

Lab 4: Reducing Crime

w203: Statistics for Data Science

Colby Carter & Jennifer Philippou

August 22, 2017

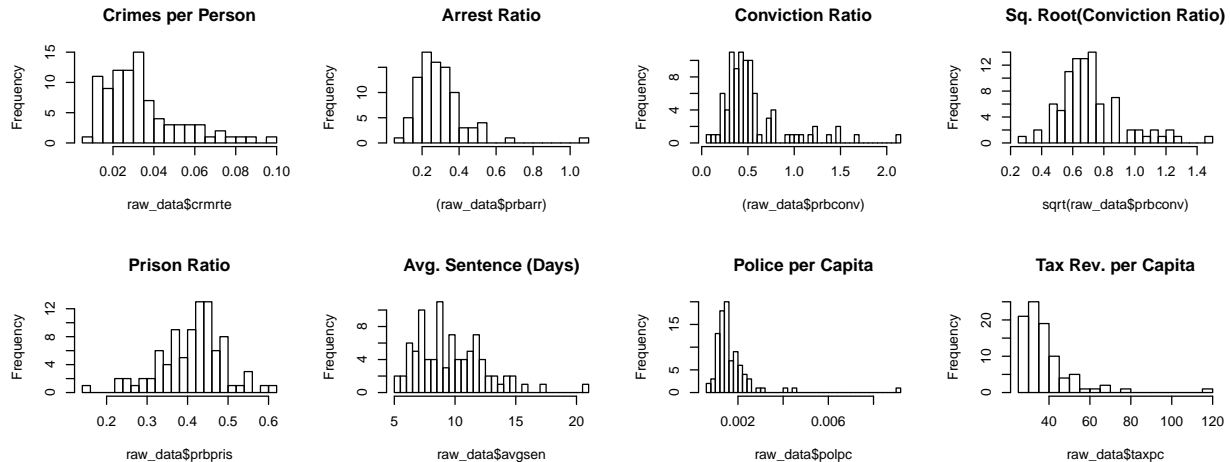
Introduction

Safety and security are fundamental human needs and society relies on the government to provide a peaceful environment. The political leaders, law enforcement agencies, and the legal system work together to bring about a halcyon community. This analysis leverages data to help local officials better understand the factors associated with crime. A stronger understanding of the current environment and relationships within the data facilitates the creation of strategic policy that will mitigate future crime and enhance the community.

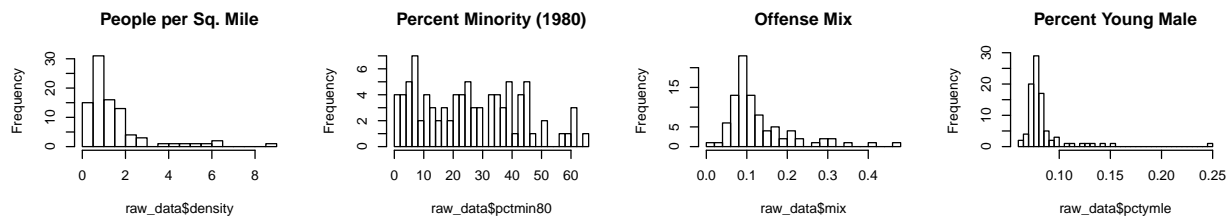
Exploratory Data Analysis

The dataset examined in this analysis include twenty-six different variables for ninety different counties. As the describe function shows, the data also has no missing values and solely represents 1987 metrics for a given point in time. The file seems to be missing a dummy variable for region so a column has been added to adjust for that. Additionally there appears to be an erroneous entry for service wage for county no. 185 as it is well above all of the other wages across the board and it has been removed. Initially the probabilities for conviction, arrest and prison that exceeded one concerned us, however further research into the dataset reveals that the metrics are actually ratios and can be included in the analysis: “PA (which is measured by the ratio of arrests to offences), probability of conviction given arrest PC (which is measured by the ratio of convictions to arrests), probability of a prison sentence given a conviction PP (measured by the proportion of total convictions resulting in prison sentences)”.

```
par(mfrow = c(2,4))
hist(raw_data$crmrte, breaks = 30, main = "Crimes per Person")
# hist(sqrt(raw_data$crmrte), breaks = 30, main = "Sq-Root of Crimes/Person")
hist((raw_data$prbarr), breaks = 30, main = "Arrest Ratio")
hist((raw_data$prbconv), breaks = 30, main = "Conviction Ratio")
hist(sqrt(raw_data$prbconv), breaks = 30, main = "Sq. Root(Conviction Ratio)")
hist(raw_data$prbpris, breaks = 30, main = "Prison Ratio")
hist(raw_data$avgsen, breaks = 30, main = "Avg. Sentence (Days)")
hist(raw_data$polpc, breaks = 30, main = "Police per Capita")
hist(raw_data$taxpc, breaks = 30, main = "Tax Rev. per Capita")
```



```
hist(raw_data$density, breaks = 30, main = "People per Sq. Mile")
hist(raw_data$pctmin80, breaks = 30, main = "Percent Minority (1980)")
hist(raw_data$mix, breaks = 30, main = "Offense Mix")
hist(raw_data$pctymle, breaks = 30, main = "Percent Young Male")
```



There are a handful of instances of very large outliers that tend to be associated with very small population densities (i.e., low population and/or high land area), which include: police per capita (0.9 compared to 0.5), tax revenue per capita (119, nearly double the next highest value), and percent young male (10 pct points higher at 25% of total population—an unusually high rate possibly explained by the presence of something like a military base). With likely explanations outside of our data, we exclude these observations from our proposed models.

```
#Remove record with top coded variable (like done in the batting averages example 13.12)
raw_data$removalFlag = ifelse(raw_data$county == "55",1, #taxpayer
                             ifelse(raw_data$county == "133",1, #pctmale
                                     ifelse(raw_data$county == "115",1, 0))) #polpc
raw_data = subset(raw_data, raw_data$removalFlag !=1)
```

For the crimes per person, probability of arrest, probability of conviction, average sentence days, offense mix, density, and tax per capita, we see varying right skewness and attempt to mitigate violations of the assumptions for linear regression by transforming these variables with the square root function (e.g., square root of the Conditional Conviction Ratio, in histogram above).

For the average weekly wage levels by industry, most distributions are fairly normal with some right skew in a couple sample distributions (e.g., see Construction and Manufacturing). We considered transformations, including taking the natural log or square root to mitigate this skew (e.g., $\text{Log}(\text{Manufacturing Wage})$ below), but given the Central Limit Theorem and our sample size above 30, we did not consider this necessary for linear regression. However, there is one case of an extreme outlier in the average wage of service workers nearly ten times higher than median county wage value, which would appear to be a data entry error likely off by an order of magnitude; we see reason to ignore this wage value error but not the observation itself and its other variables.

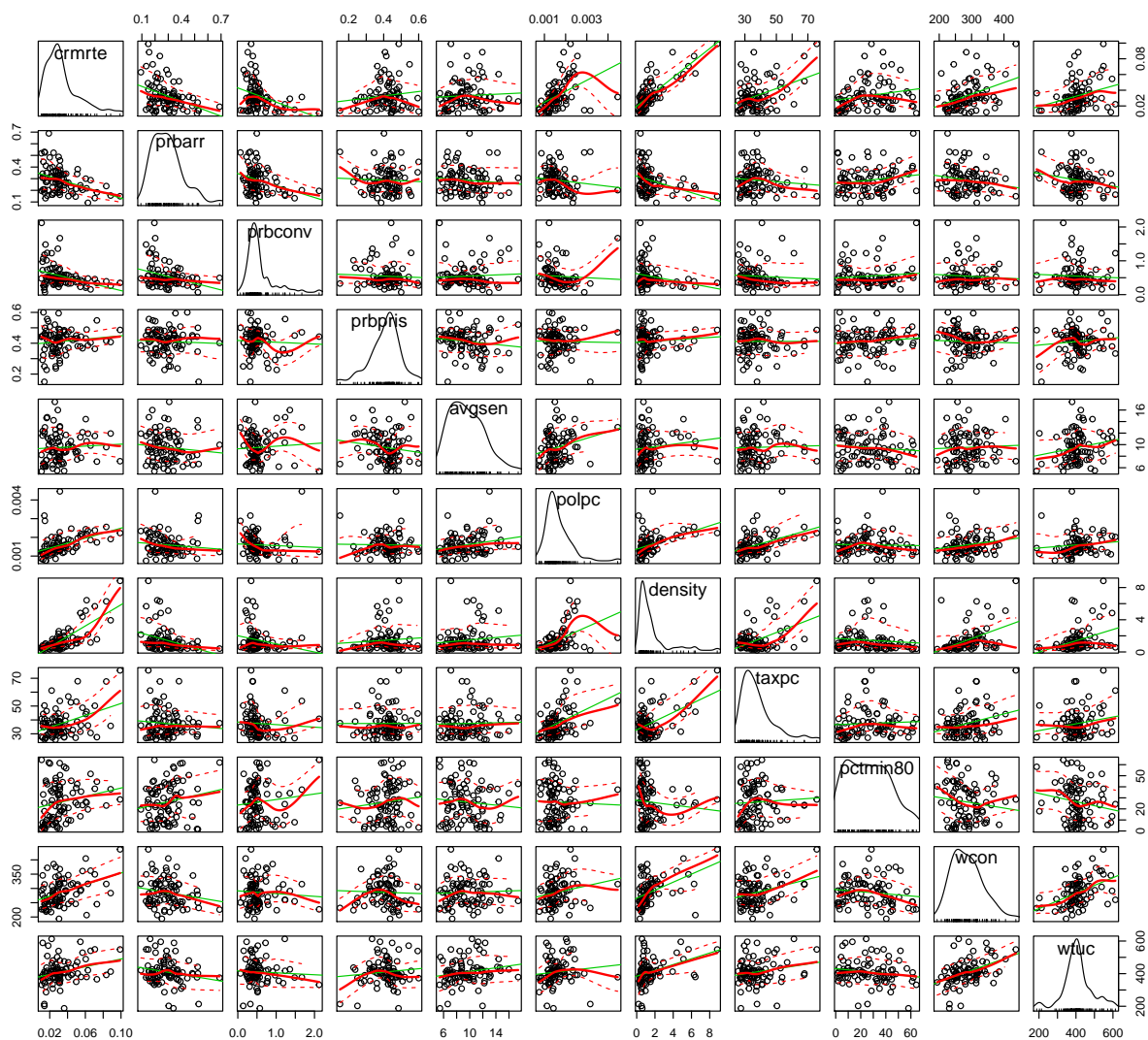
```
#Dropping service wage outlier due to likely data entry error or extremely low desity/population influence
par(mfrow = c(2,5))
hist((raw_data$wcon), breaks = 30, main = "Construction Wage")
#hist(sqrt(raw_data$wcon), breaks = 30, main = "Sq-Root of Constuction Wage")
hist((raw_data$wfed), breaks = 30, main = "Wage Federal Workers")
hist((raw_data$wfir), breaks = 30, main = "Finance/Insur/Real Estate")
hist((raw_data$wmfg), breaks = 30, main = "Manufacturing Wage")
hist(log(raw_data$wmfg), breaks = 30, main = "Log(Manufacturing Wage)")
hist((raw_data$wser), breaks = 30, main = "Wage of Service Workers")
hist((raw_data$wloc), breaks = 30, main = "Wage Local Gov.")
hist((raw_data$wsta), breaks = 30, main = "Wage State Workers")
hist((raw_data$wtrd), breaks = 30, main = "Wage Retail")
hist((raw_data$wtuc), breaks = 30, main = "Wage Trans/Util/Comm")
```



```
raw_data$wser[84] = NA
```

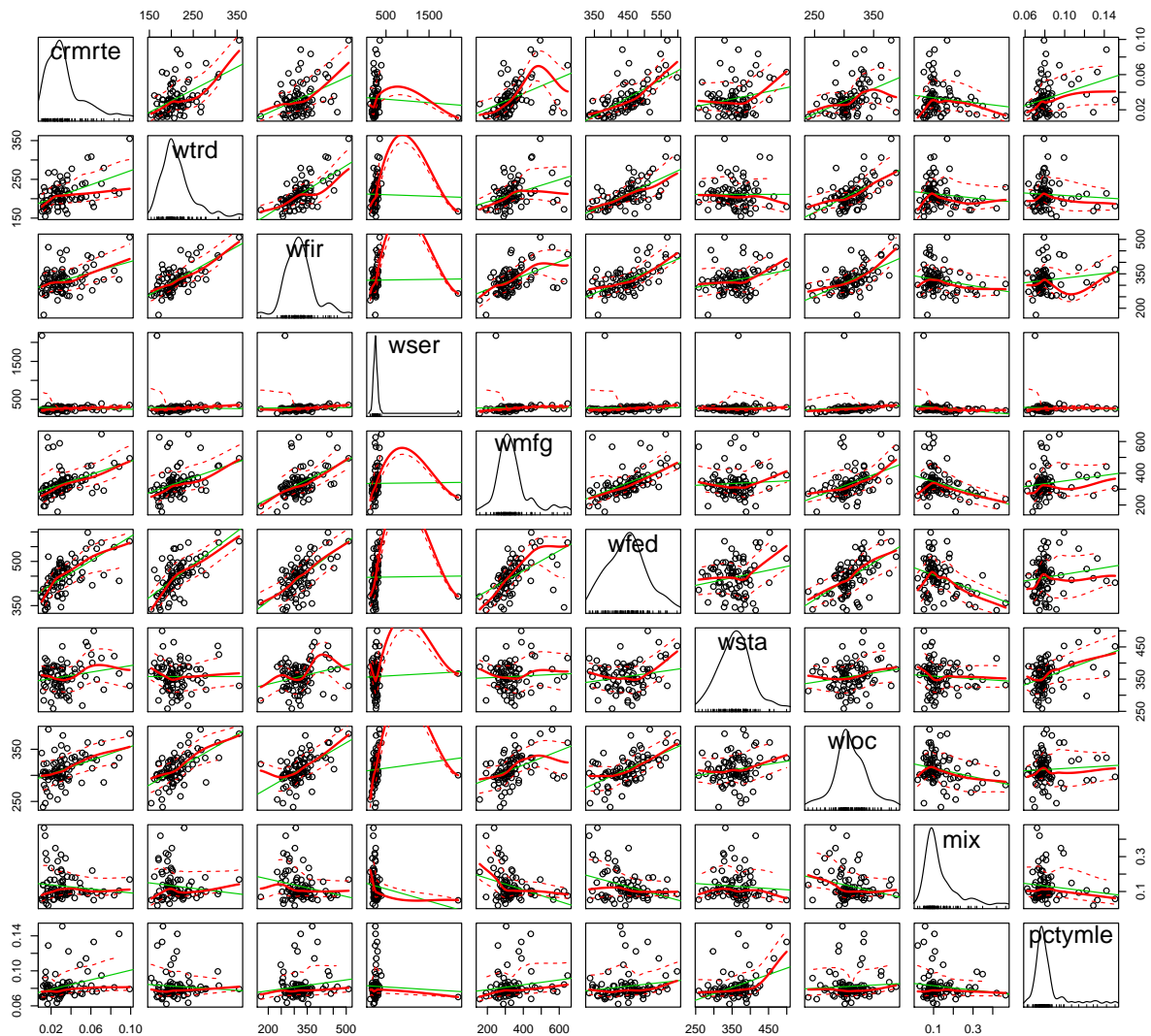
Scatter plot matrix Part I

```
scatterplotMatrix(raw_data[c(4:11,15:17)])
```



Scatterplot matrix Part II

```
scatterplotMatrix(raw_data[c(4, 18:26)])
```



Correlation between all the features:

There is a high positive correlation between the crime rate and density (0.73), and a moderate relationship between the crime rate and urban variables (0.61). There is also a moderately strong relationship between the probability of arrest and the mix (0.57). The west region and percentage minority has the lowest negative correlation (0.63), but given the binary value for west the spearman isn't super meaningful because of the number of ties in the ranking. Density is positively correlated with urban (0.8) and the federal wage (0.58). Amongst the wage values we see many moderate correlations. For example, the wage for construction has a 0.56 correlation with the wage for local gov. workers. Another example is the mild correlation between federal workers and the crime rate (0.59), density (0.58), wage for retail (0.62), wage for the service industry (0.58), wage for Finance and insurance (0.59), and the wage for local workers (0.54).

```
correlationMatrix = round(cor(raw_data[c(4:26)], use="pairwise.complete.obs"),2)
```

Correlation visual between selected variables (excluding binary variables)

```

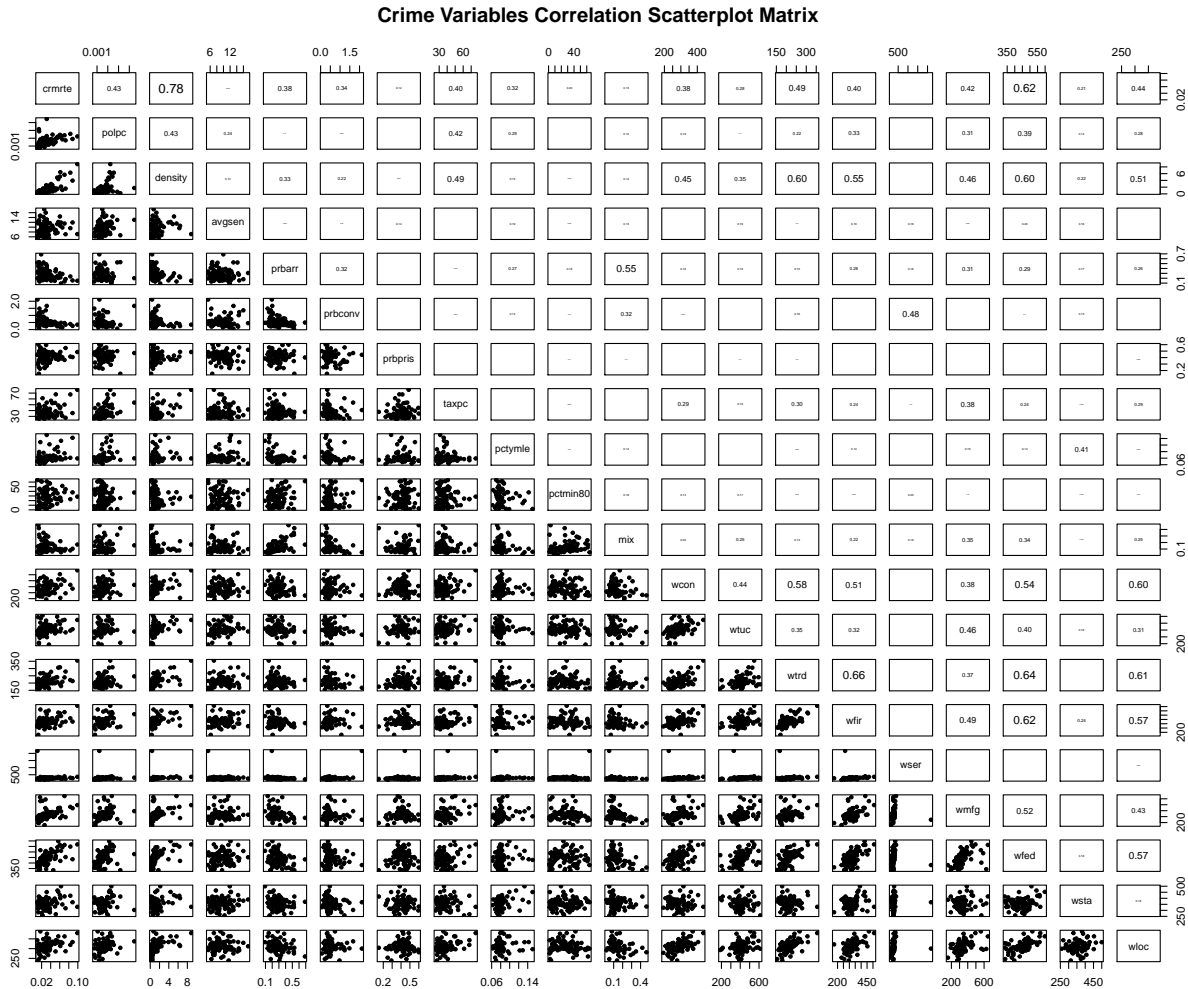
panel.cor <- function(x, y, digits = 2, prefix = "", cex.cor, ...)
{
  usr <- par("usr"); on.exit(par(usr))
  par(usr = c(0, 1, 0, 1))
  r <- abs(cor(x, y, use = "pairwise.complete.obs"))
  txt <- format(c(r, 0.123456789), digits = digits)[1]
  txt <- paste0(prefix, txt)
  if(missing(cex.cor)) cex.cor <- 0.8/strwidth(txt)
  text(0.5, 0.5, txt, cex = cex.cor * r)}

```

```

pairs(~crmte + polpc + density + avgsen + prbarr + prbconv + prbpris + taxpc + pctymle + pctmin80 + mix +

```



Identifying Key Determinants of Crime

Given our cross-section of the above county-level data, including crime statistics and likelihoods of punishment, we first seek to explain the key causal drivers of crime by hypothesizing one population model. We first conjecture which key variables may have an effect that *also* have the potential to be addressed from a policy-making position. This is the first step in an iterative process to test the robustness of these model effects and the likely biases inherent in these limited data fields and from missing, or omitted, variables that we would like to have.

Before hypothesizing our population model, we consider the types of variables in the data and which types could be translated into policy terms. First, we have fields related to police presence and the likelihood and severity of punishment for any given crime. These are our most likely levers for which we can estimate the causal relationships and compare the relative effects on crime rates, and politically can be adjusted by increasing police staffing or proposing legislation to improve crime prevention effectiveness and punishment deterrents through statute.

Next, we are given locational differences such as population density, which may offer targeted policy prescriptions or resource allocation, as well as additional effects when interacted with the aforementioned police and punishment severity levers. For example, we conjecture that there may be a difference in effect of adding police presence or increasing the severity of punishment for particular crimes in urban versus non-urban areas. We must recognize, unfortunately, that the size of the urban segment is only $n = 8$, which is likely produce too large of standard errors to make confident conclusions, but the direction of coefficients could still lead to more targeted future analysis on urban or non-urban policy.

Lastly, we have a number of economic variables including various wage levels, tax levies and demographics. While these may have explanatory power on our sample of counties, these variables do not have direct links to policy and will be incorporated later when attempting to control for non-crime factors and test for model robustness.

In expectation, we hypothesize that crime prevention proxies such as police per capita, average sentence duration, and likelihood of capture and punishment will have negative relationships with crime rates, while these effect sizes may differ in highly dense, or urban areas. We will then estimate the following population model with their raw, untransformed variables:

$$crrmrte = \beta_0 + \beta_1 polpc + \beta_2 avgscen + \beta_3 prbarr + \beta_4 prbconv + \beta_5 prbpris + \beta_6 density + \beta_7 urban + \beta_8 (polpc * urban) + \beta_9 (avgscen * urban)$$

```
model1 = lm(crrmrte ~ polpc + avgscen + prbarr + prbconv + prbpris + density + urban + (polpc * urban) +
            (avgscen * urban) + (prbarr * urban) + (prbconv * urban) + (prbpris * urban), data=raw_data)
# summary(model1)
```

From this regression output, we see we have fairly strong explanatory power of our sample, with approximately 75% of the variation in crime rate explained by the model, with several statistically significant relationships. We now test the six key assumptions of classical linear modeling:

- 1.) Linearity in parameters
- 2.) Random sample
- 3.) No perfect collinearity
- 4.) Zero conditional mean
- 5.) Homoskedasticity
- 6.) Normality of errors

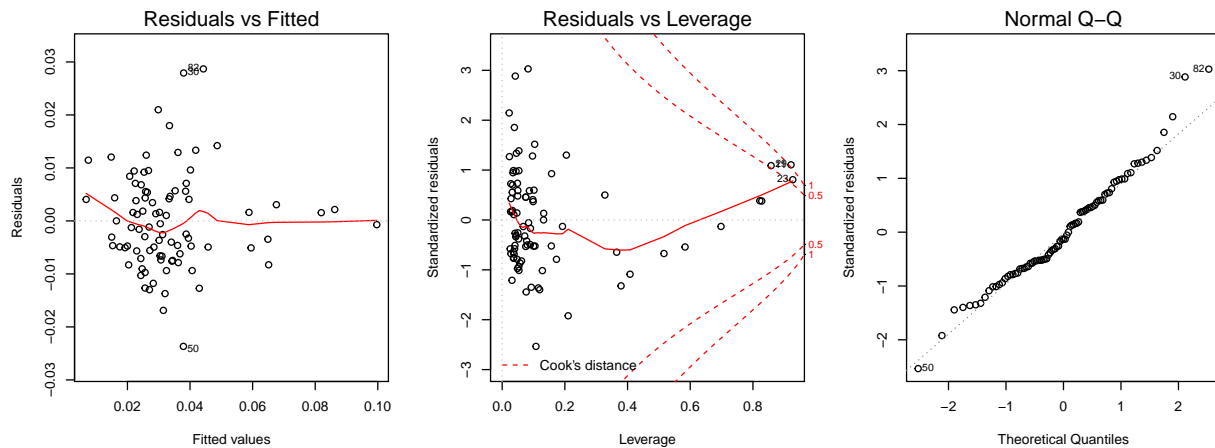
For assumptions #1-#3, we can treat each as satisfied given our population model definition: the coefficients we are estimated are strictly linear and the variables come from a cross-sectional dataset of county level data from 1987; while the sample is not randomly selected, it contains all the available cross-sections and is populated fully. Lastly, we know there is no perfect multicollinearity, as no fields can be derived from the values in another field; this is confirmed by our correlation matrix above where no variables are perfectly correlated, and the trivial fact that the model ran with no variables forced out.

We then must confirm that the expected values of our errors are all zero for any given observation and variable, as well as the homoskedasticity and normality of the errors. Looking at the first plot of residuals versus observed crime rate values, we see the fitted line of expected residuals (red) closely following zero on our residual axis, with no obvious trend as the observed crime rate values increase. Similarly, the Pearson residuals for each independent variable have expectation roughly tracking the zero line, thus not producing evidence that we do not have zero-conditional mean of the errors. However, we do see a couple data points with high leverage with Cook's distances approaching 1 (Chart: Residuals vs Leverage), likely due to the right-skew of both our depend variable and several covariates; we test a model with skewed variables transformed to more normal sample distributions shortly.

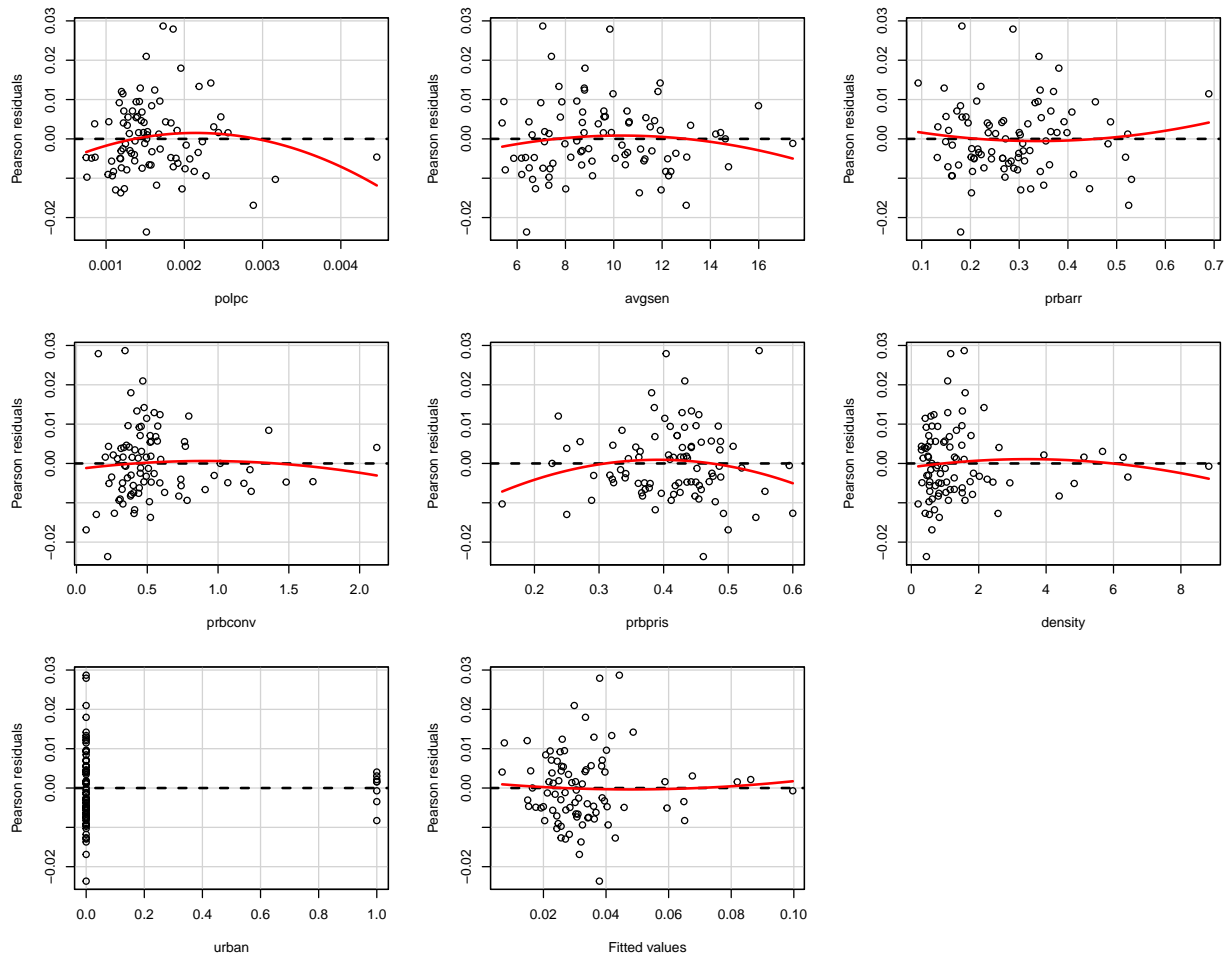
Looking at these same plots, we also do not see reason to believe there is heteroskedasticity in these data, where there would be distinct differences in variation of the errors dependent about the value of the

underlying variables; while we do see clustering of the fitted values due to skew of the underlying variables, this does not appear to lead to large differences in variance across the errors, with any trends being pulled by values at the extremes. And lastly, we can look to the Q-Q plot of the standardized residuals against the line that would represent a normal distribution and see that the residuals do appear to be close to normal, with a few points deviating from the line at the tales.

```
test_resids <- data.frame(model1$fitted.values, model1$residuals, raw_data$urban)
par(mfrow = c(1,3))
plot(model1, which = 1)
plot(model1, which = 5)
plot(model1, which = 2)
```



```
residualPlots(model1)
```

```
##          Test stat Pr(>|t|)
## polpc      -1.948   0.055
## avgsen     -0.979   0.331
## prbarr       0.684   0.496
## prbconv     -0.511   0.611
## prbpris     -1.470   0.146
## density     -1.279   0.205
## urban       0.765   0.446
## Tukey test   1.032   0.302
```

```
# acf(model1$residuals)
# pacf(model1$residuals)
```

While it would appear that interacting the *urban* indicator with our key crime covariates would absorb much of the effects of those variables, we should test the bias that omitting them would introduce as well as whether the explanatory power is significantly improved by adding them. By running the model with omitted interaction terms, we see below that the coefficients to our raw numeric variables are for the most part unchanged, while several of our interaction terms have strongly significant effects. However, by running an F-test on the two models, we do not see strong evidence that the full model has significantly more explanatory power than the restricted model with interaction terms removed ($p = .08$), so for purposes of strictly keeping key determinants of crime without sacrificing explanatory power, we would omit these terms and keep the restricted model:

```

model1_excl <- lm(crmrte ~ polpc + avgsgen + prbarr + prbconv + prbpris + density + urban, data=raw_data)
# summary(model1_excl)
stargazer(model1, model1_excl, type = "latex", report = "vcs*",
  header = FALSE,
  column.labels = c("Full Model", "Restricted Model"),
  title = "Comparison of Key Crime Determinants Only",
  star.cutoffs = c(0.05, 0.01, 0.001))

anova(model1, model1_excl)

```

Analysis of Variance Table

Model 1: $\text{crmte} \sim \text{polpc} + \text{avgsgen} + \text{prbarr} + \text{prbconv} + \text{prbpris} + \text{density} + \text{urban} + (\text{polpc} * \text{urban}) + (\text{avgsgen} * \text{urban}) + (\text{prbarr} * \text{urban}) + (\text{prbconv} * \text{urban}) + (\text{prbpris} * \text{urban})$ Model 2: $\text{crmte} \sim \text{polpc} + \text{avgsgen} + \text{prbarr} + \text{prbconv} + \text{prbpris} + \text{density} + \text{urban}$ Res.Df RSS Df Sum of Sq F Pr(>F)

	1	74	0.0072373
2	79	0.0082347	-5 -0.00099736 2.0396 0.08283 .

— Signif. codes: 0 ‘**0.001**’ ‘0.01’ ‘0.05’ ‘0.1’ ‘1’

Since we also noted the right skew of several of the independent variables, we consider whether the model is significantly improved when transforming these variables with the square root function. However, there now appears to be somewhat of a curved relationship between our residuals and observed crime rates, calling into question our assumption of zero-conditional mean of the errors. Looking at the Q-Q plot of standardized residuals, we continue to see similar approximate normality to the original model, still with some deviation from the normal curve at the extremes. Overall, the improvement does not appear to be strong enough to warrant the interpretation of square-roots of key variables:

```

model1_trans = lm(crmrte ~ sqrt(polpc) + sqrt(avgsgen) + prbarr + sqrt(prbconv) + prbpris + sqrt(density)
  + urban, data=raw_data)
summary(model1_trans)

```

```

##
## Call:
## lm(formula = crmrte ~ sqrt(polpc) + sqrt(avgsgen) + prbarr + sqrt(prbconv) +
##      prbpris + sqrt(density) + urban, data = raw_data)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -0.022552 -0.007402 -0.001698  0.006488  0.027458
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept)   0.029451   0.015055   1.956 0.053979 .
## sqrt(polpc)    0.451966   0.199223   2.269 0.026019 *
## sqrt(avgsgen) -0.003519   0.002831  -1.243 0.217553
## prbarr        -0.040003   0.012197  -3.280 0.001547 **
## sqrt(prbconv) -0.021853   0.006380  -3.425 0.000977 ***
## prbpris        0.010606   0.014024   0.756 0.451765
## sqrt(density)  0.016665   0.003943   4.227 6.31e-05 ***
## urban         0.009325   0.005978   1.560 0.122804
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 0.01025 on 79 degrees of freedom
## Multiple R-squared:  0.7081, Adjusted R-squared:  0.6822
## F-statistic: 27.37 on 7 and 79 DF,  p-value: < 2.2e-16

```

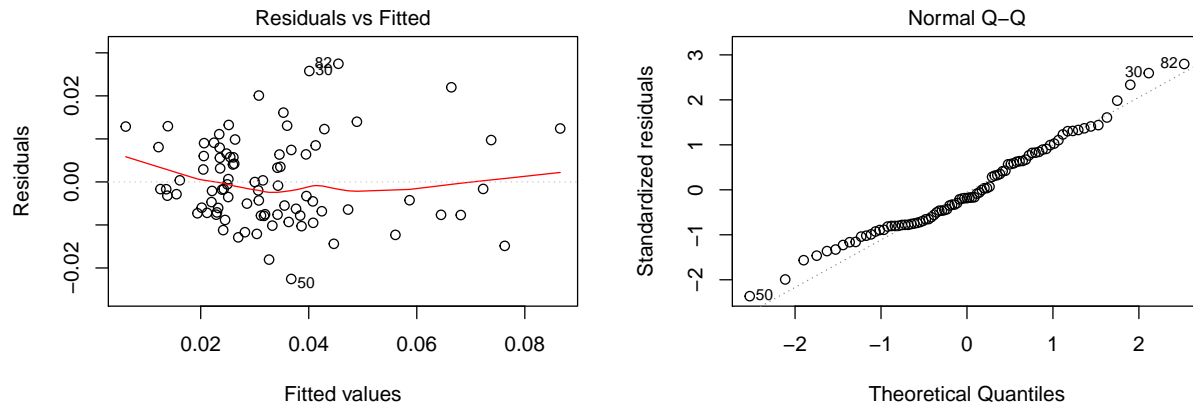
Table 1: Comparison of Key Crime Determinants Only

	<i>Dependent variable:</i>	
	crmte	
	Full Model	Restricted Model
	(1)	(2)
polpc	5.461 (2.177)*	5.346 (2.231)*
avgsen	−0.0004 (0.0004)	−0.0004 (0.0005)
prbarr	−0.044 (0.012)***	−0.046 (0.012)***
prbconv	−0.015 (0.004)***	−0.015 (0.004)***
prbpris	0.014 (0.014)	0.012 (0.014)
density	0.006 (0.002)**	0.006 (0.001)***
urban	0.590 (0.207)**	0.006 (0.007)
polpc:urban	22.181 (20.135)	
avgsen:urban	−0.015 (0.006)*	
prbarr:urban	−0.073 (0.099)	
prbconv:urban	−0.403 (0.143)**	
prbpris:urban	−0.766 (0.281)**	
Constant	0.034 (0.010)***	0.037 (0.010)***
Observations	87	87
R ²	0.745	0.710
Adjusted R ²	0.704	0.684
Residual Std. Error	0.010 (df = 74)	0.010 (df = 79)
F Statistic	18.040*** (df = 12; 74)	27.650*** (df = 7; 79)

Note:

*p<0.05; **p<0.01; ***p<0.001

```
par(mfrow = c(1,2))
plot(model1_trans, which = c(1,2))
```



Summary of Key Crime Determinants Model

Based on these three initial OLS regression models, before controlling for non-crime related control variables, we see several noteworthy effects that could influence the direction of public policy. At a high level, we see crime increase as population density increases (0.6% point increase in crime rate for each additional 100 people per square mile), and correspondingly, stronger enforcement—or higher likelihood of conviction and prison—in urban counties is associated with the largest decreases in crime rate. On the other hand, we see a moderately significant *positive* relationship between police per capita (increase in crime rate of approximately 5% points for an increase of one policeman per 100 people), which is likely reflective of a dual-causality problem: counties facing high crime rates may be responding by *then* increasing their policing presence. While we do proceed cautiously given the low sample size of our *urban* counties, the most effective allocation of resources to combat high crime rates would be improve the capabilities of law enforcement to make arrests and bring convictions to violent criminals. There appears to be less of an effect from either increasing police staff volume or sentence duration. Given these relationships, however, we need to test robustness by controlling for other factors, including demographic and economic variables.

Increasing Explanatory Power with Unbiased Controls

After determining the explanatory features that are most malleable to public policy, we look to identify the variables that are more difficult to influence, but nevertheless bare important relationships within the dataset. The distinguishing characteristics of theses explanatory features is their ability to increase the amount of variability explained while not introducing bias and violating any of the six aforementioned assumptions.

Model 1 focuses so specifically on crime and policy that it currently does not leverage much background information on each of the counties, and that is the objective of model 2. Starting with a demographic variable, introducing the tax revenue per capita to the model will enable an understanding of the relative resources available within the counties. Traditionally wealthier communities tend to have more resources and less crime and we expect to see that trend in North Carolina. Another demographic feature is the percentage of young males, the classic perpetrator of crime; naturally we anticipate a positive coefficient. Finally we have the bedrock of demographic information, the percentage of the population that's a minority, and we hypothesize that as diversity increases crime will too. A higher level characteristic of the counties is what region they belong to; however without more subject matter expertise its difficult to anticipate a trend. Outside of the demographic variables, we have added the crime mix variable which delineates the

more violent crimes (face to face) with the less interactive counterpart (other – think theft). The crime mix is difficult to control with public policy because its a facet of human nature, and didn't make it into model 1 because of that, but it could provide significant insight into crime rates. Typically cities have a reputation for more violent crimes and we expect to an interaction between the urban indicator and the mix variable. ###more on interactions...?

$$crmte = \beta_0 + \beta_1 polpc + \beta_2 avgsen + \beta_3 prbarr + \beta_4 prbconv + \beta_5 prbpris + \beta_6 density + \beta_7 urban + \beta_8 (polpc * urban) + \beta_9 (avgsen * urban)$$

```
model2 = lm(crmte ~ polpc + avgsen + prbarr + prbconv + prbpris + density + urban + (polpc * urban) +
            (avgsen * urban) + (prbarr * urban) + (prbconv * urban) + (prbpris * urban) + taxpc + west +
            pctymle + pctmin80 + mix + missingLabel, data=raw_data)
# model2
# summary(model2)
```

Adding interactions for model 2:

```
model2_rest <- lm(crmte ~ polpc + avgsen + prbarr + prbconv + prbpris + density + urban + (polpc * urban) +
                 (avgsen * urban) + (prbarr * urban) + (prbconv * urban) + (prbpris * urban) + taxpc + west +
                 central + pctymle + pctmin80 + mix + missingLabel, data=raw_data)
```

Comparing Model versions:

The control variables added to model 2 clearly enhanced the model. The output below shows an increase in adjusted R-squared from 0.70 to 0.82 and the anova test confirms the models are statistically significant so we reject the null hypothesis that the models are the same. Finally, model 2 also introduces the minority rate and it is highly statistically significant with a tiny p-value (although the power is very low). The results from the 2nd iteration of model 2 (includes the interactions) compared with the first version of model 2 are strikingly different. The model adjusted R-square moves less than 1/100th of a point and the anova test confirms the models are not significantly different.

```
#summary(model2_rest)
stargazer(model1, model2, model2_rest, type = "latex", report = "vcs",
          header = FALSE,
          title = "Comparison of Key Crime Determinants Only",
          star.cutoffs = c(0.05, 0.01, 0.001))
```

```
anova(model1, model2)
```

Analysis of Variance Table

Model 1: $crmte \sim polpc + avgsen + prbarr + prbconv + prbpris + density + urban + (polpc * urban) + (avgsen * urban) + (prbarr * urban) + (prbconv * urban) + (prbpris * urban)$ Model 2: $crmte \sim polpc + avgsen + prbarr + prbconv + prbpris + density + urban + (polpc * urban) + (avgsen * urban) + (prbarr * urban) + (prbconv * urban) + (prbpris * urban) + taxpc + west + central + pctymle + pctmin80 + mix + missingLabel$ Res.Df RSS Df Sum of Sq F Pr(>F)

1 74 0.0072373

2 68 0.0039589 6 0.0032784 9.3854 1.711e-07 *** — Signif. codes: 0 ‘’ **0.001** ’’ 0.01 ’’ 0.05 ‘:’ 0.1 ‘ ’ 1

```
anova(model2, model2_rest )
```

Analysis of Variance Table

Model 1: $crmte \sim polpc + avgsen + prbarr + prbconv + prbpris + density + urban + (polpc * urban) + (avgsen * urban) + (prbarr * urban) + (prbconv * urban) + (prbpris * urban) + taxpc + west + central + pctymle + pctmin80 + mix + missingLabel$ Model 2: $crmte \sim polpc + avgsen + prbarr + prbconv + prbpris + density + urban + (polpc * urban) + (avgsen * urban) + (prbarr * urban) + (prbconv * urban) +$

Table 2: Comparison of Key Crime Determinants Only

	<i>Dependent variable:</i>		
		crmrte	
	(1)	(2)	(3)
polpc	5.461 (2.177)*	5.897 (2.090)**	6.120 (2.684)*
avgsen	-0.0004 (0.0004)	-0.0003 (0.0004)	-0.0002 (0.0004)
prbarr	-0.044 (0.012)***	-0.049 (0.011)***	-0.049 (0.011)***
prbconv	-0.015 (0.004)***	-0.019 (0.003)***	-0.018 (0.004)***
prbpris	0.014 (0.014)	0.007 (0.011)	-0.0004 (0.012)
density	0.006 (0.002)**	0.007 (0.002)***	0.007 (0.002)***
urban	0.590 (0.207)**	0.516 (0.173)**	0.360 (0.218)
taxpc		-0.0001 (0.0001)	-0.0001 (0.0001)
west		-0.001 (0.004)	0.012 (0.014)
central		-0.002 (0.003)	-0.012 (0.009)
pctymle		0.022 (0.067)	0.093 (0.071)
pctmin80		0.0004 (0.0001)***	0.0004 (0.0001)***
mix		-0.019 (0.014)	-0.036 (0.021)
missingLabel			
polpc:urban	22.181 (20.135)	3.575 (16.267)	10.582 (23.808)
avgsen:urban	-0.015 (0.006)*	-0.013 (0.005)**	-0.013 (0.008)
prbarr:urban	-0.073 (0.099)	0.030 (0.084)	0.008 (0.121)
prbconv:urban	-0.403 (0.143)**	-0.319 (0.119)**	-0.232 (0.133)

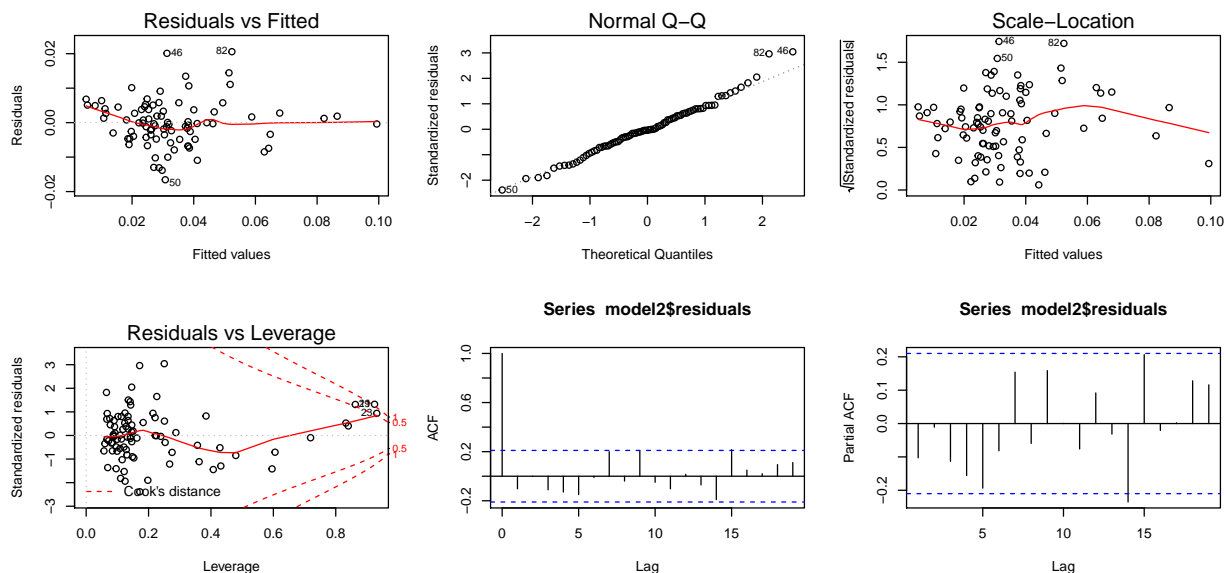
```
(prbpris * urban) + taxpc + west + central + pctymle + pctmin80 + mix + (west * polpc) + (central * polpc) + (central * mix) + (west * mix) + (urban * mix) + (urban * taxpc) + (taxpc * west) Res.Df RSS
Df Sum of Sq F Pr(>F) 1 68 0.0039589
2 61 0.0033974 7 0.00056143 1.44 0.2061
```

A quick note on the assumptions of the versions on model2:

After running the diagnostic plots, we have additional evidence to reject the version of model2 that has the interactions. Introducing interactions causes the model to have heteroskedacity (shown by a clear upward trend in the scale-location plot), reduces the normality in the Q-plot, and has 8 points with leverage above 1. Given the number of variables in the model it must be overfitting. The more parsimonious version of model 2 does not engender violated assumptions and showed results similar to model1 where the assumptions are not violated.

The results of the model comparison and the assumptions analysis both suggest that the simpler version of the control model should be the go forward model.

```
#light version of model2
par(mfrow = c(2,3))
plot(model2, which = c(1,2,3,5))
acf(model2$residuals)
pacf(model2$residuals)
```

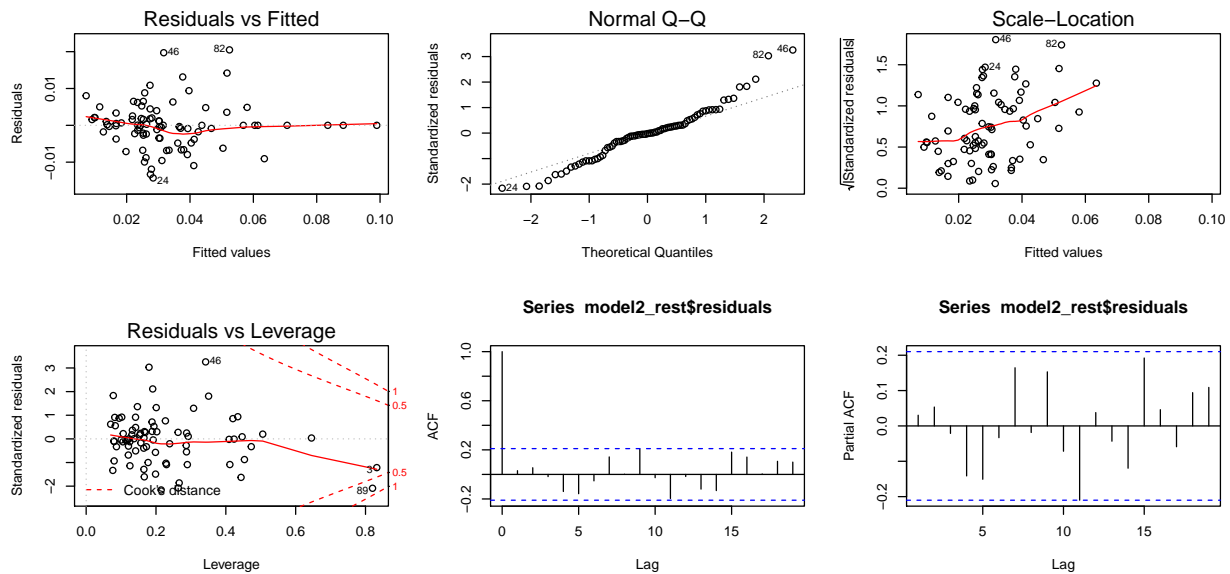


```
#model 2 (interactions)
par(mfrow = c(2,3))
plot(model2_rest, which = c(1,2,3,5))
```

```
## Warning: not plotting observations with leverage one:
## 11, 23, 28, 30, 35, 51, 55, 80
```

```
## Warning: not plotting observations with leverage one:
## 11, 23, 28, 30, 35, 51, 55, 80
```

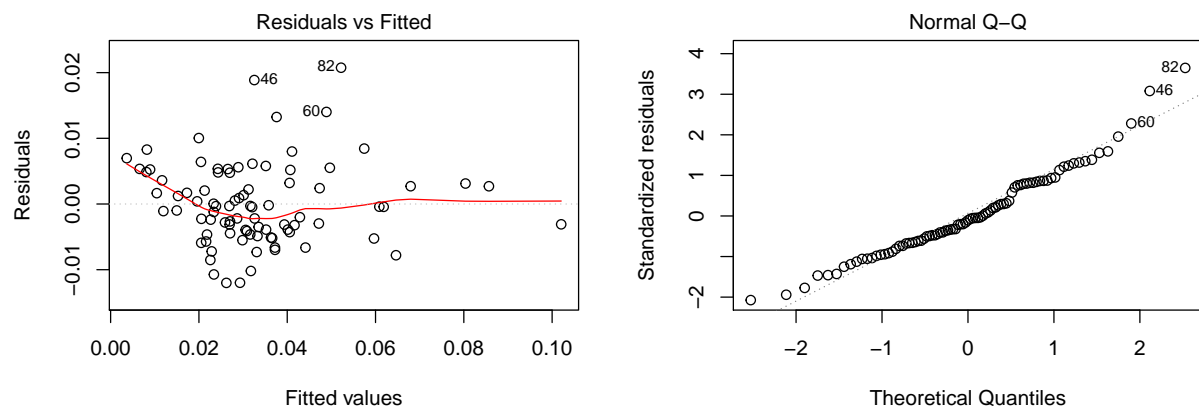
```
#residualPlots(model2)
acf(model2_rest$residuals)
pacf(model2_rest$residuals)
```

Testing Robustness with All Covariates Included

Having focused on the key determinants of crime from our available data while controlling for influential regional, demographic and economic factors, we now turn to including the remainder of our available variables in order to both maximize explanatory power of our model while testing whether the robustness of our original crime determinants hold, or whether there are strong biases generated by these other variables or omitted variables.

```
model3 <- lm(crmrte ~ polpc + avgsgen + prbarr + prbconv + prbpris + density + urban + (polpc * urban) +
              (prbarr * urban) + (prbconv * urban) + (prbpris * urban) + taxpc + west + central + pctymld +
              wcon + wtuc + wtrd + wfir + wmfgr + wfed + wsta + wloc, data=raw_data)
par(mfrow = c(1,2))
plot(model3, which = c(1,2))
```



```
stargazer(model3, model2, model1_excl, type = "latex", report = "vcs*",
  header = FALSE,
  column.labels = c("All Variables", "Key Control Variables", "Restricted Model"),
  title = "Comparison of Full Model with All Variables to Restricted Model",
  star.cutoffs = c(0.05, 0.01, 0.001))

anova(model2, model3)
```

Analysis of Variance Table

Model 1: $\text{crrmrte} \sim \text{polpc} + \text{avgsen} + \text{prbarr} + \text{prbconv} + \text{prbpris} + \text{density} + \text{urban} + (\text{polpc} * \text{urban}) + (\text{avgsen} * \text{urban}) + (\text{prbarr} * \text{urban}) + (\text{prbconv} * \text{urban}) + (\text{prbpris} * \text{urban}) + \text{taxpc} + \text{west} + \text{central} + \text{pctymle} + \text{pctmin80} + \text{mix} + \text{missingLabel}$

Model 2: $\text{crrmrte} \sim \text{polpc} + \text{avgsen} + \text{prbarr} + \text{prbconv} + \text{prbpris} + \text{density} + \text{urban} + (\text{polpc} * \text{urban}) + (\text{avgsen} * \text{urban}) + (\text{prbarr} * \text{urban}) + (\text{prbconv} * \text{urban}) + (\text{prbpris} * \text{urban}) + \text{taxpc} + \text{west} + \text{central} + \text{pctymle} + \text{pctmin80} + \text{mix} + \text{wcon} + \text{wtuc} + \text{wtrd} + \text{wfir} + \text{wmfg} + \text{wfed} + \text{wsta} + \text{wloc}$

Res.Df RSS Df Sum of Sq F Pr(>F)

1	68	0.0039589			
2	60	0.0033009	8	0.00065797	1.495 0.1782

Immediately we see a relative increase in the R^2 in the all-inclusive model but no statistical significance from the newly-added wage variables. Further, we see that our most restricted model maintains approximately its same coefficients across the key crime determinants. Performing an F-test on the all-inclusive model and model with control variables, we do not see statistical evidence that the former is significantly improved in its fit. However, we do see a minor changes to the coefficients of key crime determinants and wonder if those can be at least partly be explained by other variables in the model, leading to bias in the restricted models.

Beginning with police per capita (*polpc*), we can similarly predict that level by using our remaining covariates from the full model. Looking at the magnitudes and significance levels of the new coefficients (below), we may have expected more of from *density*, percentage young male, or the likelihood of making arrests given a crime, but these relationships appear to be relatively small. So with stable and robust estimates of the coefficients for our key determinants, we turn to variables that we are unable to include in the model and the resulting bias they may inflict on our chosen variables.

```
model_polpc <- lm(polpc ~ avgsen + prbarr + prbconv + prbpris + density + urban + taxpc + west+ central
  wcon +wtuc + wtrd + wfir+ wser +wmfg +wfed +wsta +wloc, data=raw_data)
summary(model_polpc)
```

```
##
## Call:
## lm(formula = polpc ~ avgsen + prbarr + prbconv + prbpris + density +
##      urban + taxpc + west + central + pctymle + pctmin80 + mix +
##      wcon + wtuc + wtrd + wfir + wser + wmfg + wfed + wsta + wloc,
##      data = raw_data)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -6.403e-04 -2.393e-04 -4.735e-05  1.577e-04  1.926e-03
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept) -2.379e-03  1.017e-03  -2.339  0.022462 *
## avgsen       4.133e-05  2.213e-05   1.867  0.066430 .
## prbarr      -1.679e-04  6.346e-04  -0.265  0.792132
## prbconv       2.744e-04  2.077e-04   1.321  0.191087
## prbpris      -3.515e-04  6.547e-04  -0.537  0.593181
## density       1.227e-04  7.675e-05   1.598  0.114966
## urban       -4.706e-04  3.466e-04  -1.358  0.179360
```

Table 3: Comparison of Full Model with All Variables to Restricted Model

	<i>Dependent variable:</i>		
	crmte		
	All Variables	Key Control Variables	Restricted Model
	(1)	(2)	(3)
polpc	4.778 (2.179)*	5.897 (2.090)**	5.346 (2.231)*
avgsen	-0.0003 (0.0004)	-0.0003 (0.0004)	-0.0004 (0.0005)
prbarr	-0.051 (0.011)***	-0.049 (0.011)***	-0.046 (0.012)***
prbconv	-0.018 (0.003)***	-0.019 (0.003)***	-0.015 (0.004)***
prbpris	0.003 (0.011)	0.007 (0.011)	0.012 (0.014)
density	0.006 (0.002)**	0.007 (0.002)***	0.006 (0.001)***
urban	0.440 (0.173)*	0.516 (0.173)**	0.006 (0.007)
taxpc	-0.0001 (0.0001)	-0.0001 (0.0001)	
west	0.0003 (0.004)	-0.001 (0.004)	
central	-0.002 (0.003)	-0.002 (0.003)	
pctymle	0.086 (0.074)	0.022 (0.067)	
pctmin80	0.0004 (0.0001)***	0.0004 (0.0001)***	
mix	-0.007 (0.015)	-0.019 (0.014)	
wcon	0.00002 (0.00003)		
wtuc	0.00000 (0.00001)		
wtrd	0.00002 (0.00004)		
wfir	-0.00004 (0.00002)		
wmfg	0.00000 (0.00001)		

```

## taxpc      2.102e-05  6.344e-06  3.313 0.001521 **
## west       4.427e-04  2.241e-04  1.976 0.052505 .
## central    1.672e-04  1.553e-04  1.077 0.285733
## pctymle    8.957e-03  3.668e-03  2.442 0.017395 *
## pctmin80   4.372e-06  5.405e-06  0.809 0.421667
## mix        2.753e-03  7.957e-04  3.460 0.000967 ***
## wcon       -9.920e-07  1.522e-06 -0.652 0.516993
## wtuc       -7.798e-07  8.079e-07 -0.965 0.338080
## wtrd       -2.786e-06  2.526e-06 -1.103 0.274316
## wfir       -2.606e-07  1.472e-06 -0.177 0.859980
## wser       -1.742e-07  3.025e-07 -0.576 0.566594
## wmfg        5.053e-07  7.747e-07  0.652 0.516622
## wfed        4.153e-06  1.551e-06  2.678 0.009407 **
## wsta       -4.975e-07  1.451e-06 -0.343 0.732827
## wloc        2.596e-06  2.616e-06  0.993 0.324641
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 0.0004433 on 64 degrees of freedom
## (1 observation deleted due to missingness)
## Multiple R-squared:  0.533, Adjusted R-squared:  0.3797
## F-statistic: 3.478 on 21 and 64 DF, p-value: 6.359e-05

```

Discussion of Limitations and Possible Omitted Variable Bias

Missing from this analysis are a number of variables that almost certainly influence an individual's likelihood to commit crime, as well as the estimates of our key variables in each of the above models. Namely, in addition to controls for gender, geographical region and minority mix, we would want to control unemployment, education levels (high school and college degrees), and rates of divorce or single-parent homes. While our models suggest a positive relationship between crime rates and population density, these other variables routinely have an effect on individuals' crime propensities and tend to lag for many families in more urban environments. As a result, we likely have a positive bias, or an overstatement, of the effect of population density on crime rather than the *real* driving factors like economic opportunity and family stability. Correspondingly, we see a positive relationship between crime and police per capita, suggesting that simply hiring a larger police force is not causing a reversal in the real determinants of crime. Furthermore, these models do not capture the effectiveness of government-sponsored programs, such as rehabilitation, nor the true crime rate—as some more disjointed communities may be less likely to actually report certain or many crimes. It is these communities for which we would like to see the effectiveness of current policy before advocating to allocate significant resources to, say, densely populated regions of North Carolina.

Conclusion

We conclude that while there are limitations to the explanatory power of these models due to omitted key variables, the estimates coming from these models are sufficiently robust to narrow the policy discussion as well as avoid political platitudes such as expanding police force or strengthening sentencing laws. Instead, we see opportunity to improve local law enforcement's effectiveness in bringing convictions and enforcing the laws as they are written, with particular effort toward more densely-populated regions of the state. This initial analysis also points to further study of problems afflicting these urban areas, with likely carryover to the rest of the state as well. Specifically, the strong positive relationship between crime and density is likely a *causal* one between underlying factors like educational opportunity and family structure, factors that we do not observe in these data.

Appendix: Model Summary

```
stargazer(model2, model1, model1_excl, type = "latex", report = "vcs*",
  header = FALSE,
  column.labels = c("Key Control Variables", "Urban Interactions", "Restricted Model"),
  title = "Summary of Primary Models in Crime Rate Analysis",
  star.cutoffs = c(0.05, 0.01, 0.001))
```

Table 4: Summary of Primary Models in Crime Rate Analysis

	<i>Dependent variable:</i>		
	Key Control Variables	crmte	Restricted Model
		Urban Interactions	
	(1)	(2)	(3)
polpc	5.897 (2.090)**	5.461 (2.177)*	5.346 (2.231)*
avgsen	−0.0003 (0.0004)	−0.0004 (0.0004)	−0.0004 (0.0005)
prbarr	−0.049 (0.011)***	−0.044 (0.012)***	−0.046 (0.012)***
prbconv	−0.019 (0.003)***	−0.015 (0.004)***	−0.015 (0.004)***
prbpris	0.007 (0.011)	0.014 (0.014)	0.012 (0.014)
density	0.007 (0.002)***	0.006 (0.002)**	0.006 (0.001)***
urban	0.516 (0.173)**	0.590 (0.207)**	0.006 (0.007)
taxpc	−0.0001 (0.0001)		
west	−0.001 (0.004)		
central	−0.002 (0.003)		
pctymle	0.022 (0.067)		
pctmin80	0.0004 (0.0001)***		
mix	−0.019 (0.014)		
missingLabel			
polpc:urban	3.575 (16.267)	22.181 (20.135)	
avgsen:urban	−0.013 (0.005)**	−0.015 (0.006)*	
prbarr:urban	0.030 (0.084)	−0.073 (0.099)	
prbconv:urban	−0.319 (0.119)**	−0.403 (0.143)**	