

Decision Trees and Random Forests

Problem Statement: This project involves building decision tree and random forest models to answer a number of questions. We will use the Carseats dataset that is part of the ISLR package.

Loading Required Libraries

```
library(ISLR)

library(dplyr)

library(glmnet)

library(caret)

library(rpart)

library(rpart.plot)

library(rattle)

summary(Carseats)
```

```
##      Sales      CompPrice      Income      Advertising
## Min.   : 0.000   Min.   : 77   Min.   : 21.00   Min.   : 0.000
## 1st Qu.: 5.390   1st Qu.:115   1st Qu.: 42.75   1st Qu.: 0.000
## Median : 7.490   Median :125   Median : 69.00   Median : 5.000
## Mean   : 7.496   Mean   :125   Mean   : 68.66   Mean   : 6.635
## 3rd Qu.: 9.320   3rd Qu.:135   3rd Qu.: 91.00   3rd Qu.:12.000
## Max.   :16.270   Max.   :175   Max.   :120.00   Max.   :29.000
##      Population      Price      ShelfLoc      Age      Education
## Min.   : 10.0   Min.   : 24.0   Bad   : 96   Min.   :25.00   Min.   :10.0
## 1st Qu.:139.0   1st Qu.:100.0   Good  : 85   1st Qu.:39.75   1st Qu.:12.0
## Median :272.0   Median :117.0   Medium:219   Median :54.50   Median :14.0
## Mean   :264.8   Mean   :115.8                      Mean   :53.32   Mean   :13.9
## 3rd Qu.:398.5   3rd Qu.:131.0                      3rd Qu.:66.00   3rd Qu.:16.0
## Max.   :509.0   Max.   :191.0                      Max.   :80.00   Max.   :18.0
##      Urban      US
## No :118   No :142
## Yes:282   Yes:258
##
##
##
##
```

Selecting required attributes

```
Carseats_Filtered = Carseats %>% select("Sales", "Price",
"Advertising", "Population", "Age", "Income", "Education")
summary(Carseats_Filtered)
```

```
##      Sales      Price      Advertising      Population
## Min.   : 0.000   Min.   : 24.0   Min.   : 0.000   Min.   : 10.0
## 1st Qu.: 5.390   1st Qu.:100.0   1st Qu.: 0.000   1st Qu.:139.0
## Median : 7.490   Median :117.0   Median : 5.000   Median :272.0
## Mean   : 7.496   Mean   :115.8   Mean   : 6.635   Mean   :264.8
## 3rd Qu.: 9.320   3rd Qu.:131.0   3rd Qu.:12.000   3rd Qu.:398.5
## Max.   :16.270   Max.   :191.0   Max.   :29.000   Max.   :509.0
##      Age      Income      Education
## Min.   :25.00   Min.   : 21.00   Min.   :10.0
## 1st Qu.:39.75   1st Qu.: 42.75   1st Qu.:12.0
## Median :54.50   Median : 69.00   Median :14.0
## Mean   :53.32   Mean   : 68.66   Mean   :13.9
## 3rd Qu.:66.00   3rd Qu.: 91.00   3rd Qu.:16.0
## Max.   :80.00   Max.   :120.00   Max.   :18.0
```

Developing a decision tree regression model to forecast Sales using required attributes.(“Price”, “Advertising”, “Population”, “Age”, “Income” and “Education”).

```
model_1 = rpart(Sales~.,data=Carseats_Filtered,method = 'anova')
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 1.0024099 0.06917760
## 2  0.08034146      1 0.8574847 0.9238326 0.06600774
## 3  0.06251702      2 0.7771432 0.9245128 0.06555144
## 4  0.02925241      3 0.7146262 0.8356584 0.06054659
## 5  0.02537341      4 0.6853738 0.8522120 0.06059264
## 6  0.02127094      5 0.6600003 0.8388419 0.06010782
## 7  0.02059174      6 0.6387294 0.8506578 0.06007361
## 8  0.01632010      7 0.6181377 0.8336266 0.05803310
## 9  0.01521801      8 0.6018176 0.8467105 0.05932572
##10 0.01042023      9 0.5865996 0.8676741 0.06269914
##11 0.01000559     10 0.5761793 0.8783005 0.06252426
##12 0.01000000     12 0.5561681 0.8812926 0.06252257
##
## 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)

```

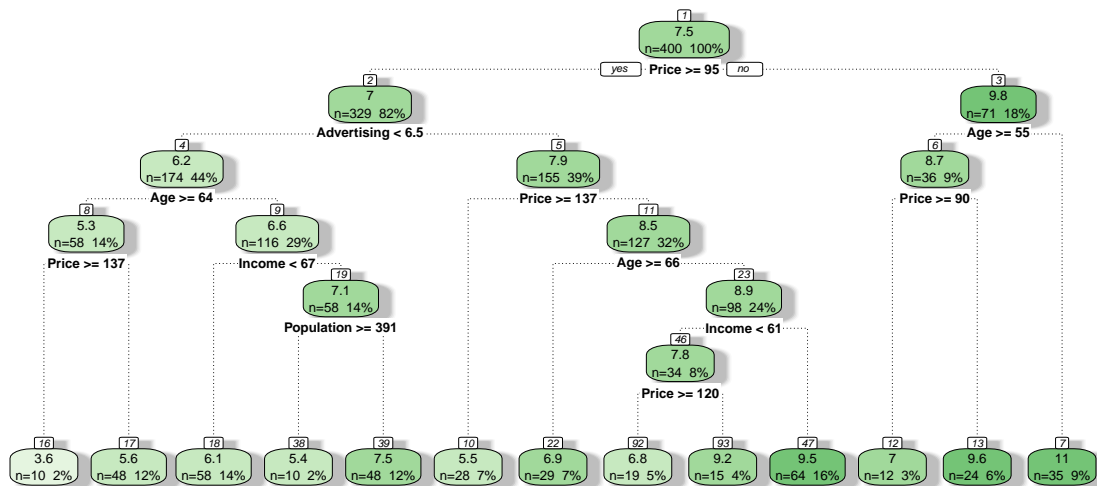
```

##      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
```

```
fancyRpartPlot(model_1)
```



Rattle 2024-Mar-07 02:10:49 r2001219

The price

attribute is utilized at the root node of the tree for the purpose of splitting.

Considering the following input: Sales=15,Price=8.71,Population=120,Advertising=0,Age=71,Income=110,Education=10.Estimating Sales for this record using the decision tree model

```
prediction_data = data.frame(Price=8.71,Population=120,Advertising=0,Age=71
,Income = 110, Education = 10)
```

```
prediction<- predict(model_1,prediction_data)
```

```
prediction
```

```
##          1
## 9.58625
```

Predicted sales value for this record is 9.58625.

Using the caret function to train a random forest (method='rf') for the same dataset.

```
set.seed(123)
model_2 <- train(Sales~.,
                 data= Carseats_Filtered,
                 method = "rf")
```

```
# Print the results
model_2
```

```
## 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    2.405819  0.2852547  1.926801
##   4    2.421577  0.2790266  1.934608
##   6    2.447373  0.2681323  1.953147
##
## RMSE was used to select the optimal model using the smallest value.
## The final value used for the model was mtry = 2.
```

Best results are obtained when mtry value is set to be 2.

Customizing 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.

```
set.seed(123)

#Cross-Validation
control <- trainControl(method = "repeatedcv",
                        repeats = 3,
                        number = 5)

#Defining search grid with mtry values of 2,3,5
search_grid <- expand.grid(mtry = c(2, 3, 5))

# Training the model using the search grid and cross-validation
model <- train(Sales~., Carseats_Filtered, method = "rf", tuneGrid = search_grid, trControl = control)
print(model)
```

```
## Random Forest
##
## 400 samples
##   6 predictor
```

```

##
## No pre-processing
## Resampling: Cross-Validated (5 fold, repeated 3 times)
## Summary of sample sizes: 320, 321, 319, 320, 320, 319, ...
## Resampling results across tuning parameters:
##
##   mtry  RMSE      Rsquared  MAE
##   2     2.405235  0.2813795  1.930855
##   3     2.401365  0.2858295  1.920612
##   5     2.417771  0.2821938  1.934886
##
## RMSE was used to select the optimal model using the smallest value.
## The final value used for the model was mtry = 3.

```