

# Gonzalez - PS405 Final Problem Sent

March 17, 2025

Question: Estimate the true conditional expectation function of `repub_change_1216`, given the provided explanatory variables.

**0.0.1 Step 1: Setting up my environment, loading the packages I will use / used, and importing the data. Also examining the data to begin thinking about creating the model.**

```
library(tidyverse)
```

```
-- Attaching core tidyverse packages ----- tidyverse 2.0.0 --
v dplyr      1.1.4      v readr      2.1.5
v forcats    1.0.0      v stringr    1.5.1
v ggplot2    3.5.1      v tibble     3.2.1
v lubridate  1.9.4      v tidyr      1.3.1
v purrr      1.0.2

-- Conflicts ----- tidyverse_conflicts() --
x dplyr::filter() masks stats::filter()
x dplyr::lag()     masks stats::lag()
i Use the conflicted package (<http://conflicted.r-lib.org/>) to force all conflicts to
```

```
library(sandwich)
```

```
library(modelsummary)
```

``modelsummary`` 2.0.0 now uses ``tinytable`` as its default table-drawing backend. Learn more at: <https://vincentarelbundock.github.io/tinytable/>

Revert to ``kableExtra`` for one session:

```
options(modelsummary_factory_default = 'kableExtra')
options(modelsummary_factory_latex = 'kableExtra')
```

```
options(modelsummary_factory_html = 'kableExtra')
```

Silence this message forever:

```
config_modelsummary(startup_message = FALSE)
```

```
library(ggeffects)
```

```
library(broom)
```

```
library(car)
```

Loading required package: carData

Attaching package: 'car'

The following object is masked from 'package:dplyr':

recode

The following object is masked from 'package:purrr':

some

```
library(kableExtra)
```

Attaching package: 'kableExtra'

The following object is masked from 'package:dplyr':

group\_rows

```
library(estimatr)
library(lmtest)
```

Loading required package: zoo

Attaching package: 'zoo'

The following objects are masked from 'package:base':

```
as.Date, as.Date.numeric
```

```
load("~/Documents/PS405-Linear-Models/final_problem_set/pres_elec16_data.RData")

#rename the datsa

election <- exam_data

# view(election)
# str(election)
```

### 0.0.2 Step 2: Clustering the data by state.

Generally, with all the county level data exists, it is probably a good idea to cluster the data by each state. Clustering the county data by each state also helps account of spatial similarities at the intra-state level, and potentially even the national level. In general, clustering will improve our accuracy, but it will also increase our standard errors why is a known byproduct of clustering.

```
#named the new column state_v2 as it is the second iteration of state data in the data  
  
election$state_v2 <- as.factor(factor(election$state))
```

### 0.0.3 Step 3: Creation of new variables.

There is a lot of demographic information within the data set, and after analyzing it for a bit, I think it would be wise to create some new variables that help us analyze our interest which is the change in republican vote share from the 2016 to 2012 elections.

The new variables I am going to create are as follows: - Change in unemployment - Change in median income - Change in poverty rate - total non-white VAP

The three new variables pertaining to change, are to tease out if economic standing had anything to do with the change in who people might vote for. The total non-white VAP is to make it easier to examine if having a larger amount of non-white VAP's is associated with a change. This new variable does raise some concern however as Hispanics have increasingly voted for republicans at a much higher rate than blacks, who they will be grouped with in this new variable.

For all these new variables, I am doing  $2016value - 2012value$ , to get the difference. For this, a positive result would mean an increase in for example median income from 2012 to 2016, and a negative result would mean a decrease.

The result in the dataframe below show a slight increase in median income, which is \$2,270, and decreases in the poverty rate, and unemployment rate.

```
# doing the math  
  
election$unempl_diff <- election$unempl_rate15 - election$unempl_rate12  
election$median_diff <- election$median_income16 - election$median_income12
```

```

election$Pov_diff <- election$poverty16 - election$poverty12

# view(election)

#storing the values to examine

mean_un <- mean(election$unempl_diff)
mean_median <- mean(election$median_diff)
mean_pov <- mean(election$Pov_diff)

#making a df to easily examine

means <- data.frame(mean_un, mean_median, mean_pov)

#print(means)

## now going to create the nonwhite variable

election$total_nonwhite_vap <- election$black_vap + election$asian_vap + election$asian_

```

#### 0.0.4 Step 5: Making models

Now that I have some new variables created, I am going to start making some models. As someone who's research interests generally fall in the lines of urban politics, race, class and political economy, I am going to focus on indicators related to those areas first. I will bring more stuff in and or make changes if those do not net any results.

For my first models, I am going to just put all of our new variables in, and see what we get.

```
#basic lm/ols model to start, without clusters.

basic_model <- lm(repub_change1216 ~ unempl_diff + median_diff + Pov_diff + total_nonwhi

# next making a model with the clusters

model_clustered <- lm_robust(repub_change1216 ~ unempl_diff + median_diff + Pov_diff + t

# Now let's compare

#modelsummary(list("Basic" = basic_model, "Clustered" = model_clustered),
#gof_omit = "AIC|BIC|Log.Lik",
#title = "Comparing Models",
#stars = TRUE,
#digits = 5)
```

### 0.0.5 Step 6: Improving the models

The models created do not explain that much. So the the next step may be to remove or add or interact variables.

```
# going to interact the variables first to see that gives us anything. Just redoing the

basic_model_interactions <- lm(repub_change1216 ~ unempl_diff * median_diff * Pov_diff *

model_clustered_interactions <- lm_robust(repub_change1216 ~ unempl_diff * median_diff *

data = election,
```

```

clusters = election$state_v2)

both <- list("Basic" = basic_model_interactions, "Clustered" = model_clustered_interactions)

#modelsummary(basic_model_interactions, stars = TRUE)

#summary(model_clustered_interactions)

# again the models pretty much explain nothing.

```

Interacting did not do anything.

Now let's bring urban population percentage into the mix.

```

basic_model_urban <- lm(repub_change1216 ~ unempl_diff + pct_urban + median_diff + Pov_diff)

# next making a model with the clusters

model_clustered_urban <- lm_robust(repub_change1216 ~ unempl_diff + pct_urban + median_diff + Pov_diff,
                                   clusters = election$state_v2)

# Now let's compare

#modelsummary(list("Basic" = basic_model_urban, "Clustered" = model_clustered_urban),
#              #gof_omit = "AIC|BIC|Log.Lik",
#              #title = "Comparing Models",
#              #stars = TRUE,
#              #digits = 5)

```

This begins to get us somewhere. The added coefficient of percentage of population that is



in an urban center, highlights that every 1 percent increase in the urban population within a county changes the amount of votes casted for the Republican presidential candidate by a decrease of 0.077. So Donald Trump made successful inroads with rural voters.

I also want to explore college educated voters, and their relationship to the change in republican vote share.

```
basic_model_college <- lm(repub_change1216 ~ unempl_diff + pct_urban + prop_noncollege +  
  
# next making a model with the clusters  
  
model_clustered_college <- lm_robust(repub_change1216 ~ unempl_diff + pct_urban + prop_n  
  
# Now let's compare  
  
#modelsummary(list("Basic" = basic_model_college, "Clustered" = model_clustered_college)  
#gof_omit = "AIC|BIC|Log.Lik",  
#title = "Comparing Models",  
#stars = TRUE,  
#digits = 5)
```

The results were pretty significant. The increase of non-college graduates in a county, and in a state, drastically increase the change in republican vote shares. The next demographic area I am going to check, and explore is occupation.

```
basic_model_occupation <- lm(repub_change1216 ~ unempl_diff + pct_urban + prop_manufactu  
  
# next making a model with the clusters  
  
model_clustered_occupation <- lm_robust(repub_change1216 ~ unempl_diff + prop_manufactur
```

```
# Now let's compare

#modelsummary(list("Basic" = basic_model_occupation, "Clustered" = model_clustered_occup
#gof_omit = "AIC|BIC|Log.Lik",
#title = "Comparing Models",
#stars = TRUE,
#digits = 5)
```

Again, we get some change. Trump made big gains with manufacturing industry, but also made big losses in agriculture and construction.

Let's re-run the latest models but with just blacks, instead of all nonwhites.

```
basic_model_occupation <- lm(repub_change1216 ~ unempl_diff + pct_urban + prop_manufactu

# next making a model with the clusters

model_clustered_occupation <- lm_robust(repub_change1216 ~ unempl_diff + prop_manufactu

# Now let's compare

#modelsummary(list("Basic" = basic_model_occupation, "Clustered" = model_clustered_occup
#gof_omit = "AIC|BIC|Log.Lik",
#title = "Comparing Models",
#stars = TRUE,
#digits = 5)
```

Nothing really changes.

### 0.0.6 Step 7: Piecing some stuff together.

So we know from our myriad of models so far that our new variables of difference in median income generally has not been helpful in any model, the same with `total_nonwhite_vap`. The difference in poverty rate has varied, but has somewhat explanatory in our models. However, throughout all of our models, the most salient variables have been the ones pertaining to education, industry, unemployment, and urban population.

Therefore, let's make a *focused* model with those variables.

```
improved_model <- lm(repub_change1216 ~ unempl_diff + prop_manufacturing + prop_agr + pr

# next making a model with the clusters

improved_model_clustered <- lm_robust(repub_change1216 ~ unempl_diff + prop_manufacturin

# Now let's compare

#modelsummary(list("Basic" = improved_model, "Clustered" = improved_model_clustered),
#gof_omit = "AIC|BIC|Log.Lik",
#title = "Comparing Models",
#stars = TRUE,
#digits = 5)
```

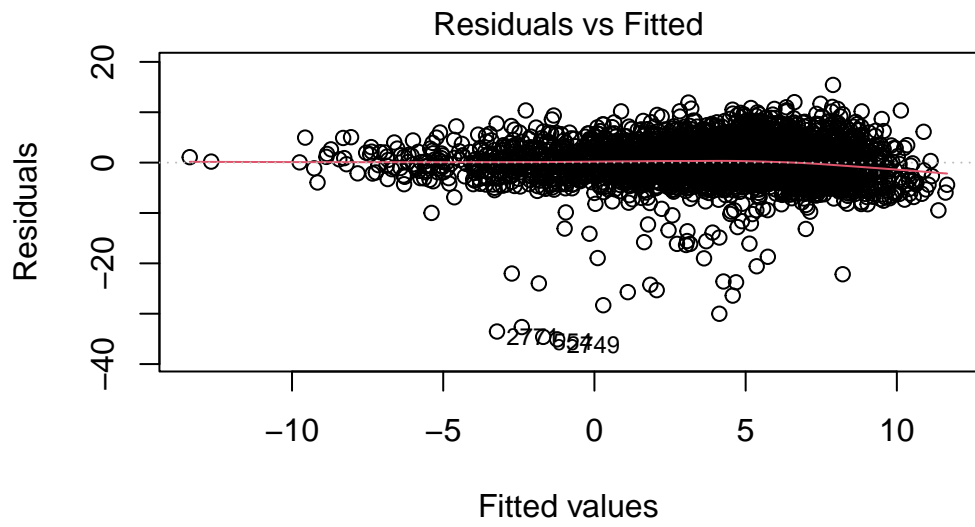
### 0.0.7 Step 8: Analyzing our model

For now, these models are what we want to go forward with.

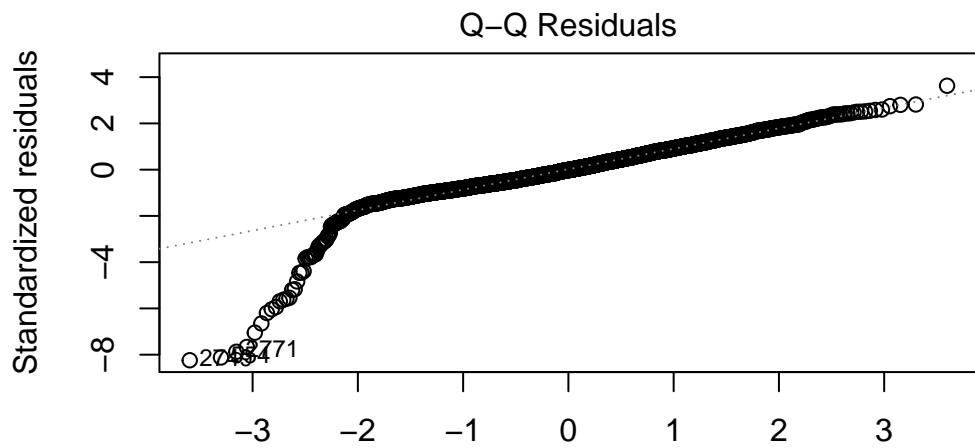
Let's analyze them to make sure they uphold with our OLS assumptions.

Checking distribution

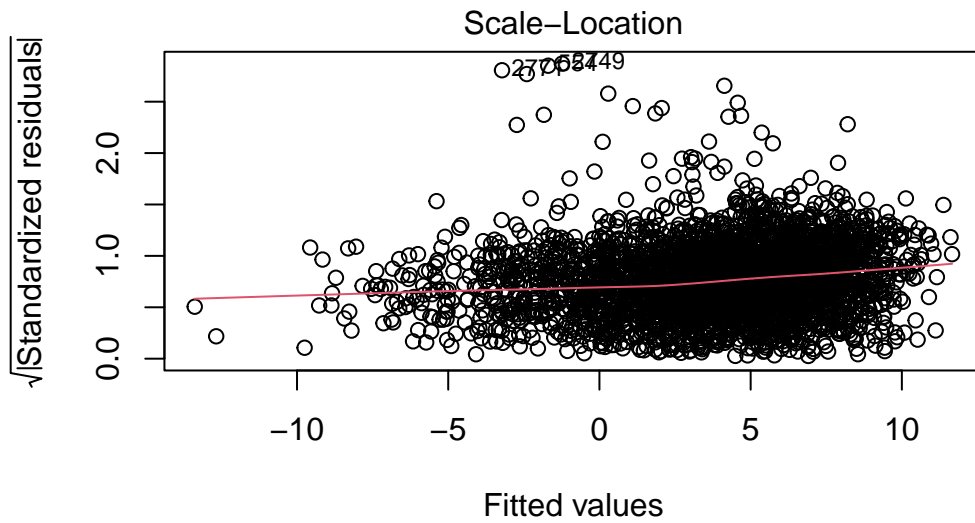
```
plot(improved_model)
```



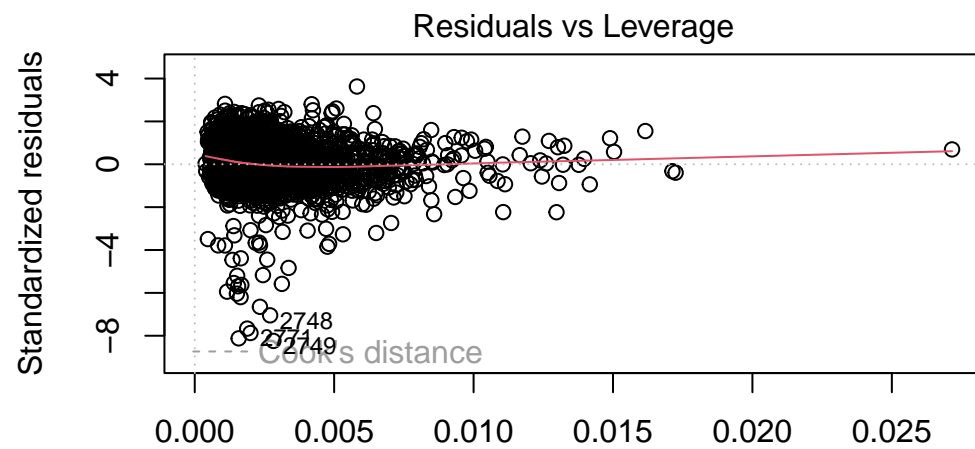
(repub\_change1216 ~ unempl\_diff + prop\_manufacturing + prop\_agr + pro



(repub\_change1216 ~ unempl\_diff + prop\_manufacturing + prop\_agr + pro

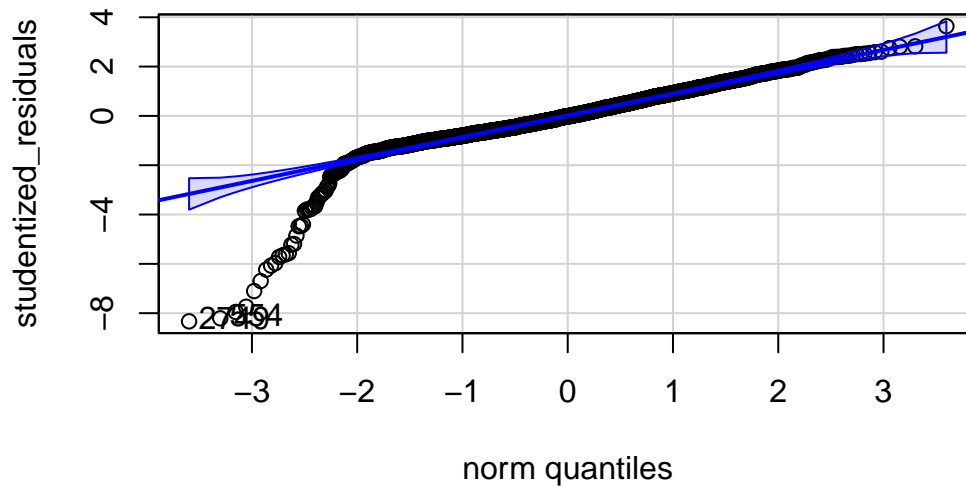


(repub\_change1216 ~ unempl\_diff + prop\_manufacturing + prop\_agr + pro



(repub\_change1216 ~ unempl\_diff + prop\_manufacturing + prop\_agr + pro

```
studentized_residuals <- rstudent(improved_model)
qqPlot(studentized_residuals)
```



```
[1] 2749 554
```

QQ plot results for the regular model is normally distributed to a point, but the lower tail has some serious outliers and is pretty extreme. This impacts some of the validity of our model. Our model is skewed negatively to the left, in statistical terms.

Let's check for homoskedasticity.

```
bptest(improved_model)
```

studentized Breusch-Pagan test

```
data: improved_model
```

```
BP = 11.868, df = 7, p-value = 0.105
```

```
bptest(improved_model_clustered)
```

studentized Breusch-Pagan test

```
data: improved_model_clustered
```

```
BP = 11.868, df = 7, p-value = 0.105
```

The models prove to be homoskedasticity, so because of this, and the size of our data, which is over 3,000 observations, the skewed tail of our QQ plot is less of a concern.

### 0.0.8 Step 8 Addressing QQ plot, non-normal residuals concern.

Even with the concern not being as high, I want to see if there is something we can do to get a normal distribution.

First, let's log the percent of people living in an urban center.

```
election$log_pct_urban <- log1p(election$pct_urban / 100) # transforming the variable,
```

Nothing really change. While my model did pass the Breusch-Pagan test, I will still use HC3 for one of the models, just to be *safe* regarding my results, since my qqplot show some concerns.

```
coeftest(improved_model, vcov = vcovHC(improved_model, type = "HC3"))
```

t test of coefficients:

	Estimate	Std. Error	t value	Pr(> t )
(Intercept)	-14.1690890	0.8357754	-16.9532	< 2.2e-16 ***
unempl_diff	0.3034771	0.0600382	5.0547	4.560e-07 ***
prop_manufacturing	9.2656043	1.4408837	6.4305	1.468e-10 ***
prop_agr	-12.2491194	1.4148551	-8.6575	< 2.2e-16 ***
prop_constr	-15.3131783	3.8001294	-4.0296	5.720e-05 ***
pct_urban	-0.0559566	0.0031915	-17.5328	< 2.2e-16 ***
prop_noncollege	27.4301236	0.9891446	27.7312	< 2.2e-16 ***
Pov_diff	0.1637868	0.0411630	3.9790	7.079e-05 ***
---				

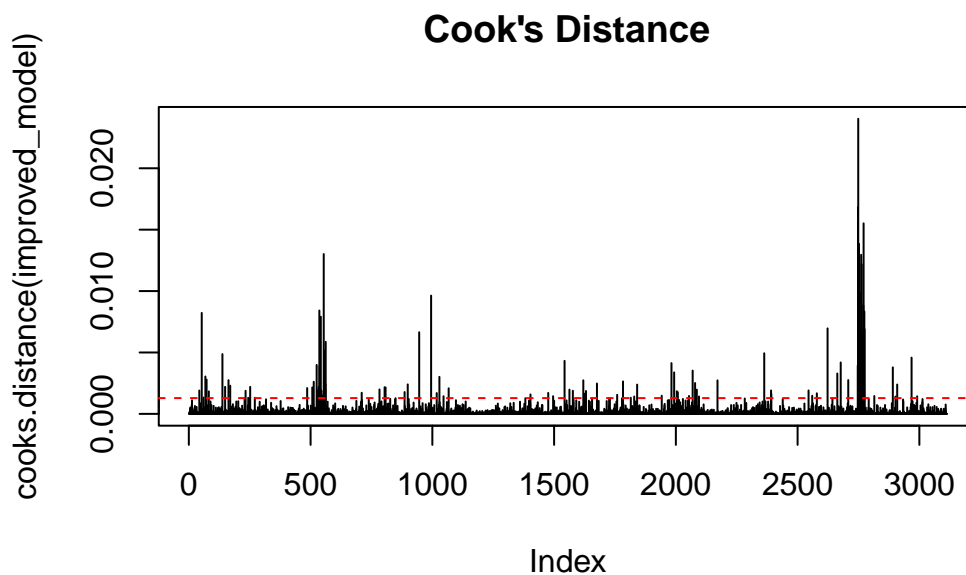
Signif. codes: 0 '\*\*\*' 0.001 '\*\*' 0.01 '\*' 0.05 '.' 0.1 ' ' 1

### 0.0.9 Step 9: Influential Points

The last thing I want to check to maybe correct the model, is to examine influential points.

```
# checking influential points
```

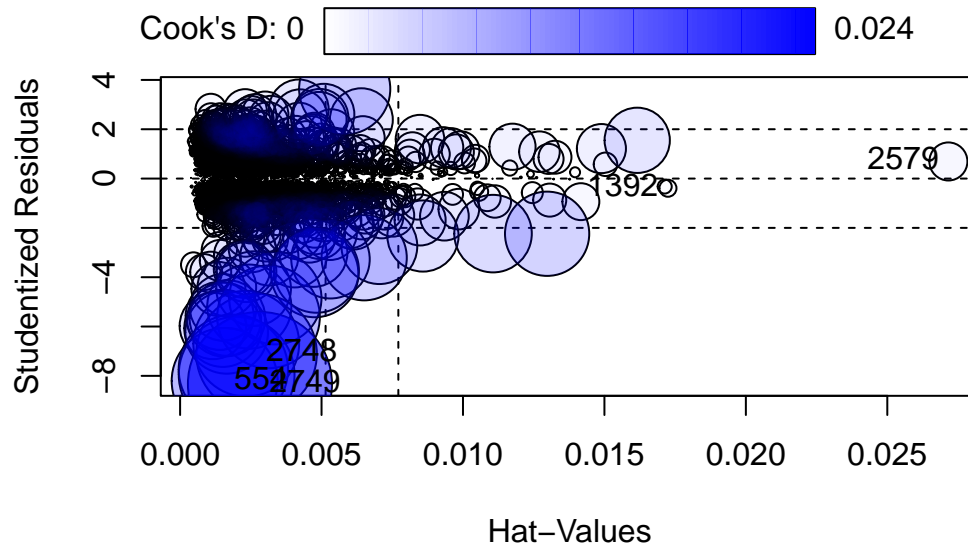
```
plot(cooks.distance(improved_model), type = "h", main = "Cook's Distance")
abline(h = 4/nrow(election), col = "red", lty = 2)
```



```
#influence plot to look at cooks distance
```

```
influencePlot(improved_model)
```





	StudRes	Hat	CookD
554	-8.2083345	0.001576079	0.0130164048
1392	-0.3920671	0.017245724	0.0003372762
2579	0.6930919	0.027159511	0.0016766591
2748	-7.1038084	0.002703059	0.0168288772
2749	-8.3316986	0.002822431	0.0240301445

# I asked ChatGPT how to get what points were influential in a specific dataframe, and g

```
influential_points <- which(cooks.distance(improved_model) > (4/nrow(election)))

#election[influential_points, ] %>%
#summarise(
#  #mean_median_income16 = mean(median_income16, na.rm = TRUE),
#  # mean_median_income12 = mean(median_income12, na.rm = TRUE),
#  # mean_poverty_2012 = mean(poverty12, na.rm = TRUE),
#  # mean_poverty_2016 = mean(poverty16, na.rm = TRUE)
#)
```

Okay, so this is sort of telling. The mean median income changes, but the poverty does not. What I am observing for this is that it may have been wrong to create variables that looked at the difference in poverty, and income or even unemployment. The reason for this is places that already had high poverty, that remained with high poverty, may have changed their vote. They could be unhappy with that fact the party they voted for last time, did not improve their material circumstance. This is not something that the rate of change would showcase.

#### 0.0.10 Step 10: Recreating models, with not new variables created.

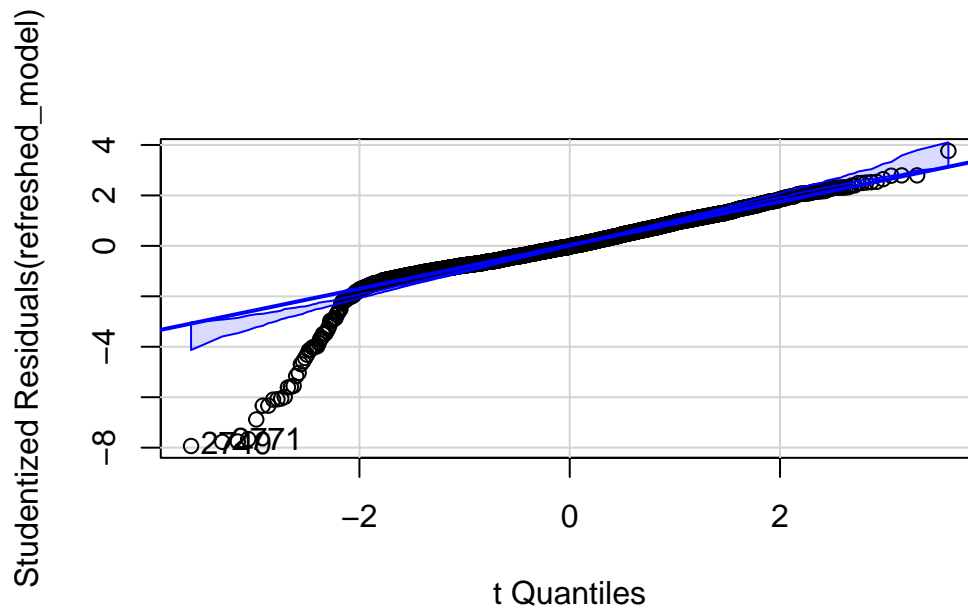
I am just going to now create two new models, with all the variables are of interest, but not combine any or make any new ones.

```
refreshed_model <- lm(repub_change1216 ~ poverty12 + poverty16 + unempl_rate12 + unempl_
refreshed_model_clustered <- lm_robust(repub_change1216 ~ poverty12 + poverty16 + unempl

#modelsummary(list("Refreshed Basic" = refreshed_model, "Refreshed Clustered" = #refresh
#gof_omit = "AIC|BIC|Log.Lik",
#title = "Comparing Models",
#stars = TRUE,
#digits = 5)
```

I will come back to the results in a bit but let's check if things are more "normal".

```
qqPlot(refreshed_model)
```



```
[1] 2749 2771
```

```
bptest(refreshed_model)
```

```
studentized Breusch-Pagan test
```

```
data: refreshed_model
```

```
BP = 23.147, df = 9, p-value = 0.005875
```

The qqPlot hardly changes, but now the model fails the homoskedastic check.

Using each unique year would violate the multi-collinearity OLS assumption, but the results do not really improve anything either. So this test was somewhat irrelevant, but it was worth checking.

### 0.0.11 Step 10: Models with seperate years

The last way we can examine the relationship(s) with the change or lack of in poverty rates, median income, and unemployment, is by examining them all individually.

```

# model for 2012

refreshed_model_2012 <- lm(repub_change1216 ~ poverty12 + unempl_rate12 + median_income12)

refreshed_model_clustered_2012 <- lm_robust(repub_change1216 ~ poverty12 + unempl_rate12 + median_income12)

#models for 2016

refreshed_model_2016 <- lm(repub_change1216 ~ poverty16 + unempl_rate15 + median_income16)

refreshed_model_clustered_2016 <- lm_robust(repub_change1216 ~ poverty16 + unempl_rate15 + median_income16)

#modelsummary(list("Refreshed Basic 2012 " = refreshed_model, "Refreshed Clustered 2016 " = refreshed_model_clustered_2016),
#gof_omit = "AIC|BIC|Log.Lik",
#title = "Comparing Models",
#stars = TRUE,
#digits = 5)

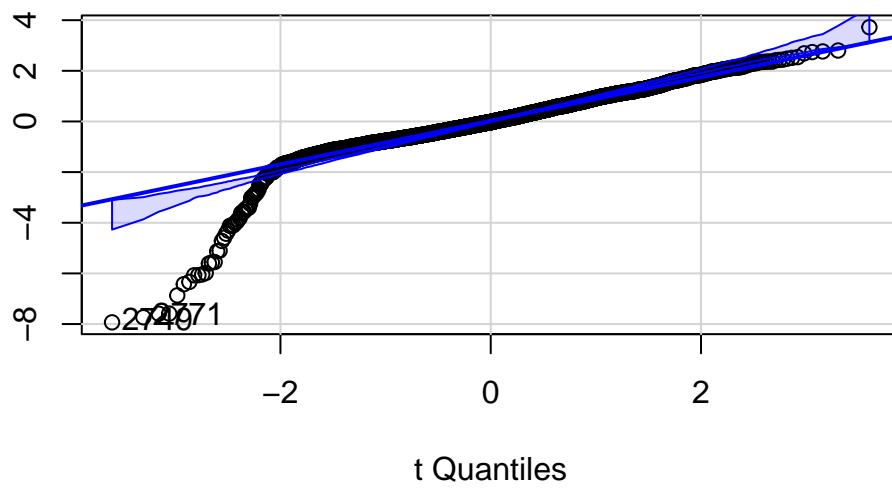
# going to leave results for now

```

Now testing assumptions

```
qqPlot(refreshed_model_2012)
```

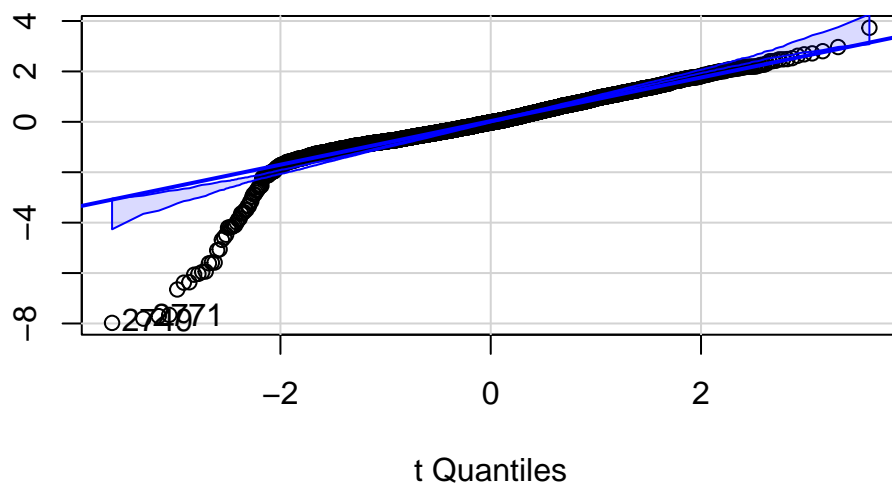
Studentized Residuals(refreshed\_model\_2012



[1] 2749 2771

```
qqPlot(refreshed_model_2016)
```

Studentized Residuals(refreshed\_model\_2016



[1] 2749 2771

```
bptest(refreshed_model_2012)
```

studentized Breusch-Pagan test

```
data: refreshed_model_2012
```

```
BP = 22.687, df = 6, p-value = 0.0009085
```

```
bptest(refreshed_model_2016)
```

studentized Breusch-Pagan test

```
data: refreshed_model_2016
```

```
BP = 22.715, df = 6, p-value = 0.0008976
```

Again, not much has changed.

### **0.0.12 Step 11: Rounding up what has been learned and observed, and some new variables.**

This section is dedicated to working through putting together all the things I would like to present, and showcase in the “final write up”.

The first thing I would like to go over, and discuss is the utility of my created variables. Those being “pov\_diff”, which speaks to the difference in poverty rate between 2016 and 2012, and “median\_diff” as well as “unempl\_diff” which also look at the difference in that same temporal period. While these are useful as potential descriptors, they cause multicollinearity issues. The reason this OLS assumption is violated is because the two poverty rates are highly correlated with each other, and the same is true for median income and unemployment rate.

```
cor(election$poverty12, election$poverty16)
```

```
[1] 0.9607138
```

```
cor(election$median_income16, election$median_income12)
```

```
[1] 0.9787943
```

```
cor(election$unempl_rate15, election$unempl_rate12)
```

```
[1] 0.8840874
```

It is a point of interest to potentially measure how change might have affected the republican vote share, but that is probably for a separate model, or could be fixed with further attempts to log variables, or even use a polynomial model. Generally however, the models with my new created variables is extreme on the left side, and is negatively skewed because its skew is to the left. The last variable I want to explore is a state-based variable, to tease out some of the results from one of the previous models. When we looked at industry, proportion of manufacturing workers was salient. Because of this, and knowing that Trump made in-roads with non-urban voters from other models as well, I am going to create a rust belt variable.

```
# defining these states the rust belt
```

```
rust_belt_states <- c("Ohio", "Michigan", "Pennsylvania", "Kentucky", "West Virginia", "
```

```
# New York is a questionable rust belt state, but I will leave it in for now.
```

```
election$rust_belt <- ifelse(election$state %in% rust_belt_states, 1, 0)
```

```
rust_belt_model <- lm(repub_change1216 ~ poverty16 + unempl_rate15 + median_income16 + p
```

```
modelsummary(rust_belt_model, stars = TRUE)
```

So the results of the rust belt variable are statistically significant. Let's make a model with the 2012 economic indicators, manufacturing, and rust belt.

```
holisitc_model <- lm(repub_change1216 ~ rust_belt + prop_manufacturing + pct_urban + pov  
modelsummary(holisitc_model, stars = TRUE)
```

So this model looks good, but unemployment rate does not seem to explain anything, the same with median\_income. So let's remove them

```
revised_model <- lm(repub_change1216 ~ rust_belt + prop_manufacturing + pct_urban + pove  
modelsummary(revised_model, stars = TRUE)
```

In this most recent model, all of our results are significant at the .0001 level. Our  $R^2$  also does not change despite removing some variables so this show they did not really explain much, or add much to our model.

The last thing I want to explore is rural-ness a bit deeper. So I am going to make a variable for count with a sub 50% urban percentage.

```
# variable creation  
election$rural <- ifelse(election$pct_urban < 50, 1, 0)
```

```
#creating new model  
#rural_model <- lm(repub_change1216 ~ rust_belt + prop_manufacturing + pct_urban + pover  
#modelsummary(rural_model, stars = TRUE)
```

,



### 0.0.13 Step 12: Putting together final materials

Based on all we have uncovered, we want to check some assumptions of our final models

Our final models are going to be as follows, and why:

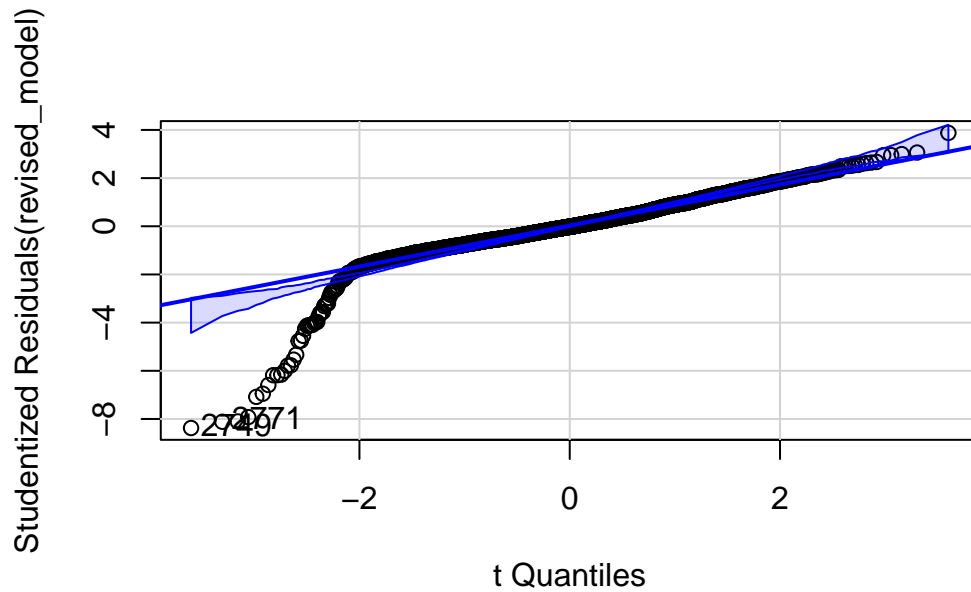
Model 1: repub\_change126, prop\_noncollege, pct\_urban, prop\_manufacturing, rust\_belt, poverty16,

```
# at this point I explored using poverty12 because that would make more sense theoreti  
  
final_rust_belt <- lm(repub_change1216 ~ poverty16 + prop_noncollege + prop_manufacturin  
  
modelsummary(final_rust_belt)
```

Why?: This model contains the most salient results, with the most explanatory power. As all these results were statistically significant.

Our model is heteroskedastic however, as scene below, so we will use HC3 standard errors. Also known as robust standard errors.

```
qqPlot(revised_model)
```



```
[1] 2749 2771
```

```
bptest(revised_model)
```

studentized Breusch-Pagan test

data: revised\_model

BP = 37.19, df = 5, p-value = 5.488e-07

```
# using hc3
```

```
final_rust_belt_hc3 <- lm_robust(repub_change1216 ~ poverty16 + rust_belt + pct_urban +
                                data = election, se_type = "HC3")
```

```
modelsummary(final_rust_belt_hc3, stars = TRUE)
```

While this model may violate the multicollinearity assumption with how its variables are

created, our original “imrpoved” model passed the BP test. So it is also of use. This can be seen below, and the model.

```
improved_model <- lm(repub_change1216 ~ unempl_diff + prop_manufacturing + prop_agr + pr  
bptest(improved_model)
```

studentized Breusch-Pagan test

data: improved\_model

BP = 11.868, df = 7, p-value = 0.105

Then lastly, we have our clustered model, of what is our best model, without the rust-belt variable as this would create issues with our state based cluster.

```
final_cluster <- lm_robust(repub_change1216 ~ unempl_diff + prop_manufacturing + prop_ag  
modelsummary(final_cluster, stars = TRUE)
```

```
# for the final presenation we should prob present this model unclustered just for refer
```

```
final_uncluster <- lm_robust(repub_change1216 ~ unempl_diff + prop_manufacturing + prop_  
modelsummary(final_uncluster, stars = TRUE)
```

Final Model Summary Table to compare is:

```

coef_names <- c(
  "(Intercept)" = "Intercept",
  "poverty16" = "Poverty Rate 2016",
  "unempl_rate15" = "Unemployment Rate 2015",
  "median_income16" = "Median Income 2016",
  "pct_urban" = "Urban Pop. Percentage",
  "prop_noncollege" = "Proportion Non-College",
  "rust_belt" = "Rust Belt States",
  "prop_agr" = "Proportion Agricultural Workers",
  "prop_manufacturing" = "Proportion Manufacturing Workers",
  "prop_constr" = "Proportion Construction Workers",
  "unempl_dff" = "Difference in 2015 and 2012 Unemployment Rate",
  "Pov_diff" = "Difference in Poverty Rate in 2016 and 2012"
)

```

```

# Create a modelsummary table comparing all models

```

```

modelsummary(
  list(
    "HC3 Robust" = final_rust_belt_hc3,
    "Unclustered" = final_uncluster,
    "Clustered" = final_cluster,
    "Improved Model" = improved_model,
    "Rust Belt" = final_rust_belt
  ),

```

```

stars = TRUE,
coef_map = coef_names
)

```

I think ultimately our HC3 model is the best. It has the highest  $R^2$ , and it also corrects for our hetero/normality issues. It also has the highest confidence for our coefficients, and best highlights the urban-rural, and manufacturing divides in American Politics.

```

#finalizing sutff

final <- modelsummary(
  final_rust_belt_hc3,
  stars = TRUE,
  coef_map = coef_names,
  title = "Table 1: HC3 Robust Regression Results",
)

df <- tidy(final_rust_belt_hc3, conf.int = TRUE)

df <- df[df$term != "(Intercept)", ]

final_plot <- ggplot(df, aes(x = estimate, y = reorder(term, estimate))) +
  geom_point(color = "blue", size = 3) + # Point for estimate
  geom_errorbarh(aes(xmin = conf.low, xmax = conf.high), height = 0.2, color = "black")
labs(
  title = "Figure 1: Plot for Final Robust Model",
  x = "Estimate",

```

```

y = "Variables"
) +
theme_minimal()

```

I just realized at this point I examined the influential points earlier in the project, but never explored removed them. So I am going to do that now for what was my “final” model I was going to present.

```

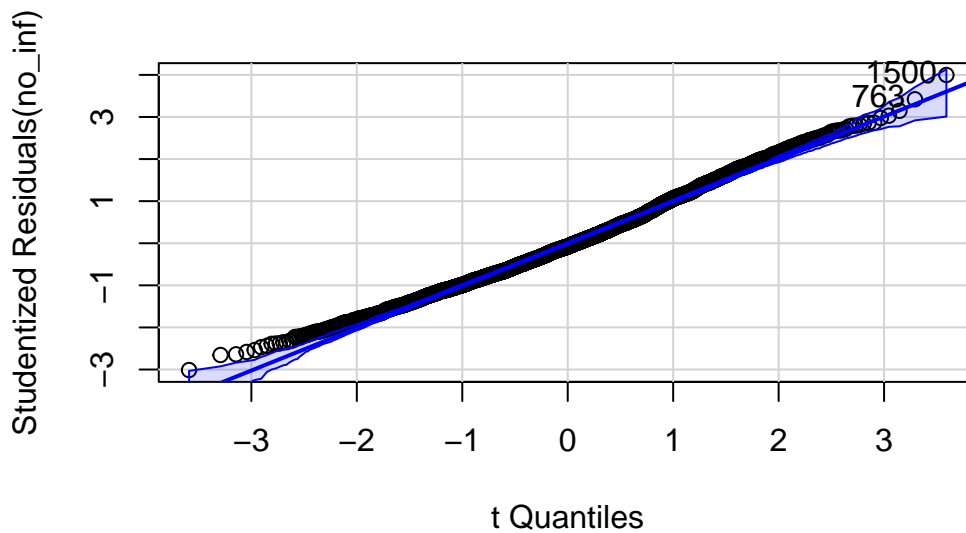
influential_points <- which(cooks.distance(final_rust_belt) > (4/nrow(election)))

election_no_influential <- election[-influential_points, ]

no_inf <- lm(repub_change1216 ~ poverty16 + rust_belt + prop_manufacturing + pct_urban +

qqPlot(no_inf)

```



763 1500

726 1440

```
bptest(no_inf)
```

studentized Breusch-Pagan test

data: no\_inf

BP = 111.85, df = 5, p-value < 2.2e-16

```
# modelsummary(no_inf)
```

```
influence <- cooks.distance(final_rust_belt)
```

```
threshold <- 4 / nrow(election)
```

```
influential_points <- which(influence > threshold)
```

```
influential_df <- election[influential_points, ]
```

```
# head(influential_df)
```

The influence points that were removed, we places that were heavy in poverty I think. Removing those, we now have a normalized qqPlot, or distribution.

```
modelsummary(  
  list("No Influential Points" = no_inf,  
        "Original Model" = final_rust_belt, "HC3" = final_rust_belt_hc3),  
  gof_map = c("r.squared", "adj.r.squared", "AIC", "BIC", "RMSE"), stars = TRUE  
)
```

Okay, so the influence points do not change the results really at all besides making our  $R^2$

higher.

#### 0.0.14 Final Write Up

This project attempted to uncover an answer to the question what caused the change in republican vote share in the 2016, and 2012 elections. There are a lot of ways to engage with this question contextually, and theoretically. The theoretical foundations for the models, and the final model for this project, beyond the statistical skills learned in the course, are taken from existing American Politics research (Fiorina 2017; Sides, Tesler, and Vavreck 2019; Hacker et al. 2024; Tesler 2016; Hacker and Pierson 2014; Brown and Mettler 2024). However, this project avoided an overly in-depth analysis with them to stay within the confines of the assignment, and focus specifically on the statistical modeling part of the problem set.

#### final

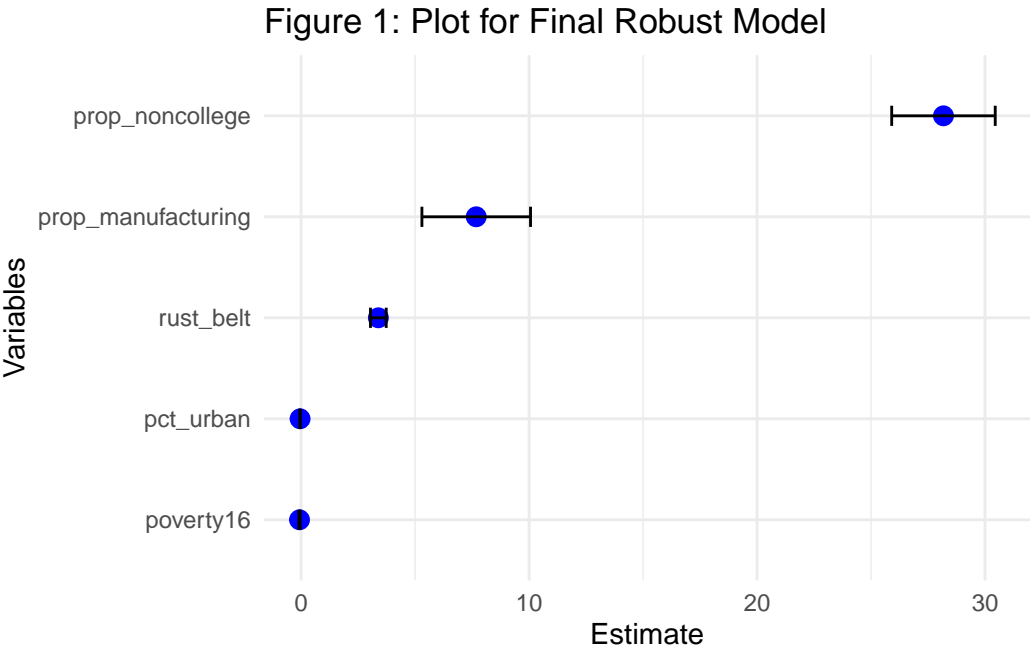
The final model this project puts forth can be seen above. This model using the variables urban population percentage, if the state is part of the rust belt or not, the 2016 poverty rate, the proportion of manufacturing workers as well as the proportion of non college graduates. Since the model struggles to satisfy the assumptions of homoskedasticity, and normally distributed residuals, this is attempted to be corrected by using robust standard errors. These variables were chosen because throughout the construction of numerous models, they were routinely the most salient. They also theoretically make sense when employing the knowledge we know about American politics, from previously mentioned literature.

The results of this model indicate key factors in the change of the republican share. Our poverty rate coefficient has a small relationship with our intercept, but this relationship is statistically significant. This relationship proves to be somewhat important, and having a larger sample size like we do likely helps confirm this. Another salient indicator, although not a tremendously large relationship, is that of the urban population percentage. When



holding other variables constant, for each 1 percentage point increase in urban population, the change in Republican vote share decreases by 0.043 percentage points. The size of the effects, and their confidence intervals can be seen below.

`final_plot`



The coefficients that were most salient, which we can see in our model table as well as our plot that showcases the estimate and the confidence intervals, are rust belt states, proportion of workers in manufacturing and proportion of non-college graduates. The proportion of non-college-educated voters has a large and positive effect. This means that counties with a higher share of non-college-educated voters saw significantly greater increases in the Republican vote share. While not as large of a relationship, the states within our rust belt variable, also have a strong, positive and statistically significant result. The same goes for areas with a high percentage of manufacturing workers, which also overlaps with rust belt states. A one percentage point increase in manufacturing workers related to a 7.682 increase in Republican vote share change from 2012 to 2016. This result was also statistically significant. These results highlight the electoral shifts made in the 2016 election, and the areas Trump had success in.

While our final model violates some OLS assumptions, as stated already, these were attempted to be corrected using robust standard errors. As demonstrated throughout the project as well, solutions were tried such as logging and transforming variables, and using diagnostic tools. The project also attempted to remove any influence points, and when doing this did produce a more *normal* distribution, however the final model presented is the model with the influence points. The reason for this is because the primary conclusions remain the same, so removing the data seemed unnecessary. The additional models without the influence points however would likely be in the appendix if this was a full article. Interactions were also explored in this project, but did not result in anything novel. Also, some of the models developed were more successful than others, and again, they would be in the final paper, or appendix were this a manuscript for an article. However, the model chosen was felt to have best overall results all things considered.

In short, I hope I have demonstrated successfully, the analytical research process used to develop a linear regression model, that aims to answer this problem set's question. To do this, I employed tools learned throughout the quarter, such as exploring the transformation of variables, running diagnostics to explore my models, and general coding skills for data presentation and visualization. I also hope that I have been successful in creating a useful model, that helps presents answers to the question itself. Ultimately, after all of this, this project finds that Trump made in-roads with citizens who were not college degree holders in the rust belt, more than anything else. And the nature of American federalism winning these unique area, was worth more than making national in-roads, which is why he won the election but lost the popular vote.

#### **0.0.15 Addendum**

As mentioned periodically in the body, this project employed Large-Language Models, specifically Open AI's ChatGPT at various points in this project. It did so when exploring the utility of log-transforming the variables. As when doing so by myself, I ran into some issues

with some of the rows being “-inf”, or N/A, which led to models not running. I also asked for some help on how to store the influence points in a data-frame, to get a rough idea what the influence points were. I did this twice. Then a few times I asked questions related to wording, defining terms, or just attempting to get a better understanding of something I was trying to talk about, or use. I also had a lot of troubling render this project, so I had to use ChatGPT to debug some LaTeX and quarto issues. Overall, ChatGPT was helpful, but generally, I only find it helpful if I have pre-existing context, or know exactly what I am looking for it. If I am troubleshooting something completely blind, or have no context, it does not really seem to be all that helpful. If not, more of a burden.

	(1)
(Intercept)	−10.061*** (1.551)
poverty16	−0.169*** (0.022)
unempl_rate15	−0.032 (0.051)
median_income16	0.000*** (0.000)
pct_urban	−0.039*** (0.003)
prop_noncollege	27.020*** (1.244)
rust_belt	3.685*** (0.192)
Num.Obs.	3111
R2	0.460
R2 Adj.	0.459
AIC	17 608.4
BIC	17 656.7
Log.Lik.	−8796.188
F	441.340
RMSE	4.09
+ p <0.1, * p <0.05, ** p <0.01, *** p <0.001	

	(1)
(Intercept)	−9.185*** (1.547)
rust_belt	3.347*** (0.198)
prop_manufacturing	7.503*** (1.166)
pct_urban	−0.041*** (0.003)
poverty16	−0.158*** (0.022)
unempl_rate15	−0.020 (0.050)
median_income16	0.000*** (0.000)
prop_noncollege	24.589*** (1.293)
Num.Obs.	3111
R2	0.467
R2 Adj.	0.466
AIC	17 569.1
BIC	17 623.5
Log.Lik.	−8775.558
F	389.137
RMSE	4.06

+ p <0.1, \* p <0.05, \*\* p <0.01,  
\*\*\* p <0.001

	(1)
(Intercept)	−17.252*** (0.858)
rust_belt	3.388*** (0.195)
prop_manufacturing	7.682*** (1.172)
pct_urban	−0.045*** (0.003)
poverty16	−0.070*** (0.013)
prop_noncollege	28.177*** (1.138)
Num.Obs.	3111
R2	0.461
R2 Adj.	0.460
AIC	17 605.3
BIC	17 647.6
Log.Lik.	−8795.632
F	530.190
RMSE	4.09
+ p <0.1, * p <0.05, ** p <0.01, *** p <0.001	

	(1)
(Intercept)	−17.252 (0.858)
poverty16	−0.070 (0.013)
prop_noncollege	28.177 (1.138)
prop_manufacturing	7.682 (1.172)
rust_belt	3.388 (0.195)
pct_urban	−0.045 (0.003)
Num.Obs.	3111
R2	0.461
R2 Adj.	0.460
AIC	17 605.3
BIC	17 647.6
Log.Lik.	−8795.632
F	530.190
RMSE	4.09

	(1)
(Intercept)	−17.252*** (0.830)
poverty16	−0.070*** (0.014)
rust_belt	3.388*** (0.176)
pct_urban	−0.045*** (0.003)
prop_manufacturing	7.682*** (1.215)
prop_noncollege	28.177*** (1.158)
Num.Obs.	3111
R2	0.461
R2 Adj.	0.460
AIC	17 605.3
BIC	17 647.6
RMSE	4.09
+ p <0.1, * p <0.05, ** p <0.01, *** p <0.001	



	(1)
(Intercept)	−14.169*** (1.987)
unempl_diff	0.303 (0.244)
prop_manufacturing	9.266** (3.571)
prop_agr	−12.249** (3.940)
prop_constr	−15.313** (5.917)
pct_urban	−0.056*** (0.007)
prop_noncollege	27.430*** (2.724)
Pov_diff	0.164** (0.061)
Num.Obs.	3111
R2	0.415
R2 Adj.	0.414
AIC	17 860.9
BIC	17 915.3
RMSE	4.26
Std.Errors	by: state_v2
+ p <0.1, * p <0.05, ** p <0.01, *** p <0.001	

	(1)
(Intercept)	−14.169*** (0.834)
unempl_diff	0.303*** (0.060)
prop_manufacturing	9.266*** (1.438)
prop_agr	−12.249*** (1.411)
prop_constr	−15.313*** (3.792)
pct_urban	−0.056*** (0.003)
prop_noncollege	27.430*** (0.987)
Pov_diff	0.164*** (0.041)
Num.Obs.	3111
R2	0.415
R2 Adj.	0.414
AIC	17 860.9
BIC	17 915.3
RMSE	4.26

+ p <0.1, \* p <0.05, \*\* p <0.01,  
\*\*\* p <0.001

	HC3 Robust	Unclustered	Clustered	Improved Mod
Intercept	−17.252*** (0.830)	−14.169*** (0.834)	−14.169*** (1.987)	−14.169*** (0.914)
Poverty Rate 2016	−0.070*** (0.014)			
Urban Pop. Percentage	−0.045*** (0.003)	−0.056*** (0.003)	−0.056*** (0.007)	−0.056*** (0.003)
Proportion Non-College	28.177*** (1.158)	27.430*** (0.987)	27.430*** (2.724)	27.430*** (1.090)
Rust Belt States	3.388*** (0.176)			
Proportion Agricultural Workers		−12.249*** (1.411)	−12.249** (3.940)	−12.249*** (1.446)
Proportion Manufacturing Workers	7.682*** (1.215)	9.266*** (1.438)	9.266** (3.571)	9.266*** (1.373)
Proportion Construction Workers		−15.313*** (3.792)	−15.313** (5.917)	−15.313*** (3.730)
Difference in Poverty Rate in 2016 and 2012		0.164*** (0.041)	0.164** (0.061)	0.164*** (0.042)
Num.Obs.	3111	3111	3111	3111
R2	0.461	0.415	0.415	0.415
R2 Adj.	0.460	0.414	0.414	0.414
AIC	17 605.3	17 860.9	17 860.9	17 860.9
BIC	17 647.6	17 915.3	17 915.3	17 915.3
Log.Lik.				−8921.461
F				314.607
RMSE	4.09	4.26	4.26	4.26
Std.Errors			by: state_v2	

+ p <0.1, \* p <0.05, \*\* p <0.01, \*\*\* p <0.001

	No Influential Points	Original Model	HC3
(Intercept)	−16.560*** (0.681)	−17.252*** (0.858)	−17.252*** (0.830)
poverty16	−0.103*** (0.011)	−0.070*** (0.013)	−0.070*** (0.014)
rust_belt	3.105*** (0.152)	3.388*** (0.195)	3.388*** (0.176)
prop_manufacturing	7.756*** (0.930)	7.682*** (1.172)	7.682*** (1.215)
pct_urban	−0.042*** (0.002)	−0.045*** (0.003)	−0.045*** (0.003)
prop_noncollege	28.202*** (0.909)	28.177*** (1.138)	28.177*** (1.158)
R2	0.574	0.461	0.461
R2 Adj.	0.573	0.460	0.460

+ p < 0.1, \* p < 0.05, \*\* p < 0.01, \*\*\* p < 0.001

Table 1: Table 1: HC3 Robust Regression Results

	(1)
Intercept	−17.252*** (0.830)
Poverty Rate 2016	−0.070*** (0.014)
Urban Pop. Percentage	−0.045*** (0.003)
Proportion Non-College	28.177*** (1.158)
Rust Belt States	3.388*** (0.176)
Proportion Manufacturing Workers	7.682*** (1.215)
Num.Obs.	3111
R2	0.461
R2 Adj.	0.460
AIC	17 605.3
BIC	17 647.6
RMSE	4.09

+ p <0.1, \* p <0.05, \*\* p <0.01, \*\*\* p <0.001

## References

- Brown, Trevor E., and Suzanne Mettler. 2024. “Sequential Polarization: The Development of the Rural-Urban Political Divide, 1976–2020.” *Perspectives on Politics* 22 (3): 630–58. <https://doi.org/10.1017/S1537592723002918>.
- Fiorina, Morris P. 2017. *Unstable Majorities: Polarization, Party Sorting, and Political Stalemate*. Hoover Institution Press Publication, no. 685. Stanford, California: Hoover Institution Press, Stanford University.
- Hacker, Jacob S., Amelia Malpas, Paul Pierson, and Sam Zacher. 2024. “Bridging the Blue Divide: The Democrats’ New Metro Coalition and the Unexpected Prominence of Redistribution.” *Perspectives on Politics* 22 (3): 609–29. <https://doi.org/10.1017/S1537592723002931>.
- Hacker, Jacob S., and Paul Pierson. 2014. “After the ‘Master Theory’: Downs, Schattschneider, and the Rebirth of Policy-Focused Analysis.” *Perspectives on Politics* 12 (3): 643–62. <https://doi.org/10.1017/S1537592714001637>.
- Sides, John, Michael Tesler, and Lynn Vavreck. 2019. *Identity Crisis: The 2016 Presidential Campaign and the Battle for the Meaning of America*. New edition. Princeton: Princeton University Press.
- Tesler, Michael. 2016. *Post-Racial or Most-Racial? Race and Politics in the Obama Era*. Chicago Studies in American Politics. Chicago London: The University of Chicago Press.