/veterans;
77 run;
```

```sas
73  proc sgscatter data=election;
74      matrix &x deel/ diagonal=(histogram normal);
75  run;

76
77  title;

78  proc corr data=election nosimple;
79      var &x;
80  run;

81
82  *#####################################;
83  *### STEP 2: Data spliting 50%:50% ###;
84  *#####################################;

85
86  proc surveyselect data=election out=train_test_split method=srs samprate=0.5
            outall seed=123 noprint;
87      samplingunit y;
88  run;

89
90  data train_set;
91      set train_test_split;
92      where Selected=1;
93  run;

94  data test_set;
95      set train_test_split;
96      where Selected=0;
97  run;

98
99  /* Standardize data! */
100 proc stdize data=train_set method=mean out=cent_train;
101     var &x;
102 run;

103 proc stdize data=test_set method=mean out=cent_test;
104     var &x;
105 run;

106
107 proc print data=cent_test;
108 run;

109
110 data cent_train;
111     set cent_train;
112     intwhiteincome=white*income;
113     intwhitehousehold=white*household;
114     intwhitegender=white*gender;
115     intwhiteeducation=white*education;
116     intwhitefirms=white*firms;
117     intwhiteveterans=white*veterans;
118     intincomehousehold=income*household;
119     intincomegender=income*gender;
120     intincomeeducation=income*education;
121     intincomefirms=income*firms;
122     intincomeveterans=income*veterans;
123     inthouseholdgender=household*gender;
124     inthouseholdeducation=household*education;
125     inthouseholdfirms=household*firms;
126     inthouseholdveterans=household*veterans;
127     intgendereducation=gender*education;
128     intgenderfirms=gender*firms;
129     intgenderveterans=gender*veterans;
130     inteducationfirms=education*firms;
131     inteducationveterans=education*veterans;
132     intfirmsveterans=firms*veterans;
133 run;

131
132 data cent_test;
133     set cent_test;
134     intwhiteincome=white*income;
135     intwhitehousehold=white*household;
136     intwhitegender=white*gender;
137     intwhiteeducation=white*education;
138     intwhitefirms=white*firms;
139     intwhiteveterans=white*veterans;
140     intincomehousehold=income*household;
141     intincomegender=income*gender;
142     intincomeeducation=income*education;
143     intincomefirms=income*firms;
144     intincomeveterans=income*veterans;
145     inthouseholdgender=household*gender;
        inthouseholdeducation=household*education;
        inthouseholdfirms=household*firms;
```

```
146        inthouseholdveterans=household*veterans;
147        intgendereducation=gender*education;
148        intgenderfirms=gender*firms;
149        intgenderveterans=gender*veterans;
150        inteducationfirms=education*firms;
151        inteducationveterans=education*veterans;
152        intfirmsveterans=firms*veterans;
153   run;

154   %let int_x = intwhiteincome intwhitehousehold intwhitegender intwhiteeducation intwhitefirms intwhiteveterans

156           intincomehousehold intincomegender intincomeeducation intincomefirms intincomeveterans

158           inthouseholdgender inthouseholdeducation inthouseholdfirms inthouseholdveterans

160           intgendereducation intgenderfirms intgenderveterans

162           inteducationfirms inteducationveterans

163           intfirmsveterans;
164   *#####################################;
165   *### STEP 3: Forward selection    ###;
166   *#####################################;

167   /* 3.1 Simple Linear regression */
168   /* t:0.0376 F:61.83 */
169   proc reg data=train_set;
170       model y=income / vif;
171       title "Simple linear regression";
172       run;
173   quit;

174   title;

176   /* First extra predictor*/
177   /* income need to stay in the model */
178   /* 3.2 determine the second fixed variable via t-value and p-value */
179   /* t:28.40 F:449.92 */
180   proc reg data=cent_train;
181       model y=income white;
182       run;
183   quit;

184   /* t:-9.36 F:76.41 */
185   proc reg data=cent_train;
186       model y=income household;
187       run;
188   quit;

189
190   /* t:-5.10 F:44.42   */
191   proc reg data=cent_train;
192       model y=income gender;
193       run;
        quit;

194
195   /* t:3.22 F:36.28 */
196   proc reg data=cent_train;
197       model y=income education;
198       run;
199   quit;

200   /* t:-15.92 F:162.64 */
201   proc reg data=cent_train;
202       model y=income firms;
203       run;
204   quit;

205
206   /* t:-18.27 F:204.34 */
207   proc reg data=cent_train;
208       model y=income veterans;
209       run;
        quit;

210
211   /* we choose white as a candidate varible, now check for it */
212   /* r_square:0.3626 vif < 10*/
213   proc reg data=cent_train;
214       model y=income white / vif;
215       run;
216   quit;

217   /* 3.3 Add white to model, determine the third fixed variable via t-value and p-value */
218   /*t:-2.28 F:302.47 */
```

```
219  proc reg data=cent_train;
220      model y=income white household;
221      run;
222  quit;
223
     /*t:-5.05 F:315.56 */
224  proc reg data=cent_train;
225      model y=income white gender;
226      run;
227  quit;
228
229  /*t:-1.25 F:300.57 */
230  proc reg data=cent_train;
231      model y=income white education;
232      run;
233  quit;
234  /*t:-14.67 F:412.37 */
235  proc reg data=cent_train;
236      model y=income white firms;
237      run;
238  quit;
239
     /*t:-17.04 F:451.66 */
240  proc reg data=cent_train;
241      model y=income white veterans;
242      run;
243  quit;
244
245  /* we choose veterans as a candidate varible, now check for it */
246  /* r_square:0.4605 vif<10*/
247  proc reg data=cent_train;
248      model y=income white veterans / vif;
249      run;
250  quit;
251  /* 3.3 Add veterans to model, determine the fourth fixed variable via t-value and p-value */
252  /*t:0.09 F:338.54 p:0.9249*/
253  proc reg data=cent_train;
254      model y=income white veterans household;
255      run;
256  quit;
257  /*t:-1.15 F:339.15 p:0.2485*/
258  proc reg data=cent_train;
259      model y=income white veterans gender;
260      run;
261  quit;
262  /*t:0.52 F:338.66 p:0.6029*/
263  proc reg data=cent_train;
264      model y=income white veterans education;
265      run;
266  quit;
267
268  /*t:-3.33 F:343.69 */
269  proc reg data=cent_train;
270      model y=income white veterans firms;
271      run;
272  quit;
273
274  /* we choose firms as a candidate varible, now check for it */
275  /* r_square:0.4639 vif<10*/
276  proc reg data=cent_train;
277      model y=income white veterans firms/vif;
278      run;
279  quit;
280  /* 3.5 Add firms, determine the fifth fixed variable via t-value f-value and p-value (reject */
281  /*r:0.22 F:274.79 p:0.8236*/
282  proc reg data=cent_train;
283      model y=income white veterans firms household;
284      run;
285  quit;
286  /*r:-0.94 F:275.11 p:0.3458*/
287  proc reg data=cent_train;
288      model y=income white veterans firms gender;
289      run;
290  quit;
291  /*r:-0.78 F:275.01 p:0.4330*/
```

```sas
292   proc reg data=cent_train;
293       model y=income white veterans firms education;
294       run;
295   quit;
296
297   /* 3.6 fix model with income white veterans and firms, r_square:0.4639 vif < 10 */
298   /* Add one interaction variable */
299   proc reg data=cent_train;
300       model y=income white veterans firms intwhiteincome;
301       run;
302   quit;
303
304   proc reg data=cent_train;
305       model y=income white veterans firms intwhitehousehold;
306       run;
307   quit;
308
309   proc reg data=cent_train;
310       model y=income white veterans firms intwhitegender;
311       run;
312   quit;
313
314   proc reg data=cent_train;
315       model y=income white veterans firms intwhiteeducation;
316       run;
317   quit;
318
319   proc reg data=cent_train;
320       model y=income white veterans firms intwhitefirms;
321       run;
322   quit;
323
324   proc reg data=cent_train;
325       model y=income white veterans firms intwhiteveterans;
326       run;
327   quit;
328
329   proc reg data=cent_train;
330       model y=income white veterans firms intincomehousehold;
331       run;
332   quit;
333
334   proc reg data=cent_train;
335       model y=income white veterans firms intincomegender;
336       run;
337   quit;
338
339   proc reg data=cent_train;
340       model y=income white veterans firms intincomeeducation;
341       run;
342   quit;
343
344   proc reg data=cent_train;
345       model y=income white veterans firms intincomefirms;
346       run;
347   quit;
348
349   proc reg data=cent_train;
350       model y=income white veterans firms intincomeveterans;
351       run;
352   quit;
353
354   proc reg data=cent_train;
355       model y=income white veterans firms inthouseholdgender;
356       run;
357   quit;
358
359   proc reg data=cent_train;
360       model y=white income veterans firms inthouseholdeducation;
361       run;
362   quit;
363
364   proc reg data=cent_train;
365       model y=income white veterans firms inthouseholdfirms;
366       run;
367   quit;
368
369   proc reg data=cent_train;
370       model y=income white veterans firms inthouseholdveterans;
371       run;
372   quit;
```

```
365   proc reg data=cent_train;
366       model y=income white veterans firms intgendereducation;
367       run;
368   quit;
369
370   proc reg data=cent_train;
371       model y=income white veterans firms intgenderfirms;
372       run;
      quit;
373
374   proc reg data=cent_train;
375       model y=income white veterans firms intgenderveterans;
376       run;
377   quit;
378
379   proc reg data=cent_train;
380       model y=income white veterans firms inteducationfirms;
381       run;
      quit;
382
383   proc reg data=cent_train;
384       model y=income white veterans firms inteducationveterans;
385       run;
386   quit;
387
388   proc reg data=cent_train;
389       model y=income white veterans firms intfirmsveterans;
390       run;
      quit;
391
392   /* we choose inthouseholdeducation as a candidate interaction varible, now check for it */
393   /* r_square:0.5173 vif<10 */
394   proc reg data=cent_train;
395       model y=income white veterans firms inthouseholdeducation/ vif;
396       run;
397       /* 3.7 add inthouseholdeducation, determine the second interaction variable*/
398   proc reg data=cent_train;
399       model y=income white veterans firms inthouseholdeducation intwhiteincome;
400       run;
401   quit;
402
403   proc reg data=cent_train;
404       model y=income white veterans firms inthouseholdeducation intwhitehousehold;
405       run;
      quit;
406
407   proc reg data=cent_train;
408       model y=income white veterans firms inthouseholdeducation intwhitegender;
409       run;
410   quit;
411
412   proc reg data=cent_train;
413       model y=income white veterans firms inthouseholdeducation intwhiteeducation;
414       run;
      quit;
415
416   proc reg data=cent_train;
417       model y=income white veterans firms inthouseholdeducation intwhitefirms;
418       run;
      quit;
419
420   proc reg data=cent_train;
421       model y=income white veterans firms inthouseholdeducation intwhiteveterans;
422       run;
423   quit;
424
425   proc reg data=cent_train;
426       model y=income white veterans firms inthouseholdeducation intincomehousehold;
427       run;
      quit;
428
429   proc reg data=cent_train;
430       model y=income white veterans firms inthouseholdeducation intincomegender;
431       run;
432   quit;
433
434   proc reg data=cent_train;
435       model y=income white veterans firms inthouseholdeducation intincomeeducation;
436       run;
      quit;
437
```

```sas
438   proc reg data=cent_train;
439       model y=income white veterans firms inthouseholdeducation intincomefirms;
440       run;
441   quit;
442
443   proc reg data=cent_train;
444       model y=income white veterans firms inthouseholdeducation intincomeveterans;
445       run;
      quit;
446
447   proc reg data=cent_train;
448       model y=income white veterans firms inthouseholdeducation inthouseholdgender;
449       run;
450   quit;
451
452   proc reg data=cent_train;
453       model y=income white veterans firms inthouseholdeducation inthouseholdfirms;
454       run;
      quit;
455
456   proc reg data=cent_train;
457       model y=income white veterans firms inthouseholdeducation inthouseholdveterans;
458       run;
459   quit;
460
461   proc reg data=cent_train;
462       model y=income white veterans firms inthouseholdeducation intgendereducation;
463       run;
      quit;
464
465   proc reg data=cent_train;
466       model y=income white veterans firms inthouseholdeducation intgenderfirms;
467       run;
468   quit;
469
470   proc reg data=cent_train;
471       model y=income white veterans firms inthouseholdeducation intgenderveterans;
472       run;
      quit;
473
474   proc reg data=cent_train;
475       model y=income white veterans firms inthouseholdeducation inteducationfirms;
476       run;
      quit;
477
478   proc reg data=cent_train;
479       model y=income white veterans firms inthouseholdeducation inteducationveterans;
480       run;
481   quit;
482
483   proc reg data=cent_train;
484       model y=income white veterans firms inthouseholdeducation intfirmsveterans;
485       run;
      quit;
486
487   proc reg data=cent_train;
488       model y=income white veterans Firms inthouseholdeducation intincomeeducation;
489       run;
490
491   proc reg data=cent_train;
492       model y=income white veterans firms inthouseholdeducation intincomeeducation;
493       run;
494
495       /* we choose intwhiteeducation as a candidate interaction varible, now check for it */
          /*r:0.5300 vif<10*/
496   proc reg data=cent_train;
497       model y=income white veterans firms inthouseholdeducation intwhiteeducation /
498           vif;
499       run;
500
          /* 3.8 add intwhiteeducation, determine the third interaction variable*/
501   proc reg data=cent_train;
502       model y=income white veterans firms inthouseholdeducation intwhiteeducation
503           intwhiteincome;
504       run;
505   quit;
506
507   proc reg data=cent_train;
508       model y=income white veterans firms inthouseholdeducation intwhiteeducation
509           intwhitehousehold;
          run;
510   quit;
```

```
511
512   proc reg data=cent_train;
513       model y=income white veterans firms inthouseholdeducation intwhiteeducation
514          intwhitegender;
515       run;
516   quit;
517   proc reg data=cent_train;
518       model y=income white veterans firms inthouseholdeducation intwhiteeducation
519          intwhitefirms;
520       run;
521   quit;
522
523   proc reg data=cent_train;
524       model y=income white veterans firms inthouseholdeducation intwhiteeducation
525          intwhiteveterans;
526       run;
527   quit;
528   proc reg data=cent_train;
529       model y=income white veterans firms inthouseholdeducation intwhiteeducation
530          intincomehousehold;
531       run;
532   quit;
533   proc reg data=cent_train;
534       model y=income white veterans firms inthouseholdeducation intwhiteeducation
535          intincomegender;
536       run;
537   quit;
538
539   proc reg data=cent_train;
540       model y=income white veterans firms inthouseholdeducation intwhiteeducation
541          intincomeeducation;
542       run;
543   quit;
544   proc reg data=cent_train;
545       model y=income white veterans firms inthouseholdeducation intwhiteeducation
546          intincomefirms;
547       run;
548   quit;
549   proc reg data=cent_train;
550       model y=income white veterans firms inthouseholdeducation intwhiteeducation
551          intincomeveterans;
552       run;
553   quit;
554
555   proc reg data=cent_train;
556       model y=income white veterans firms inthouseholdeducation intwhiteeducation
557          inthouseholdgender;
558       run;
559   quit;
560   proc reg data=cent_train;
561       model y=income white veterans firms inthouseholdeducation intwhiteeducation
562          inthouseholdfirms;
563       run;
564   quit;
565
566   proc reg data=cent_train;
567       model y=income white veterans firms inthouseholdeducation intwhiteeducation
568          inthouseholdveterans;
569       run;
570   quit;
571   proc reg data=cent_train;
572       model y=income white veterans firms inthouseholdeducation intwhiteeducation
573          intgendereducation;
574       run;
575   quit;
576   proc reg data=cent_train;
577       model y=income white veterans firms inthouseholdeducation intwhiteeducation
578          intgenderfirms;
579       run;
580   quit;
581   proc reg data=cent_train;
582       model y=income white veterans firms inthouseholdeducation intwhiteeducation
583          intgenderveterans;
```

```
584        run;
585    quit;
586
587    proc reg data=cent_train;
588        model y=income white veterans firms inthouseholdeducation intwhiteeducation
               inteducationfirms;
589        run;
590    quit;
591
592    proc reg data=cent_train;
593        model y=income white veterans firms inthouseholdeducation intwhiteeducation
               inteducationveterans;
594        run;
595
596    quit;
597
598    proc reg data=cent_train;
599        model y=income white veterans firms inthouseholdeducation intwhiteeducation
               intfirmsveterans;
600        run;
601    quit;
602
603    /* we choose intincomefirms as a candidate interaction varible, now check for it */
604    /*r_square:0.5345 vif<10*/
605    proc reg data=cent_train;
606        model y=white income veterans firms inthouseholdeducation intwhiteeducation
               intincomefirms / vif;
607        run;
608
609        /* 3.9 add intincomefirms, determine the fourth interaction variable*/
610    proc reg data=cent_train;
611        model y=income white veterans firms inthouseholdeducation intwhiteeducation
               intincomefirms intwhiteincome;
612        run;
613    quit;
614
615    proc reg data=cent_train;
616        model y=income white veterans firms inthouseholdeducation intwhiteeducation
               intincomefirms intwhitehousehold;
617        run;
618
619    quit;
620
621    proc reg data=cent_train;
622        model y=income white veterans firms inthouseholdeducation intwhiteeducation
               intincomefirms intwhitegender;
623        run;
624    quit;
625
626    proc reg data=cent_train;
627        model y=income white veterans firms inthouseholdeducation intwhiteeducation
               intincomefirms intwhitefirms;
628        run;
629    quit;
630
631    proc reg data=cent_train;
632        model y=income white veterans firms inthouseholdeducation intwhiteeducation
               intincomefirms intwhiteveterans;
633        run;
634
635    quit;
636
637    proc reg data=cent_train;
638        model y=income white veterans firms inthouseholdeducation intwhiteeducation
               intincomefirms intincomehousehold;
639        run;
640    quit;
641
642    proc reg data=cent_train;
643        model y=income white veterans firms inthouseholdeducation intwhiteeducation
               intincomefirms intincomegender;
644        run;
645    quit;
646
647    proc reg data=cent_train;
648        model y=income white veterans firms inthouseholdeducation intwhiteeducation
               intincomefirms intincomeeducation;
649        run;
650
651    quit;
652
653    proc reg data=cent_train;
654        model y=income white veterans firms inthouseholdeducation intwhiteeducation
               intincomefirms intincomeveterans;
655        run;
656    quit;
```

```
657
658   proc reg data=cent_train;
659       model y=income white veterans firms inthouseholdeducation intwhiteeducation
660           intincomefirms inthouseholdgender;
661       run;
662   quit;

663   proc reg data=cent_train;
664       model y=income white veterans firms inthouseholdeducation intwhiteeducation
665           intincomefirms inthouseholdfirms;
666       run;
667   quit;
668
669   proc reg data=cent_train;
670       model y=income white veterans firms inthouseholdeducation intwhiteeducation
671           intincomefirms inthouseholdveterans;
672       run;
673   quit;

674   proc reg data=cent_train;
675       model y=income white veterans firms inthouseholdeducation intwhiteeducation
676           intincomefirms intgendereducation;
677       run;
678   quit;

679   proc reg data=cent_train;
680       model y=income white veterans firms inthouseholdeducation intwhiteeducation
681           intincomefirms intgenderfirms;
682       run;
683   quit;
684
685   proc reg data=cent_train;
686       model y=income white veterans firms inthouseholdeducation intwhiteeducation
687           intincomefirms intgenderveterans;
688       run;
689   quit;

690   proc reg data=cent_train;
691       model y=income white veterans firms inthouseholdeducation intwhiteeducation
692           intincomefirms inteducationfirms;
693       run;
694   quit;
695
696   proc reg data=cent_train;
697       model y=income white veterans firms inthouseholdeducation intwhiteeducation
698           intincomefirms inteducationveterans;
699       run;
700   quit;

701   proc reg data=cent_train;
702       model y=income white veterans firms inthouseholdeducation intwhiteeducation
703           intincomefirms intfirmsveterans;
704       run;
705   quit;

706   proc reg data=cent_train;
707       model y=income white veterans Firms inthouseholdeducation intwhiteeducation
708           intincomefirms intincomeeducation;
709       run;
710
711   proc reg data=cent_train;
712       model y=income white veterans firms inthouseholdeducation intwhiteeducation
713           intincomefirms intincomeeducation;
714       run;
715       /* we choose intwhitehousehold as a candidate interaction varible, now check for it */
716       /*r:0.5412 vif<10*/
717   proc reg data=cent_train;
718       model y=income white veterans firms inthouseholdeducation intwhiteeducation
719           intincomefirms intwhitehousehold/ vif;
720       run;
721       /* 3.10 add intwhitehousehold, determine the fifth interaction variable*/
722   proc reg data=cent_train;
723       model y=income white veterans firms inthouseholdeducation intwhiteeducation
724           intincomefirms intwhitehousehold intwhiteincome;
725       run;
726   quit;
727
728   proc reg data=cent_train;
729       model y=income white veterans firms inthouseholdeducation intwhiteeducation
            intincomefirms intwhitehousehold intwhitegender;
```

```
730        run;
731   quit;
732
733   proc reg data=cent_train;
734       model y=income white veterans firms inthouseholdeducation intwhiteeducation
735           intincomefirms intwhitehousehold intwhitefirms;
736       run;
      quit;
737
738   proc reg data=cent_train;
739       model y=income white veterans firms inthouseholdeducation intwhiteeducation
740           intincomefirms intwhitehousehold intwhiteveterans;
741       run;
742   quit;
743
744   proc reg data=cent_train;
745       model y=income white veterans firms inthouseholdeducation intwhiteeducation
746           intincomefirms intwhitehousehold intincomehousehold;
        run;
747   quit;
748
749   proc reg data=cent_train;
750       model y=income white veterans firms inthouseholdeducation intwhiteeducation
751           intincomefirms intwhitehousehold intincomegender;
752       run;
      quit;
753
754   proc reg data=cent_train;
755       model y=income white veterans firms inthouseholdeducation intwhiteeducation
756           intincomefirms intwhitehousehold intincomeeducation;
757       run;
758   quit;
759
760   proc reg data=cent_train;
761       model y=income white veterans firms inthouseholdeducation intwhiteeducation
762           intincomefirms intwhitehousehold intincomeveterans;
763       run;
      quit;
764
765   proc reg data=cent_train;
766       model y=income white veterans firms inthouseholdeducation intwhiteeducation
767           intincomefirms intwhitehousehold inthouseholdgender;
768       run;
769   quit;
770   proc reg data=cent_train;
771       model y=income white veterans firms inthouseholdeducation intwhiteeducation
772           intincomefirms intwhitehousehold inthouseholdfirms;
773       run;
774   quit;
775
776   proc reg data=cent_train;
777       model y=income white veterans firms inthouseholdeducation intwhiteeducation
778           intincomefirms intwhitehousehold inthouseholdveterans;
779       run;
      quit;
780
781   proc reg data=cent_train;
782       model y=income white veterans firms inthouseholdeducation intwhiteeducation
783           intincomefirms intwhitehousehold intgendereducation;
784       run;
      quit;
785
786   proc reg data=cent_train;
787       model y=income white veterans firms inthouseholdeducation intwhiteeducation
788           intincomefirms intwhitehousehold intgenderfirms;
789       run;
790   quit;
791
792   proc reg data=cent_train;
793       model y=income white veterans firms inthouseholdeducation intwhiteeducation
794           intincomefirms intwhitehousehold intgenderveterans;
795       run;
      quit;
796
797   proc reg data=cent_train;
798       model y=income white veterans firms inthouseholdeducation intwhiteeducation
799           intincomefirms intwhitehousehold inteducationfirms;
800       run;
      quit;
801
802   proc reg data=cent_train;
```

```sas
803        model y=income white veterans firms inthouseholdeducation intwhiteeducation
804            intincomefirms intwhitehousehold inteducationveterans;
805        run;
806  quit;
807
808  proc reg data=cent_train;
809        model y=income white veterans firms inthouseholdeducation intwhiteeducation
810            intincomefirms intwhitehousehold intfirmsveterans;
811        run;
     quit;
812
813  proc reg data=cent_train;
814        model y=white Income veterans Firms inthouseholdeducation intwhiteeducation
815            intincomefirms intwhitehousehold intincomeeducation;
816        run;
817
     /* we choose intincomehousehold as a candidate interaction varible, now check for it */
818  proc reg data=cent_train;
819        model y=white income veterans firms inthouseholdeducation intwhiteeducation
820            intincomefirms intwhitehousehold intincomehousehold/ vif;
821        run;
822
     /* 3.11 add intincomehousehold, determine the sisth interaction variable*/
823  proc reg data=cent_train;
824        model y=income white veterans firms inthouseholdeducation intwhiteeducation
825            intincomefirms intwhitehousehold intincomehousehold intwhiteincome;
826        run;
827  quit;
828
829  proc reg data=cent_train;
830        model y=income white veterans firms inthouseholdeducation intwhiteeducation
831            intincomefirms intwhitehousehold intincomehousehold intwhitegender;
832        run;
     quit;
833
834  proc reg data=cent_train;
835        model y=income white veterans firms inthouseholdeducation intwhiteeducation
836            intincomefirms intwhitehousehold intincomehousehold intwhitefirms;
837        run;
838  quit;
839
840  proc reg data=cent_train;
841        model y=income white veterans firms inthouseholdeducation intwhiteeducation
842            intincomefirms intwhitehousehold intincomehousehold intwhiteveterans;
843        run;
     quit;
844
845  proc reg data=cent_train;
846        model y=income white veterans firms inthouseholdeducation intwhiteeducation
847            intincomefirms intwhitehousehold intincomehousehold intincomegender;
848        run;
     quit;
849
850  proc reg data=cent_train;
851        model y=income white veterans firms inthouseholdeducation intwhiteeducation
852            intincomefirms intwhitehousehold intincomehousehold intincomeeducation;
853        run;
854  quit;
855
856  proc reg data=cent_train;
857        model y=income white veterans firms inthouseholdeducation intwhiteeducation
858            intincomefirms intwhitehousehold intincomehousehold intincomeveterans;
859        run;
     quit;
860
861  proc reg data=cent_train;
862        model y=income white veterans firms inthouseholdeducation intwhiteeducation
863            intincomefirms intwhitehousehold intincomehousehold inthouseholdgender;
864        run;
     quit;
865
866  proc reg data=cent_train;
867        model y=income white veterans firms inthouseholdeducation intwhiteeducation
868            intincomefirms intwhitehousehold intincomehousehold inthouseholdfirms;
869        run;
870  quit;
871
872  proc reg data=cent_train;
873        model y=income white veterans firms inthouseholdeducation intwhiteeducation
874            intincomefirms intwhitehousehold intincomehousehold inthouseholdveterans;
     run;
875  quit;
```

```
876
877   proc reg data=cent_train;
878       model y=income white veterans firms inthouseholdeducation intwhiteeducation
879           intincomefirms intwhitehousehold intincomehousehold intgendereducation;
880       run;
881   quit;
882
883   proc reg data=cent_train;
884       model y=income white veterans firms inthouseholdeducation intwhiteeducation
885           intincomefirms intwhitehousehold intincomehousehold intgenderfirms;
886       run;
887   quit;
888
889   proc reg data=cent_train;
890       model y=income white veterans firms inthouseholdeducation intwhiteeducation
891           intincomefirms intwhitehousehold intincomehousehold intgenderveterans;
892       run;
893   quit;
894
895   proc reg data=cent_train;
896       model y=income white veterans firms inthouseholdeducation intwhiteeducation
897           intincomefirms intwhitehousehold intincomehousehold inteducationfirms;
898       run;
899   quit;
900
901   proc reg data=cent_train;
902       model y=income white veterans firms inthouseholdeducation intwhiteeducation
903           intincomefirms intwhitehousehold intincomehousehold inteducationveterans;
904       run;
905   quit;
906
907   proc reg data=cent_train;
908       model y=income white veterans firms inthouseholdeducation intwhiteeducation
909           intincomefirms intwhitehousehold intincomehousehold intfirmsveterans;
910       run;
911   quit;
912
913   /* we choose  intincomegender as a candidate interaction varible, now check for it */
914   /*r:0.5531 vif<10*/
915   proc reg data=cent_train;
916       model y=white income veterans firms inthouseholdeducation intwhiteeducation
917           intincomefirms intwhitehousehold intincomehousehold intincomegender / vif;
918       run;
919
920       /* 3.12 add intincomegender, determine the seventh interaction variable*/
921   proc reg data=cent_train;
922       model y=income white veterans firms inthouseholdeducation intwhiteeducation
923           intincomefirms intwhitehousehold intincomehousehold intincomegender
924           intwhiteincome;
925       run;
926   quit;
927
928   proc reg data=cent_train;
929       model y=income white veterans firms inthouseholdeducation intwhiteeducation
930           intincomefirms intwhitehousehold intincomehousehold intincomegender
931           intwhitegender;
932       run;
933   quit;
934
935   proc reg data=cent_train;
936       model y=income white veterans firms inthouseholdeducation intwhiteeducation
937           intincomefirms intwhitehousehold intincomehousehold intincomegender
938           intwhitefirms;
939       run;
940   quit;
941
942   proc reg data=cent_train;
943       model y=income white veterans firms inthouseholdeducation intwhiteeducation
944           intincomefirms intwhitehousehold intincomehousehold intincomegender
945           intwhiteveterans;
946       run;
947   quit;
948
      proc reg data=cent_train;
          model y=income white veterans firms inthouseholdeducation intwhiteeducation
              intincomefirms intwhitehousehold intincomehousehold intincomegender
              intincomeeducation;
          run;
      quit;

      proc reg data=cent_train;
          model y=income white veterans firms inthouseholdeducation intwhiteeducation
```

```
949          intincomefirms intwhitehousehold intincomehousehold intincomegender
950          intincomeveterans;
951      run;
952  quit;
953
954  proc reg data=cent_train;
955      model y=income white veterans firms inthouseholdeducation intwhiteeducation
956          intincomefirms intwhitehousehold intincomehousehold intincomegender
957          inthouseholdgender;
958      run;
959  quit;
960
961  proc reg data=cent_train;
962      model y=income white veterans firms inthouseholdeducation intwhiteeducation
963          intincomefirms intwhitehousehold intincomehousehold intincomegender
964          inthouseholdfirms;
965      run;
966  quit;
967
968  proc reg data=cent_train;
969      model y=income white veterans firms inthouseholdeducation intwhiteeducation
970          intincomefirms intwhitehousehold intincomehousehold intincomegender
971          inthouseholdveterans;
972      run;
973  quit;
974
975  proc reg data=cent_train;
976      model y=income white veterans firms inthouseholdeducation intwhiteeducation
977          intincomefirms intwhitehousehold intincomehousehold intincomegender
978          intgendereducation;
979      run;
980  quit;
981
982  proc reg data=cent_train;
983      model y=income white veterans firms inthouseholdeducation intwhiteeducation
984          intincomefirms intwhitehousehold intincomehousehold intincomegender
985          intgenderfirms;
986      run;
987  quit;
988
989  proc reg data=cent_train;
990      model y=income white veterans firms inthouseholdeducation intwhiteeducation
991          intincomefirms intwhitehousehold intincomehousehold intincomegender
992          intgenderveterans;
993      run;
994  quit;
995
996  proc reg data=cent_train;
997      model y=income white veterans firms inthouseholdeducation intwhiteeducation
998          intincomefirms intwhitehousehold intincomehousehold intincomegender
999          inteducationfirms;
1000     run;
1001 quit;
1002
1003 proc reg data=cent_train;
1004     model y=income white veterans firms inthouseholdeducation intwhiteeducation
1005         intincomefirms intwhitehousehold intincomehousehold intincomegender
1006         inteducationveterans;
1007     run;
1008 quit;
1009
1010 proc reg data=cent_train;
1011     model y=income white veterans firms inthouseholdeducation intwhiteeducation
1012         intincomefirms intwhitehousehold intincomehousehold intincomegender
1013         intfirmsveterans;
1014     run;
1015 quit;
1016
1017 /* we choose inteducationveterans as a candidate interaction varible, now check for it */
1018 /*r_square:0.5555 vif<10 */
1019 proc reg data=cent_train;
1020     model y=white income veterans firms inthouseholdeducation intwhiteeducation
1021         intincomefirms intwhitehousehold intincomehousehold intincomegender
         inteducationveterans/ vif;
     run;

     /* 3.13 add inteducationveterans, determine the eighth interaction variable*/
proc reg data=cent_train;
     model y=income white veterans firms inthouseholdeducation intwhiteeducation
         intincomefirms intwhitehousehold intincomehousehold intincomegender
         inteducationveterans intwhiteincome;
     run;
```

```
1022  quit;
1023
1024  proc reg data=cent_train;
1025      model y=income white veterans firms inthouseholdeducation intwhiteeducation
1026          intincomefirms intwhitehousehold intincomehousehold intincomegender
1027          inteducationveterans intwhitegender;
1028      run;
1029  quit;
1030  proc reg data=cent_train;
1031      model y=income white veterans firms inthouseholdeducation intwhiteeducation
1032          intincomefirms intwhitehousehold intincomehousehold intincomegender
1033          inteducationveterans intwhitefirms;
1034      run;
1035  quit;
1036
1037  proc reg data=cent_train;
1038      model y=income white veterans firms inthouseholdeducation intwhiteeducation
1039          intincomefirms intwhitehousehold intincomehousehold intincomegender
1040          inteducationveterans intwhiteveterans;
1041      run;
      quit;
1042
1043  proc reg data=cent_train;
1044      model y=income white veterans firms inthouseholdeducation intwhiteeducation
1045          intincomefirms intwhitehousehold intincomehousehold intincomegender
1046          inteducationveterans intincomeeducation;
1047      run;
      quit;
1048
1049  proc reg data=cent_train;
1050      model y=income white veterans firms inthouseholdeducation intwhiteeducation
1051          intincomefirms intwhitehousehold intincomehousehold intincomegender
1052          inteducationveterans intincomeveterans;
1053      run;
      quit;
1054
1055  proc reg data=cent_train;
1056      model y=income white veterans firms inthouseholdeducation intwhiteeducation
1057          intincomefirms intwhitehousehold intincomehousehold intincomegender
1058          inteducationveterans inthouseholdgender;
1059      run;
      quit;
1060
1061  proc reg data=cent_train;
1062      model y=income white veterans firms inthouseholdeducation intwhiteeducation
1063          intincomefirms intwhitehousehold intincomehousehold intincomegender
1064          inteducationveterans inthouseholdfirms;
1065      run;
1066  quit;
1067
1068  proc reg data=cent_train;
1069      model y=income white veterans firms inthouseholdeducation intwhiteeducation
1070          intincomefirms intwhitehousehold intincomehousehold intincomegender
1071          inteducationveterans inthouseholdveterans;
1072  quit;
1073
1074  proc reg data=cent_train;
1075      model y=income white veterans firms inthouseholdeducation intwhiteeducation
1076          intincomefirms intwhitehousehold intincomehousehold intincomegender
1077          inteducationveterans intgendereducation;
1078      run;
      quit;
1079
1080  proc reg data=cent_train;
1081      model y=income white veterans firms inthouseholdeducation intwhiteeducation
1082          intincomefirms intwhitehousehold intincomehousehold intincomegender
1083          inteducationveterans intgenderfirms;
1084      run;
      quit;
1085
1086  proc reg data=cent_train;
1087      model y=income white veterans firms inthouseholdeducation intwhiteeducation
1088          intincomefirms intwhitehousehold intincomehousehold intincomegender
1089          inteducationveterans intgenderveterans;
1090      run;
1091  quit;
1092
1093  proc reg data=cent_train;
1094      model y=income white veterans firms inthouseholdeducation intwhiteeducation
          intincomefirms intwhitehousehold intincomehousehold intincomegender
```

```
1095            inteducationveterans inteducationfirms;
1096        run;
1097  quit;
1098
1099  proc reg data=cent_train;
1100        model y=income white veterans firms inthouseholdeducation intwhiteeducation
1101            intincomefirms intwhitehousehold intincomehousehold intincomegender
1102            inteducationveterans intfirmsveterans;
1103  quit;
1104
1105  /* we choose intwhitefirms as a candidate interaction varible, now check for it */
1106  /*r_square:0.5577 vif<10*/
1107  proc reg data=cent_train;
1108        model y=white income veterans firms inthouseholdeducation intwhiteeducation
1109            intincomefirms intwhitehousehold intincomehousehold intincomegender
1110            inteducationveterans intwhitefirms/ vif;
1111        run;
1112
1113        /* 3.14 add intwhitefirms, determine the ninth interaction variable*/
1114  proc reg data=cent_train;
1115        model y=income white veterans firms inthouseholdeducation intwhiteeducation
1116            intincomefirms intwhitehousehold intincomehousehold intincomegender
1117            inteducationveterans intwhitefirms intwhiteincome;
1118        run;
1119  quit;
1120
1121  proc reg data=cent_train;
1122        model y=income white veterans firms inthouseholdeducation intwhiteeducation
1123            intincomefirms intwhitehousehold intincomehousehold intincomegender
1124            inteducationveterans intwhitefirms intwhitegender;
1125        run;
1126  quit;
1127
1128  proc reg data=cent_train;
1129        model y=income white veterans firms inthouseholdeducation intwhiteeducation
1130            intincomefirms intwhitehousehold intincomehousehold intincomegender
1131            inteducationveterans intwhitefirms intwhiteveterans;
1132        run;
1133  quit;
1134
1135  proc reg data=cent_train;
1136        model y=income white veterans firms inthouseholdeducation intwhiteeducation
1137            intincomefirms intwhitehousehold intincomehousehold intincomegender
1138            inteducationveterans intwhitefirms intincomeeducation;
1139        run;
1140  quit;
1141
1142  proc reg data=cent_train;
1143        model y=income white veterans firms inthouseholdeducation intwhiteeducation
1144            intincomefirms intwhitehousehold intincomehousehold intincomegender
1145            inteducationveterans intwhitefirms intincomeveterans;
1146        run;
1147  quit;
1148
1149  proc reg data=cent_train;
1150        model y=income white veterans firms inthouseholdeducation intwhiteeducation
1151            intincomefirms intwhitehousehold intincomehousehold intincomegender
1152            inteducationveterans intwhitefirms inthouseholdgender;
1153        run;
1154  quit;
1155
1156  proc reg data=cent_train;
1157        model y=income white veterans firms inthouseholdeducation intwhiteeducation
1158            intincomefirms intwhitehousehold intincomehousehold intincomegender
1159            inteducationveterans intwhitefirms inthouseholdfirms;
1160        run;
1161  quit;
1162
1163  proc reg data=cent_train;
1164        model y=income white veterans firms inthouseholdeducation intwhiteeducation
1165            intincomefirms intwhitehousehold intincomehousehold intincomegender
1166            inteducationveterans intwhitefirms inthouseholdveterans;
1167        run;
      quit;
```

```
      proc reg data=cent_train;
          model y=income white veterans firms inthouseholdeducation intwhiteeducation
              intincomefirms intwhitehousehold intincomehousehold intincomegender
              inteducationveterans intwhitefirms intgendereducation;
          run;
      quit;
```

```sas
1168
1169   proc reg data=cent_train;
1170       model y=income white veterans firms inthouseholdeducation intwhiteeducation
1171           intincomefirms intwhitehousehold intincomehousehold intincomegender
1172           inteducationveterans intwhitefirms intgenderfirms;
1173       run;
1174   quit;

1175   proc reg data=cent_train;
1176       model y=income white veterans firms inthouseholdeducation intwhiteeducation
1177           intincomefirms intwhitehousehold intincomehousehold intincomegender
1178           inteducationveterans intwhitefirms intgenderveterans;
1179       run;
1180   quit;

1181   proc reg data=cent_train;
1182       model y=income white veterans firms inthouseholdeducation intwhiteeducation
1183           intincomefirms intwhitehousehold intincomehousehold intincomegender
1184           inteducationveterans intwhitefirms inteducationfirms;
1185       run;
1186   quit;

1187   proc reg data=cent_train;
1188       model y=income white veterans firms inthouseholdeducation intwhiteeducation
1189           intincomefirms intwhitehousehold intincomehousehold intincomegender
1190           inteducationveterans intwhitefirms intfirmsveterans;
1191       run;
1192   quit;

1193
1194   /* we choose intfirmsveterans as a candidate interaction varible, now check for it */
1195   /*r_square:0.5589 vif<10 */
1196   proc reg data=cent_train;
1197       model y=white income veterans firms inthouseholdeducation intwhiteeducation
1198           intincomefirms intwhitehousehold intincomehousehold intincomegender
1199           inteducationveterans intwhitefirms intfirmsveterans/ vif;
1200       run;

1201       /* 3.15 add intfirmsveterans, determine the tenth interaction variable*/
1202   proc reg data=cent_train;
1203       model y=income white veterans firms inthouseholdeducation intwhiteeducation
1204           intincomefirms intwhitehousehold intincomehousehold intincomegender
1205           inteducationveterans intwhitefirms intfirmsveterans intwhiteincome;
1206       run;
1207   quit;

1208   proc reg data=cent_train;
1209       model y=income white veterans firms inthouseholdeducation intwhiteeducation
1210           intincomefirms intwhitehousehold intincomehousehold intincomegender
1211           inteducationveterans intwhitefirms intfirmsveterans intwhitegender;
1212       run;
1213   quit;

1214   proc reg data=cent_train;
1215       model y=income white veterans firms inthouseholdeducation intwhiteeducation
1216           intincomefirms intwhitehousehold intincomehousehold intincomegender
1217           inteducationveterans intwhitefirms intfirmsveterans intwhiteveterans;
1218       run;
1219   quit;

1220   proc reg data=cent_train;
1221       model y=income white veterans firms inthouseholdeducation intwhiteeducation
1222           intincomefirms intwhitehousehold intincomehousehold intincomegender
1223           inteducationveterans intwhitefirms intfirmsveterans intincomeeducation;
1224       run;
1225   quit;

1226   proc reg data=cent_train;
1227       model y=income white veterans firms inthouseholdeducation intwhiteeducation
1228           intincomefirms intwhitehousehold intincomehousehold intincomegender
1229           inteducationveterans intwhitefirms intfirmsveterans intincomeveterans;
1230       run;
1231   quit;

1232
1233   proc reg data=cent_train;
1234       model y=income white veterans firms inthouseholdeducation intwhiteeducation
1235           intincomefirms intwhitehousehold intincomehousehold intincomegender
1236           inteducationveterans intwhitefirms intfirmsveterans inthouseholdgender;
1237       run;
1238   quit;

1239   proc reg data=cent_train;
1240       model y=income white veterans firms inthouseholdeducation intwhiteeducation
```

```
1241            intincomefirms intwhitehousehold intincomehousehold intincomegender
1242            inteducationveterans intwhitefirms intfirmsveterans inthouseholdfirms;
1243        run;
1244    quit;
1245
1246    proc reg data=cent_train;
1247        model y=income white veterans firms inthouseholdeducation intwhiteeducation
1248            intincomefirms intwhitehousehold intincomehousehold intincomegender
1249            inteducationveterans intwhitefirms intfirmsveterans inthouseholdveterans;
1250        run;
1251    quit;
1252
1253    proc reg data=cent_train;
1254        model y=income white veterans firms inthouseholdeducation intwhiteeducation
1255            intincomefirms intwhitehousehold intincomehousehold intincomegender
1256            inteducationveterans intwhitefirms intfirmsveterans intgendereducation;
1257        run;
1258    quit;
1259
1260    proc reg data=cent_train;
1261        model y=income white veterans firms inthouseholdeducation intwhiteeducation
1262            intincomefirms intwhitehousehold intincomehousehold intincomegender
1263            inteducationveterans intwhitefirms intfirmsveterans intgenderfirms;
1264        run;
1265    quit;
1266
1267    proc reg data=cent_train;
1268        model y=income white veterans firms inthouseholdeducation intwhiteeducation
1269            intincomefirms intwhitehousehold intincomehousehold intincomegender
1270            inteducationveterans intwhitefirms intfirmsveterans intgenderveterans;
1271        run;
1272    quit;
1273
1274    proc reg data=cent_train;
1275        model y=income white veterans firms inthouseholdeducation intwhiteeducation
1276            intincomefirms intwhitehousehold intincomehousehold intincomegender
1277            inteducationveterans intwhitefirms intfirmsveterans inteducationfirms;
1278        run;
1279    quit;
1280
1281    /* all the p value > 0.05 and r_square is 0.5589, so the final model is  */
1282    proc reg data=cent_train;
1283        model y=white income veterans firms inthouseholdeducation intwhiteeducation
1284            intincomefirms intwhitehousehold intincomehousehold intincomegender
1285            inteducationveterans intwhitefirms intfirmsveterans/ vif;
1286        run;
1287        *####################################;
1288        *### STEP 4: Final Model verify    ###;
1289        *####################################;
1290        %let reg_x= white income veterans firms inthouseholdeducation intwhiteeducation intincomefirms intwhitehouseho
1291
1292        /*1. Check multicollinearity*/
1293    proc reg data=cent_train;
1294        model y=&reg_x / vif;
1295        run;
1296    quit;
1297
1298    /*2. Check the assumptions */
1299    /* Normality: verify qqplot of (studentised) residuals */
1300    proc reg data=cent_train;
1301        model y=&reg_x;
1302        output out=resid r=rman p=pman student=student;
1303        run;
1304    quit;
1305
1306    /* Linearity: verify plot of (studentised) residuals vs predicted values*/
1307    /* Homoscedasticity: verify (squared) residuals vs predicted values */
1308    data resid2;
1309        set resid;
1310        rman2=rman**2;
1311    run;
1312
1313    proc sgplot data=resid2;
```

Lines continue:

```
1305    proc sgplot data=resid2;
1306        scatter x=pman y=rman2;
1307        refline 0 / axis=y lineattrs=(color=red);
1308    run;
1309
1310    /*3. Check for outliers */
1311    proc reg data=cent_train noprint;
1312        model y=&reg_x / r;
1313        output out=cookdis cookd=cdist;
1313        run;
```

```sas
1314  quit;
1315
1316  data cookdis2;
1317      set cookdis;
1318      n=_n_;
1319  run;
1320
      title "Cook's distance threshold 4/n=0.00255";
1321
1322  proc sgplot data=cookdis2;
1323      scatter x=n y=cdist;
1324      refline 0.0025 / axis=y lineattr=(color=red);
1325  run;
1326
      title;
1327
1328  /*Remove outlier according to Cook's distance threshold 4/n=0.00255;*/
1329  data cent_training_outlier_removed;
1330      set cookdis2;
1331      where cdist<0.00255;
1332
1333      /* 93 outlier removed*/
      run;
1334
1335  /*Build model after removing outlier*/
1336  /*r_square:0.5714,intwhitehousehold(p=0.4253) and intfirmsveterans(p=0.5833)*/
1337  proc reg data=cent_training_outlier_removed;
1338      model y=&reg_x / vif;
1339      run;
1340  quit;
1341
1342  /*Remove intwhitehousehold and intfirmsveterans and check outlier again*/
1343  /*Only 2 outliers and they are very close to the transhold, r_square:0.5714*/
1344
      %let reg_x= white income veterans firms inthouseholdeducation intwhiteeducation intincomefirms intincomehousehold
1345
1346  proc reg data=cent_training_outlier_removed noprint;
1347      model y=&reg_x / r;
1348      output out=cookdis cookd=cdist;
1349      run;
1350  quit;
1351
1352  data cookdis2;
1353      set cookdis;
1354      n=_n_;
      run;
1355
1356  title "Cook's distance threshold 4/n=0.00255";
1357
1358  proc sgplot data=cookdis2;
1359      scatter x=n y=cdist;
1360      refline 0.0025 / axis=y lineattr=(color=red);
      run;
1361
1362  title;
1363
1364  /*Export the dataset for Logistic Regression*/
1365  proc export data=cent_training_outlier_removed (keep=y county &reg_x)
1366          outfile="/folders/myfolders/StatmodProject2020/cent_train.csv" dbms=csv replace;
1367  run;
1368
1369  proc export data=cent_test (keep=y county &reg_x)
1370          outfile="/folders/myfolders/StatmodProject2020/cent_test.csv" dbms=csv replace;
      run;
1371
1372  *####################################;
1373  *### STEP 5: Testing Final Model   ###;
1374  *####################################;
1375
1376  proc reg data=cent_training_outlier_removed outest=train_estimate noprint;
1377      model y=&reg_x;
      run;
1378
1379  proc score data=cent_test score=train_estimate out=test_result type=parms
1380          predict;
1381      var y &reg_x;
1382  run;
1383
1384  ods listing gpath=&dir;
1385  ods graphics on;
      ods select FitPlot;
1386
```

```sas
proc reg data=test_result;
    model y=model1;
    plot model1*y / pred conf;
    run;
    ods graphics off;
    ods listing;
```

```sas
proc reg data=test_result;
    model y=model1;
    plot model1*y / pred conf;
    run;
    ods graphics off;
    ods listing;
```

# US_election_final

## DU Feihong, Tang Lin, Zhang Shengmin

### 2020/12/11

```r
#import data#
mypath = '.'
setwd(mypath)
cent_train<-read.csv('cent_train.csv')
cent_test<-read.csv('cent_test.csv')
whole_dataset<- rbind(cent_train,cent_test)
str(cent_train)
str(cent_test)
plot(cent_train$y,xlab = 'series',ylab = 'vote_percent')

#--------------------------------------------------------------
#6.categorize y#

cent_train$cat_y <-ifelse(cent_train$y>0.5,1,0)
cent_test$cat_y <-ifelse(cent_test$y>0.5,1,0)


#fit with categorized y# state what we have found when compare two models # #
fit.raw.train <-lm(y ~ white+income+veterans+firms+inthouseholdeducation+
                     intwhiteeducation+intincomefirms+intincomehousehold+
                     intincomegender+inteducationveterans+intwhitefirms,
                 data=cent_train)



fit.cat.train <- glm(cat_y ~ white+income+veterans+firms+inthouseholdeducation+
                     intwhiteeducation+intincomefirms+intincomehousehold+
                     intincomegender+inteducationveterans+intwhitefirms,
                 family = binomial,data=cent_train)
fit.raw.train.summary<-summary(fit.raw.train)
fit.cat.train.summary<-summary(fit.cat.train)

beta1 <- fit.cat.train$coefficients[2]
con.odds.white <- exp(beta1)
plot(effects::effect('white', fit.cat.train))


#7.----------------------------------------------------------------------
glm.pred.raw <- predict(fit.raw.train, newdata = cent_test, type = "response")
glm.pred.cat <- ifelse(glm.pred.raw >0.5,1,0)
glm.pred.cat<-as.numeric(glm.pred.cat)
```

```r
#test result#
p_vs_a<-as.matrix(table(actual=cent_test$cat_y,predict=glm.pred.cat))
n = sum(p_vs_a) # number of instances
nc = nrow(p_vs_a) # number of classes
diag = diag(p_vs_a) # number of correctly classified instances per class
rowsums = apply(p_vs_a, 1, sum) # number of instances per class
colsums = apply(p_vs_a, 2, sum) # number of predictions per class
p = rowsums / n # distribution of instances over the actual classes
q = colsums / n # distribution of instances over the predicted classes
accuracy = sum(diag) / n
precision = diag / colsums
recall = diag / rowsums

#overlap------------------------------------------------------------
overlapEst(cent_test$cat_y, glm.pred.cat)
#difference-------------------------------------
t.test(glm.pred.cat,cent_test$cat_y,paired = T)
#variance----------------------------------
var.test(glm.pred.cat, cent_test$cat_y, alternative = "two.sided")
########################################################################
fit.raw.whole <-lm(y ~ white+income+veterans+firms+inthouseholdeducation+
                    intwhiteeducation+intincomefirms+intincomehousehold+
                    intincomegender+inteducationveterans+intwhitefirms,
                 data=whole_dataset)
glm.pred.whole <- predict(fit.raw.whole, newdata = whole_dataset, type = "response")
glm.pred.whole<-as.numeric(glm.pred.whole)
glm.pred.whole.cat <- ifelse(glm.pred.whole >0.52,'win',
                             ifelse(glm.pred.whole<0.48,'lose','undecided'))
whole_dataset$result<-glm.pred.whole.cat


undecided <- whole_dataset%>%
  subset(result == 'undecided')

factor(whole_dataset$county)
factor(undecided$county)
table(undecided$county,undecided$result)
```