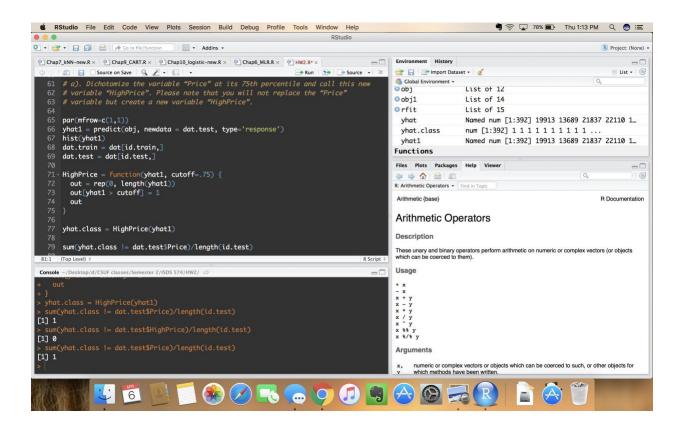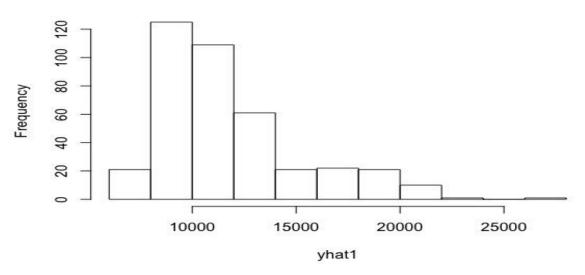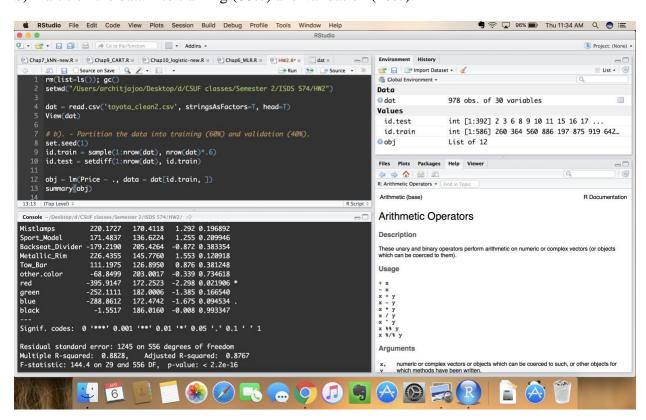## HOMEWORK 2

**a)** Dichotomize the variable "Price" at its 75<sup>th</sup> percentile and call this new variable "HighPrice". Please note that you will not replace the "Price" variable but create a new variable "HighPrice".
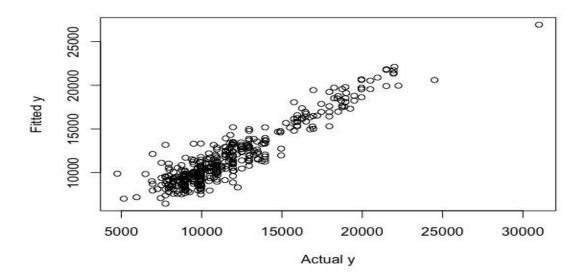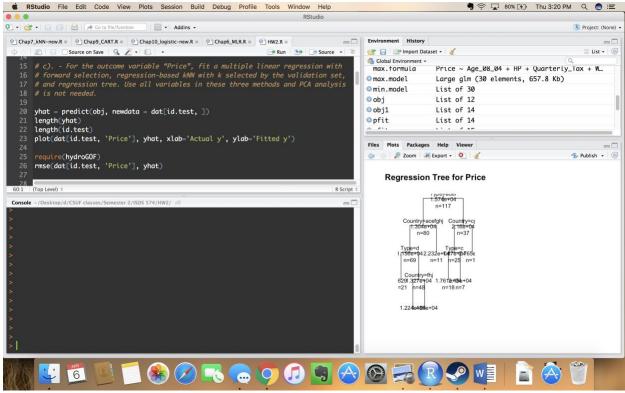
## Histogram of yhat1



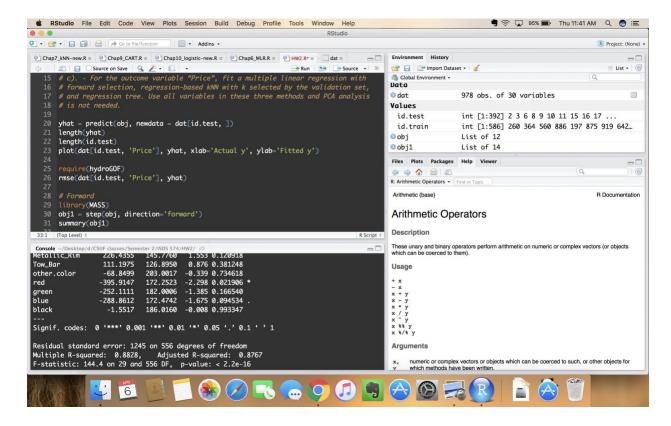**b)** Partition the data into training (60%) and validation (40%).



**c)** For the outcome variable "Price", fit a multiple linear regression with   forward selection, regression-based kNN with k selected by the validation set, and regression tree. Use all variables in these three methods and PCA analysis is not needed.
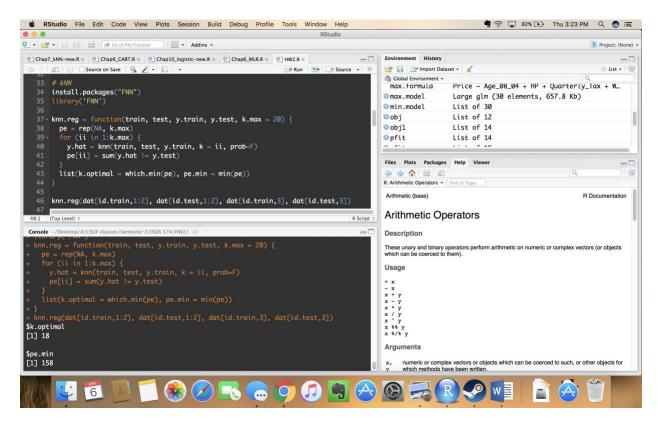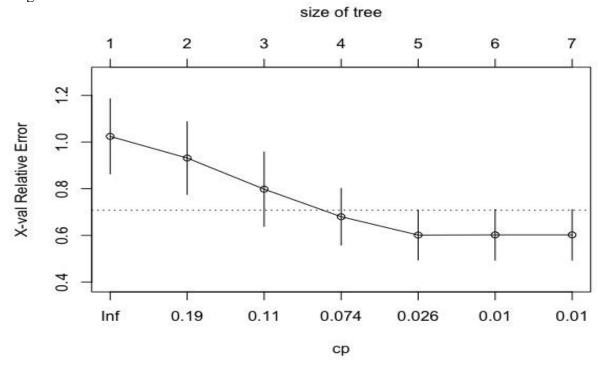
Forward selection:

kNN:

Regression tree:

**Regression Tree for Price**



Type=adei
1.574e+04
n=117

Country=acefghj
1.304e+04
n=80

Country=cgj
2.16e+04
n=37

Type=d
1.156e+04
n=69

2.232e+04
n=11

Type=c
1.87e+04
n=25

2.765e+0
n=12

7629
n=21

Country=fhj
1.327e+04
n=48

1.761e+04
n=18

2.15e+04
n=7

1.224e+04    1.485e+04
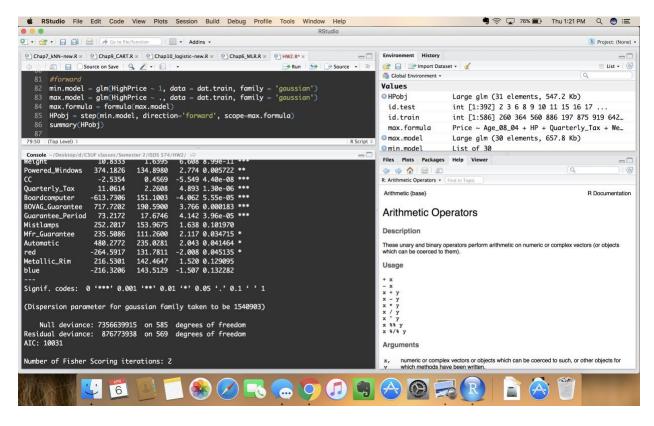
**d)** Compare the above three models obtained and pick up the best one. Explain why you think it is the best one. Note that I already took out the "grey" as the reference group for color. You don't have to do any data cleaning in this homework.
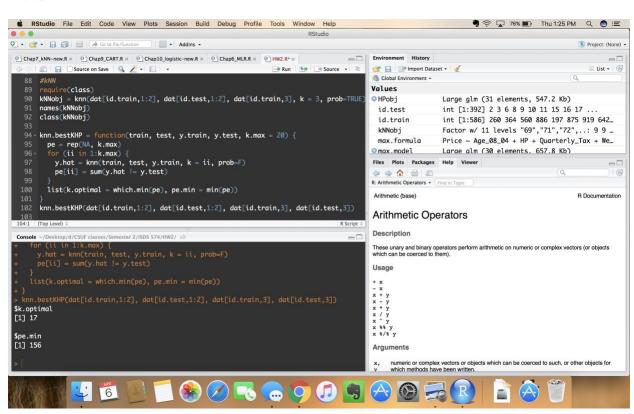
Multiple Linear Regression through forward selection is the best model because it gives lowest RMSE.

**e)** For the outcome variable "HighPrice", fit a logistic regression with forward regression, classification-based kNN with k selected by the validation set, and classification tree.
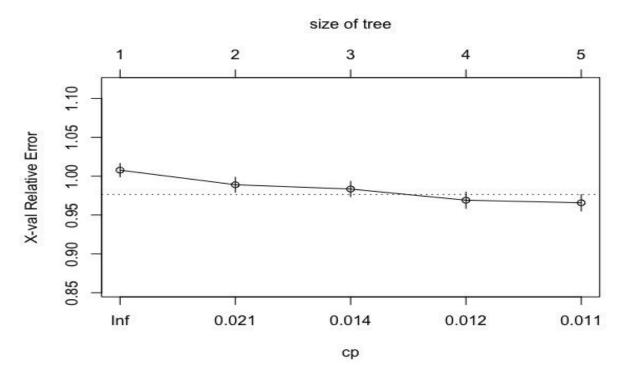
Forward:

kNN:

Classification tree:

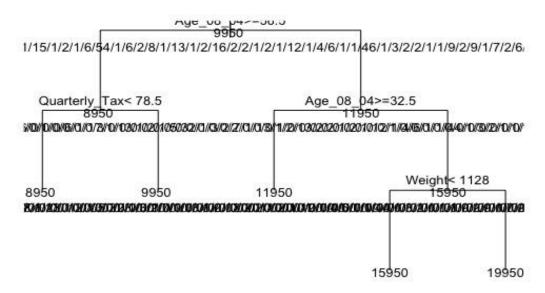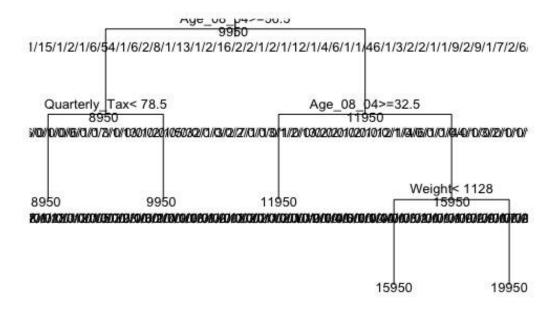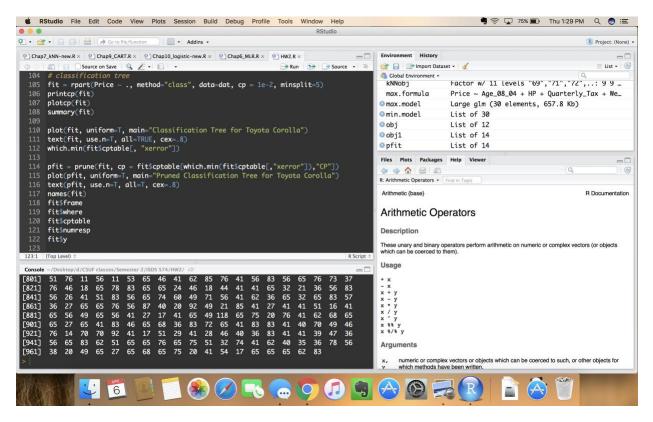## Classification Tree for Toyota Corolla

Age_08_04>=36.5
9950
1/15/1/2/1/6/54/1/6/2/8/1/13/1/2/16/2/2/1/2/1/12/1/4/6/1/1/46/1/3/2/2/1/1/9/2/9/1/7/2/6

Quarterly_Tax< 78.5
8950

Age_08_04>=32.5
11950

8950

9950

11950

Weight< 1128
15950

15950

19950

# Pruned Classification Tree for Toyota Corolla

Age_08_04>=36.5
9950
1/15/1/2/1/6/54/1/6/2/8/1/13/1/2/16/2/2/1/2/1/12/1/4/6/1/1/46/1/3/2/2/1/1/9/2/9/1/7/2/6

Quarterly_Tax< 78.5
8950

Age_08_04>=32.5
11950

8950

9950

11950

Weight< 1128
15950

15950

19950

. **f)** Compare the models obtained and pick up the best one. Explain why you think it is the best one.

Classification tree is the best model because it gives lowest misclassification rate.