

Predicting Unemployment

Dirk Broadhead, Yuxin Chen, Jordan Okada, Audrey Woodwell

17 April 2020

1 Introduction

1.1 Unemployment

Widespread unemployment can severely harm the many aspects of the economy. All of the unemployed individuals lose their source of income and businesses are often unable to produce as many goods or services. Additionally, the purchasing power of the unemployed individuals is decreased, which often negatively impacts other businesses. By knowing the unemployment rate in a given community, civic leaders and policy makers can help determine the overall health and growth of the economy. A low unemployment rate is typically better because it implies that individuals who are searching for a job are more likely to find employment. The purpose of our study is to attempt to understand the relationship between unemployment and several of the other population indicators regularly collected by the United States government. Our study focuses specifically on the data collected from the state of Utah.

1.2 American Community Survey

The United States of America population census takes place every 10 years. However, taking a count of every individual in an entire year is a costly undertaking, so instead a sample survey of the population, known as the American Community Survey (ACS), is conducted every month by the Census Bureau. This survey is sent to a random sample of 3.5 million addresses across the United States. In this way, communities can get current estimates of various statistics for the population every year, rather than relying solely on the information from the last census. The ACS is also used by both local and national leaders to determine the need for further funding, programs, or other projects. For our analysis, we use data from the ACS in 2017 (*The Importance of the American Community Survey and the 2020 Census* (2020)).

1.3 Census Tracts

”Census tracts are relatively permanent small-area geographic divisions of a county or statistically equivalent entity defined for the tabulation and presentation of data” (*Census Tracts for the 2020*

Census-Final Criteria (2018)) We decided to use census tracts rather than cities or counties because a census tract allowed us to closely examine the population of a small geographic area within a county or city. Using data from smaller areas gives us more information to use in our analysis. For example, there are 588 census tract in Utah, which means we have more than 500 observations in our data set. However, there are only 29 counties in Utah, an extremely significant decrease in the number and diversity of the data set.

2 Data

Variable Name	Definition
Income	The average household income
Percent Men	Percent of Men in the total population. The complement is the percent of women.
Poverty Rate	Percent of population in poverty
Child Poverty Rate	Percent of children in poverty
Total Population	Number of people in each census tract
Mean Commute	Average work commute in a census tract
Proportion of Voting Age Citizens	Percent of population above voting age
Income Error	Median household income error
Income Per Capita	The average income per person of a census tract

Figure 1: Variables Used in the Study

The data in the ACS contains information from across the United States. We first separated out the 588 census tracts that are located in Utah. Figure 1 shows the initial nine explanatory variables we began with. We used these variables to explain the unemployment rate by census tract in Utah.

There were three census tracts with no reported values other than the location of the tract (county and state). For the purpose of this analysis we removed these observations because they added no additional information.

Before beginning our analysis, we separated 20% of the observations (approximately 118 data points) from the main data set. 80% of the data was used for the creation of various inference and prediction models. The small portion of the observations that were removed was used later to test the validity of our potential models.

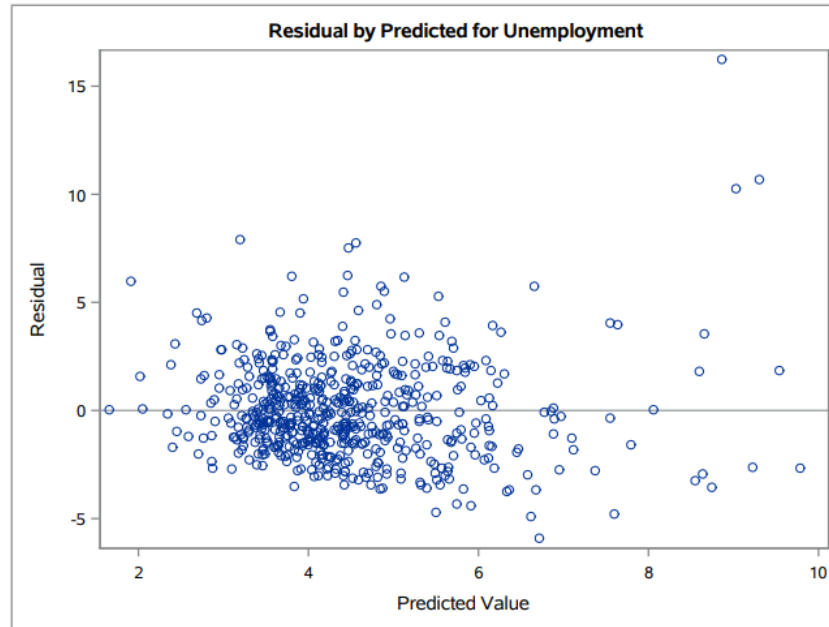


Figure 2: Residual Plot of unemployment

P-value for Brown-Forsythe test of constant variance in Residual vs. Predicted

Obs	t_BF	BF_pvalue
1	4.24265	.000026692

Figure 3: Brow-Forsythe Test of Constant Variance

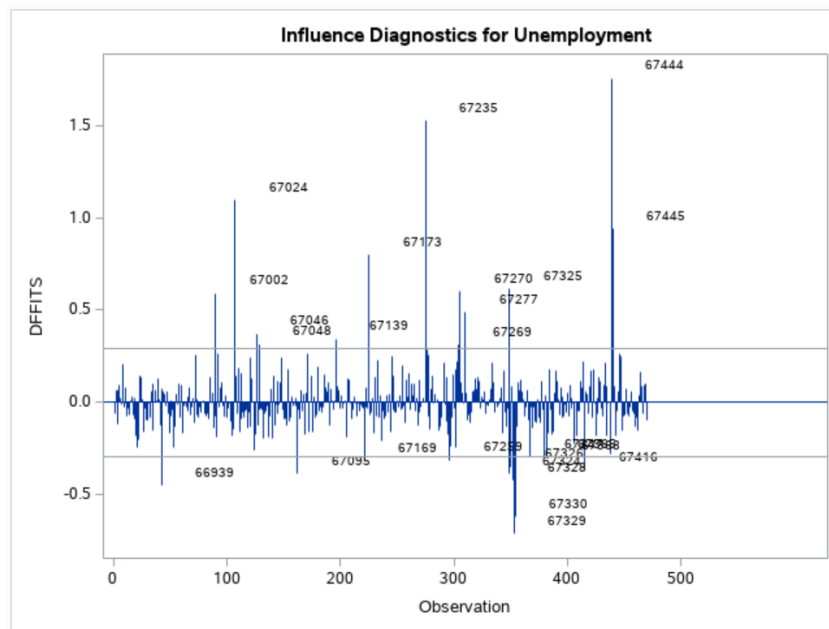
The necessary conditions of a linear model appear to be violated. While the distribution of unemployment is approximately normal, there is evidence of slight non-constant variance. Figure 2 shows the distribution of the residuals against the predicted values, and Figure 3 is the p-value for the Brown Forsythe test of constant variance. Both of these outputs indicate that the residuals for unemployment have a heteroscedastic distribution. Taking a log transformation of unemployment sufficiently reduced the non-constant variance. The significant p-value for the Brown-Forsythe test of constant variance on the log transformed data is given below in Figure 4.

Obs	t_BF	BF_pvalue
1	1.33008	0.18419

Figure 4: Brown-Forsythe Test of Constant Variance

After performing a transformation on unemployment, we tested for multicollinearity in the independent variables. Total population is the sum of the number of men and the number of women

and is an exact linear combination of those two variables. The variables for men, women, and total population had very high variance inflation factors (3195, 2895, and 12036 respectively). The condition index also reflected this multicollinearity. Additionally, the number of voting age citizens had a very high variance inflation factor and was highly collinear with men, women, and total population. Rather than deleting several variables, we transformed the way the variables were presented. We encoded the number of men as a percentage of total population. This automatically includes the data about the total number of women, because the percent of women in the total population is the complement of the percent of men. We also expressed the number of voting age citizens as a percent of the total population. Total population was left alone. This adequately reduced the variance inflation factors to be well below 10. The variance inflation factors for the other variables were also low, suggesting that none of the other variables were highly collinear.



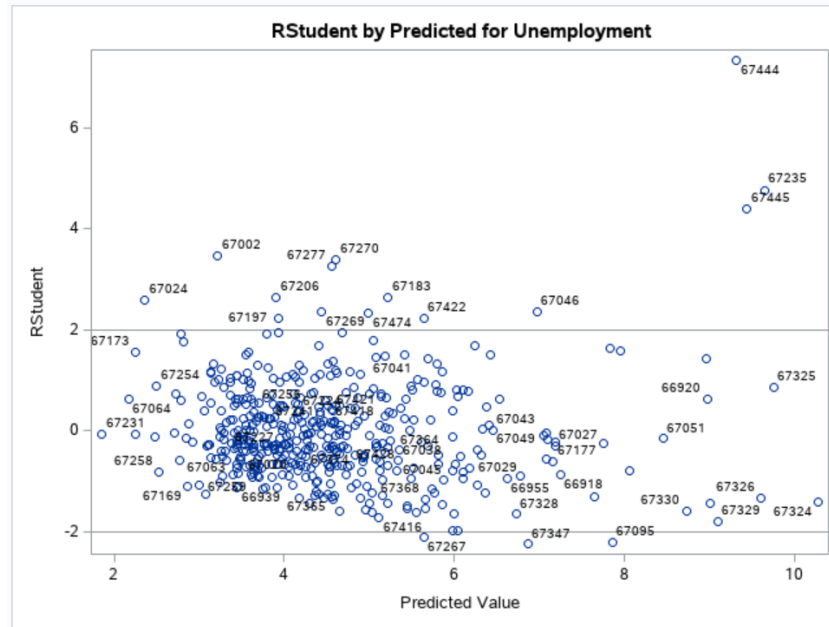


Figure 6: Studentized Residuals for Unemployment

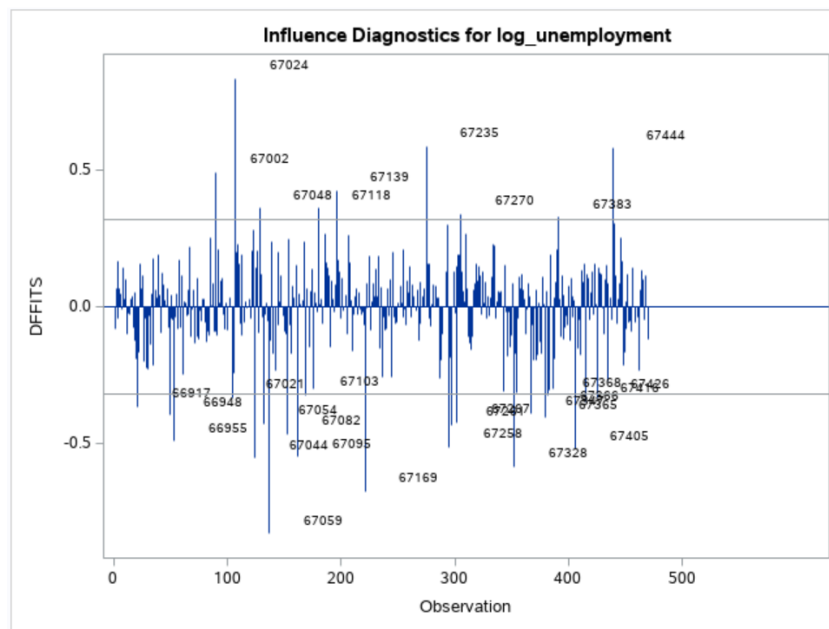


Figure 7: DFFITS for Log Transformed Unemployment

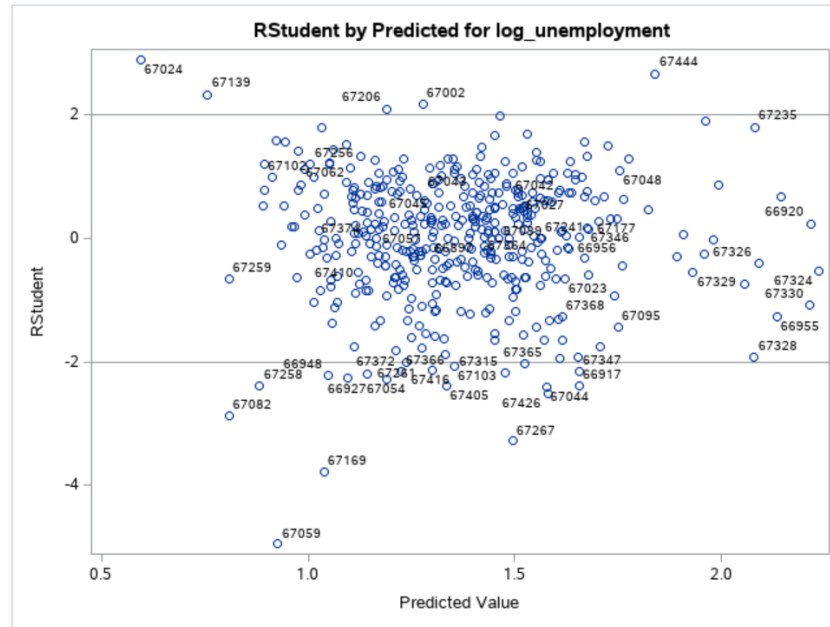


Figure 8: Studentized Residuals for Log Transformed Unemployment

There were several influential points that can be seen in the DFFITS plot in Figure 5. Additionally there were several outliers that can be seen in the plot of studentized residuals for unemployment in Figure 6. These influential points and outliers exist because most of our explanatory variables are right-skewed. Rather than removing so many observations, we transformed most of the explanatory variables. We used a log-transformation on most of the variables; with the exception of percent of men, percent of voting age citizens, and mean commute, which were almost exactly normally distributed to begin with. The distribution of the explanatory variables following the transformation was much closer to normal. And, while the transformation did not fully remove all of the influential points, it did sufficiently mitigate their effects. Figures 7 and 8 show an example of the reduced effects of both the influential points and the outliers following the transformation.

3 Model Testing

To determine the best model for prediction and inference using ordinary least squares regression, we used a variety of variable selection techniques. Both stepwise selection and backward selection suggest using a model with only two variables: poverty and income per capita. The best model based on the CP, adjusted R^2 , AIC, and SBC in all possible regressions also included poverty and income per capita but suggested the addition of income error.

Our remedial measures did transform the data enough to satisfy the assumptions of a linear model. However, we also attempted to evaluate our data using a regression tree to see if a nonparametric method created better predictions. The tree had two branches and suggested income and the child poverty level were the two most important variables when making predictions.

4 Comparing Potential Models

In addition to the basic models suggested above we also tested to see if several interaction terms were significant. First, we tested to see if poverty and income per capita had any interaction. We tested to see if either poverty or income per capita had higher order interactions, specifically if they interact with themselves. The model with potential interaction terms that we tested was as follows:

$$\hat{Y} = \beta_0 + \beta_1(\text{income per capita}) + \beta_2(\text{poverty}) + \beta_{1,2}(\text{income per capita})(\text{poverty}) + \beta_3(\text{income per capita})^2 + \beta_4(\text{poverty})^2$$

The interactions were insignificant: $(\text{income per capita})(\text{poverty})$ had a p-value of $p = 0.8243$, $(\text{income per capita})^2$ had a p-value of $p = 0.7611$, and $(\text{poverty})^2$ had a p-value of $p = 0.1151$. Because these interaction terms did not have a significant effect on unemployment, they will not be included in the model.

In order to determine which of the models performed the best, we calculated the ability of each model to make predictions on the separated 20% of the data. Using the test, set we calculated the mean square prediction rate (MSPR) for four linear regression models: one that contained no predictors, one that contained poverty and income per capita as predictors, a third that contained poverty, income per capita, and income error as predictors, and a fourth that contained all of the initial variables. We also calculated the MSPR for predictions made by the regression tree.

Model	MSPR
Tree	0.3657
2 Variables	0.3301
3 Variables	0.3293
Full	0.3282
None	0.3672

Figure 9: Mean Squared Prediction Rate for Potential Models

As seen in Figure 9, the model with no predictor variables had the highest MSPR. The regression tree predicted only slightly better than the empty model, but there is not a significant difference. The full model and the models with two and three variables had very similar prediction error rates, that were approximately 0.329 (± 0.001). All three of these models decreased the prediction error by approximately 10%.

Because the models including two variables, three variables, and all nine initial variables all performed very similarly, we could justifiably choose any of of them for the final reported model. However, under the assumption that a simpler model is a better model, we will use the two variable model.

5 Final Model and Assumptions

The equation of our final model is given as:

$$\hat{y} = 5.13718 - 0.40302(\text{Income Per Capita}) + 0.14494(\text{Poverty})$$

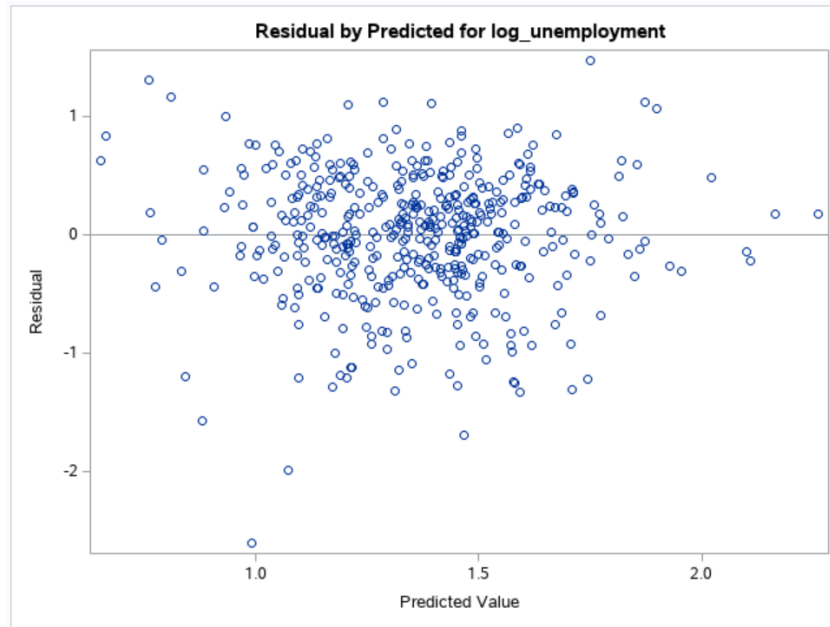


Figure 10: Residuals by Predicted Values

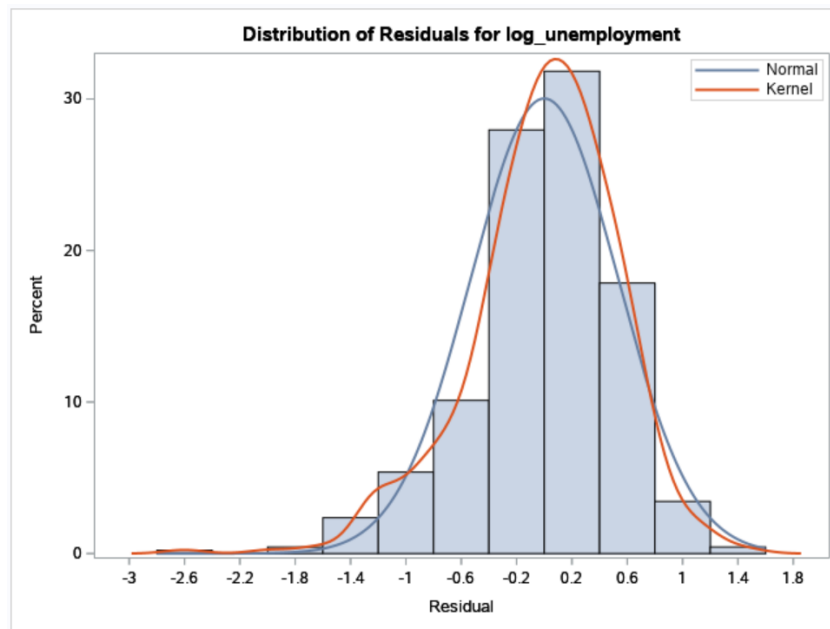


Figure 11: Distribution of Residuals

As seen in Figure 10 the distribution of residuals is homoscedastic, and Figure 11 shows that the distribution of the residuals is approximately normal. We can conclude that the assumptions for a linear model are satisfied.

Income per capita and unemployment have a negative relationship. After accounting for the effect of poverty for every unit increase in income per capita, the unemployment rate decreases by 0.40302. The influence of poverty on unemployment is positive. While holding the effect of income per capita constant for every unit increase in the poverty level, the unemployment rate increases by 0.14494.

6 Conclusion

Our predictions regarding the unemployment rate were not highly accurate. We are, however, able to infer several things about what influences the unemployment rate. We can conclude that the poverty rate is positively correlated with the unemployment rate; both rates increase together. Income per capita on the other hand is negatively related to the unemployment rate, meaning that as one of the two increases, the other decreases. Intuitively, the aforementioned relationships make sense. An increase in the poverty level means people are making less money and are more likely to be unemployed. However, as the average amount of money people in a census tract make increases it is likely that more individuals are actually employed.

We only studied the census tracts that are located in Utah. However, the same information is collected annually across all of the United States. Further study might attempt to generalize our techniques, to predict the unemployment of the United States as a whole, rather than just an individual state.

It would also be beneficial to see if our prediction of unemployment could be improved by adding different, potentially more applicable, predictor variables. Some of these variables might be a measure of the quality of K-12 schools, the average educational attainment in the labor force, or the average age of the labor force. Additionally future studies might use a different, more applicable, study to collect data. The ACS is beneficial because it contains a massive amount of data, it happens yearly, and the information is free to the public. However, it could be beneficial in future studies to collect information with the purpose of studying unemployment.

References

- Census tracts for the 2020 census-final criteria.* (2018, October 30). U.S. Census Bureau. Retrieved 2020-03-30, from <https://www.federalregister.gov/documents/2018/11/13/2018-24567/census-tracts-for-the-2020-census-final-criteria>
- The importance of the american community survey and the 2020 census.* (2020, February). Retrieved 2020-03-28, from <https://www.census.gov/programs-surveys/acs/about/acs-and-census.html>

7 Appendix SAS Code

```
161
162     FILENAME REFFILE '/home/u42026342/STAT5100/Final Paper/utahCensus17.csv';
163
164 PROC IMPORT DATAFILE=REFFILE
165 DBMS=CSV
166 OUT=WORK.census;
167 GETNAMES=YES;
168 RUN;
169
170 /*Initial Check of Data */
171 proc sgscatter data=census;
172 matrix unemployment income men women TotalPop VotingAgeCitizen Poverty ChildPoverty
173 MeanCommute IncomeErr IncomeperCap / markerattrs = (symbol=circlefilled size=2pt);
174 title1 'Scatter Plot Matrix of Variables';
175 run;
176
177 proc univariate data=census;
178 var unemployment income men women TotalPop VotingAgeCitizen Poverty ChildPoverty
179 MeanCommute IncomeErr IncomeperCap;
180 hist unemployment income men women TotalPop VotingAgeCitizen Poverty ChildPoverty
181 MeanCommute IncomeErr IncomeperCap;
182 run;
183
184 /*Fixing Multi Collinearity, Influential Point and Outliers, and Training set */
185 data census1; set census1;
186 percentmen= men/TotalPop;
187 percentVoters=VotingAgeCitizen/TotalPop;
188
189 proc surveyselect data=census1 seed=1337 out=census rate=0.2 outall;
190 /* Withhold 20% for validation */
191 run;
192
193 data train;
194 set census;
195 if Selected=0;
196 run;
197 data test;
198 set census;
199 if Selected=1;
200 run;
201 data train; set train;
202 data train; set train;
203 log_unemployment=log(unemployment);
204 log_men=log(men);
```

```

205 log_women=log(women);
206 log_childpov=log(ChildPoverty);
207 log_population=log(TotalPop);
208 log_poverty=log(Poverty);
209 log_income=log(income);
210 log_VotingAgeCitizen=log(VotingAgeCitizen);
211 log_IncomeErr =log(IncomeErr);
212 log_IncomeperCap=log(IncomeperCap);
213 run;
214
215 proc reg data=train plots(label unpack)=(cooksd Rstudentbyleverage dffits dfbetas);
216 id var1;
217 model unemployment = income men women TotalPop VotingAgeCitizen Poverty ChildPoverty
218 MeanCommute IncomeErr IncomeperCap;
219 output out=out2 r=resid p=pred;
220 title1 'Initial Regression Model Predicting Unemployment';
221 title2 'Also Influential Points and Outliers';
222 run;
223 %resid_num_diag(dataset=out2, datavar=resid, label='Residual', predvar=pred,
224 predlabel='Predicted');
225
226 proc univariate data=train;
227 variable log_unemployment log_income percentment TotalPop percentVoters log_poverty log_
228 MeanCommute log_IncomeErr log_IncomeperCap;
229 histogram log_unemployment log_income log_men log_women log_population log_VotingAgeCiti
230 MeanCommute log_IncomeErr log_IncomeperCap;
231 run;
232
233 proc reg data=train plots(label unpack)=(cooksd Rstudentbyleverage dffits dfbetas);
234 id VAR1;
235 model log_unemployment = income men women TotalPop VotingAgeCitizen Poverty ChildPoverty
236 MeanCommute IncomeErr IncomeperCap; vif collin;
237 output out=out4 r=resid p=pred;
238 title1 'Initial Regression Model Predicting lpg_Unemployment';
239 run;
240 %resid_num_diag(dataset=out4, datavar=resid, label='Residual', predvar=pred,
241 predlabel='Predicted');
242
243 proc reg data=train;
244 id VAR1;
245 model log_unemployment = log_income percentment log_population percentvoters log_poverty
246 MeanCommute log_IncomeErr log_IncomeperCap / vif collin;
247 output out=out3 r=resid p=pred;
248 title1 'Initial Regression Model Predicting Unemployment';
249 run;

```

```

250 %resid_num_diag(dataset=out3, datavar=resid, label='Residual', predvar=pred,
251 predlabel='Predicted');
252
253 /* Variable Selection */
254 proc reg data=train;
255 model log_unemployment = log_income percentment log_population percentvoters log_poverty
256 MeanCommute log_IncomeErr log_IncomeperCap / selection=cp adjrsq aic sbc;
257 title1 'All possible regression - Variable Selection';
258 run;
259
260 proc reg data=train;
261 model log_unemployment = log_income percentment log_population percentvoters log_poverty
262 MeanCommute log_IncomeErr log_IncomeperCap / selection=stepwise slentry=0.1 slstay=0.1;
263 title1 'Stepwise Selection - Variable Selection';
264 run;
265
266 proc reg data=train;
267 model log_unemployment = log_income percentment log_population percentvoters log_poverty
268 MeanCommute log_IncomeErr log_IncomeperCap / selection=backward slstay=0.1;
269 title1 'Backward Selection - Variable Selection';
270
271 proc reg data= train plots(unpack)=diagnostics;
272 model log_unemployment = log_IncomeperCap log_poverty;
273 output out=out1 r=resid p=pred;
274 title1 'Potential Final Model';
275 run;
276 %resid_num_diag(dataset=out1, datavar=resid, label='Residual', predvar=pred,
277 predlabel='Predicted');
278
279 /* TESTING INTERACTIONS */
280 data train; set train;
281 incomexpoverty= log_incomeperCap*log_poverty;
282 incomesquare = (log_incomeperCap)**2;
283 povertysquare= (log_poverty)**2;
284
285 proc reg data=train;
286 model log_unemployment = log_IncomeperCap log_poverty incomexpoverty incomesquare povert
287 run;
288
289 /* Regression Tree */
290 proc hpsplit data=train seed=15531;
291 model log_unemployment = log_income percentment log_population percentvoters log_poverty
292 MeanCommute log_IncomeErr log_IncomeperCap;
293 code file='/home/u42026342/STAT5100/Final Paper/regressiontree.sas/'; /* This saves the
294 run;

```

295

296 /* Validation */

297 data test; set test;

298 log_unemployment=log(unemployment);

299 log_men=log(men);

300 log_women=log(women);

301 log_childpov=log(ChildPoverty);

302 log_population=log(TotalPop);

303 log_poverty=log(Poverty);

304 log_income=log(income);

305 log_VotingAgeCitizen=log(VotingAgeCitizen);

306 log_IncomeErr =log(IncomeErr);

307 log_IncomeperCap=log(IncomeperCap);

308 run;

309

310 proc reg data= train noprint;

311 model log_unemployment = log_IncomeperCap log_poverty;

312 store VarSelectModel;

313 run;

314

315 proc reg data= train noprint;

316 model log_unemployment = log_IncomeperCap log_poverty log_IncomeErr;

317 store VarSelectModel3;

318 run;

319

320 proc reg data=train noprint;

321 model log_unemployment = log_income percentment percentVoters log_population log_poverty

322 MeanCommute log_IncomeErr log_IncomeperCap;

323 store FullModel;

324 run;

325

326 /* Regression Tree */

327 proc hpsplit data=train seed=15531 noprint;

328 model log_unemployment = log_income percentment log_population percentvoters log_poverty

329 MeanCommute log_IncomeErr log_IncomeperCap;

330 code file='/home/u42026342/STAT5100/Final Paper/regressiontree.sas/'; /* This saves the

331 run;

332

333 proc reg data=train noprint;

334 model log_unemployment = ;

335 store emptyModel;

336 run;

337

338 proc plm restore=VarSelectModel;

339 score data=test out=newTest predicted;

```

340 run;
341
342 proc plm restore=VarSelectModel3;
343 score data=test out=newTest3var predicted;
344 run;
345
346 proc plm restore=FullModel;
347 score data=test out=newTest2 predicted;
348 run;
349
350 proc plm restore=emptyModel;
351 score data=test out=newTest3 predicted;
352 run;
353
354 data newTest; set newTest;
355 ASE = (log_unemployment - predicted)**2;
356 run;
357
358 data newTest3var; set newTest3var;
359 ASE = (log_unemployment - predicted)**2;
360 run;
361
362 data newTest2; set newTest2;
363 ASE = (log_unemployment - predicted)**2;
364 run;
365
366 data newTest3; set newTest3;
367 ASE = (log_unemployment - predicted)**2;
368 run;
369
370 data scored;
371 set test;
372 %include '/home/u42026342/STAT5100/Final Paper/regressiontree.sas/';
373 run;
374
375 data testTree;
376 set scored;
377 ASE = (log_unemployment - P_log_unemployment)**2;
378 run;
379 proc means data = testTree;
380 var ASE;
381 title1 'From Regression Tree';
382 run;
383
384 proc means data= newTest;

```

```

385 var ASE;
386 title1 'With 2 Variables';
387 title2 'From Variable Selection';
388 run;
389
390 proc means data= newTest3var;
391 var ASE;
392 title1 'With 3 Variables';
393 title2 'From Variable Selection';
394 run;
395
396 proc means data= newTest2;
397 var ASE;
398 title1 'Model with All variables';
399 run;
400
401 proc means data= newTest3;
402 var ASE;
403 title1 'Model with no Variables';
404 run;

```