# HW 3

*Team 2*

*April 10, 2019*

## Contents

## Overview

In this homework assignment, you will explore, analyze and model a data set containing information on crime for various neighborhoods of a major city. Each record has a response variable indicating whether or not the crime rate is above the median crime rate (1) or not (0). Below is a short description of the variables of interest in the data set:

1. `zn`: proportion of residential land zoned for large lots (over 25000 square feet) (predictor variable)
2. `indus`: proportion of non-retail business acres per suburb (predictor variable)
3. `chas`: a dummy var. for whether the suburb borders the Charles River (1) or not (0) (predictor variable)
4. `nox`: nitrogen oxides concentration (parts per 10 million) (predictor variable)
5. `rm`: average number of rooms per dwelling (predictor variable)
6. `age`: proportion of owner-occupied units built prior to 1940 (predictor variable)
7. `dis`: weighted mean of distances to five Boston employment centers (predictor variable)
8. `rad`: index of accessibility to radial highways (predictor variable)
9. `tax`: full-value property-tax rate per $10,000 (predictor variable)???
10. `ptratio`: pupil-teacher ratio by town (predictor variable)
11. `black`: 1000(Bk - 0.63)2 where Bk is the proportion of blacks by town (predictor variable)
12. `lstat`: lower status of the population (percent) (predictor variable)
13. `medv`: median value of owner-occupied homes in $1000s (predictor variable)
14. `target`: whether the crime rate is above the median crime rate (1) or not (0) (response variable)

## Objective

Your objective is to build a binary logistic regression model on the training data set to predict whether the neighborhood will be at risk for high crime levels. You will provide classifications and probabilities for the evaluation data set using your binary logistic regression model. You can only use the variables given to you (or variables that you derive from the variables provided). Below is a short description of the variables of interest in the data set:

## Dependencies

For this project, we used Rstudio, ggplot2, and corrplot.

```
#install.packages('corrplot')

require(ggplot2)
require(corrplot)
require(dplyr)
require(randomForest)
```

# Data Exploration

First, we read the data as a csv then performed some simple statistical calculations so that we could explore the data. Below we can see a sample of the data output as it was read from the csv.
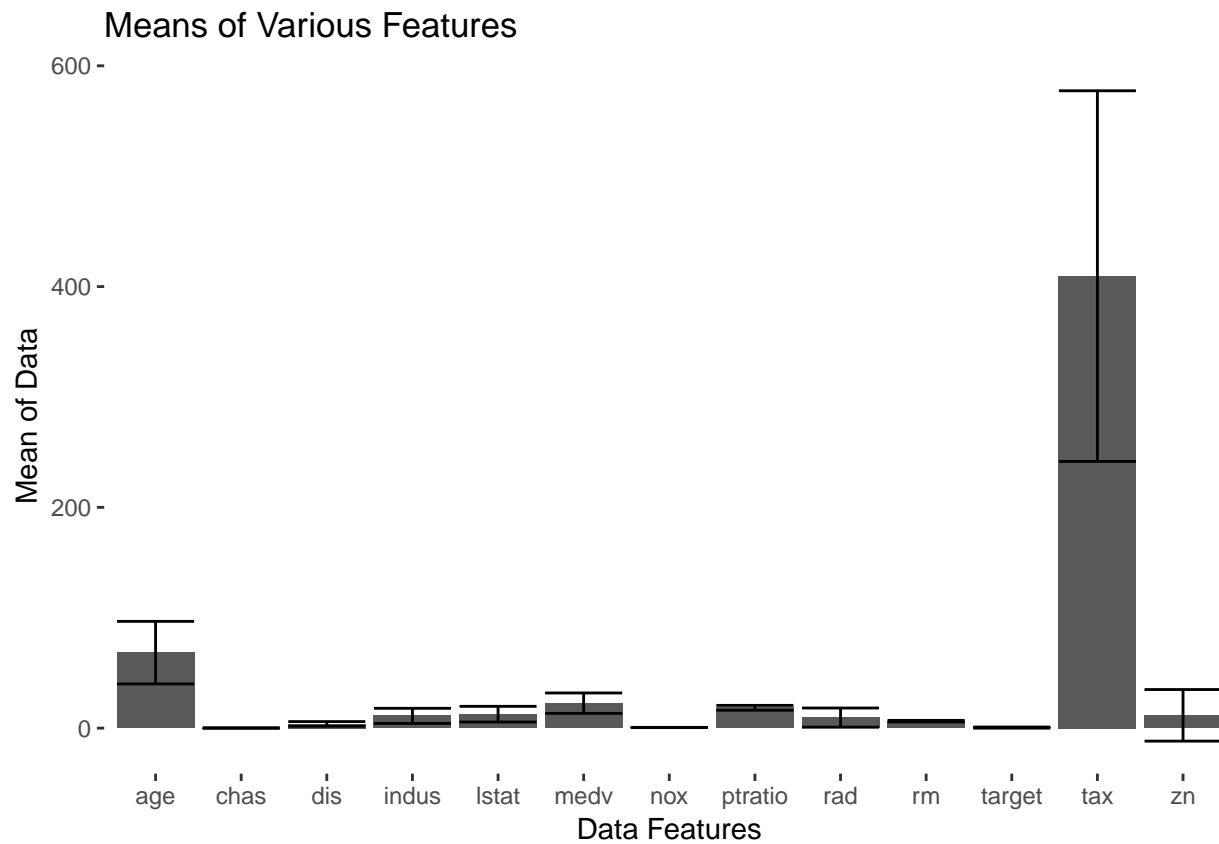
| zn | indus | chas | nox | rm | age | dis | rad | tax | ptratio | lstat | medv | target |
|----|-------|------|-----|----|-----|-----|-----|-----|---------|-------|------|--------|
| 0 | 19.58 | 0 | 0.605 | 7.929 | 96.2 | 2.0459 | 5 | 403 | 14.7 | 3.70 | 50.0 | 1 |
| 0 | 19.58 | 1 | 0.871 | 5.403 | 100.0 | 1.3216 | 5 | 403 | 14.7 | 26.82 | 13.4 | 1 |
| 0 | 18.10 | 0 | 0.740 | 6.485 | 100.0 | 1.9784 | 24 | 666 | 20.2 | 18.85 | 15.4 | 1 |
| 30 | 4.93 | 0 | 0.428 | 6.393 | 7.8 | 7.0355 | 6 | 300 | 16.6 | 5.19 | 23.7 | 0 |
| 0 | 2.46 | 0 | 0.488 | 7.155 | 92.2 | 2.7006 | 3 | 193 | 17.8 | 4.82 | 37.9 | 0 |

## Summary Statistics

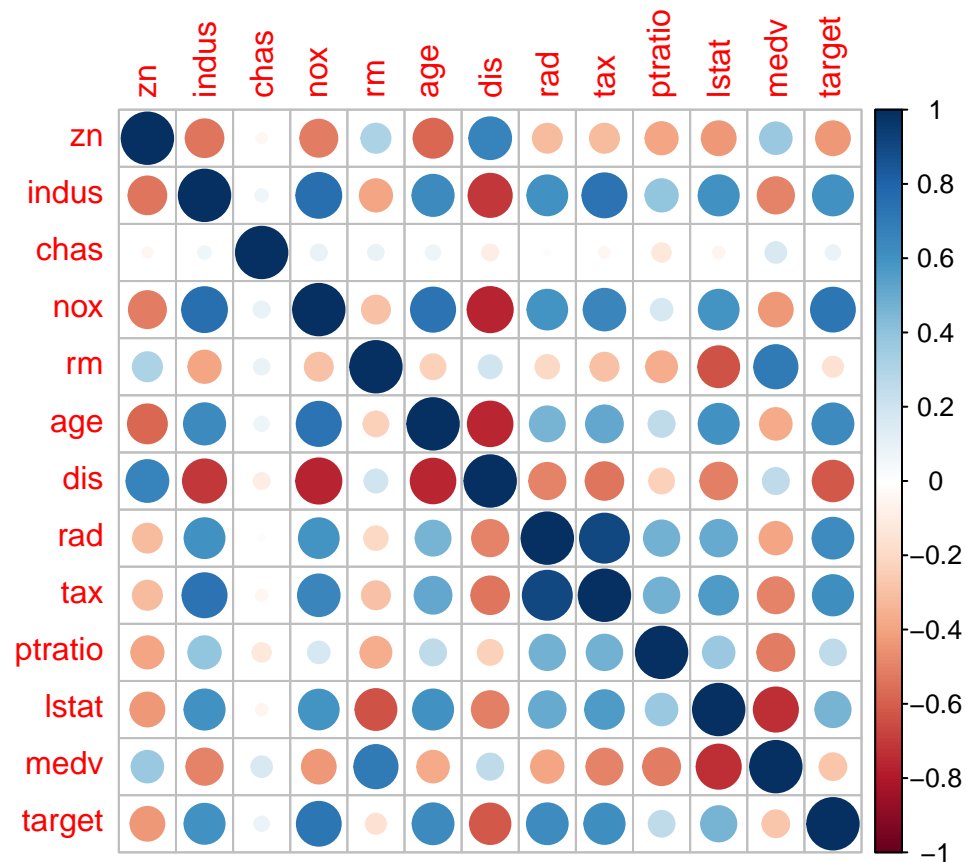We then calculated the mean and standard deviation for each data vector:

| | means | sds |
|---|-------|-----|
| zn | 11.5772532 | 23.3646511 |
| indus | 11.1050215 | 6.8458549 |
| chas | 0.0708155 | 0.2567920 |
| nox | 0.5543105 | 0.1166667 |
| rm | 6.2906738 | 0.7048513 |
| age | 68.3675966 | 28.3213784 |
| dis | 3.7956929 | 2.1069496 |
| rad | 9.5300429 | 8.6859272 |
| tax | 409.5021459 | 167.9000887 |
| ptratio | 18.3984979 | 2.1968447 |
| lstat | 12.6314592 | 7.1018907 |
| medv | 22.5892704 | 9.2396814 |
| target | 0.4914163 | 0.5004636 |

Below is a bar chart that illutrates the average and standard deviation for each of our data vectors. As we can see, the `tax` vector is a totally different magnitude than the rest. Models involving this vector will benefit from normalization or scaling.

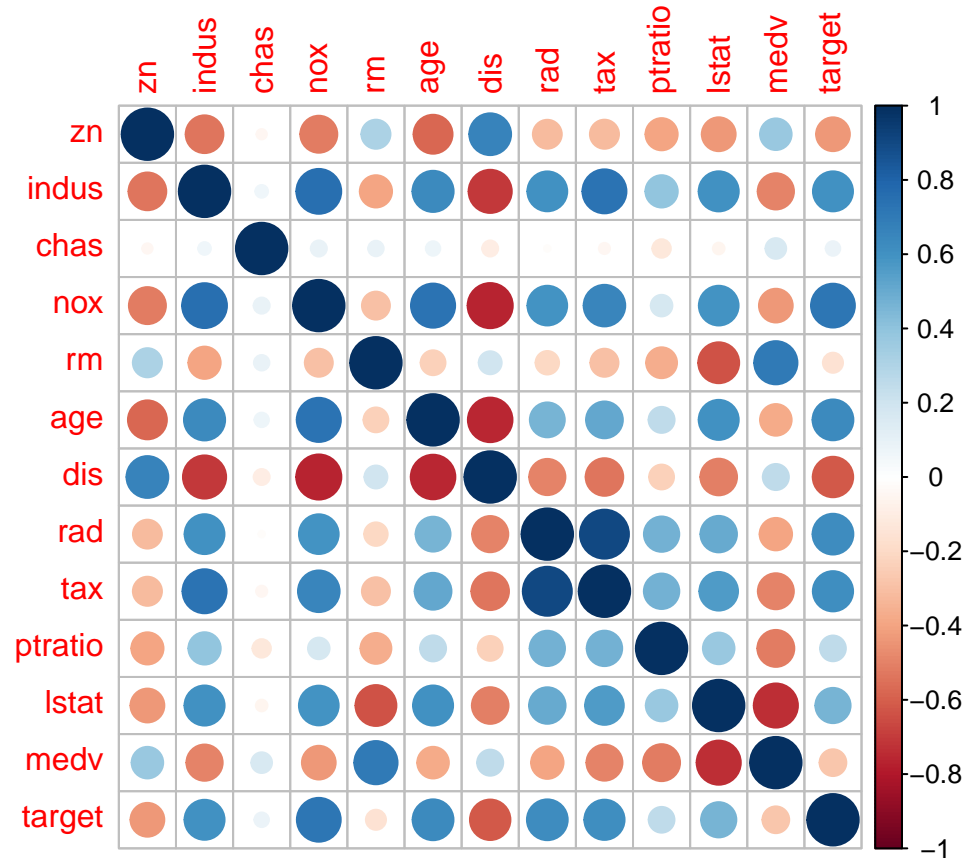## Means of Various Features



## Correlation

We can see our correlation matrix below. A dark blue circle represents a strong positive relationship and a dark red circle represents a strong negative relationship between two variables. We can see that `indus`, `nox`, `target`, and `dis` have the most colinearity. Likewise, these vectors are the best predictors for the target value. Note that this plot only includes rows tha have data in each column.

**EDIT QUESTION: SEEING AS THERE ARE NO NA VALUES, CAN WE JUST MENTION THAT FIRST THEN DO ONE CORRPLOT? DOING TWO FOR MISSING ROW VALUES SEEMS REDUNDANT**

We can compare the plot above to the one below, which includes rows without all of the data present. The availability of data does not significantly affect the results.
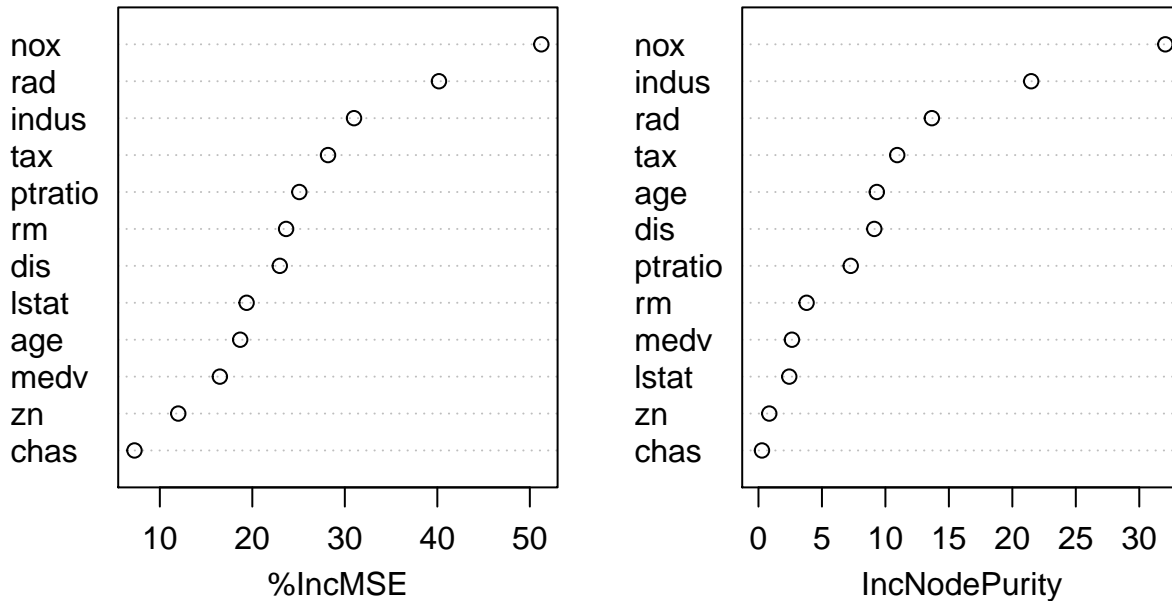
We can explore how many `NAs` are in each column to see if we need to impute any of the variables:

| zn | indus | chas | nox | rm | age | dis | rad | tax | ptratio | lstat | medv | target |
|-----|-------|------|-----|-----|-----|-----|-----|-----|---------|-------|------|--------|
| 466 | 466 | 466 | 466 | 466 | 466 | 466 | 466 | 466 | 466 | 466 | 466 | 466 |

As we can see, each data vector has the same number of entries, 466. Imputation will not be necessary. Finally, we can use the `randomforest` package to verify our assumptions from the correlation plot.

fit



We verified our assumptions above using 1000 random forests. The `nox`, `rad`, `indus`, and `tax` have the most effect. While `dis`is strongly colinear, it has less effect on the target. This is likely due to it encoding information stored redundantly in another vector.

# Data Preparation

Describe how you have transformed the data by changing the original variables or creating new variables. If you did transform the data or create new variables, discuss why you did this. Here are some possible transformations.

- [x] Fix missing Values (None!)
- [ ] Create Flags for missing variable
- [ ] Bin Data
- [ ] Transforms (BoxCox, etc)
- [ ] Combine Variables ?

# Build Models

- [ ] 3 binary logistic models
- [ ] forward, stepwise, random forest, etc
- [ ] Inferences
- [ ] Coefficients

# Select Models

- [ ] Use Log Likelihood, AIC, ROC curve,
- [ ] Evaluate Training Set
- [ ] Accuracy, Error, Precision, Sensitivity, Specificity, F1 score, AUC, conf matrix (hint: use assignment 2, and check outthis link )
- [ ] Make predictions with test set and interpret