# Assignment 2

shiva gadila

2023-11-11

QA1. What is the key idea behind bagging? Can bagging deal both with high variance(overfitting) and high bias (underfitting)?

ANSWER:- Bagging aims to reduce the variance (overfitting) of a machine-learning model. The core concept involves training multiple models on different random subsets of the dataset. This diversity in training data helps the models capture different aspects of the underlying patterns in the data. In regression or classification tasks, the models are merged by calculating their average (for regression) or through a majority vote (for classification). This process helps in mitigating errors introduced by each individual model, leading to a reduction in the final prediction's variability. Bagging is particularly effective in addressing high variance (overfitting) by mitigating the impact of outliers and noise in individual models. However, it may not directly deal with high bias (underfitting). To achieve a balanced model, ensemble methods like bagging work well when combined with diverse models. By incorporating different perspectives from various models, the ensemble can often strike a better balance between bias and variance.
Bagging reduces variance but may not directly handle high bias; combining it with boosting, which emphasizes misclassified instances, forms a potent strategy for addressing both high variance and high bias in machine learning models.

QA2. Why bagging models are computationally more efficient when compared to boosting models with the same number of weak learners?

ANSWER:- Bagging models demonstrate enhanced computational efficiency due to their ability to train weak learners independently in parallel. Each learner works on a random subset of the data, allowing them to be trained in isolation. This parallelization significantly speeds up the training process compared to the sequential training required by boosting models.

In bagging, the training of each weak learner is like a simultaneous, independent expedition into the data landscape. This parallel approach contrasts with boosting, where the training journey is more like a relay race. Boosting's sequential training method means each runner (weak learner) relies on the insights handed over by the previous one, limiting the speed of the entire process.

Moreover, boosting models often demand more rounds of refinement to reach their optimal solution, amplifying the computational load. In essence, while bagging models efficiently explore the data world in parallel, boosting models engage in a sequential relay,

making bagging computationally more effective, especially when considering both parallelization and iteration aspects.

QA3. James is thinking of creating an ensemble mode to predict whether a given stock will go up or down in the next week. He has trained several decision tree models but each model is not performing any better than a random model. The models are also very similar to each other. Do you think creating an ensemble model by combining these tree models can boost the performance? Discuss your answer.

ANSWER :-If James is facing a scenario where each individual decision tree model performs no better than a random model and the models are strikingly similar, the prospect of creating an ensemble model might not substantially improve performance.

Ensemble models like bagging and boosting thrive on aggregating diverse weak models to construct a potent predictor. However, if the individual models are both subpar and remarkably alike, the resulting ensemble model may not deliver the desired performance boost. In such cases, it might be more effective to first focus on improving the individual decision tree models. This could involve tweaking hyperparameters, exploring alternative approaches, or acquiring more data for model training.

Alternatively, James could consider a different ensemble method, like stacking, where the decision tree models' outputs serve as inputs for another model better suited for handling weak models. It's crucial to note that the success of any ensemble method hinges on the diversity and quality of the individual models it incorporates. Before diving into ensemble strategies, refining the foundation—the individual decision tree models—should be a priority.

QA4. Consider the following Table that classifies some objects into two classes of edible (+) and non- edible (-), based on some characteristics such as the object color, size and shape. What would be the Information gain for splitting the dataset based on the "Size" attribute?

```r
# PARENT
parent <- -((8/16) * log2(8/16) + (8/16) * log2(8/16))
print(parent)

## [1] 1

# SMALL
small <- -((6/8) * log2(6/8) + (2/8) * log2(2/8))
print(small)

## [1] 0.8112781

# LARGE
large <- -((3/8) * log2(3/8) + (5/8) * log2(5/8))
print(large)

## [1] 0.954434

# Calculate the mean entropy for SMALL and LARGE
values <- c(small, large)
```

```
mean_entropy <- mean(values)
print(mean_entropy)

## [1] 0.8828561

# Calculate Information Gain
info_gain <- parent - mean_entropy
print(info_gain)

## [1] 0.1171439
```

QA5. Why is it important that the m parameter (number of attributes available at each split) to be optimally set in random forest models? Discuss the implications of setting this parameter too small or too large In random forest models, setting the m parameter, which determines the number of attributes available at each split, is crucial for optimal performance. Here's a discussion on the implications of setting this parameter too small or too large:

Optimal Setting: The ideal value for the m parameter varies depending on the specific dataset and the problem being addressed. Finding the right balance is essential to ensure that the random forest model can effectively capture patterns and relationships in the data.

Setting m Too Small: If the m parameter is set too low, it leads to a reduction in the diversity among individual decision trees in the forest. This similarity may result in a limited ability to capture the full range of patterns and relationships in the data. Consequently, the model is prone to significant bias or underfitting, leading to a decrease in overall predictive performance.

Setting m Too Large: Conversely, if the m parameter is set too high, it encourages individual decision trees to be overly dissimilar. This diversity may lead to a forest that struggles to generalize well to new data, as the trees may become too specialized to the training set. This scenario increases the risk of overfitting, where the model performs well on the training data but fails to generalize effectively to unseen data. Finding the optimal setting for the m parameter is essential as it determines the delicate balance between diversity and similarity among decision trees in a random forest. This setting plays a direct role in shaping the model's capacity to generalize effectively and provide accurate predictions for new, unseen data.

PART 2

```
# Loading required libraries
library(ISLR)
library(dplyr)

##
## Attaching package: 'dplyr'

## The following objects are masked from 'package:stats':
##
##     filter, lag
```

```
## The following objects are masked from 'package:base':
##
##     intersect, setdiff, setequal, union

library(glmnet)

## Loading required package: Matrix

## Loaded glmnet 4.1-8

library(caret)

## Loading required package: ggplot2

## Loading required package: lattice

library(rpart)
library(rpart.plot)
```
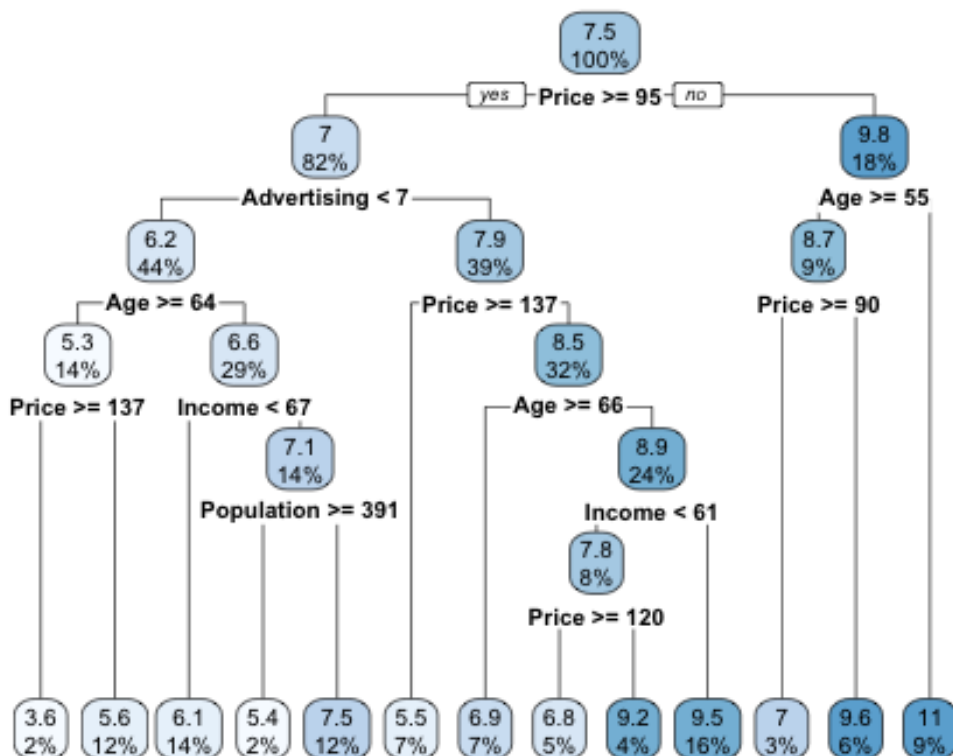
##Question 1 . Build a decision tree regression model to predict Sales based on all other attributes ("Price", "Advertising", "Population", "Age", "Income" and "Education"). Which attribute is used at the top of the tree (the root node) for splitting? Hint: you can either plot () and text() functions or use the summary() function to see the decision tree rule

```
#LOADING THE DATASET CARSEATS:
Carseats_Filtered <- Carseats%>%select("Sales","Price","Advertising","Populat
ion","Age","Income","Education")

model_1 <- rpart(Sales~., data=Carseats_Filtered, method = 'anova')
rpart.plot(model_1)
```

```
summary(model_1)

## Call:
## rpart(formula = Sales ~ ., data = Carseats_Filtered, method = "anova")
##   n= 400
##
##            CP nsplit rel error    xerror       xstd
## 1  0.14251535      0 1.0000000 1.0077985 0.06966818
## 2  0.08034146      1 0.8574847 0.8877077 0.06124215
## 3  0.06251702      2 0.7771432 0.8628614 0.06042133
## 4  0.02925241      3 0.7146262 0.8212436 0.05558607
## 5  0.02537341      4 0.6853738 0.8467304 0.05627570
## 6  0.02127094      5 0.6600003 0.8250220 0.05478347
## 7  0.02059174      6 0.6387294 0.8244607 0.05416011
## 8  0.01632010      7 0.6181377 0.8357625 0.05509159
## 9  0.01521801      8 0.6018176 0.8238084 0.05443783
## 10 0.01042023      9 0.5865996 0.8208132 0.05716180
## 11 0.01000559     10 0.5761793 0.8326996 0.05575377
## 12 0.01000000     12 0.5561681 0.8415278 0.05616440
##
## Variable importance
##        Price Advertising         Age      Income  Population   Education
##           49          18          16           8           6           3
```

```
## 
## Node number 1: 400 observations,    complexity param=0.1425153
##   mean=7.496325, MSE=7.955687
##   left son=2 (329 obs) right son=3 (71 obs)
##   Primary splits:
##       Price       < 94.5  to the right, improve=0.14251530, (0 missing)
##       Advertising < 7.5   to the left,  improve=0.07303226, (0 missing)
##       Age         < 61.5  to the right, improve=0.07120203, (0 missing)
##       Income      < 61.5  to the left,  improve=0.02840494, (0 missing)
##       Population  < 174.5 to the left,  improve=0.01077467, (0 missing)
## 
## Node number 2: 329 observations,    complexity param=0.08034146
##   mean=7.001672, MSE=6.815199
##   left son=4 (174 obs) right son=5 (155 obs)
##   Primary splits:
##       Advertising < 6.5   to the left,  improve=0.11402580, (0 missing)
##       Price       < 136.5 to the right, improve=0.08411056, (0 missing)
##       Age         < 63.5  to the right, improve=0.08091745, (0 missing)
##       Income      < 60.5  to the left,  improve=0.03394126, (0 missing)
##       Population  < 23    to the left,  improve=0.01831455, (0 missing)
##   Surrogate splits:
##       Population < 223   to the left,  agree=0.599, adj=0.148, (0 split)
##       Education  < 10.5  to the right, agree=0.565, adj=0.077, (0 split)
##       Age        < 53.5  to the right, agree=0.547, adj=0.039, (0 split)
##       Income     < 114.5 to the left,  agree=0.547, adj=0.039, (0 split)
##       Price      < 106.5 to the right, agree=0.544, adj=0.032, (0 split)
## 
## Node number 3: 71 observations,    complexity param=0.02537341
##   mean=9.788451, MSE=6.852836
##   left son=6 (36 obs) right son=7 (35 obs)
##   Primary splits:
##       Age        < 54.5  to the right, improve=0.16595410, (0 missing)
##       Price      < 75.5  to the right, improve=0.08365773, (0 missing)
##       Income     < 30.5  to the left,  improve=0.03322169, (0 missing)
##       Education  < 10.5  to the right, improve=0.03019634, (0 missing)
##       Population < 268.5 to the left,  improve=0.02383306, (0 missing)
##   Surrogate splits:
##       Advertising < 4.5   to the right, agree=0.606, adj=0.200, (0 split)
##       Price       < 73    to the right, agree=0.592, adj=0.171, (0 split)
##       Population  < 272.5 to the left,  agree=0.592, adj=0.171, (0 split)
##       Income      < 79.5  to the right, agree=0.592, adj=0.171, (0 split)
##       Education   < 11.5  to the left,  agree=0.577, adj=0.143, (0 split)
## 
## Node number 4: 174 observations,    complexity param=0.02127094
##   mean=6.169655, MSE=4.942347
##   left son=8 (58 obs) right son=9 (116 obs)
##   Primary splits:
##       Age        < 63.5  to the right, improve=0.078712160, (0 missing)
##       Price      < 130.5 to the right, improve=0.048919280, (0 missing)
##       Population < 26.5  to the left,  improve=0.030421540, (0 missing)
```

```
##         Income       < 67.5  to the left,  improve=0.027749670, (0 missing)
##         Advertising < 0.5   to the left,  improve=0.006795377, (0 missing)
##    Surrogate splits:
##         Income       < 22.5  to the left,  agree=0.678, adj=0.034, (0 split)
##         Price        < 96.5  to the left,  agree=0.672, adj=0.017, (0 split)
##         Population < 26.5  to the left,  agree=0.672, adj=0.017, (0 split)
##
## Node number 5: 155 observations,    complexity param=0.06251702
##   mean=7.935677, MSE=7.268151
##   left son=10 (28 obs) right son=11 (127 obs)
##   Primary splits:
##         Price        < 136.5 to the right, improve=0.17659580, (0 missing)
##         Age          < 73.5  to the right, improve=0.08000201, (0 missing)
##         Income       < 60.5  to the left,  improve=0.05360755, (0 missing)
##         Advertising < 13.5  to the left,  improve=0.03920507, (0 missing)
##         Population  < 399   to the left,  improve=0.01037956, (0 missing)
##    Surrogate splits:
##         Advertising < 24.5  to the right, agree=0.826, adj=0.036, (0 split)
##
## Node number 6: 36 observations,    complexity param=0.0163201
##   mean=8.736944, MSE=4.961043
##   left son=12 (12 obs) right son=13 (24 obs)
##   Primary splits:
##         Price        < 89.5  to the right, improve=0.29079360, (0 missing)
##         Income       < 39.5  to the left,  improve=0.19043350, (0 missing)
##         Advertising < 11.5  to the left,  improve=0.17891930, (0 missing)
##         Age          < 75.5  to the right, improve=0.04316067, (0 missing)
##         Education   < 14.5  to the left,  improve=0.03411396, (0 missing)
##    Surrogate splits:
##         Advertising < 16.5  to the right, agree=0.722, adj=0.167, (0 split)
##         Income       < 37.5  to the left,  agree=0.722, adj=0.167, (0 split)
##         Age          < 56.5  to the left,  agree=0.694, adj=0.083, (0 split)
##
## Node number 7: 35 observations
##   mean=10.87, MSE=6.491674
##
## Node number 8: 58 observations,    complexity param=0.01042023
##   mean=5.287586, MSE=3.93708
##   left son=16 (10 obs) right son=17 (48 obs)
##   Primary splits:
##         Price       < 137   to the right, improve=0.14521540, (0 missing)
##         Education  < 15.5  to the right, improve=0.07995394, (0 missing)
##         Income      < 35.5  to the left,  improve=0.04206708, (0 missing)
##         Age         < 79.5  to the left,  improve=0.02799057, (0 missing)
##         Population < 52.5  to the left,  improve=0.01914342, (0 missing)
##
## Node number 9: 116 observations,    complexity param=0.01000559
##   mean=6.61069, MSE=4.861446
##   left son=18 (58 obs) right son=19 (58 obs)
##    Primary splits:
```

```
##          Income      < 67      to the left,   improve=0.05085914, (0 missing)
##          Population < 392    to the right, improve=0.04476721, (0 missing)
##          Price       < 127    to the right, improve=0.04210762, (0 missing)
##          Age         < 37.5   to the right, improve=0.02858424, (0 missing)
##          Education  < 14.5   to the left,   improve=0.01187387, (0 missing)
##    Surrogate splits:
##          Education  < 12.5   to the right, agree=0.586, adj=0.172, (0 split)
##          Age         < 58.5   to the left,   agree=0.578, adj=0.155, (0 split)
##          Price       < 144.5 to the left,   agree=0.569, adj=0.138, (0 split)
##          Population < 479    to the right, agree=0.560, adj=0.121, (0 split)
##          Advertising < 2.5    to the right, agree=0.543, adj=0.086, (0 split)
##
## Node number 10: 28 observations
##    mean=5.522857, MSE=5.084213
##
## Node number 11: 127 observations,     complexity param=0.02925241
##    mean=8.467638, MSE=6.183142
##    left son=22 (29 obs) right son=23 (98 obs)
##    Primary splits:
##          Age         < 65.5   to the right, improve=0.11854590, (0 missing)
##          Income      < 51.5   to the left,   improve=0.08076060, (0 missing)
##          Advertising < 13.5   to the left,   improve=0.04801701, (0 missing)
##          Education  < 11.5   to the right, improve=0.02471512, (0 missing)
##          Population < 479    to the left,   improve=0.01908657, (0 missing)
##
## Node number 12: 12 observations
##    mean=7.038333, MSE=2.886964
##
## Node number 13: 24 observations
##    mean=9.58625, MSE=3.834123
##
## Node number 16: 10 observations
##    mean=3.631, MSE=5.690169
##
## Node number 17: 48 observations
##    mean=5.632708, MSE=2.88102
##
## Node number 18: 58 observations
##    mean=6.113448, MSE=3.739109
##
## Node number 19: 58 observations,     complexity param=0.01000559
##    mean=7.107931, MSE=5.489285
##    left son=38 (10 obs) right son=39 (48 obs)
##    Primary splits:
##          Population < 390.5 to the right, improve=0.10993270, (0 missing)
##          Price       < 124.5 to the right, improve=0.07534567, (0 missing)
##          Advertising < 0.5    to the left,   improve=0.07060488, (0 missing)
##          Age         < 45.5   to the right, improve=0.04611510, (0 missing)
##          Education  < 11.5   to the right, improve=0.03722944, (0 missing)
##
```

```
## Node number 22: 29 observations
##    mean=6.893793, MSE=6.08343
##
## Node number 23: 98 observations,     complexity param=0.02059174
##    mean=8.933367, MSE=5.262759
##    left son=46 (34 obs) right son=47 (64 obs)
##    Primary splits:
##        Income      < 60.5  to the left,  improve=0.12705480, (0 missing)
##        Advertising < 13.5  to the left,  improve=0.07114001, (0 missing)
##        Price       < 118.5 to the right, improve=0.06932216, (0 missing)
##        Education   < 11.5  to the right, improve=0.03377416, (0 missing)
##        Age         < 49.5  to the right, improve=0.02289004, (0 missing)
##    Surrogate splits:
##        Education < 17.5  to the right, agree=0.663, adj=0.029, (0 split)
##
## Node number 38: 10 observations
##    mean=5.406, MSE=2.508524
##
## Node number 39: 48 observations
##    mean=7.4625, MSE=5.381106
##
## Node number 46: 34 observations,     complexity param=0.01521801
##    mean=7.811471, MSE=4.756548
##    left son=92 (19 obs) right son=93 (15 obs)
##    Primary splits:
##        Price       < 119.5 to the right, improve=0.29945020, (0 missing)
##        Advertising < 11.5  to the left,  improve=0.14268440, (0 missing)
##        Income      < 40.5  to the right, improve=0.12781140, (0 missing)
##        Population  < 152   to the left,  improve=0.03601768, (0 missing)
##        Age         < 49.5  to the right, improve=0.02748814, (0 missing)
##    Surrogate splits:
##        Education   < 12.5  to the right, agree=0.676, adj=0.267, (0 split)
##        Advertising < 7.5   to the right, agree=0.647, adj=0.200, (0 split)
##        Age         < 53.5  to the left,  agree=0.647, adj=0.200, (0 split)
##        Population  < 240   to the right, agree=0.618, adj=0.133, (0 split)
##        Income      < 41.5  to the right, agree=0.618, adj=0.133, (0 split)
##
## Node number 47: 64 observations
##    mean=9.529375, MSE=4.5078
##
## Node number 92: 19 observations
##    mean=6.751053, MSE=3.378915
##
## Node number 93: 15 observations
##    mean=9.154667, MSE=3.273025
```

##The price attribute is used as the root node for spitting as it has the highest value in value importance (using summary function) and also we can see the decision tree plot where the price attribute is the first root node.

```
summary(model_1)

## Call:
## rpart(formula = Sales ~ ., data = Carseats_Filtered, method = "anova")
##   n= 400
##
##            CP nsplit rel error    xerror      xstd
## 1  0.14251535      0 1.0000000 1.0077985 0.06966818
## 2  0.08034146      1 0.8574847 0.8877077 0.06124215
## 3  0.06251702      2 0.7771432 0.8628614 0.06042133
## 4  0.02925241      3 0.7146262 0.8212436 0.05558607
## 5  0.02537341      4 0.6853738 0.8467304 0.05627570
## 6  0.02127094      5 0.6600003 0.8250220 0.05478347
## 7  0.02059174      6 0.6387294 0.8244607 0.05416011
## 8  0.01632010      7 0.6181377 0.8357625 0.05509159
## 9  0.01521801      8 0.6018176 0.8238084 0.05443783
## 10 0.01042023      9 0.5865996 0.8208132 0.05716180
## 11 0.01000559     10 0.5761793 0.8326996 0.05575377
## 12 0.01000000     12 0.5561681 0.8415278 0.05616440
##
## Variable importance
##       Price Advertising         Age      Income  Population   Education
##          49          18          16           8           6           3
##
## Node number 1: 400 observations,    complexity param=0.1425153
##   mean=7.496325, MSE=7.955687
##   left son=2 (329 obs) right son=3 (71 obs)
##   Primary splits:
##       Price       < 94.5  to the right, improve=0.14251530, (0 missing)
##       Advertising < 7.5   to the left,  improve=0.07303226, (0 missing)
##       Age         < 61.5  to the right, improve=0.07120203, (0 missing)
##       Income      < 61.5  to the left,  improve=0.02840494, (0 missing)
##       Population  < 174.5 to the left,  improve=0.01077467, (0 missing)
##
## Node number 2: 329 observations,    complexity param=0.08034146
##   mean=7.001672, MSE=6.815199
##   left son=4 (174 obs) right son=5 (155 obs)
##   Primary splits:
##       Advertising < 6.5   to the left,  improve=0.11402580, (0 missing)
##       Price       < 136.5 to the right, improve=0.08411056, (0 missing)
##       Age         < 63.5  to the right, improve=0.08091745, (0 missing)
##       Income      < 60.5  to the left,  improve=0.03394126, (0 missing)
##       Population  < 23    to the left,  improve=0.01831455, (0 missing)
##   Surrogate splits:
##       Population < 223   to the left,  agree=0.599, adj=0.148, (0 split)
##       Education  < 10.5  to the right, agree=0.565, adj=0.077, (0 split)
##       Age        < 53.5  to the right, agree=0.547, adj=0.039, (0 split)
##       Income     < 114.5 to the left,  agree=0.547, adj=0.039, (0 split)
##       Price      < 106.5 to the right, agree=0.544, adj=0.032, (0 split)
##
```

```
## Node number 3: 71 observations,    complexity param=0.02537341
##   mean=9.788451, MSE=6.852836
##   left son=6 (36 obs) right son=7 (35 obs)
##   Primary splits:
##       Age         < 54.5  to the right, improve=0.16595410, (0 missing)
##       Price       < 75.5  to the right, improve=0.08365773, (0 missing)
##       Income      < 30.5  to the left,  improve=0.03322169, (0 missing)
##       Education   < 10.5  to the right, improve=0.03019634, (0 missing)
##       Population  < 268.5 to the left,  improve=0.02383306, (0 missing)
##   Surrogate splits:
##       Advertising < 4.5   to the right, agree=0.606, adj=0.200, (0 split)
##       Price       < 73    to the right, agree=0.592, adj=0.171, (0 split)
##       Population  < 272.5 to the left,  agree=0.592, adj=0.171, (0 split)
##       Income      < 79.5  to the right, agree=0.592, adj=0.171, (0 split)
##       Education   < 11.5  to the left,  agree=0.577, adj=0.143, (0 split)
##
## Node number 4: 174 observations,    complexity param=0.02127094
##   mean=6.169655, MSE=4.942347
##   left son=8 (58 obs) right son=9 (116 obs)
##   Primary splits:
##       Age         < 63.5  to the right, improve=0.078712160, (0 missing)
##       Price       < 130.5 to the right, improve=0.048919280, (0 missing)
##       Population  < 26.5  to the left,  improve=0.030421540, (0 missing)
##       Income      < 67.5  to the left,  improve=0.027749670, (0 missing)
##       Advertising < 0.5   to the left,  improve=0.006795377, (0 missing)
##   Surrogate splits:
##       Income      < 22.5  to the left,  agree=0.678, adj=0.034, (0 split)
##       Price       < 96.5  to the left,  agree=0.672, adj=0.017, (0 split)
##       Population  < 26.5  to the left,  agree=0.672, adj=0.017, (0 split)
##
## Node number 5: 155 observations,    complexity param=0.06251702
##   mean=7.935677, MSE=7.268151
##   left son=10 (28 obs) right son=11 (127 obs)
##   Primary splits:
##       Price       < 136.5 to the right, improve=0.17659580, (0 missing)
##       Age         < 73.5  to the right, improve=0.08000201, (0 missing)
##       Income      < 60.5  to the left,  improve=0.05360755, (0 missing)
##       Advertising < 13.5  to the left,  improve=0.03920507, (0 missing)
##       Population  < 399   to the left,  improve=0.01037956, (0 missing)
##   Surrogate splits:
##       Advertising < 24.5  to the right, agree=0.826, adj=0.036, (0 split)
##
## Node number 6: 36 observations,    complexity param=0.0163201
##   mean=8.736944, MSE=4.961043
##   left son=12 (12 obs) right son=13 (24 obs)
##   Primary splits:
##       Price       < 89.5  to the right, improve=0.29079360, (0 missing)
##       Income      < 39.5  to the left,  improve=0.19043350, (0 missing)
##       Advertising < 11.5  to the left,  improve=0.17891930, (0 missing)
##       Age         < 75.5  to the right, improve=0.04316067, (0 missing)
```

```
##         Education   < 14.5  to the left,   improve=0.03411396, (0 missing)
##     Surrogate splits:
##         Advertising < 16.5  to the right, agree=0.722, adj=0.167, (0 split)
##         Income      < 37.5  to the left,  agree=0.722, adj=0.167, (0 split)
##         Age         < 56.5  to the left,  agree=0.694, adj=0.083, (0 split)
##
## Node number 7: 35 observations
##   mean=10.87, MSE=6.491674
##
## Node number 8: 58 observations,    complexity param=0.01042023
##   mean=5.287586, MSE=3.93708
##   left son=16 (10 obs) right son=17 (48 obs)
##   Primary splits:
##         Price      < 137  to the right, improve=0.14521540, (0 missing)
##         Education  < 15.5 to the right, improve=0.07995394, (0 missing)
##         Income     < 35.5 to the left,  improve=0.04206708, (0 missing)
##         Age        < 79.5 to the left,  improve=0.02799057, (0 missing)
##         Population < 52.5 to the left,  improve=0.01914342, (0 missing)
##
## Node number 9: 116 observations,    complexity param=0.01000559
##   mean=6.61069, MSE=4.861446
##   left son=18 (58 obs) right son=19 (58 obs)
##   Primary splits:
##         Income     < 67   to the left,  improve=0.05085914, (0 missing)
##         Population < 392  to the right, improve=0.04476721, (0 missing)
##         Price      < 127  to the right, improve=0.04210762, (0 missing)
##         Age        < 37.5 to the right, improve=0.02858424, (0 missing)
##         Education  < 14.5 to the left,  improve=0.01187387, (0 missing)
##     Surrogate splits:
##         Education   < 12.5  to the right, agree=0.586, adj=0.172, (0 split)
##         Age         < 58.5  to the left,  agree=0.578, adj=0.155, (0 split)
##         Price       < 144.5 to the left,  agree=0.569, adj=0.138, (0 split)
##         Population  < 479   to the right, agree=0.560, adj=0.121, (0 split)
##         Advertising < 2.5   to the right, agree=0.543, adj=0.086, (0 split)
##
## Node number 10: 28 observations
##   mean=5.522857, MSE=5.084213
##
## Node number 11: 127 observations,    complexity param=0.02925241
##   mean=8.467638, MSE=6.183142
##   left son=22 (29 obs) right son=23 (98 obs)
##   Primary splits:
##         Age         < 65.5 to the right, improve=0.11854590, (0 missing)
##         Income      < 51.5 to the left,  improve=0.08076060, (0 missing)
##         Advertising < 13.5 to the left,  improve=0.04801701, (0 missing)
##         Education   < 11.5 to the right, improve=0.02471512, (0 missing)
##         Population  < 479  to the left,  improve=0.01908657, (0 missing)
##
## Node number 12: 12 observations
##   mean=7.038333, MSE=2.886964
```

```
## 
## Node number 13: 24 observations
##    mean=9.58625, MSE=3.834123
## 
## Node number 16: 10 observations
##    mean=3.631, MSE=5.690169
## 
## Node number 17: 48 observations
##    mean=5.632708, MSE=2.88102
## 
## Node number 18: 58 observations
##    mean=6.113448, MSE=3.739109
## 
## Node number 19: 58 observations,    complexity param=0.01000559
##    mean=7.107931, MSE=5.489285
##    left son=38 (10 obs) right son=39 (48 obs)
##    Primary splits:
##        Population  < 390.5 to the right, improve=0.10993270, (0 missing)
##        Price       < 124.5 to the right, improve=0.07534567, (0 missing)
##        Advertising < 0.5   to the left,  improve=0.07060488, (0 missing)
##        Age         < 45.5  to the right, improve=0.04611510, (0 missing)
##        Education   < 11.5  to the right, improve=0.03722944, (0 missing)
## 
## Node number 22: 29 observations
##    mean=6.893793, MSE=6.08343
## 
## Node number 23: 98 observations,    complexity param=0.02059174
##    mean=8.933367, MSE=5.262759
##    left son=46 (34 obs) right son=47 (64 obs)
##    Primary splits:
##        Income      < 60.5  to the left,  improve=0.12705480, (0 missing)
##        Advertising < 13.5  to the left,  improve=0.07114001, (0 missing)
##        Price       < 118.5 to the right, improve=0.06932216, (0 missing)
##        Education   < 11.5  to the right, improve=0.03377416, (0 missing)
##        Age         < 49.5  to the right, improve=0.02289004, (0 missing)
##    Surrogate splits:
##        Education < 17.5  to the right, agree=0.663, adj=0.029, (0 split)
## 
## Node number 38: 10 observations
##    mean=5.406, MSE=2.508524
## 
## Node number 39: 48 observations
##    mean=7.4625, MSE=5.381106
## 
## Node number 46: 34 observations,    complexity param=0.01521801
##    mean=7.811471, MSE=4.756548
##    left son=92 (19 obs) right son=93 (15 obs)
##    Primary splits:
##        Price       < 119.5 to the right, improve=0.29945020, (0 missing)
##        Advertising < 11.5  to the left,  improve=0.14268440, (0 missing)
```

```
##         Income      < 40.5   to the right, improve=0.12781140, (0 missing)
##         Population  < 152    to the left,  improve=0.03601768, (0 missing)
##         Age         < 49.5   to the right, improve=0.02748814, (0 missing)
##    Surrogate splits:
##         Education   < 12.5   to the right, agree=0.676, adj=0.267, (0 split)
##         Advertising < 7.5    to the right, agree=0.647, adj=0.200, (0 split)
##         Age         < 53.5   to the left,  agree=0.647, adj=0.200, (0 split)
##         Population  < 240    to the right, agree=0.618, adj=0.133, (0 split)
##         Income      < 41.5   to the right, agree=0.618, adj=0.133, (0 split)
##
## Node number 47: 64 observations
##    mean=9.529375, MSE=4.5078
##
## Node number 92: 19 observations
##    mean=6.751053, MSE=3.378915
##
## Node number 93: 15 observations
##    mean=9.154667, MSE=3.273025
```

## Queston 2

Consider the following input: Sales=9, Price=6.54, Population=124, Advertising=0, Age=76, Income= 110, Education=10 What will be the estimated Sales for this record using the decision tree model?

```
Sales <- c(9)
Price <- c(6.54)
Population <- c(124)
Advertising <- c(0)
Age <- c(76)
Income <- c(110)
Education <- c(10)
Test <- data.frame(Sales,Price,Population,Advertising,Age,Income,Education)
```

Now that our test set is prepared for evaluation within our model, the next step involves forecasting sales based on the provided data.

```
Pred_sales_2 <- predict(model_1, Test)
Pred_sales_2

##       1
## 9.58625
```

As per the prediction from our model using the predict function, the decision tree suggests that sales for this given record are estimated to be approximately 9.58625 units.

## Question 3

Use the caret function to train a random forest (method='rf') for the same dataset. Use the caret default settings. By default, caret will examine the "mtry" values of 2,4, and 6.Recall

that mtry is the number of attributes available for splitting at each splitting node. Which mtry value gives the best performance?
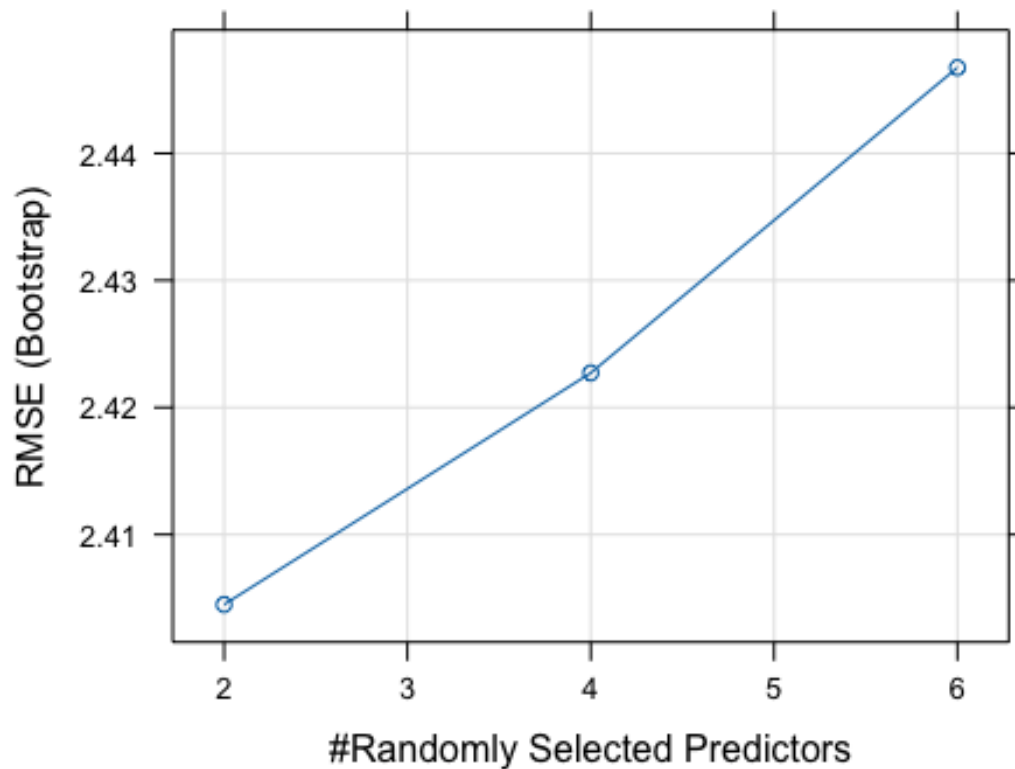
```
set.seed(123)
Model_forest_caret <- train(Sales~., data = Carseats_Filtered, method = 'rf')

summary(Model_forest_caret)

##                  Length Class      Mode
## call                 4  -none-     call
## type                 1  -none-     character
## predicted          400  -none-     numeric
## mse                500  -none-     numeric
## rsq                500  -none-     numeric
## oob.times          400  -none-     numeric
## importance           6  -none-     numeric
## importanceSD         0  -none-     NULL
## localImportance      0  -none-     NULL
## proximity            0  -none-     NULL
## ntree                1  -none-     numeric
## mtry                 1  -none-     numeric
## forest              11  -none-     list
## coefs                0  -none-     NULL
## y                  400  -none-     numeric
## test                 0  -none-     NULL
## inbag                0  -none-     NULL
## xNames               6  -none-     character
## problemType          1  -none-     character
## tuneValue            1  data.frame list
## obsLevels            1  -none-     logical
## param                0  -none-     list

print(Model_forest_caret)

## Random Forest
##
## 400 samples
##   6 predictor
##
## No pre-processing
## Resampling: Bootstrapped (25 reps)
## Summary of sample sizes: 400, 400, 400, 400, 400, 400, ...
## Resampling results across tuning parameters:
##
##   mtry  RMSE      Rsquared   MAE
##   2     2.404492  0.2860284  1.925868
##   4     2.422719  0.2782951  1.935483
##   6     2.446748  0.2686232  1.952884
##
## RMSE was used to select the optimal model using the smallest value.
## The final value used for the model was mtry = 2.
```

```
plot(Model_forest_caret)
```



#Since **2** mtry has the lowest RMSE, this value is the best fit for mtry.

## ##QUESTION 4

Customize the search grid by checking the model's performance for mtry values of 2, 3 and 5 using 3 repeats of 5-fold cross validation.

```
control <- trainControl(method="repeatedcv", number=5, repeats=3, search="gri
d")
tunegrid <- expand.grid(.mtry=c(2,3,5))
rf_gridsearch <- train(Sales~., data=Carseats_Filtered, method="rf", tuneGrid
=tunegrid,trControl=control)
print(rf_gridsearch)

## Random Forest
##
## 400 samples
##   6 predictor
##
## No pre-processing
## Resampling: Cross-Validated (5 fold, repeated 3 times)
## Summary of sample sizes: 321, 320, 320, 320, 319, 320, ...
```

```
## Resampling results across tuning parameters:
##
##   mtry  RMSE       Rsquared   MAE
##   2     2.388490   0.2902905  1.902942
##   3     2.390502   0.2898689  1.899672
##   5     2.402758   0.2869045  1.905036
##
## RMSE was used to select the optimal model using the smallest value.
## The final value used for the model was mtry = 2.
```

```
# Create a plot to visualize the results of the grid search
plot(rf_gridsearch)
```