

Lab 3: Reducing Crime (DRAFT: Stage 1)

N. Akkineni, K. Hanna, A. Thorp

November 27, 2018

Contents

Introduction (Stage 1: Draft Report)	1
Exploratory Data Analysis	2
Data Summary	2
Data Clean Up	3
Concerns about the data	3
Misabeled Region	4
Univariate Analysis	4
Key Variable	4
Explanatory Variables	5
Correlation Matrix	8
Model Analysis	10
Model1 - Minimum Specification	11
Model2 - Optimal Specification	14
Model3 - Optimal-2 Specification Best-Fit model	17
Model4 - Using all variables Specification	20
Omitted Variables	23
Conclusion	26
Appendix A: Codebook	27

Introduction (Stage 1: Draft Report)

Our team has been hired to provide research for a local political campaign, which would like to understand the determinants of crime rates and to provide policy suggestions that are applicable to local government. We examine the provided dataset to determine if a model with causal interpretation is feasible. After examining the data, we detail four regression models and find that estimators related to variables used to operationalize the concept of certainty of punishment are directionally consistent and statistically significant. From this we draw a limited policy recommendation to adopt a model of community policing to improve trust and information flow to law enforcement. However, our policy recommendations are limited because omitted variable bias confounds our estimators. Should local officials desire more robust conclusions, we recommend involving data scientists in the data collection process to improve our ability to draw causal inference from our modeling process and thus be able to make robust policy recommendations.

```
library(knitr)
library(kableExtra)
suppressMessages(library(car))
suppressMessages(library(stargazer))
suppressMessages(library(lmtest))
library(corrplot)
```

```
## corrplot 0.84 loaded

crime <- read.csv('crime_v2.csv',header = TRUE, sep = ",")

# Convert columns to factors and logical.
crime$west <- as.logical(crime$west)
crime$central <- as.logical(crime$central)
crime$urban <- as.logical(crime$urban)

codebook <- read.csv('codebook.csv')

# Create a variable for counties not in west or central.
crime$other <- !(crime$west | crime$central)

# Average of all weekly wage variables.
crime$avgwage = (crime$wcon + crime$wtuc + crime$wtrd + crime$wfir +
                 crime$wser + crime$wmfg + crime$wfed + crime$wsta +
                 crime$wloc)/9

# Used for displaying subsets of variables.
columns_pretty_sort = c("county", "year", "crm rte", "prbarr", "prbconv",
                        "prbpris", "avgsen", "mix", "polpc", "density",
                        "taxpc", "other", "central", "west", "urban",
                        "pctmin80", "pctymle", "avgwage", "wcon", "wtuc", "wtrd",
                        "wfir", "wser", "wmfg", "wfed", "wsta", "wloc")
```

Exploratory Data Analysis

Data Summary

We were provided with a dataset that includes crime statistics from the North Carolina Department of Corrections' prison and probation files, demographic statistics taken from the decennial census, police data derived from FBI police agency data and wage data from the North Carolina Employment Security commission. In all we were provided with 25 variables and 90 counties.

Some of the values in this dataset were calculated from other datasets and we found some characteristics in the dataset which may bring its veracity in to question and we have addressed those below and in our analyses.

Table 1: Crime Data Codebook

Variable	Label
county	county identifier
year	1987
crm rte	crimes committed per person
prbarr	'probability' of arrest
prbconv	'probability' of conviction
prbpris	'probability' of prison sentence
avgsen	avg. sentence, days
polpc	police per capita
density	people per sq. mile

Table 1: Crime Data Codebook (*continued*)

Variable	Label
taxpc	tax revenue per capita
west	=1 if in western N.C.
central	=1 if in central N.C.
urban	=1 if in SMSA
pctmin80	perc. Minority, 1980
wcon	weekly wage, construction
wtuc	wkly wge, trns, util, commun
wtrd	wkly wge whlesle, retail, trade
wfir	wkly wge, fin, ins, real est
wser	wkly wge, service industry
wmfg	wkly wge, manufacturing
wfed	wkly wge fed employees
wsta	wkly wge state employees
wloc	wkly wge local gov emps
mix	offense mix: face-to-face/other
pctymle	percent young male

Data Clean Up

Null Rows

The dataset contained an apostrophe 6 rows after the data which caused the csv reader to create 6 invalid rows. We removed these rows as they contain no data.

```
# Delete the 6 empty observations at the end, including the row with the apostrophe.
crime <- crime[!is.na(crime$county) & !is.na(crime$crmrte), ]
```

```
# Convert prbconv to numeric
crime$prbconv = as.numeric(as.character(crime$prbconv))
```

We found two identical observations for county 193. There is no logical reason to have two identical observations in this cross-sectional dataset, so we feel removing one of these two observations can only improve the quality of our analysis.

```
# county 193 is duplicated, remove one
crime = crime[!duplicated(crime), ]
```

Concerns about the data

prbarr (Probability of Arrest)

We found that county 115 contained a value of 1.09 in prbarr (probability of arrest) which is not possible. We believe this to be a labeling error, as it is a ratio used to approximate a probability, rather than a true probability.

prbconv (Probability of Conviction)

We found 10 observations with values greater than 1, which, again, is not a possible value for probability. The documentation in the codebook specifies that “(t)he probability of conviction is proxied by the ratio of

convictions to arrests”, which leaves some ambiguity, however it is plausible to have values greater than 1 as a single arrest can result in multiple convictions and persons can be convicted in absentia. Similar to the prbarr, we believe this to be a labeling error, as it is a ratio used to approximate a probability, rather than a true probability.

Omitted Counties Creating Bias

The dataset contained 90 counties, and there are 100 counties in North Carolina. The observation id labeled ‘county’ in our data set appears to contain FIPS codes, if this assumption is correct the following are the missing counties: Camden County, Carteret County, Clay County, Gates County, Graham County, Iredell County, Jones County, Mitchell County, Tyrrell County and Yancey County

These missing counties will introduce a slight clustering bias into our analyses at a minimum, and possibly a more significant bias if they were omitted based on specific criteria whether deliberate or not.

Mislabeled Region

County 87 is both part of the West region and Central region. It’s possible to straddle the border however, we’d expect more instances if it was done in this manner, thus we expect this is mislabeled and should belong in only one region. This is the only case where this occurs.

Univariate Analysis

```
#hist(log(crime$crmrte))

#hist(log(crime$density), breaks = 20)
#hist(crime$density)
#hist(log(crime$pctmin80))
#hist(poly(log(crime$pctmin80),2))
#hist(log(crime$pctmle)^2)
#hist(poly(log(crime$taxpc), 4))
```

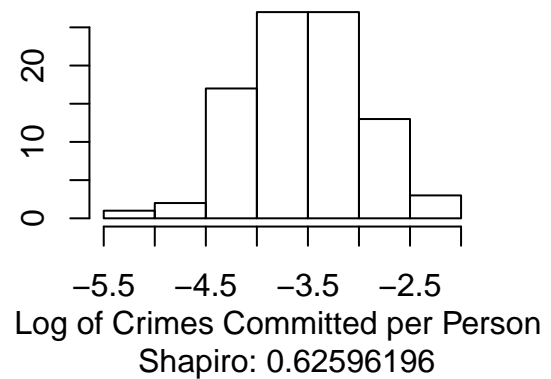
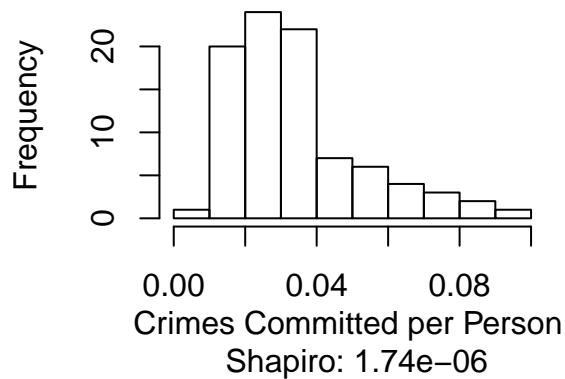
Key Variable

```
quick_uni_analysis = function(variable, description, roundto = 8) {
  hist(variable, xlab = paste(tools::toTitleCase(description),
    paste('\n Shapiro:',
      round(as.numeric(shapiro.test(variable)[2]), roundto)
    )), main = "")
  hist(log(variable),
    xlab = tools::toTitleCase(paste('Log of', description,
      paste('\n Shapiro:', ... = round(as.numeric(shapiro.test(log(variable))[2]), roundto)
    )), main = "", ylab = '')
}
```

Crimes Committed Per Person

Campaign Significance: The political campaign which hired us is interested in policy prescriptions derived from causal analysis of crime rates. This is the key variable our models will attempt to explain.

```
par(mfrow=c(1,2))
quick_uni_analysis(crime$crm rte, 'crimes committed per person')
```



Crimes committed per capita has a fairly strong positive skew, applying a natural log transformation creates a more symmetrical distribution and results in a Shapiro-Wilk test p-value that we cannot reject. The transformed variable is preferable for modelling.

```
crm rte.outliers = boxplot(crime$crm rte, plot = FALSE)$out
length(crm rte.outliers)
```

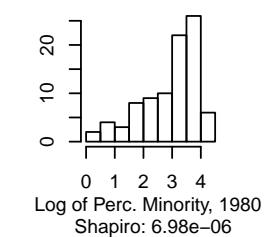
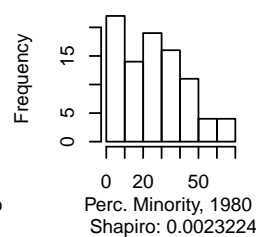
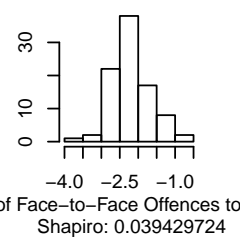
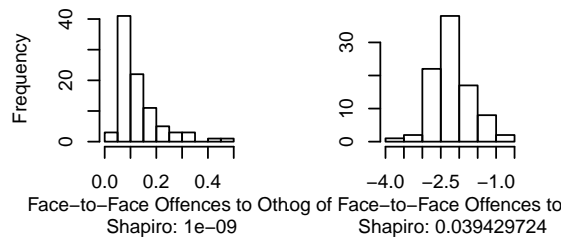
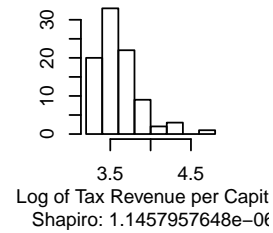
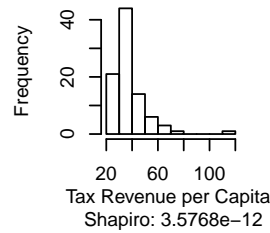
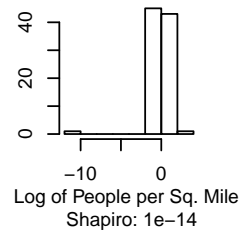
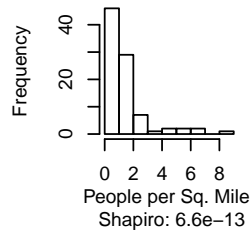
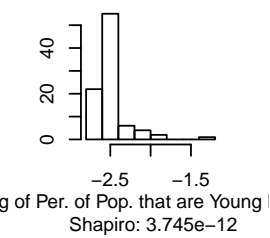
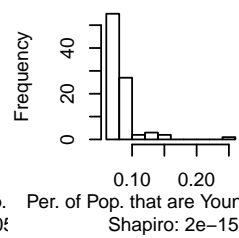
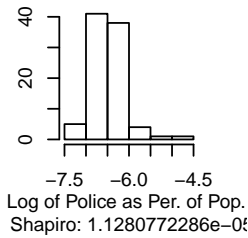
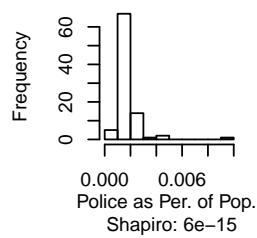
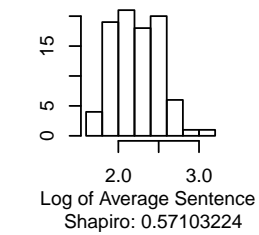
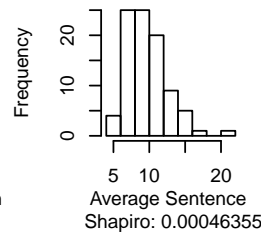
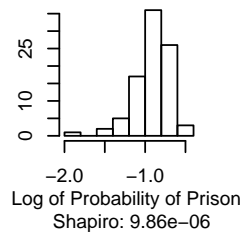
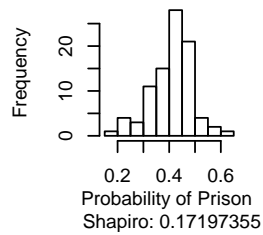
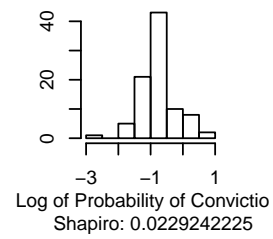
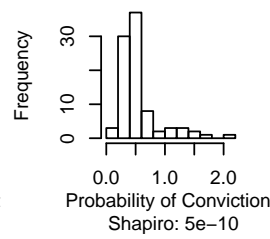
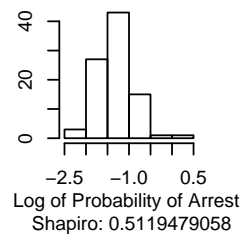
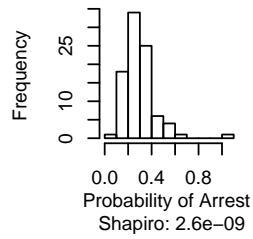
```
## [1] 5
```

Crime rate has 5 outliers, none that are extreme or causing concern.

Explanatory Variables

Diagrams of Key Variables With and Without Log Transformations

```
par(mfrow=c(5,4))
quick_uni_analysis(crime$prbarr, 'Probability of Arrest', roundto = 10)
quick_uni_analysis(crime$prbconv, 'Probability of Conviction', roundto = 10)
quick_uni_analysis(crime$prbpris, 'Probability of Prison')
quick_uni_analysis(crime$avg sen, 'Average Sentence')
quick_uni_analysis(crime$polpc, 'Police as Per. of Pop.', roundto = 15)
quick_uni_analysis(crime$pctymle, 'Per. of Pop. That Are Young Males', roundto = 15)
quick_uni_analysis(crime$den sity, 'people per sq. mile', roundto = 14)
quick_uni_analysis(crime$taxpc, 'tax revenue per capita', roundto = 16)
quick_uni_analysis(crime$mix, 'Face-to-face offences to other', roundto = 9)
quick_uni_analysis(crime$pctmin80, 'perc. minority, 1980')
```



```
par(mfrow=c(1,1)) #Reset
```

Probability of Arrest

Probability of Arrest has a positive skew, applying a natural log transformation creates a more symmetrical distribution and results in a Shapiro-Wilk test p-value that we cannot reject the null hypothesis of normality. The transformed variable is preferable for modelling.

Probability of Conviction

The log transform is preferable - both for interpretation and for better adhering to modeling assumptions. However, even the logged version fails a Shapiro-Wilk normality test. Something to keep in mind.

Probability of Prison

From an interpretation standpoint, the logged version is preferable, although from an modeling assumption standpoint, the unlogged version is preferable.

Average Sentence

The logged version is preferable from both an interpretation and modeling assumption standpoint.

Police as a Percentage of Population

Both logged and un-logged versions of police as a percentage of the population are non-normal. Neither is inherently preferable from a modeling assumptions standpoint.

Both logged and un-logged versions of the percent of population that is young and male are non-normal. Neither is inherently preferable from a modeling assumptions standpoint. ⁽¹⁾ The Journal of Law & Economic <https://www.jstor.org/stable/10.1086/666614>

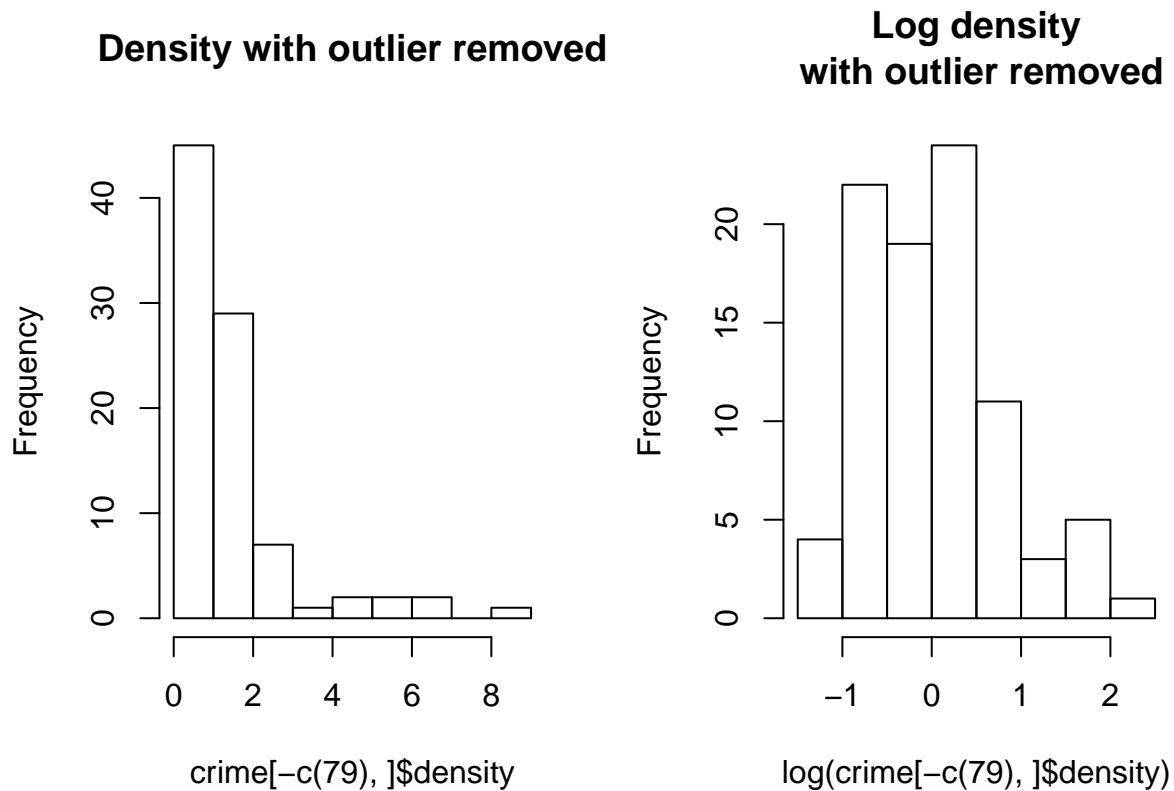
Percentage of Population That Are Young Males

This variable has a strong positive skew, using a natural log transformation results in a distribution that is still skewed, however, it is closer to normal.

People per Square Mile

There is one fairly extreme outlier for county 173, with 0.000023 people per square mile. This will affect our modeling, specifically the cooks distance. Without that outlier, this variable is more normal with a log transformation.

```
par(mfrow=c(1,2))
hist(crime[-c(79), ]$density, main = "Density with outlier removed")
hist(log(crime[-c(79), ]$density), main = "Log density\n with outlier removed")
```



Tax Revenue per Capita

Tax revenue per capita has a strong positive skew, using a natural log transformation results in a distribution that is still skewed, however, it is closer to normal.

Face-to-face offences to Other (Offence Mix)

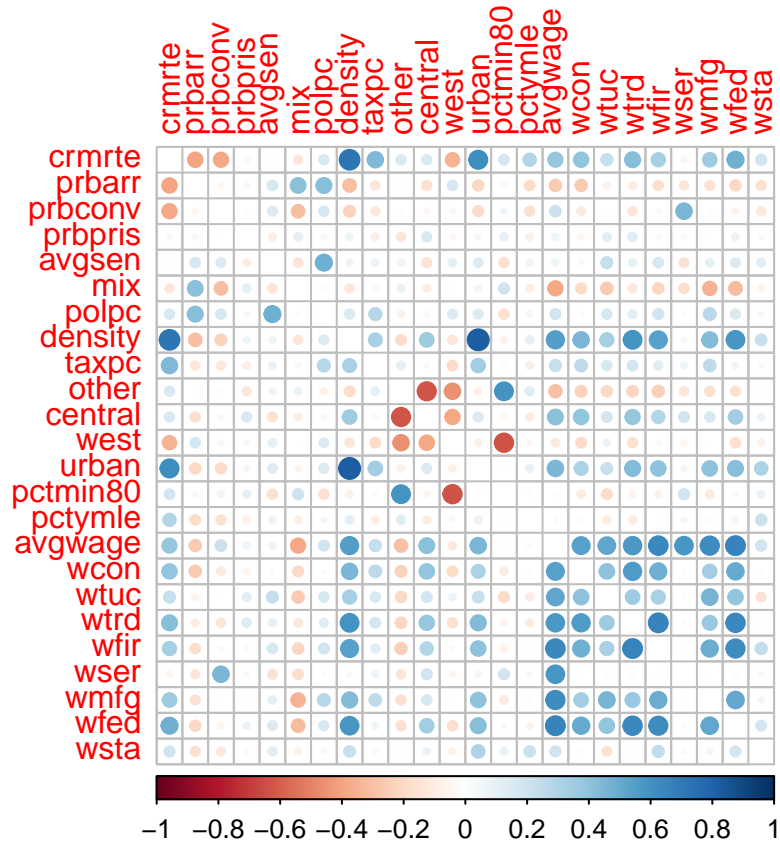
This is a ratio of face-to-face crimes to all other crimes. Face-to-face crimes include violent crimes and those with a higher probability of violence.

Campaign Significance: Violent crimes create fear and fear is a strong motivator for voters.

The mix of face-to-face crimes to other crimes has a positive skew, applying a natural log transformation creates a more symmetrical distribution. However, the resulting Shapiro-Wilk test would still reject the null hypothesis of normality. That said, the log transformation is preferable for modelling. ## Relationships

Correlation Matrix

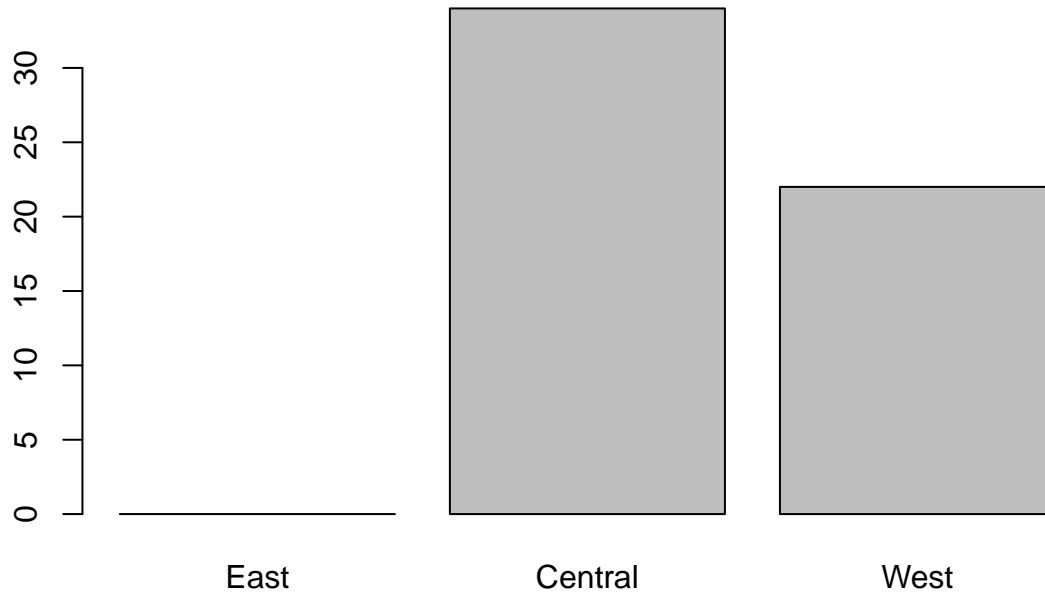
```
corrmatrix <- cor(crime[,columns_pretty_sort[3:26]])
corrplot(corrmatrix, cl.pos = "b", diag = FALSE)
```

Using the correlation matrix above we can see there are some strong correlations between Crime Rate (crmrte) and Population Density (density), Crime Rate and Urban, Other (likely east) counties and Percentage Minority from 1980 (pctmin80) and West and Percentage Minority, interestingly opposite correlations in those last two. We also see strong correlation with Population Density and the Average Wage (avgwage) and several wages, specifically trade, financial, insurance real estate and federal employees.

```
barplot(c(sum(crime$east), sum(crime$central), sum(crime$west)),
        names.arg = c("East", "Central", "West"), main = "County locations")
```

County locations



```
crime$region <- ifelse(crime$west == 1, "west",
                      ifelse(crime$central == 1, "central",
                            ifelse(crime$east == 1, "east", "other")))
region = aggregate(density ~ region, data = crime, mean)
region$polpc = aggregate(polpc ~ region, data = crime, mean)[2]
region$crmrte = aggregate(crmrte ~ region, data = crime, mean)[2]
colnames(region) = c("Region", "Density", "Police per Cap", "Crime Rate")

#region
kable(region)
```

Region	Density	Police per Cap	Crime Rate
central	2.047960	0.001637046	0.03699627
west	1.074319	0.001970259	0.02209975

```
#kable(region, "latex", longtable = TRUE, booktabs = TRUE, caption = "Regions") #>%
#   kable_styling(full_width = TRUE, latex_options = c("HOLD_position", "striped", "repeat_header"), r
```

Model Analysis

Transforming the variables and storing them in the data frame

```
crime$log_crmrte <- log(crime$crmrte)
crime$log_prbarr <- log(crime$prbarr)
crime$log_prbconv <- log(crime$prbconv)
crime$log_prbpris <- log(crime$prbpris)
```

```

crime$log_avgsen <- log(crime$avgsen)
crime$log_polpc <- log(crime$polpc)
crime$log_density <- log(crime$density)
crime$log_taxpc <- log(crime$taxpc)
crime$log_pctmin80 <- log(crime$pctmin80)
crime$log_mix <- log(crime$mix)
crime$log_pctymle <- log(crime$pctymle)

```

Model1 - Minimum Specification

We anticipate that the crime rate depends on certainty of punishment and severity of punishment. As such, our base model includes variables that attempt to operationalize those concepts. The probability of arrest and probability of punishment operationalize certainty of punishment, while probability of prison and sentence length operationalize severity of punishment.

$$\log(\text{crmte}) = \beta_0 + \beta_1 \log(\text{prbarr}) + \beta_2 \log(\text{prbconv}) + \beta_3 \log(\text{prbpris}) + \beta_4 \log(\text{avgsen}) + u$$

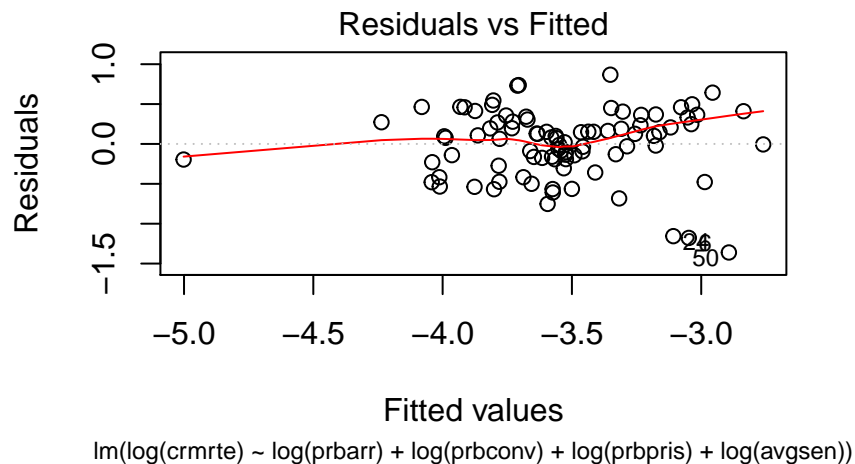
```

model1 = lm(log(crmte) ~ log(prbarr) + log(prbconv) +
            log(prbpris) + log(avgsen), data = crime)

```

Assumption - Linear model

```
plot(model1, which=1, cex.sub=0.75)
```



From the residuals vs fitted plot, we don't see a non-linear relationship, so we treat this as a valid assumption.

Assumption - Random Sampling

We do not have enough information about the sample selection to investigate this claim. We know that our dataset includes 90 counties and that North Carolina has only 100 counties, so our sample is close to approximating the population. However, we do not know how the counties in our dataset were selected. This uncertainty leads us to be more cautious in interpreting results from our models.

Assumption - Multicollinearity

```
# Find the maximum Correlation
cor_model1 <- data.matrix(subset(crime,
                                select = c("log_crmrte", "log_prbarr", "log_prbconv", "log_prbpris", "log_avgsen")))
cor_model1 = cor(cor_model1)
diag(cor_model1) = 0 # Zero the diagonal as we are uninterested in those.
max(cor_model1)
```

```
## [1] 0.06960729
```

```
vif(model1)
```

```
## log(prbarr) log(prbconv) log(prbpris) log(avgsen)
## 1.043435 1.042909 1.016372 1.015728
```

Based on pairwise correlation and variance inflation factors, we do not detect evidence of multicollinearity negatively impacting our specification.

Assumption - 4 Exogeneity (Zero Conditional Mean)

From the residuals vs fitted plot, the red line is very influenced by the outliers on the ends, though it remains largely centered on zero, so we treat this as most-likely a valid assumption.

Additional Assumptions

Homoscedasticity

From the same residuals vs fitted plot, we see it is very scattered with extreme outliers. So, it is not easy to determine Homoscedasticity from this plot only.

```
bptest(model1)
```

```
##
## studentized Breusch-Pagan test
##
## data: model1
## BP = 10.702, df = 4, p-value = 0.03012
```

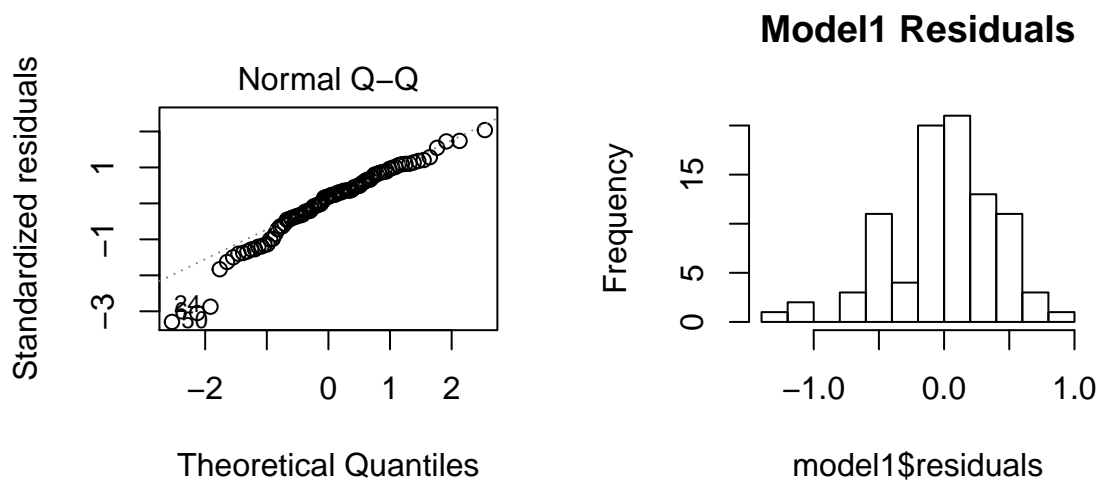
```
ncvTest(model1)
```

```
## Non-constant Variance Score Test
## Variance formula: ~ fitted.values
## Chisquare = 7.467969, Df = 1, p = 0.0062806
```

Both tests are showing small p-values showing that we have to reject the hypothesis. Homoscedasticity does not appear to be a valid assumption here indicating that our standard errors may not be used for inference. We will need to return to this in future analyses to examine our ability to use robust standard errors or other compensating techniques in order to make inference from our models.

Normality of Residuals

```
plot(model1, which=2)
hist(model1$residuals, main="Model1 Residuals")
```

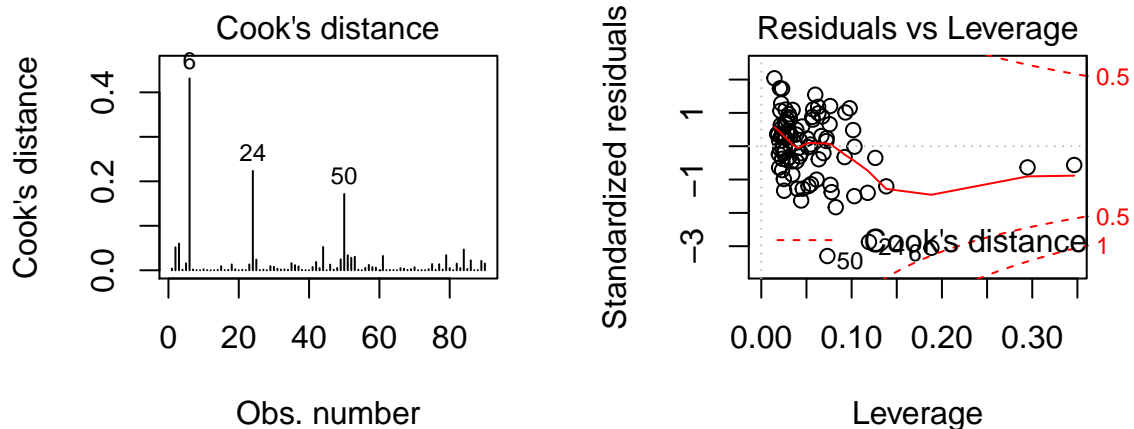


$y \sim \log(\text{prbarr}) + \log(\text{prbconv}) + \log(\text{prbp})$
Other

than a few outliers, the distribution is relatively normal for our given sample size. We treat this as a valid assumption.

Cook's Distance:

```
plot(model1, which=4)
plot(model1, which=5)
```



$y \sim \log(\text{prbarr}) + \log(\text{prbconv}) + \log(\text{prbp})$
There

are some influential values however cook's distance is within the bounds.

Calculating AIC

```
AIC(model1)
```

```
## [1] 109.9382
```

The AIC for this model is 110.0643.

Model2 - Optimal Specification

In addition to the explanatory variables introduced in our #Model1, we have decided to include the following variables in the model.

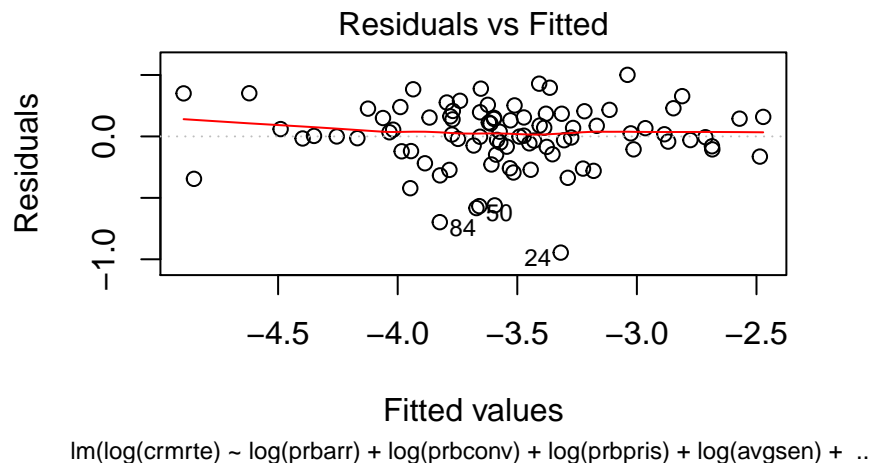
$$\begin{aligned} \log(\text{crmrte}) = & \beta_0 + \beta_1 \log(\text{prbarr}) + \beta_2 \log(\text{prbconv}) + \beta_3 \log(\text{prbpris}) + \\ & \beta_4 \log(\text{avgsen}) + \beta_5 \log(\text{polpc}) + \beta_6 \log(\text{taxpc}) + \\ & \beta_7 \log(\text{density}) + \beta_8 \log(\text{pctymle}) + \beta_9 \log(\text{pctmin80}) + u \end{aligned}$$

```
model2 = lm(log(crmrte) ~ log(prbarr) + log(prbconv) + log(prbpris) + log(avgsen)
+ log(polpc) + log(taxpc) + log(density) + log(pctymle) + log(pctmin80), data = crime)
```

Checking if our Assumptions are valid in the model-

Assumption - Linear model

```
plot(model2, which=1, cex.sub=0.75)
```



From the residuals vs fitted plot, we don't see a non-linear relationship, so we treat this as a valid assumption.

Assumption - Random Sampling

Our views are identical to the previous specification.

Assumption - Multicollinearity

```
# Find the maximum Correlation
cor_model2 <- data.matrix(subset(crime,
  select = c("log_crmrte", "log_prbarr", "log_prbconv", "log_prbpris",
    "log_avgsen", "log_polpc", "log_taxpc", "log_density", "log_pctymle", "log_pctmin80")))
cor_model2 = cor(cor_model2)
diag(cor_model2) = 0 # Zero the diagonal as we are uninterested in those.
max(cor_model2)
```

```
## [1] 0.4936425
```

```
vif(model2)
```

```
##   log(prbarr)  log(prbconv)  log(prbpris)  log(avgsen)  log(polpc)
##   1.469188    1.297357    1.344846    1.332245    1.587111
##   log(taxpc)  log(density)  log(pctymle) log(pctmin80)
##   1.465737    1.595347    1.365754    1.105714
```

Based on pairwise correlation and variance inflation factors, we do not detect evidence of multicollinearity negatively impacting our specification.

Assumption - Exogeneity (Zero Conditional Mean)

From the residuals vs fitted plot, the red line is influenced by outliers on either end, however the zero conditional mean appears to hold.

Additional Assumptions - Homoscedasticity

From the same residuals Vs fitted plot, we see that it is scattered with extreme outliers. It is not easy to determine homoscedasticity from this plot only. Running some additional tests

```
bptest(model2)
```

```
##
## studentized Breusch-Pagan test
##
## data:  model2
## BP = 17.355, df = 9, p-value = 0.04344
```

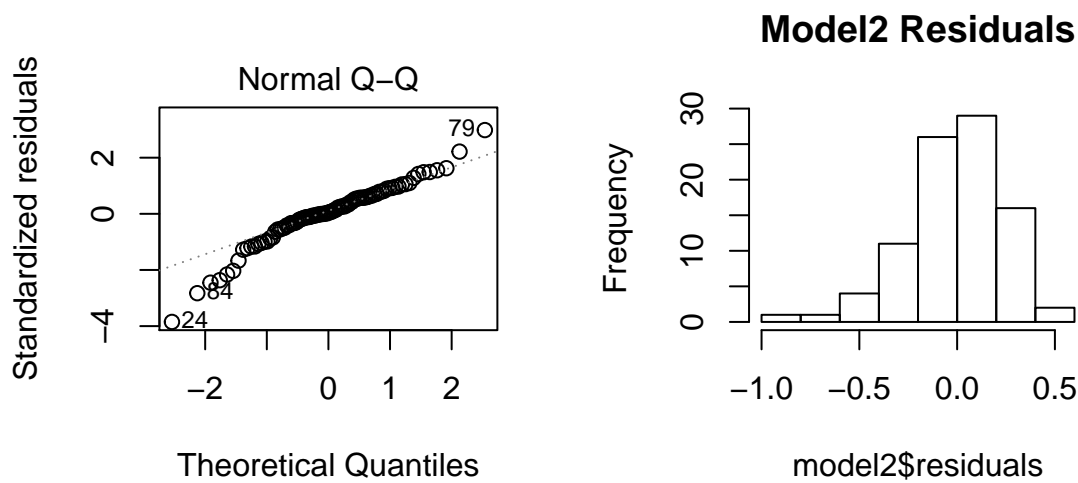
```
ncvTest(model2)
```

```
## Non-constant Variance Score Test
## Variance formula: ~ fitted.values
## Chisquare = 0.5980037, Df = 1, p = 0.43934
```

Both tests are showing large p-values showing that we fail to reject the null hypothesis.

```
#####Normality of Residuals
```

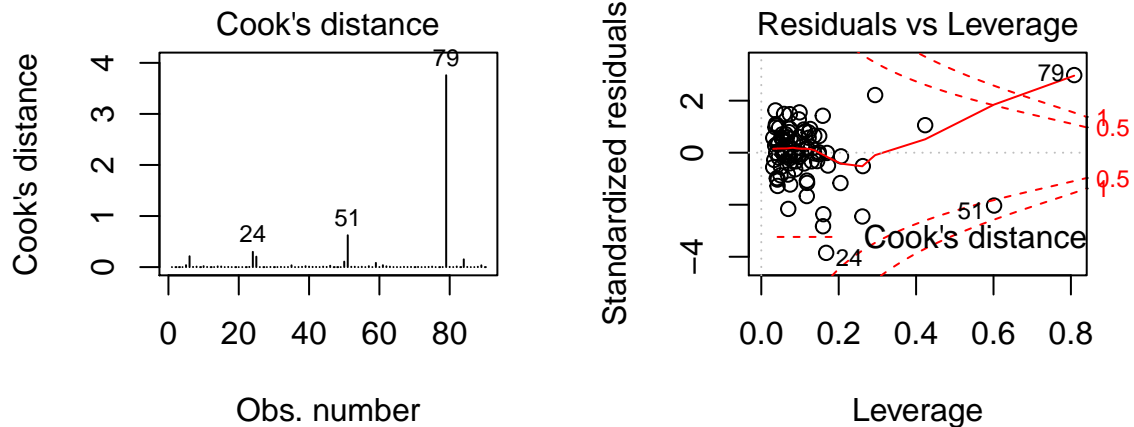
```
plot(model2, which=2)
hist(model2$residuals,main="Model2 Residuals")
```



Other than a few outliers, the distribution is relatively normal for our given sample size, so we treat it as a valid assumption.

####Cook's Distance:

```
plot(model2, which=4)
plot(model2, which=5)
```



There are some influential values [24,51,79] however cook's distance is within bounds.

```
AIC(model2)
```

```
## [1] 30.5825
```

The AIC for this model is 30.5825

Model3 - Optimal-2 Specification Best-Fit model

After observing Model2 - We found that tax and percent male doesn't influence crime extensively, so we remove those in our best-fit model

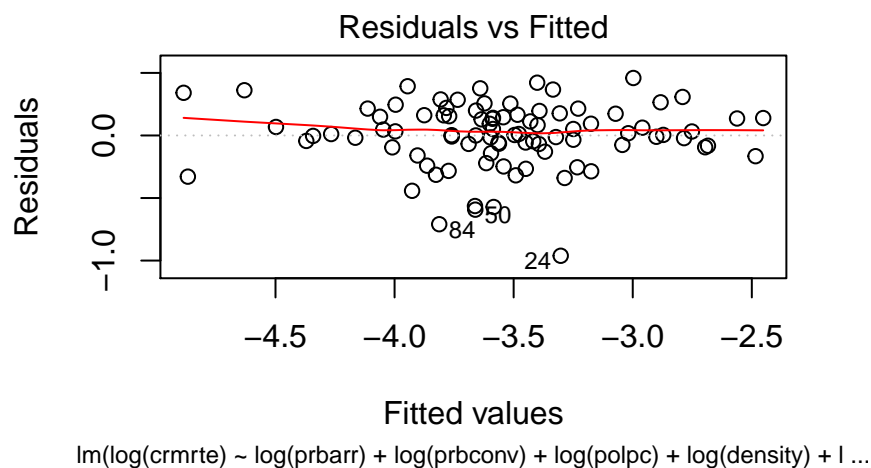
$$\log(\text{crmrte}) = \beta_0 + \beta_1 \log(\text{prbarr}) + \beta_2 \log(\text{prbconv}) + \beta_3 \log(\text{prbpris}) + \beta_4 \log(\text{avgsen}) + \beta_5 \log(\text{polpc}) + \beta_6 \log(\text{density}) + \beta_7 \log(\text{pctmin80}) + u$$

```
model3 = lm(log(crmrte) ~ log(prbarr) + log(prbconv) + log(polpc) +  
            log(density) + log(pctmin80) + log(prbpris) + log(avgsen) , data = crime)
```

Checking if our Assumptions are valid in the model-

Assumption - Linear model

```
plot(model3, which=1, cex.sub=0.75)
```



From the residuals vs fitted plot, we don't see a non-linear relationship, so we treat this as a valid assumption.

Assumption - Random Sampling

Our views are identical to the previous model.

Assumption - Multicollinearity

```
# Find the maximum Correlation  
cor_model3 <- data.matrix(subset(crime,  
                                select = c("log_crmrte", "log_prbarr", "log_prbconv",  
                                            "log_avgsen", "log_polpc", "log_density", "log_pctmin80")))  
cor_model3 = cor(cor_model3)  
diag(cor_model3) = 0 # Zero the diagonal as we are uninterested in those.  
max(cor_model3)
```

```
## [1] 0.4936425
```

```
vif(model3)
```

```
##    log(prbarr)  log(prbconv)    log(polpc)  log(density) log(pctmin80)
##      1.264585      1.111820      1.237717      1.571169      1.044778
##  log(prbpris)    log(avgsen)
##      1.318298      1.321456
```

Based on pairwise correlation and variance inflation factors, we do not detect evidence of multicollinearity negatively impacting our specification.

Assumption - Exogeneity (Zero Conditional Mean)

From the residuals vs fitted plot, the red line is very influenced by the outlier on either end, however the zero conditional mean appears to hold.

Additional Assumptions - Homoscedasticity

From the same residuals vs fitted plot, we see it is scattered with extreme outliers. So, it is not easy to determine Homoscedasticity from this plot only.

```
bptest(model3)
```

```
##
##  studentized Breusch-Pagan test
##
## data:  model3
## BP = 8.4026, df = 7, p-value = 0.2984
```

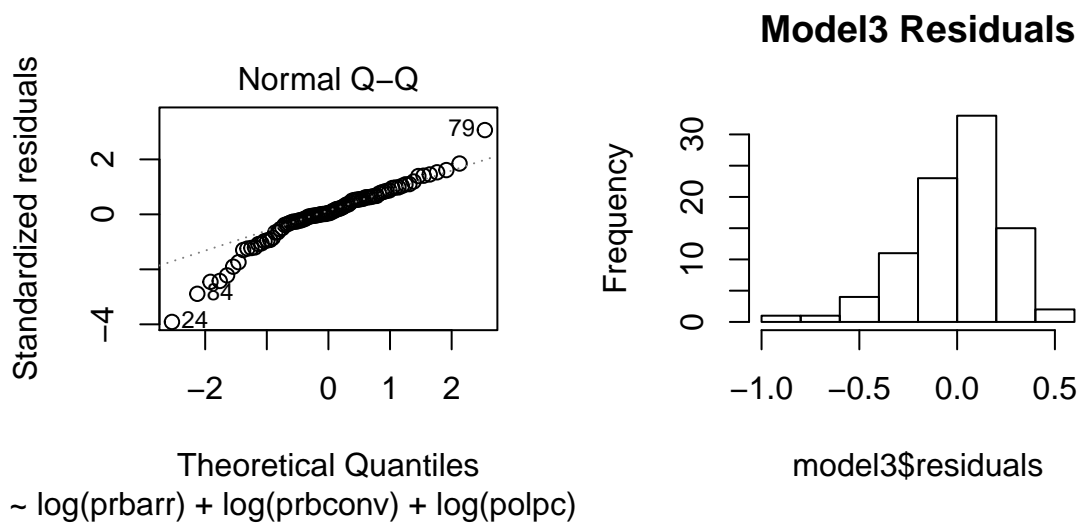
```
ncvTest(model3)
```

```
## Non-constant Variance Score Test
## Variance formula: ~ fitted.values
## Chisquare = 0.5987105, Df = 1, p = 0.43907
```

Both tests are showing insignificant p-values showing that we fail to reject the null hypothesis of homoskedasticity.

Normality of Residuals

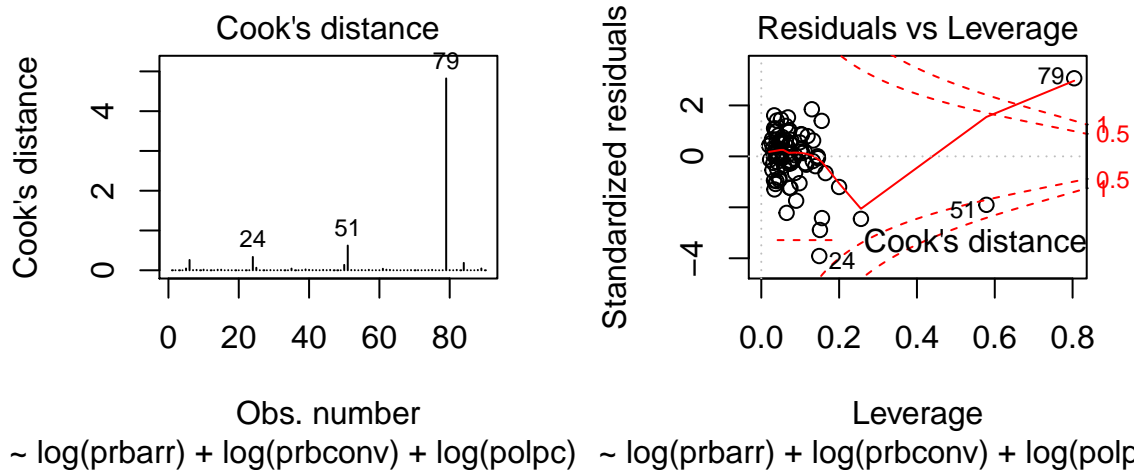
```
plot(model3, which=2)
hist(model3$residuals, main="Model3 Residuals")
```



Other than a few outliers, the distribution is relatively normal for our given sample size, so we treat this as a valid assumption.

Cook's Distance:

```
plot(model3, which=4)
plot(model3, which=5)
```



There are some influential values [6,51,79] however cook's distance is within the bounds.

Calculating AIC

```
AIC(model3)
```

```
## [1] 26.9529
```

The AIC for this model is 26.9529

Model4 - Using all variables Specification

The team wanted to check the robustness of prior model specifications by creating a model with maximal inclusion from our dataset to see if estimated relationships shift.

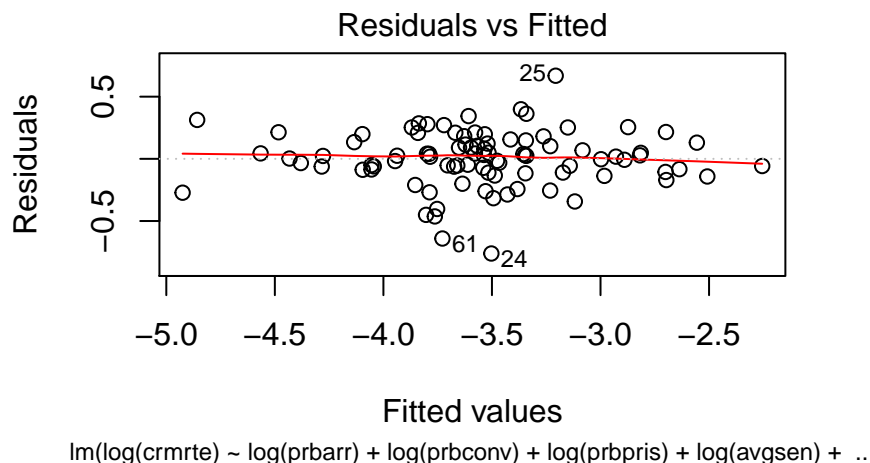
$$\begin{aligned} \log(\text{crmte}) = & \beta_0 + \beta_1 \log(\text{prbarr}) + \beta_2 \log(\text{prbconv}) + \beta_3 \log(\text{prbpris}) + \\ & \beta_4 \log(\text{avgse}) + \beta_5 \log(\text{polpc}) + \beta_6 \log(\text{taxpc}) + \beta_7 \log(\text{west}) + \\ & \beta_8 \log(\text{central}) + \beta_9 \log(\text{urban}) + \beta_{10} \log(\text{pctmin80}) + \beta_{11} \log(\text{wcon}) + \beta_{12} \log(\text{wtuc}) + \\ & \beta_{13} \log(\text{wtrd}) + \beta_{14} \log(\text{wfir}) + \beta_{15} \log(\text{wser}) + \beta_{16} \log(\text{wmfg}) + \beta_{17} \log(\text{wfed}) + \\ & \beta_{18} \log(\text{wsta}) + \beta_{19} \log(\text{wloc}) + \beta_{20} \log(\text{mix}) + \beta_{21} \log(\text{density}) + \beta_{22} \log(\text{pctymle}) + u \end{aligned}$$

```
model4 = lm(log(crmte) ~ log(prbarr) + log(prbconv) + log(prbpris) +  
  log(avgse) + log(polpc) + log(density) + log(taxpc) +  
  west + central + urban + log(pctmin80) + wcon + wtuc +  
  wtrd + wfir + wser + wmfg + wfed + wsta + wloc +  
  log(mix) + log(pctymle), data = crime)
```

Checking if our Assumptions are valid in the model-

Assumption - Linear model

```
plot(model4, which=1, cex.sub=0.75)
```



Assumption - Random Sampling

Our views are identical to the previous model.

Assumption - Multicollinearity

```
vif(model4)
```

##	log(prbarr)	log(prbconv)	log(prbpris)	log(avgse)	log(polpc)
##	1.812595	2.069457	1.536031	1.683955	2.538936

```
## log(density)    log(taxpc)        west        central        urban
##      3.126722      2.087337      4.293576      1.921333      2.073034
## log(pctmin80)      wcon          wtuc          wtrd          wfir
##      3.798204      2.214858      1.798198      3.233880      2.841757
##           wser          wmfq          wfed          wsta          wloc
##      1.367321      2.047091      3.344744      1.797655      2.276723
##      log(mix)    log(pctymle)
##      2.251470      1.708311
```

Computing VIF indicates that there is no multicollinearity.

Assumption - Exogeneity (Zero Conditional Mean)

From the residuals vs fitted plot, the red line is very influenced by the outliers on either end, but the zero conditional mean appears to hold.

Additional Assumptions - Homoscedasticity

From the same residuals Vs fitted plot, we see it is scattered with extreme outliers. So, it is not easy to determine Homoscedasticity from this plot only.

```
bptest(model4)
```

```
##
## studentized Breusch-Pagan test
##
## data:  model4
## BP = 50.118, df = 22, p-value = 0.0005653
```

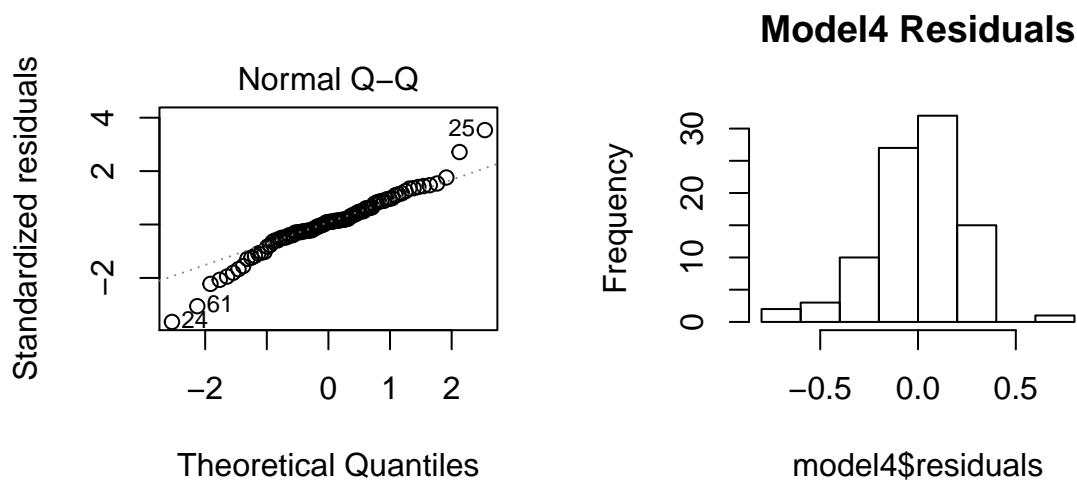
```
ncvTest(model4)
```

```
## Non-constant Variance Score Test
## Variance formula: ~ fitted.values
## Chisquare = 0.1546225, Df = 1, p = 0.69416
```

Both tests are showing large p-values showing that we fail to accept the null hypothesis.

Normality of Residuals

```
plot(model4, which=2)
hist(model4$residuals, main="Model4 Residuals")
```



$\sim \log(\text{prbarr}) + \log(\text{prbconv}) + \log(\text{prbpris})$

than a few outliers, the distribution is relatively normal for our given sample size.

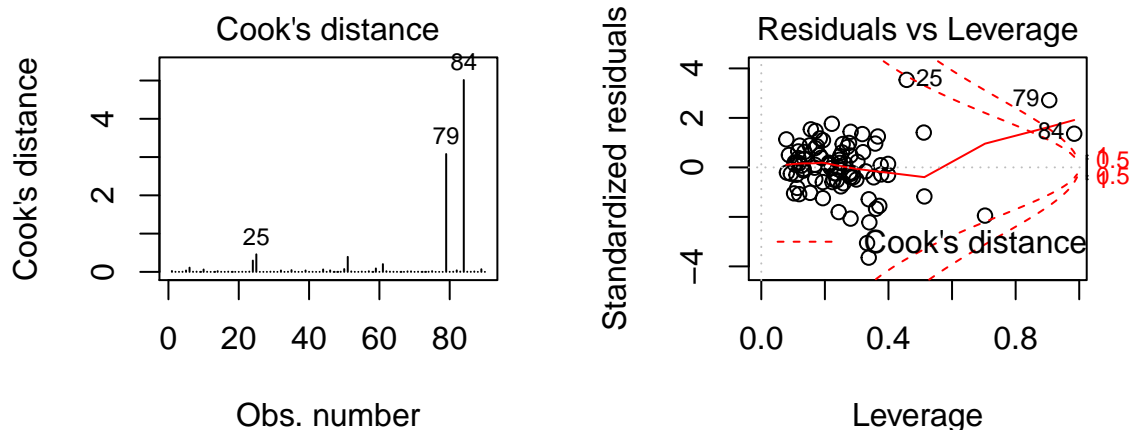
Other

Cook's Distance:

```
plot(model4, which=4)
plot(model4, which=5)
```

```
## Warning in sqrt(crit * p * (1 - hh)/hh): NaNs produced
```

```
## Warning in sqrt(crit * p * (1 - hh)/hh): NaNs produced
```



$\sim \log(\text{prbarr}) + \log(\text{prbconv}) + \log(\text{prbpris}) \sim \log(\text{prbarr}) + \log(\text{prbconv}) + \log(\text{prbpris})$

There are some influential values [25,79,84] however cook's distance is within bounds. ##### Calculating AIC

```
AIC(model4)
```

```
## [1] 31.94569
```

The AIC for this model is 31.94569

```
stargazer(model1, model2,model3, model4, type = "latex",
  report = "vc", # Don't report errors, since we haven't covered them
  title = "Linear Models Predicting Crime Rate",
  keep.stat = c("rsq", "n"),
  omit.table.layout = "n",
  column.labels=c("Not good","Good","Better","Using All"),
  dep.var.caption = "Measuring Crime Rate",
  dep.var.labels = "Crime Rate")
```

% Table created by stargazer v.5.2.2 by Marek Hlavac, Harvard University. E-mail: hlavac at fas.harvard.edu
 % Date and time: Sat, Dec 08, 2018 - 3:13:17 PM

Model1 -

Only 41.6% of crime rate is being explained by our model. prbarr, prbconv are shown to have a negative effect on crime, while the other independent variables prbpris and avgsen have a positive effect on crime.

We do not want to overinterpret the model as assumptions of homoskedasticity appear to be violated. If we accept the model at face value, we would see that the estimators for the probability of arrest and probability of conviction are statistically significantly different from zero.

Model 2 -

78.4% of crime rate is being explained by this model. Polpc, Density and Pctmin80 appear to have a positive effect on crime. While other independent variables are estimated as having a negative effect on crime.

Model 3 -

78.3% of crime rate is being explained by this model.

This says that crime rate decreases if more people are arrested, convicted, sentenced in prison

Model 4 -

83.5% of crime rate is being explained by this model Using all variables, we got a better AIC compared to Model 1 and very close to Model 2- but it is not a best fit model

Best-Fit Model

Based on our analysis - we found that our Model3 is the best fit model with low AIC value 26.9529.

Omitted Variables

In order to make valid policy recommendations, we need confidence that our estimated coefficients for policy-relevant variables are unbiased, statistically significant, and practically significant. Statistical software makes it quite easy to determine if there is a relationship between a given variable and the dependent variable that is statistically significantly different from zero - an area of analysis that we will expand upon in follow-ups to this piece. Practical significance of our estimates requires just one extra step to interpret the meaning of the estimate for each variable under consideration. Accounting for elements which could bias our estimates is more difficult and, to some degree, not a solvable problem.

We only have observational data available. Moreover, we are not able to design or even infer experiments for our data generating process. As such, we are left to reason about counterfactuals, rather than conduct

Table 2: Linear Models Predicting Crime Rate

	Measuring Crime Rate			
	Not good	Crime Rate		Using All
		Good	Better	
	(1)	(2)	(3)	(4)
log(prbarr)	-0.724	-0.522	-0.506	-0.533
log(prbconv)	-0.473	-0.403	-0.391	-0.298
log(prbpris)	0.160	-0.332	-0.320	-0.301
log(avgsen)	0.076	-0.235	-0.232	-0.308
log(polpc)		0.553	0.527	0.480
log(taxpc)		-0.068		0.018
west				0.052
central				-0.077
urban				0.071
log(density)		0.164	0.162	0.139
log(pctymle)		-0.066		0.101
log(pctmin80)		0.266	0.261	0.255
wcon				0.001
wtuc				0.0001
wtrd				0.001
wfir				-0.001
wser				-0.0004
wmfg				-0.0001
wfed				0.001
wsta				-0.00003
wloc				0.0001
log(mix)				0.062
Constant	-4.868	-1.417	-1.615	-1.984
Observations	90	90	90	90
R ²	0.416	0.784	0.783	0.835

experiments to verify the implications of our model. Additionally, we have a flawed data collection process, which we also have no ability to correct for. Our desired population variables are by-in-large not included in the dataset we were provided. Some of these desired variables are practically or ethically unobservable. Others were operationalized in a flawed manner, with a negative impact on our ability to model relationships with a causal interpretation. We address some of these issues here.

Our ideal model of the causes of the crime rate would be something like:

$$\begin{aligned} \text{crime_rate} = & \beta_0 + \beta_1 \text{crtty_punish} + \beta_2 \text{svrty_punish} + \beta_3 \text{poverty_rate} + \\ & \beta_4 \text{educ} + \beta_5 \text{social_cohesion} + \beta_6 \text{weapon_availability} + \beta_7 \text{real_wage} + \\ & \beta_8 \text{low_skill_unemployment_rate} + \beta_9 \text{age_15_to_30_proportion_population} + \\ & \beta_{10} \text{percent_of_population_previously_committed_crime} + \\ & \beta_{11} \text{percent_of_population_previously_imprisoned} + \dots + \text{error} \end{aligned}$$

Unfortunately, we are unable to observe virtually all of these concepts.

Some concepts have been operationalized in our dataset. For example, certainty of punishment has been operationalized through three variables: 1) the percent of the population which are police, 2) the proportion of arrests to crimes, and 3) the proportion of convictions to arrest. This is among the most effective operationalizations in this dataset. Severity of punishment is also operationalized through 1) the proportion of convictions that result in a prison sentence and 2) the average length of a prison sentence. Nominal wages are operationalized in the dataset with average wages for certain industry groupings. None of poverty rate, education, social cohesion, weapon availability, cost of living, or the low skill unemployment rate are operationalized within this dataset.

Moreover, certain variables that are included in our dataset are likely correlated with many of our desired variables, but actually measure something distinct - introducing the possibility for model estimates based on those variables to be biased and thus misleading. For example, the `pctmin80` variable measures the percent of a county that was minority in 1980 - 7 years prior to our other observations. Setting the time divergence aside and extrapolating from national trends in the U.S. in the 1980s, the percentage of a county that belongs to a minority class could be the result of *red lining*, a segregating practice which led to poverty and low-quality education, both positively correlated with crime rate. It may also exhibit a parabolic relation with social cohesion. If we were to include `pctmin80` in our regression, we would expect the model estimate to be biased as we have not adjusted for the impacts of education, poverty, or social cohesion. Examining the impact of education alone on the estimator for `pctmin80` - as education was likely negatively correlated with `pctmin80`, and we expect educated to be negatively related to the crime rate, the model's estimate of the impact of the percent of a county which was minority in 1980 would be upwardly biased. In other words, the estimator for `pctmin80` in the underspecified model would imply a much larger relationship between `pctmin80` and crime rate than actually exists.

Similarly, our dataset contains a variable `density` which is likely correlated with two of our desired but unobserved explanatory variables: social cohesion and poverty. In practice, in the U.S. in the 1980s, we would expect social cohesion to be negatively correlated with density, while poverty may be positively correlated with density. We expect the beta for social cohesion to crime rate to be negative, while the beta for poverty to crime rate is expected to be positive. The impact of both of these omitted variables is that the model's estimate for density is likely upwardly biased. As with `pctmin80`, the model would again overestimate the impact of density on crime rate.

Our ability to interpret the variable `polpc` in our dataset is also compromised by omitted variable bias. While we understand the idea that increased police presence should increase the certainty of punishment (more likely to be detected and more likely to be caught) *ceteris paribus*, in our current dataset, we do not have the ability to use `polpc` in this way. We are unable to observe the counterfactual of the same location with the same characteristics at the same point in time having more or less police. Rather, the variable in our dataset is the current level of police as a percent of the population. Given that we expect local governments to respond to increased crime by highering more police, our model is more likely to reflect that higher crime rate locations

also have higher police concentrations. Given an alternate work environment where we could retrieve more data, we might think about attempting to compensate for this by locating police concentration and crime rate statistics for previous years, then using them to create variables for the percentage point change in police concentration, which we could use to explain a newly created variable for the percentage point change in crime rate for a given location. However, in their current single point in time forms, our model is likely to estimate the relationship between police percentage and crime rate as positive, thus providing a misleading estimate for the relationship we would actually like to observe.

Finally, our dataset contains several variables with nominal wages for certain industries. Including these in our model is likely to be somewhat misleading, producing biased estimators because these measures are not adjusted for cost of living. Said in other terms, each of the nominal wage indicators is likely positively correlated with our desired explanatory variable - real wages. Conceptually, we expect the relationship between real wages and crime rate to be negative, while the relationship between real wages and nominal wages is positive. As such our model's estimator for wages is likely to understate the impact of wages on crime rate. As such, these nominal wage variables are an imperfect proxy for the desired variable real wages

Conclusion

We examined several models of crime rate and found a directionally consistent, statistically significant negative relationship for the probability of arrest and the probability of conviction on crime rate. As such, policies adopted should focus on increasing the certainty of punishment for committing crimes. One such policy could focus on improving information flow from local communities to police and judicial officials. A good model to build off of is community policing, where police focus on developing ties to the local community to build trust and thereby promote flow of needed information.

That said, our ability to draw policy prescriptions from our models is limited due to notable omitted variable bias, which leads our model's estimators to be biased. These omitted variable biases are not possible to overcome while limited to the current data collection process. Should more work requiring causal inference be desired on these relationships in the future, we would seek input into the data collecting process in order to correct for some of our omitted variable biases.

Appendix A: Codebook