

Lab 3: Reducing Crime (DRAFT)

W203 Statistics

Luke Evans, Daniel Rasband, and Yulia Zamriy

April 1, 2018

An analysis of Crime in North Carolina to support Policy Decisions

Abstract...

TBA

Introduction

Crime is expected to be a significant issue during the upcoming election in North Carolina. Using statistical techniques, this report intends to provide data driven insights into the determinants of crime in the state. A mixture of both long- and short-term policy suggestions will be included to address the factors that exacerbate crime, and to capitalize on those factors that may inhibit it.

Exploratory Data Analysis

The data utilized to conduct this statistical analysis generally relates to 1987, with a single variable from 1980 (percent minority). Data is provided for most counties in North Carolina, and can be further grouped by region (West, Central and Other). Granularity below the county level is not available.

While averages and rates are presented for many variables, the absolute numbers, for example of population, are not. This can generate some challenges when discussing practical significance.

Data Cleaning

Our initial exploration of the data has revealed several notable features. The information below provides the dimensions of the raw data: 25 variables and 97 rows.

```
setwd("/home/yulia/Documents/MIDS/mids-w203-lab3/")
crime <- read.csv("crime_v2.csv", stringsAsFactors = FALSE)
dim(crime)
```

```
## [1] 97 25
```

There are a total of 91 observations in the dataset; 6 rows are completely devoid of data and can be excluded. It should be noted that there are 100 counties in North Carolina; therefore this dataset contains data for 91% of them. It is not possible to tell if the excluded counties are randomly excluded or share specific features that may bias this data set.

Counties range in population numbers from 15,000 people to over 1 million. The data provided has many ratios and averages, but without the actual numbers relating to those numbers, it can be hard to draw practical significance from conclusions as each county will be considered equal to any other. As electoral representation in general does not follow population density, there may be advantages to analyzing data at a county level only, but it should be considered depending on the inference that is being generated.

```
crime <- na.omit(crime)
```

From the summary, the probability of conviction dimension, prbconv, is of a data format which is not consistent with analysis in R. This can be corrected, and the result stored in a new column

```
crime$prbconv <- as.numeric(as.character(crime$prbconv))
```

The variables are made up of ratios, specifically the probability of arrest, conviction and prison sentence, the percent minority, young male, police and tax revenue per capita, and the ratio (mix) of face to face crimes to other types. The mean of several variables are provided for each county: a series of weekly wages in different business segments, and prison sentences in days. Finally, an indicator of the location of the county in the state is also provided indicating west, central, where other can be identified by difference. The `urban` variable also indicates whether the county is a "Standard Metropolitan Statistical Area." Below is a summary of variables including some summary statistics:

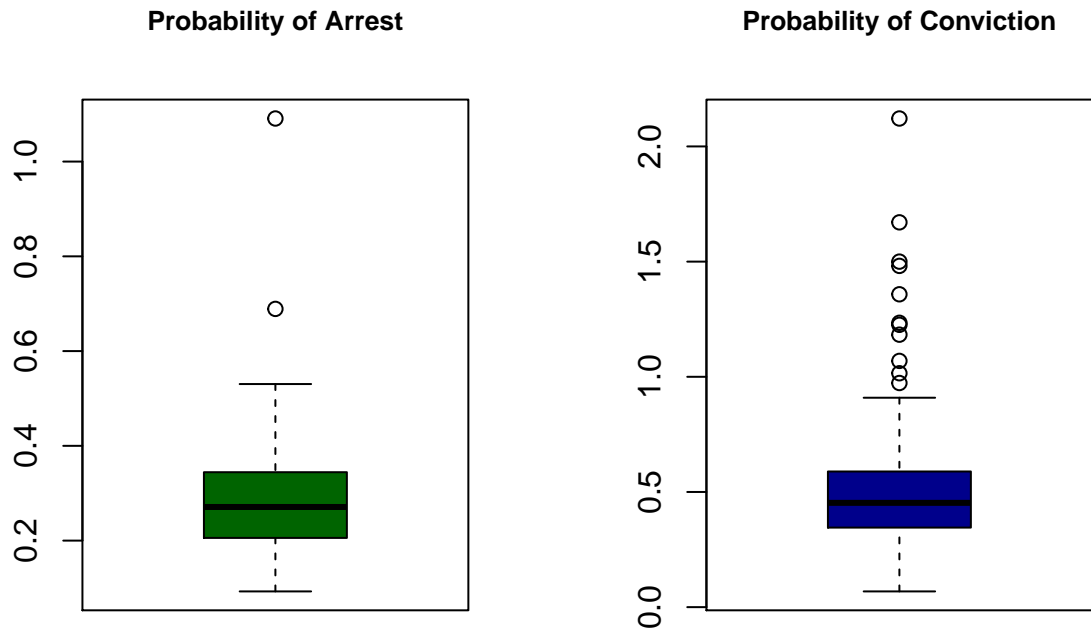
```
crime_summary <- data.frame(t(mapply(summary, crime)))
crime_summary <- crime_summary[,c("Min.", "Mean", "Max.")]
crime_summary$Min. <- round(crime_summary$Min., 5)
crime_summary$Mean <- round(crime_summary$Mean, 4)
crime_summary$Max. <- round(crime_summary$Max., 4)
kable(crime_summary, booktabs = TRUE) %>%
  kable_styling(font_size = 7)
```

	Min.	Mean	Max.
county	1.00000	101.6154	197.0000
year	87.00000	87.0000	87.0000
crmrte	0.00553	0.0334	0.0990
prbarr	0.09277	0.2949	1.0909
prbconv	0.06838	0.5513	2.1212
prbpris	0.15000	0.4108	0.6000
avgsen	5.38000	9.6468	20.7000
polpc	0.00075	0.0017	0.0091
density	0.00002	1.4288	8.8277
taxpc	25.69287	38.0551	119.7615
west	0.00000	0.2527	1.0000
central	0.00000	0.3736	1.0000
urban	0.00000	0.0879	1.0000
pctmin80	1.28365	25.4955	64.3482
wcon	193.64316	285.3585	436.7666
wtuc	187.61726	411.6680	613.2261
wtrd	154.20900	211.5529	354.6761
wfir	170.94017	322.0982	509.4655
wser	133.04306	275.5642	2177.0681
wmfg	157.41000	335.5887	646.8500
wfed	326.10001	442.9007	597.9500
wsta	258.32999	357.5220	499.5900
wloc	239.17000	312.6808	388.0900
mix	0.01961	0.1288	0.4651
pctymle	0.06216	0.0840	0.2487

From the above table it can be seen that in several counties, the probability of arrest or the probability of conviction variables are greater than one, indicating that more arrests were carried out than crimes committed, or more convictions than those arrested.

```
par(mfrow=c(1,2))
boxplot(crime$prbarr,
  col = "darkgreen",
  cex.main = 0.8,
  main = "Probability of Arrest")
```

```
boxplot(crime$prbconv,
        col = "darkblue",
        cex.main = 0.8,
        main = "Probability of Conviction")
```

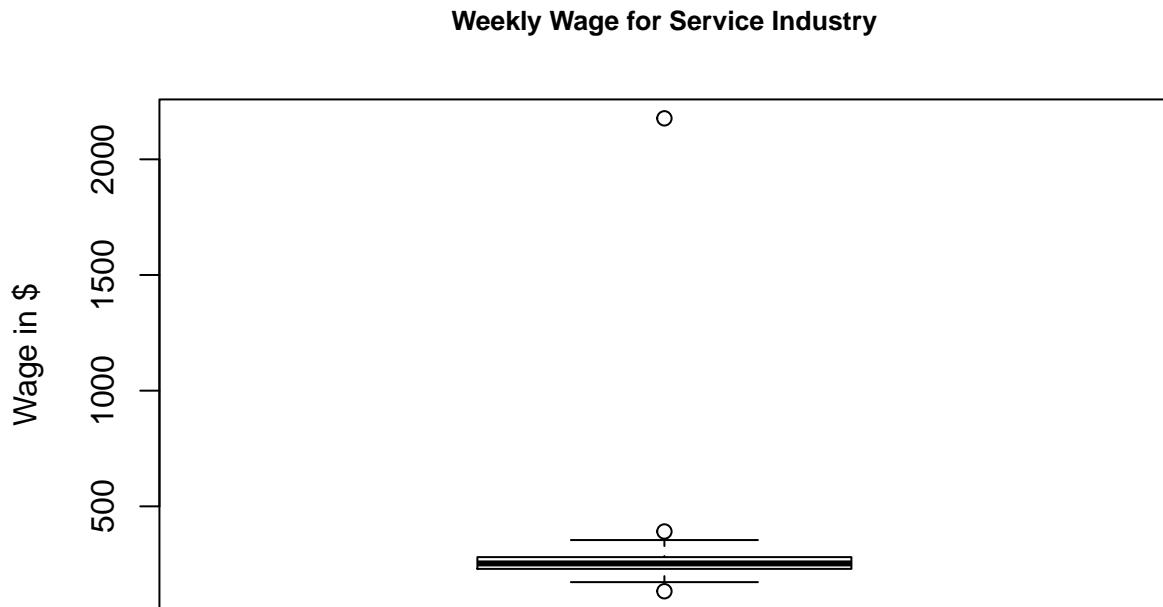


In case of probability of arrests, there is only one observation where the value is above 1, and it is significantly higher than the next closest value. Is it possible that there are more arrests than offences in one county? As this is time-limited data covering a single year, it is possible that crimes committed in the previous year and not recorded as a 1987 crime actually generated an arrest in 1987. Similarly, convictions may also have occurred in 1987, with the arrest relating to that conviction occurring in a prior period.

Additionally, the table identifies some unusual features in some of the variables, including some significant outliers. Some of these outliers clearly appear to be inconsistent with the data and will be mentioned and corrected here; others may be more subtle and will be discussed as they are considered in models.

In the series of variables noting the weekly wages in a county, there is an exceptional value in the in one of the counties, as seen in the below boxplot.

```
boxplot(crime$wser,
        cex.main = 0.8,
        main = "Weekly Wage for Service Industry",
        ylab = 'Wage in $')
```



This one value is not only over 9 standard deviations from the mean (as seen below) of wser wages, but greater than any other weekly wage value in the state.

```
(max(crime$wser) - mean(crime$wser)) / sd(crime$wser)
```

```
## [1] 9.21935
```

The observation will need to be maintained, and therefore only the service weekly wage value will be replaced by an imputed value.

After developing a predictive model using the total of average weekly wages, with which the wages of the service sector is strongly correlated, the value of \$211 per week is not dissimilar from the mean of \$254 per week and therefore use of the mean as imputed value appears reasonable. A new field is populated so that we do not lose the original values.

```
crime$wser_imp <- ifelse(crime$wser > 2000, mean(crime[crime$wser < 2000,]$wser), crime$wser)
summary(crime$wser)
```

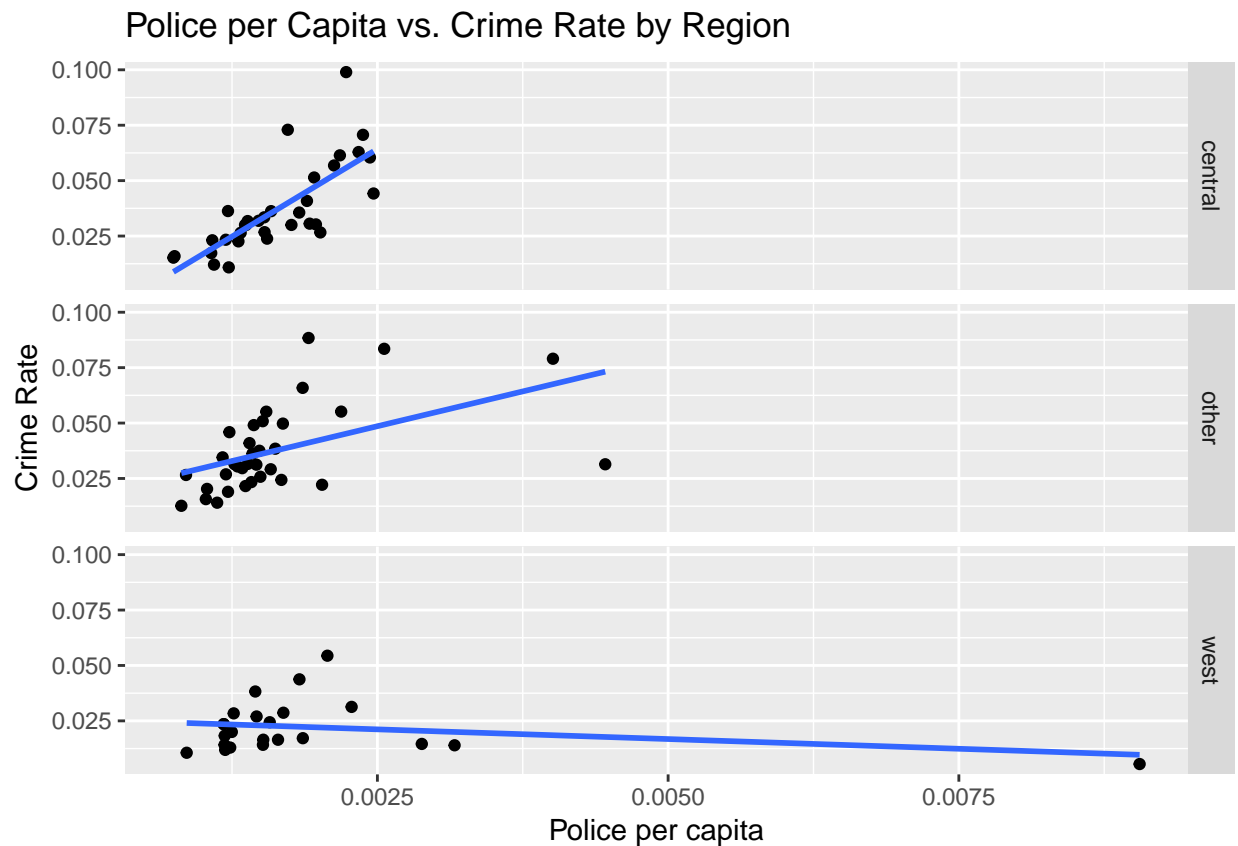
```
##      Min. 1st Qu.  Median    Mean 3rd Qu.    Max.
##  133.0   229.7   253.2   275.6   280.5  2177.1
```

```
summary(crime$wser_imp)
```

```
##      Min. 1st Qu.  Median    Mean 3rd Qu.    Max.
##  133.0   229.7   253.2   254.4   277.2   391.3
```

The variable for police per capita (polpc) also has a notable outlier. This has immediately generated some incongruent results with the rest of the dataset when segmented by region, as seen in the regression plots below.

```
crime$region <- ifelse(crime$west == 1, "west", ifelse(crime$central == 1, "central", "other"))
ggplot(crime, aes(polpc, crmrte)) +
  geom_point() +
  facet_grid(region~.) +
  geom_smooth(method = "lm", se = FALSE) +
  xlab("Police per capita") +
  ylab("Crime Rate") +
  ggtitle("Police per Capita vs. Crime Rate by Region")
```



It is clear that the impact of this observation is significant to the trend of police per capita on crime rate. Additionally, according to governing.com, police per population in Washington DC (where we would expect the highest concentration of police force) is 0.0065, significantly lower than our outlier point. Based on this analysis, we decided to recode the outlier with the mean of *polpc* in the West region:

```
crime$polpc_imp <-
  ifelse(crime$polpc == max(crime$polpc), mean(crime[crime$west == 1 & crime$polpc < 0.009,]$polpc), crime$polpc)
summary(crime$polpc)
```

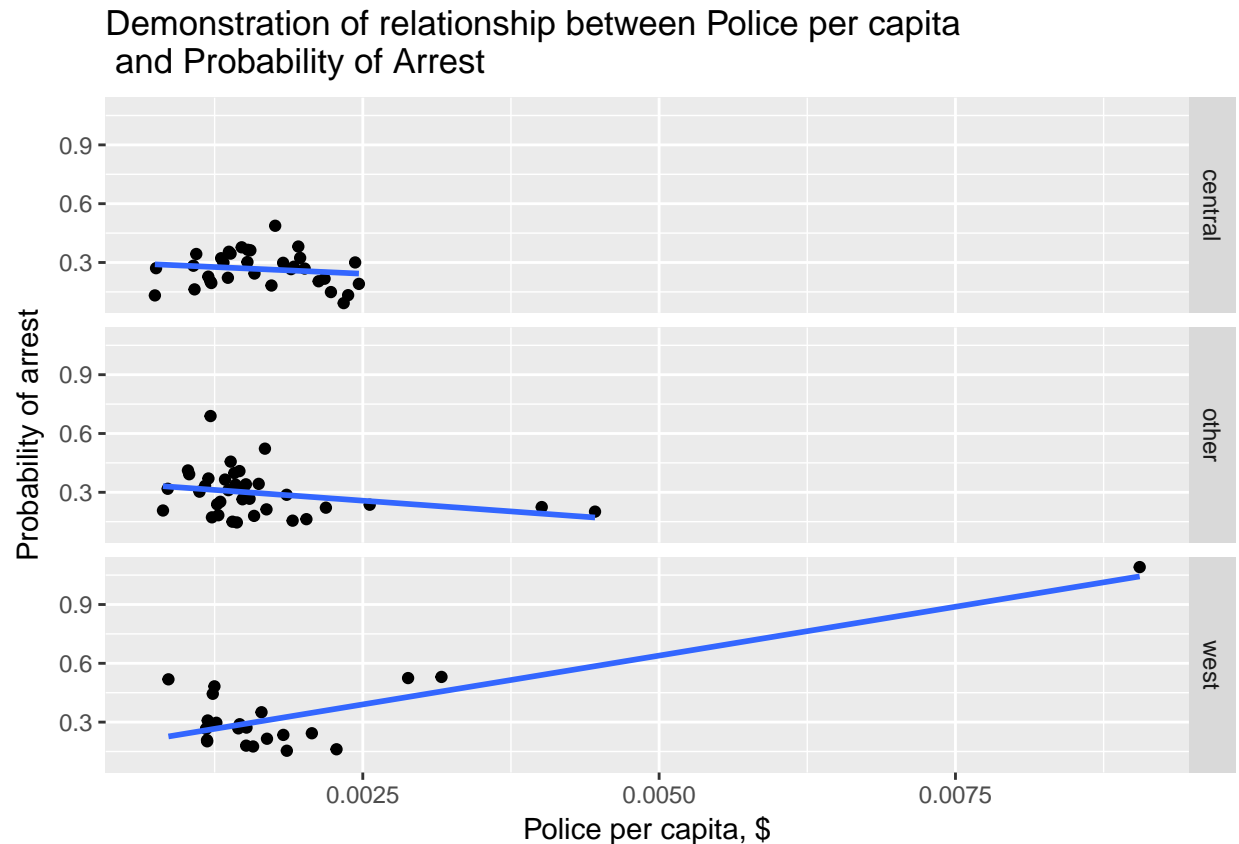
```
##      Min.   1st Qu.   Median     Mean   3rd Qu.    Max.
## 0.0007459 0.0012308 0.0014853 0.0017022 0.0018768 0.0090543
```

```
summary(crime$polpc_imp)
```

```
##      Min.   1st Qu.   Median     Mean   3rd Qu.    Max.
## 0.0007459 0.0012308 0.0014853 0.0016204 0.0018583 0.0044592
```

However, with more in-depth analysis we discovered that the outlier for *polpc* belongs to the same observation as the outlier for *prbarr*.

```
ggplot(crime, aes(polpc, prbarr)) +
  geom_point() +
  facet_grid(region~.) +
  geom_smooth(method = "lm", se = FALSE) +
  xlab("Police per capita, $") +
  ylab("Probability of arrest") +
  ggtitle("Demonstration of relationship between Police per capita \n and Probability of Arrest")
```



Moreover, the correlation between two variables changes from positive to negative if we exclude the observation with the outlier:

```
cat("Correlation with the outlier included:", cor(crime$polpc, crime$prbarr), "\n")

## Correlation with the outlier included: 0.4264409

cat("Correlation with the outlier excluded:", cor(crime[-51,]$polpc, crime[-51,]$prbarr), "\n")

## Correlation with the outlier excluded: -0.1241811
```

Hence, we decided that should create new variable for *prbarr* as well and impute mean for the probability of arrests instead of the original value of above one:

```
crime$prbarr_imp <-
  ifelse(crime$prbarr > 1, mean(crime[crime$west == 1 & crime$prbarr < 1,]$prbarr), crime$prbarr)
summary(crime$prbarr)
```

```
##      Min. 1st Qu.  Median    Mean 3rd Qu.    Max.
## 0.09277 0.20568 0.27095 0.29492 0.34438 1.09091
```

```
summary(crime$prbarr_imp)
```

```
##      Min. 1st Qu.  Median    Mean 3rd Qu.    Max.
## 0.09277 0.20568 0.27095 0.28622 0.34323 0.68902
```

Though other variables appear to have exceptional values or outliers (particularly probability of arrest and the percent young male), none are as clear. These outliers will be addressed during the development of the models as appropriate and with due consideration for the practical significance and the leverage and influence they have on the models developed.

Correlations

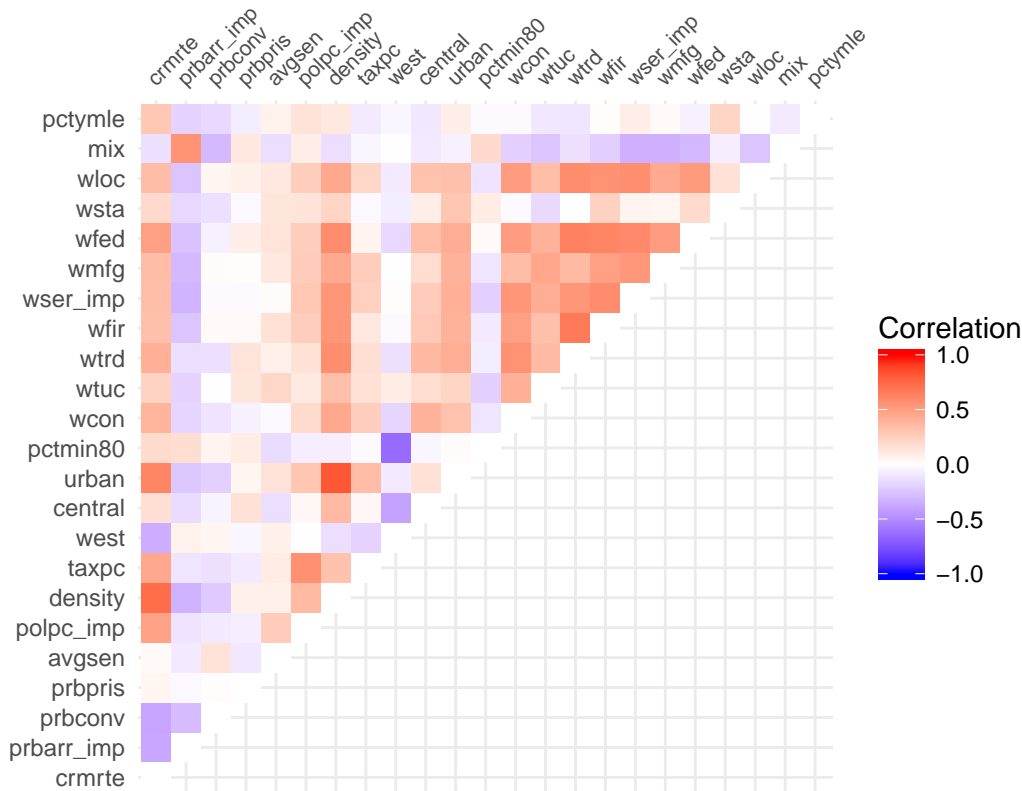
To conclude our initial data exploration, we have developed an easy-to-reference correlation heatmap for quick identification of positive or negative correlations between variables in the data set.

```
ind_variables <- c( 'crmrtte',
  'prbarr_imp', 'prbconv', 'prbpris', 'avgsgen',
  'polpc_imp', 'density', 'taxpc', 'west', 'central', 'urban', 'pctmin80', 'wcon',
  'wtuc', 'wtrd', 'wfir', 'wser_imp', 'wmfg', 'wfed', 'wsta', 'wloc', 'mix',
  'pctymle'
)

cor_mat <- round(cor(crime[,ind_variables]),2)
get_upper_tri <- function(cor_mat){
  cor_mat[lower.tri(cor_mat)]<- NA
  return(cor_mat)
}
cor_mat_upper <- get_upper_tri(cor_mat)
cor_mat_upper2 <- melt(cor_mat_upper, na.rm = TRUE)
cor_mat_upper2[cor_mat_upper2$value == 1,]$value <- 0

ggplot(data = cor_mat_upper2, aes(Var1, Var2, fill = value)) +
  geom_tile() +
  scale_fill_gradient2(low = "blue", high = "red", mid = "white",
    midpoint = 0, limit = c(-1,1), space = "Lab",
    name = "Correlation") +
  theme_minimal() +
  scale_x_discrete(position = "top") +
  theme(axis.text.x = element_text(angle = 45, vjust = 1, size = 8, hjust = 0),
    axis.title.x=element_blank(),
    axis.title.y=element_blank()) +
  coord_fixed() +
  ggtitle("Correlation Matrix")
```

Correlation Matrix



Summary of variables

The table below summarizes all variables in the dataset, and includes the expected impact of each on the dependent variable, crime rate, along with the actual correlation. Also included, as a framework for the analysis, is an assessment of the rapidity at which policy could be enacted and be effective. The support and lobbying for judges with perspectives that would support policies relating to custodial terms and the length of those could be implemented quickly. However developing incentives and strategies to reduce population density will take a longer time to generate results

```
var_labels <- c('crimes committed per person',
'probability of arrest', 'probability of conviction',
'probability of prison sentence', 'avg. sentence, days',
'police per capita', 'people per sq. mile', 'tax revenue per capita',
'=1 if in western N.C.', '=1 if in central N.C.', '=1 if in SMSA',
'perc. minority, 1980', 'weekly wage, construction',
'wkly wge, trns, util, commun', 'wkly wge, whlesle, retail trade',
'wkly wge, fin, ins, real est', 'wkly wge, service industry',
'wkly wge, manufacturing', 'wkly wge, fed employees',
'wkly wge, state employees', 'wkly wge, local gov emps',
'offense mix: face-to-face/other', 'percent young male'
)
impact <- c("Dependent",
"Negative" , "Negative", "Negative", "Negative",
"Negative", "Positive", "Negative",
"Unclear", "Unclear", "Unclear", "Unclear",
```



```

"Negative", "Negative", "Negative",
"Negative", "Negative", "Negative", "Negative",
"Negative", "Negative", "Unclear", "Positive")
control <- c("NA", "Medium Term", "Medium Term", "Short Term", "Short Term",
"Medium Term", "Long Term", "Long Term",
"No", "No", "No", "Long Term",
"Medium Term", "Medium Term", "Medium Term",
"Medium Term", "Medium Term", "Medium Term", "Medium Term",
"Short Term", "Medium Term", "No", "Long Term")
cor_w_crimerate <- round(cor(crime[,ind_variables])[1,],2)
desc <- data.frame(ind_variables, var_labels, impact, cor_w_crimerate, control,
row.names = NULL)
colnames(desc) <- c("Explanatory Variables",
"Explanation",
"Expected Impact on Crime Rate",
"Correlation w/ Crime Rate",
"Policy Impact Timeframe")

kable(desc, booktabs = TRUE, align = c("llccc")) %>%
kable_styling(latex_options = c("scale_down"),
full_width = FALSE) %>%
row_spec(0, bold = TRUE) %>%
column_spec(1, width = "8em") %>%
column_spec(3, width = "10em") %>%
column_spec(4, width = "8em") %>%
column_spec(5, width = "9em")

```

Explanatory Variables	Explanation	Expected Impact on Crime Rate	Correlation w/ Crime Rate	Policy Impact Timeframe
crmrte	crimes committed per person	Dependent	1.00	NA
prbarr_imp	probability of arrest	Negative	-0.38	Medium Term
prbconv	probability of conviction	Negative	-0.39	Medium Term
prbpris	probability of prison sentence	Negative	0.05	Short Term
avgsen	avg. sentence, days	Negative	0.03	Short Term
polpc_imp	police per capita	Negative	0.48	Medium Term
density	people per sq. mile	Positive	0.73	Long Term
taxpc	tax revenue per capita	Negative	0.45	Long Term
west	=1 if in western N.C.	Unclear	-0.35	No
central	=1 if in central N.C.	Unclear	0.17	No
urban	=1 if in SMSA	Unclear	0.62	No
pctmin80	perc. minority, 1980	Unclear	0.19	Long Term
wcon	weekly wage, construction	Negative	0.39	Medium Term
wtuc	wkly wge, trns, util, commun	Negative	0.23	Medium Term
wtrd	wkly wge, whlesle, retail trade	Negative	0.41	Medium Term
wfir	wkly wge, fin, ins, real est	Negative	0.33	Medium Term
wser_imp	wkly wge, service industry	Negative	0.34	Medium Term
wmfg	wkly wge, manufacturing	Negative	0.35	Medium Term
wfed	wkly wge, fed employees	Negative	0.49	Medium Term
wsta	wkly wge, state employees	Negative	0.20	Short Term
wloc	wkly wge, local gov emps	Negative	0.35	Medium Term
mix	offense mix: face-to-face/other	Unclear	-0.13	No
pctymle	percent young male	Positive	0.29	Long Term

The Model Building Process

Overview

As we are moving into model building section of the report, let's outline our objective: identify the impact of causal variables on crime rate to build crime-fighting policies. What are the causal variables of interest in this case? We hypothesise that in this dataset there are two variables that cause crime rate to increase/decrease: probability of arrest and probability of conviction. The third probability variable, *prbpris*, has a weak correlation with crime rate. Most likely this is due to the fact that prison sentence is far enough from the act of a crime to be ineffective in altering criminal behavior.

Our first model will be developed with these two variables along with two control variables that will help us to get unbiased estimates of our main variables of interest (explained in the appropriate section). Our second model will expand on the first one. We will add variables that help us improve the fit of the model without interacting significantly with our main causal effects. The added variables also make sense in term of interpretability. The third model will contain all provided variables (except county and year as they are constants). This model will be used to demonstrate that our model # 2 is robust. The last part of this report will focus on residuals of all three models.

Dependent variable

Our main dependent variable is crime rate (*crm rte*), which is defined as "Crimes committed per person".

After careful consideration, in order for us to understand the impact of our main causal effects (probability of arrest and probability of conviction) onto crime rate, we decided to transform our dependent variable by taking a natural log:

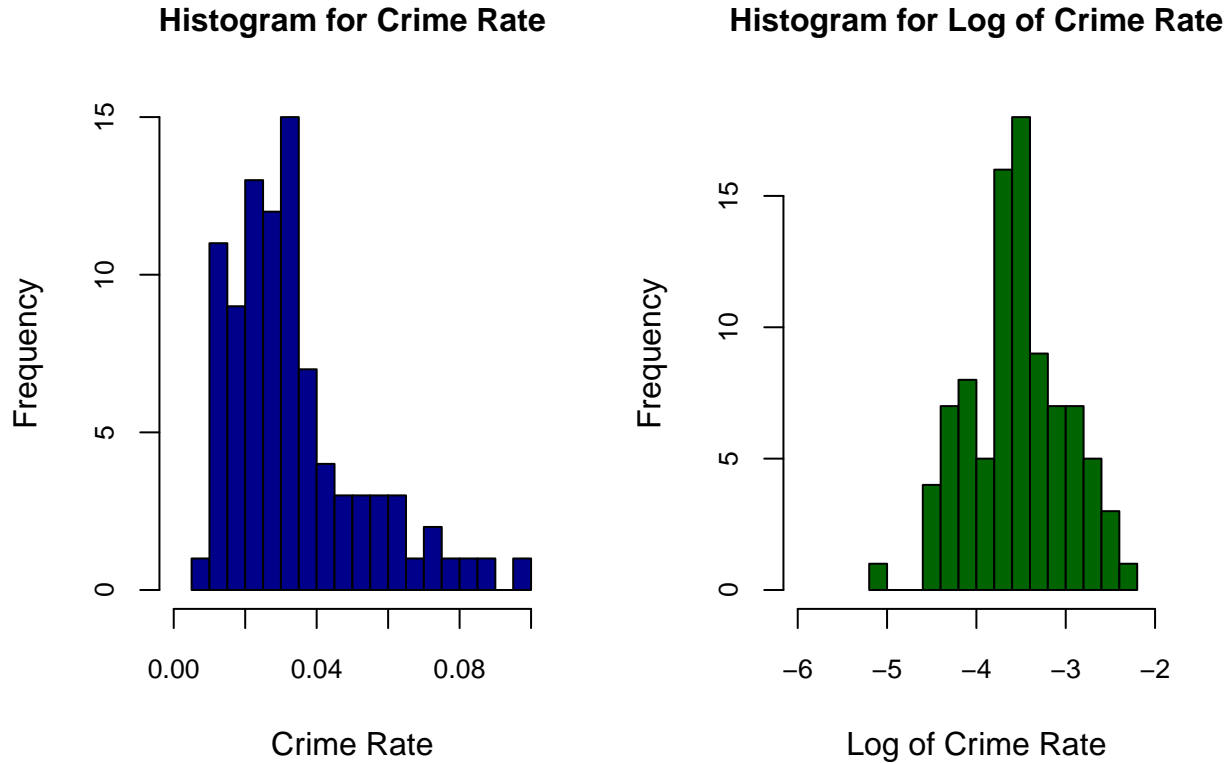
Since our variable is a ratio (crimer per person), hypothetically it can vary between 0 and 1 (though it's highly unlikely to find a county with such a high crime rate). This makes it not very suitable for OLS because this method can predict values outside 0 to 1 range. Natural log will help us only with part of the problem (avoiding negative values in prediction of actual crime rate). Caveat: in our dataset crime rate variable is never equal to zero. Hence, transformation is straight forward. However, since zero is a real possible value, we would need to watch out for those values while transforming crime rate in different datasets.

This would allow us to interpret the coefficients of our predictive factors as semi-elasticities: if probability of arrest goes up by x points, then the crime rate decreases by $100 \cdot x\%$ (assuming our stated hypothesis is true and the probability of arrest *prbarr* has a negative effect). If we were to keep the variable as is, we would interpret the coefficient for *prbarr* as: if probability of arrest goes up by x points, then the crime rate decreases by y crimes per person. However, this interpretation does not allow us to judge the practical significance of the effect (is that y big or small?).

Let's take a look at histograms for *crm rte* (as it is and transformed):

```
par(mfrow=c(1,2))
hist(crime$crm rte,
     breaks = 15,
     xlim = c(0,0.1),
     col = "darkblue",
     cex.main = 1,
     cex.axis = 0.8,
     xlab = "Crime Rate",
     main = "Histogram for Crime Rate")
hist(log(crime$crm rte),
     breaks = 15,
     xlim = c(-6,-2),
     cex.main = 1,
     cex.axis = 0.8,
```

```
xlab = "Log of Crime Rate",
col = "darkgreen",
main = "Histogram for Log of Crime Rate")
```



Based on the above charts, *crmrte* is skewed towards the right tail (a number of counties have large crime rates). The log of *crmrte*, on the other hand, looks normally distributed. This definition of the dependent variables will help us build a model with a better fit.

Main control variables

Our primary focus in this analysis is two variables: *prbarr* and *prbconv*. These two variables, the probability of arrest and the probability of conviction respectively, have relatively high correlation with crime rate and have potential to be influenced by political action. We will try to understand how probability of arrest *prbarr* and probability of conviction *prbconv* impact crime rate. If they are strong causal factors, we can define policies that influence these two factors and, hence, help us lower crime rates across North Carolina.

Earlier in this report, we hypothesised that these two variables will have negative impact on our dependent variable: the higher the probabilities of arrest and conviction, the lower the crime rate. Before building a model with these two variables, however, we want to make a case of including two more variables in our first model: *density* and *west*.

First, let's consider crime rate by region (we recoded the third region as "other" for analysis purposes):

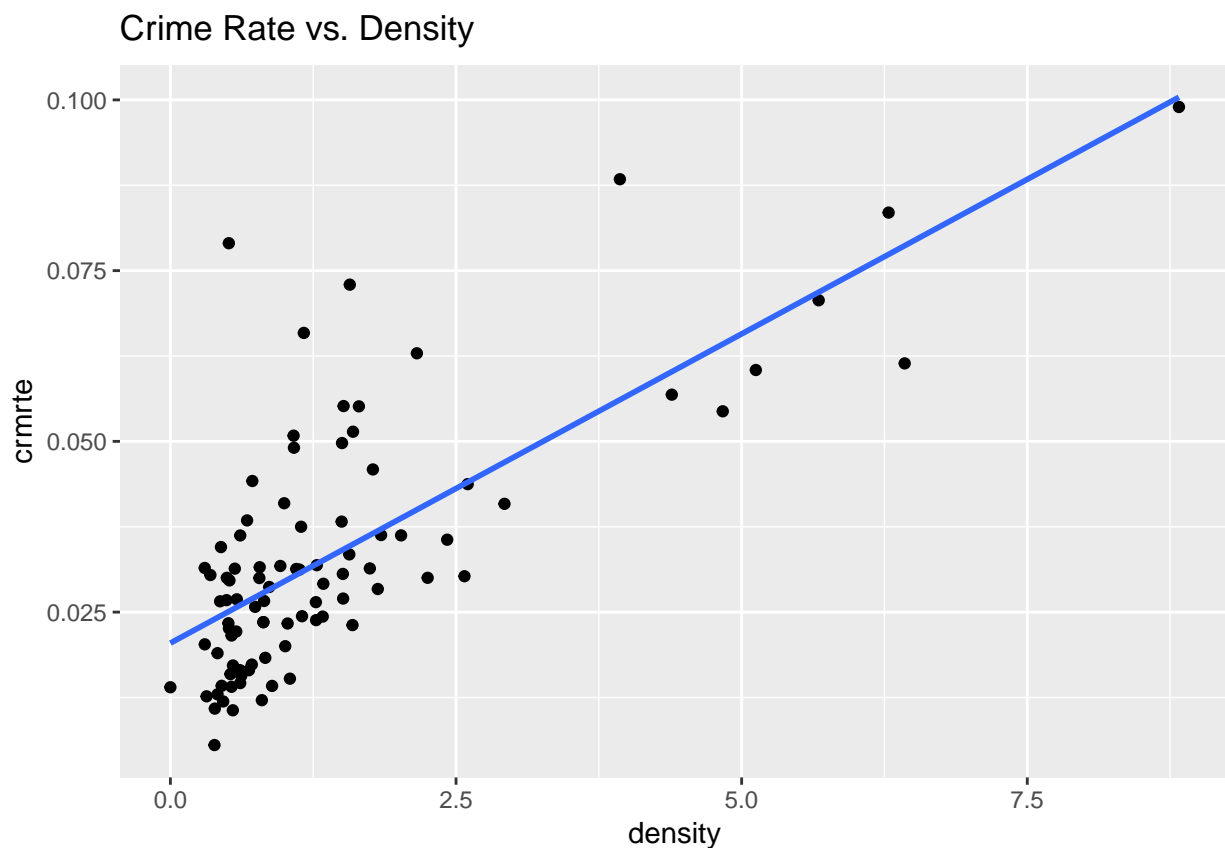
```
crime$region <- ifelse(crime$west == 1, "west",
                      ifelse(crime$central == 1, "central", "other"))
aggregate(crmrte ~ region, data = crime, mean)
```

```
##      region      crrmte
## 1 central 0.03699627
## 2  other 0.03739491
## 3   west 0.02216183
```

Based on the table above, crime rate in the West region is lower than in the Central and Other regions. We therefore need to control for regionality in order to get an unbiased read on the two selected probability variables.

On the other hand, density has the highest correlation with crime rate (0.73). And the chart below shows clear support for a strong linear relationship between the two variables:

```
ggplot(crime, aes(density, crrmte)) +
  geom_point() +
  geom_smooth(method = "lm", se = FALSE) +
  ggtitle("Crime Rate vs. Density")
```



We also know that *west* and *density* have a different relationship with *crrmte*; even though crime rate is the lowest in the West, density is the highest in the Central region. Hence, we need both *west* and *density* in our initial model to get unbiased estimates of *prbarr* and *prbconv*.

```
aggregate(density ~ region, data = crime, mean)
```

```
##      region  density
## 1 central 2.047960
## 2  other 1.085503
## 3   west 1.062994
```

Note: we tested *central* and *urban* in our models and they were not significant predictors for crime rate.

Model #1

Our first model contains four variables: *density*, *west*, *prbarr*, *prbconv*. The coefficients for these variables are calculated here:

```
model1.ind_vars <- c("density", "west", "prbarr_imp", "prbconv")
model1.formula <- as.formula(paste("log(crmrte) ~", paste(model1.ind_vars, collapse = " + "), sep = " "))
model1 <- lm(model1.formula, data = crime)
coef1 <- data.frame(model1$coefficients)
coef1$model1.coefficients <- round(coef1$model1.coefficients, 3)
colnames(coef1) <- "Model 1 Coefficients"
coef1
```

```
##           Model 1 Coefficients
## (Intercept)          -2.779
## density              0.137
## west                 -0.394
## prbarr_imp           -1.688
## prbconv              -0.686
```

Model 1 Coefficient Interpretation

- *density* 0.15: for each person per square mile increase in density, crime rate increases by 0.15% when everything else stays the same.
- *west* -0.36: crime rate in the West is 0.36% lower than in Central and Other regions on average (and controlling for all other included variables).
- *prbarr* -1.27: for each point increase in probability of arrest, crime rate decreases by 1.27%.
- *prbconv* -0.55: for each point increase in probability of arrest crime rate decreases by 0.55%

Our model is consistent with our initial hypothesis: both probability variables have a negative impact on crime rate. Moreover, a one-point change in probability of arrest has a larger impact on crime rate than one-point change in probability of conviction. This confirms our hypothesis that probability of arrest has a stronger effect on crime rate because it's closer to the act of crime (being arrested is easier to relate to than being convicted).

The adjusted R^2 for this model is 67.7%:

```
summary(model1)$adj.r.squared
```

```
## [1] 0.6716899
```

All of the coefficients are highly statistically significant when we look at heteroskedastic-robust errors:

```
coeftest(model1, vcov = vcovHC, level = 0.05)
```

```
##
## t test of coefficients:
##
##           Estimate Std. Error  t value  Pr(>|t|)
## (Intercept) -2.778965   0.202167 -13.7459 < 2.2e-16 ***
## density      0.136989   0.026966  5.0801 2.167e-06 ***
## west        -0.393927   0.074816 -5.2653 1.018e-06 ***
## prbarr_imp   -1.688143   0.372034 -4.5376 1.834e-05 ***
## prbconv      -0.685746   0.144098 -4.7589 7.789e-06 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
```

Note: we will analyze the residuals later on, after we develop all three models.

Model #2

We tested most of other variables in the dataset that potentially could be related to crime rate (using correlations). At the end we decided to add the following variables: *polpc*, *pctmin80*, and an interaction of *west* and *polpc*. Police per capita is not only highly correlated with crime rate, but it does seem to have direct link to crime. What we see from correlation analysis is counter intuitive at first glance: the higher police per capita, the higher crime rate. Logically, we would expect that increasing police presence would decrease crime rate over time. However, this dataset is panel data at one moment in time, not a time series. Hence, counties with higher police per capita require more police presence. Therefore, this variable is necessary for control purposes and it improves the fit of the model. Percent minorities also help with model fit. It is hard to hypothesize why correlation with crime rate is positive. Do poorer counties have larger minority populations? In this case, personal income would be a confounding variable that we don't have. Or counties with more minorities have more gangs (on the ethnic basis)? This would also be a confounding variable. In any case, percent minorities will be used as representative of omitted factors. Before adding these variables, though, we have decided to transform the *polpc* variable by taking its log. We perform this transformation for two reasons: there is a stronger correlation between the log of police per capita and the log of the crime rate variable than that of the raw values or of the log of the crime rate and the raw police per capita. This tells us that there is a better correlation between percent changes in the two variables than the raw changes. Performing this transformation also improves the quality of our regression model. The transformation is performed here:

```
crime$polpc_imp.ln <- log(crime$polpc_imp)
```

Correlation is strongest between the logs of both variables:

```
c(
  cor(crime$polpc_imp.ln, log(crime$crmrte)),
  cor(crime$polpc, log(crime$crmrte)),
  cor(crime$polpc, crime$crmrte)
)
```

```
## [1] 0.5099462 0.0126232 0.1698849
```

We also combine *west* and the transformed *polpc* (*polpc_imp.ln*) by multiplying them. Why do we add this interaction? Because the police per capita in the West region has much lower correlation with crime rate than police per capita in the other regions, as is shown here:

```
cat("Correlation in the West:",
    cor(log(crime[crime$region=="west",]$crmrte), crime[crime$region=="west",]$polpc_imp.ln), "\n")
```

```
## Correlation in the West: 0.1739071
```

```
cat("Correlation in the Central:",
    cor(log(crime[crime$region=="central",]$crmrte), crime[crime$region=="central",]$polpc_imp.ln), "\n")
```

```
## Correlation in the Central: 0.8012123
```

```
cat("Correlation in the Other:",
    cor(log(crime[crime$region=="other",]$crmrte), crime[crime$region=="other",]$polpc_imp.ln), "\n")
```

```
## Correlation in the Other: 0.5901029
```

Now, to the actual model:

```
model2.ind_vars <- c("density", "west", "prbarr_imp", "prbconv", "polpc_imp.ln",
                    "pctmin80", "west*polpc_imp.ln")
model2.formula <- as.formula(paste("log(crmrte) ~", paste(model2.ind_vars, collapse = " + "), sep = " "))
model2 <- lm(model2.formula, data = crime)
coef2 <- data.frame(model2$coefficients)
```

```
coef2$model2.coefficients <- round(coef2$model2.coefficients,3)
colnames(coef2) <- "Model 2 Coefficients"
coef2
```

```
##                Model 2 Coefficients
## (Intercept)                1.228
## density                    0.089
## west                      -4.243
## prbarr_imp                -2.024
## prbconv                   -0.743
## polpc_imp.ln              0.636
## pctmin80                  0.010
## west:polpc_imp.ln        -0.630
```

Model 2 Coefficient Interpretation

Some changes can be seen with the the introduction of police per capita, percent minority in 1980, and police per capita in the West region variables:

- *density* (Before: 0.15, After: 0.10) - The effect of density has decreased slightly.
- *west* (Before: -0.36, After: -3.80) - The effect of a county being in the West region has become much stronger, likely due to controlling for the effect of police per capita in the West.
- *prbarr* (Before: -1.27, After: -1.43) - The probability of arrest has a slightly stronger effect.
- *prbconv* (Before: -0.55, After: -0.74) - The effect of the probability of conviction has also increased substantially.
- *polpc_imp.ln* (0.60) - This coefficient indicates that a 60% increase in police per capita results in a 1% increase in crime rate.
- *pctmin80* (0.007) - This coefficient indicates that a 0.7% increase in the minority population means a 1% increase in crime per capita.
- *west * polpc_imp.ln* (-0.56) - A 56% decrease in police per capita in the West region results in a 1% increase in crime rate for those counties. This second model remains consistent with our initial hypothesis, and actually improves the strength of both probability variables slightly. The overall predictive strength of the model has also increased:

The adjusted R^2 for this model is 79.5%, which is 11.8ppt higher than our first model:

```
summary(model2)$adj.r.squared
```

```
## [1] 0.8009718
```

All of the coefficients are statistically significant—with the exception of the interaction between West and police per capita—when we look at heteroskedastic-robust errors:

```
coeftest(model2, vcov = vcovHC, level = 0.05)
```

```
##
## t test of coefficients:
##
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept)    1.2282470  0.8071212   1.5218 0.1318694
## density        0.0889342  0.0247597   3.5919 0.0005551 ***
## west          -4.2425746  1.2678084  -3.3464 0.0012317 **
## prbarr_imp    -2.0237368  0.3274119  -6.1810 2.275e-08 ***
## prbconv       -0.7426217  0.1185968  -6.2617 1.602e-08 ***
## polpc_imp.ln   0.6360198  0.1208091   5.2647 1.083e-06 ***
## pctmin80       0.0099089  0.0021410   4.6282 1.347e-05 ***
## west:polpc_imp.ln -0.6295556  0.1970078  -3.1956 0.0019741 **
```

```
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
```

Model #3 - All Variables

Finally, our last model includes all variables, including our imputed variables. We transform all wage variables to match the logarithmic transformation of the crime rate.

```
# Transform wage variables:
wage.vars <- c("wcon", "wtuc", "wtrd", "wfir", "wser_imp", "wmfg", "wfed", "wsta", "wloc")
wage.vars.ln <- mapply(function(var.name) paste(var.name, '.ln', sep=''), wage.vars)
crime[, wage.vars.ln] <- log(crime[, wage.vars])

# Create model 3:
model3.ind_vars <- c("prbarr_imp", "prbconv", "prbpris", "avgsen",
                    "polpc_imp.ln", "density", "taxpc", "west", "central",
                    "urban", "pctmin80", "wcon.ln", "wtuc.ln", "wtrd.ln",
                    "wfir.ln", "wser_imp.ln", "wmfg.ln", "wfed.ln", "wsta.ln",
                    "wloc.ln", "mix", "pctymle")
model3.formula <- as.formula(paste("log(crmrte) ~ ", paste(model3.ind_vars, collapse = ' + '), sep = ' '))
model3 <- lm(model3.formula, data = crime)
coef3 <- data.frame(model3$coefficients)
coef3$model3.coefficients <- round(coef3$model3.coefficients, 3)
colnames(coef3) <- "Model 3 Coefficients"
coef3
```

```
##           Model 3 Coefficients
## (Intercept)           -2.880
## prbarr_imp            -1.704
## prbconv              -0.654
## prbpris              -0.120
## avgsen               -0.024
## polpc_imp.ln          0.515
## density               0.107
## taxpc                 0.001
## west                 -0.200
## central              -0.161
## urban                -0.120
## pctmin80              0.008
## wcon.ln               0.299
## wtuc.ln               0.163
## wtrd.ln               0.252
## wfir.ln              -0.164
## wser_imp.ln          -0.526
## wmfg.ln              -0.042
## wfed.ln               0.795
## wsta.ln              -0.324
## wloc.ln               0.096
## mix                  -0.595
## pctymle               2.484
```

Despite adding a lot more variables, the R^2 of the all-inclusive regression model went up only to 82.6% (from 79.5% in model 2):


```
summary(model3)$adj.r.squared
```

```
## [1] 0.8233825
```

That means that additional variables did not contribute to the model fit.

Now let's compare all three models:

```
se.model1 <- sqrt(diag(vcovHC(model1)))
se.model2 <- sqrt(diag(vcovHC(model2)))
se.model3 <- sqrt(diag(vcovHC(model3)))
stargazer(model1, model2, model3,
  type = "text", omit.stat = "f",
  se = list(se.model1, se.model2, se.model3),
  star.cutoffs = c(0.05, 0.01, 0.001),
  font.size = "tiny",
  no.space = FALSE,
  align = TRUE)
```

```
##
## =====
##                               Dependent variable:
##                               -----
##                               log(crmrte)
##                               (1)         (2)         (3)
## -----
## density          0.137***      0.089***      0.107
##                  (0.027)      (0.025)      (0.056)
##
## taxpc                                0.001
##                                      (0.007)
##
## west          -0.394***      -4.243***      -0.200
##                  (0.075)      (1.268)      (0.116)
##
## central                                -0.161*
##                                      (0.076)
##
## urban                                -0.120
##                                      (0.224)
##
## prbarr_imp      -1.688***      -2.024***      -1.704***
##                  (0.372)      (0.327)      (0.308)
##
## prbconv         -0.686***      -0.743***      -0.654***
##                  (0.144)      (0.119)      (0.122)
##
## prbpris                                -0.120
##                                      (0.423)
##
## avgsen                                -0.024
##                                      (0.015)
##
## polpc_imp.ln          0.636***      0.515**
##                  (0.121)      (0.189)
```

```

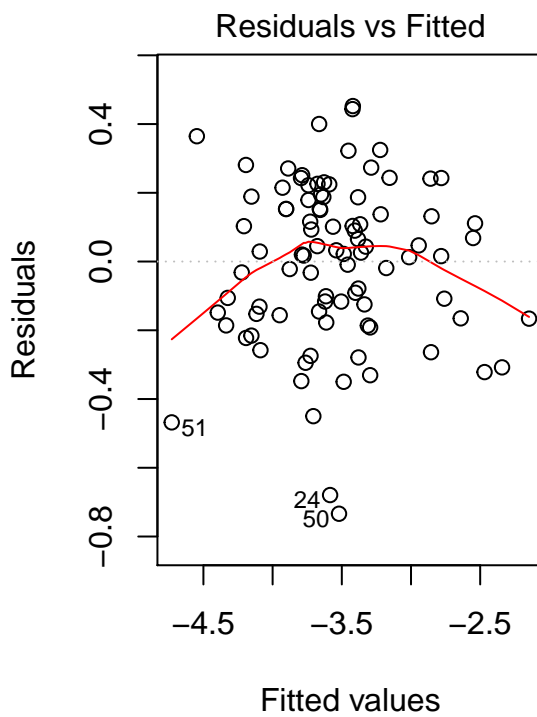
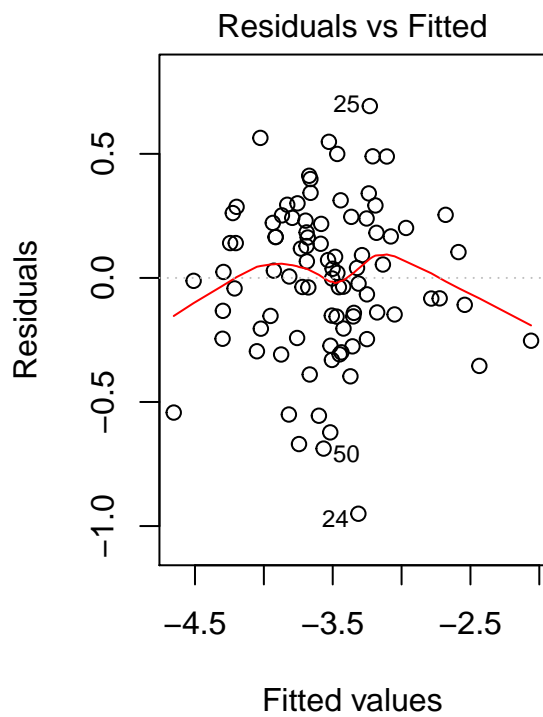
##
## pctmin80                0.010***      0.008**
##                        (0.002)      (0.003)
##
## west:polpc_imp.ln      -0.630**
##                        (0.197)
##
## wcon.ln                0.299
##                        (0.235)
##
## wtuc.ln                0.163
##                        (0.287)
##
## wtrd.ln                0.252
##                        (0.314)
##
## wfir.ln               -0.164
##                        (0.341)
##
## wser_imp.ln           -0.526
##                        (0.310)
##
## wmfg.ln               -0.042
##                        (0.164)
##
## wfed.ln                0.795
##                        (0.420)
##
## wsta.ln               -0.324
##                        (0.323)
##
## wloc.ln                0.096
##                        (0.643)
##
## mix                   -0.595
##                        (0.536)
##
## pctymle                2.484
##                        (1.325)
##
## Constant              -2.779***      1.228      -2.880
##                        (0.202)      (0.807)      (4.196)
## -----
## Observations           91           91           91
## R2                     0.686         0.816         0.867
## Adjusted R2            0.672         0.801         0.823
## Residual Std. Error 0.313 (df = 86) 0.244 (df = 83) 0.230 (df = 68)
## =====
## Note:                  *p<0.05; **p<0.01; ***p<0.001

```

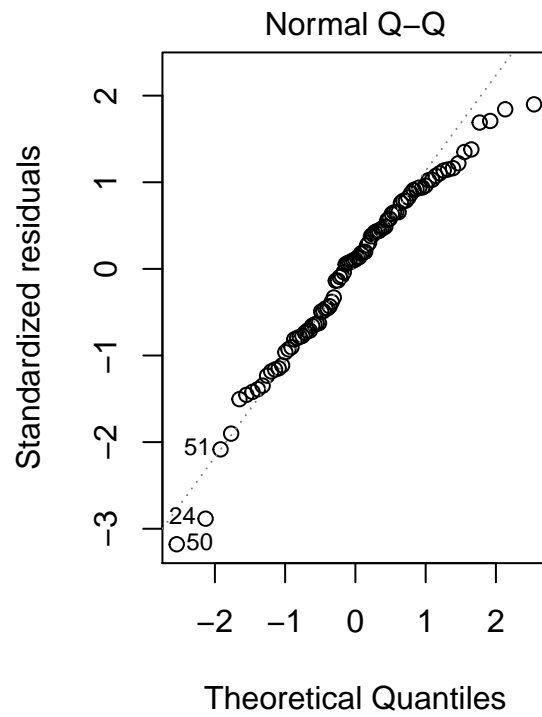
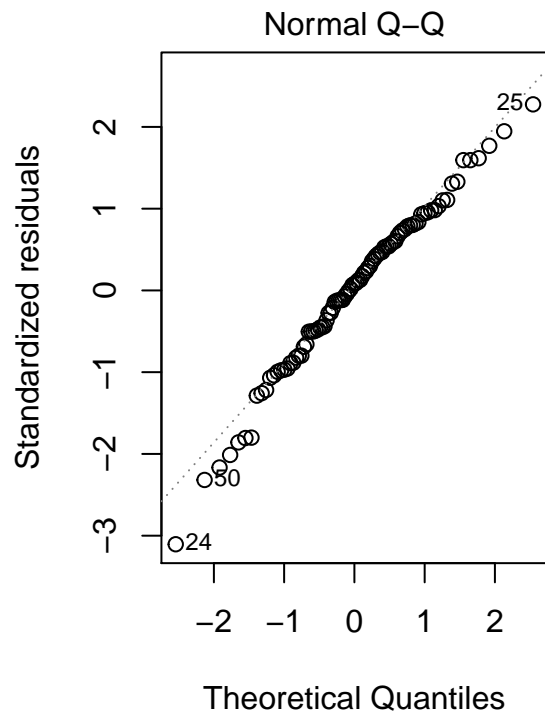
(work in progress - will need to look harder at the impact on the coefficients for the independent variables we have selected)

Residual Analysis

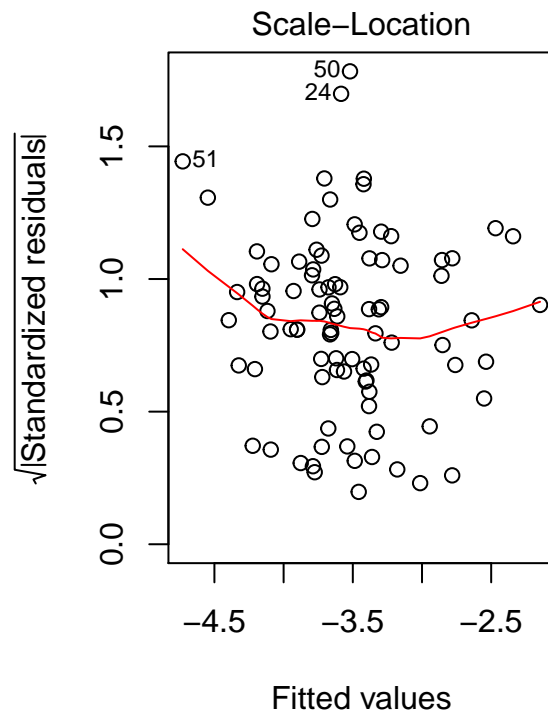
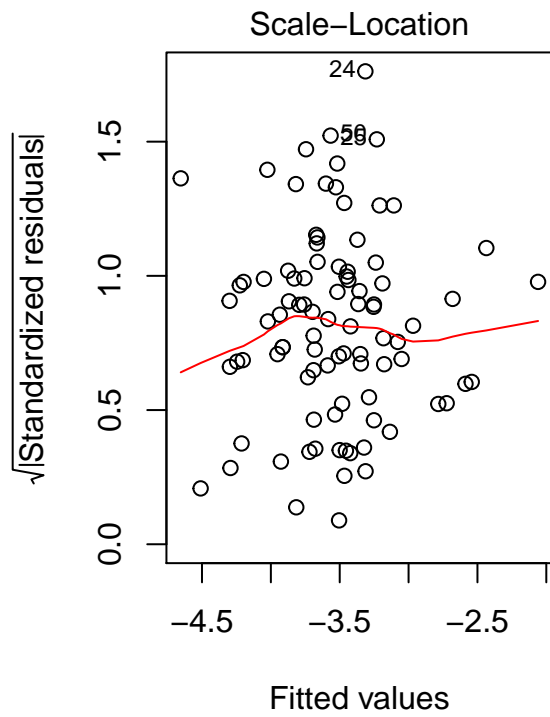
```
par(mfrow=c(1,2))
plot(model1, which = 1)
plot(model2, which = 1)
```



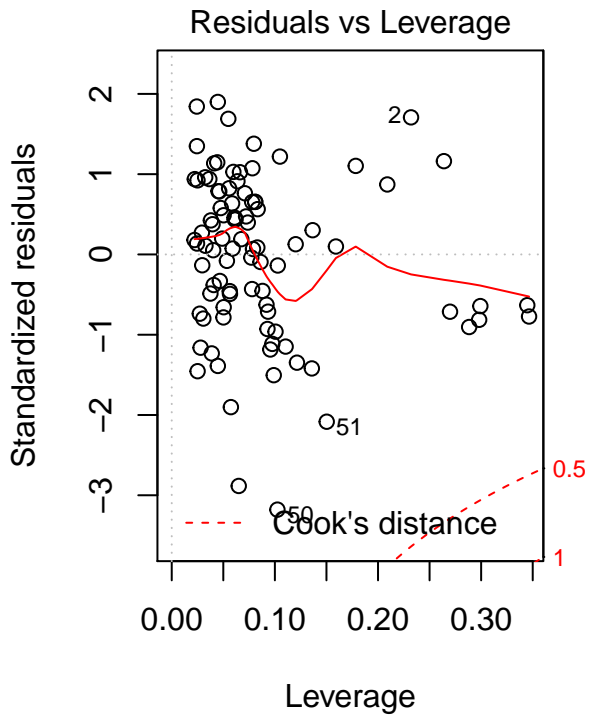
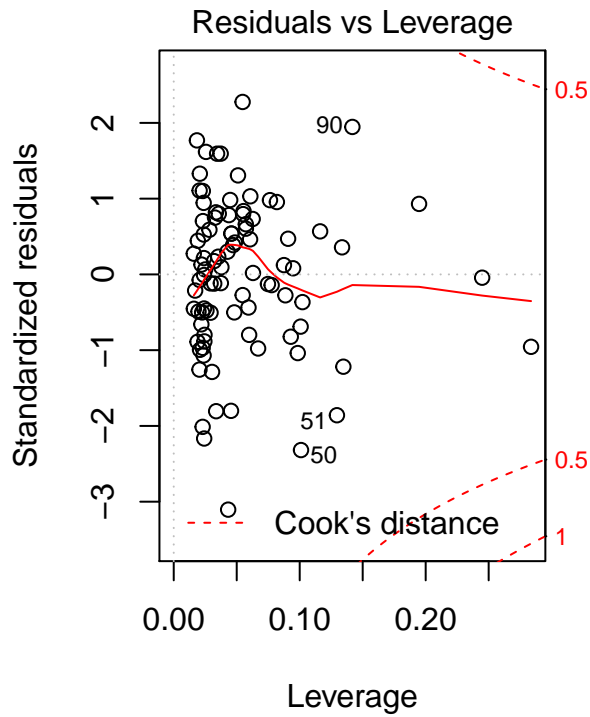
```
par(mfrow=c(1,2))
plot(model1, which = 2)
plot(model2, which = 2)
```



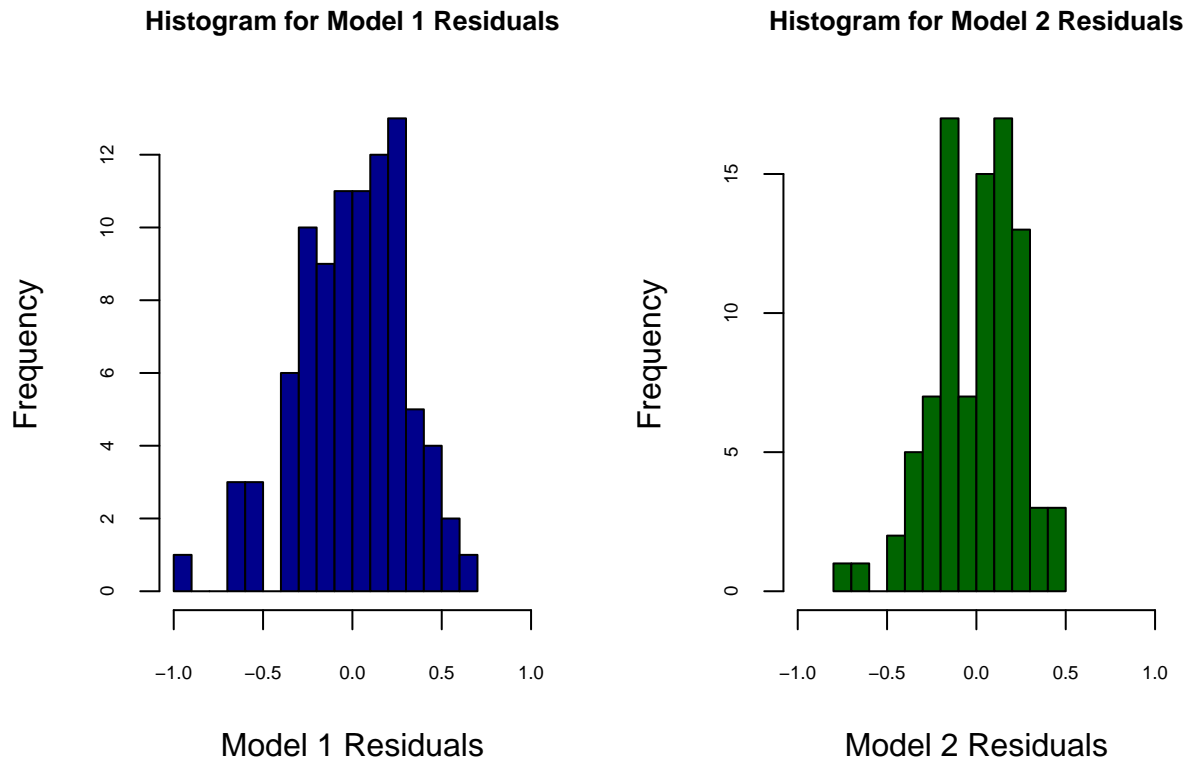
```
par(mfrow=c(1,2))
plot(model11, which = 3)
plot(model12, which = 3)
```



```
par(mfrow=c(1,2))
plot(model11, which = 5)
plot(model12, which = 5)
```



```
par(mfrow=c(1,2))
hist(model1$residuals,
      breaks = 15,
      col = "darkblue",
      xlim = c(-1,1),
      cex.main = 0.8,
      cex.axis = 0.6,
      xlab = "Model 1 Residuals",
      main = "Histogram for Model 1 Residuals")
hist(model2$residuals,
      breaks = 15,
      col = "darkgreen",
      xlim = c(-1,1),
      cex.main = 0.8,
      cex.axis = 0.6,
      xlab = "Model 2 Residuals",
      main = "Histogram for Model 2 Residuals")
```



Omitted Variable Analysis

Drug and alcohol abuse levels

The presence of drug and alcohol problems in a community is a significant contributing factor to crime rates in many areas. There is an expectation that less affluent communities in urban areas would be most impacted, which may explain some of the higher rates of crime in higher density populated counties and bias the coefficients of density and the Urban areas in models. These coefficients may have factors related to drug and alcohol abuse included in them

Recidivism

There have been several studies that suggest that someone who has committed a crime in the past is more likely to commit crimes in the *future*¹. The proportion of people with prior convictions in a county could be an additional driver that would impact crime rates. It is unclear where this may bias, but may be unlikely to be in the most affluent areas with the higher wages and taxes per capita.

¹The Offending, Crime and Justice Survey (2003); RECIDIVISM AMONG FEDERAL OFFENDERS: A COMPREHENSIVE OVERVIEW

Unemployment levels

Unemployment is likely to have an impact on crime rates. Unemployment may be higher in the young and minorities, and therefore some of the coefficients relating to those variables may be overstated.

Education levels

The level of education in a county could be an indicator of some crimes. There may be covariance between lower levels of education and lower wages

Strength of community

Strong community ties, generally in rural areas, can have a suppressing effect on crime.

Gang presence**Unreported crime**

The stigma of some crimes for victims within a community, the feeling that nothing will be done to catch the perpetrators or perhaps vigilante justice may lead to crimes in some areas being under reported. Sexual assaults specifically can be difficult for victims to report for fear of community isolation or reprisals in smaller communities. The presence of gangs, undocumented immigrants or local judicial services being overwhelmed and unavailable may be a cause in some more urban areas. The presence of unreported crimes will impact the probability of arrest as the total number of crimes that have occurred is understated