

Data 621 - Homework 3

Group 4 Layla Quinones, Ian Costello, Dmitriy Burtsev & Esteban Aramayo

11/7/2021

Overview

General Objective

For this assignment, we will be exploring, analyzing, and modeling data related to crime statistics for various areas of a major U.S. city. The primary objective is to understand how, or if, variable indicate whether crime in a particular area will be above or below the median crime rate for the entire city. The models will be binary logistic regression using combinations or constructions of the variables provided.

About the Data

- `zn`: proportion of residential land zoned for large lots (over 25000 square feet) (predictor variable)
- `indus`: proportion of non-retail business acres per suburb (predictor variable)
- `chas`: a dummy var. for whether the suburb borders the Charles River (1) or not (0) (predictor variable)
- `nox`: nitrogen oxides concentration (parts per 10 million) (predictor variable)
- `rm`: average number of rooms per dwelling (predictor variable)
- `age`: proportion of owner-occupied units built prior to 1940 (predictor variable)
- `dis`: weighted mean of distances to five Boston employment centers (predictor variable)
- `rad`: index of accessibility to radial highways (predictor variable)
- `tax`: full-value property-tax rate per \$10,000 (predictor variable)
- `ptratio`: pupil-teacher ratio by town (predictor variable)
- `lstat`: lower status of the population (percent) (predictor variable)
- `medv`: median value of owner-occupied homes in \$1000s (predictor variable)
- `target`: whether the crime rate is above the median crime rate (1) or not (0) (response variable)

Libraries Used

We use pretty standard packages for this assignment, including the ever-useful `tidyverse`, `ggplot2`, and `caret`. New additions for this assignment include `VIM`, `DataExplorer`, and `broom`. (**Code Appendix 1.1**)

Data Exploration

As usual, our data are stored on GitHub at our team's main repository for easy access across team members (**Code Appendix 2.2**). The variables are data type doubles, except for two variables: `chas` and `target`. While these are integer data types, they will be treated as categorical factors. Taking a peek into the data, we can get a sense of the distribution and structure of the data set (**Code Appendix 2.3**). More visuals will be required to look deeper, however.

	vars	n	mean	sd	median	trimmed	mad	min	max	range	skew	kurtosis	se
zn	1	466	11.5772532	23.3646511	0.00000	5.3542781	0.0000000	0.0000	100.0000	100.0000	2.1768152	3.8135765	1.0823466
indus	2	466	11.1050215	6.8458549	9.69000	10.9082353	9.3403800	0.4600	27.7400	27.2800	0.2885450	-1.2432132	0.3171281
chas	3	466	0.0708155	0.2567920	0.00000	0.0000000	0.0000000	0.0000	1.0000	1.0000	3.3354899	9.1451313	0.0118957
nox	4	466	0.5543105	0.1166667	0.53800	0.5442684	0.1334340	0.3890	0.8710	0.4820	0.7463281	-0.0357736	0.0054045
rm	5	466	6.2906738	0.7048513	6.21000	6.2570615	0.5166861	3.8630	8.7800	4.9170	0.4793202	1.5424378	0.0326516
age	6	466	68.3675966	28.3213784	77.15000	70.9553476	30.0226500	2.9000	100.0000	97.1000	-0.5777075	-1.0098814	1.3119625
dis	7	466	3.7956929	2.1069496	3.19095	3.5443647	1.9144814	1.1296	12.1265	10.9969	0.9988926	0.4719679	0.0976026
rad	8	466	9.5300429	8.6859272	5.00000	8.6978610	1.4826000	1.0000	24.0000	23.0000	1.0102788	-0.8619110	0.4023678
tax	9	466	409.5021459	167.9000887	334.50000	401.5080214	104.5233000	187.0000	711.0000	524.0000	0.6593136	-1.1480456	7.7778214
ptratio	10	466	18.3984979	2.1968447	18.90000	18.5970588	1.9273800	12.6000	22.0000	9.4000	-0.7542681	-0.4003627	0.1017669
lstat	11	466	12.6314592	7.1018907	11.35000	11.8809626	7.0720020	1.7300	37.9700	36.2400	0.9055864	0.5033688	0.3289887
medv	12	466	22.5892704	9.2396814	21.20000	21.6304813	6.0045300	5.0000	50.0000	45.0000	1.0766920	1.3737825	0.4280200
target	13	466	0.4914163	0.5004636	0.00000	0.4893048	0.0000000	0.0000	1.0000	1.0000	0.0342293	-2.0031131	0.0231835

```
## Rows: 466
## Columns: 13
## $ zn      <dbl> 0, 0, 0, 30, 0, 0, 0, 0, 0, 80, 22, 0, 0, 22, 0, 0, 100, 20, 0~
## $ indus   <dbl> 19.58, 19.58, 18.10, 4.93, 2.46, 8.56, 18.10, 18.10, 5.19, 3.6~
## $ chas    <int> 0, 1, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, ~
## $ nox     <dbl> 0.605, 0.871, 0.740, 0.428, 0.488, 0.520, 0.693, 0.693, 0.515, ~
## $ rm      <dbl> 7.929, 5.403, 6.485, 6.393, 7.155, 6.781, 5.453, 4.519, 6.316, ~
## $ age     <dbl> 96.2, 100.0, 100.0, 7.8, 92.2, 71.3, 100.0, 100.0, 38.1, 19.1, ~
## $ dis     <dbl> 2.0459, 1.3216, 1.9784, 7.0355, 2.7006, 2.8561, 1.4896, 1.6582~
## $ rad     <int> 5, 5, 24, 6, 3, 5, 24, 24, 5, 1, 7, 5, 24, 7, 3, 3, 5, 5, 24, ~
## $ tax     <int> 403, 403, 666, 300, 193, 384, 666, 666, 224, 315, 330, 398, 66~
## $ ptratio <dbl> 14.7, 14.7, 20.2, 16.6, 17.8, 20.9, 20.2, 20.2, 20.2, 16.4, 19~
## $ lstat   <dbl> 3.70, 26.82, 18.85, 5.19, 4.82, 7.67, 30.59, 36.98, 5.68, 9.25~
## $ medv    <dbl> 50.0, 13.4, 15.4, 23.7, 37.9, 26.5, 5.0, 7.0, 22.2, 20.9, 24.8~
## $ target  <int> 1, 1, 1, 0, 0, 0, 1, 1, 0, 0, 0, 0, 1, 1, 0, 0, 1, 1, 1, 0, ~
```

Table 1. Glimpse of data structure (Code Appendix 2.4)

Missing Data Checks

Following standard procedure, we check for any missing data in the set (**Code Appendix 2.5**). It appears there are no missing values as the following figures demonstrate. Again, we can see that most of the variables appear to be continuous. From the description of the predictors in the overview section of this document, we know that some of them can be treated as discrete and/or categorical. No columns with missing values were detected and all rows are complete.

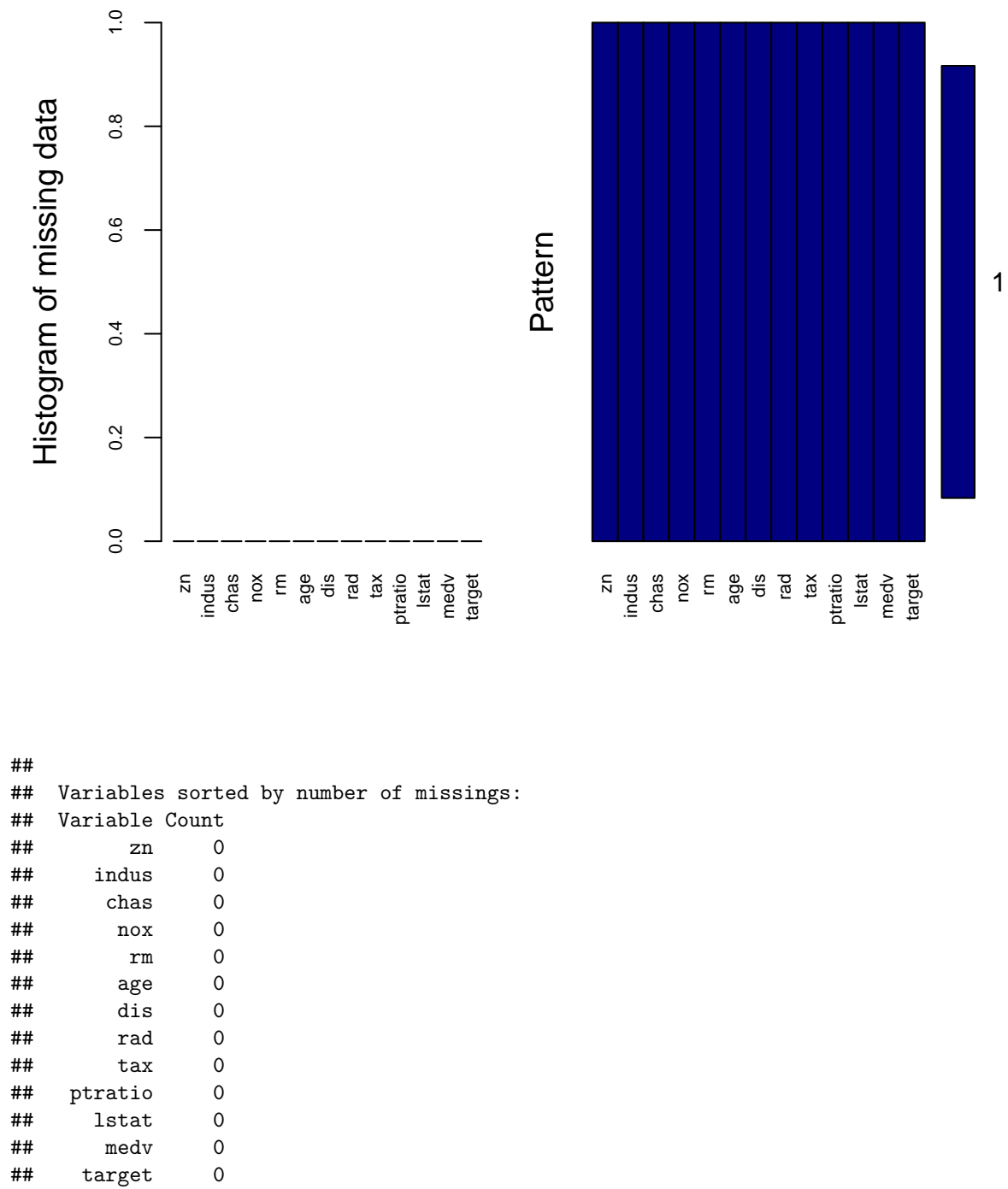


Figure 1. Plot missing values with VIM package

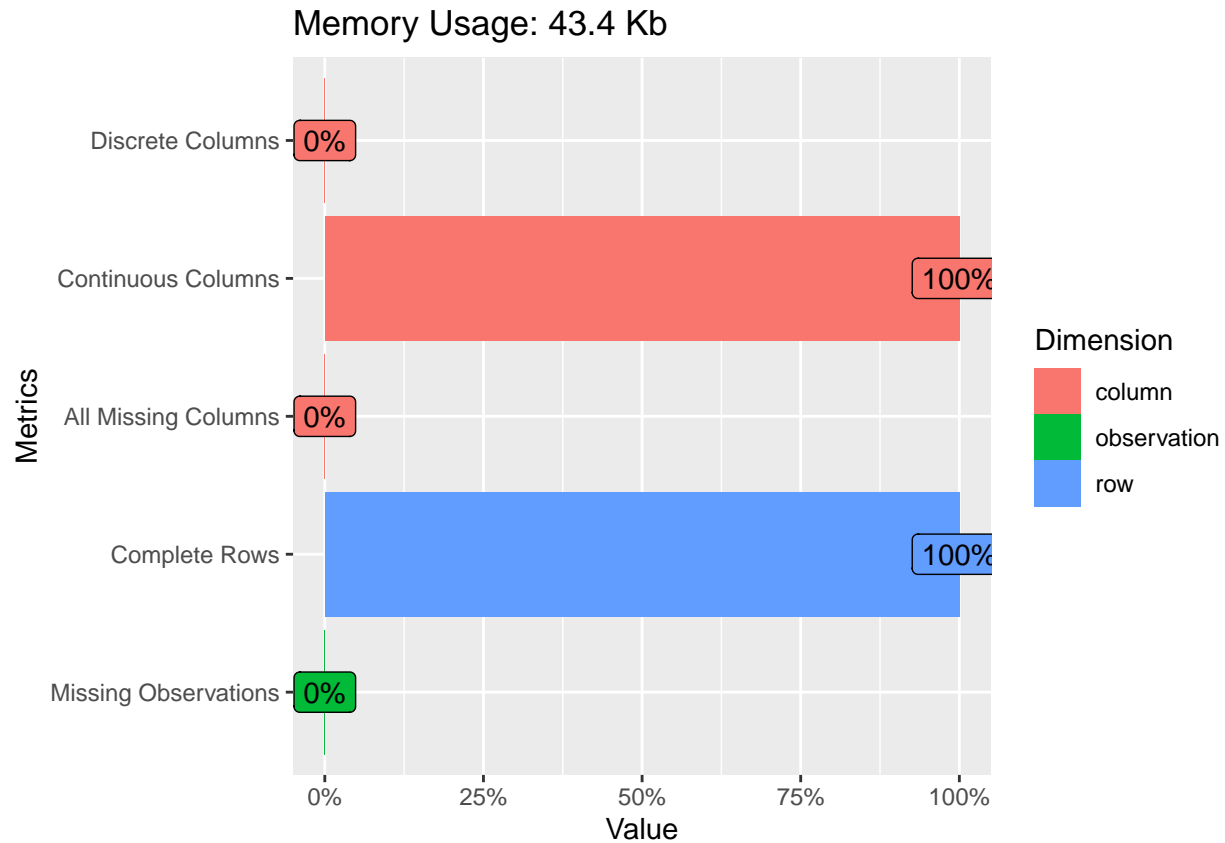


Figure 2. Plot missing values with DataExplorer package

Feature Histograms

For each of the variables, these histograms in figure 3 (**Code Appendix 2.6**) provide a nice overview of each feature, its variation, and paths for potential transformations later on for model construction. Histograms are a quick way to see the shape of the distributions for each feature. Of note are the normally distributed variables, median home value and number of rooms per dwelling. The remaining features appear quite skewed, especially land zoning, distance from employment, tax value, and distance to radial highways. We can also begin to see the affect of outliers that we'll have to account for.

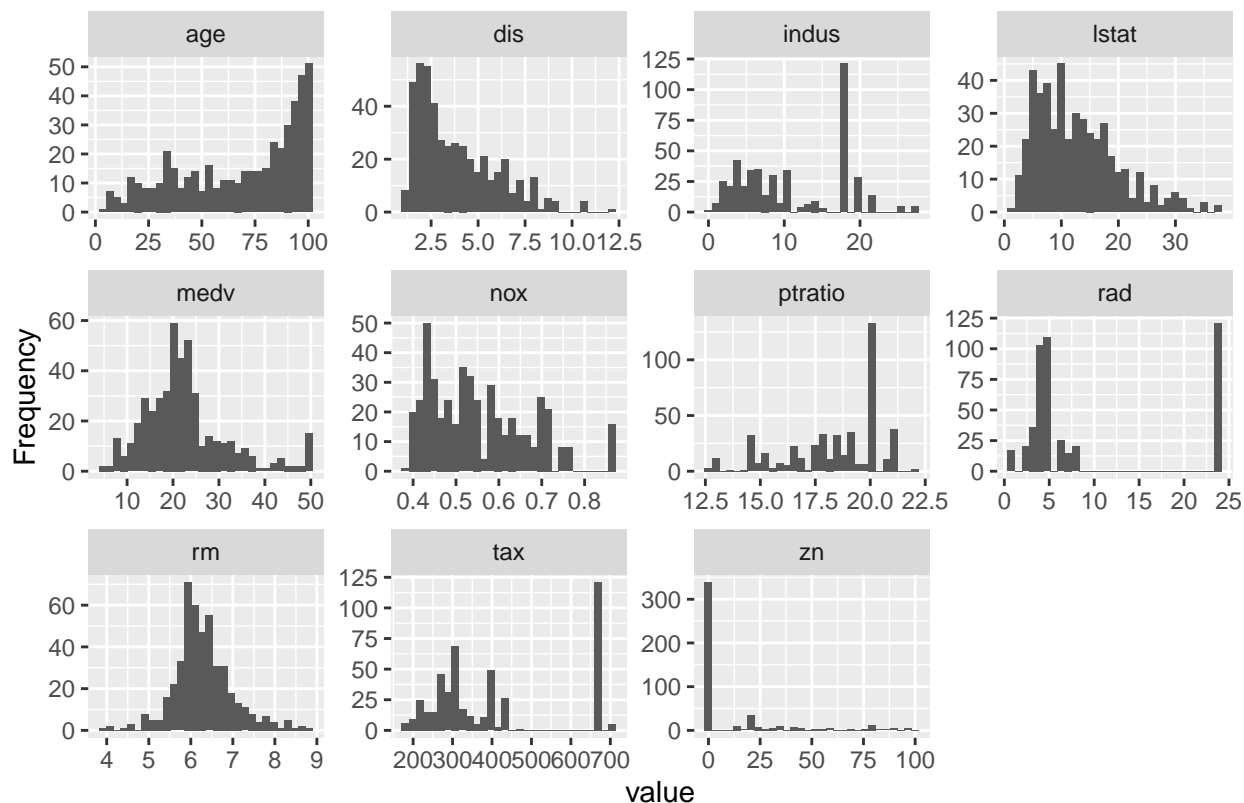


Figure 3. Feature histograms

Feature Boxplots

This box plot visualization (Figure 4) (**Code Appendix 2.7**) gives us an idea of the outliers we have in each variable, but does not give us a good sense of the distribution. We can use the histograms (Figure 3) above to interpret shape. We apply a log re-scaling to better compare the values across variables using a common scale and use notches to compare groups. If the notches of two boxes do not overlap, then this suggests that the medians are significantly different.

For the features we see significant outliers, we can decide to either throw out that variable out altogether and not consider it in our models or impute the outliers with median values. Before deciding on a course of action, we'll look at a few other things.

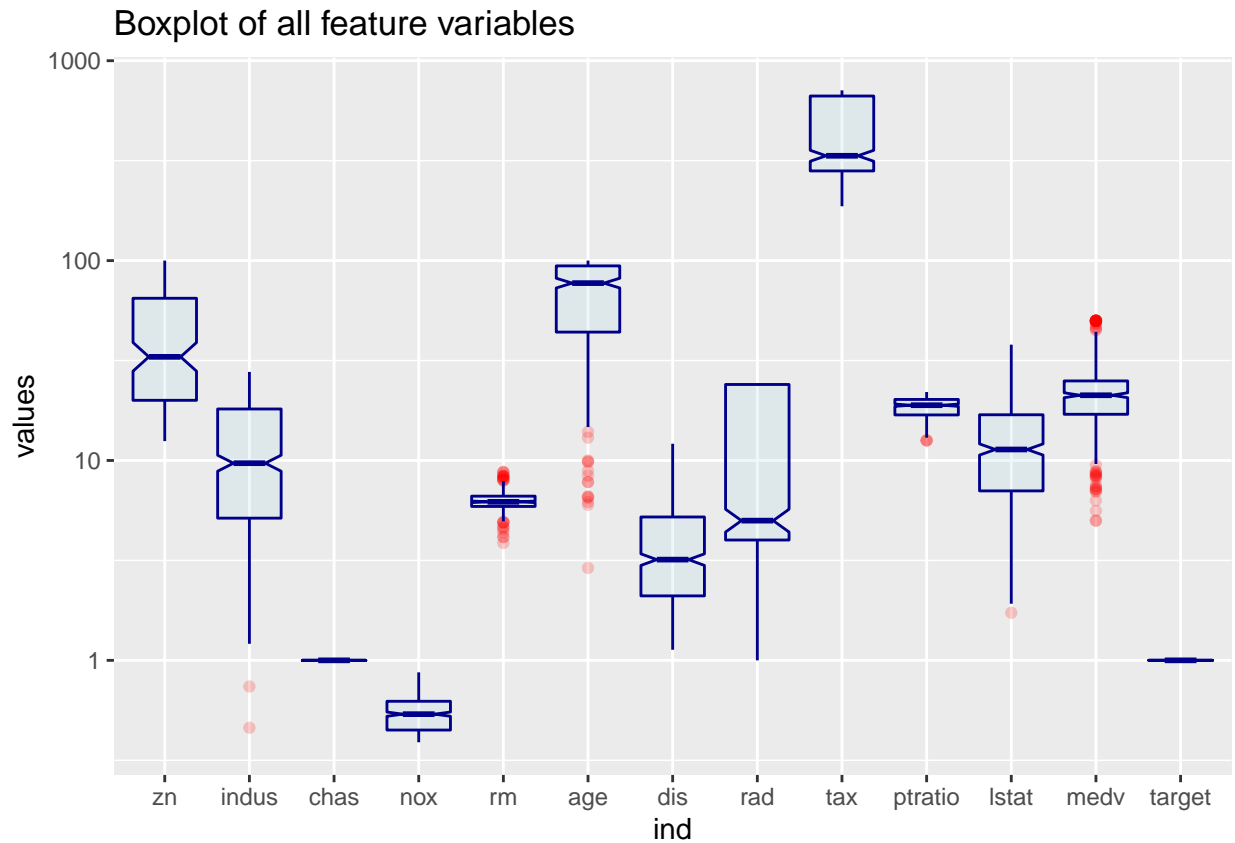


Figure 4. Feature box plots

Feature QQ Plots

The Quantile-Quantile, or QQ plots (figures 5 and 6) (**Code Appendix 2.8**) are used to visualize the deviation of the predictors compared to the normal distribution. It is rather normal to see tails on either side of the QQ as outliers will deviate significantly from the normal distribution. Too much and the feature will not be a good one to use as a predictor.

QQ Plots

Consistent with our other analysis of the features, **zn**, **rad**, and **tax** are not following the normal distribution enough to be helpful in our models, and with exception of the “chas” predictor, all other predictors will need to be transformed for linear regression. (**Code Appendix 2.8.1**)

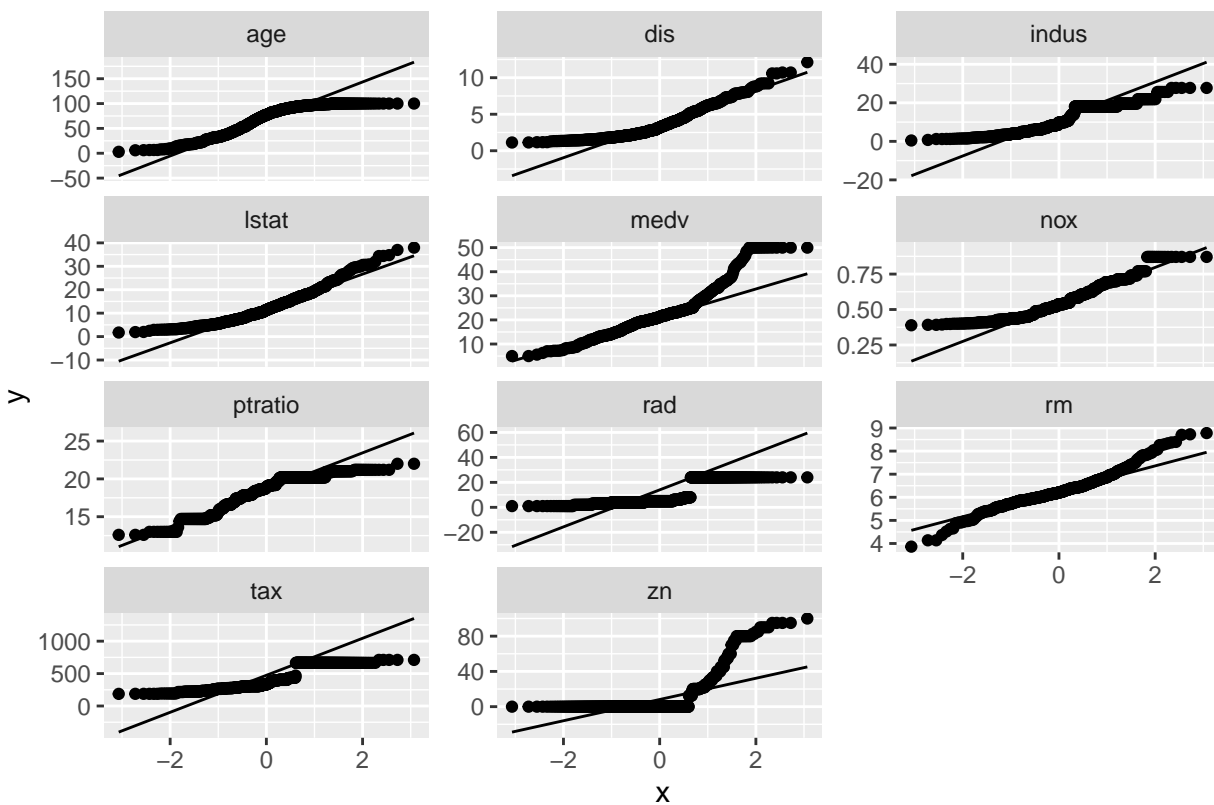


Figure 5. Feature QQ plot

Log QQ Plots

With the log transformation, the distributions look better now. So, as part of the data preparation we will transform the necessary predictors before we use them for the models. (**Code Appendix 2.8.2**)

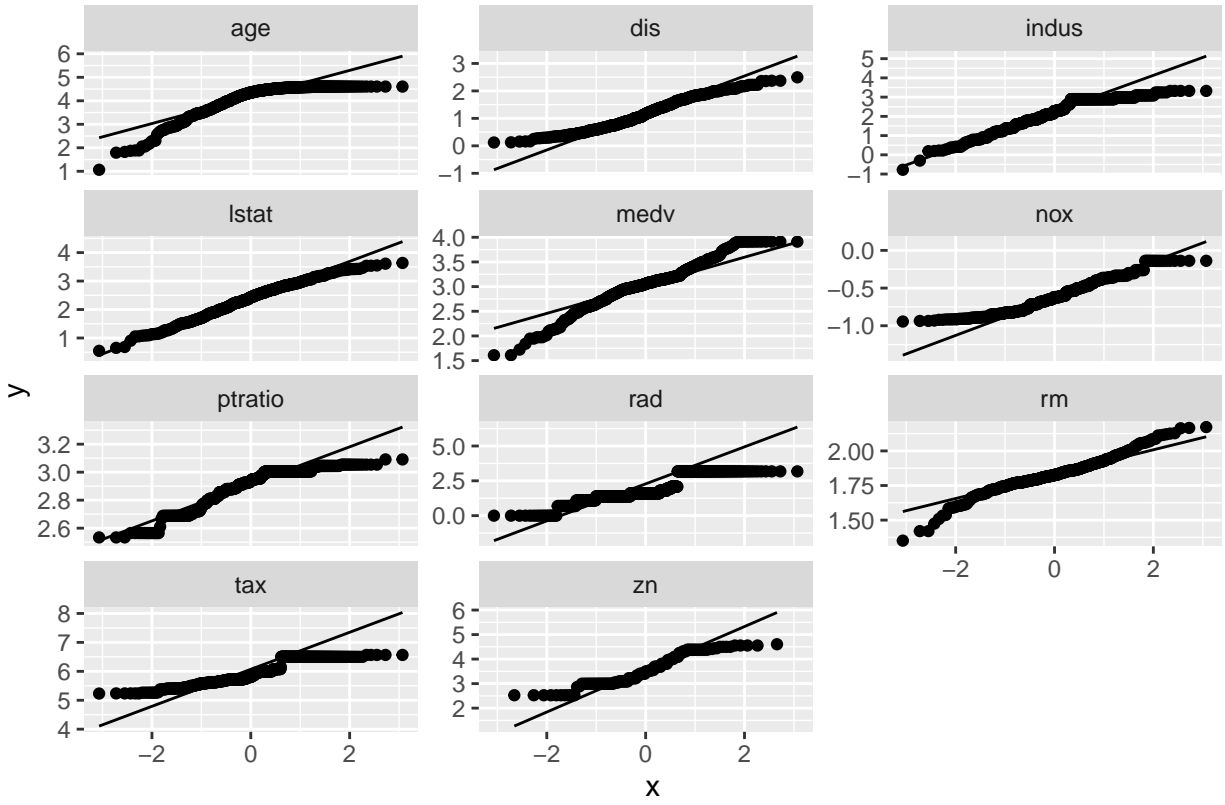


Figure 6. Feature QQ plots, log transformation

Correlation

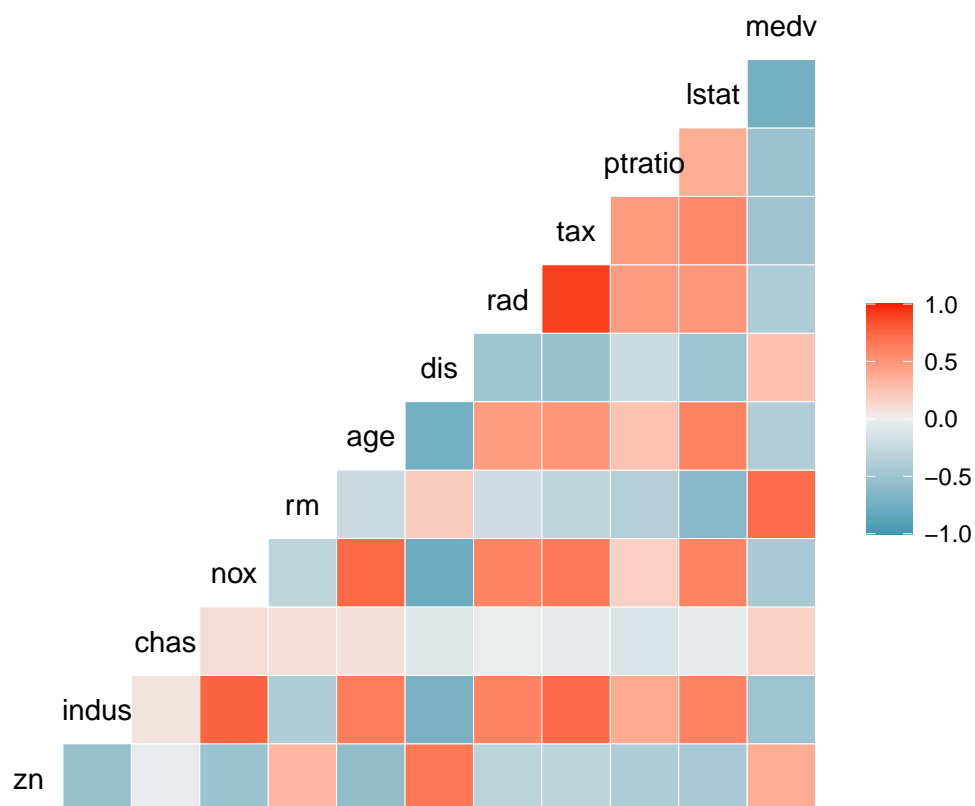


Figure 7. Correlation matrix

It is important to check for features which may also be correlated. Simply, having multiple features relate to themselves can cause overfitting, reduced p values, and strange variances in the data. To avoid this, we exclude one or more of the variables. In the correlation matrix (Figure 7) (**Code Appendix 2.9**), we see that `tax` is very intertwined with two other variables besides the target, showing up as bright red. We'll take care when constructing our models not to use those.

Var1	Var2	Correlation
rad	tax	0.91
indus	nox	0.76
nox	age	0.74
indus	tax	0.73
nox	target	0.73*
rm	medv	0.71
age	target	0.63*
rad	target	0.03*
tax	target	0.61*

Table 3. Correlation results

Data Preparation

Outlier Value Analysis

Coming back to the outliers, let's analyze the variables with outliers and their relation to the **target** variable. (Code Appendix 3.2.1)

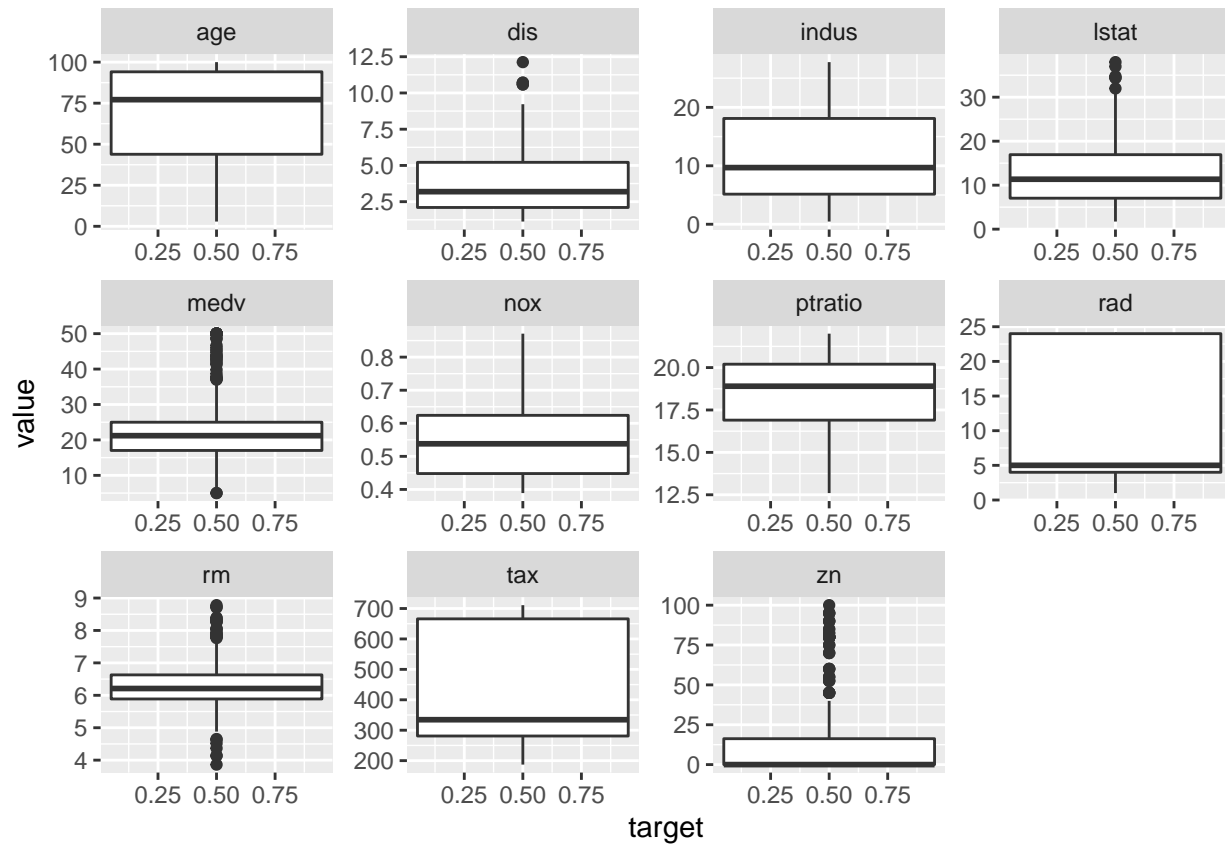


Figure 8. Feature box plots with outliers

Outlier imputation

The boxplots above confirm that there are obvious outliers for variables **dis**, **indus**, **lstat**, **medv**, **ptratio**, **rm**, **tax**, and **zn**. These outliers need to be imputed to prevent them from skewing the results of the modeling (Code Appendix 3.2.2). We use the median values to impute outliers. After imputing the variables with outliers, let's analyze them again with respect to their relation to the **target** variable (Code Appendix 3.2.3). This is to ensure that outliers have been eliminated or at least minimized as much as possible.

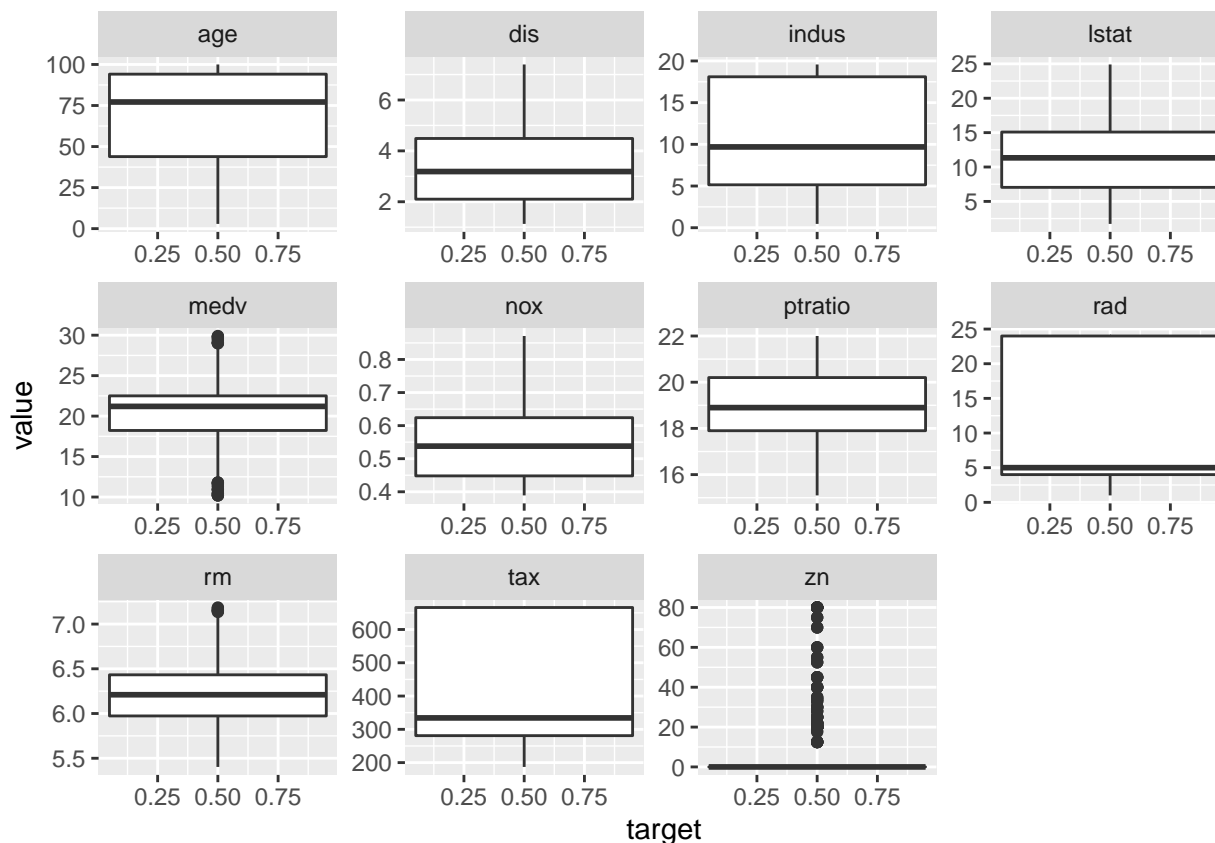


Figure 9. Feature box plots with imputer outliers

Model Building

Model 1: Kitchen Sink

AIC Value = 171.94

This model contains all values except for `tax` which was removed from consideration in any model due to its multicollinearity with other features. No transformations were performed on this model. The rule of thumb for these types of regressions are the AIC value. Lower AIC values indicate a strong model. **(Code Appendix 4.2)**

Model 2: Kitchen Sink Transformed in Part

AIC Value = 173.86

For this model, we include all features, squaring `age` and log transforming `lstat` in hopes of producing a lower AIC score. **(Code Appendix 4.3)**

Model 3: Kitchen Sink Transformed More

AIC Value = 168.1

Now we see the AIC score starting to decrease, we log `zn`, `dis`, and keep `age` squared, though we are unsure if this last transformation has any impact. **(Code Appendix 4.4)**

Model 4: More Transformations

AIC Value = 161.9

This model lowers the AIC score and reduces the residual deviance, logging and scaling `zn`, squaring `age`, and multiplying `rad` by `rm`. (**Code Appendix 4.5**)

Model 5: Variable Importance

AIC Value = 161.9

Keeping model 4's setup, except for deleting `indus` leads to higher residual deviance and the same AIC score. This could be a final contender for our model. We learn that `indus` and `zn` are not very important to this model. (**Code Appendix 4.6**)

Model 6: Narrow Variable Importance

AIC Value = 161.1

Taking what we've learned from model 5, we delete `indus` and multiply `ptratio` by `nox` and push the AIC even lower. (**Code Appendix 4.7**)

Model Selection

Before proceeding to final model selection, we will test for accuracy and error. As a final check we will look at the ROC plots. (**Code Appendix 5.1**)

Model Error

For each model we use the `predict()` function and check the error and r^2 values (**Code Appendix 5.2.1**). The RMSE is decreasing with each model, possibly meaning that the fit is better with each (**Code Appendix 5.2.2**). It is difficult to tell as the training sample is rather small. Model 6's r^2 value may mean that by deleting `indus` from the model it lost some information that would lower the r^2 .

```
##           modelOne  modelTwo modelThree modelFour modelFive  modelSix
## RMSE      0.2592275 0.2572563  0.2389191 0.2306940 0.2308036 0.2170661
## Rsquared  0.7319253 0.7368322  0.7733806 0.7891527 0.7883487 0.8132272
## MAE       0.1256217 0.1251212  0.1143695 0.1054849 0.1094327 0.1011464
```

Confusion Matrix and Accuracy Measurement

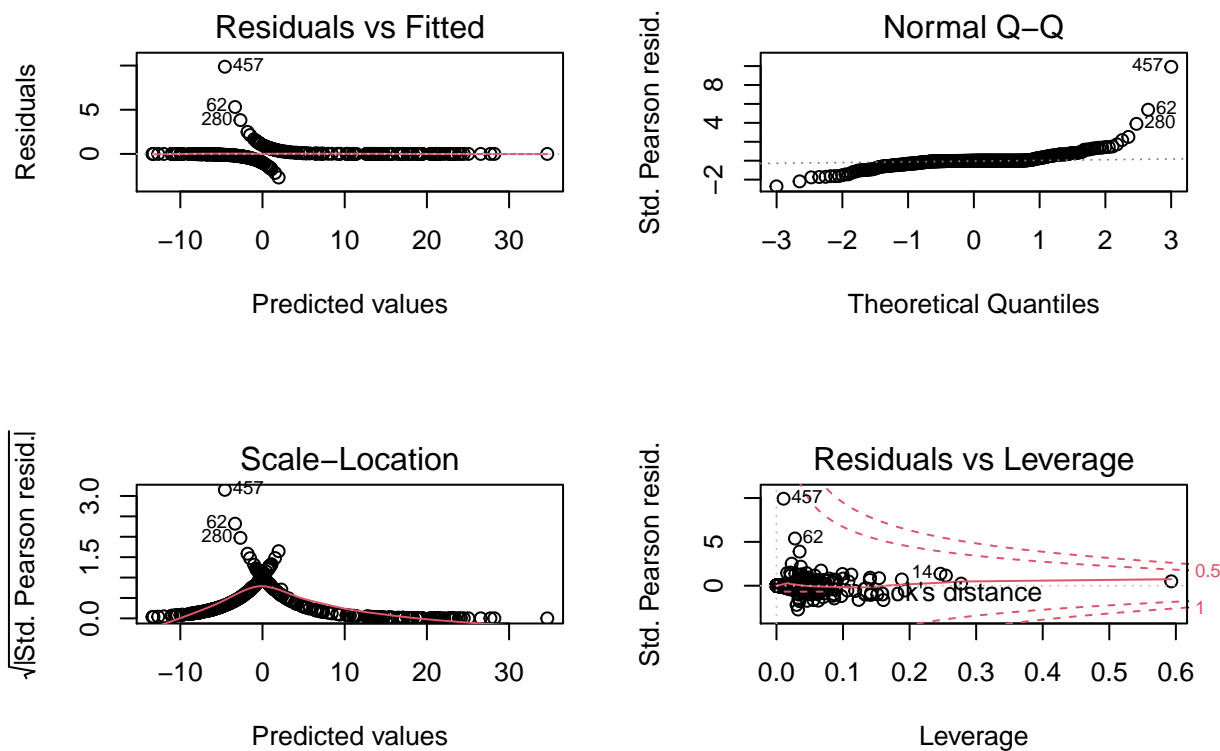
In terms of accuracy, it appears that models 5 and 6 are stand outs with the highest accuracy scores. (**Code Appendix 5.3 & 5.4**)

```
##   cOne.overall cTwo.overall cThree.overall cFour.overall cFive.overall
## 1    0.9139785    0.9247312    0.9354839    0.9462366    0.9354839
##   cSix.overall
## 1    0.9569892
```

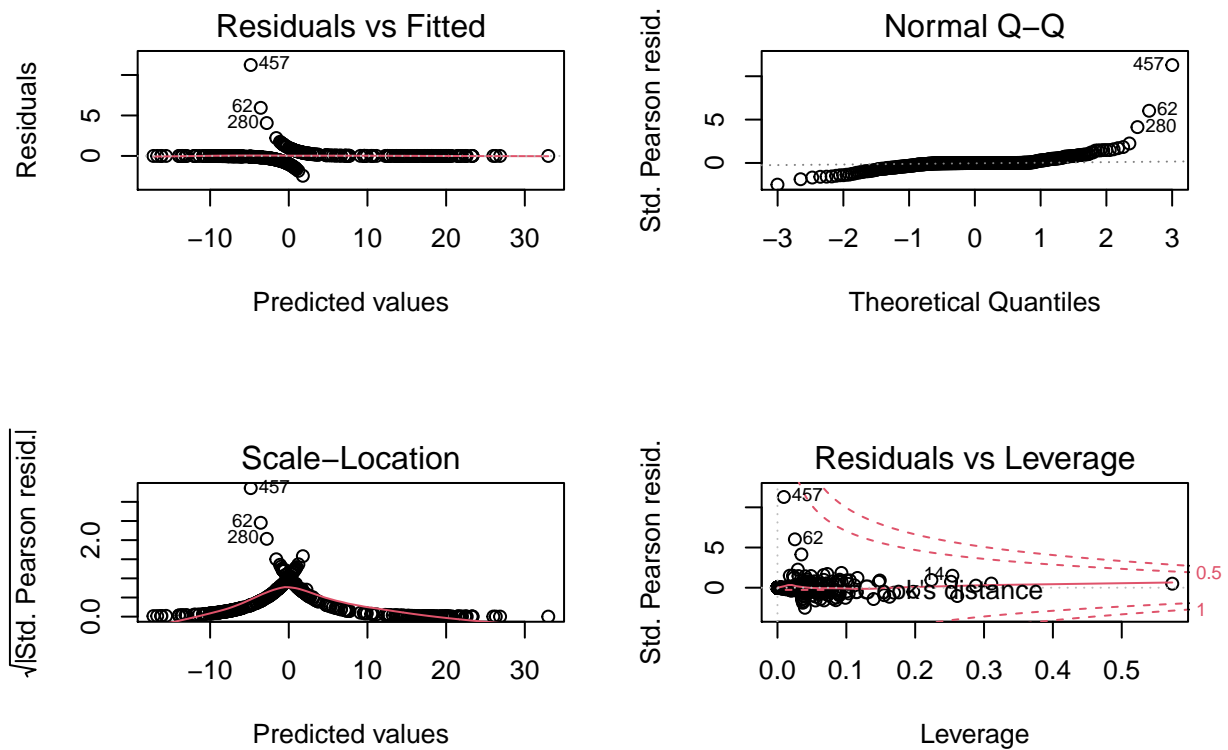
Final Model Plots

In reviewing the final plots for these models, we check again the residuals and QQ plots for each. Both models look extremely close for each of the graphs, with very slight differences. (Code Appendix 5.5.1 & 5.5.2)

Model 5



Model 6



ROC

The ROC plot in figure 10 (**Code Appendix 5.6**) seems to provide the tie breaker for us. Model 6 is 0.1 higher and will be our final model to use for our predictions.

```
##
## Call:
## roc.default(response = train$target, predictor = modelFive$fitted.values, percent = TRUE, plot =
##
## Data: modelFive$fitted.values in 198 controls (train$target 0) < 175 cases (train$target 1).
## Area under the curve: 97.54%
```

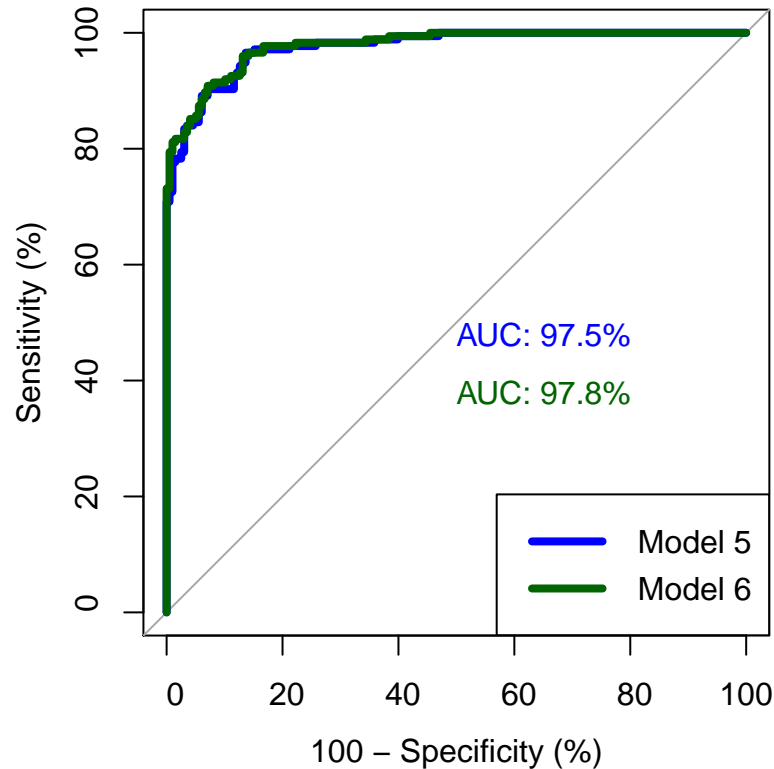


Figure 10. ROC plot with models 5 and 6

Conclusion

We conclude by writing our predictions to the `test` data in a CSV using model 6. To do this, we took the average from the predicted values and assigned either 1 or 0 if the value fell above or below the average. (Code Appendix 5.7)

Code Appendix

1.1 Libraries Used

We use pretty standard packages for this assignment, including the ever-useful `tidyverse`, `ggplot2`, and `caret`. New additions for this assignment include `VIM`, `DataExplorer`, and `broom`.

```
{r, warning = FALSE, message = FALSE, echo=FALSE} library(tidyverse) library(ggplot2) library(VIM)
library(GGally) library(caret) library(broom) library(kableExtra) library(tidymodels) library(DataExplorer)
library(psych) library(pROC)
```

2.1 Data Exploration

2.2 Data Import

```
{r, warning = FALSE, message = FALSE, echo=FALSE} rawTrain <- read.csv("https://raw.githubusercontent.com/MsQCompSci/Data621Group4/main/HW3/crime-training-data_modified.csv", header = TRUE, stringsAsFactors = FALSE) rawTest <- read.csv("https://raw.githubusercontent.com/MsQCompSci/Data621Group4/main/HW3/crime-evaluation-data_modified.csv")
```

2.3 Summary Stats

```
{r, warning = FALSE, message = FALSE, echo=FALSE} glimpse(rawTrain)
```

2.4 Table 1. Glimpse of data structure**

```
{r, warning = FALSE, message = FALSE, echo=FALSE} kable(describe(rawTrain), booktabs = TRUE) %>% kable_styling(latex_options = "scale_down")
```

2.5 Missing Data Checks

2.5.1 Figure 1. Plot missing values with VIM package

```
{r, warning = FALSE, message = FALSE, echo=FALSE} #plot missing values using VIM package  
aggr(rawTrain, col=c('navyblue','red'), numbers=TRUE, sortVars=TRUE, labels=names(rawTrain),  
cex.axis=.7, gap=3, ylab=c("Histogram of missing data","Pattern"))
```

2.5.2 Figure 2. Plot missing values with DataExplorer package

```
{r, warning = FALSE, message = FALSE, echo=FALSE} DataExplorer::plot_intro(rawTrain)
```

2.6 Feature Histograms Figure 3. Feature histograms

```
{r, warning = FALSE, message = FALSE, echo=FALSE} DataExplorer::plot_histogram(rawTrain)
```

2.7 Feature Boxplots Figure 4. Feature box plots

```
{r, warning = FALSE, message = FALSE, echo=FALSE} ggplot(stack(rawTrain), aes(x = ind, y = values))  
+ geom_boxplot(color = "darkblue", fill = "lightblue", alpha = 0.2, outlier.color = "red", outlier.fill = "red",  
outlier.alpha = 0.2, notch = TRUE) + labs(title = "Boxplot of all feature variables") + scale_y_log10()
```

2.8 Feature QQ Plots

2.8.1 QQ Plots Figure 5. Feature QQ plot

```
{r, warning = FALSE, message = FALSE, echo=FALSE} qq_train_data <- rawTrain[, c("age", "dis",  
"indus", "lstat", "medv", "nox", "ptratio", "rad", "rm", "tax", "zn")]
```

```
DataExplorer::plot_qq(qq_train_data, nrow = 4L, ncol = 3L)
```


2.8.2 Log QQ Plots Figure 6. Feature QQ plots, log transformation

```
{r, warning = FALSE, message = FALSE, echo=FALSE} log_qq_train_data <- DataExplorer::update_columns(qq_train_data, ind = names(qq_train_data), what = log)
DataExplorer::plot_qq(log_qq_train_data, nrow = 4L, ncol = 3L)
```

2.9 Correlation Figure 7. Correlation matrix

```
{r, warning = FALSE, message = FALSE, echo=FALSE} #correlation matrix for predictors ggcorr(rawTrain%>% select(zn:medv))
```

3.1 Data Preparation

3.2 Outlier Value Analysis

3.2.1 Figure 8. Feature box plots with outliers

```
{r, warning = FALSE, message = FALSE, echo=FALSE} pred_vs_target <- gather(rawTrain, variable, value, -c(chas,target))
ggplot(pred_vs_target, aes(x = target, y = value)) + geom_boxplot() + facet_wrap(~variable, scales = 'free')
```

3.2.2 Outlier imputation

```
{r, warning = FALSE, message = FALSE, echo=FALSE} rawTrain_prepmed <- rawTrain %>% mutate( dis = ifelse(dis > 7.5, median(dis), dis), indus = ifelse(indus > 21, median(indus), indus), lstat = ifelse(lstat > 25, median(lstat), lstat), medv = ifelse(medv > 30 | medv < 10, median(medv), medv), ptratio = ifelse(ptratio < 15.0, median(ptratio), ptratio), rm = ifelse(rm > 7.2 | rm < 5.4, median(rm), rm), tax = ifelse(tax > 700.0, median(tax), tax), zn = ifelse(zn > 80, median(zn), zn) )
```

3.2.3 Figure 9. Feature box plots with imputed outliers

```
{r, warning = FALSE, message = FALSE, echo=FALSE} pred_vs_target <- gather(rawTrain_prepmed, variable, value, -c(chas,target))
ggplot(pred_vs_target, aes(x = target, y = value)) + geom_boxplot() + facet_wrap(~variable, scales = 'free')
```

4.1 Model Building

```
{r, warning = FALSE, message = FALSE, echo=FALSE} dt <- createDataPartition(rawTrain_prepmed$target, p = .8, list = FALSE, times = 1)
train <- rawTrain[dt,] test <- rawTrain[-dt,]
```

4.2 Model 1: Kitchen Sink

```
#remove Tax due to high correlation with other variables
modelOne <- glm(target ~ zn + indus + chas + nox + rm + age + dis + rad + ptratio + lstat + medv , data = train)
summary(modelOne)
```

```
##
## Call:
## glm(formula = target ~ zn + indus + chas + nox + rm + age + dis +
##      rad + ptratio + lstat + medv, family = "binomial", data = train)
##
## Deviance Residuals:
##      Min       1Q   Median       3Q      Max
## -1.94713  -0.15956  -0.00142   0.00257   2.65592
##
## Coefficients:
##              Estimate Std. Error z value Pr(>|z|)
## (Intercept) -43.51032    7.52713  -5.780 7.45e-09 ***
## zn          -0.09509    0.04283  -2.220 0.026415 *
## indus       -0.13017    0.05126  -2.539 0.011101 *
## chas         1.01891    0.87812   1.160 0.245913
## nox         52.30104    9.25460   5.651 1.59e-08 ***
## rm          -0.66619    0.80062  -0.832 0.405355
## age          0.04045    0.01562   2.589 0.009629 **
## dis          1.01464    0.27237   3.725 0.000195 ***
## rad          0.58932    0.15916   3.703 0.000213 ***
## ptratio      0.29704    0.13673   2.172 0.029824 *
## lstat        0.08466    0.06129   1.381 0.167147
## medv         0.21103    0.07651   2.758 0.005814 **
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## (Dispersion parameter for binomial family taken to be 1)
##
##      Null deviance: 515.67  on 372  degrees of freedom
## Residual deviance: 154.29  on 361  degrees of freedom
## AIC: 178.29
##
## Number of Fisher Scoring iterations: 9
```

4.3 Model 2: Kitchen Sink Transformed in Part

```
#remove Tax squared age and log lstat
modelTwo <- glm(target ~ zn + indus + chas + nox + rm + age^2 + dis + rad + ptratio + log2(lstat) + medv , data = train)
summary(modelTwo)
```

```
##
## Call:
## glm(formula = target ~ zn + indus + chas + nox + rm + age^2 +
##      dis + rad + ptratio + log2(lstat) + medv, family = "binomial",
##      data = train)
##
```

```
## Deviance Residuals:
##      Min       1Q   Median       3Q      Max
## -1.94183  -0.15525  -0.00158   0.00262   2.96601
##
## Coefficients:
##              Estimate Std. Error z value Pr(>|z|)
## (Intercept) -44.10925     8.12209  -5.431 5.61e-08 ***
## zn          -0.09195     0.04212  -2.183 0.029044 *
## indus       -0.12396     0.05088  -2.436 0.014839 *
## chas         1.09660     0.89052   1.231 0.218171
## nox         52.03778     9.27240   5.612 2.00e-08 ***
## rm          -0.76466     0.82253  -0.930 0.352553
## age          0.04377     0.01558   2.809 0.004967 **
## dis          1.02821     0.27097   3.795 0.000148 ***
## rad          0.59060     0.15700   3.762 0.000169 ***
## ptratio      0.30841     0.13756   2.242 0.024965 *
## log2(lstat)  0.48983     0.56545   0.866 0.386339
## medv         0.21596     0.07660   2.819 0.004811 **
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## (Dispersion parameter for binomial family taken to be 1)
##
##      Null deviance: 515.67  on 372  degrees of freedom
## Residual deviance: 155.45  on 361  degrees of freedom
## AIC: 179.45
##
## Number of Fisher Scoring iterations: 9
```

#This one has a litter lower AIC

4.4 Model 3: Kitchen Sink Transformed More

```
#log10(zn + 1), log10(dis) and deleted log2(lstat) - not significant
modelThree <- glm(target ~ log10(zn + 1) + indus + chas + nox + rm + age^2 + log10(dis) + rad + ptratio
summary(modelThree)
```

```
##
## Call:
## glm(formula = target ~ log10(zn + 1) + indus + chas + nox + rm +
##      age^2 + log10(dis) + rad + ptratio + medv, family = "binomial",
##      data = train)
##
## Deviance Residuals:
##      Min       1Q   Median       3Q      Max
## -1.8942  -0.1635  -0.0119   0.0019   3.2496
##
## Coefficients:
##              Estimate Std. Error z value Pr(>|z|)
## (Intercept) -47.44387     8.10602  -5.853 4.83e-09 ***
## log10(zn + 1) -1.13423     0.58075  -1.953 0.050813 .
```

```
## indus      -0.07491    0.05175   -1.448 0.147708
## chas       0.97674    0.88319    1.106 0.268759
## nox       56.57968    9.30018    6.084 1.17e-09 ***
## rm        -1.24766    0.77840   -1.603 0.108968
## age        0.05521    0.01515    3.645 0.000267 ***
## log10(dis) 11.52593    2.63543    4.373 1.22e-05 ***
## rad        0.63877    0.16354    3.906 9.39e-05 ***
## ptratio    0.36363    0.14786    2.459 0.013925 *
## medv       0.24444    0.08127    3.008 0.002631 **
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## (Dispersion parameter for binomial family taken to be 1)
##
## Null deviance: 515.67  on 372  degrees of freedom
## Residual deviance: 149.60  on 362  degrees of freedom
## AIC: 171.6
##
## Number of Fisher Scoring iterations: 9
```

#AIC is lower again (not sure if age² is helpful)

4.5 Model 4: More Transformations

```
#combine rad and rm (multiplied) - they seemed to correspond in their distributions
modelFour<- glm(target ~ log10(zn + 1) + indus + chas + nox + age^2 + log10(dis) + rad*rm + ptratio +
summary(modelFour)
```

```
##
## Call:
## glm(formula = target ~ log10(zn + 1) + indus + chas + nox + age^2 +
##      log10(dis) + rad * rm + ptratio + medv, family = "binomial",
##      data = train)
##
## Deviance Residuals:
##      Min       1Q   Median       3Q      Max
## -1.8465  -0.1084  -0.0075   0.0079   3.1947
##
## Coefficients:
##              Estimate Std. Error z value Pr(>|z|)
## (Intercept)  -36.85625    8.10598  -4.547 5.45e-06 ***
## log10(zn + 1)  -1.17363    0.62901  -1.866 0.062063 .
## indus         -0.09321    0.05289  -1.762 0.077990 .
## chas           0.85010    0.91197   0.932 0.351254
## nox          59.56105    9.73564   6.118 9.49e-10 ***
## age           0.05902    0.01603   3.683 0.000231 ***
## log10(dis)    11.69841    2.71228   4.313 1.61e-05 ***
## rad          -1.60615    0.56149  -2.861 0.004230 **
## rm           -3.45779    1.16251  -2.974 0.002935 **
## ptratio       0.39201    0.15071   2.601 0.009292 **
## medv          0.27377    0.08725   3.138 0.001703 **
```

```
## rad:rm          0.36543    0.10880    3.359 0.000783 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## (Dispersion parameter for binomial family taken to be 1)
##
##      Null deviance: 515.67  on 372  degrees of freedom
## Residual deviance: 142.42  on 361  degrees of freedom
## AIC: 166.42
##
## Number of Fisher Scoring iterations: 9
```

#AIC is lower #Not sure what the rationale is for this working but it lowered the AIC number and Resid

4.6 Model 5: Variable Importance

```
#delete indus
modelFive<-glm(target ~ log10(zn+1)+ nox + age^2 + log10(dis) + rad*rm + ptratio + medv, data = train,
summary(modelFive)
```

```
##
## Call:
## glm(formula = target ~ log10(zn + 1) + nox + age^2 + log10(dis) +
##      rad * rm + ptratio + medv, family = "binomial", data = train)
##
## Deviance Residuals:
##      Min       1Q   Median       3Q      Max
## -2.04661  -0.13593  -0.00786   0.00694   3.02945
##
## Coefficients:
##              Estimate Std. Error z value Pr(>|z|)
## (Intercept)  -33.93967    7.64288  -4.441 8.97e-06 ***
## log10(zn + 1)  -1.36012    0.60023  -2.266 0.023452 *
## nox           53.87017    8.88189   6.065 1.32e-09 ***
## age            0.06060    0.01563   3.877 0.000106 ***
## log10(dis)     12.83739    2.68341   4.784 1.72e-06 ***
## rad           -1.60839    0.58723  -2.739 0.006164 **
## rm            -3.69867    1.17934  -3.136 0.001711 **
## ptratio        0.34089    0.14908   2.287 0.022212 *
## medv           0.30914    0.08738   3.538 0.000403 ***
## rad:rm         0.37353    0.11273   3.314 0.000921 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## (Dispersion parameter for binomial family taken to be 1)
##
##      Null deviance: 515.67  on 372  degrees of freedom
## Residual deviance: 146.14  on 363  degrees of freedom
## AIC: 166.14
##
## Number of Fisher Scoring iterations: 9
```

```
#AIC is higher #resiudal deviance is lower
# I looked at the histograms and looked for complementary shapes to decide what to multiply
```

4.7 Model 6: Narrow Variable Importance

```
#multiply ptratio*nox (remove squared from age)
modelSix<- glm(target ~ log10(zn + 1) + age + ptratio*nox + log10(dis) + rad*rm + medv, data = train,
summary(modelSix)
```

```
##
## Call:
## glm(formula = target ~ log10(zn + 1) + age + ptratio * nox +
##      log10(dis) + rad * rm + medv, family = "binomial", data = train)
##
## Deviance Residuals:
##      Min       1Q   Median       3Q      Max
## -1.97208  -0.12192  -0.00215   0.00749   3.11302
##
## Coefficients:
##              Estimate Std. Error z value Pr(>|z|)
## (Intercept)  -84.88763   27.49948  -3.087 0.002023 **
## log10(zn + 1)  -1.55446    0.64317  -2.417 0.015654 *
## age           0.06447    0.01653   3.899 9.64e-05 ***
## ptratio       3.08655    1.38630   2.226 0.025983 *
## nox          141.95535   46.78231   3.034 0.002410 **
## log10(dis)    12.28015    2.72109   4.513 6.39e-06 ***
## rad          -1.52131    0.58228  -2.613 0.008984 **
## rm           -3.58744    1.20264  -2.983 0.002855 **
## medv         0.31120    0.09094   3.422 0.000621 ***
## ptratio:nox   -4.82724    2.42437  -1.991 0.046467 *
## rad:rm        0.36078    0.11232   3.212 0.001317 **
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## (Dispersion parameter for binomial family taken to be 1)
##
##      Null deviance: 515.67  on 372  degrees of freedom
## Residual deviance: 141.42  on 362  degrees of freedom
## AIC: 163.42
##
## Number of Fisher Scoring iterations: 9
```

```
#AIC is lower
```

5.1 Model Selection

Before proceeding to final model selection, we will test for accuracy and error. As a final check we will look at the ROC plots.

5.2 Model Error

5.2.1

```
{r, warning = FALSE, message = FALSE} #Make predictions predOne = predict(modelOne,test, type = "response") predTwo = predict(modelTwo,test, type = "response") predThree = predict(modelThree,test, type = "response") predFour = predict(modelFour,test, type = "response") predFive = predict(modelFive,test, type = "response") predSix = predict(modelSix,test, type = "response")
```

##5.2.2

```
{r, warning = FALSE, message = FALSE} #Error Measures data.frame(modelOne = postResample(pred = predOne, obs = testtarget), modelTwo = postResample(pred = predTwo, obs = testtarget), modelThree = postResample(pred = predThree, obs = testtarget), modelFour = postResample(pred = predFour, obs = testtarget), modelFive = postResample(pred = predFive, obs = testtarget), modelSix = postResample(pred = predSix, obs = testtarget))
```

5.3 Confusion Matrix and Accuracy Measurment

```
{r, warning = FALSE, message = FALSE, echo=FALSE} Extract Accuracy
```

5.3.1 Model 1

```
resultsFitOne <- ifelse(predOne > 0.5,1,0) resultsFitOne <- as.factor(resultsFitOne)
cOne <- confusionMatrix(as.factor(testtarget),resultsFitOne)accOne <- as.data.frame(cOneoverall)[1]
accOne<- accOne %>% slice(1)
```

5.3.2 Model 2

```
resultsFitTwo <- ifelse(predTwo > 0.5,1,0) resultsFitTwo <- as.factor(resultsFitTwo)
cTwo <- confusionMatrix(resultsFitTwo, as.factor(testtarget))accTwo <- as.data.frame(cTwooverall)[1]
accTwo<- accTwo %>% slice(1)
```

5.3.3 Model 3

```
resultsFitThree<- ifelse(predThree > 0.5,1,0) resultsFitThree <- as.factor(resultsFitThree)
cThree <- confusionMatrix(resultsFitThree, as.factor(testtarget))accThree <- as.data.frame(cThreeoverall)[1]
accThree<- accThree%>% slice(1)
```

5.3.4 Model 4

```
resultsFitFour<- ifelse(predFour > 0.5,1,0) resultsFitFour <- as.factor(resultsFitFour)
cFour <- confusionMatrix(resultsFitFour, as.factor(testtarget))accFour <- as.data.frame(cFouroverall)[1]
accFour<- accFour%>% slice(1)
```

5.3.5 Model 5

```
resultsFitFive<- ifelse(predFive > 0.5,1,0) resultsFitFive <- as.factor(resultsFitFive)
cFive <- confusionMatrix(resultsFitFive, as.factor(testtarget))accFive <- as.data.frame(cFiveoverall)[1]
accFive<- accFive%>% slice(1)
```

5.3.6 Model 6

```
resultsFitSix<- ifelse(predSix > 0.5,1,0) resultsFitSix <- as.factor(resultsFitSix)
cSix<- confusionMatrix(resultsFitSix, as.factor(testtarget))accSix <- as.data.frame(cSix$overall)[1]
accSix<- accSix%>% slice(1)
```

5.4 Accuracy Results

```
{r, warning = FALSE, message = FALSE, echo=FALSE} #create a table with accuracies data.frame(c(accOne,
accTwo, accThree, accFour,accFive, accSix))
```

5.5 Final Model Plots

5.5.1 Model 5

```
{r, warning = FALSE, message = FALSE, echo=FALSE} par(mfrow = c(2,2)) plot(modelFive)
```

5.5.2 Model 6

```
{r, warning = FALSE, message = FALSE, echo=FALSE} par(mfrow = c(2,2)) plot(modelSix)
```

5.6 ROC Figure 10. ROC plot with models 5 and 6

```
{r, warning = FALSE, message = FALSE, echo=FALSE} par(pty = "s") roc(traintarget, modelFivefitted.values,
plot = TRUE, legacy.axes = TRUE, percent=TRUE, col="blue", lwd=4, print.auc = TRUE)
plot.roc(traintarget, modelSixfitted.values, percent=TRUE, col="dark green", lwd=4, print.auc=TRUE,
add=TRUE, print.auc.y=40) legend("bottomright", legend=c("Model 5", "Model 6"), col=c("blue", "dark
green"), lwd=4)
```

5.7 Conclusion

```
{r echo=FALSE, message=FALSE, warning=FALSE} rawTesttarget_prob <- predict(modelSix, newdata =
rawTest)mean_test <- mean(rawTesttarget_prob) rawTesttarget_pred <- ifelse(rawTesttarget_prob
>= mean_test, 1, 0) rawTest %>% write.csv(., "crime_pred.csv", row.names = F)
```