

Lab3 Draft, w203: Statistics for Data Science

Avinash Chandrasekaran, Deepak Nagaraj, Saurav Datta

March 31, 2018

1. Introduction

Our team has been hired to provide research for a political campaign. The campaign has obtained a dataset of crime statistics for a selection of counties in North Carolina. Our task is to examine the data to help the campaign understand the determinants of crime and to generate policy suggestions that are applicable to local government.

The data provided consists of 25 variables and 91 different observations collected in a given year. Moreover the dataset obtained is a single cross-section of data collected from variety of different sources. For the analysis made in this research, we will assume that the data collected from different counties in NC were randomly sampled.

Our primary analysis of data will include ordinary least squares regressions to make casual estimates and we will clearly explain how omitted variables may affect our conclusions. We begin our research by conducting exploratory analysis of the dataset to gain a better understanding of the variables.

2. Data Input

Let us read the data and have a first look.

```
# Read the csv file
crime_data_raw = read.csv("crime_v2.csv")
```

Empty rows

There appears to be 6 rows of NA's across all variables. We can simply use `na.omit()`, because the number of all-NA rows matches the count on all the variables.

Column formatting

We also notice that 'prbconv' is a factor while the rest of the variables are numeric. West, central and urban are really categorical variables and not integers: we will treat them as such.

```
# Remove NA rows
crime_data = na.omit(crime_data_raw)

# convert factor to numeric for variable prbconv
crime_data$prbconv = as.numeric(levels(crime_data$prbconv)[crime_data$prbconv])
```

Unused variables

County and Year variables just represent the different counties and the year the data was collected. Year is always 87. Hence, we can safely remove these from the dataset for further analysis.

```
crime_data = crime_data %>% dplyr::select(-c(year, county))
```

Duplicate records

We also noticed a duplicate record (record #89) in the dataset. As this could potentially affect our regression analysis, we will remove the duplicate record.

```
duplicated(crime_data)[duplicated(crime_data) == TRUE]
```

```
## [1] TRUE
```

```
crime_data = distinct(crime_data)
```

3. Exploratory Data Analysis

We will start with an explanatory note on transformations. Any skew in the original data may cause the residuals not to follow normal distribution. If this happens, it violates an assumption of the LS regression model: we will not be able to draw inferences from our model. Hence it is important to ensure our residuals to follow normal distribution as much as possible, and to transform our predictors if that helps.

We will now try to get a sense of each variable in the dataset.

Single variable analysis

There are a total of 90 observations across 23 different variables. We will now explore each of the variables collected in the data.

```
# Utility function to describe a column variable
f_describe_col = function(col, do_log = FALSE, plot_model = FALSE, do_sqrt = FALSE) {
  y = log(crime_data$crmrte)
  print(summary(col))
  par(mfrow = c(2, 2))
  if (is.numeric(col)) {
    hist(col, main = "Histogram")
    boxplot(col, main = "Box plot")
  }
  if (do_log == TRUE) {
    x = log(col)
    hist(x, main = "Histogram, log")
  } else if (do_sqrt == TRUE) {
    x = sqrt(col)
    hist(x, main = "Histogram, sqrt")
  } else {
    x = col
  }
  if (is.numeric(col)) {
    print(signif(cor(x, y), 3))
  }
  m = lm(y ~ x)
  print(summary(m))
  plot(x, y, main = "Cor. with crime rate")
  if (is.numeric(col))
```

```

    abline(m, col = "blue")
  if (plot_model == TRUE) {
    plot(m)
    if (do_log == TRUE) {
      m = lm(y ~ x)
      plot(m)
    }
  }
}

```

Wage distribution: Services

We will start with wage distribution in services industry.

There is a large outlier for wages in service industry, *wser*. Observation #84 has very large influence as shown by Cook's distance.

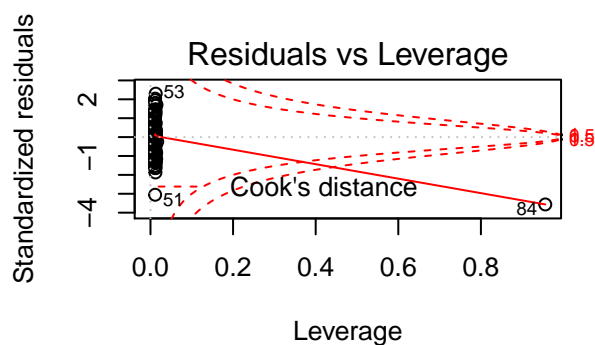
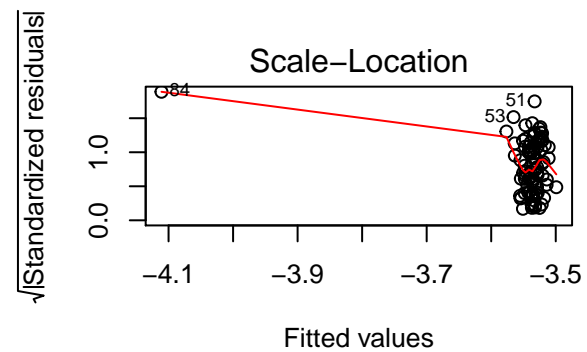
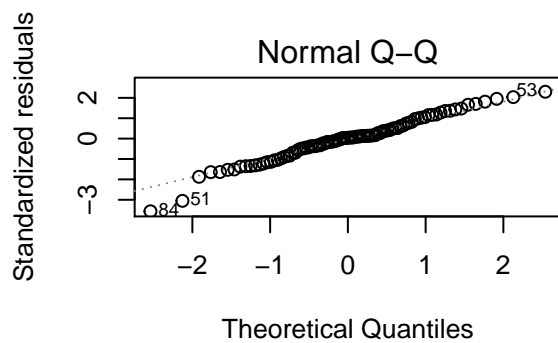
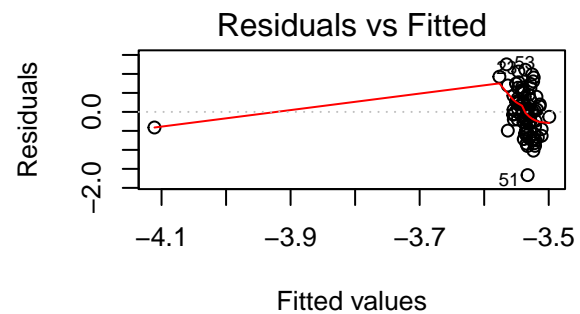
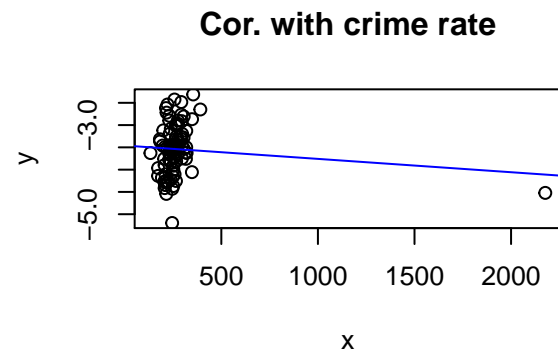
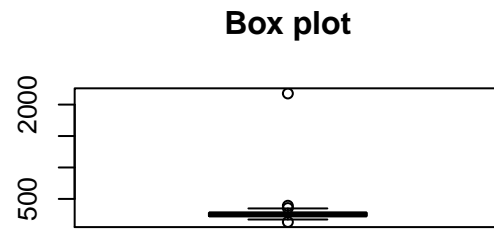
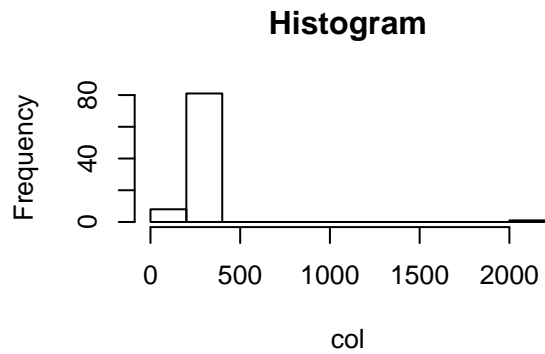
```
f_describe_col(crime_data$wser, plot_model = TRUE)
```

```

##      Min. 1st Qu.  Median    Mean 3rd Qu.    Max.
##    133.0   229.3   253.1   275.3   277.6  2177.1

## [1] -0.113
##
## Call:
## lm(formula = y ~ x)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -1.66428 -0.35914  0.02214  0.32067  1.25238
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept) -3.4593123  0.0964071 -35.882  <2e-16 ***
## x           -0.0002993  0.0002802  -1.068    0.288
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 0.5483 on 88 degrees of freedom
## Multiple R-squared:  0.0128, Adjusted R-squared:  0.00158
## F-statistic: 1.141 on 1 and 88 DF,  p-value: 0.2884

```



Let us look at the outlier observation:

```
crime_data %>% slice(84) %>% select(everything())
```

```
## # A tibble: 1 x 23
```

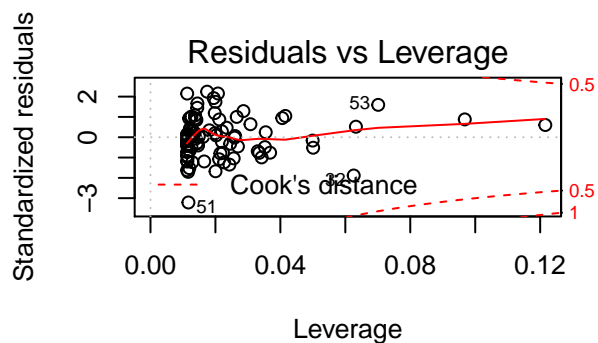
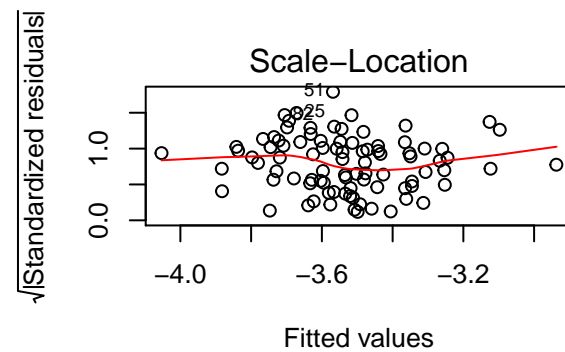
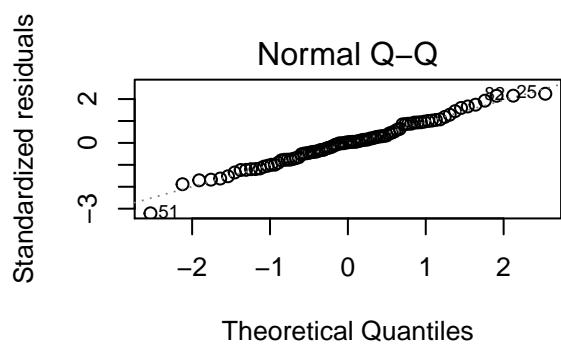
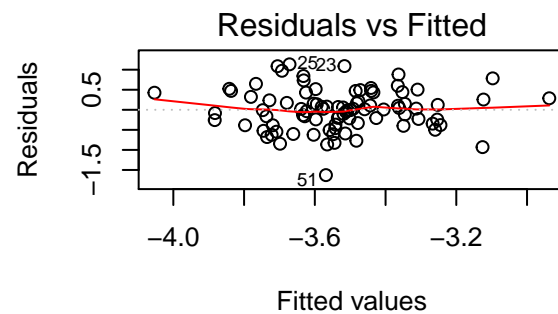
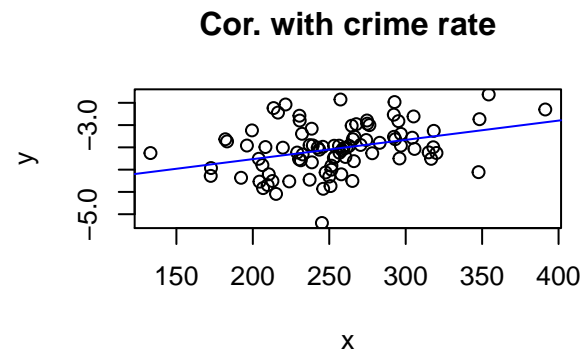
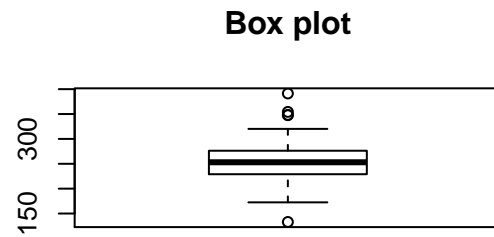
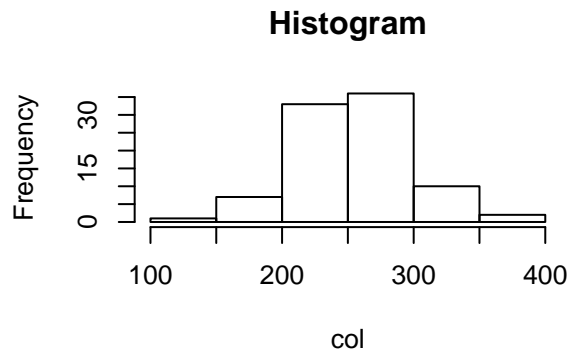
```
##   crmrte prbarr prbconv prbpris avgsen   polpc density taxpc  west central
##   <dbl> <dbl>   <dbl>   <dbl> <dbl>   <dbl>   <dbl> <dbl> <int>   <int>
## 1 0.0109 0.195    2.12   0.443   5.38 0.00122 0.389 40.8    0       1
## # ... with 13 more variables: urban <int>, pctmin80 <dbl>, wcon <dbl>,
## #   wtuc <dbl>, wtrd <dbl>, wfir <dbl>, wser <dbl>, wmfg <dbl>,
## #   wfed <dbl>, wsta <dbl>, wloc <dbl>, mix <dbl>, pctymle <dbl>
```

This value is way too large. Let us remove it from the dataset and replot.

```
crime_data = crime_data %>% slice(-84)
f_describe_col(crime_data$wser, plot_model = TRUE)
```

```
##      Min. 1st Qu.  Median      Mean 3rd Qu.      Max.
##    133.0  229.0   253.0   254.0   276.3   391.3

## [1] 0.352
##
## Call:
## lm(formula = y ~ x)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -1.62840 -0.34568  0.01331  0.32162  1.13342
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept) -4.628358   0.317962 -14.556  < 2e-16 ***
## x             0.004322   0.001234   3.503 0.000729 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 0.5101 on 87 degrees of freedom
## Multiple R-squared:  0.1236, Adjusted R-squared:  0.1135
## F-statistic: 12.27 on 1 and 87 DF,  p-value: 0.0007285
```



This looks better. We will come back to this variable later, when we consider all wages.

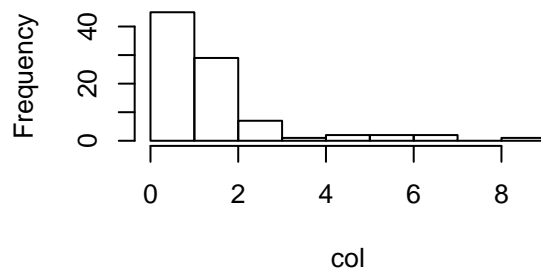
Population density

```
f_describe_col(crime_data$density, do_log = TRUE, plot_model = TRUE)
```

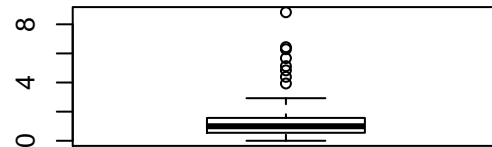
```
##      Min. 1st Qu.  Median    Mean 3rd Qu.    Max.
## 0.00002 0.54786 0.99623 1.44743 1.57028 8.82765

## [1] 0.491
##
## Call:
## lm(formula = y ~ x)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -1.49715 -0.25485  0.02209  0.30714  1.33192
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept) -3.51602    0.05039  -69.773 < 2e-16 ***
## x             0.19301    0.03671   5.257 1.03e-06 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 0.4747 on 87 degrees of freedom
## Multiple R-squared:  0.2411, Adjusted R-squared:  0.2324
## F-statistic: 27.64 on 1 and 87 DF,  p-value: 1.033e-06
```

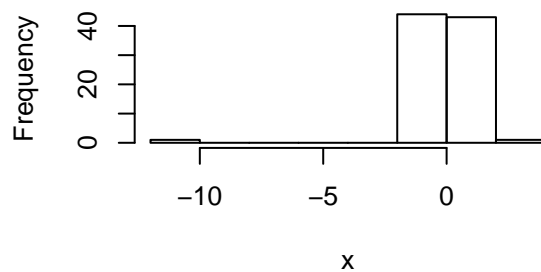
Histogram



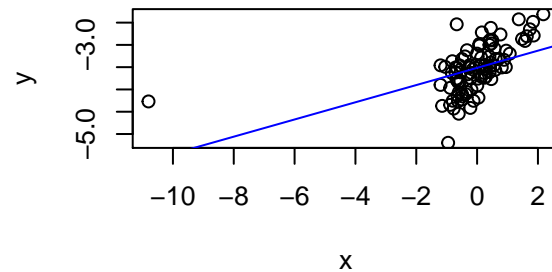
Box plot

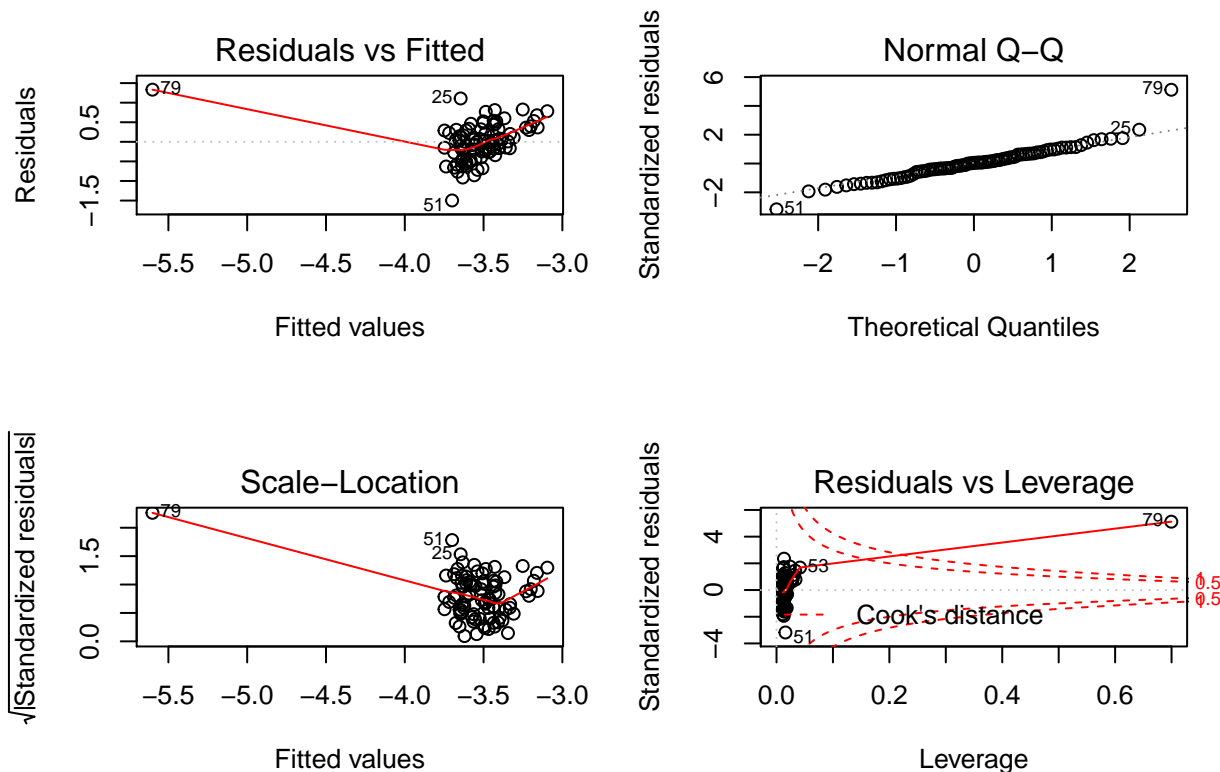


Histogram, log



Cor. with crime rate





We see that density is a highly skewed distribution. Most observations are from counties with low population density. However, observation #79 has Cook's distance beyond 1, meaning extreme leverage:

```
crime_data %>% slice(79) %>% select(everything())
```

```
## # A tibble: 1 x 23
##   crmrte prbarr prbconv prbpris avgsen  polpc density taxpc  west central
##   <dbl> <dbl> <dbl> <dbl> <dbl> <dbl> <dbl> <dbl> <int> <int>
## 1 0.0140 0.530 0.328 0.150 6.64 0.00316 2.03e-5 37.7 1 0
## # ... with 13 more variables: urban <int>, pctmin80 <dbl>, wcon <dbl>,
## #   wtuc <dbl>, wtrd <dbl>, wfir <dbl>, wser <dbl>, wmfg <dbl>,
## #   wfed <dbl>, wsta <dbl>, wloc <dbl>, mix <dbl>, pctymle <dbl>
```

The density is 2E-5, which is extremely low. With so few people, the observation may not add a lot of meaning. It also affects the model. We will remove this observation from the dataset and replot.

```
crime_data = crime_data %>% slice(-79)
f_describe_col(crime_data$density, do_log = TRUE, plot_model = TRUE)
```

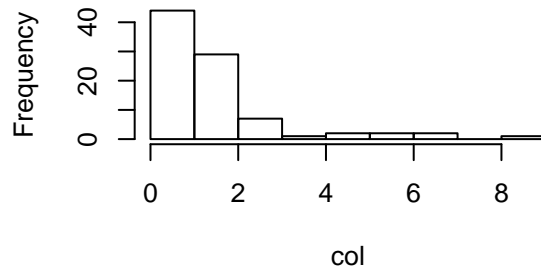
```
##   Min. 1st Qu.  Median    Mean 3rd Qu.    Max.
## 0.3006 0.5599 1.0008 1.4639 1.5762 8.8277

## [1] 0.677
##
## Call:
## lm(formula = y ~ x)
##
## Residuals:
##   Min       1Q   Median       3Q      Max
## -1.1980 -0.2924 -0.0115  0.2752  1.3262
##
## Coefficients:
```

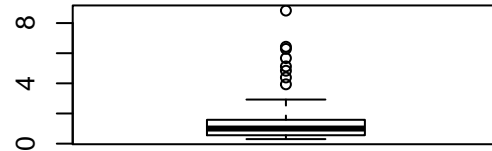


```
##           Estimate Std. Error t value Pr(>|t|)
## (Intercept) -3.54417    0.04262 -83.161  < 2e-16 ***
## x           0.47755    0.05600   8.528  4.5e-13 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 0.3991 on 86 degrees of freedom
## Multiple R-squared:  0.4582, Adjusted R-squared:  0.4519
## F-statistic: 72.72 on 1 and 86 DF,  p-value: 4.504e-13
```

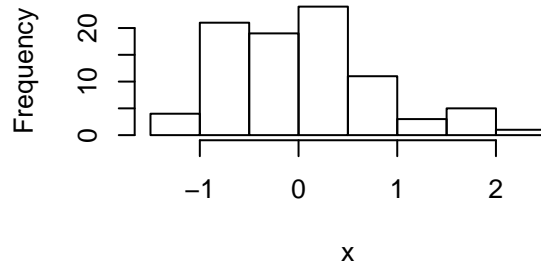
Histogram



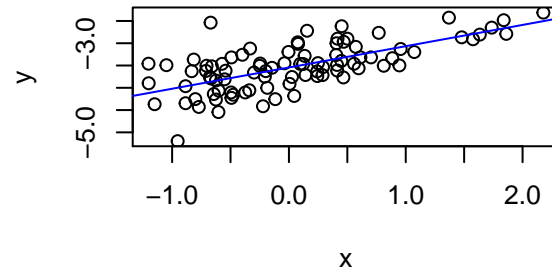
Box plot

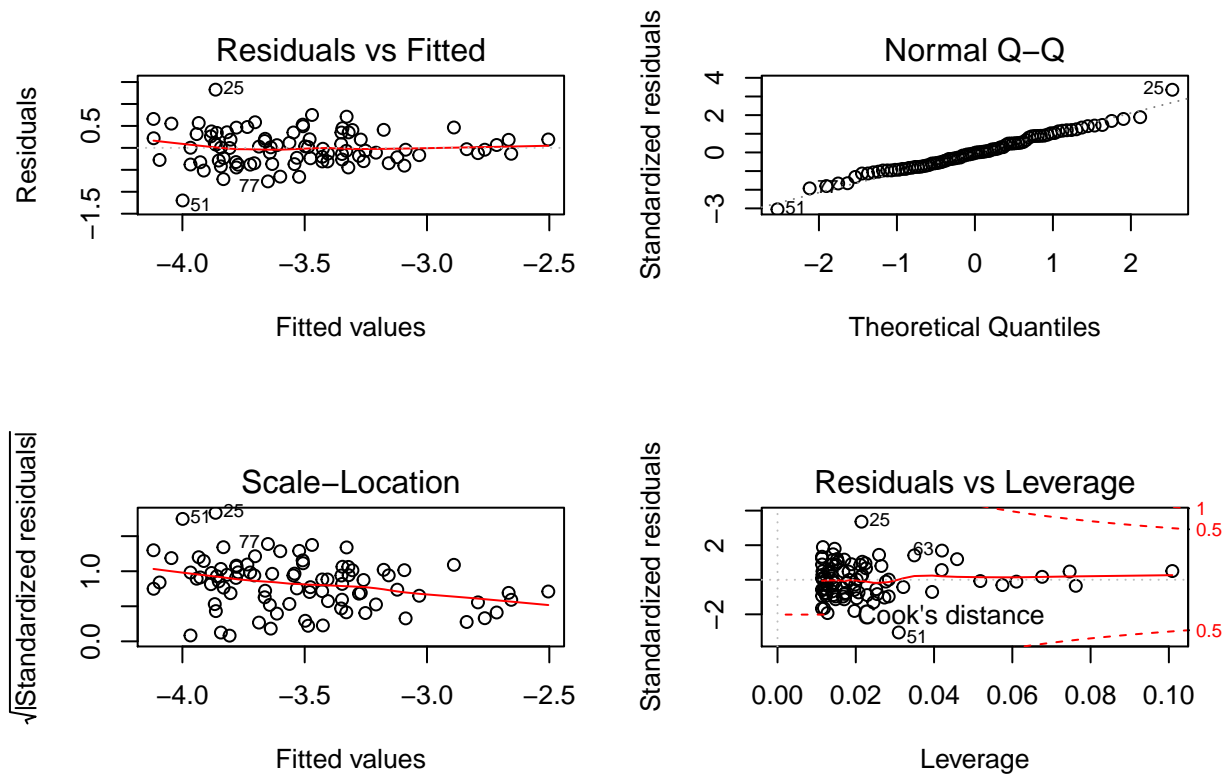


Histogram, log



Cor. with crime rate





The histogram of density shows quite a bit positive skew. The log transformation shows a more promising normal distribution. There are no more outliers with large leverage as measured by Cook's distance.

We see high positive correlation with crime rate. We will surely consider this variable in our model.

Police per capita

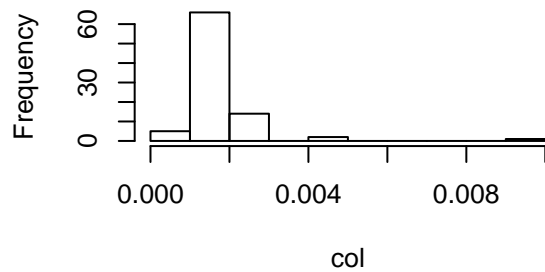
```
f_describe_col(crime_data$polpc, do_log = TRUE, plot_model = TRUE)
```

```
##      Min.   1st Qu.   Median     Mean  3rd Qu.    Max.
## 0.0007459 0.0012448 0.0014897 0.0016970 0.0018679 0.0090543

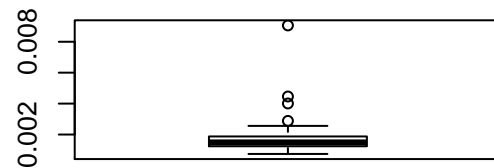
## [1] 0.316
##
## Call:
## lm(formula = y ~ x)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -2.48377 -0.26110  0.04185  0.29879  1.04446
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept)  -0.5494     0.9651  -0.569  0.57066
## x              0.4599     0.1491   3.085  0.00274 **
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 0.5144 on 86 degrees of freedom
```

```
## Multiple R-squared:  0.09966,    Adjusted R-squared:  0.08919
## F-statistic: 9.519 on 1 and 86 DF,  p-value: 0.002735
```

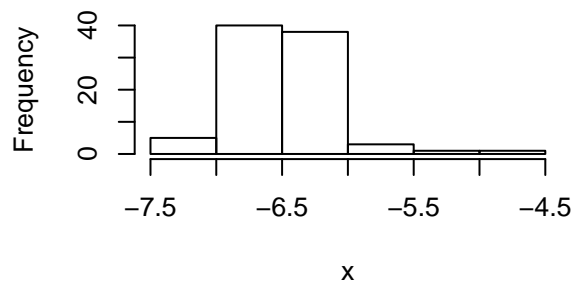
Histogram



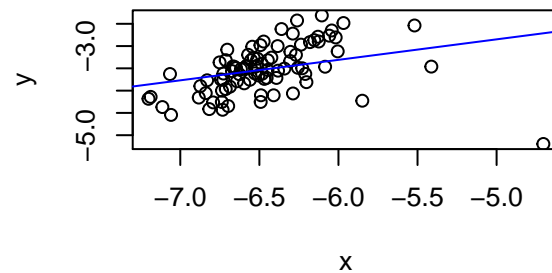
Box plot



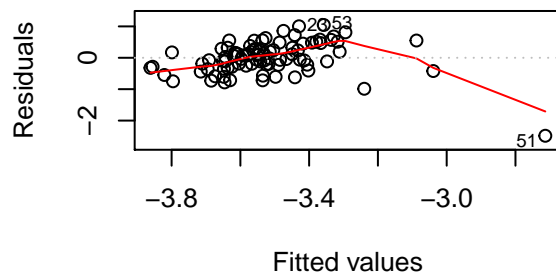
Histogram, log



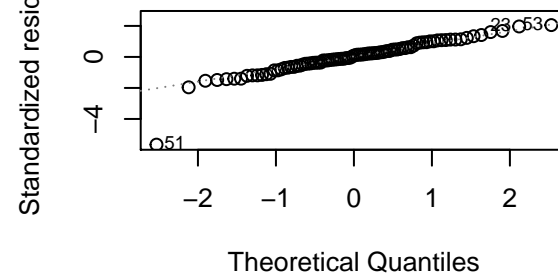
Cor. with crime rate



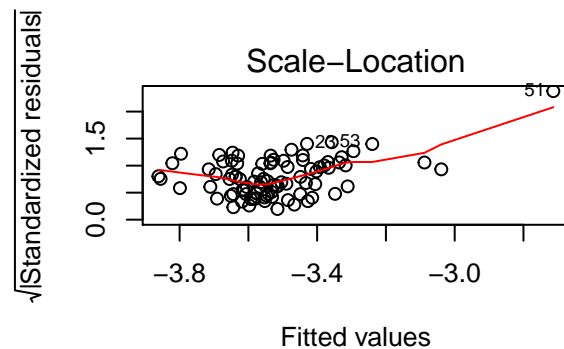
Residuals vs Fitted



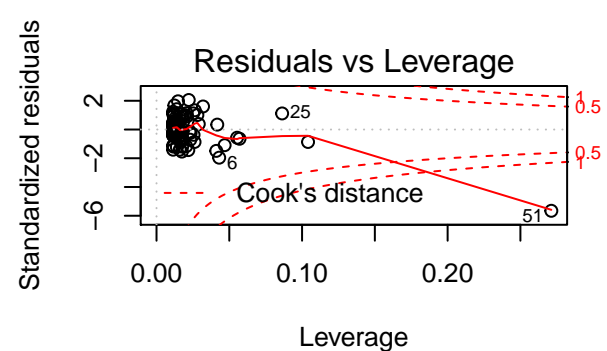
Normal Q-Q



Scale-Location



Residuals vs Leverage



Police per capita has positive skew. Taking log helps, but we still see a very large outlier. Fitting a model, we see that observation #51 has Cook's distance beyond 1. This is a lot of leverage. Let us look at the

observation:

```
crime_data %>% slice(51) %>% select(everything())
```

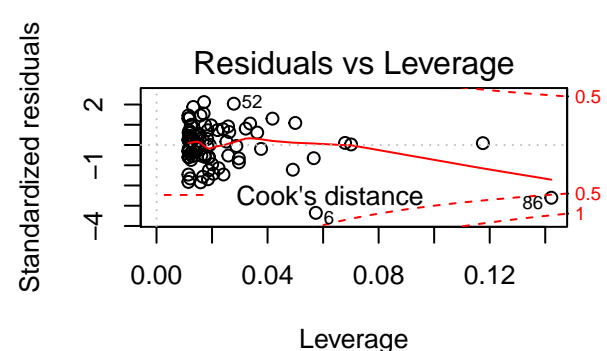
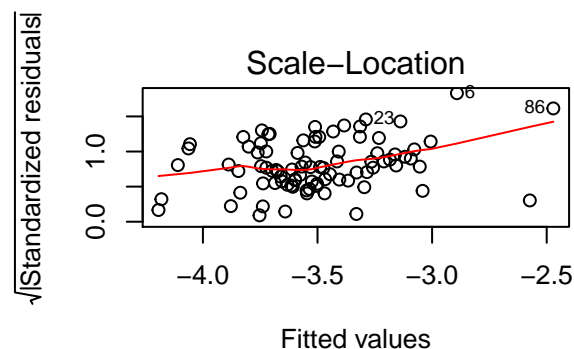
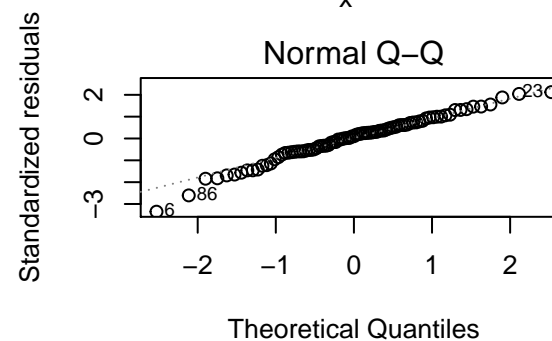
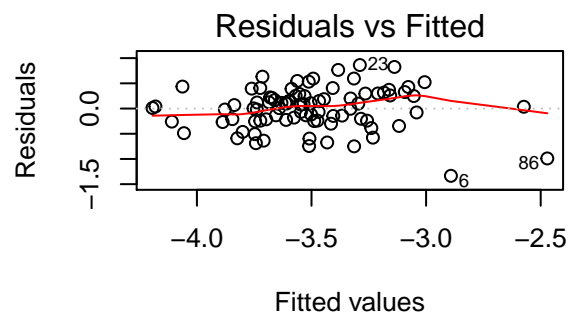
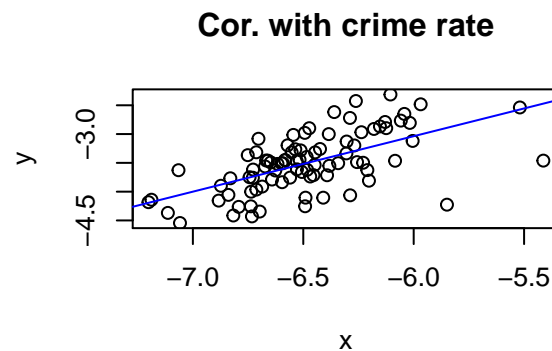
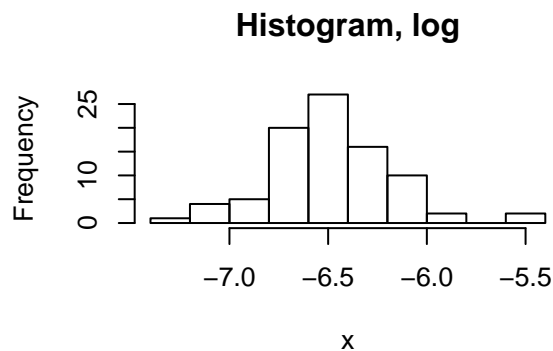
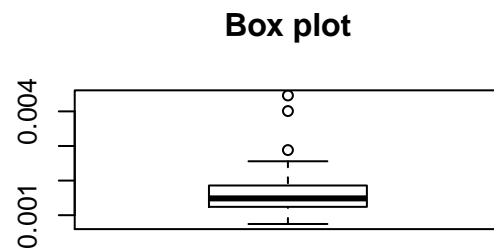
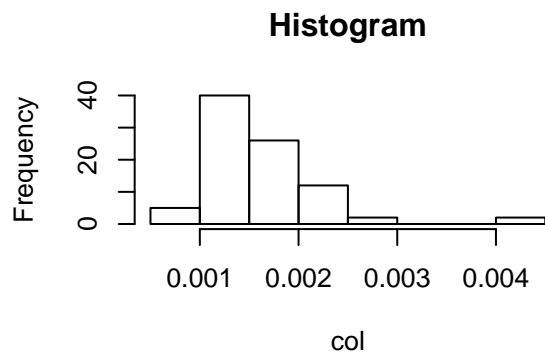
```
## # A tibble: 1 x 23
##   crmrte prbarr prbconv prbpris avgsen  polpc density taxpc  west
##   <dbl> <dbl> <dbl> <dbl> <dbl> <dbl> <dbl> <dbl> <int>
## 1 0.00553 1.09 1.50 0.500 20.7 0.00905 0.386 28.2 1
## # ... with 14 more variables: central <int>, urban <int>, pctmin80 <dbl>,
## #   wcon <dbl>, wtuc <dbl>, wtrd <dbl>, wfir <dbl>, wser <dbl>,
## #   wmfg <dbl>, wfed <dbl>, wsta <dbl>, wloc <dbl>, mix <dbl>,
## #   pctymle <dbl>
```

The crime rate, at 0.005, is extremely low for police per capita close to the maximum. It is questionable whether this observation is accurate. We will remove this reading from the dataset and replot.

```
crime_data = crime_data %>% slice(-51)
f_describe_col(crime_data$polpc, do_log = TRUE, plot_model = TRUE)
```

```
##      Min.    1st Qu.      Median        Mean     3rd Qu.      Max.
## 0.0007459 0.0012413 0.0014853 0.0016124 0.0018583 0.0044592

## [1] 0.603
##
## Call:
## lm(formula = y ~ x)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -1.33462 -0.23989  0.04081  0.25531  0.86296
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept)   2.7439     0.8985   3.054 0.00302 **
## x              0.9635     0.1384   6.961 6.63e-10 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 0.4101 on 85 degrees of freedom
## Multiple R-squared:  0.3631, Adjusted R-squared:  0.3556
## F-statistic: 48.45 on 1 and 85 DF,  p-value: 6.627e-10
```



```
crime_data$log_polpc = log(crime_data$polpc)
```

The distribution looks better now. We see quite strong positive correlation of 0.6 with crime rate: high number of police per capita is associated with high crime rate. It is probably a cause, rather than a result. More police may have been deployed to deal with higher amount of crime.

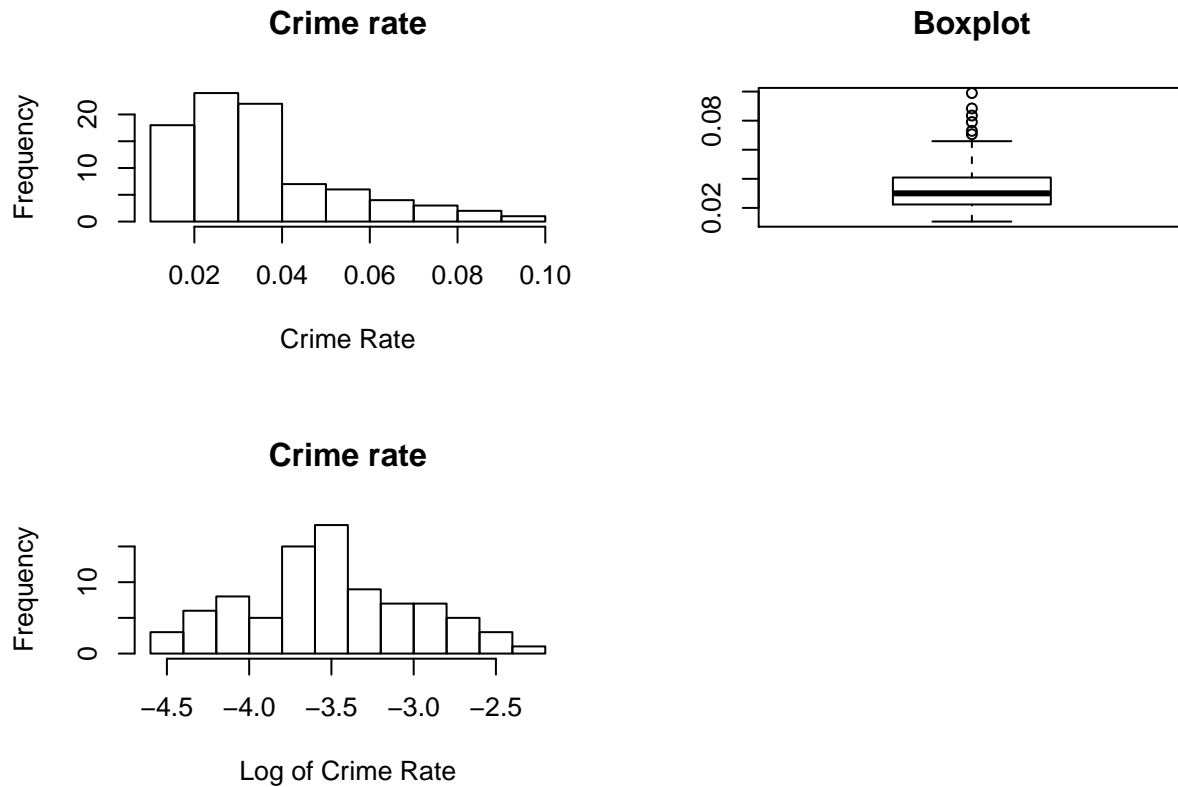
For our first model, we will *not* include this variable.

Crime rate

```
summary(crime_data$crmrte)
```

```
##      Min. 1st Qu.  Median    Mean 3rd Qu.    Max.
## 0.01062 0.02238 0.03002 0.03432 0.04090 0.09897
```

```
par(mfrow = c(2, 2))
hist(crime_data$crmrte, main = "Crime rate", xlab = "Crime Rate")
boxplot(crime_data$crmrte, main = "Boxplot")
hist(log(crime_data$crmrte), main = "Crime rate", xlab = "Log of Crime Rate")
```



The crime rate variable is the key dependent variable of interest. Looking at the histogram, the distribution is positively skewed to the left. We can take the log transformation which makes the variable appear more normally distributed.

Probability of arrest

```
f_describe_col(crime_data$prbarr, plot_model = TRUE)
```

```
##      Min. 1st Qu.  Median    Mean 3rd Qu.    Max.
## 0.09277 0.20568 0.27095 0.28454 0.34323 0.68902
```

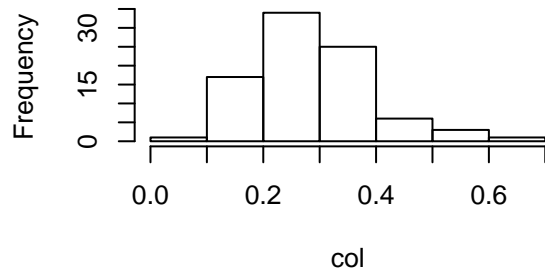
```
## [1] -0.374
```

```
##
```

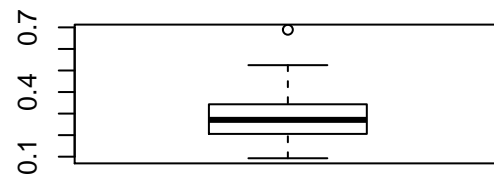
```
## Call:
```

```
## lm(formula = y ~ x)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -1.00436 -0.28274  0.02566  0.27730  0.94751
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept)  -2.9935     0.1463  -20.464  < 2e-16 ***
## x             -1.7911     0.4817   -3.718  0.000359 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 0.4765 on 85 degrees of freedom
## Multiple R-squared:  0.1399, Adjusted R-squared:  0.1298
## F-statistic: 13.82 on 1 and 85 DF,  p-value: 0.000359
```

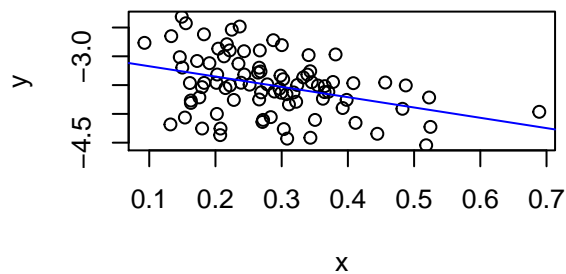
Histogram



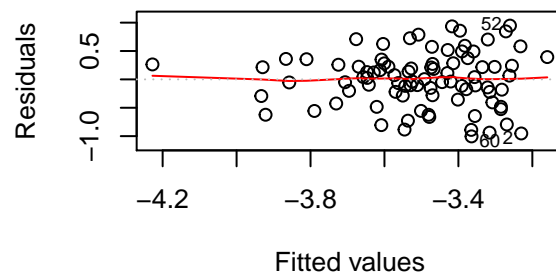
Box plot

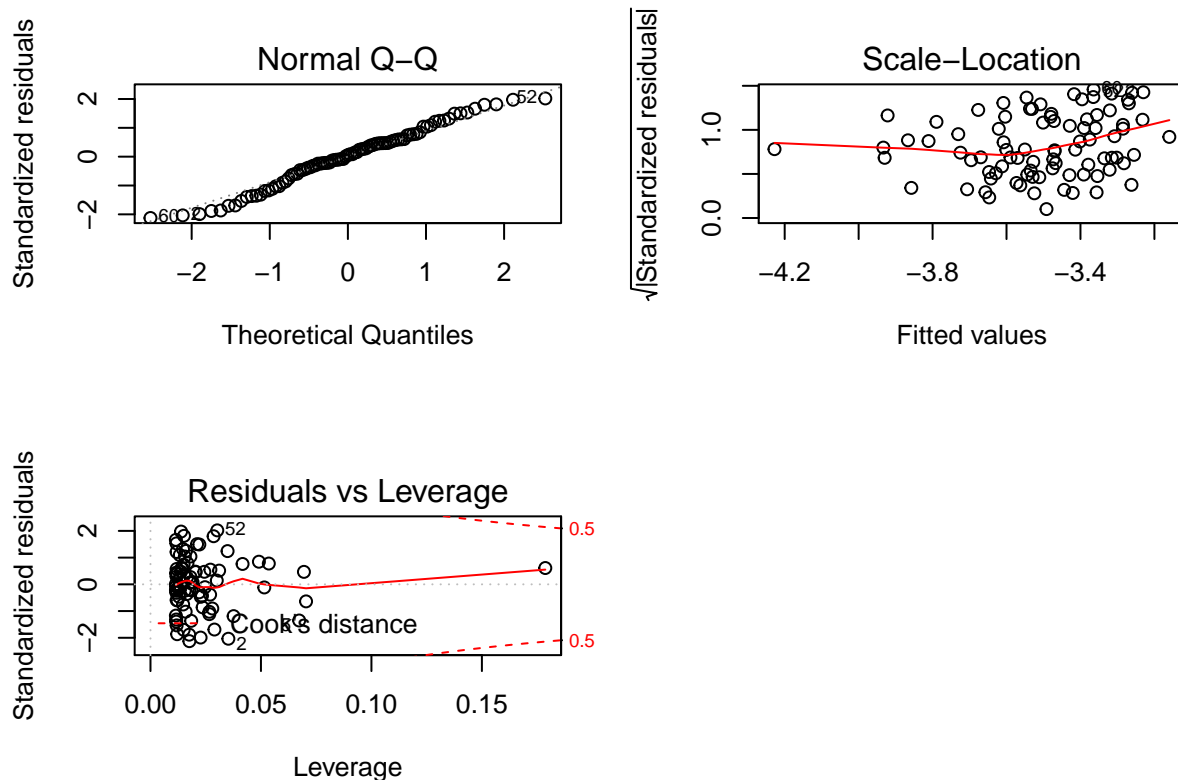


Cor. with crime rate



Residuals vs Fitted





The plot looks fairly normal; there is only one outlier.

There is fairly negative correlation of -0.37: as probability of arrests increases, crime rate goes down. It may be that arrests are a deterrent.

We will include *prbarr* in our model.

Probability of conviction

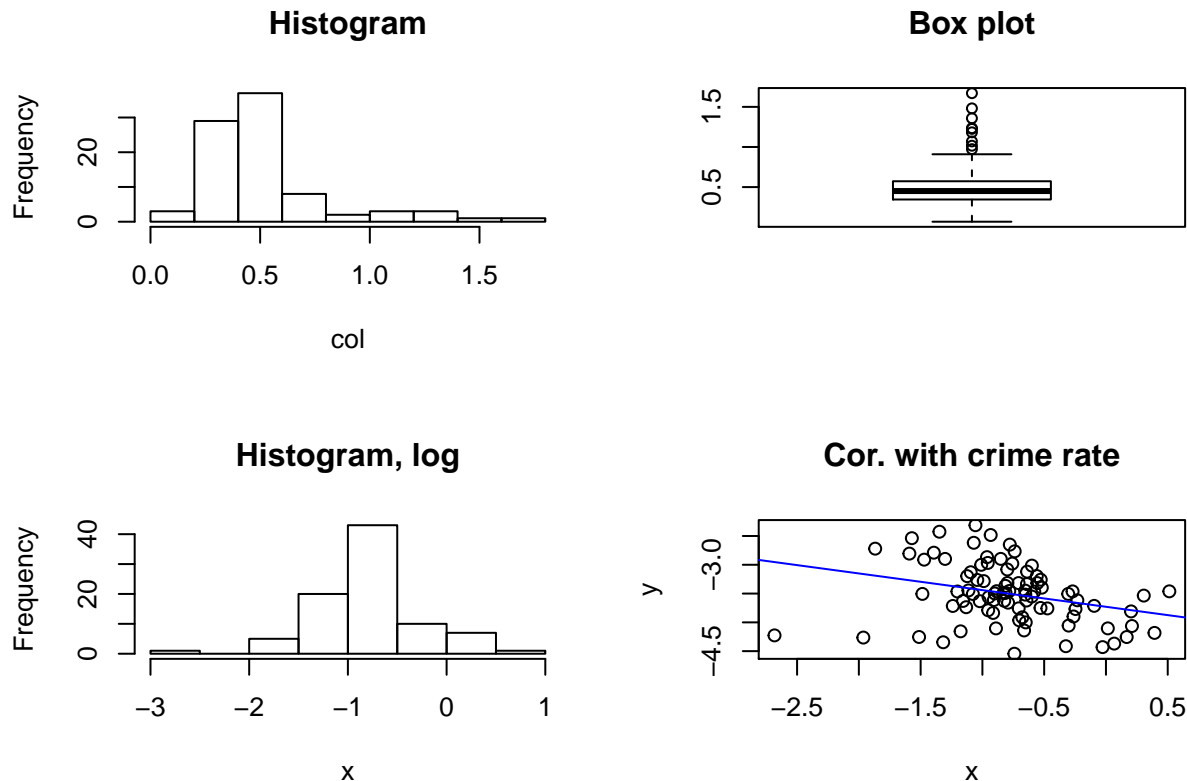
```
f_describe_col(crime_data$prbconv, do_log = TRUE)

##      Min. 1st Qu.  Median    Mean 3rd Qu.    Max.
## 0.06838 0.34541 0.45057 0.52446 0.57269 1.67052

## [1] -0.299
##
## Call:
## lm(formula = y ~ x)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -1.27319 -0.30410  0.00397  0.34472  1.11079
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept) -3.72938    0.09431  -39.55 < 2e-16 ***
## x           -0.28936    0.10012   -2.89  0.00489 **
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
```



```
##
## Residual standard error: 0.4903 on 85 degrees of freedom
## Multiple R-squared:  0.08948,    Adjusted R-squared:  0.07877
## F-statistic: 8.353 on 1 and 85 DF,  p-value: 0.004885
```



```
crime_data$log_prbconv = log(crime_data$prbconv)
```

This variable has quite a bit of left skew. It also has many outliers after the 3rd quartile. There are a few beyond 1 as well. Again, this is because we are not looking at a real probability but a ratio of convictions to arrests. It is possible, although perhaps uncommon, that a suspect is arrested once but convicted on multiple charges.

Taking a log transform improves the skew, although the spread is still quite a bit. There are no outliers with large influence as measured by Cook's distance.

There is moderate negative correlation with crime rate of -0.3. As convictions go up, crime rate goes down. Since we have already considered *prbarr*, let us check if *prbconv* has high correlation with *prbarr*:

```
print(cor(crime_data$prbarr, crime_data$prbconv))
```

```
## [1] -0.296224
```

```
print(cor(crime_data$prbarr, crime_data$log_prbconv))
```

```
## [1] -0.2855235
```

Not much. We will include $\log_{prbconv}$ in our model.

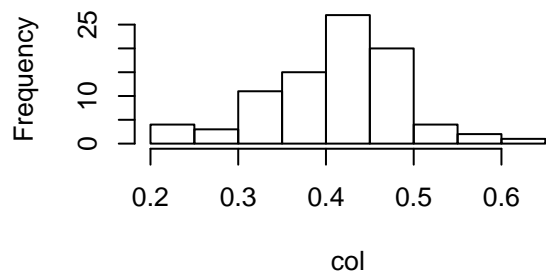
Probability of prison sentence

```
f_describe_col(crime_data$prbpris)
```

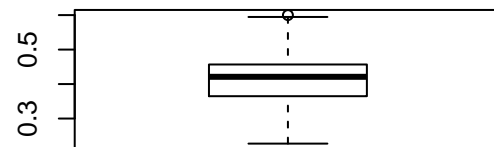
```
##      Min. 1st Qu.  Median    Mean 3rd Qu.    Max.
## 0.2273 0.3648 0.4211 0.4122 0.4568 0.6000

## [1] 0.0206
##
## Call:
## lm(formula = y ~ x)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -1.04583 -0.28459 -0.00612  0.30505  1.17991
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept)  -3.5599     0.3039  -11.71  <2e-16 ***
## x              0.1377     0.7250   0.19    0.85
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 0.5137 on 85 degrees of freedom
## Multiple R-squared:  0.0004245, Adjusted R-squared:  -0.01134
## F-statistic: 0.0361 on 1 and 85 DF, p-value: 0.8498
```

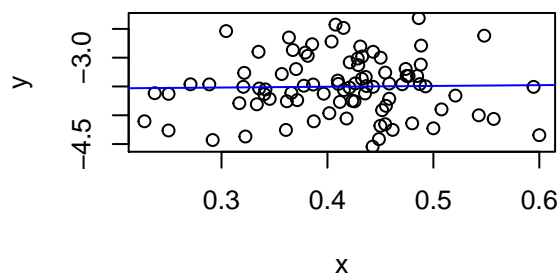
Histogram



Box plot



Cor. with crime rate



This histogram plot looks fairly normal and we don't observe any outliers. However, correlation is almost nonexistent wrt crime rate.

We will *not* consider this variable in our model.

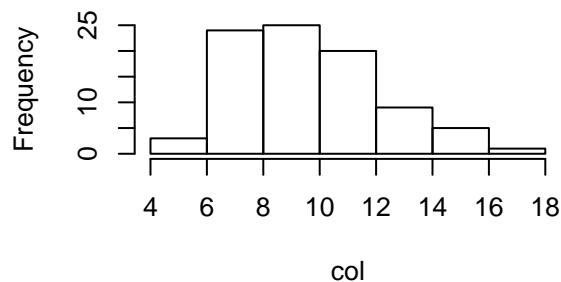
Average sentence duration

```
f_describe_col(crime_data$avgsen, do_log = TRUE)
```

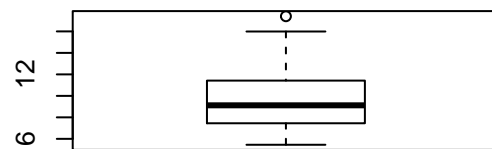
```
##      Min. 1st Qu.  Median      Mean 3rd Qu.      Max.
##    5.450  7.450   9.120   9.647  11.420  17.410

## [1] 0.0741
##
## Call:
## lm(formula = y ~ x)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -1.02352 -0.29899 -0.00019  0.29845  1.22870
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept)  -3.8245     0.4725  -8.095 3.65e-12 ***
## x              0.1440     0.2102   0.685  0.495
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 0.5124 on 85 degrees of freedom
## Multiple R-squared:  0.005488,    Adjusted R-squared:  -0.006212
## F-statistic: 0.469 on 1 and 85 DF,  p-value: 0.4953
```

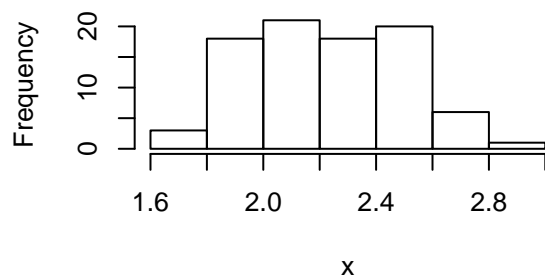
Histogram



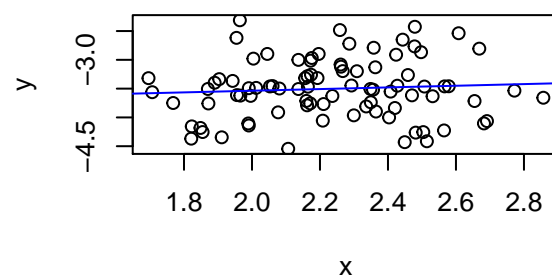
Box plot



Histogram, log



Cor. with crime rate



The average sentence in days looks slightly positive skewed, which we can correct with a log transform. But correlation is absent with respect to crime rate. It is interesting because we would expect that longer sentences would deter crime.

Perhaps we can use this data to make a policy recommendation to reduce sentences over long periods of time, or to be more lenient in pardoning criminals already serving long sentences.

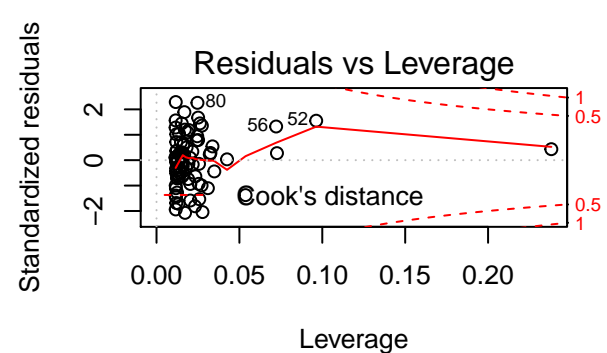
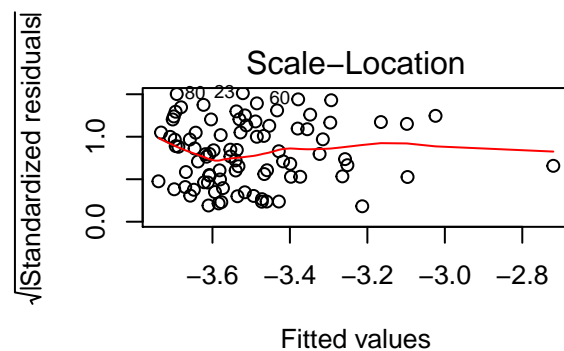
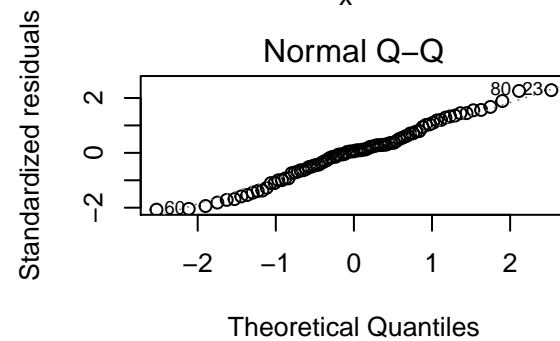
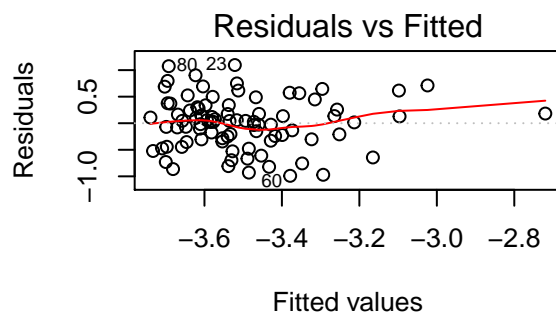
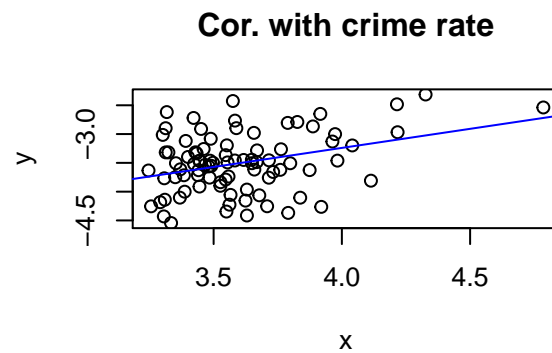
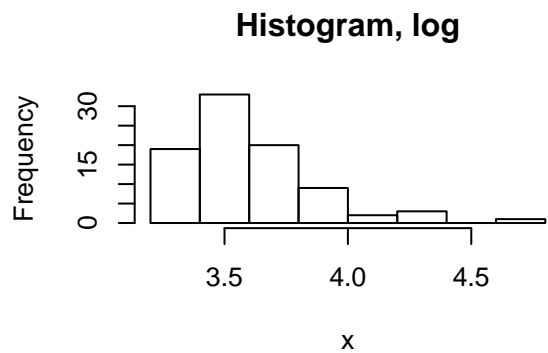
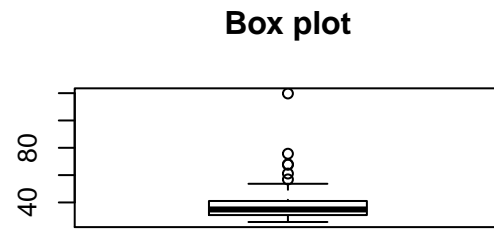
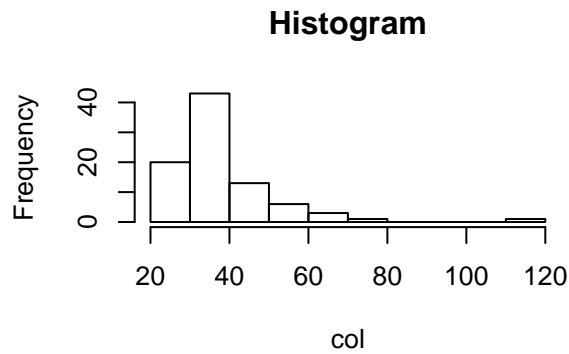
We will *not* consider this variable in our model.

Tax revenue per capita

```
f_describe_col(crime_data$taxpc, do_log = TRUE, plot_model = TRUE)
```

```
##      Min. 1st Qu.  Median    Mean 3rd Qu.    Max.
##    25.69  30.77   34.87   38.25  41.08  119.76

## [1] 0.347
##
## Call:
## lm(formula = y ~ x)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -0.99014 -0.30321  0.02651  0.29081  1.09548
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept)  -5.8892     0.7009  -8.403 8.72e-13 ***
## x              0.6623     0.1940   3.414 0.000983 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 0.4818 on 85 degrees of freedom
## Multiple R-squared:  0.1206, Adjusted R-squared:  0.1102
## F-statistic: 11.65 on 1 and 85 DF,  p-value: 0.0009833
```



Tax revenue also shows positive skew, with one outlier indicating high tax revenue per capita (>100). It does not show a lot of leverage, however, so we will keep the value.

We also see considerable positive correlation with crime rate. It may be that tax revenue is a proxy for wealth, and high amount of wealth attracts crime.

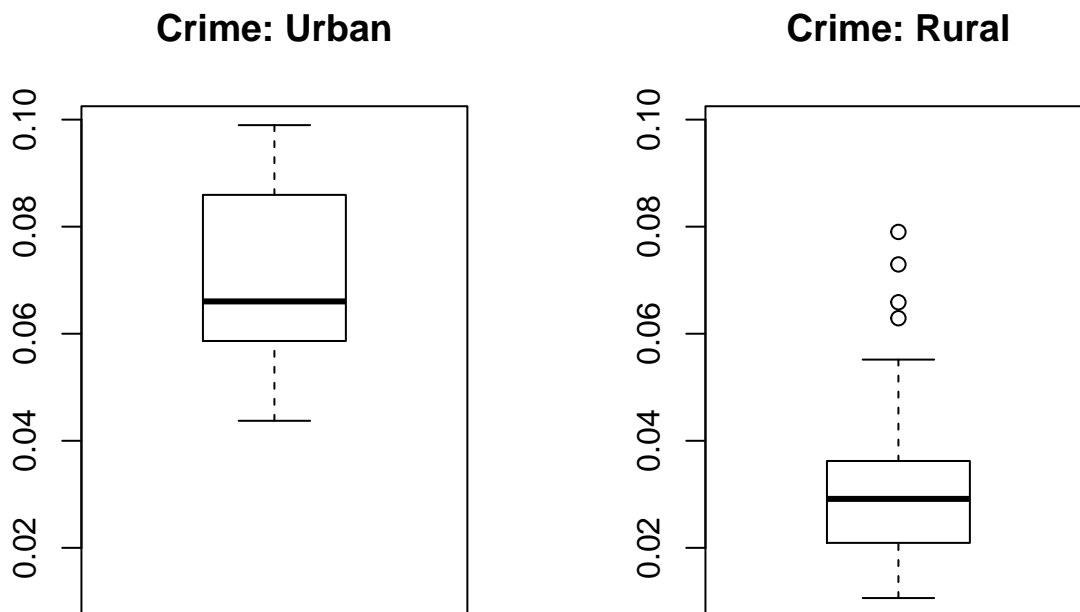
We will *not* include this variable in a first model.

Urban population

```
print(length(crime_data$urban[crime_data$urban == 1]))

## [1] 8

urban_crime_data = crime_data %>% filter(urban == 1) %>% dplyr::select(-urban)
rural_crime_data = crime_data %>% filter(urban == 0) %>% dplyr::select(-urban)
par(mfrow = c(1, 2))
lmts = range(urban_crime_data$crmrte, rural_crime_data$crmrte)
boxplot(urban_crime_data$crmrte, main = "Crime: Urban", ylim = lmts)
boxplot(rural_crime_data$crmrte, main = "Crime: Rural", ylim = lmts)
```



It is worth noting that there are only 8 observations classified urban in this dataset. Median crime rate in urban regions is double that of rural regions.

Let us fit a model and see if our variable is salient.

```
print(cor(crime_data$urban, log(crime_data$crmrte)))

## [1] 0.513772

m = lm(log(crmrte) ~ factor(urban), data = crime_data)
print(summary(m))

##
## Call:
## lm(formula = log(crmrte) ~ factor(urban), data = crime_data)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -0.95860 -0.23686  0.03449  0.26254  1.04801
##
## Coefficients:
```

```
##               Estimate Std. Error t value Pr(>|t|)
## (Intercept)   -3.5861     0.0496 -72.307 < 2e-16 ***
## factor(urban)1  0.9030     0.1636   5.521 3.59e-07 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 0.4408 on 85 degrees of freedom
## Multiple R-squared:  0.264, Adjusted R-squared:  0.2553
## F-statistic: 30.48 on 1 and 85 DF, p-value: 3.593e-07
```

We do see a strong correlation between observations classified “urban” and crime rate, and the same is reflected by the low p-value in the model summary.

Let us check if there is correlation between “urban” and “density”:

```
cor(crime_data$density, crime_data$urban)
```

```
## [1] 0.8218822
```

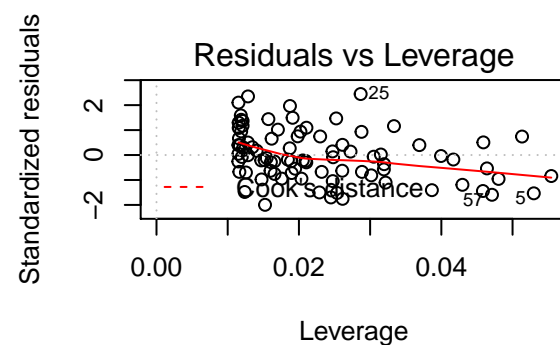
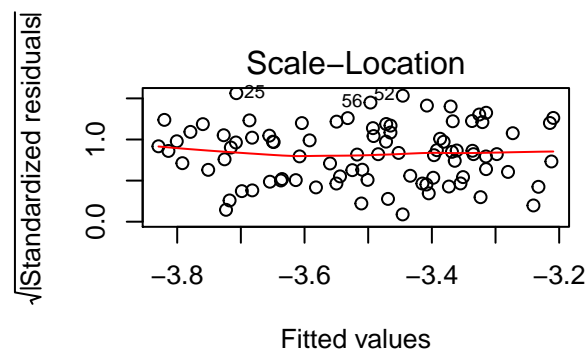
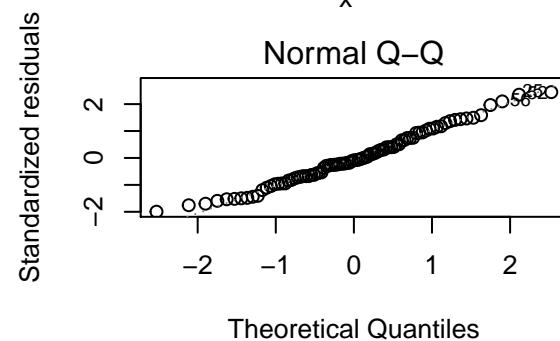
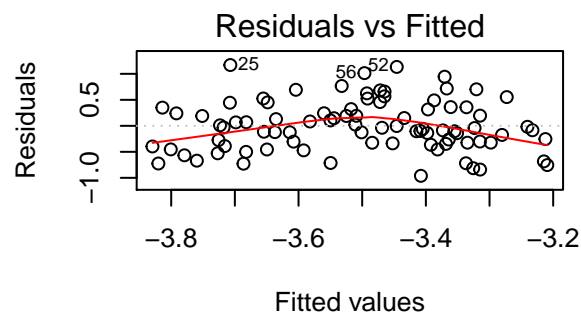
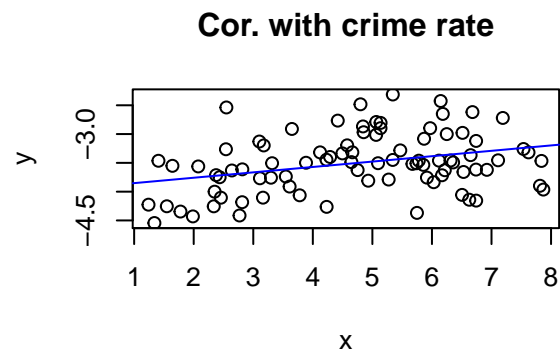
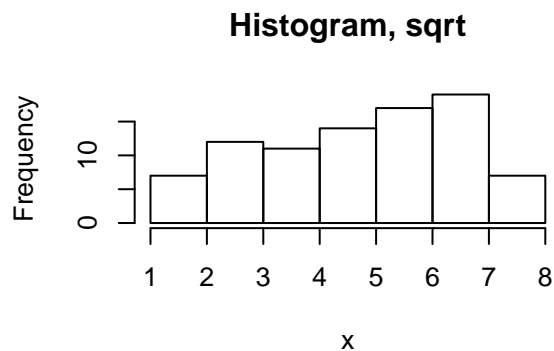
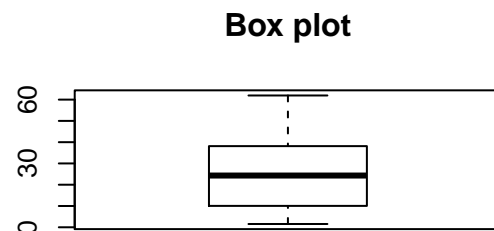
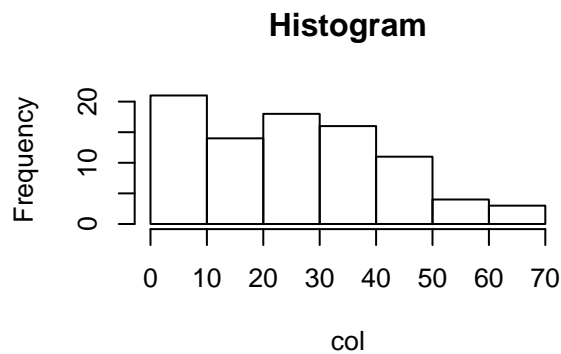
This is quite high, so we run a risk of multicollinearity.

Therefore, and since we have already selected density (with an additional advantage of more number of observations), we will *not* include this variable in our model.

Percent minority

```
f_describe_col(crime_data$pctmin80, plot_model = TRUE, do_sqrt = TRUE)
```

```
##      Min. 1st Qu.  Median    Mean 3rd Qu.    Max.
##    1.541 10.044  24.312  25.553  38.142  61.942
## [1] 0.329
##
## Call:
## lm(formula = y ~ x)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -0.96104 -0.33142 -0.04262  0.33590  1.16900
##
## Coefficients:
##               Estimate Std. Error t value Pr(>|t|)
## (Intercept) -3.94553     0.14741 -26.765 < 2e-16 ***
## x            0.09355     0.02916   3.208  0.00189 **
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 0.4853 on 85 degrees of freedom
## Multiple R-squared:  0.108, Adjusted R-squared:  0.0975
## F-statistic: 10.29 on 1 and 85 DF, p-value: 0.001886
```



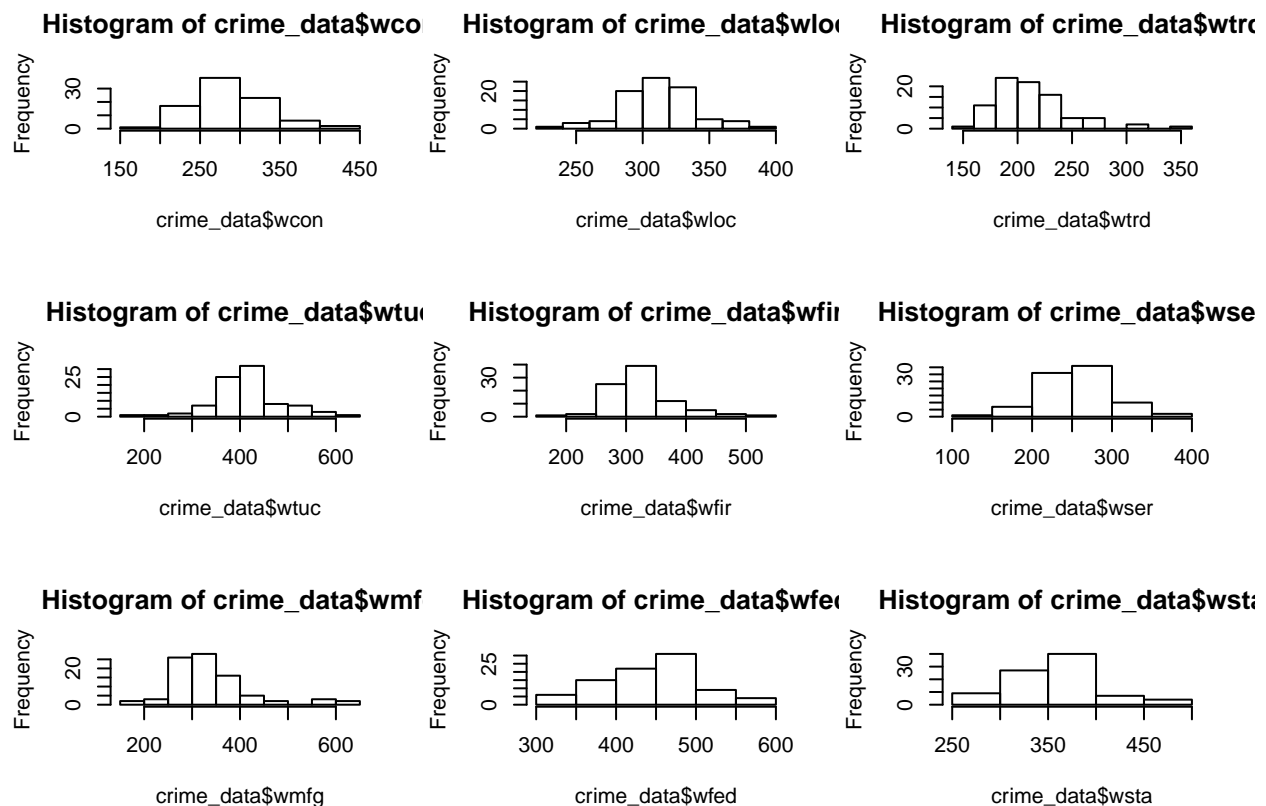
```
crime_data$sqrt_pctmin80 = sqrt(crime_data$pctmin80)
```

Minority percentage has positive skew, but no outliers. Taking square root reshapes the distribution nicely. There is a fair amount of positive correlation with crime rate (0.27). It may be that as minorities increase, there is loss of social homogeneity and/or hate crime.

We will include this (transformed) variable in our model.

Wage distribution

```
par(mfrow = c(3, 3))
hist(crime_data$wcon)
hist(crime_data$wloc)
hist(crime_data$wtrd)
hist(crime_data$wtuc)
hist(crime_data$wfir)
hist(crime_data$wser)
hist(crime_data$wmfg)
hist(crime_data$wfed)
hist(crime_data$wsta)
```



Most of the wage variables conform to normal distributions. We do not have to worry about transformations.

Let us look which of them have high correlation with crime rate, considering all those with $R > 0.25$ (arbitrarily).

```
wage_cols = c("wcon", "wloc", "wtrd", "wtuc", "wfir", "wser", "wmfg", "wfed", "wsta")
cor(log(crime_data$crmrte), crime_data[, wage_cols])
```

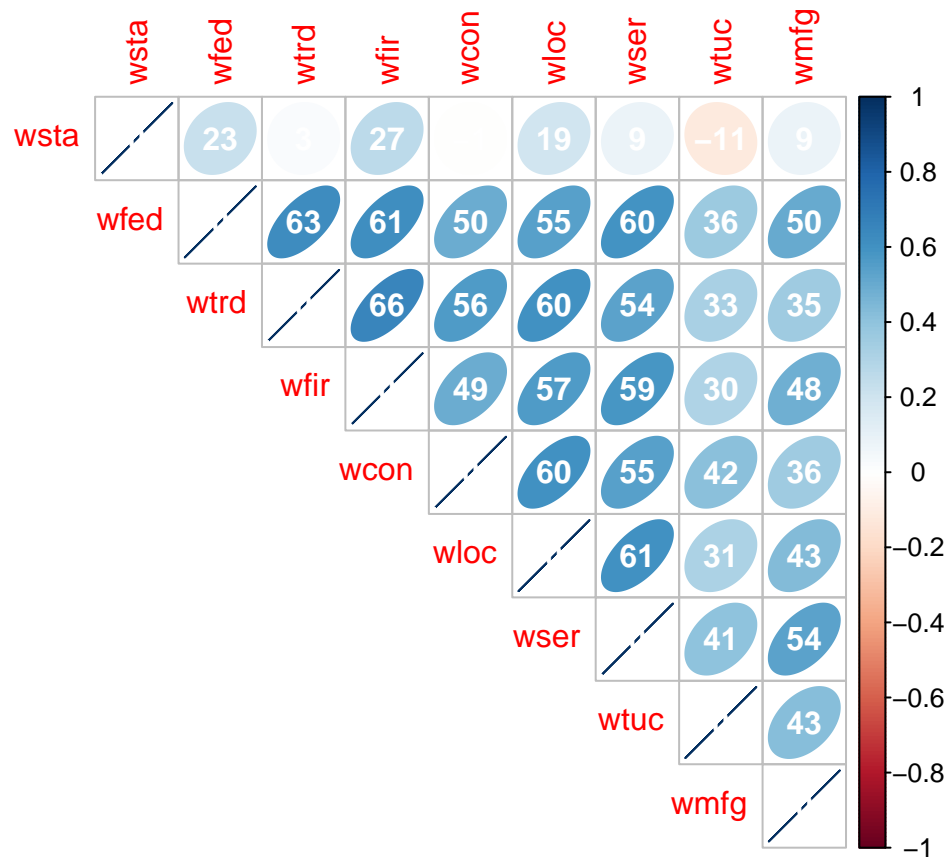
```
##           wcon      wloc      wtrd      wtuc      wfir      wser      wmfg
## [1,] 0.3296627 0.4173795 0.396715 0.2092351 0.2975497 0.3528975 0.3441131
##           wfed      wsta
## [1,] 0.5306463 0.2000234
```

This eliminates wsta and wtuc, but we are still left with 7 categories.

- wfed (0.50)
- wcon (0.37)
- wtrd (0.37)
- wser (0.34)
- wloc (0.28)
- wmfg (0.28)
- wfir (0.27)

As a different approach, let us check if the wages have high correlation among them. This will allow us to eliminate possible multi-collinearity.

```
corrplot(cor(crime_data[, wage_cols]), type = "upper", diag = TRUE, addCoef.col = "white",
         addCoefasPercent = TRUE, order = "hclust", method = "ellipse")
```



Indeed, a lot of the wage categories above have a high degree of correlation among them, but all are less than 70. We cannot eliminate any wage categories this way.

As a third approach, let us check for variance inflation instead:

```
m = lm(log(crmrte) ~ wfed + wcon + wtrd + wser + wloc + wmfg + wfir + wsta + wtuc,
       data = crime_data)
print(vif(m))
```

```
##      wfed      wcon      wtrd      wser      wloc      wmfg      wfir      wsta
## 2.304902 1.961994 2.541301 2.248561 2.247945 1.692487 2.452187 1.274484
##      wtuc
## 1.431049
```

Again, no VIF is above 5. This procedure also does not eliminate any wage categories.

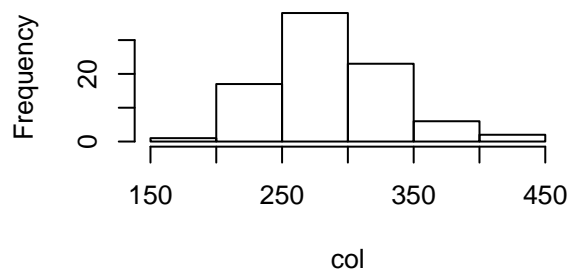
We can take a call, and choose *wcon* as a proxy for a first model. We can include other wages in a second model.

```
f_describe_col(crime_data$wcon, plot_model = TRUE)
```

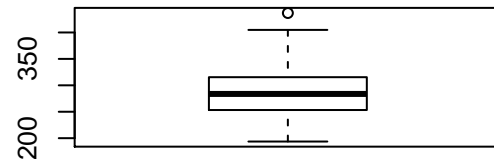
```
##      Min. 1st Qu.  Median    Mean 3rd Qu.    Max.
##    193.6  253.3   283.7   287.6   315.4   436.8

## [1] 0.33
##
## Call:
## lm(formula = y ~ x)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -1.11943 -0.29602 -0.04366  0.29642  1.09103
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept) -4.534364   0.324539  -13.972  < 2e-16 ***
## x              0.003586   0.001114   3.219   0.00182 **
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 0.4851 on 85 degrees of freedom
## Multiple R-squared:  0.1087, Adjusted R-squared:  0.09819
## F-statistic: 10.36 on 1 and 85 DF,  p-value: 0.00182
```

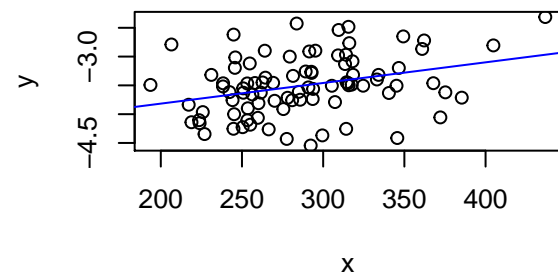
Histogram



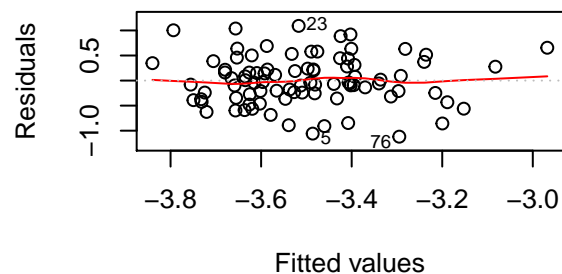
Box plot

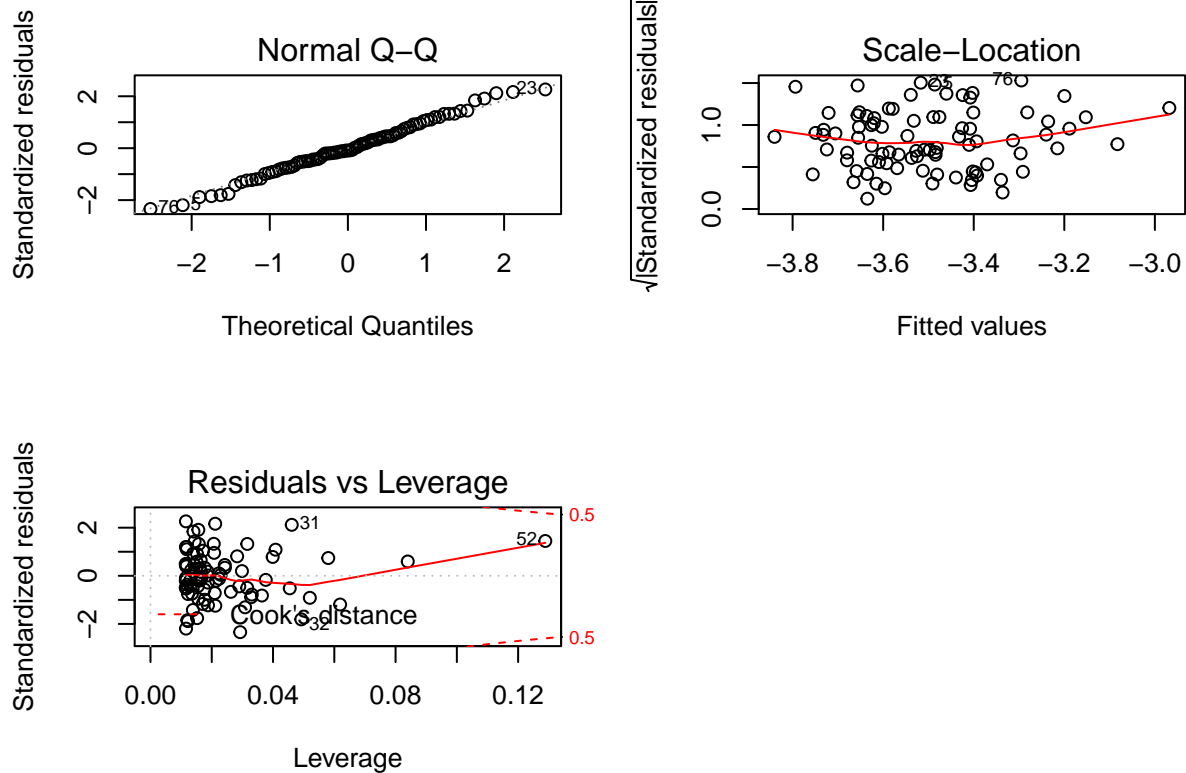


Cor. with crime rate



Residuals vs Fitted



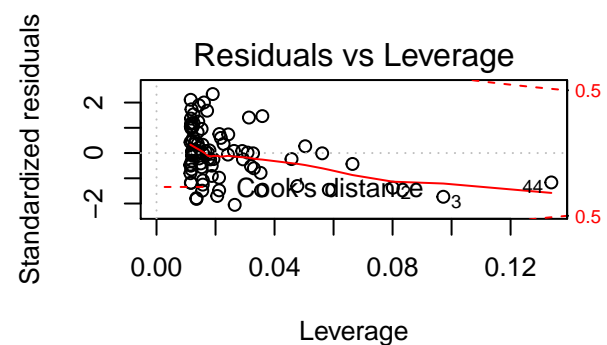
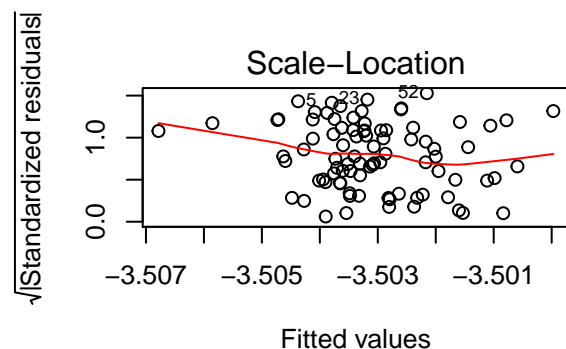
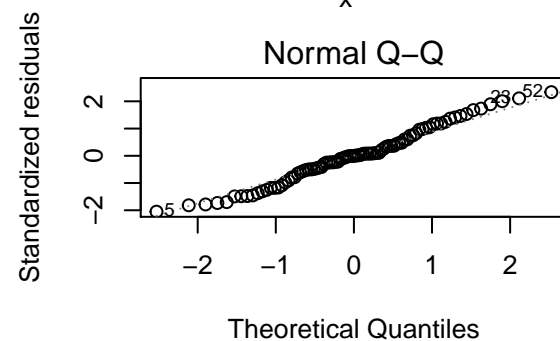
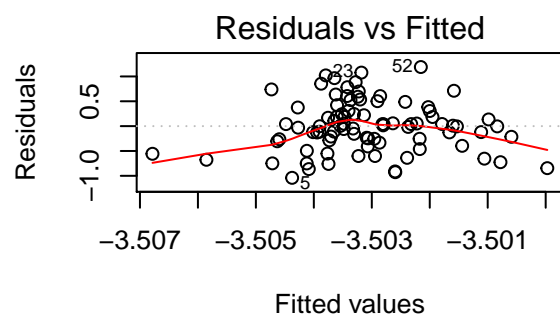
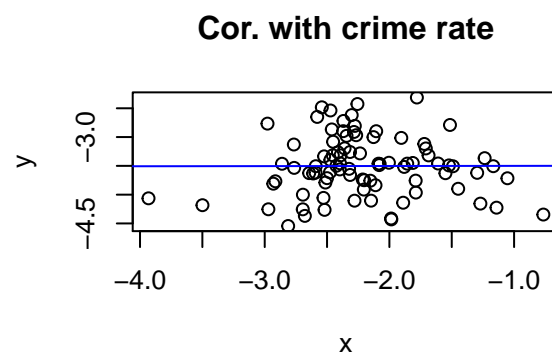
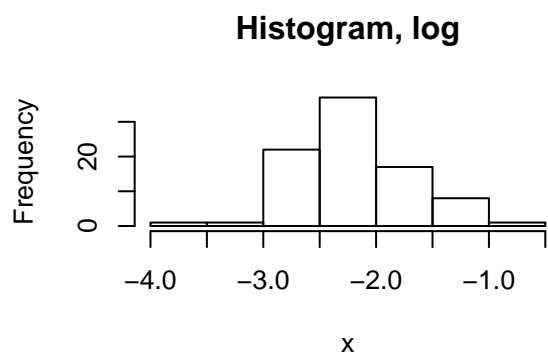
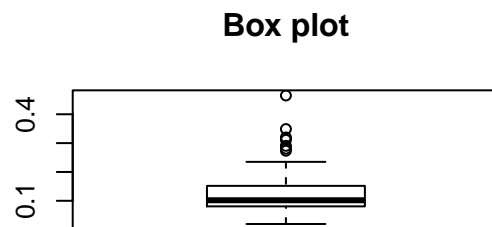
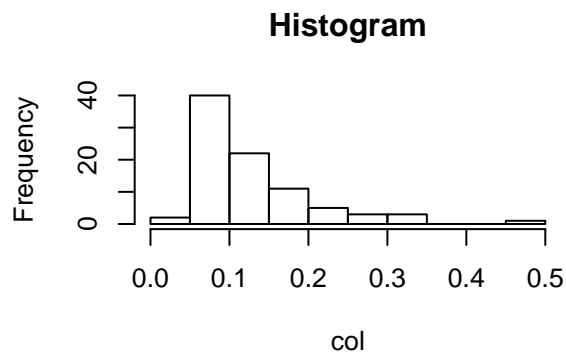


Offense Mix

```
f_describe_col(crime_data$mix, do_log = TRUE, plot_model = TRUE)

##      Min. 1st Qu.  Median    Mean 3rd Qu.    Max.
## 0.01961 0.08073 0.10186 0.12695 0.15175 0.46512

## [1] 0.00224
##
## Call:
## lm(formula = y ~ x)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -1.04034 -0.29549 -0.00194  0.30604  1.18917
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept) -3.498325   0.236725 -14.778  <2e-16 ***
## x             0.002151   0.104257   0.021   0.984
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 0.5138 on 85 degrees of freedom
## Multiple R-squared:  5.01e-06, Adjusted R-squared: -0.01176
## F-statistic: 0.0004259 on 1 and 85 DF, p-value: 0.9836
```



Offense mix does not seem to have any correlation with crime rate. The distribution is skewed, but a log transform fixes it. Outliers exist, but none have leverage as detected by Cook's distance.

We will *not* include offense mix in our models.

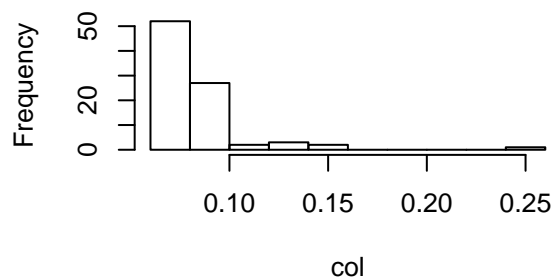
Percent of young males

```
f_describe_col(crime_data$pctymle, do_sqrt = TRUE, plot_model = TRUE)
```

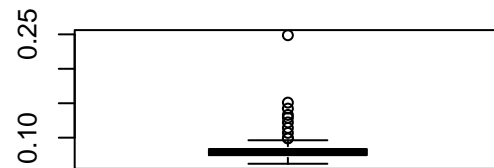
```
##      Min. 1st Qu.  Median    Mean 3rd Qu.    Max.
## 0.06216 0.07464 0.07787 0.08443 0.08355 0.24871

## [1] 0.281
##
## Call:
## lm(formula = y ~ x)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -0.94503 -0.29623  0.03136  0.28149  1.22093
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept)  -4.7301     0.4573  -10.344 < 2e-16 ***
## x              4.2514     1.5737   2.701  0.00833 **
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 0.4931 on 85 degrees of freedom
## Multiple R-squared:  0.07907,    Adjusted R-squared:  0.06823
## F-statistic: 7.298 on 1 and 85 DF,  p-value: 0.008332
```

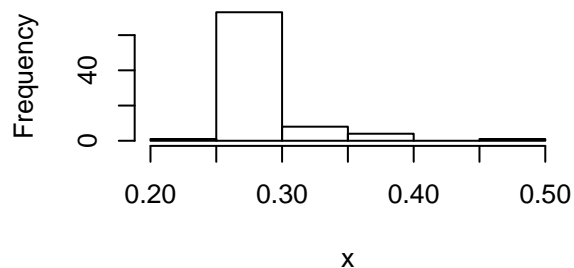
Histogram



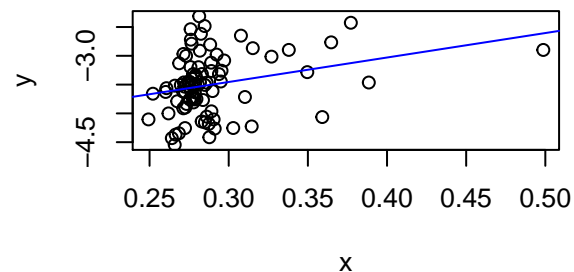
Box plot

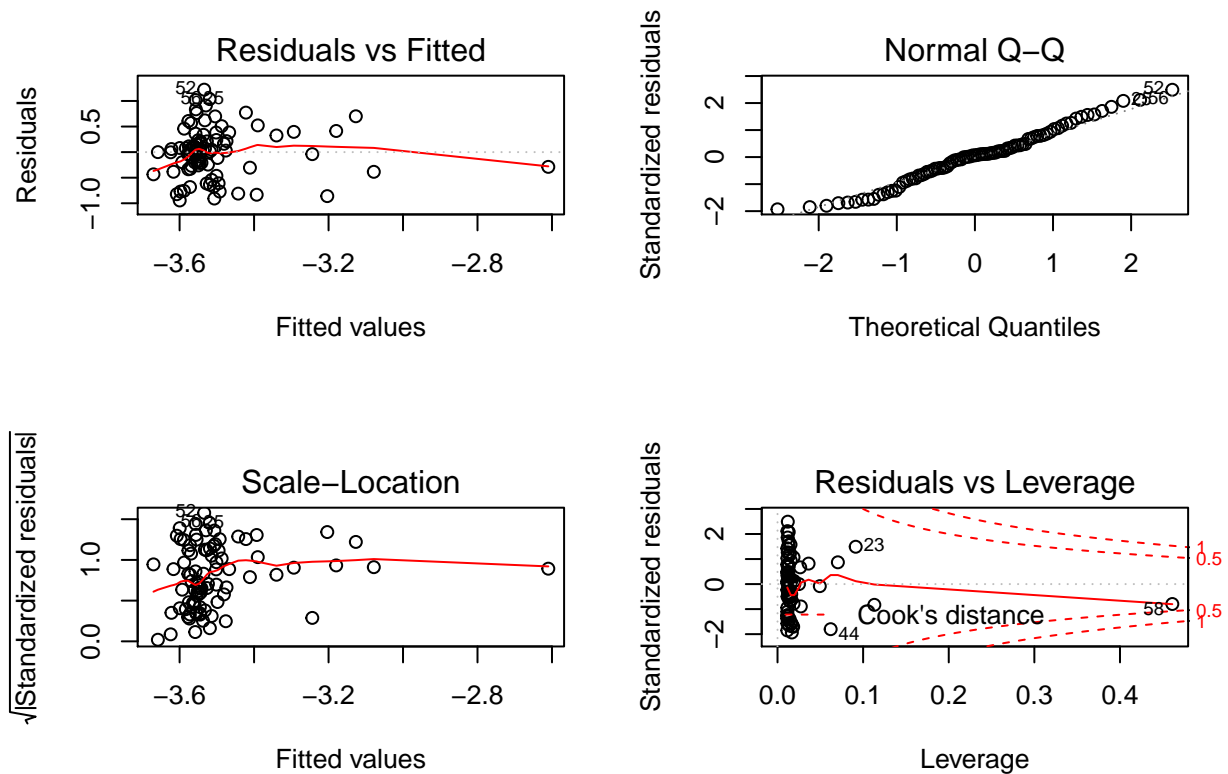


Histogram, sqrt



Cor. with crime rate





```
crime_data$sqrt_pctymle = sqrt(crime_data$pctymle)
```

We see moderate positive correlation with higher percentage of young males. There is positive skew, which we correct by taking a square root. Boxplot shows outliers, but none has outsized influence (Cook's distance < 0.5).

We will include this variable in our model.

Categorical variables

We have the following categorical variables in the dataset:

- Direction: west, central, other
- Urban or rural

We will use these to come up with separate models, based on different factors, later in this analysis.

Summary of variables

Here is a summary table of variables used in our models.

Variable	Transform?	Model1?	Model2?	Model3?	Remarks
county	N/A				Unused
year	N/A				Unused
prbarr		Y	Y	Y	
prbconv	log	Y	Y	Y	
prbpris				Y	No corr. found
avgsgen				Y	No corr. found

Variable	Transform?	Model1?	Model2?	Model3?	Remarks
polpc	log		Y	Y	Effect, not cause
density	log	Y	Y	Y	Causal
taxpc	log		Y	Y	Omit var: wealth
west	N/A				Categ, sep. model
central	N/A				Categ, sep. model
urban					Cor. with density
pctmin80	sqrt	Y	Y	Y	Causal
wcon		Y	Y	Y	Proxy for wages
wtuc				Y	Low corr. found
wtrd			Y	Y	
wfir			Y	Y	
wser			Y	Y	
wmfg			Y	Y	
wfed			Y	Y	
wsta				Y	Low corr. found
wloc			Y	Y	
mix	log			Y	No corr. found
pctymle	sqrt	Y	Y	Y	Causal

Data Transformation

Based on the univariate analysis performed above, we can opt to take the following transformations of the variables to make better analysis and judgement calls:

Log transformation of the Crime Rate, Police per Capita, Density per sq. mile, Tax revenue per capita And finally scaling the percent (percent young male) and probabilities (arrest, conviction and prison sentence) to be between 0-100

```
crime_data$log_crmrte = log(crime_data$crmrate)
crime_data$log_density = log(crime_data$density)
crime_data$log_polpc = log(crime_data$polpc)
crime_data$log_taxpc = log(crime_data$taxpc)
crime_data$adj_pctymle = crime_data$pctymle * 100
crime_data$adj_prbarr = crime_data$prbarr * 100
crime_data$adj_prbconv = crime_data$prbconv * 100
crime_data$adj_prbpris = crime_data$prbpris * 100
```

A final summary table of our dataset with all transformation and data cleansing performed is displayed below:

```
stargazer(crime_data, title = "Descriptive Statistics", digits = 1)
```

% Table created by stargazer v.5.2 by Marek Hlavac, Harvard University. E-mail: hlavac at fas.harvard.edu
 % Date and time: Sat, Mar 31, 2018 - 21:40:02

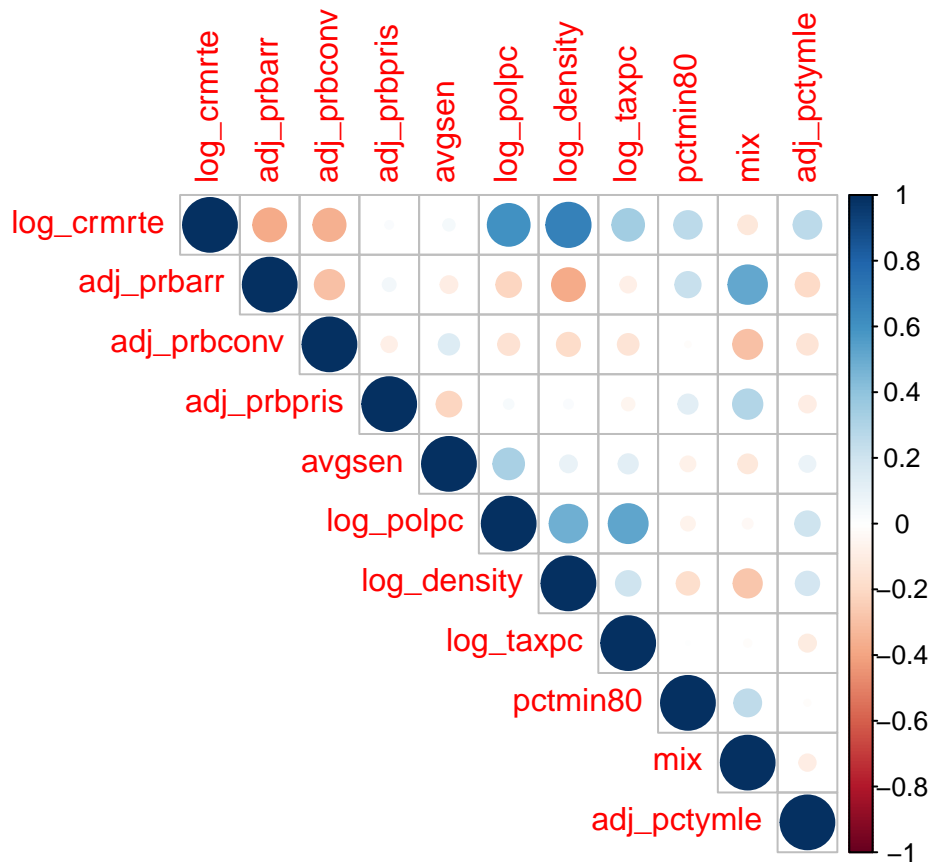
Table 2: Descriptive Statistics

Statistic	N	Mean	St. Dev.	Min	Max
-----------	---	------	----------	-----	-----

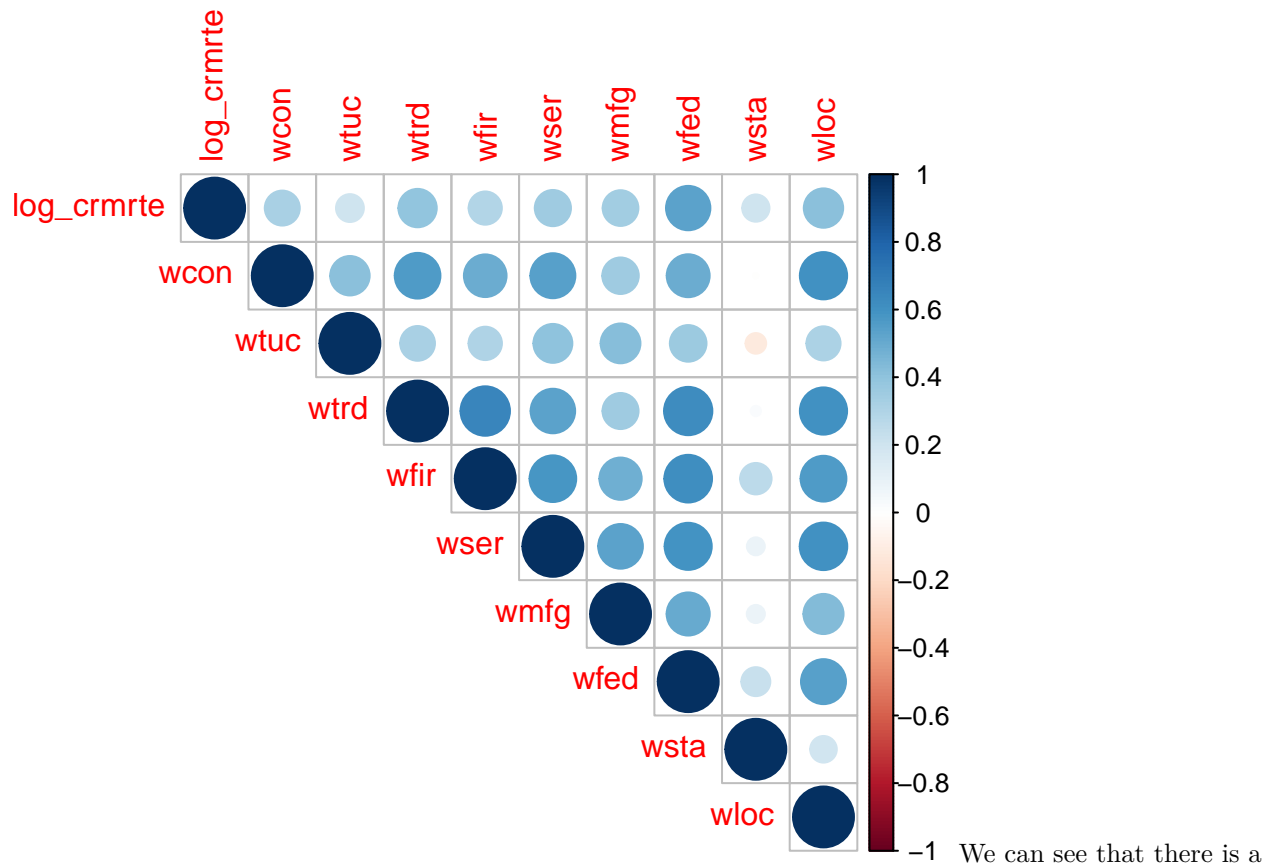
Bi-variate Analysis

The correlation plot between the different variables is as follows:


```
corrplot(cor(crime_data[, c("log_crmrte", "adj_prbarr", "adj_prbconv", "adj_prbpris",
"avgsen", "log_polpc", "log_density", "log_taxpc", "pctmin80", "mix", "adj_pctymle")]),
type = "upper")
```



```
corrplot(cor(crime_data[, c("log_crmrte", "wcon", "wtuc", "wtrd", "wfir", "wser",
"wmfg", "wfed", "wsta", "wloc")]), type = "upper")
```



high positive correlation between:

- log of crime rate vs. log of policy per capita, log of tax revenue per capita, log of density and percent young male
- log of crime rate vs. most of the wage variables

And there is a high negative correlation between:

- log of crime rate vs. probability of arrests and conviction

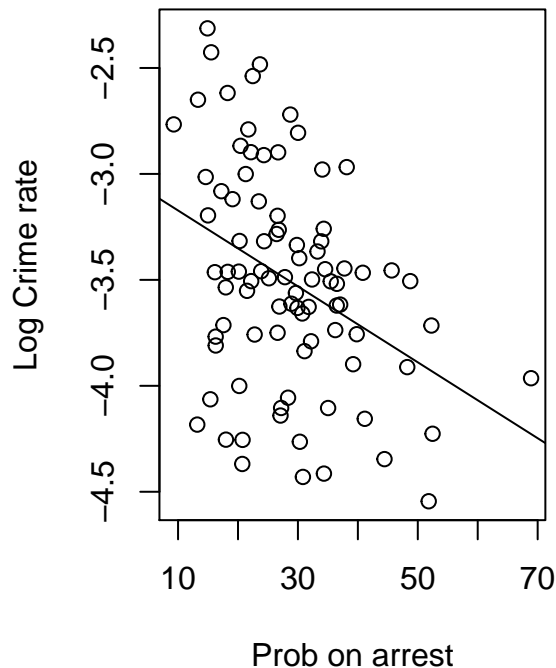
The positive correlation observed makes sense for the following reasons:

- 1) More densely populated regions tends to observe more crimes
- 2) More wealthy areas (more wages and taxes) tend to have more crimes
- 3) More crimes leads to more police presence in a particular county to monitor and reduce crime rate

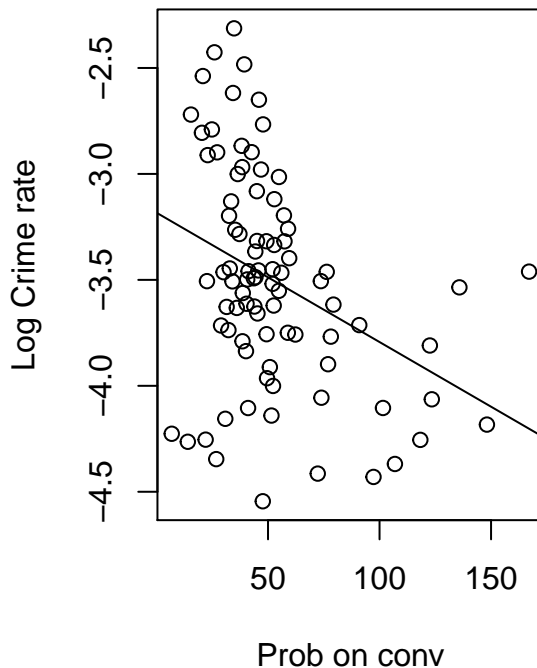
The negative correlations can be further observed using:

```
par(mfrow = c(1, 2))
plot(crime_data$adj_prbarr, crime_data$log_crmrte, main = "Probability of arrest",
     ylab = "Log Crime rate", xlab = "Prob on arrest")
abline(lm(crime_data$log_crmrte ~ crime_data$adj_prbarr))
plot(crime_data$adj_prbconv, crime_data$log_crmrte, main = "Probability of conviction vs. crime rate",
     ylab = "Log Crime rate", xlab = "Prob on conv")
abline(lm(crime_data$log_crmrte ~ crime_data$adj_prbconv))
```

Probability of arrest



Probability of conviction vs. crime



As seen above, as the probability of arrests and conviction go down, there are more criminals on the loose which leads to higher crime rates observed

TODO: Talk about other possible correlations here?

TODO: Discuss other interesting bi-variate analysis?

3. Model Specification and Assumptions

In our earlier analysis, we observed some key relationships between crime rate and other variables presented. Some of these variables had high positive correlation to crime rate while some others exhibited strong negative correlation.

For our first simple model, we will choose a subset of these variables that we believe are most important determinants of crime rate.

Model 1

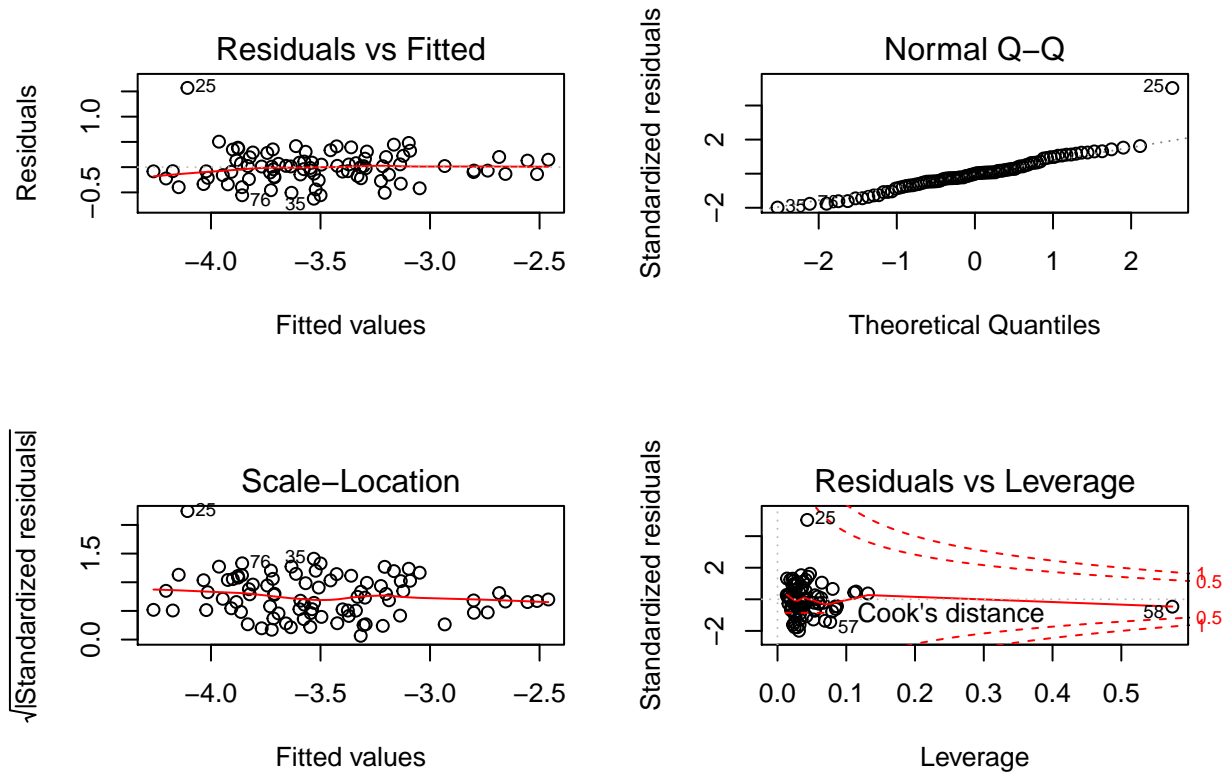
$$\log(\text{CrimeRate}) = \beta_0 + \beta_1 \log(\text{Density}) + \beta_2(\text{YoungMale}) + \beta_3(\text{Minority}) + u$$

It is common knowledge that areas with higher density have more crime. Therefore we include that factor in our model. Similarly we hypothesized that crime rate is high among minority and young male population, so we round off our model with that factored in as well.

```
model1 = lm(log(crmrte) ~ (log_density) + pctymle + pctmin80, data = crime_data)
model1$coefficients
```

```
## (Intercept) log_density    pctymle    pctmin80
## -4.08397802  0.48206791   2.91099248  0.01203591
```

```
par(mfrow = c(2, 2))
plot(model1)
```



```
AIC(model1)
```

```
## [1] 54.12082
```

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

```
## [1] 0.6230685
```

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

```
## [1] 0.6094445
```

Model 2

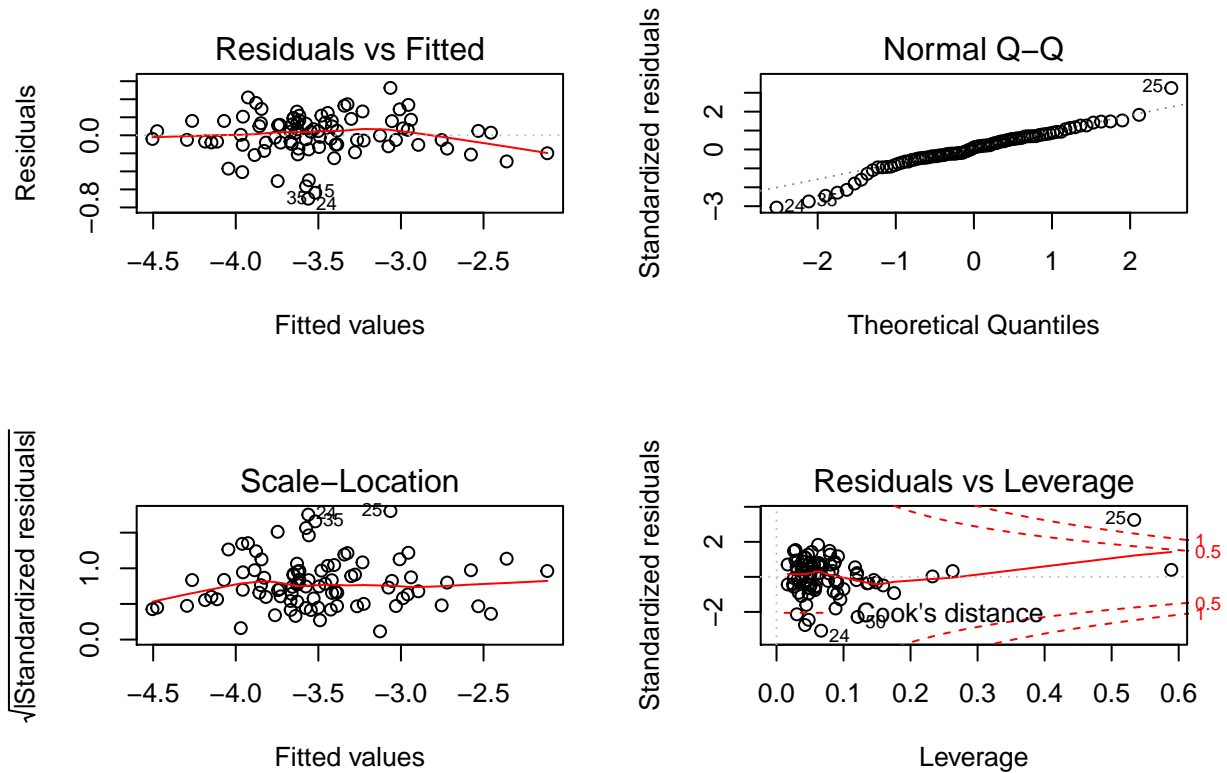
high probability of arrests and conviction act as deterrents to crime.

$$\log(\text{CrimeRate}) = \beta_0 + \beta_1 \log(\text{Density}) + \beta_2 (\text{YoungMale}) + \beta_3 (\text{Minority}) + \beta_4 (\text{Conviction}) + \beta_5 (\text{Arrest}) + \beta_6 (\text{Tax}) + u$$

```
model2 = lm(log(crmrte) ~ (log_density) + pctymle + pctmin80 + adj_prbarr + adj_prbconv +
  taxpc, data = crime_data)
model2$coefficients
```

```
## (Intercept) log_density    pctymle    pctmin80  adj_prbarr
## -3.624003297 0.344827949  1.888649001  0.013369932 -0.015748832
## adj_prbconv    taxpc
## -0.005324798  0.008561615
```

```
par(mfrow = c(2, 2))
plot(model12)
```



```
AIC(model12)
```

```
## [1] 4.918033
```

```
summary(model12)$r.squared
```

```
## [1] 0.8001537
```

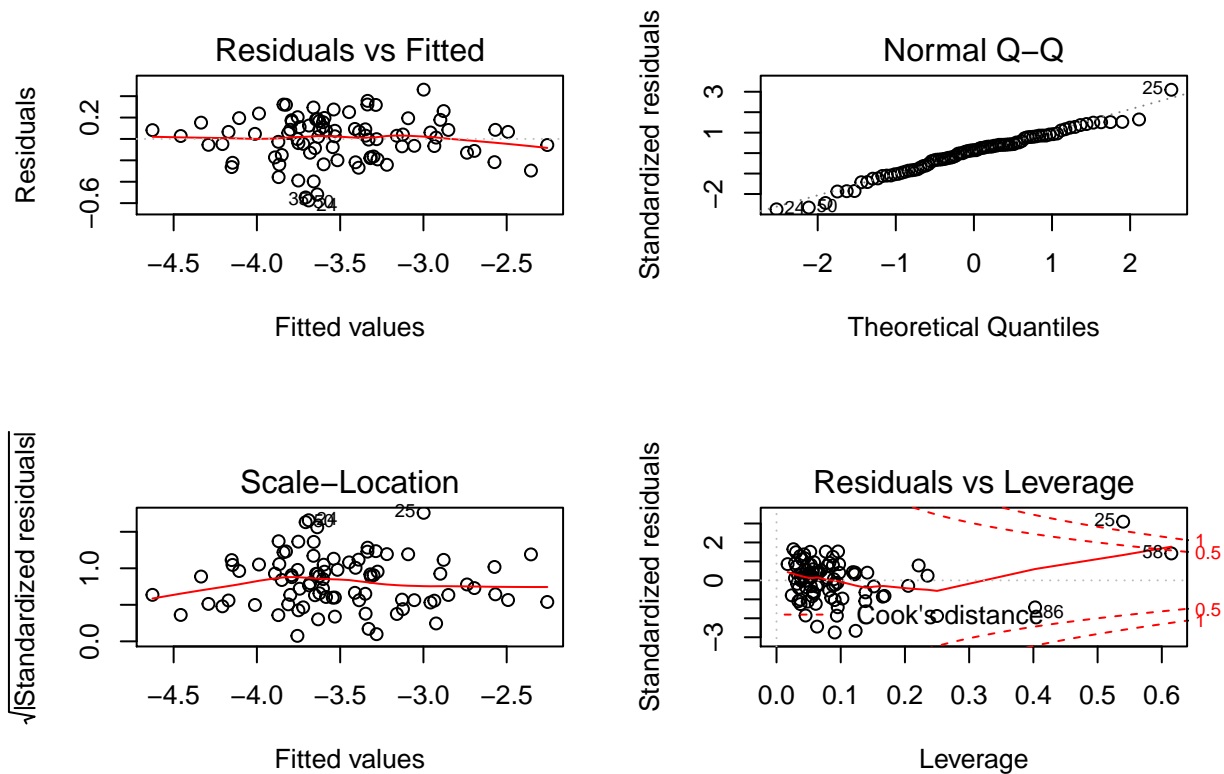
Model 3

```
everything
```

```
model3 = lm(log(crmrte) ~ (log_density) + pctymle + pctmin80 + adj_prbarr + adj_prbconv +
  taxpc + log_polpc, data = crime_data)
model3$coefficients
```

```
## (Intercept) log_density      pctymle      pctmin80  adj_prbarr
## -0.862515848  0.283874000  0.914663878  0.013211261 -0.015982962
## adj_prbconv      taxpc      log_polpc
## -0.005392071  0.004087450  0.384090018
```

```
par(mfrow = c(2, 2))
plot(model3)
```



```
AIC(model3)
```

```
## [1] -7.458543
```

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

```
## [1] 0.8305936
```

```
stargazer(model11, model12, model13)
```

```
% Table created by stargazer v.5.2 by Marek Hlavac, Harvard University. E-mail: hlavac at fas.harvard.edu
% Date and time: Sat, Mar 31, 2018 - 21:40:08
```

5. Discussion of omitted variables (Identify what you think are the 5-10 most important omitted variables that bias results you care about.)

Education

Unemployment

Poverty

Table 3:

	<i>Dependent variable:</i>		
	log(crmrte)		
	(1)	(2)	(3)
log_density	0.482*** (0.047)	0.345*** (0.039)	0.284*** (0.039)
pctymle	2.911* (1.476)	1.889 (1.142)	0.915 (1.089)
pctmin80	0.012*** (0.002)	0.013*** (0.002)	0.013*** (0.001)
adj_prbarr		-0.016*** (0.003)	-0.016*** (0.003)
adj_prbconv		-0.005*** (0.001)	-0.005*** (0.001)
taxpc		0.009*** (0.002)	0.004* (0.002)
log_polpc			0.384*** (0.102)
Constant	-4.084*** (0.139)	-3.624*** (0.201)	-0.863 (0.756)
Observations	87	87	87
R ²	0.623	0.800	0.831
Adjusted R ²	0.609	0.785	0.816
Residual Std. Error	0.319 (df = 83)	0.237 (df = 80)	0.219 (df = 79)
F Statistic	45.733*** (df = 3; 83)	53.385*** (df = 6; 80)	55.333*** (df = 7; 79)

Note:

*p<0.1; **p<0.05; ***p<0.01