

Lab2: Data Handling

2025-02-01

This exercise involves the Boston housing data set. Write R codes to answer the following questions.

1. To begin, load in the Boston data set. The Boston data set is part of the MASS library in R. Per capita crime rate is the response variable. Print the first few observations of the dataset

How many rows are in this data set? How many columns? What do the rows and columns represent (check R help)?

```
library(MASS)
head(Boston)
```

```
##      crim zn indus chas   nox    rm  age    dis rad tax ptratio  black lstat
## 1 0.00632 18  2.31    0 0.538 6.575 65.2 4.0900   1 296    15.3 396.90  4.98
## 2 0.02731  0  7.07    0 0.469 6.421 78.9 4.9671   2 242    17.8 396.90  9.14
## 3 0.02729  0  7.07    0 0.469 7.185 61.1 4.9671   2 242    17.8 392.83  4.03
## 4 0.03237  0  2.18    0 0.458 6.998 45.8 6.0622   3 222    18.7 394.63  2.94
## 5 0.06905  0  2.18    0 0.458 7.147 54.2 6.0622   3 222    18.7 396.90  5.33
## 6 0.02985  0  2.18    0 0.458 6.430 58.7 6.0622   3 222    18.7 394.12  5.21
##   medv
## 1 24.0
## 2 21.6
## 3 34.7
## 4 33.4
## 5 36.2
## 6 28.7
```

```
dim(Boston)
```

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

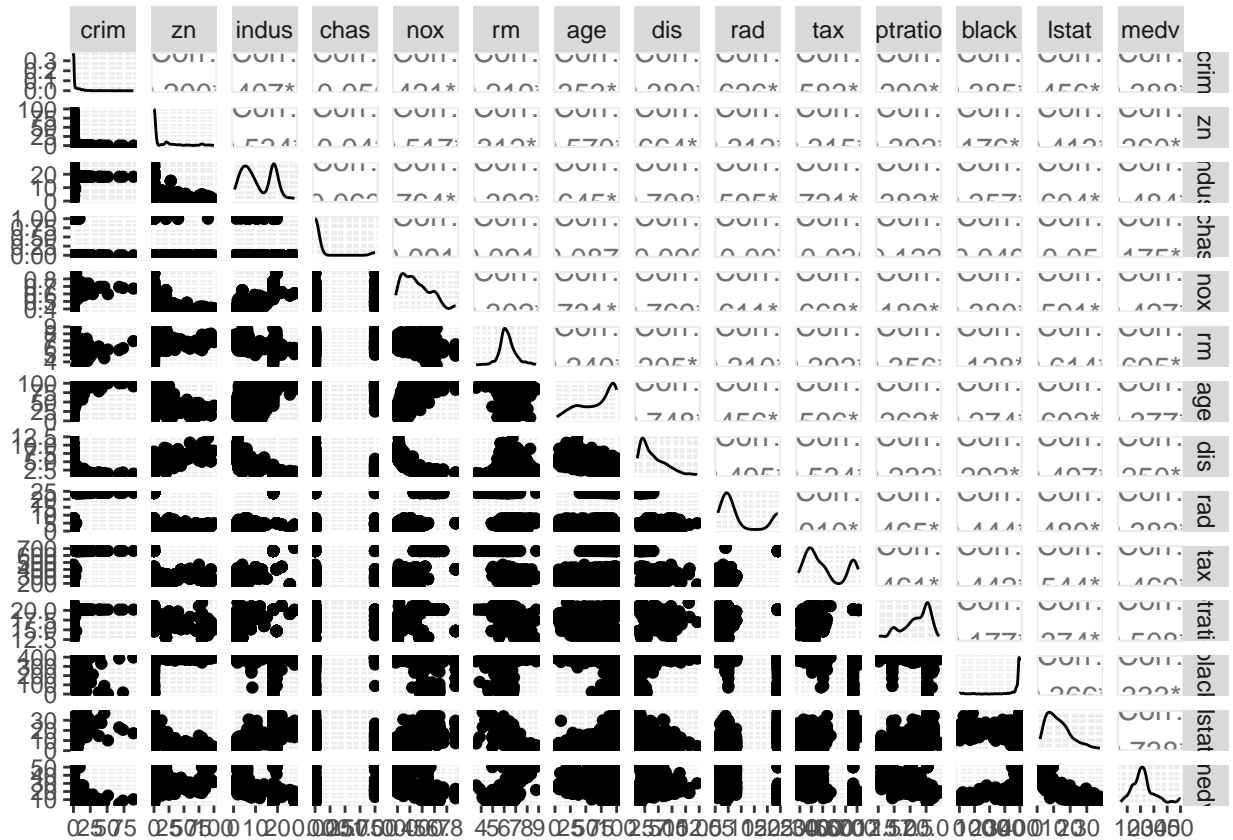
- There are 506 rows and 14 columns.
- Each row corresponds to a specific neighborhood or district within the Boston area.
- Each column represents a different attribute related to housing, crime rates, and economic factors.

2. Make some pairwise scatterplots of the predictors (columns) in this data set. Describe your findings.

```
library(ggplot2)
library(GGally)
```

```
## Registered S3 method overwritten by 'GGally':
##   method from
##   +.gg      ggplot2
```

```
ggpairs(Boston)
```



As percentage of lower-income households (lstat) increases, median home values (medv) tend to decrease. More rooms in a house are associated with a lower percentage of low-income residents and higher median home values.

Higher tax rates are associated with higher pupil-teacher ratios. Higher industrialization correlates with higher nitrogen oxides concentration.

Some variables, such as medv and lstat, or nox vs distances to employment centres, show a curved pattern rather than a straight-line trend, suggesting a non-linear relationship.

3. Are any of the predictors associated with the per capita crime rate? If so, explain the relationship.

```
correlations = cor(Boston)
sorted_correlations = sort(correlations["crim", ], decreasing = TRUE)
```

4. Find the summary of each predictor. Check the top five observations with the highest per capita crime rate.

```
summary(Boston)
```

```
##          crim          zn          indus          chas
## Min.   : 0.00632   Min.   : 0.00   Min.   : 0.46   Min.   :0.00000
## 1st Qu.: 0.08205   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
```

```
sorted_data = sort(Boston[,1], decreasing = TRUE, index.return = TRUE)
data_5 = sorted_data$ix[1:5]
data_5_Boston = Boston[data_5,]
```

5. How many of the census tracts (observations) in this data set bound the Charles river?

```
sum(Boston$chas == 1)
```

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

6. What is the median pupil-teacher ratio among the towns in this data set?

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

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

7. Which census tract of Boston has the lowest median value of owner-occupied homes? Print the corresponding row.

```
Boston[which.min(Boston$medv),]
```

```
##      crim zn indus chas   nox   rm age   dis rad tax ptratio black lstat
## 399 38.3518  0  18.1    0 0.693 5.453 100 1.4896  24 666    20.2 396.9 30.59
##      medv
## 399      5
```

8. In this dataset, how many of the census tracts average more than eight rooms per dwelling? Find the summary of the census tracts that average more than eight rooms per dwelling.

```
sum(Boston$rm > 8)
```

```
## [1] 13
```

```
summary(Boston[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
```

9. Save this dataset in an Excel or CSV file.

```
write.csv(Boston, file = "Boston.csv", row.names = FALSE)

library(openxlsx)
write.xlsx(x = Boston, file = 'Boston_data.xlsx')
```