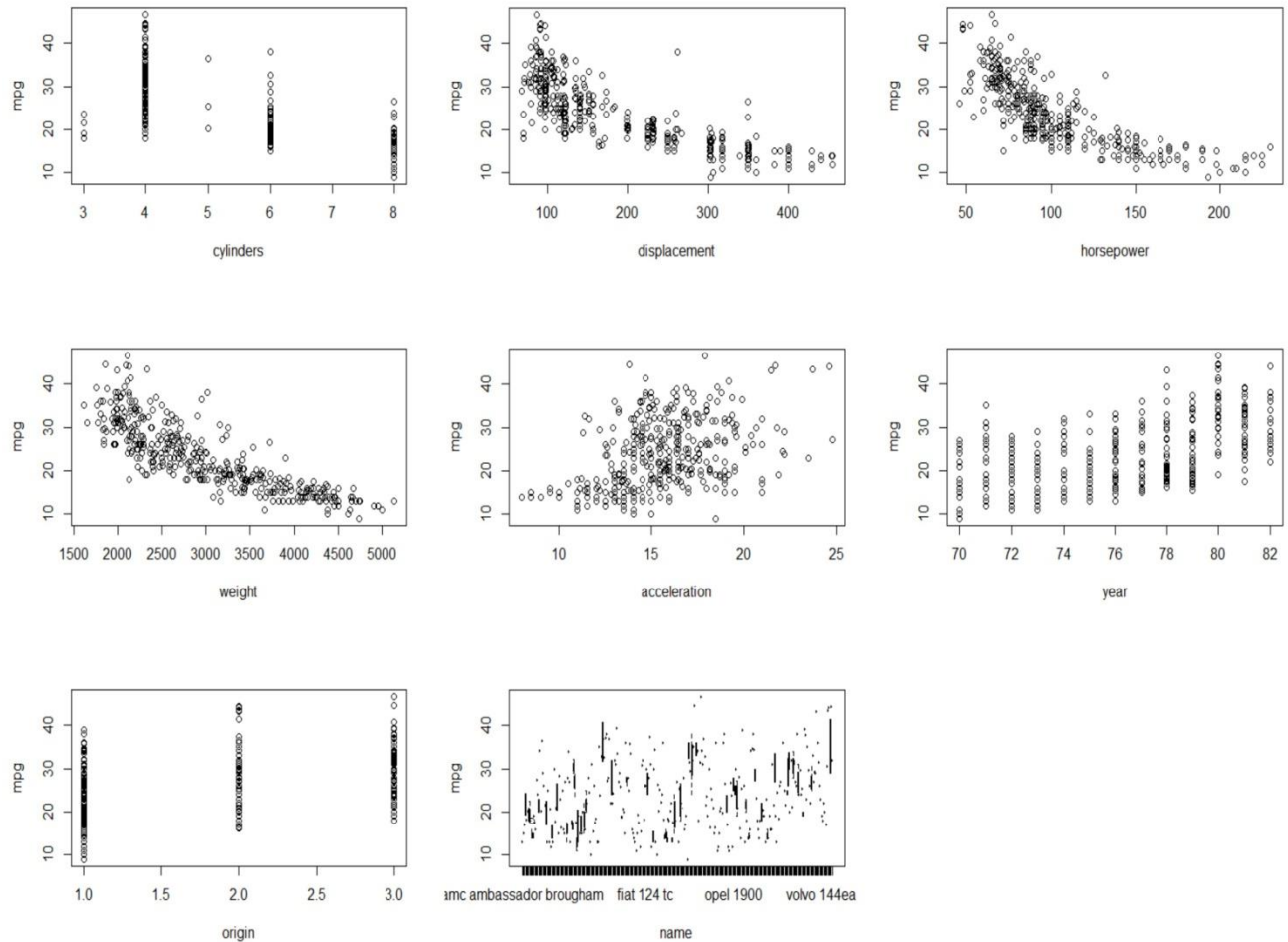Statistical Data Mining I

# Homework 1
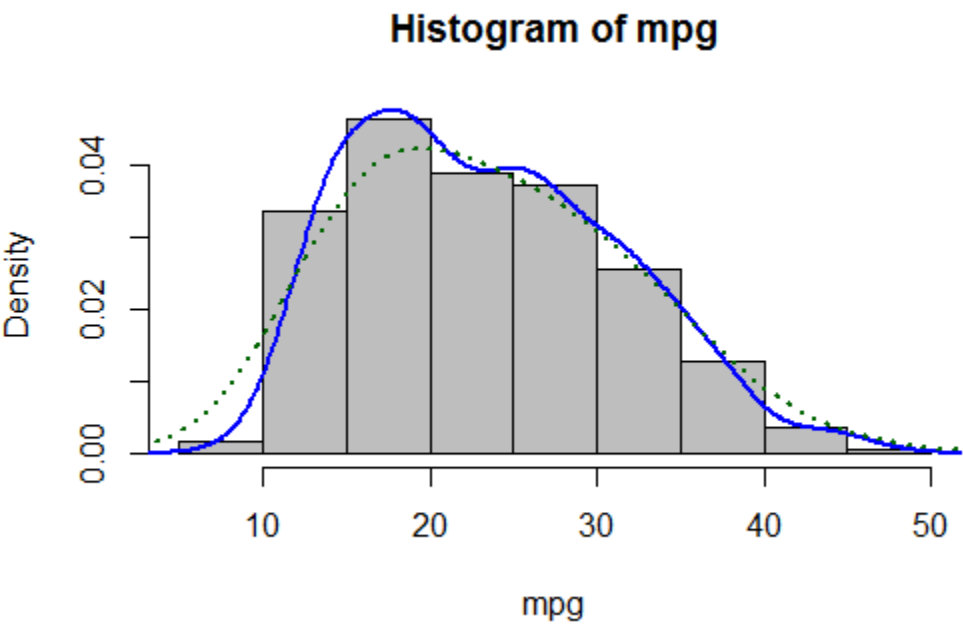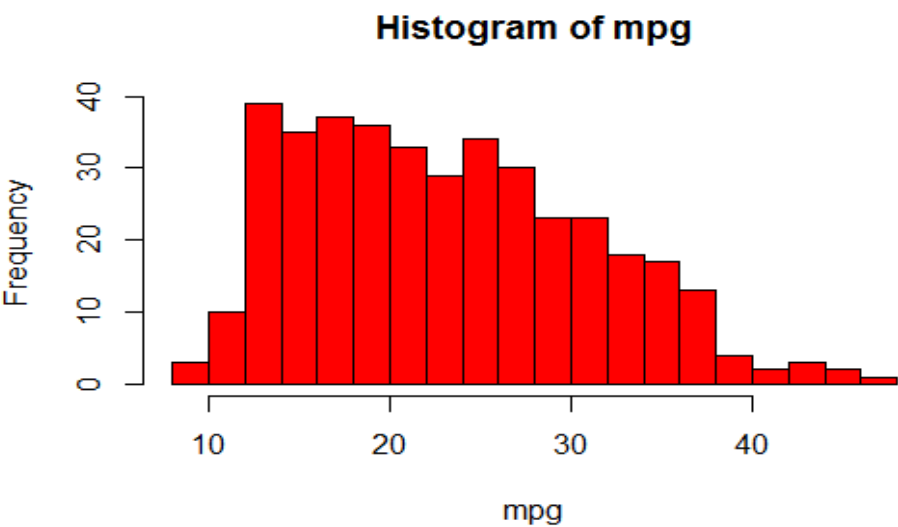
Student Number:71
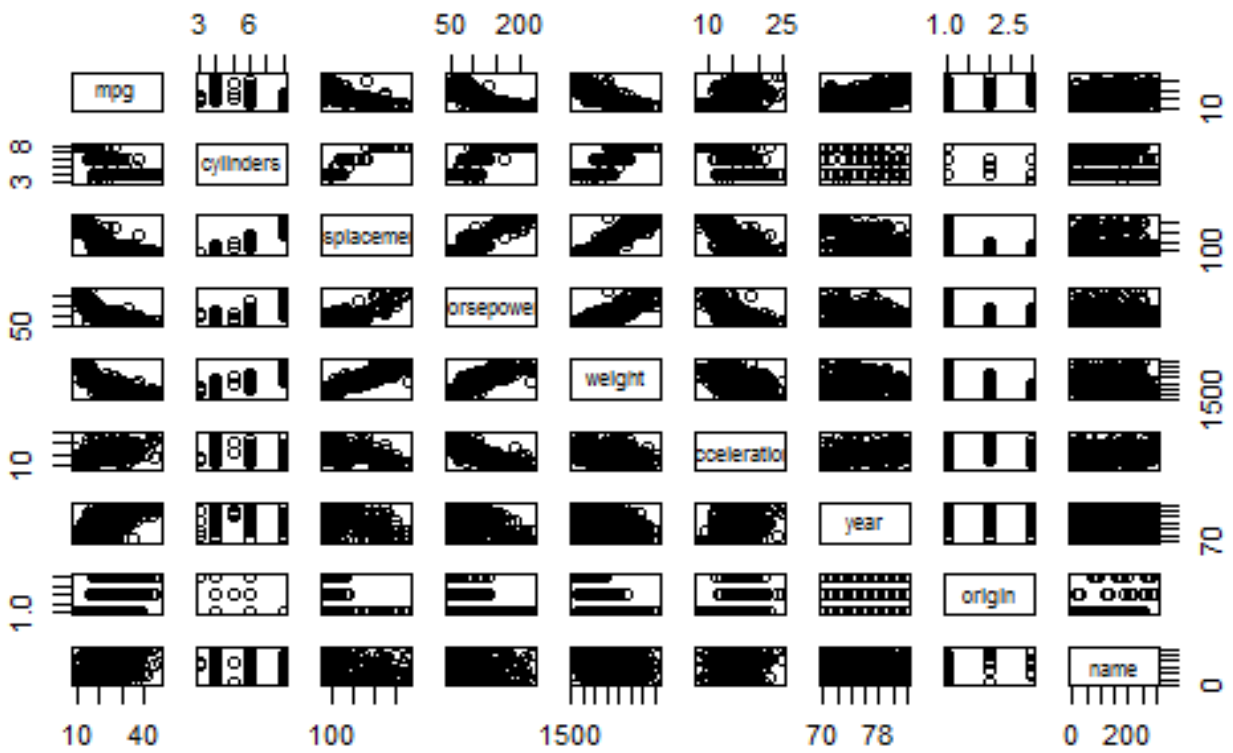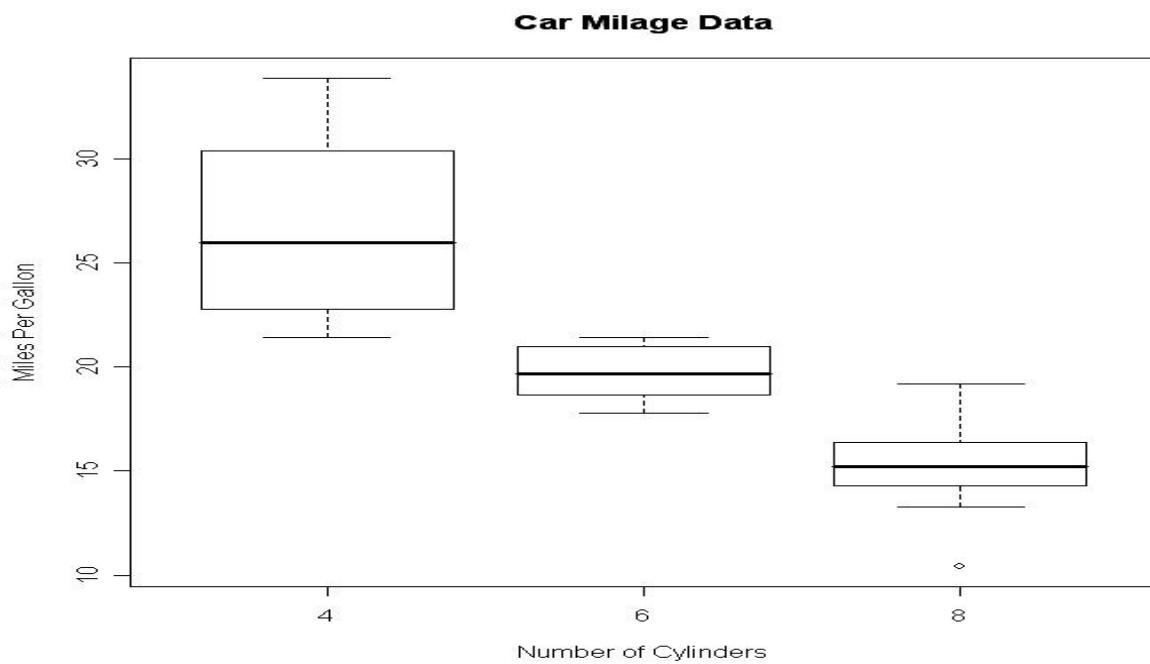
Problem 1:



i)    We seem to get more mileage per gallon on a 4 cylinder vehicle than the others.
ii)   Weight, displacement and horsepower seem to have an inverse effect with mpg. We
      see an overall increase in mpg over the years. Almost doubled in one decade.
iii)   Japanese cars have higher mpg than US or European cars.

Exploratory Data Analysis:
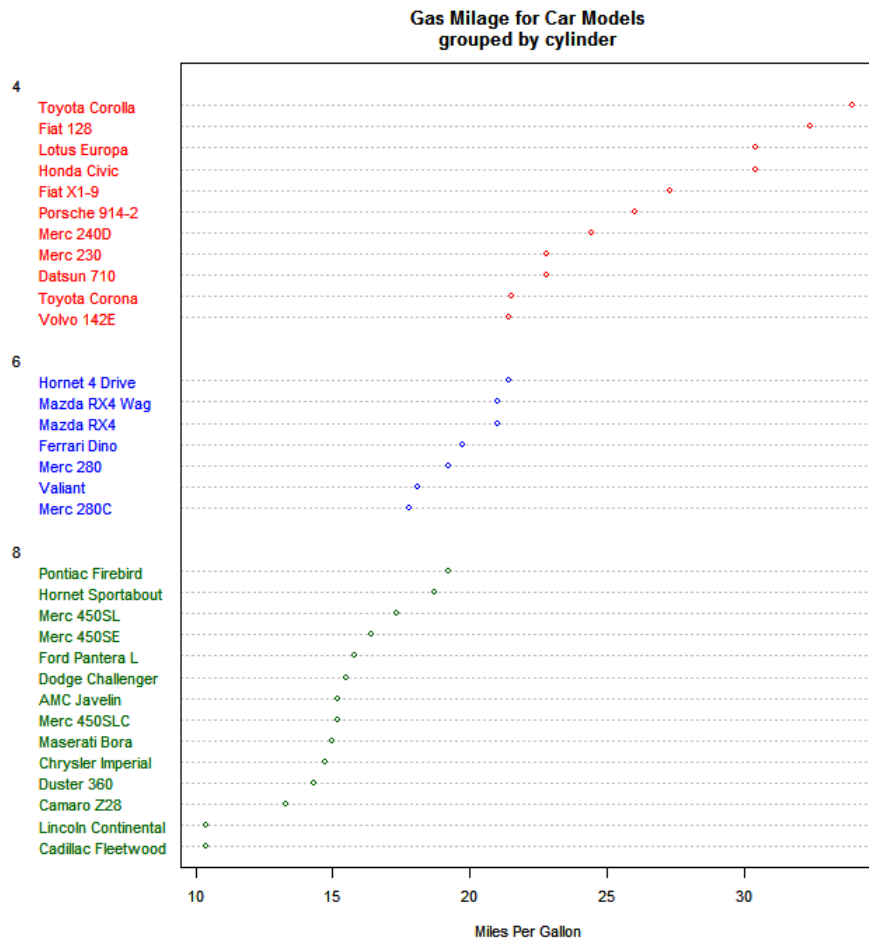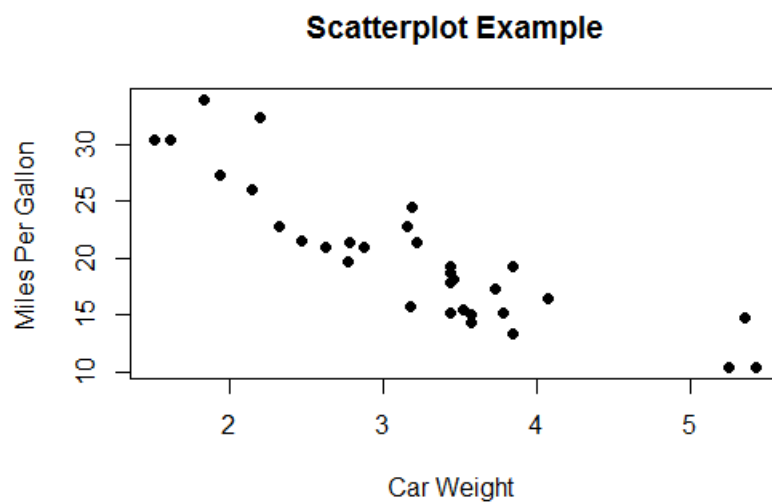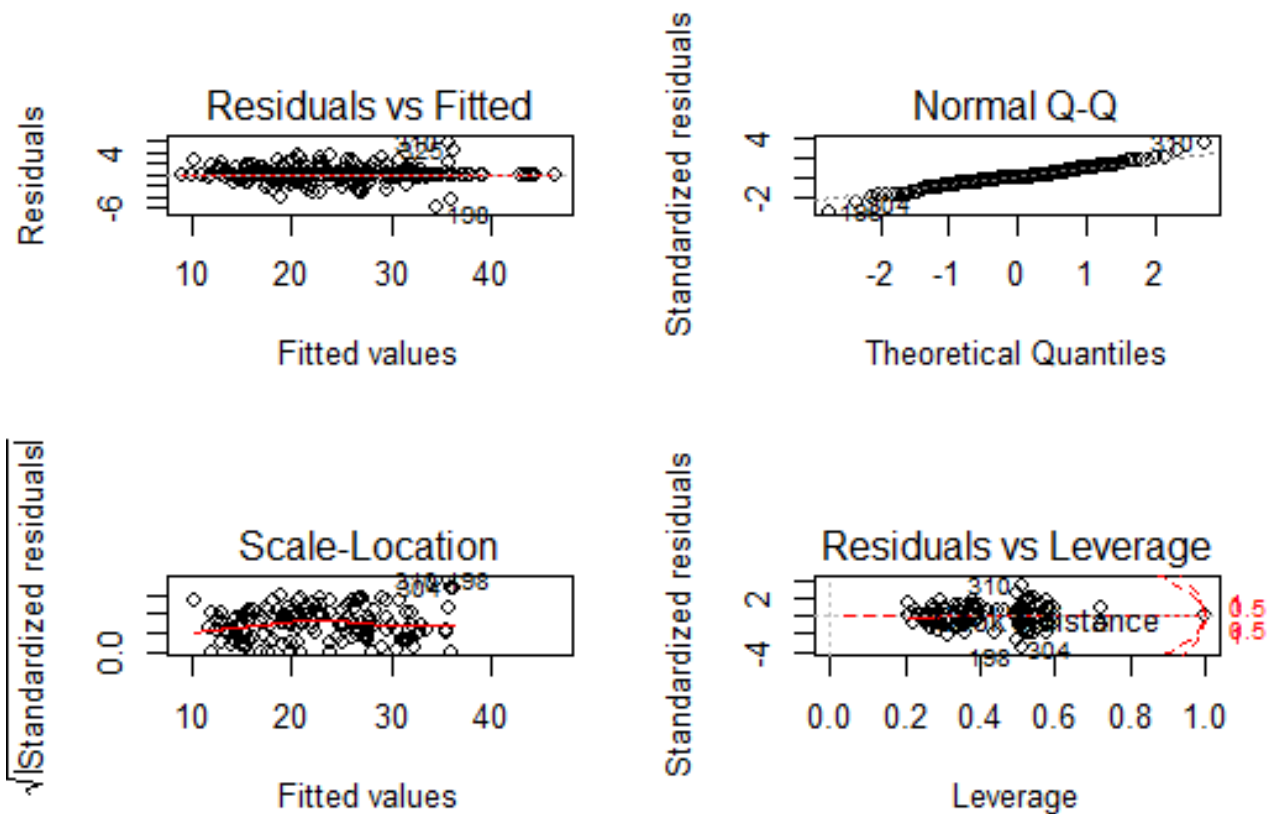
**Histogram of mpg**



**Histogram of mpg**

B)Box Plot



Car Milage Data

C)Dot Plot

**Gas Milage for Car Models grouped by cylinder**
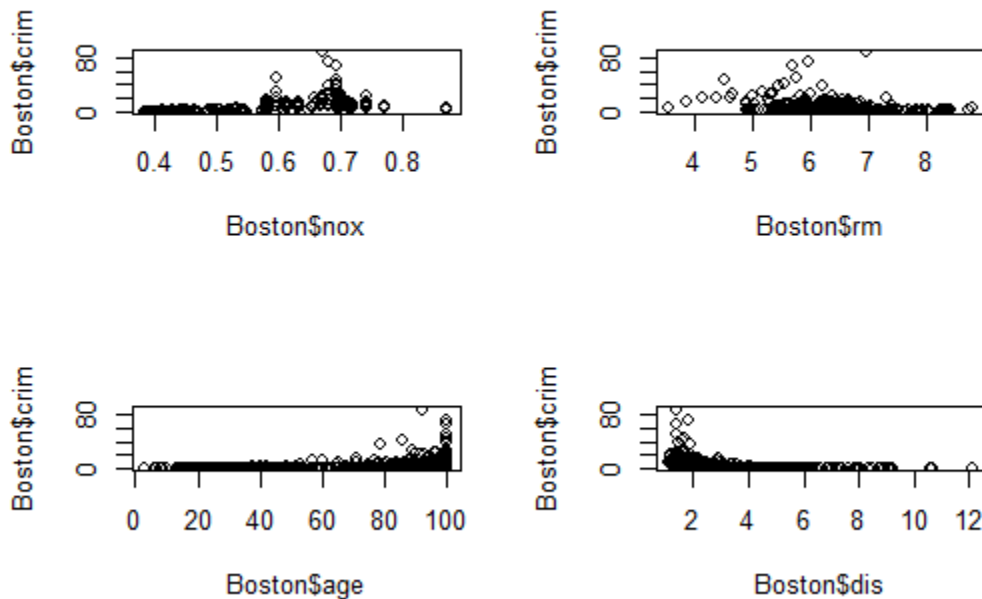


4) Scatter Plot:

**Scatterplot Example**

Problem 2:



A) Indeed multiple predictors appear to be statistically significant. Displacement, weight, year, and origin are included there.

B) The year coefficient suggests a positive correlation between increasing year, and increasing mpg, in other words, mpg gets better with newer cars. This makes sense, as advances are made and engines become more efficient.

C) From the p-values, we can see that the interaction between displacement and weight is statistically signifcant, while the interactiion between cylinders and displacement is not..

Problem 3:

    i)       Linear regression is does not help us here, partly, because the pixels for different samples do not align properly.

    ii)     The linear regression does better on the test data than on the training data.

    iii)    Nearest neighbor results are quite reasonable. The training error results are reduced by the fact that there is one direct hit.

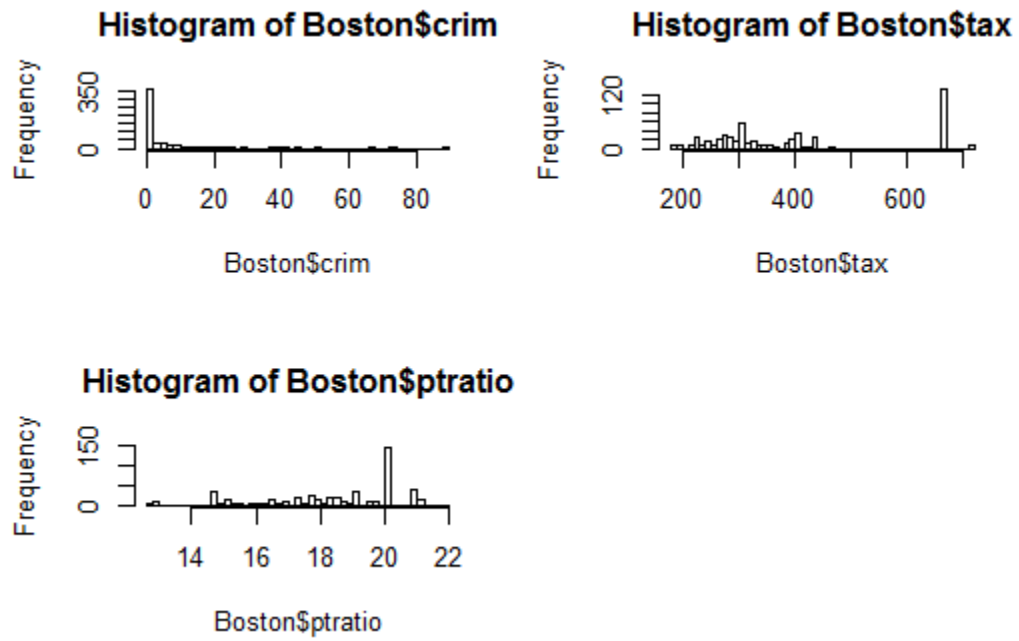    iv)    The amount of error increases as the number of neighbors is increased.

Problem 4:

a) Pairwise Scatterplots



b) There may be a relationship between crim and nox, rm, age, dis, lstat and medv.

c)

**Histogram of Boston$crim**



**Histogram of Boston$tax**



**Histogram of Boston$ptratio**



d) 64 suburbs average more than seven rooms per dwelling.
   13 suburbs average more than eight rooms per dwelling.