

# HW 3

*Team 2*

*April 10, 2019*

## Contents

<b>Overview</b>	<b>1</b>
Objective . . . . .	2
Dependencies . . . . .	2
<b>Data Exploration</b>	<b>2</b>
Summary Statistics . . . . .	2
Correlation . . . . .	5
<b>Data Preparation</b>	<b>7</b>
Variable Exploration . . . . .	7
<b>rad</b> . . . . .	9
<b>indus</b> . . . . .	11
<b>tax</b> . . . . .	13
Variable Transformation . . . . .	15
<b>Build Models</b>	<b>23</b>
<b>Select Models</b>	<b>23</b>

## 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

Replication of our work requires the following packages in Rstudio:

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

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

## 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.

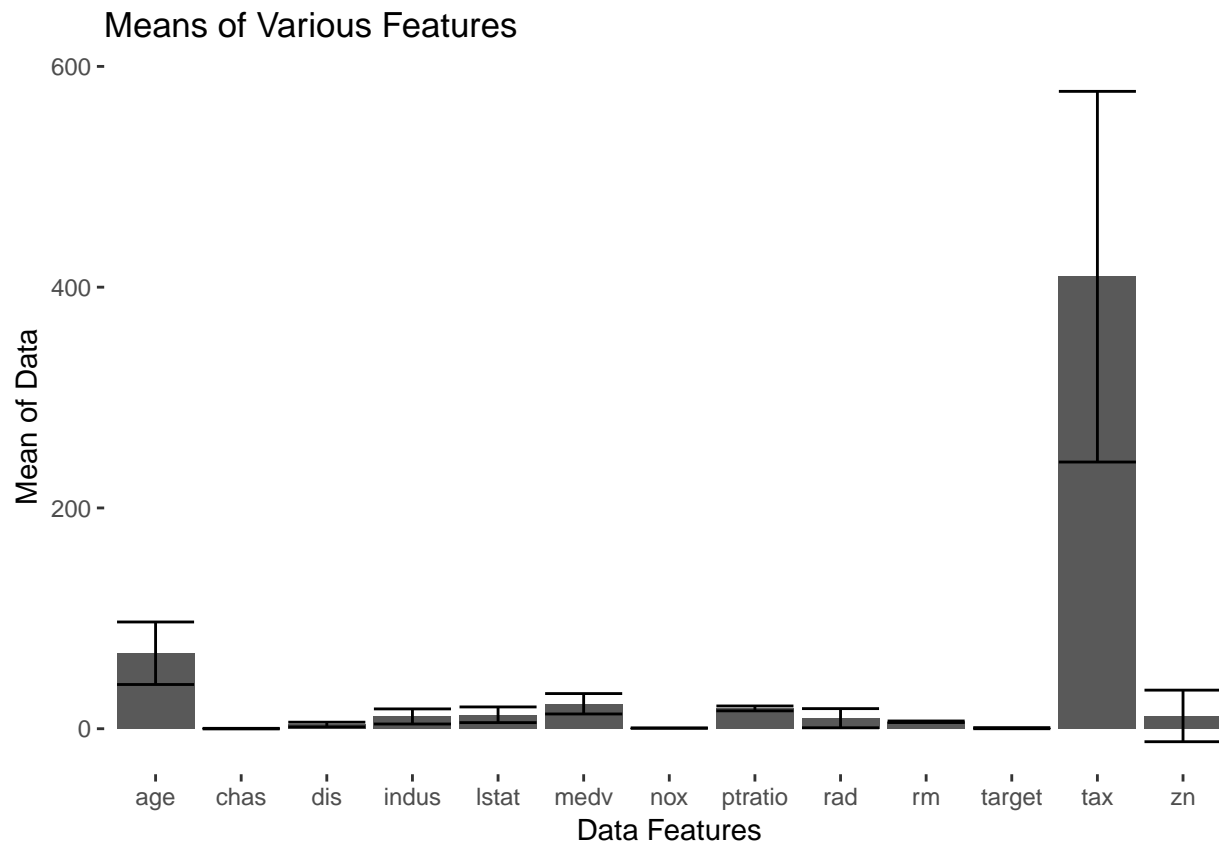
zn	indus	chas	nox	rm	age	dis	rad	tax	ptratio	lstat	medv	target
0	19.58	0	0.605	7.929	96.2	2.0459	5	403	14.7	3.70	50.0	1
0	19.58	1	0.871	5.403	100.0	1.3216	5	403	14.7	26.82	13.4	1
0	18.10	0	0.740	6.485	100.0	1.9784	24	666	20.2	18.85	15.4	1
30	4.93	0	0.428	6.393	7.8	7.0355	6	300	16.6	5.19	23.7	0
0	2.46	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

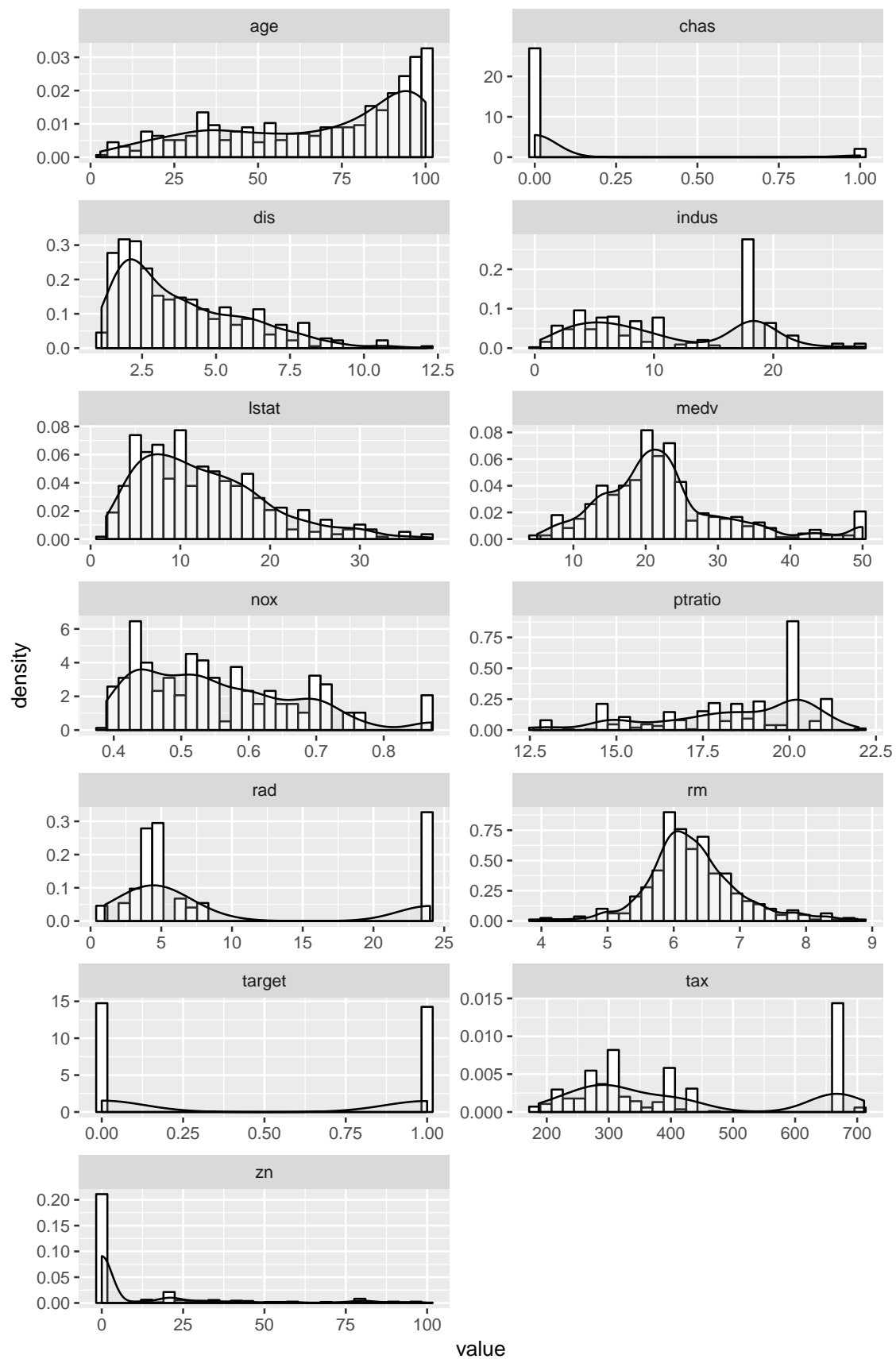
Below is a bar chart that illustrates the average and standard deviation for each of our data vectors. As we can see, the **tax** vector is a totally different magnitude than the rest. Models involving this vector will benefit from normalization or scaling.



## Histogram

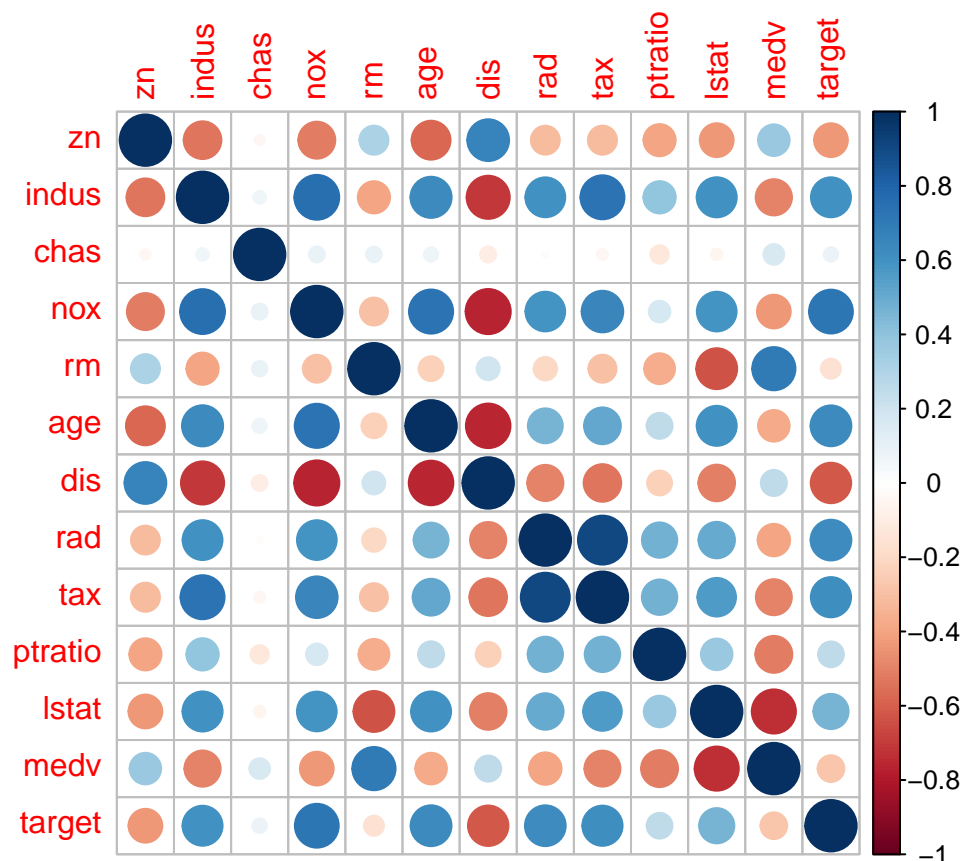
The following histograms help visualize the spread and skewness of the raw data.

```
ggplot(data = gather(training), mapping = aes(x = value)) +
  geom_histogram(aes(y=..density..), colour="black", fill="white")+
  geom_density(alpha=.2, fill="lightgrey")+
  facet_wrap(~key, ncol = 2, scales = 'free')
```



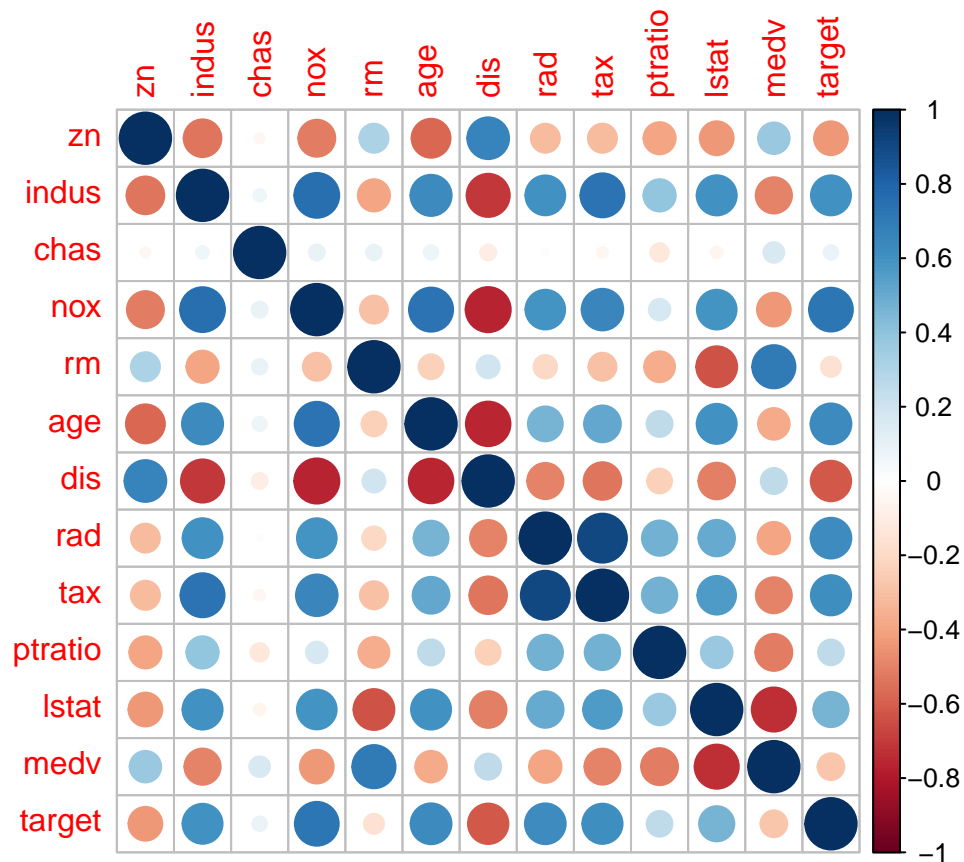
## Correlation

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 -jm**

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.

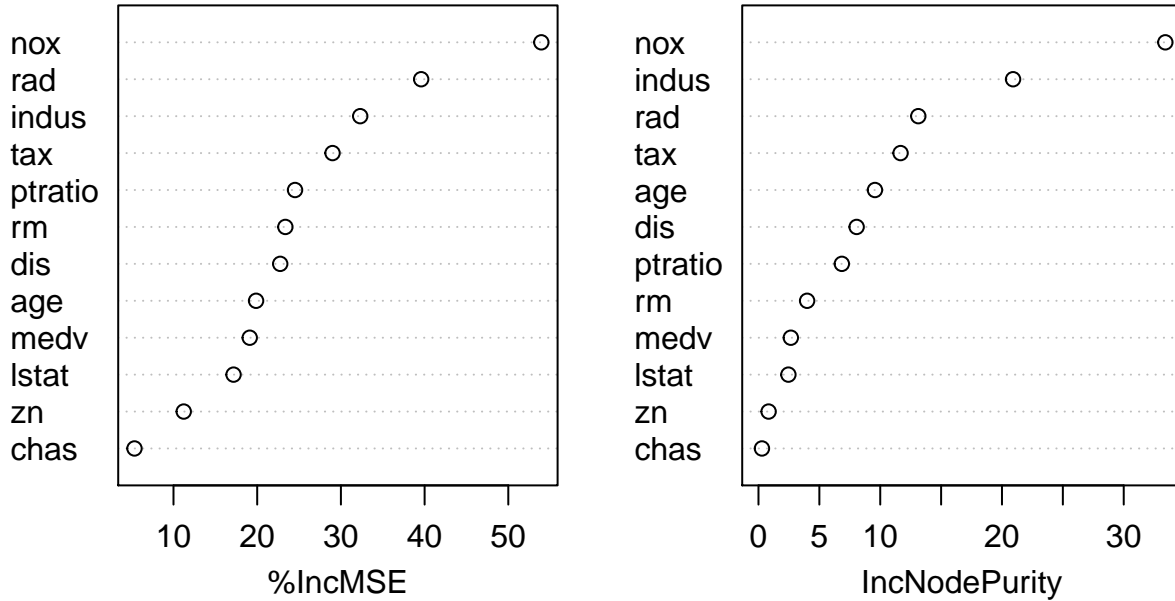


We can explore how many NAs are in each column to see if we need to impute any of the variables:

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

In the exploration section, we identified that there were no missing values in the dataset would affect the outcome of our model. We additionally found that the following four variables had the strongest correlation with our target goal:

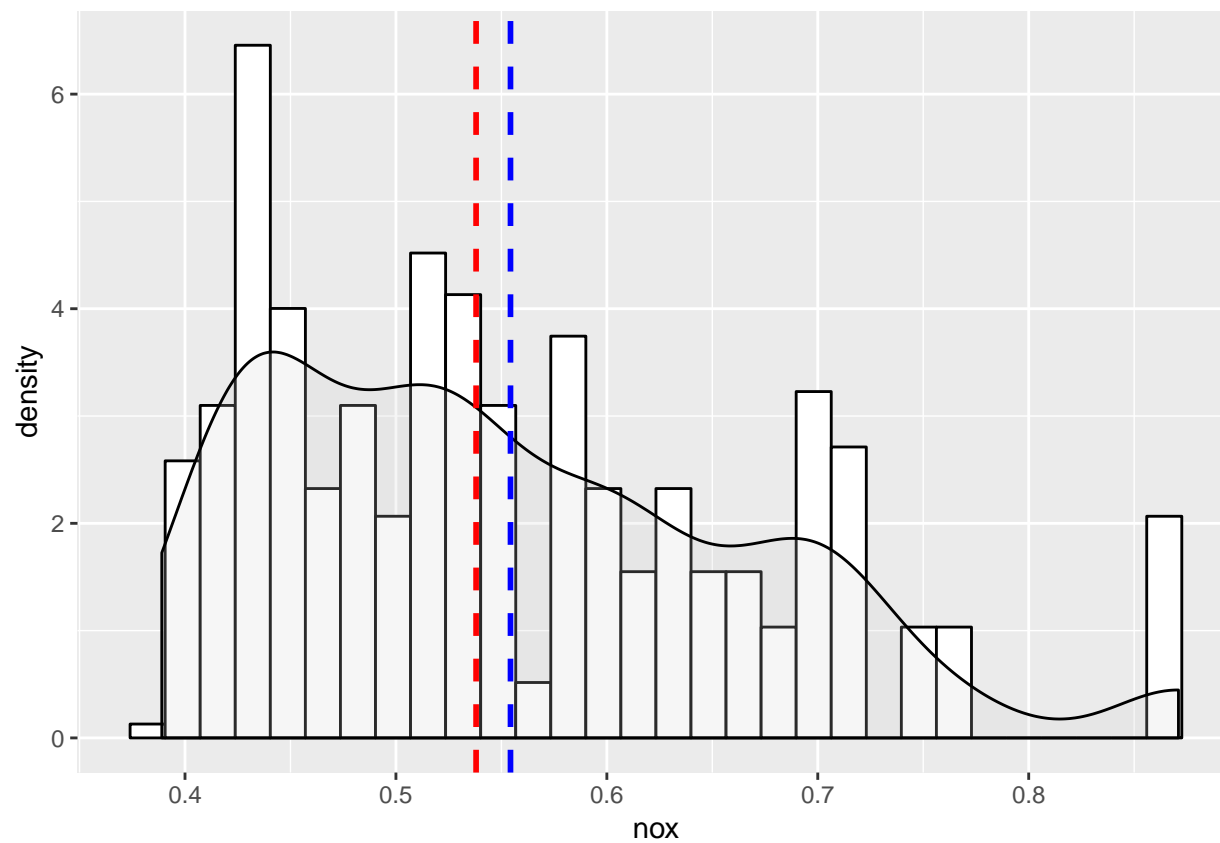
1. `nox`
2. `rad`
3. `indus`
4. `tax`

In the following section, we will analyze and transform these variables to use in the development of our model:

## Variable Exploration

### `nox` variable

This variable represents nitrogen oxides concentration (parts per 10 million). The plots below show the `nox` data is multimodal, skewed right, and centered around 0.55. The red line shows the median, whereas the blue line depicts the mean value for this variable.

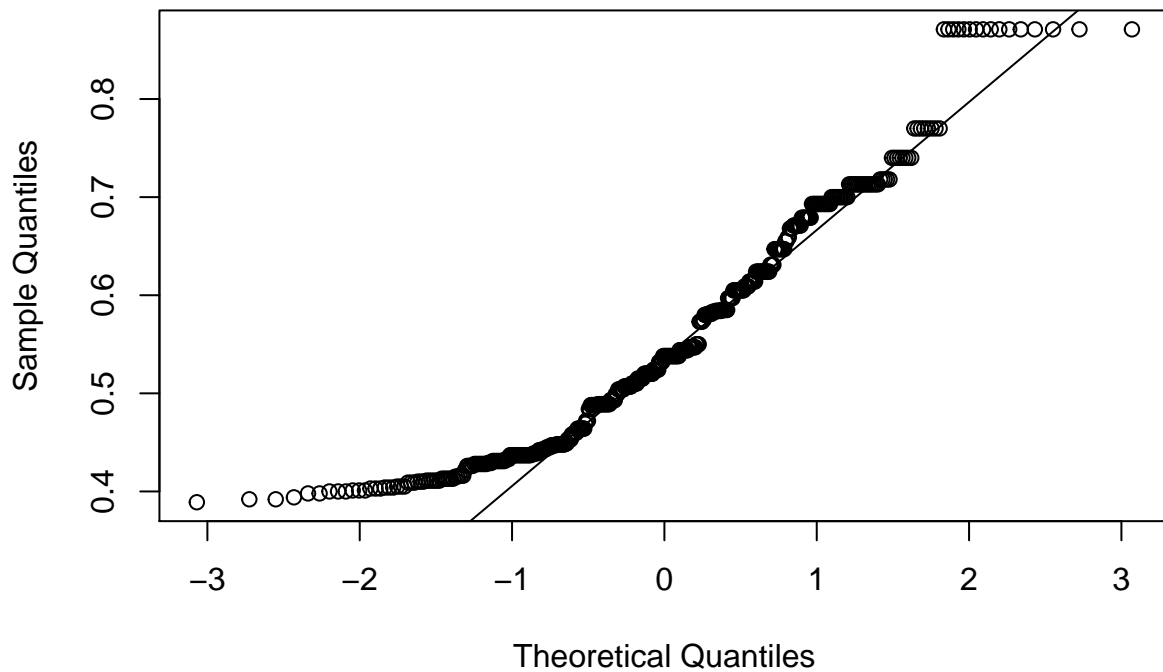


```
##      Min. 1st Qu.  Median    Mean 3rd Qu.    Max.
## 0.3890 0.4480 0.5380 0.5543 0.6240 0.8710
```

The qqplot below confirms that this variable does not follow a normal distribution.



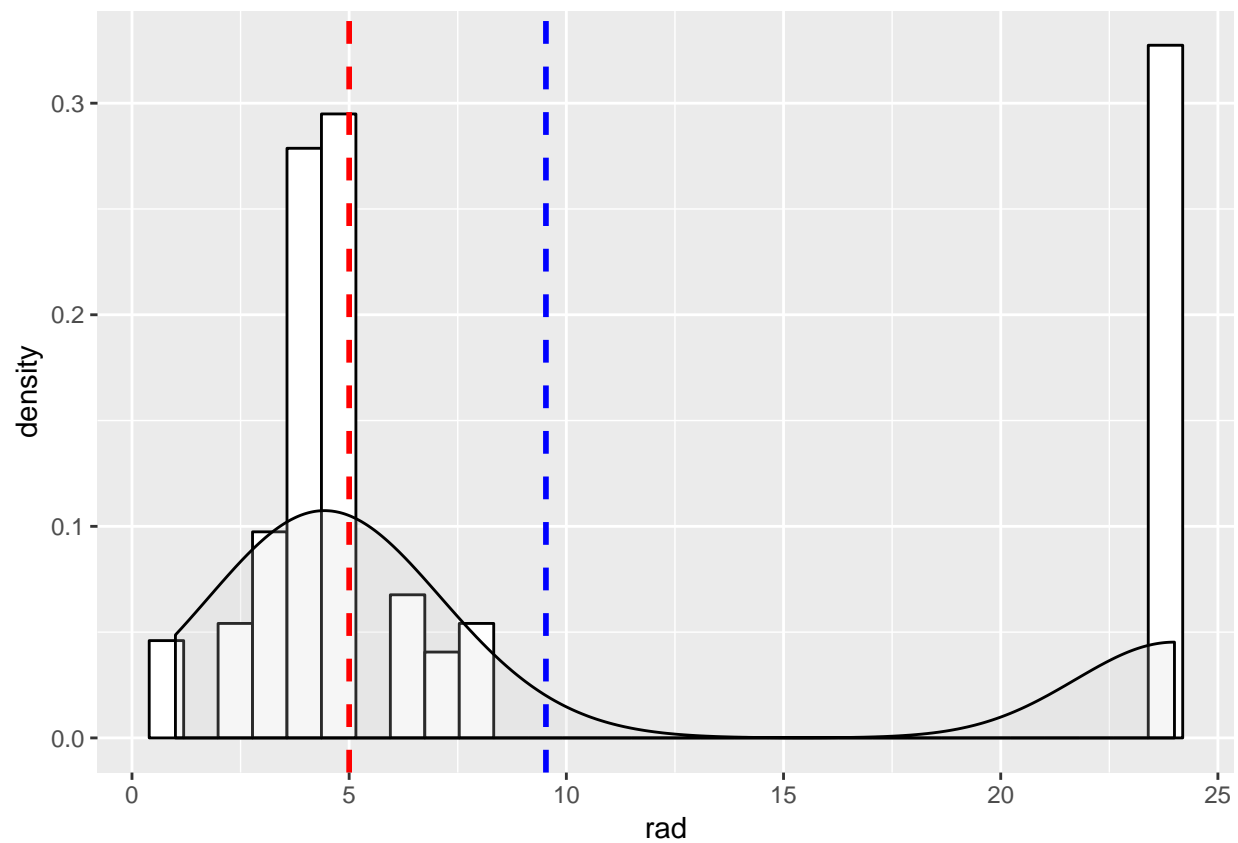
## Normal Q-Q Plot



```
stud_teach_ratio <- training$ptratio  
stud_teach_ratio <- cut(stud_teach_ratio, breaks = 3, labels=c("Small", "Medium", "Large"))
```

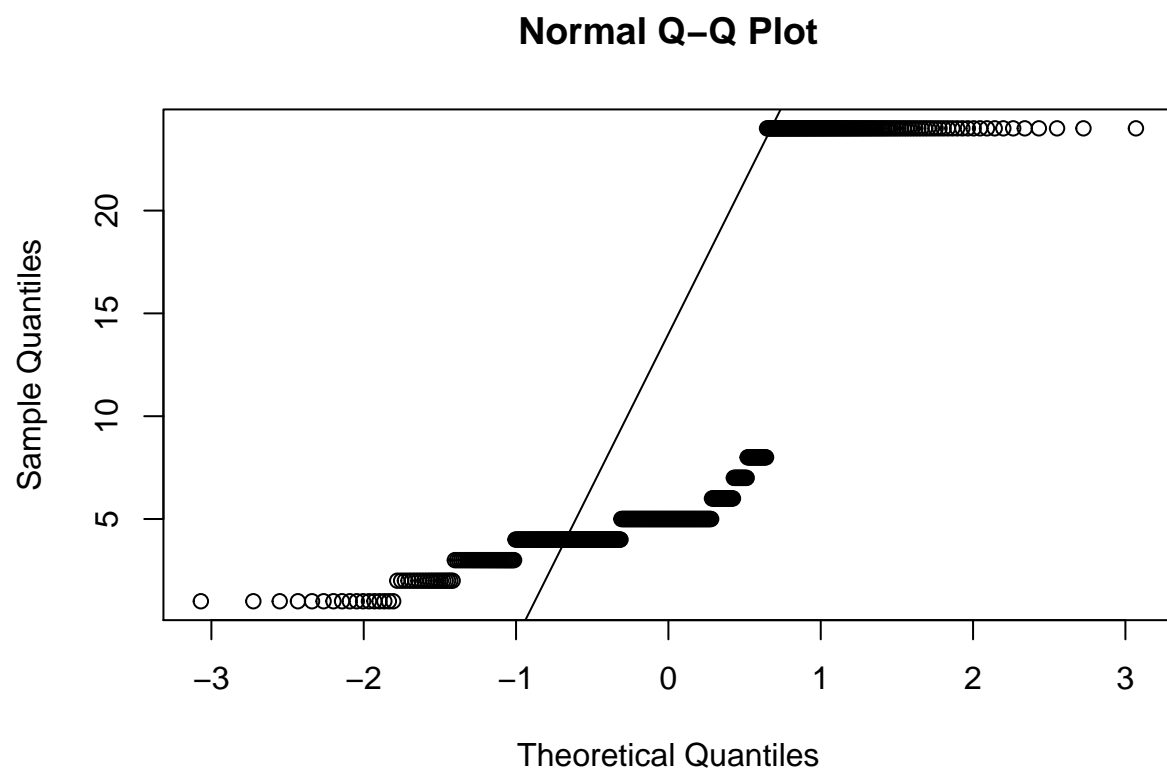
### rad

This is an index variable that represents accessibility to radial highways. The plots below show the **rad** data is bimodal and centered around 5. This data almost follows a normal distribution, however there is an extreme outlier that skews the mean to the right. The red line shows the median, whereas the blue line depicts the mean value for this variable.



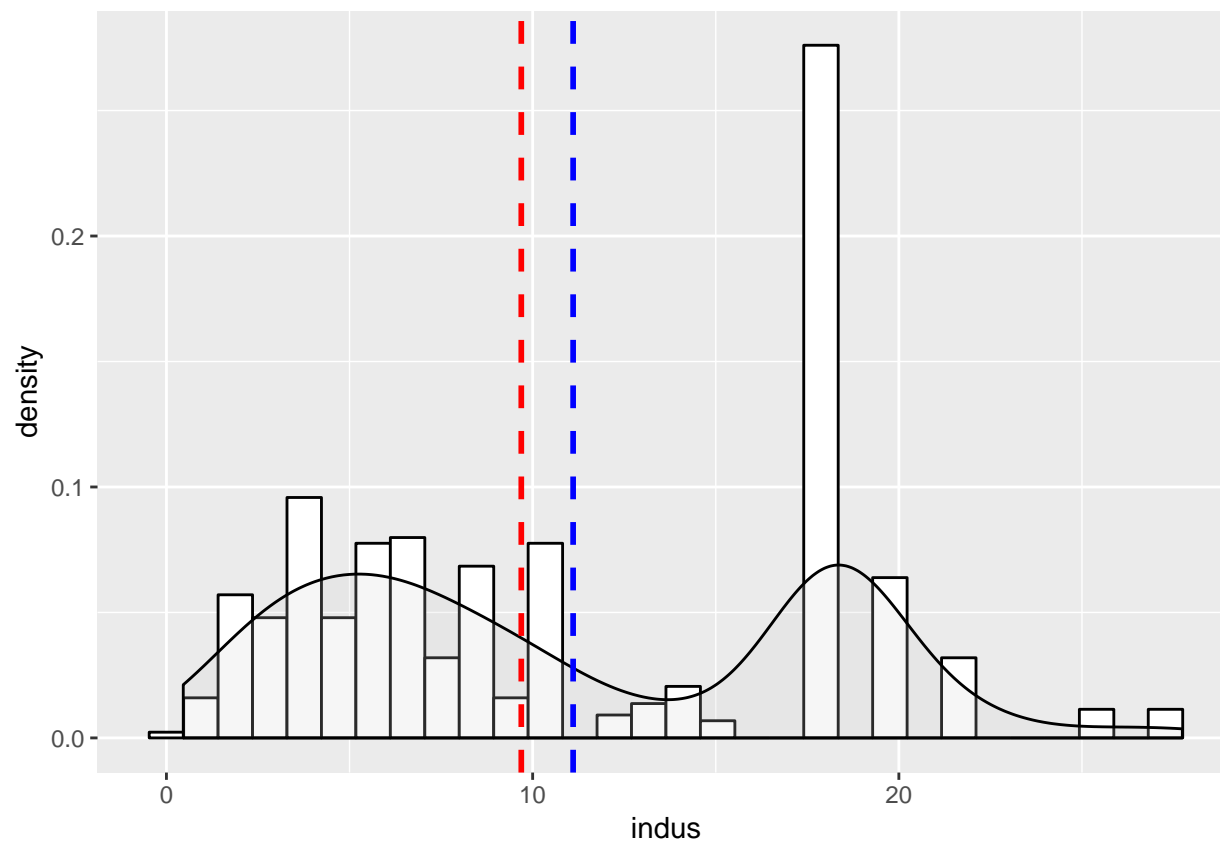
##	Min.	1st Qu.	Median	Mean	3rd Qu.	Max.
##	1.00	4.00	5.00	9.53	24.00	24.00

The `qqplot` below confirms that this variable does not follow a normal distribution. TODO: Make this a binary variable +/- the median.



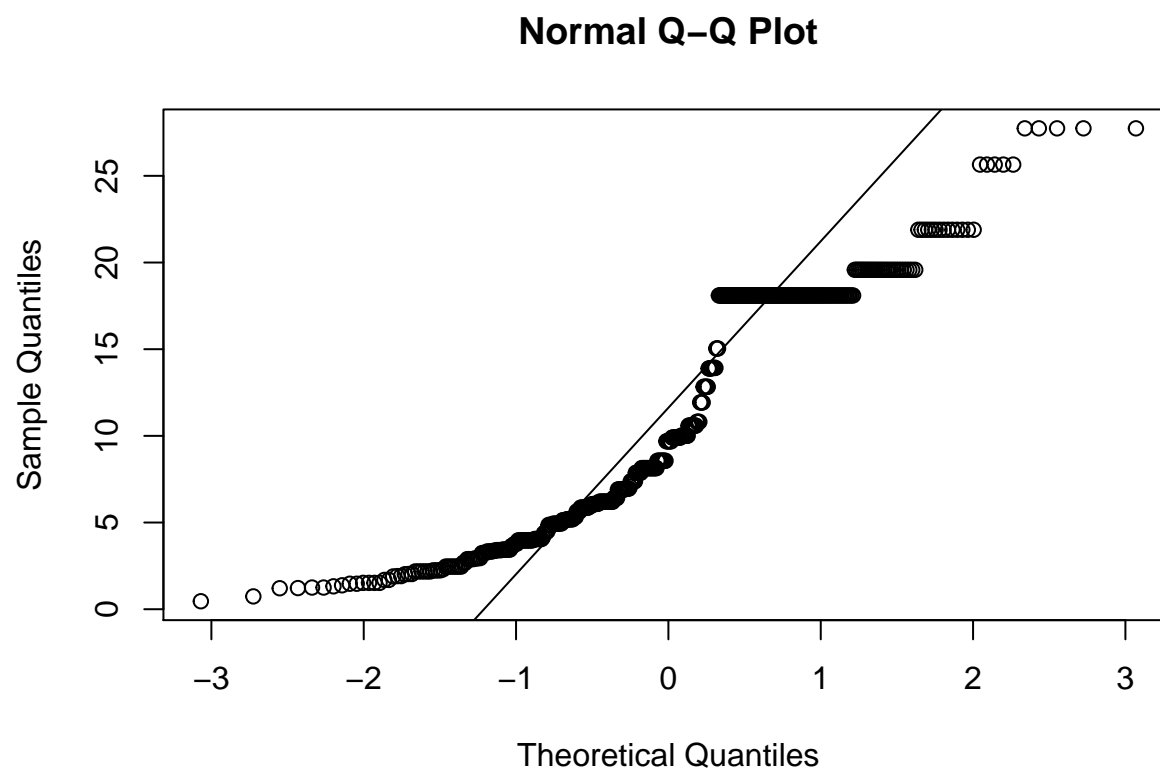
#### **indus**

This variable represents the proportion of non-retail business acres per suburb. The plots below show the **indus** data is bimodal, skewed right, and centered around 10. The red line shows the median, whereas the blue line depicts the mean value for this variable.



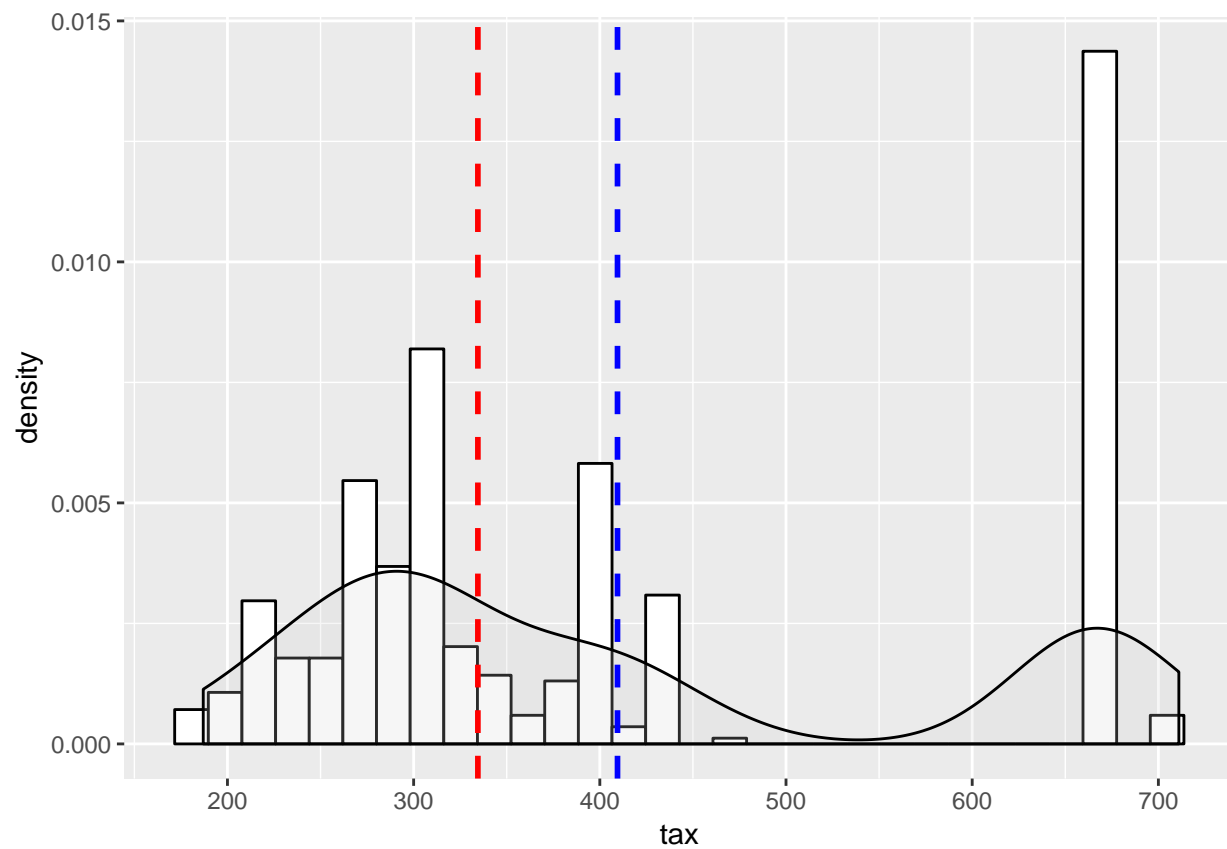
```
##      Min. 1st Qu.  Median    Mean 3rd Qu.    Max.
##  0.460   5.145   9.690  11.105  18.100  27.740
```

The `qqplot` below confirms that this variable does not follow a normal distribution. TODO: maybe this could be a bifucated variable too. See the two different normal-ish distributions?



#### **tax**

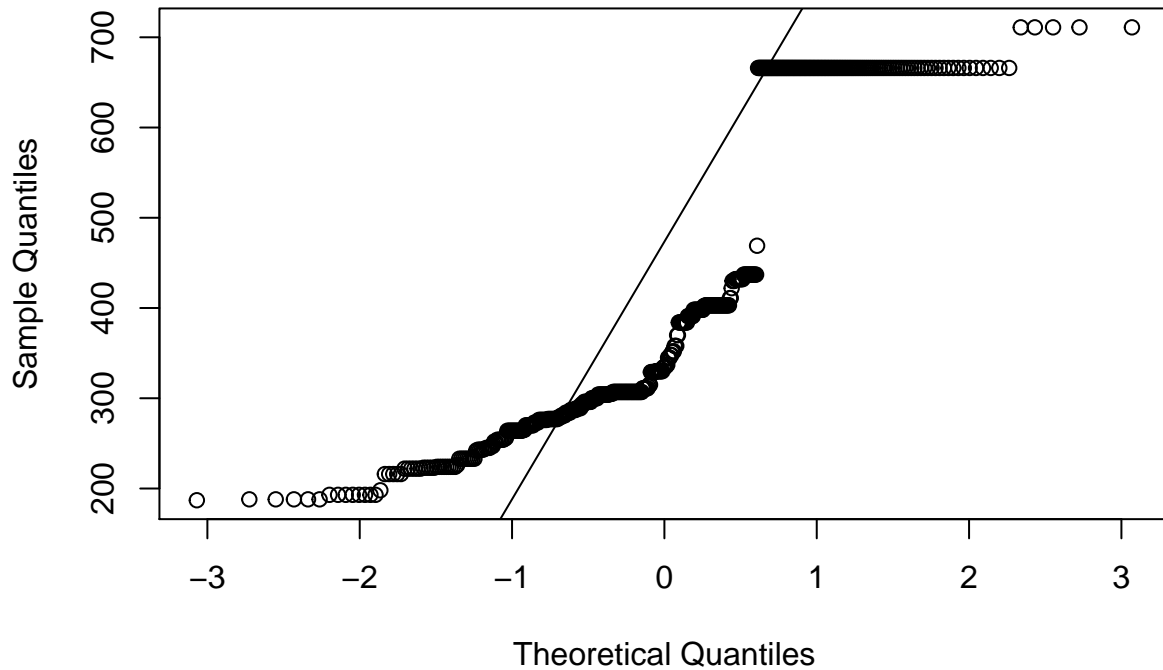
This variable represents full-value property-tax rate per \$10,000. The values stored in this variable are significantly larger than the ones previously explored. The plots below show that the `tax` data is bimodal, skewed right, and centered around 330. The red line shows the median, whereas the blue line depicts the mean value for this variable.



```
##      Min. 1st Qu.  Median      Mean 3rd Qu.      Max.
##    187.0   281.0   334.5   409.5   666.0   711.0
```

The qqplot below confirms that this variable does not follow a normal distribution.

## Normal Q-Q Plot



## Variable Transformation

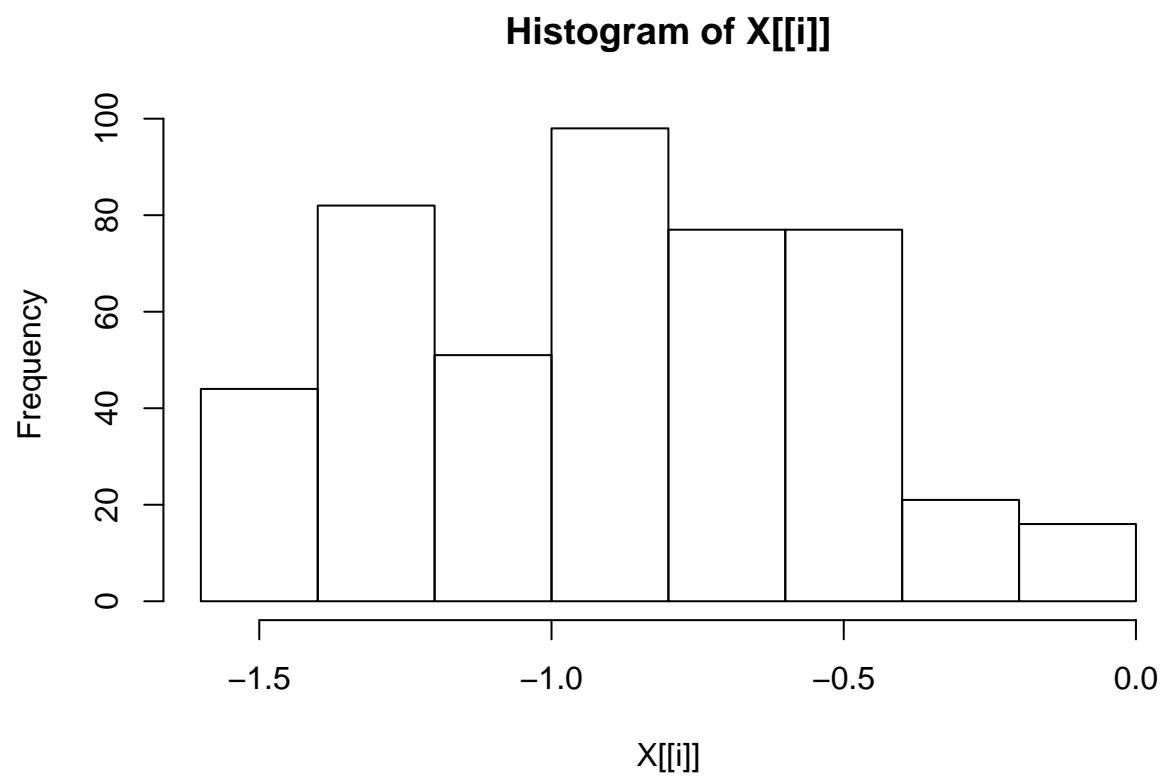
Our goal is to transform the following data

```
data <- training %>% select(nox, rad, indus, tax)
data %>% slice(1:5) %>% kable() %>% kable_styling()
```

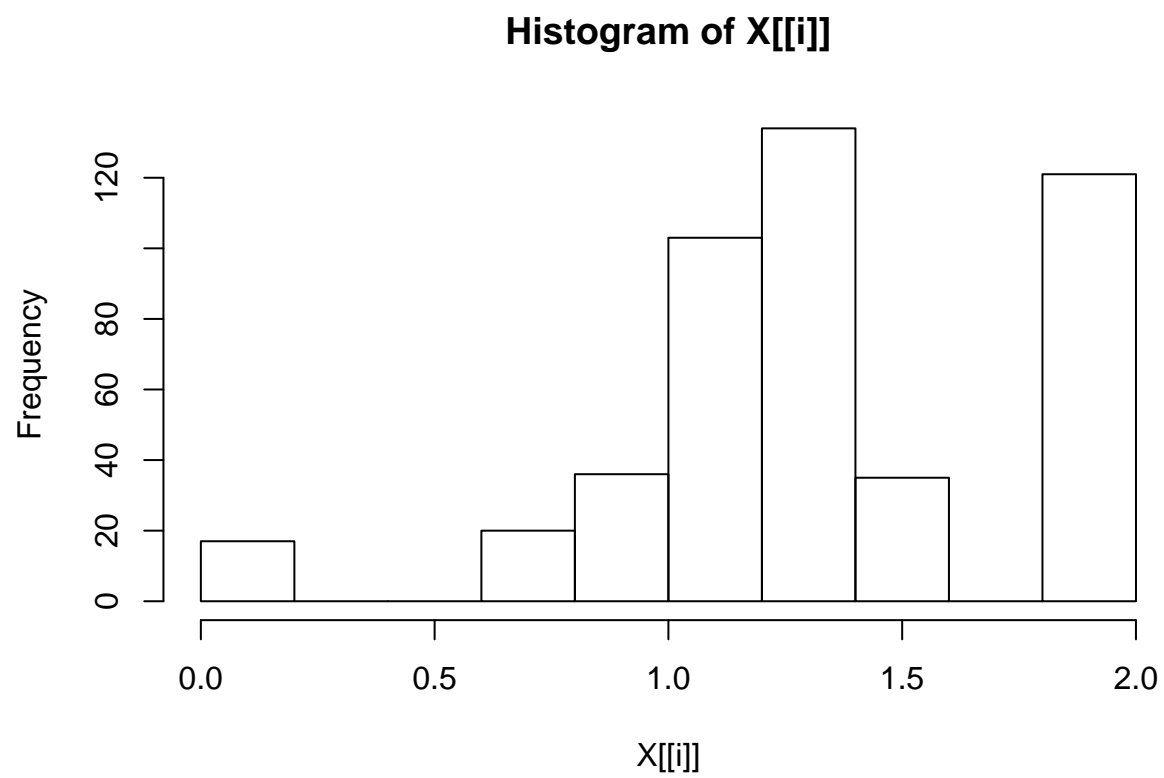
nox	rad	indus	tax
0.605	5	19.58	403
0.871	5	19.58	403
0.740	24	18.10	666
0.428	6	4.93	300
0.488	3	2.46	193

We can fix this with a Box-Cox transformation, using the `forecast` package in R. NOTE: Nox looks pretty good. The others definitely need to be transformed in some other way (or not at all), since the results are pretty not-normal.

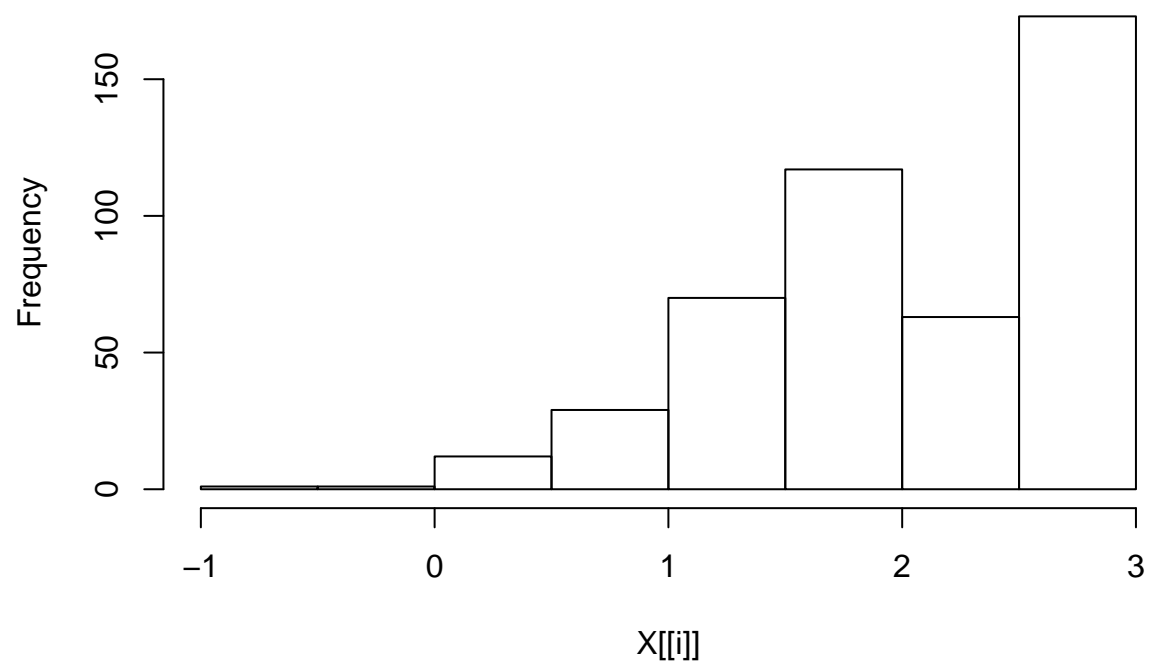
```
data2 <- data
transform <- function(x){
  x <- BoxCox(x, lambda = BoxCox.lambda(x))
}
data2 <- as.data.frame(sapply(data, transform))
sapply(data2, hist)
```

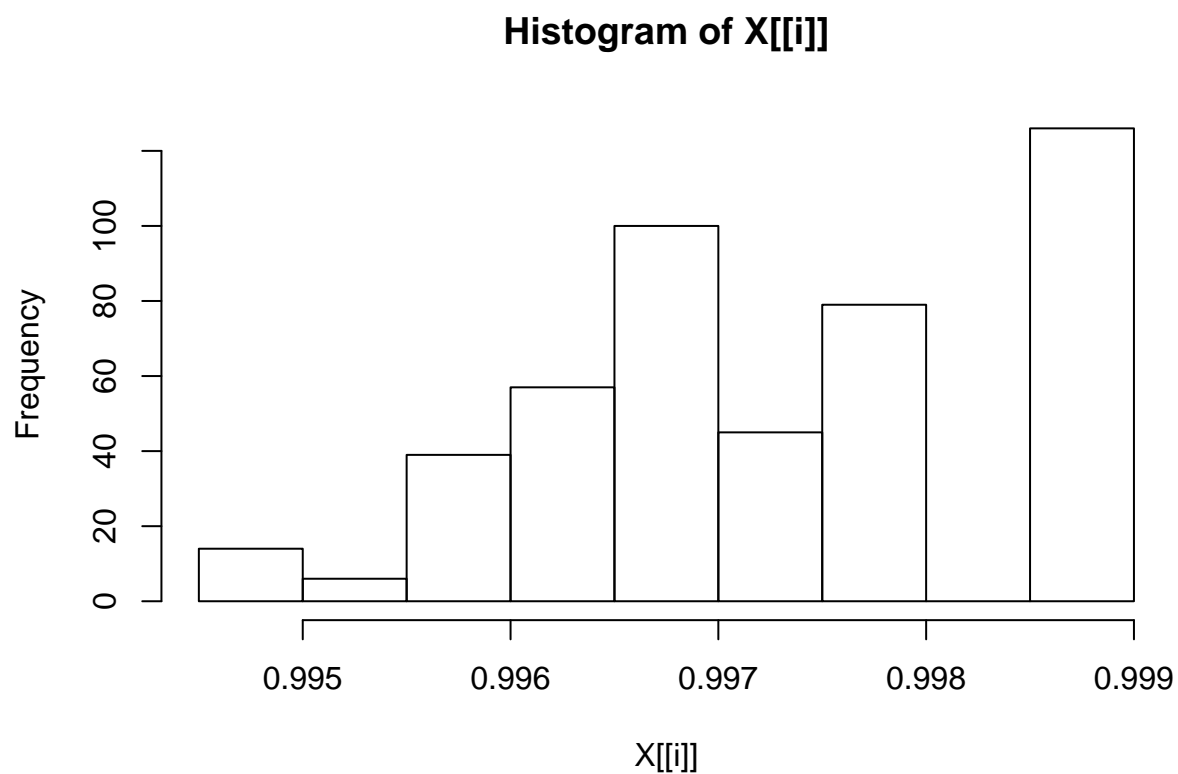






**Histogram of  $X[[i]]$**

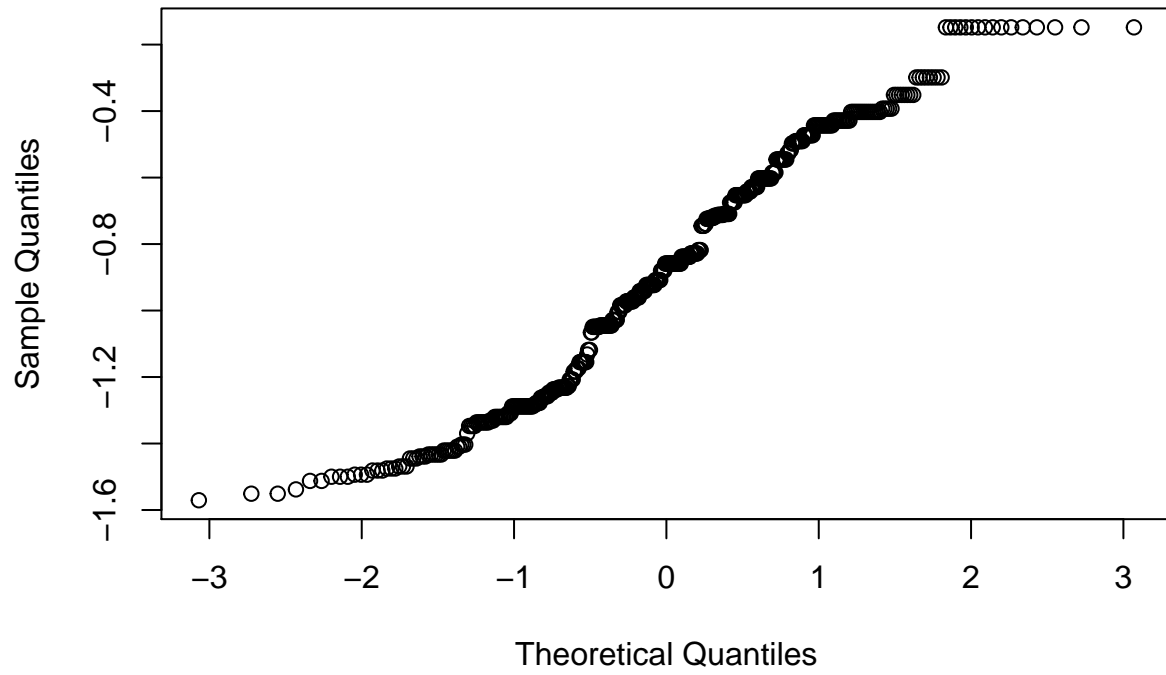


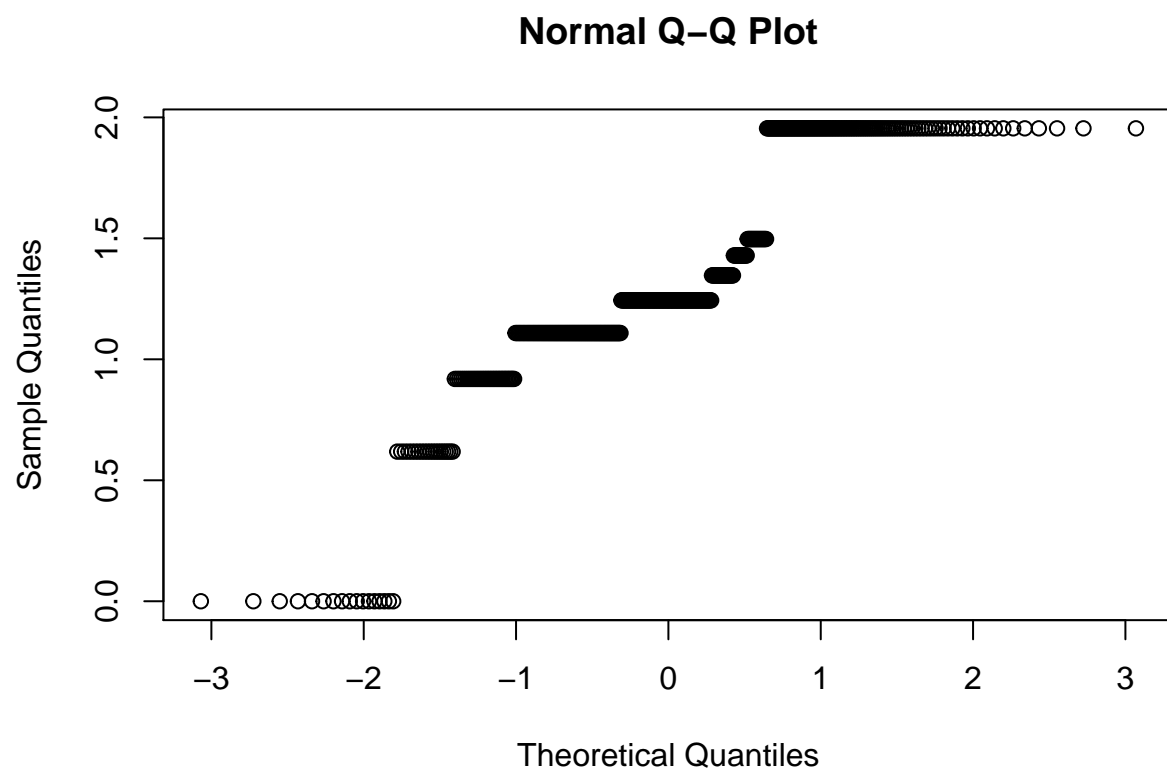


```
##          nox      rad      indus      tax
## breaks  Numeric,9 Numeric,11 Numeric,9 Numeric,10
## counts  Integer,8 Integer,10 Integer,8 Integer,9
## density  Numeric,8 Numeric,10 Numeric,8 Numeric,9
## mids     Numeric,8 Numeric,10 Numeric,8 Numeric,9
## xname    "X[[i]]"  "X[[i]]"  "X[[i]]"  "X[[i]]"
## equidist TRUE      TRUE      TRUE      TRUE
```

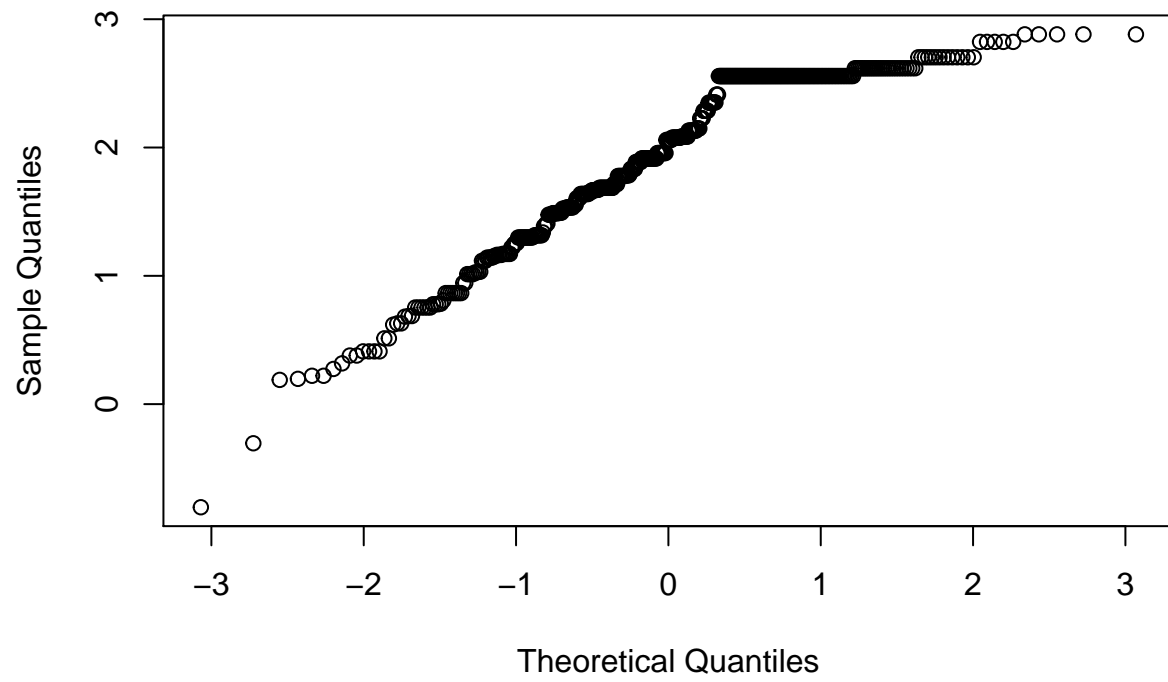
```
sapply(data2, qqnorm)
```

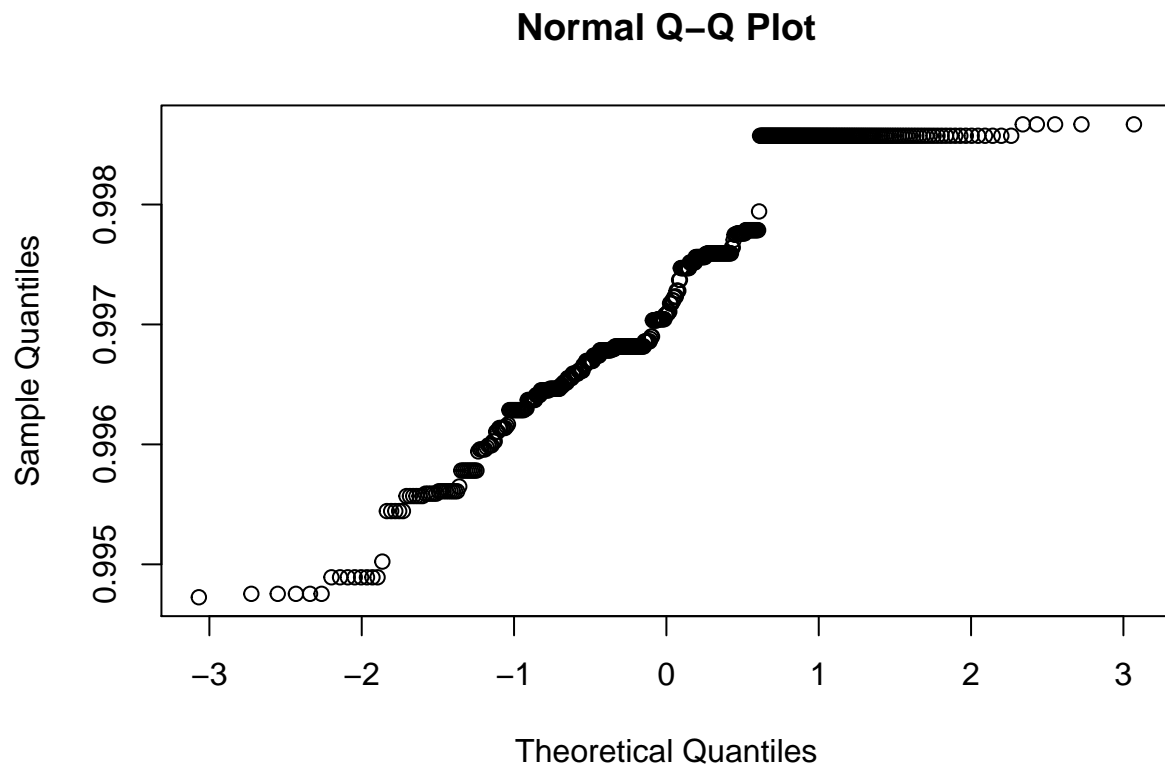
Normal Q-Q Plot





Normal Q-Q Plot





```
##   nox      rad      indus      tax
## x Numeric,466 Numeric,466 Numeric,466 Numeric,466
## y Numeric,466 Numeric,466 Numeric,466 Numeric,466
```

## 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 out [this link](#) )
- [ ] Make predictions with test set and interpret