

Conceptual

1.

(a)	Obs.	X1	X2	X3	Distance(0, 0, 0)	Y
	1	0	3	0	3	Red
	2	2	0	0	2	Red
	3	0	1	3	$\sqrt{10} \sim 3.2$	Red
	4	0	1	2	$\sqrt{5} \sim 2.2$	Green
	5	-1	0	1	$\sqrt{2} \sim 1.4$	Green
	6	1	1	1	$\sqrt{3} \sim 1.7$	Red

(b) Green. Observation #5 is the closest neighbor for K = 1.

(c) Red. Observations #2, 5, 6 are the closest neighbors for K = 3. 2 is Red,
5 is Green, and 6 is Red.

(d) Small. A small K would be flexible for a non-linear decision boundary,
whereas a large K would try to fit a more linear boundary because it takes more points into consideration.

hw1_solution

Applied

1

(a)

```
library(readr)
college <- read.csv("C:/Users/jzz0121/Desktop/STAT6000/datasets/College.csv")
```

(b)

```
fix(college)
rownames(college) <- college[, 1]
college <- college[, -1]
fix(college)
```

(c)

```
# i.
summary(college)
```

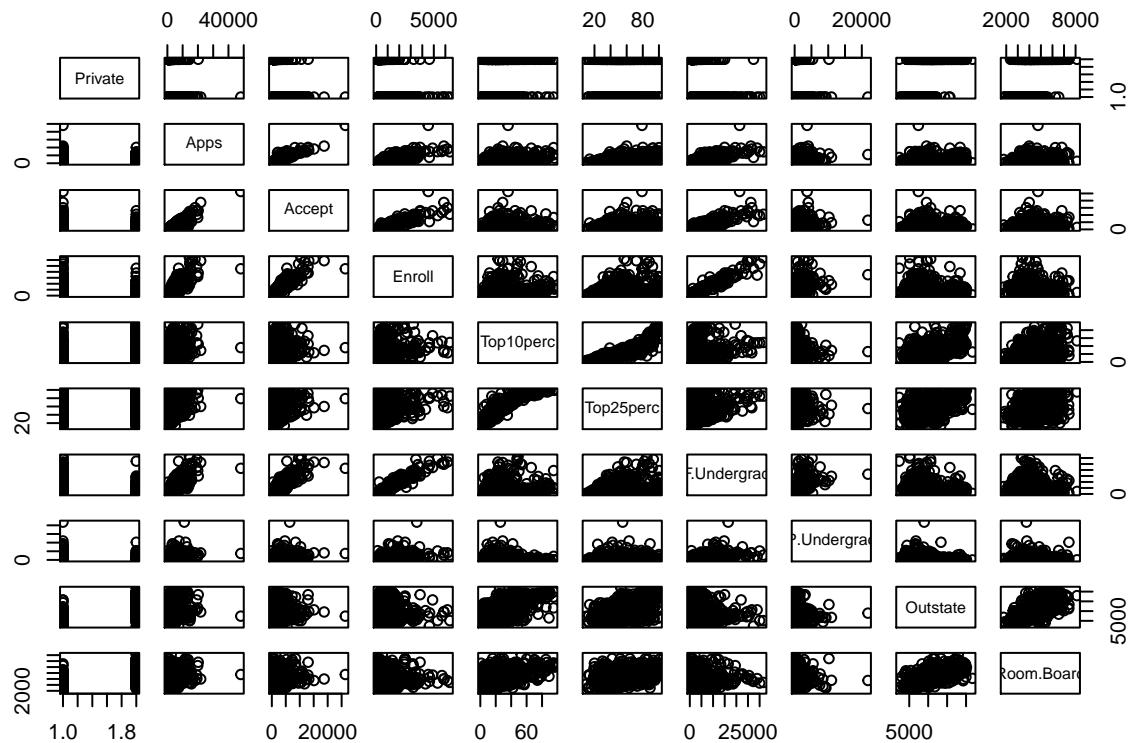
```
##   Private      Apps      Accept      Enroll    Top10perc
##   No :212  Min.   : 81  Min.   : 72  Min.   : 35  Min.   : 1.00
##   Yes:565  1st Qu.: 776 1st Qu.: 604 1st Qu.: 242 1st Qu.:15.00
##               Median :1558  Median :1110  Median :434  Median :23.00
##               Mean   :3002  Mean   :2019  Mean   :780  Mean   :27.56
##               3rd Qu.:3624 3rd Qu.:2424 3rd Qu.:902 3rd Qu.:35.00
##               Max.  :48094 Max.  :26330 Max.  :6392  Max.  :96.00
##   Top25perc    F.Undergrad    P.Undergrad      Outstate
##   Min.   : 9.0  Min.   :139  Min.   : 1.0  Min.   :2340
##   1st Qu.: 41.0 1st Qu.:992  1st Qu.: 95.0  1st Qu.:7320
##   Median : 54.0  Median :1707  Median : 353.0  Median :9990
##   Mean   : 55.8  Mean   :3700  Mean   : 855.3  Mean   :10441
##   3rd Qu.: 69.0 3rd Qu.:4005 3rd Qu.: 967.0  3rd Qu.:12925
##   Max.   :100.0  Max.   :31643  Max.   :21836.0  Max.   :21700
##   Room.Board      Books      Personal      PhD
##   Min.   :1780  Min.   : 96.0  Min.   : 250  Min.   :  8.00
##   1st Qu.:3597  1st Qu.:470.0  1st Qu.: 850  1st Qu.: 62.00
##   Median :4200  Median :500.0  Median :1200  Median : 75.00
##   Mean   :4358  Mean   :549.4  Mean   :1341  Mean   : 72.66
##   3rd Qu.:5050  3rd Qu.:600.0  3rd Qu.:1700  3rd Qu.: 85.00
##   Max.   :8124  Max.   :2340.0  Max.   :6800  Max.   :103.00
##   Terminal      S.F.Ratio      perc.alumni      Expend
##   Min.   : 24.0  Min.   : 2.50  Min.   : 0.00  Min.   : 3186
##   1st Qu.: 71.0  1st Qu.:11.50  1st Qu.:13.00  1st Qu.: 6751
##   Median : 82.0  Median :13.60  Median :21.00  Median : 8377
##   Mean   : 79.7  Mean   :14.09  Mean   :22.74  Mean   : 9660
##   3rd Qu.: 92.0  3rd Qu.:16.50  3rd Qu.:31.00  3rd Qu.:10830
```

```

##   Max.    :100.0   Max.    :39.80   Max.    :64.00   Max.    :56233
##   Grad.Rate
##   Min.    : 10.00
##   1st Qu.: 53.00
##   Median  : 65.00
##   Mean    : 65.46
##   3rd Qu.: 78.00
##   Max.    :118.00

# ii.
pairs(college[,1:10])

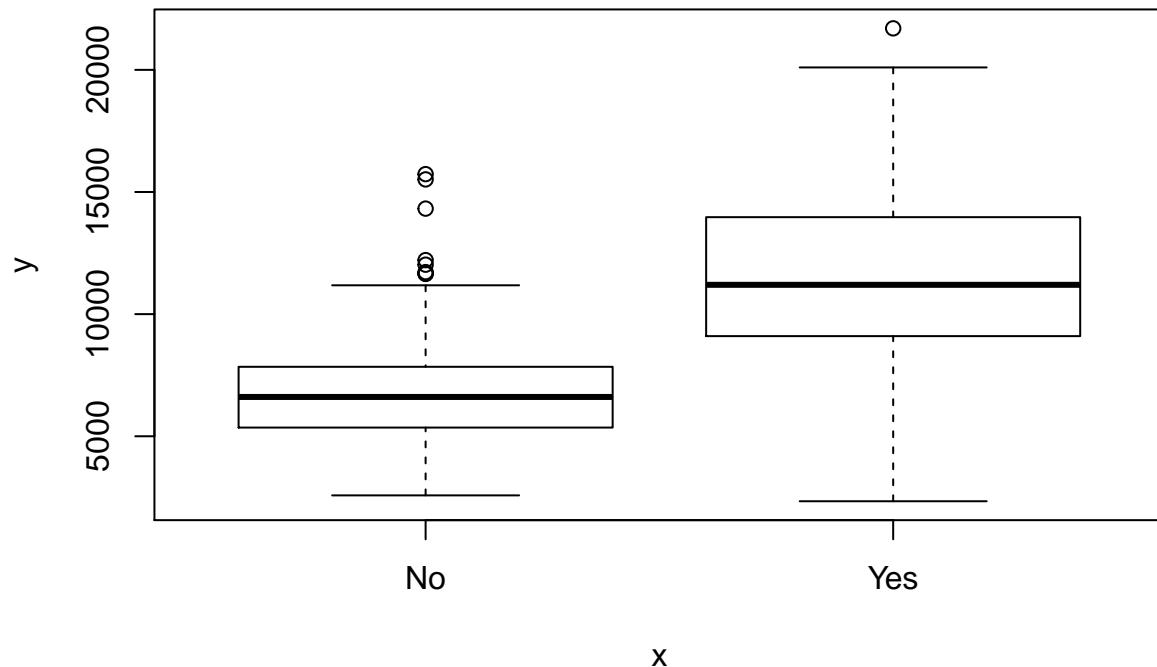
```



```

# iii.
plot(college$Private, college$Outstate)

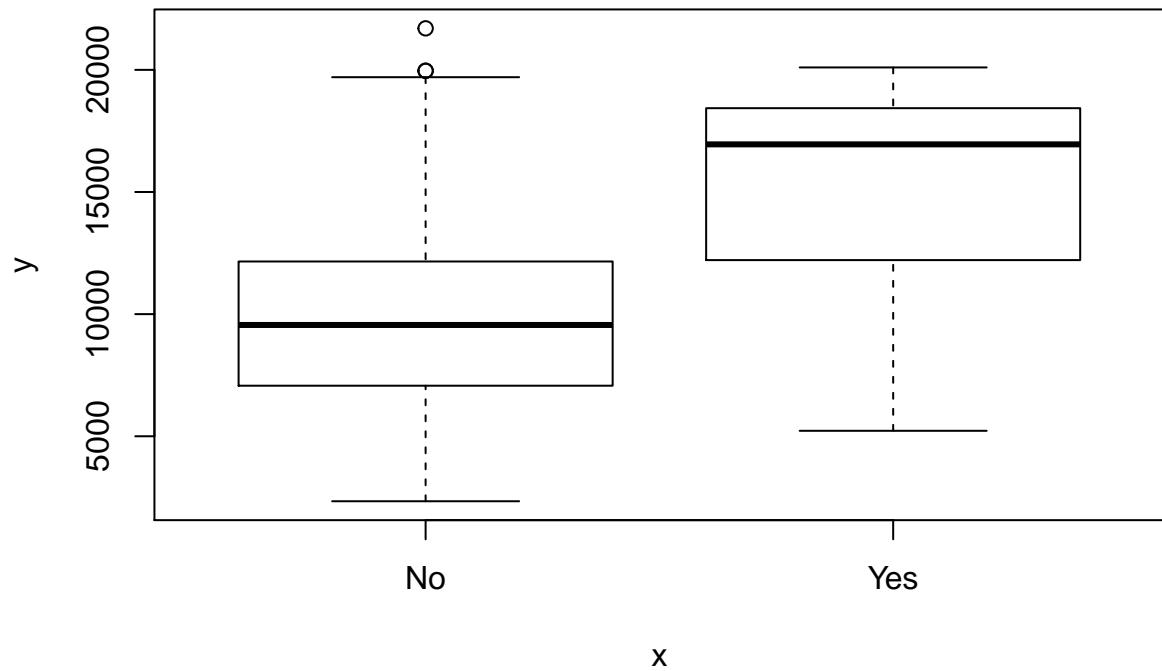
```



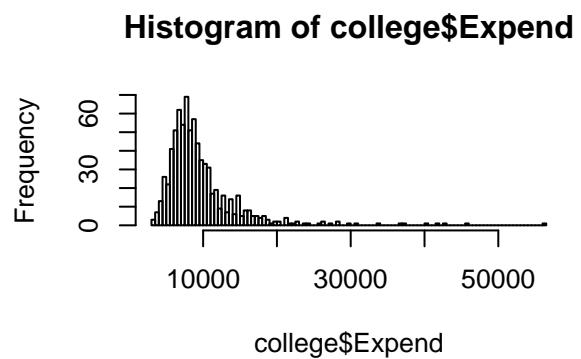
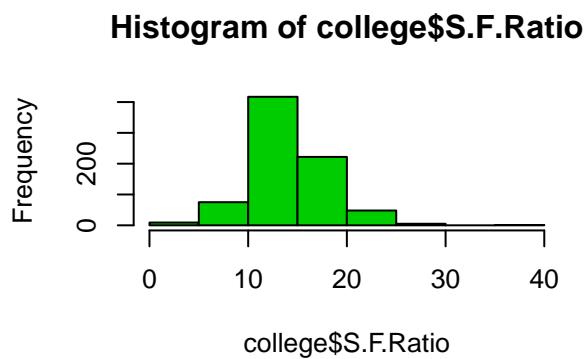
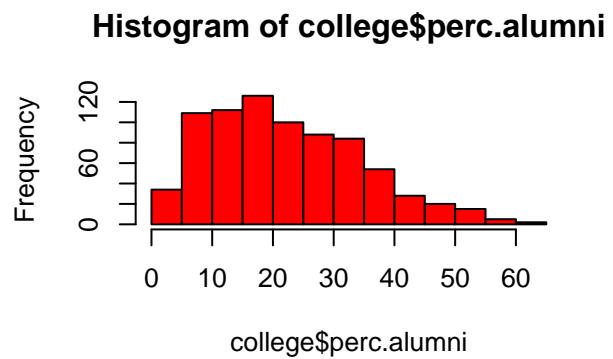
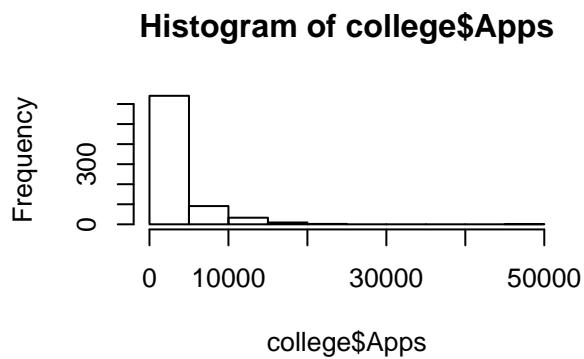
```
# iv.
Elite <- rep("No", nrow(college))
Elite[college$Top10perc > 50] <- "Yes"
Elite <- as.factor(Elite)
college <- data.frame(college, Elite)
summary(college$Elite)

##  No Yes
## 699 78

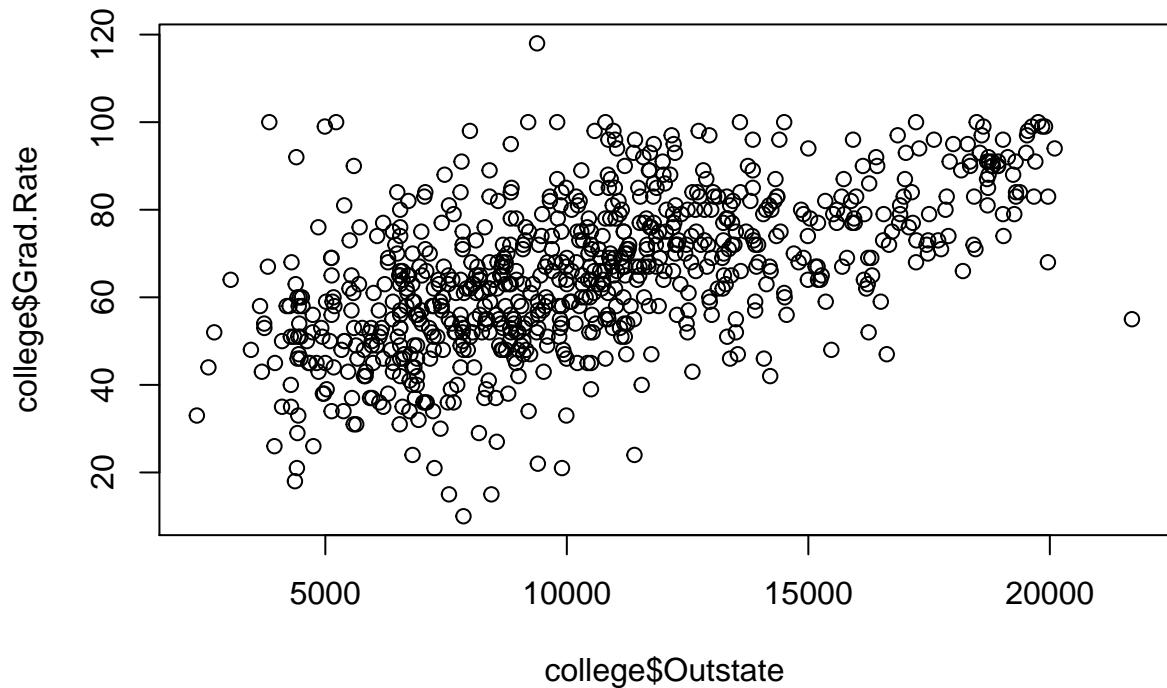
plot(college$Elite, college$Outstate)
```



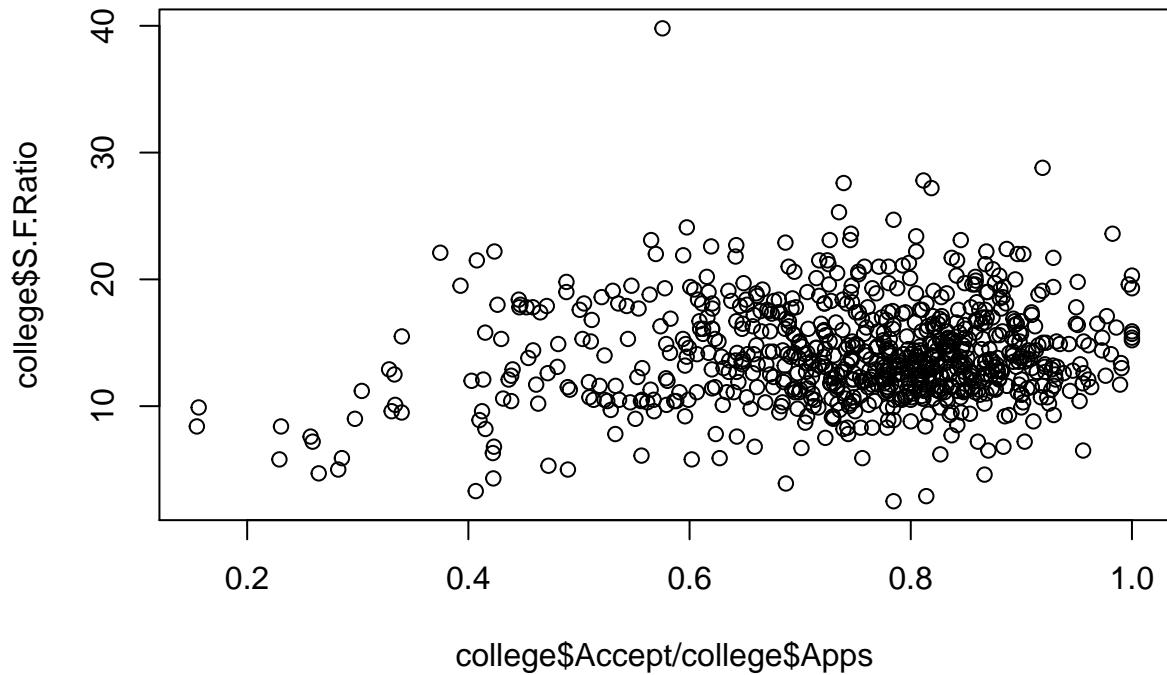
```
# v.  
par(mfrow=c(2,2))  
hist(college$Apps)  
hist(college$perc.alumni, col=2)  
hist(college$S.F.Ratio, col=3, breaks=10)  
hist(college$Expend, breaks=100)
```



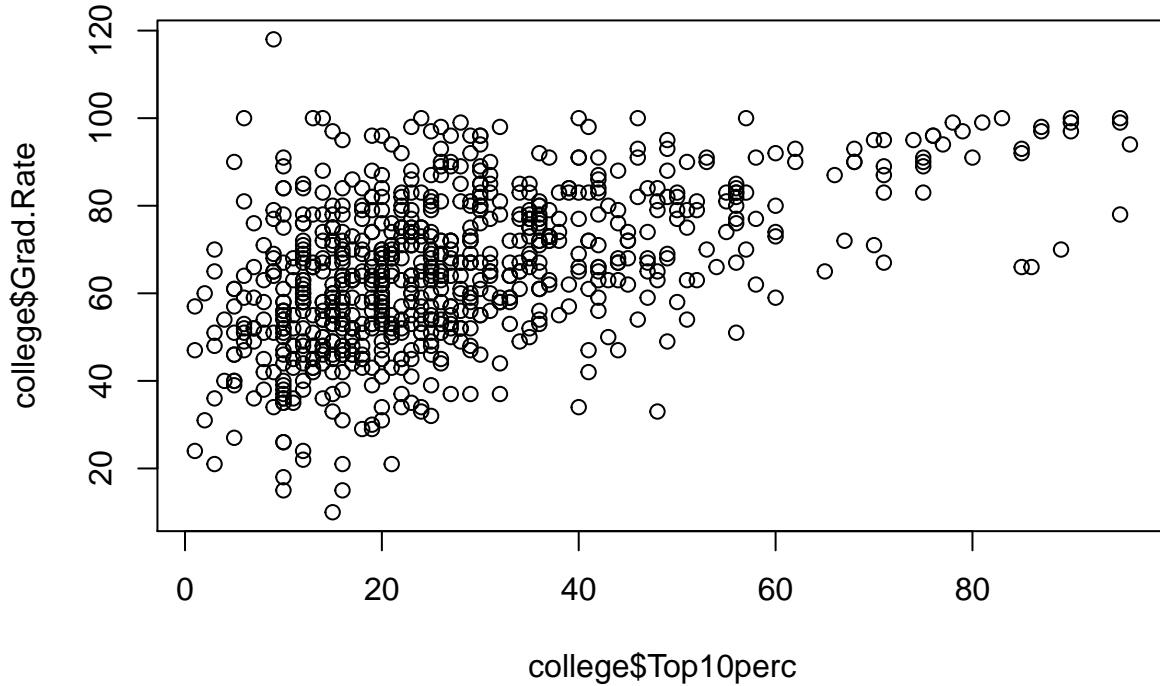
```
# vi.  
par(mfrow=c(1,1))  
# High tuition correlates to high graduation rate.  
plot(college$Outstate, college$Grad.Rate)
```



```
# Colleges with low acceptance rate tend to have low S:F ratio.  
plot(college$Accept / college$Apps, college$S.F.Ratio)
```



```
# Colleges with the most students from top 10% perc don't necessarily have  
# the highest graduation rate. Also, rate > 100 is erroneous!  
plot(college$Top10perc, college$Grad.Rate)
```



2

```

## mark ? as NA
Auto <- read.csv("C:/Users/jzz0121/Desktop/STAT6000/datasets/Auto.csv", header=T, na.strings="?")
summary(Auto) ## notice there are NA in horsepower

##      mpg      cylinders      displacement      horsepower
## Min.   : 9.00  Min.   :3.000  Min.   :68.0  Min.   :46.0
## 1st Qu.:17.50  1st Qu.:4.000  1st Qu.:104.0  1st Qu.:75.0
## Median :23.00  Median :4.000  Median :146.0  Median :93.5
## Mean   :23.52  Mean   :5.458  Mean   :193.5  Mean   :104.5
## 3rd Qu.:29.00  3rd Qu.:8.000  3rd Qu.:262.0  3rd Qu.:126.0
## Max.   :46.60  Max.   :8.000  Max.   :455.0  Max.   :230.0
##                               NA's   :5
##      weight      acceleration      year      origin
## Min.   :1613  Min.   :8.00  Min.   :70.00  Min.   :1.000
## 1st Qu.:2223  1st Qu.:13.80  1st Qu.:73.00  1st Qu.:1.000
## Median :2800  Median :15.50  Median :76.00  Median :1.000
## Mean   :2970  Mean   :15.56  Mean   :75.99  Mean   :1.574
## 3rd Qu.:3609  3rd Qu.:17.10  3rd Qu.:79.00  3rd Qu.:2.000
## Max.   :5140  Max.   :24.80  Max.   :82.00  Max.   :3.000
##
##      name
## ford pinto   : 6

```

```

##   amc matador    : 5
##   ford maverick : 5
##   toyota corolla: 5
##   amc gremlin    : 4
##   amc hornet     : 4
##   (Other)         :368

Auto <- na.omit(Auto)
dim(Auto)

## [1] 392   9

summary(Auto)

##      mpg          cylinders      displacement      horsepower
##  Min.   :9.00   Min.   :3.000   Min.   :68.0   Min.   :46.0
##  1st Qu.:17.00  1st Qu.:4.000  1st Qu.:105.0  1st Qu.:75.0
##  Median :22.75  Median :4.000  Median :151.0  Median :93.5
##  Mean   :23.45  Mean   :5.472  Mean   :194.4  Mean   :104.5
##  3rd Qu.:29.00  3rd Qu.:8.000  3rd Qu.:275.8  3rd Qu.:126.0
##  Max.   :46.60  Max.   :8.000  Max.   :455.0  Max.   :230.0
##
##      weight        acceleration       year        origin
##  Min.   :1613   Min.   :8.00   Min.   :70.00  Min.   :1.000
##  1st Qu.:2225  1st Qu.:13.78  1st Qu.:73.00  1st Qu.:1.000
##  Median :2804  Median :15.50  Median :76.00  Median :1.000
##  Mean   :2978  Mean   :15.54  Mean   :75.98  Mean   :1.577
##  3rd Qu.:3615  3rd Qu.:17.02  3rd Qu.:79.00  3rd Qu.:2.000
##  Max.   :5140  Max.   :24.80  Max.   :82.00  Max.   :3.000
##
##      name
##  amc matador    : 5
##  ford pinto     : 5
##  toyota corolla : 5
##  amc gremlin    : 4
##  amc hornet     : 4
##  chevrolet chevette: 4
##  (Other)         :365

```

(a)

quantitative: mpg, cylinders, displacement, horsepower, weight, acceleration, year.

qualitative: name, origin

(b)

```
# apply the range function to the first seven columns of Auto
sapply(Auto[, 1:7], range)
```

```
##      mpg cylinders displacement horsepower weight acceleration year
##  [1,] 9.0        3           68        46    1613        8.0      70
##  [2,] 46.6       8           455       230    5140       24.8      82
```

(c)

```
sapply(Auto[, 1:7], mean)

##          mpg      cylinders displacement horsepower      weight
##  23.445918      5.471939    194.411990     104.469388  2977.584184
## acceleration      year
##      15.541327    75.979592
```

```
sapply(Auto[, 1:7], sd)

##          mpg      cylinders displacement horsepower      weight
##    7.805007      1.705783    104.644004     38.491160   849.402560
## acceleration      year
##      2.758864      3.683737
```

(d)

```
newAuto <- Auto[-(10:85),]

sapply(newAuto[, 1:7], range)

##          mpg cylinders displacement horsepower weight acceleration year
## [1,] 11.0         3           68          46   1649        8.5     70
## [2,] 46.6         8           455         230   4997       24.8     82
```

```
sapply(newAuto[, 1:7], mean)

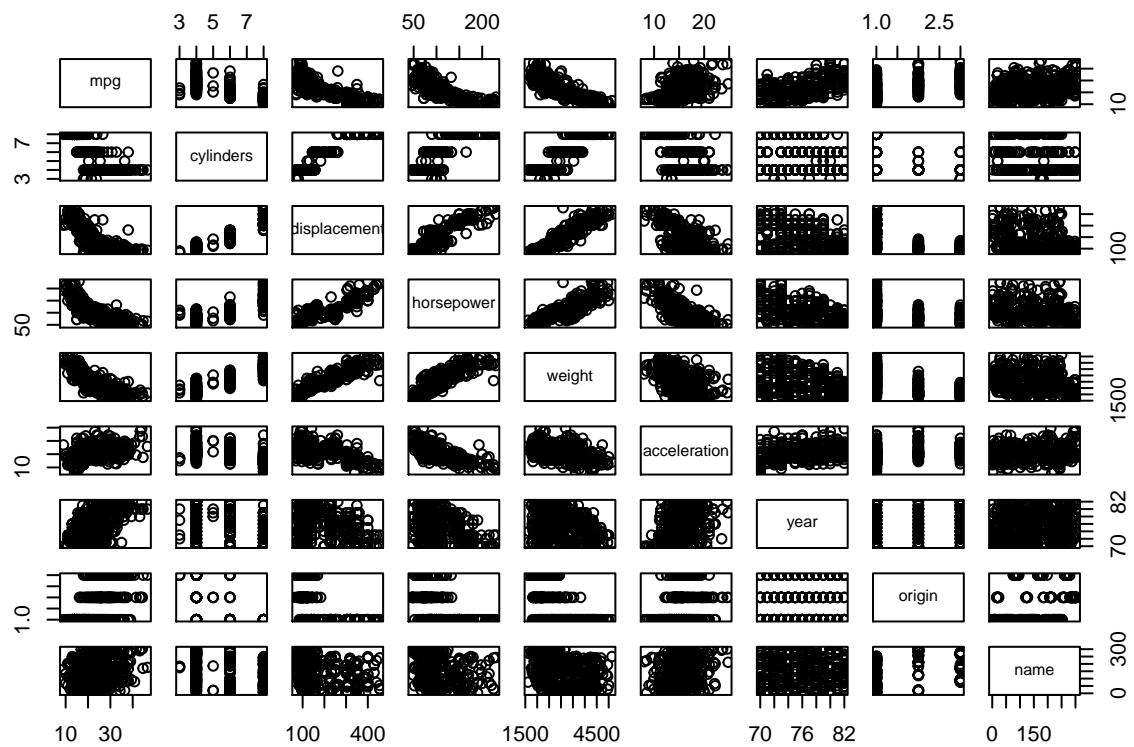
##          mpg      cylinders displacement horsepower      weight
##  24.404430      5.373418    187.240506     100.721519  2935.971519
## acceleration      year
##      15.726899    77.145570
```

```
sapply(newAuto[, 1:7], sd)

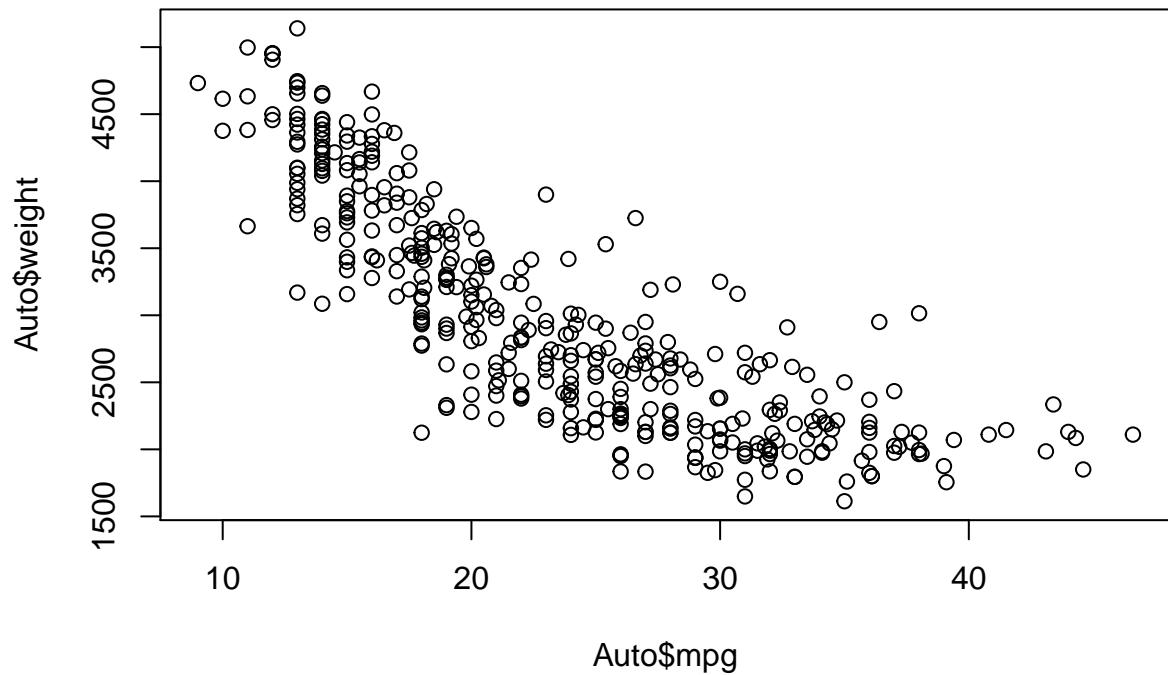
##          mpg      cylinders displacement horsepower      weight
##    7.867283      1.654179    99.678367     35.708853   811.300208
## acceleration      year
##      2.693721      3.106217
```

(e)

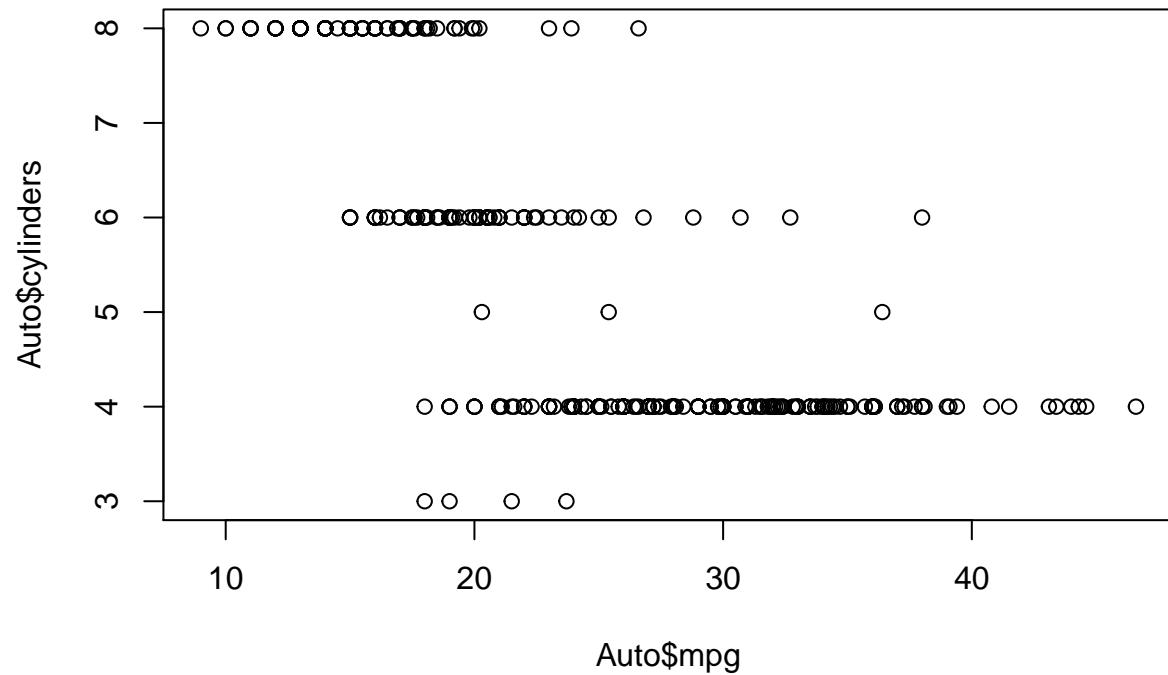
```
pairs(Auto)
```



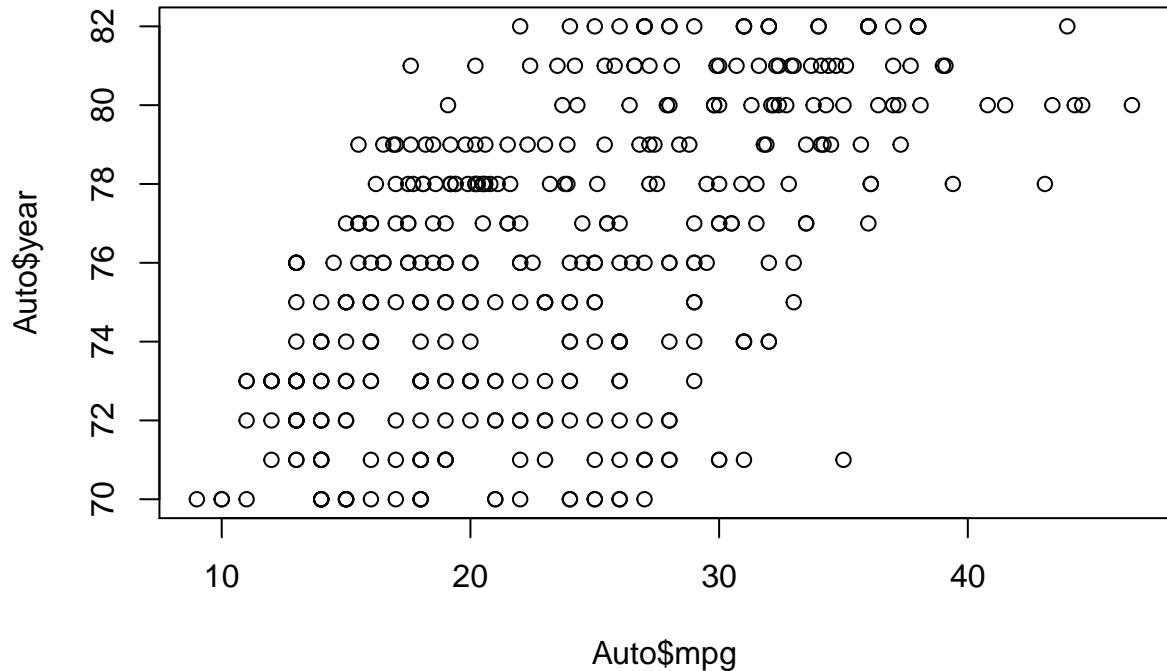
```
# Heavier weight correlates with lower mpg.
plot(Auto$mpg, Auto$weight)
```



```
# More cylinders, less mpg.  
plot(Auto$mpg, Auto$cylinders)
```



```
# Cars become more efficient over time.  
plot(Auto$mpg, Auto$year)
```



- (f) See pairs plot in (e). All of the predictors show some correlation with mpg. The name predictor has too little observations per name though, so using this as a predictor is likely to result in overfitting the data and will not generalize well.

3

(a) 506 rows, 14 columns.

14 features, 506 housing values in Boston suburbs

```
library(MASS)
# ?Boston
dim(Boston)
```

```
## [1] 506 14
```

(b)

X correlates with: a, b, c

crim: age, dis, rad, tax, ptratio

zn: indus, nox, age, lstat

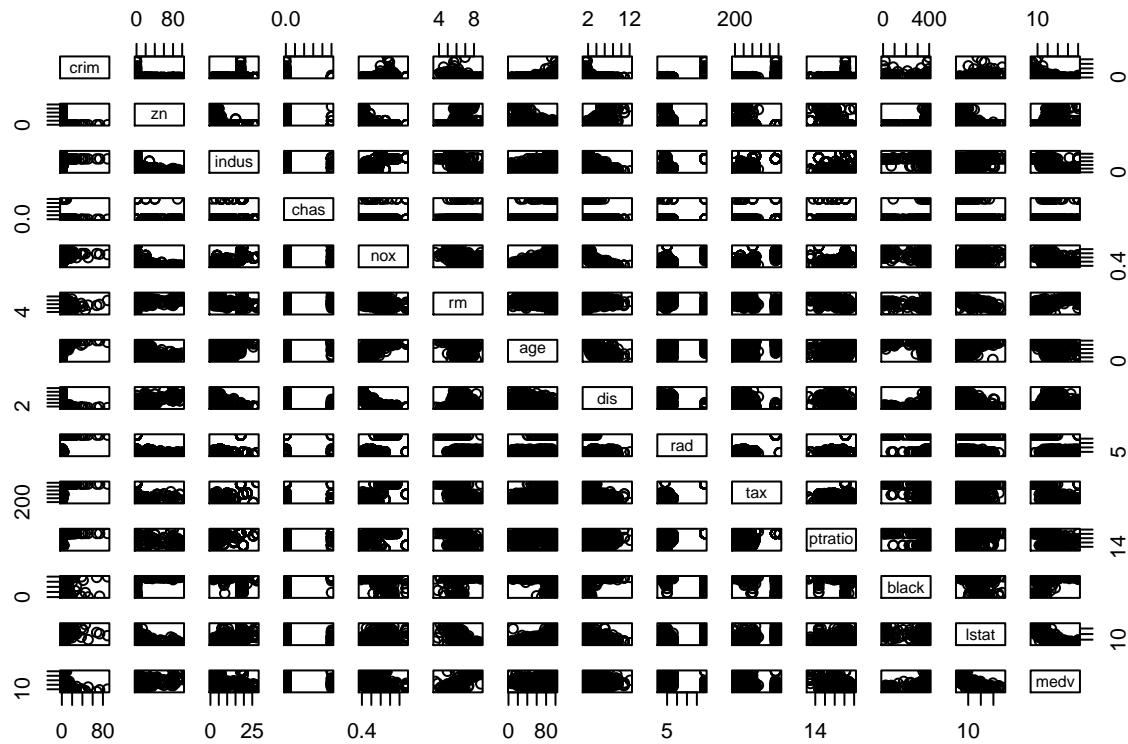
indus: age, dis

nox: age, dis

dis: lstat

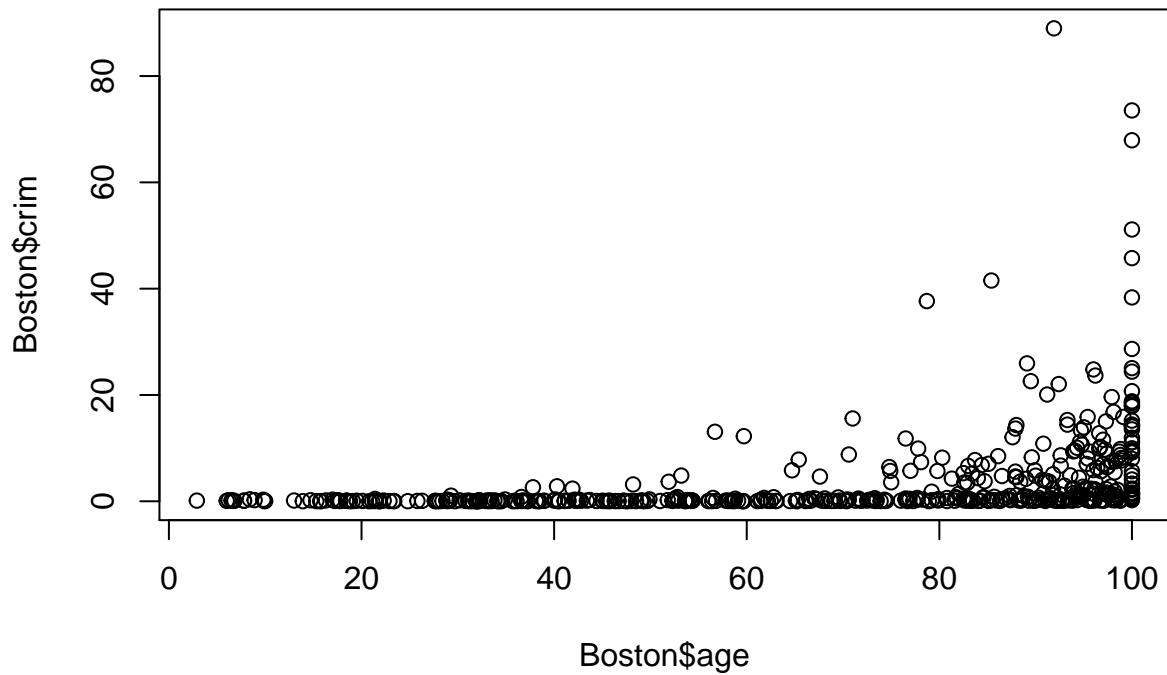
lstat: medv

```
pairs(Boston)
```

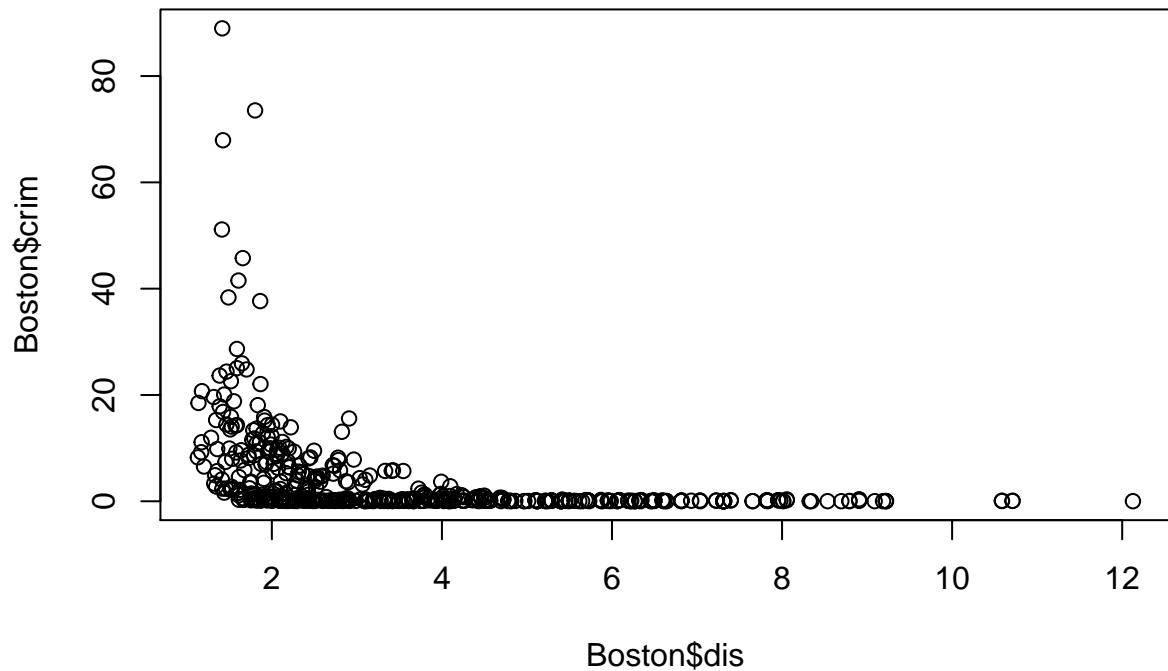


(c)

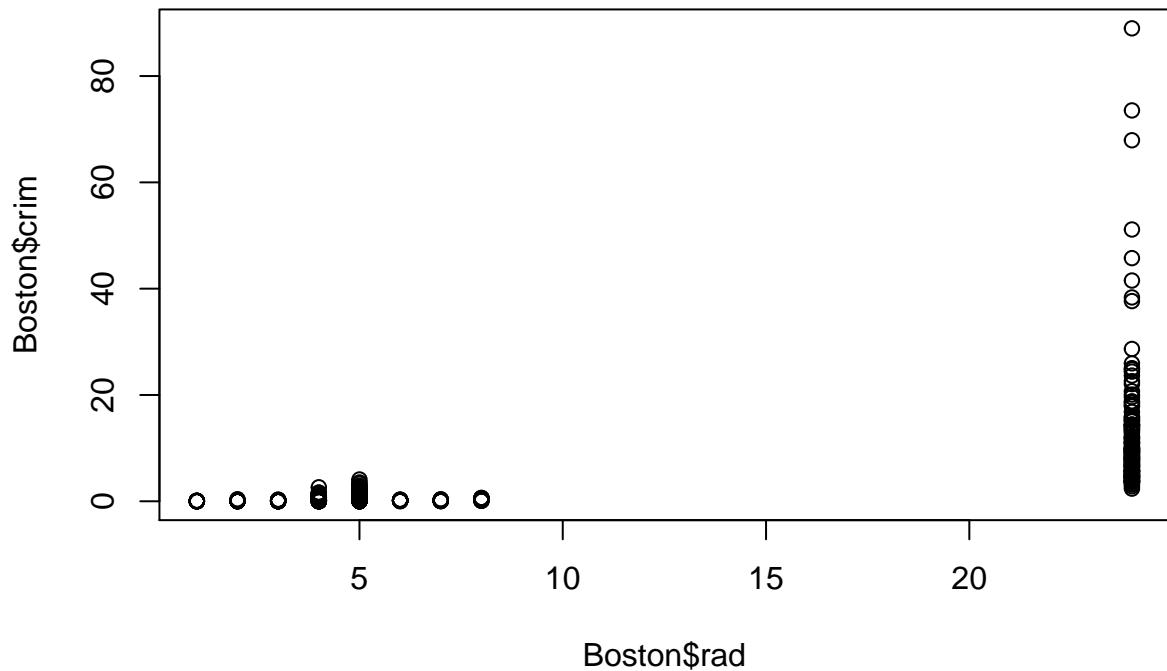
```
# Older homes, more crime  
plot(Boston$age, Boston$crim)
```



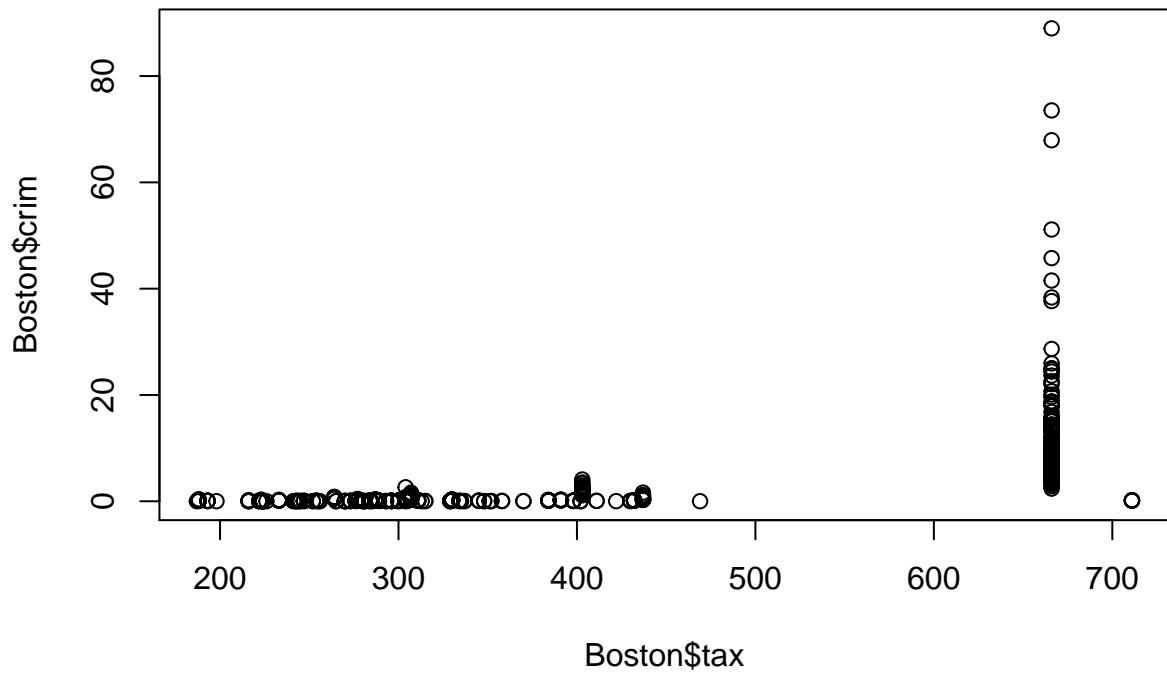
```
# Closer to work-area, more crime  
plot(Boston$dis, Boston$crim)
```



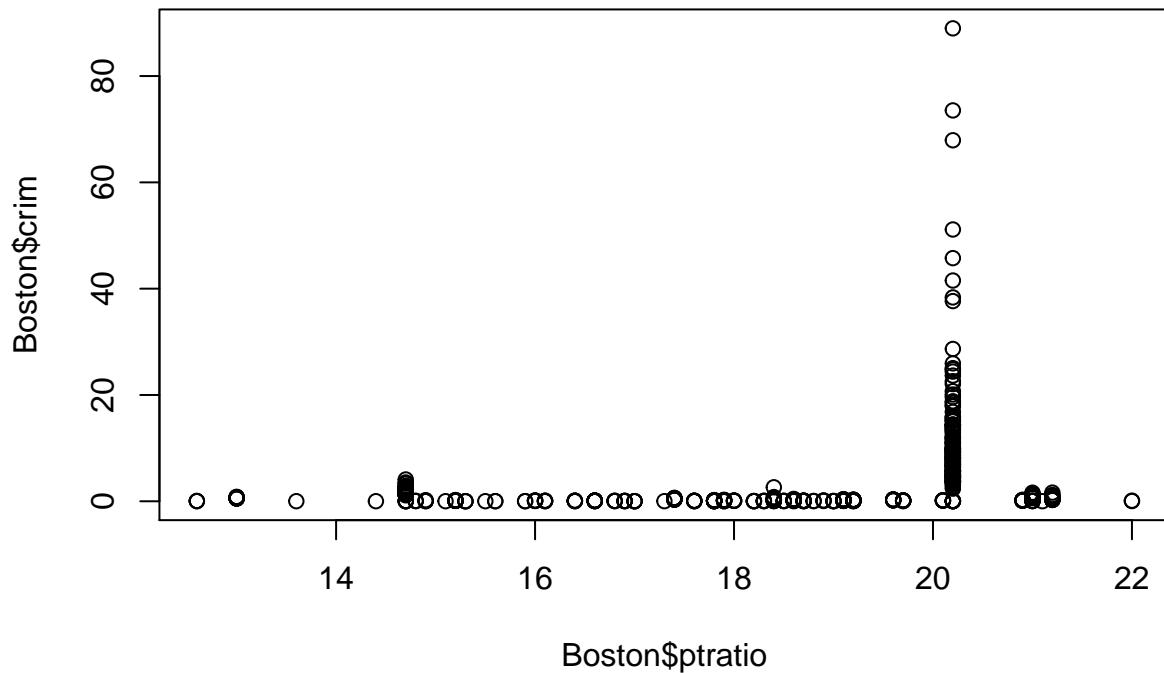
```
# Higher index of accessibility to radial highways, more crime  
plot(Boston$rad, Boston$crim)
```



```
# Higher tax rate, more crime  
plot(Boston$tax, Boston$crim)
```

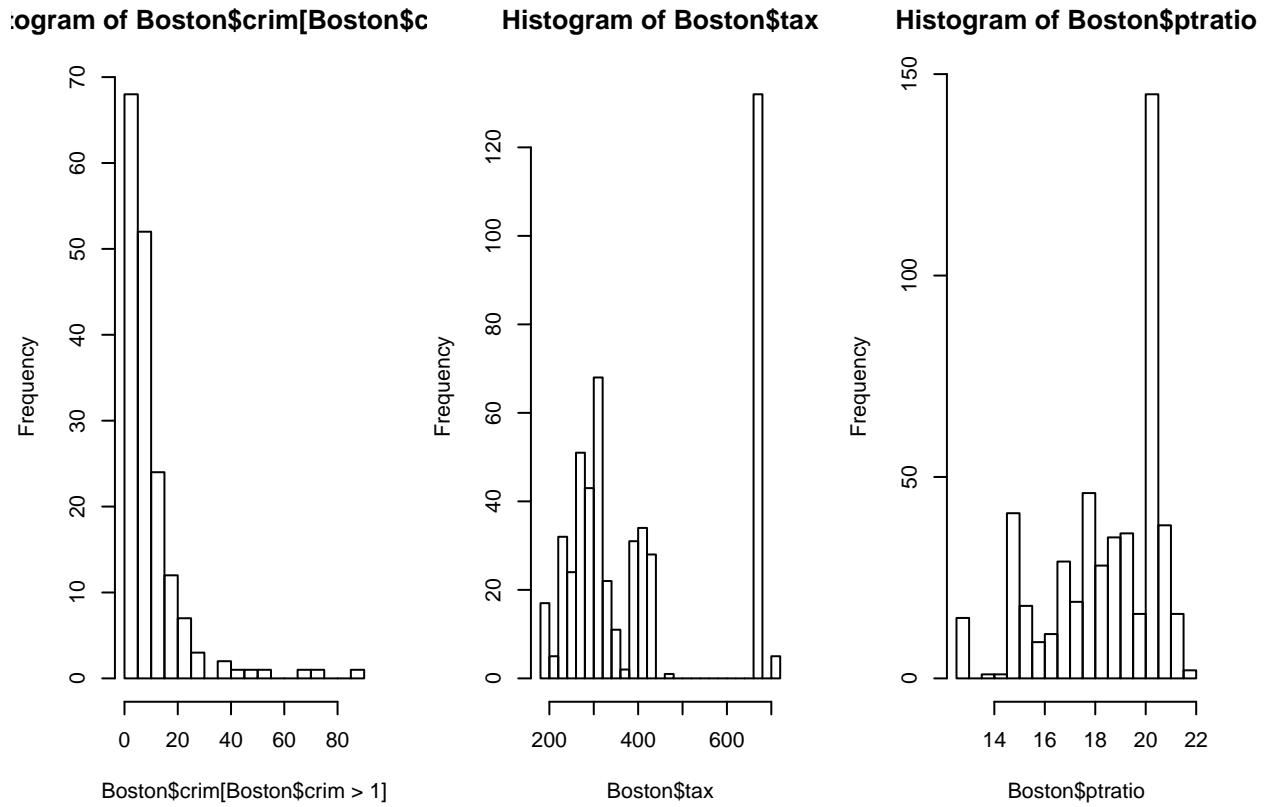


```
# Higher pupil:teacher ratio, more crime  
plot(Boston$ptratio, Boston$crim)
```



(d)

```
par(mfrow=c(1,3))
hist(Boston$crim[Boston$crim>1], breaks=25)
# most cities have low crime rates, but there is a long tail: 18 suburbs appear
# to have a crime rate > 20, reaching to above 80
hist(Boston$tax, breaks=25)
# there is a large divide between suburbs with low tax rates and a peak at 660-680
hist(Boston$pstatus, breaks=25)
```



```
# a skew towards high ratios, but no particularly high ratios
```

(e)

```
# 35 suburbs
dim(subset(Boston, chas == 1))
```

```
## [1] 35 14
```

```
table(Boston$chas)
```

```
##
##      0      1
## 471   35
```

(f)

```
median(Boston$ptratio)
```

```
## [1] 19.05
```

(g)

```

subset(Boston, medv == min(Boston$medv))

##      crim zn indus chas   nox     rm age     dis rad tax ptratio black
## 399 38.3518 0 18.1    0 0.693 5.453 100 1.4896 24 666    20.2 396.90
## 406 67.9208 0 18.1    0 0.693 5.683 100 1.4254 24 666    20.2 384.97
## lstat medv
## 399 30.59    5
## 406 22.98    5

```

```
summary(Boston)
```

```

##      crim             zn            indus            chas
##  Min.   : 0.00632   Min.   : 0.00   Min.   : 0.46   Min.   :0.00000
##  1st Qu.: 0.08204   1st Qu.: 0.00   1st Qu.: 5.19   1st Qu.:0.00000
##  Median : 0.25651   Median : 0.00   Median : 9.69   Median :0.00000
##  Mean   : 3.61352   Mean   : 11.36  Mean   :11.14   Mean   :0.06917
##  3rd Qu.: 3.67708   3rd Qu.: 12.50  3rd Qu.:18.10   3rd Qu.:0.00000
##  Max.   :88.97620   Max.   :100.00  Max.   :27.74   Max.   :1.00000
##      nox              rm            age             dis
##  Min.   :0.3850   Min.   :3.561   Min.   : 2.90   Min.   : 1.130
##  1st Qu.:0.4490   1st Qu.:5.886   1st Qu.: 45.02  1st Qu.: 2.100
##  Median :0.5380   Median :6.208   Median : 77.50  Median : 3.207
##  Mean   :0.5547   Mean   :6.285   Mean   : 68.57  Mean   : 3.795
##  3rd Qu.:0.6240   3rd Qu.:6.623   3rd Qu.: 94.08  3rd Qu.: 5.188
##  Max.   :0.8710   Max.   :8.780   Max.   :100.00  Max.   :12.127
##      rad              tax            ptratio          black
##  Min.   : 1.000   Min.   :187.0   Min.   :12.60   Min.   : 0.32
##  1st Qu.: 4.000   1st Qu.:279.0   1st Qu.:17.40   1st Qu.:375.38
##  Median : 5.000   Median :330.0   Median :19.05   Median :391.44
##  Mean   : 9.549   Mean   :408.2   Mean   :18.46   Mean   :356.67
##  3rd Qu.:24.000   3rd Qu.:666.0   3rd Qu.:20.20   3rd Qu.:396.23
##  Max.   :24.000   Max.   :711.0   Max.   :22.00   Max.   :396.90
##      lstat            medv
##  Min.   : 1.73   Min.   : 5.00
##  1st Qu.: 6.95   1st Qu.:17.02
##  Median :11.36   Median :21.20
##  Mean   :12.65   Mean   :22.53
##  3rd Qu.:16.95   3rd Qu.:25.00
##  Max.   :37.97   Max.   :50.00

```

```
# Not the best place to live, but certainly not the worst.
```

(h)

```
dim(subset(Boston, rm > 7))
```

```
## [1] 64 14
```

```
# 64
dim(subset(Boston, rm > 8))
```

```

## [1] 13 14

# 13
summary(subset(Boston, rm > 8))

##      crim          zn          indus          chas
##  Min.   :0.02009  Min.   : 0.00  Min.   : 2.680  Min.   :0.0000
##  1st Qu.:0.33147  1st Qu.: 0.00  1st Qu.: 3.970  1st Qu.:0.0000
##  Median :0.52014  Median : 0.00  Median : 6.200  Median :0.0000
##  Mean   :0.71879  Mean   :13.62  Mean   : 7.078  Mean   :0.1538
##  3rd Qu.:0.57834  3rd Qu.:20.00  3rd Qu.: 6.200  3rd Qu.:0.0000
##  Max.   :3.47428  Max.   :95.00  Max.   :19.580  Max.   :1.0000
##      nox          rm          age          dis
##  Min.   :0.4161  Min.   :8.034  Min.   : 8.40  Min.   :1.801
##  1st Qu.:0.5040  1st Qu.:8.247  1st Qu.:70.40  1st Qu.:2.288
##  Median :0.5070  Median :8.297  Median :78.30  Median :2.894
##  Mean   :0.5392  Mean   :8.349  Mean   :71.54  Mean   :3.430
##  3rd Qu.:0.6050  3rd Qu.:8.398  3rd Qu.:86.50  3rd Qu.:3.652
##  Max.   :0.7180  Max.   :8.780  Max.   :93.90  Max.   :8.907
##      rad          tax          ptratio        black
##  Min.   : 2.000  Min.   :224.0  Min.   :13.00  Min.   :354.6
##  1st Qu.: 5.000  1st Qu.:264.0  1st Qu.:14.70  1st Qu.:384.5
##  Median : 7.000  Median :307.0  Median :17.40  Median :386.9
##  Mean   : 7.462  Mean   :325.1  Mean   :16.36  Mean   :385.2
##  3rd Qu.: 8.000  3rd Qu.:307.0  3rd Qu.:17.40  3rd Qu.:389.7
##  Max.   :24.000  Max.   :666.0  Max.   :20.20  Max.   :396.9
##      lstat         medv
##  Min.   :2.47  Min.   :21.9
##  1st Qu.:3.32  1st Qu.:41.7
##  Median :4.14  Median :48.3
##  Mean   :4.31  Mean   :44.2
##  3rd Qu.:5.12  3rd Qu.:50.0
##  Max.   :7.44  Max.   :50.0

# compare with summary(Boston)
# relatively lower crime (comparing range), lower lstat (comparing range)

```