

16/4/18

## \* Classification and Regression Trees :- (CART)

- Build a tree by splitting on variables
- To predict an outcome for an observation, follow the splits and at the end, predict the most frequent outcome.
- There are different ways to control the splits {<sup>less</sup>/<sub>more</sub>}
  - (1) By setting a lower bound for the number of points in each subset.  
In R it is called minbucket
    - The smaller it is, the more splits will be generated.
    - If too small - overfitting will occur
    - If too large - model will be simple with poor accuracy.

- Stevens = read.csv("stevens.csv")

- library(caret)

spl = sample.split(stevens\$Reverso, SplitRatio = 0.7)

Train = subset(stevens, spl == TRUE)

Test = subset(stevens, spl == FALSE)

- install.packages("rpart")

install.packages("rpart.plot")

Creating CART

- StevensTree = rpart(Reverso ~ Circuit + ... + Uncert,

data = Train,

method = "class", minbucket = 25)

- `prp (StevensTree)` To plot the tree

[CART is a series of decision rules which can easily be explained]

Predictions with CART

- `PredictCART = predict(StevensTree, newdata = Test, type = "class")`

Creating Confusion matrix

- `table(Test$Reverse, PredictCART)`

	0	1
0	41	36
1	22	71

$$\text{Accuracy} = \frac{(TN + TP)}{N}$$
$$(41 + 71) / (41 + 36 + 22 + 71) = 65.9\%$$

Generating ROC curve

- `library(ROCR)`

`PredictRoc = predict(StevensTree, newdata = Test)`

Plotting ROC curves

- `pred = prediction(PredictRoc [2], Test$Reverse)`

`perf = performance(pred, "tpr", "fpr")`

$\uparrow$                        $\uparrow$   
true +ve rate    false +ve rate

`plot(pred)`

`plot(pred, colorize = TRUE)`

Generating AUC

- `performance(pred, "auc") @ y.values`

## \* Random Forests :

- Designed to improve prediction accuracy of CART
- Builds large number of CART trees.
- Each tree can split on only a random subset of the variables.
- Each tree is built from a "bagged" / "bootstrapped" sample of the data.

### Parameters:

(1) Minimum number of observations in a subset (minbucket parameter from cart)

- In R, controlled by nodesize parameter

- Smaller nodesize leads to bigger tree  
may take longer in R.

(2) Number of trees to build

- In R, called the ntree parameter

- install package ("randomForest")

Creating Random Forest

- Stevens Forest = randomForest (Reverse ~ Circuit + ... + Unew)

data = Train,

nodesize = 25,

ntree = 200)

[Trees can be used for Classification problems &  
Regression problems]

[When we do a Classification problem, the outcome should  
be a factor. So, convert Reverse to factor in Train &  
Test set]

- Train & Reverse = id. factors (Train & Reverse)
- Test & Reverse = id. factors (Test & Reverse)

Predictions with Random Forest

- Predict Forest = predict (Stevens Forest, newdata = Test)

Creating Confusion Matrix

- Table (Test & Reverse, Predict Forest)

0	40	27
1	19	74

$$\text{Accuracy} = \frac{(40 + 74)}{(40 + 27 + 19 + 74)} \\ = 0.67\%$$

#### \* CART

- Interpretable "
- White Box "
- Transparent "

#### Random Forests

- Black Box "
- Unclear "
- High Accuracy
- Low Interpretability "

#### \* Optimal Classification Tree : (OCT)

- Competitive with Black Box and is interpretable
- OCT - tree with parallel splits (one variable/split).
- OCT-H - tree with hyperplane splits (multiple variables/split)

- a method to predict misclassification parameter in CART method.
- We first split the given training set into  $K$  equal pieces.
- Use  $K-1$  folds to estimate a model and test model on remaining one fold ("validation set")
- Repeat the above step for each  $K$  folds
- When we use cross-validation in R, we'll use a parameter called Complexity Parameter ( $cp$ ).
- Smaller  $cp$  leads to bigger tree (might overfit)
- install.packages ("caret")
- install.packages ("e1071")

defining how many folds

- numFolds = trainControl (method = "cv", number = 10)

Getting the possible value for CP parameter:

- cpGrid = expand.grid (cp = seq(0.01, 0.9, 0.01))  
 $(cp = (0:10) * 0.01)$

Performing Cross Validation

- train (Reruse ~ Circuit + ... Uncat, data = Train, method = "rpart")
  - trainControl = numFolds,
  - tuneGrid = cpGrid)

- the best cp will be the obj (Type 'tr' to get it)

Creating the cart model & with cp

- $\text{StevensTreeCV} = \text{rpart}(\text{Paverse} \sim \text{Circuit} + \dots + \text{District},$   
data = Train,  
method = "class",  
cp = 0.18)

$$\text{Accuracy} = 72.35\%$$

## +> Practicing CART Models :

Finding Accuracy for Baseline Model :

- Claims = read.csv("ClaimsData.csv")

Now split the data in Train and Test set using  
sample.split

- ClaimsTrain, ClaimsTest

↓ Table 1 (63)

- Table(ClaimsTest & bucket 200, ClaimsTest & bucket 2008)

$$\text{Accuracy} = 0.6828$$

To find PenaltyAve

- PenaltyMatrix = matrix(c(0, 1, 2, 3, 4, 2, 0, 1, 2, 3, 4, 2, 0, 1, 2, 3, 6, 4, 2, 0, 1, 8, 6, 4, 2, 0), byrow = TRUE, nrow = 5)

- sum(as.matrix(Table1) \* PenaltyMatrix) / nrow(ClaimsTest)

$$\text{Penalty ave} = 0.74$$

## Finding Accuracy with CART model :

- library (export)
  - library (export, plot)
    - ! all disease as input
  - ClaimTree = export (bucket 2009 ~ age + ... + bucket 2008 +  
reinforcement 2008, data = ClaimTree  
method = "class", CP = 0.0005)

prep(ClaimToee)



- Table (Claim Test & basket 2029 , Product Test)

$$\text{Accuracy} = 0.71$$

- sum (rel. matrix (Table 1) \* Penalty Matrix) / max (Chemi Test)

$$\text{Penalty per} = 0.757$$

- Claims Trend = Impact( bucket 2009 ~ ago + -----)

data = ChainTrain, method = "class"

$cp = 0.00005$ , params = list (loss = Penalty Metric)

To minimize the worst type of errors it decreases the penalty air.

## \* Regression Trees :

- Trees can also be used for Regression
- the output at each leaf of the tree is no longer a category, but a number.
- Just like classification trees, regression trees can capture non-linearities that linear regression can't.

boston - dataset

- plot (boston & LON, boston & LAT)  
points (boston & LON [boston & CHAS == 1],  
boston & LAT [boston & CHAS == 1], col = "blue",  
pch = 19)
- latlontree = rpart (MEDV ~ LAT + LON, data = boston)  
plot (latlontree)  
points boston & LON  
fittedvalues = predict (latlontree)  
points (boston & LON [fittedvalues >= 21.2],  
boston & LAT [fittedvalues >= 21.2], col = "blue",  
pch = 19)

## \* RSS :

- Residual sum of squares / Sum of square differences

$$RSS = \sum_{i=1}^n (y_i - f(x_i))^2$$

$$\sum_{\text{leaves}} (\text{RSS at each leaf}) + \lambda S$$

$$CP = \frac{\lambda}{RSS(\text{no split})}$$

## \* Text Analytics :-

- To build Text Analytics model we use Bag of words method
- Bag of words - Counts the number of times each word occurred in a string.

Cleaning up Irregularities - Apple, APPLE, @apple

Stemming - Algorithmic process of performing reduction

Eg: argue argued argues arguing

argue  
↳(stem)

- tweets = read.csv("tweets.csv", stringsAsFactors = FALSE)

- tweets & Negative = is.negative(tweets\$Avg <= -1)

↳ creates a new variable with  
TRUE if tweets Avg is  $\leq -1$   
FALSE if tweets Avg is  $> -1$

Preprocessing Text Data (Bag of Words Approach)

- install.packages("tm")

- install.packages("SnowballC")

- corpus = Corpus(VectorSource(tweets\$Tweet))

↳ A collection of documents

corpus[1] - Shows the first Tweet

stopwords("english") [1:10]

- corpus = stem-map (corpus, stopWords)
- corpus = stem-map (corpus, removePunctuation)
- corpus = stem-map (corpus, removeWords,  
c("apple", stopWords("english")))

- corpus = stem-map (corpus, stemDocument)

Extracting the word frequencies

- frequencies = DocumentTermMatrix (corpus)  $\rightarrow \text{dtm} : 3289 \text{ words}$

inspect (frequencies [0:000:1000, 500:515])

firstFreqTerms (frequencies, docfreq = 20)

$\uparrow$  to find the most popular terms  
here, the words that appear atleast 20 times

- sparse = removeSparseTerms (frequencies, 0.995)

$\uparrow$   
Sparsity Threshold

sparse  $\rightarrow \text{dtm} : 209 \text{ words}$

Converting sparse matrix to dataframe

- tweetsSparse = as.data.frame (dtm.matrix (sparse))

- colnames (tweetsSparse) = make.names (colnames  
(tweetsSparse))

$\uparrow$   
convert variable names to appropriate names

Adding the dependent variable to Dataset

- tweetsSparse & Negative = tweets & Negative

Splitting the data into Train and Test sets

- set.seed(123)

split = sample::split(tweetsSparse ~ negative,  
splitRatio = 0.7)

trainSparse = subset(tweetsSparse, split == TRUE)

testSparse = subset(tweetsSparse, split == FALSE)

Use CART to build predictive model

- tweetCART = rpart(Negative ~ ., data = trainSparse,  
method = "class")

ppr(tweetCART)

- predictCART = predict(tweetCART, newdata = testSparse,  
type method = "class")

Accuracy

- table(testSparse ~ negative, predictCART)  $\rightarrow 0.88$

Using baseline model for Accuracy

- table(testSparse ~ negative)  $\begin{array}{cc} \text{FALSE} & \text{TRUE} \\ 300 & 99 \end{array} \rightarrow 300/399 = 0.75$

Using Random forests to build Accuracy

- tweetRF = randomForest(Negative ~ ., data = trainSparse)

predictRF = predict(tweetRF, newdata = testSparse)

table(testSparse ~ negative, predictRF)  $\rightarrow 0.885$

- `graph ("cat", "dogs and cats", fixed = TRUE) → TRUE`
- `graph ("cat", "dogs and rats", fixed = TRUE) → FALSE`

Eg :

```
wikiWords2[!HTTP == ifelse(graph ("http", wikiWords,  
fixed = TRUE), 1, 0)]
```

## \* Clustering :-

- Netflix Algorithm - Cinematch
- Collaborative Filtering - making predictions for a user by collecting information from other users.

Content Filtering - making predictions ~~at~~ <sup>by</sup> number of applications comparison b/w content of the item and a user profile.

- Clustering ~ unsupervised learning methods
  - ↳ methods - Hierarchical clustering  
K-means clustering

## ⇒ Distance between points :

Euclidean Distance (ED)

$$d_{ij} = \sqrt{(x_{i1} - x_{j1})^2 + (x_{i2} - x_{j2})^2 + \dots + (x_{ik} - x_{jk})^2}$$

k = number of independent variables

## Distance between clusters :

Minimum distance - distance between the points that are closest in that cluster.

Centroid distance - distance between centroid of clusters  
↑  
avg of all data points  
in that cluster

Naturalize Data = Subtracting the mean of data  
Standard deviation

### \* Hierarchical Clustering:

Dendrogram - the display of hierarchical cluster process

Read data from text file

- movies = read.table ("movies.txt", header = FALSE,  
sep = "|", quote = "\")

Adding col names to movies

- colnames (movies) = c ("ID", "Title", "ReleaseDate", ...  
"Western")

To remove variables from data-set

- movies \$ ID = NULL
- movies \$ ReleaseDate = NULL

To remove duplicates in data-set

- movies = unique (movies)

### Hierarchical Clustering

To find distance b/w data points

- distances = dist (movies [1:20], method = "euclidean")

Cluster the data points

- clusterMovies = hclust (distances, method = "ward")

Plotting dendograms

- plot (clusterMovies)

Plotting Data points to its clusters

Dividing into

- clusterGroups = cutree (clusterCentroids, k=10) 10 clusters

To know the percentage of movies in each genre

- tapply (movies & Actions, clusterGroups, mean)

- subset (movies, title > "Men in Black (1997)")  $\rightarrow 257$   
clusterGroups [257]  $\rightarrow 2$  (Cluster 2)

- cluster 2 = subset (movies, clusterGroup == 2)

cluster 2 \$ Title [1:10]  $\rightarrow$  will show the Top 10 movies  
in cluster 2

#### \* K-means Clustering :

(1) Specify desired number of clusters

(2) Randomly assign each datapoint to a cluster

(3) Compute cluster centroids

(4) Reassign each point to the closest cluster centroid

(5) Re-compute cluster centroids.

(6) Repeat 4 & 5 until no improvements are made.

## \* Image Segmentation :-

Applications :

- Medical Imaging - Locate & Headache tumor
- Object Detection - photos , facial recognition.
- Recognition Tasks - fingerprint , iris

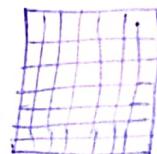
Methods :

- Clustering methods
- Edge detection
- Region-growing methods

## \* Clustering method :

Intensity vector of size  $n = 7 \times 7$

Pairwise distances  $\frac{n(n-1)}{2}$



7x7 pixels

Example :

Reading a matrix data

- flower = read.csv ("flower.csv", header = FALSE)
- flowerMatrix = as.matrix (flower)  $\leftarrow$  Converting to matrix
- flowerVector = as.vector (flowerMatrix)  $\leftarrow$  Converting to vector

Hierarchical Clustering

To find distance

- distance = dist (flowerVector, method = "euclidean")

- clusterIntensity = hcclus<sup>t</sup>.distance, method = "ward")

Plotting clusters in dendrogram

- plot (clusterIntensity)

- rect.hcclus (clusterIntensity, K=3, border = "red")

Getting clusters

- flowerClusters = cutree (clusterIntensity, K=3)

To get mean intensities of 3 clusters

- tapply (flowerVector, flowerClusters, mean)

Output the image

- dims (flowerClusters) = c(50, 50) ← Converting vector  
its matrix

- image (flowerClusters, ~~axes~~ = FALSE) ← o/p image

~~image~~ <sup>Matrix</sup> → flowerClusters, ~~axes~~ = FALSE, ~~col~~ =

- image (flowerMatrix, ~~axes~~ = FALSE, <sup>I/P image</sup>,  
col = gray (rep(c(0, 1), length = 256)))

### Example 2 :

- healthy = read.csv ("healthy.csv", header = FALSE)  
healthyMatrix = as.matrix (healthy)  
healthyVector = as.vector (healthyMatrix) ←

### Specifying number of clusters

- K = 5

set.seed(1)

k means cluster

- KMC = kmeans (healthyVector, centers = K,  
iter.max = 1000)

### Extracting the image from cluster

- healthyCluster = KMC\$cluster

[From KMC we can get the means (centers) of each cluster and  
size of each cluster  
↑  
size(KMC)]

### Output Image

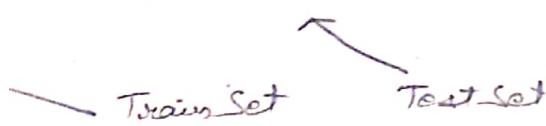
- dims (healthyClusters) = c (nrow(healthyMatrix),  
ncol(healthyMatrix))  
image (healthyClusters, axes = FALSE, col = rainbow (k))

## Detecting Tumors via HR from Example 2

- tumor4 = read.csv("tumor.csv", header = FALSE)

tumorMatrix = as.matrix(tumor4)

tumorVector = as.vector(tumorMatrix)



- install.packages("flexclust")

Conversion of data for prediction

- KMC.Kcca = as.kcca(KMC, healthyVector)

tumorCluster = predict(KMC.Kcca, newdata = tumorVector)

- dim(tumorCluster) = c(nrow(tumorMatrix), ncol(tumorMatrix))

image(tumorCluster, axes = FALSE, col = rainbow(10))

## \* Visualizations :-

Data visualization - mapping of data properties to visual properties

ggplot → Data - dataframes

Aesthetic mapping - color, shape, scale,  
x-y axis, subsets

Geometric objects - points, lines, barplots,  
bars, polygons

- `WHO = read.csv("WHO.csv")`

- `install.packages("ggplot2")`

- `scatterplot = ggplot(WHO, aes(x=GNP, y=Fertility))`

- `scatterplot + geom_point()` ← point graph

- `scatterplot + geom_line()` ← line graph

`geom_point(color = "blue", size = 3, shape = 17)`  
↑  
Triangle = 17

- `scatterplot + geom_point() + ggtitle("GNP vs FR")`

Saving the plot into pdf

- `pdf("MyPlot.pdf")`

`print(Fertility GNP plot)`

`dev.off()`

Three variable plot

- ggplot(who, aes(x = GNI, y = Fertility Rate, col = Region))  
+ geom\_point()

1 scatterplot1 =

- ggplot(who, aes(x = log(Fertility Rate), y = Under 15))  
+ geom\_point()

- model = lm(Under 15 ~ log(Fertility Rate), data = who)

scatterplot1 + stat\_smooth(method = "lm")

, level = 0.99

, se = FALSE

, color = "orange")

→ Heat Maps :

3 variables - x-axis, y-axis, shades of colors

- mut = read.csv("mut.csv", stringsAsFactors = FALSE)

Converting Date

- mut\$date = strptime(mut\$date, format = "%m/%d/%y  
%H:%M")

mut\$weekday = weekdays(mut\$date)

mut\$Hour = mut\$date %>% hour

Converting table to dataframe

- weekdayCounts = as.data.frame(table(mut\$weekday))

- ggplot(weekdayCounts, aes(x = var1, y = Freq))

↑ scatterplot2 + geom\_line(aes(group = 1))

To get the weekdays in chronological order

`days = c("Sunday",  
"Monday",  
:  
"Saturday")`

```
- scatterplot2 + xlab("Day of week")  
+ ylab("Total motor thefts")
```

## Converting Data from table to data-frames

- Day Hours Counts = as. data-frame (table (mvt & weekday, mvt & hour))

- Day  $\Delta$  Counts & Hours = abs. measure (abs. character (Day  $\Delta$  Count  
& hours))

- ggplot(DayHawkCounts, aes(x = Hours, y = Frequency)) +  
 geom\_line(aes(group = Variable))  
 , col = Variable

## Plotting Heatmap

- ~~graph~~ / Day Hairs Counts,  $\text{des}(\text{x} = \text{Hairs}, \text{y} = \text{Var},)$

+ geom - tile (and  $\text{fill} = \text{Floor}$ )

+ scale-filt-gradient (names = "Total MV Thggl", low = "red", high = "red")

+ shows (axis-title - y = element - blank ( ))

2 to get rid of our new name

## \* MAP Plotting :

- install - package ("maps")
- install - package ("ggmap")

Loading a map

- chicago = get\_map (location = "chicago", zoom = 11)  
ggmap(chicago)

Plotting on map

- ggmap(chicago) + geom\_point (data = mtc [1:100,],  
aes (x = Longitude, y = Latitude))

- LatLongCounts = as.data.frame (table (round (mtc\$longitude, 2),  
round (mtc\$latitude, 2)))  
↑  
creating a data frame with latitude & longitude  
and rounding off with 2 terms.

Converting a factor variable to numeric variable

LatLongCounts\$Long = as.numeric(as.character(LatLongCounts\$Long))

LatLongCounts\$Lat = as.numeric(as.character(LatLongCounts\$Lat))

- ggmap(chicago) + geom\_point (data = LatLongCounts,  
aes (x = Long, y = Lat, color = Fixeg,  
size = Fixeg))

+ scale\_color\_gradient (low = "yellow", high = "red")

(OR)

- ggmap(chicago) + geom\_tile (data = LatLongCounts,  
aes (x = Long, y = Lat, alpha = Fixeg),  
fill = "red")

→ HeatMap on the united states : (m - 10, a - 12)

- murders = read\_csv ("murders.csv")

Loading a map (United States)

- stateMap = map\_data ("state")

Plotting the US map

- ggplot (stateMap, aes (x = long, y = lat, group = group))

+ geom\_polygon (fill = "white", color = "black")

Making the state variable same in murders & stateMap df

- murders & regions = stateover (murders & state)

Merging the 2 df

- murderMap = merge (stateMap, murders, by = "regions")

- ggplot (murderMap, aes (x = long, y = lat, group = group),  
fill = MurderRate )

+ geom\_polygon (color = "black")

+ scale\_fill\_gradient (low = "black", high = "red",  
guide = "legend")

→ Bar graph :

- intl = read.csv("intl.csv")
- ggplot (intl, aes (x = Regions, y = Percentage of Intl))
  - + geom\_bar (stat = "identity")
  - + geom\_text (aes (label = Percentage of Intl))
- intl = transforms (intl, Region = reorder(Regions,
  - ↑  
ordering based on the  
percentage of Intl values
  - Percent of Intl))
  - ↑  
decreasing order
- intl + Percent of Intl = intl + Percent of Intl + 100
  - ↑ changing y-axis

Updated ggplot :

- ggplot (intl, aes (x = Regions, y = Percentage of Intl))
  - + geom\_bar (stat = "identity", fill = "darkblue")
  - + geom\_text (aes (label = Percentage of Intl), vjust = -0.4)
  - + ylab ("Percent of International students")
  - + theme (axis.title.x = element\_blank ()),  
axis.text.x = element\_text (angle = 45,  
vjust = 1))