

Clustering

Intra-cluster distances
are minimized

Inter cluster distances
are maximized

A grouping of data objects, such that the objects within a group are similar (or related) to one another and different from (or unrelated) to the objects in other groups

→ It is a classic case of unsupervised learning → as there is no right answer, no $y = f(x)$ situation.

If I have 1 million records and say 100 columns, it becomes very difficult to describe the data or in other words, getting coherent pieces / homogeneous groups. So purpose of clustering is to find homogeneous groups which are close to each other and can help us in describing the data.

→ It is an important step for preprocessing of data before running models
→ Also helps in identifying outliers

Outliers — objects that do not belong to any cluster or forms clusters of very small cardinality

Outlier analysis — when we are interested in discovering outliers (interesting cases) for some study, not clusters

Clustering results are used:

- As a standalone tool to get insight into data distribution
 - visualization of clusters may unveil important information
- As a preprocessing step for other algorithms
 - Efficient indexing or comparison often relies on clustering

Applications of clustering

- Image processing - cluster images based on their visual content
- document clustering - where documents are clustered based on various words, length of docs, themes, topics, to be create better user experiences, if this person reads this book is considered, to read another book also
- Ex, segmenting for efficient marketing of car based on a survey, on a scale of 1 to 10, we analyse the attitude of driver towards driving and come up with 3 groups to target
 - Speed oriented
 - Comfort oriented
 - Safety oriented

Supervised → classification - you begin with known classes
 (there is right answer) e.g. will you enroll in this program or not?

Unsupervised → clustering - you do not know any classes beforehand
 (there is no right answer) e.g. What can be the target programs based on needs
 classification is kind of individual's problem, clustering is more of group problem

Attributes (dimensions)

- Real value attributes / variables / Number where we can do all sort of mathematical functions, e.g. salary, height
- Binary attribute - Yes/No, Male/Female, True/False
- Nominal attribute - Categorical - extension of binary attribute
 - Default/Not default
 - e.g. religion (Christian, Muslim, Buddhist, Hindu etc)
- Ordinal/Ranked attributes, e.g. Military ranks (soldier, sergeant, lieutenant, captain etc)

Clustering algorithm works best when all variables are real valued.

→ To know whether 2 rows/people/object are close to each other or not, we need to measure difference b/w them. 2 important conditions are

* distance b/w 2 rows can't be -ve [Row 1 - - -]

* distance b/w 2 identical rows is 0 [Row 2 - - -
(say Row 1 and Row 1)]

The distance $d(x, y)$ between two objects x and y is a metric if

- $d(i, j) \geq 0$ (non-negativity)
- $d(i, i) = 0$ (isolation)

- $d(i, j) = d(j, i)$ (symmetry)

- $d(i, j) \leq d(i, h) + d(h, j)$ (triangular inequality)

distance is measured differently for real, boolean, categorical and ordinal - for ex, distance for Yes/No/categorical variable can not be measured like real and it becomes more like a compromise

Sometimes different weights may be associated with different variables based on applications and data semantic

for ex, if age is more significant than income, we can give it weightage

fundamental matrix for clustering

Data Matrix

attributes / dimensions

$x_{11}, x_{12}, \dots, x_{1d}$

$x_{n1}, x_{n2}, \dots, x_{nd}$

tables/objects

Distance Matrix

objects

0

$d(2, 1) = 0$

$d(3, 1) \quad d(3, 2) = 0$

0

$d(n, 1) \quad d(n, 2)$

for FA/PCA, fundamental input is correlation matrix

Dice coefficient
is also used

Distance functions for Binary vectors

e.g. to find groups (clusters) of matching pairs, series of binary questions are asked (response in form of 0 or 1)

Jaccard Similarity b/w binary vectors X and Y

$$JSim(X, Y) = \frac{X \cap Y}{X \cup Y}$$

Many a times 1 and 0 does not give equal ~~also~~ amount of information like do you like football - 1 means clear info but 0 does not give clear info. If not football then what. However, for male/female 1 and 0 gives equal information.

→ Hence 0 is not considered in Jaccard calculation

Jaccard distance b/w binary vectors X and Y

$$Jdist(X, Y) = 1 - JSim(X, Y) \quad |Q_1| Q_2 | Q_3 | Q_4 | Q_5 | Q_6$$

$$JSim = 1/6 \quad (\text{Both } 1) \quad |x| \quad |1| |0| |0| |1| |1|$$

$$Jdist = 5/6 \quad (\text{At Least 1}) \quad |y| \quad |0| |1| |1| |0| |0|$$

→ Do not consider double 0's

Distance function for real valued vector

Minkowski distance - L_p norms - generic.

$$L_p(x, y) = \left[|x_1 - y_1|^p + |x_2 - y_2|^p + \dots + |x_d - y_d|^p \right]^{1/p} = \left[\sum_{i=1}^d |x_i - y_i|^p \right]^{1/p}$$

$p = +\infty$ integer

$d = \text{no of columns/variables}$

$x_1, y_1 \& x_2, y_2$ are corresponding values of variables in row x and y .

(Root) $\sqrt{x_1}$ | column 1 | column 2 | ... | column d

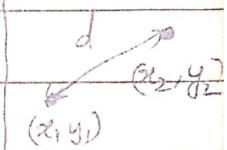
* Binary variables are not always symmetric i.e. 1 & 0 might not give always similar information.

Manhattan (or city block) distance - if $p=1$, L_1

$$L_1(x, y) = |x_1 - y_1| + |x_2 - y_2| + \dots + |x_d - y_d|$$

Euclidean distance - if $p=2$, L_2

$$L_2(x, y) = \sqrt{|x_1 - y_1|^2 + |x_2 - y_2|^2 + \dots + |x_d - y_d|^2}$$



$$\text{distance b/w 2 points} = d = \sqrt{(x_1 - x_2)^2 + (y_1 - y_2)^2}$$

Typically in clustering we use euclidean distance.

Sometimes we use manhattan distance.

* For different versions of models (to have maximum efficiency), we use different values of p and see which one works best

Weightage distance -

$$d(x, y) = \sqrt{(w_1|x_1 - y_1|^p + w_2|x_2 - y_2|^p + \dots + w_d|x_d - y_d|^p)}$$

* In general we will have mixed variables (binary as well as real). 2 ways of clustering

- Hierarchical clustering

Single linkage

Complete linkage

Ward's procedure

Advanced linkage

- K-means clustering (partitioning clustering)

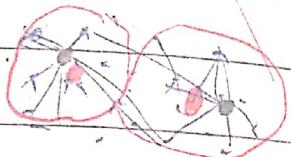
K-Means

K stands for how many clusters so if $K=3$, 3 means clustering, Given a set of X of n points in a d -dimensional space and an integer K .

Generally 2 to 7 is a healthy number of clusters.

you start with an arbitrary number and eventually discover how many clusters are good enough so that overall distance is minimized.

- say for all those points we decide 2 clusters
- Random centers are picked and distance of each point is calculated from that center (euclidean or jacard)
- Based on minimum distances, 2 clusters are formed.
- New center for each cluster is chosen based on average of all points in that cluster (Centroid)
- Repeat the above process till we don't see any improvement in variation reduction & i.e. Convergence (no further reduction in distance)
- we may get clusters of different densities as well as different sizes.
- Any outliers may force the centroid to change and bootstrap out of outliers can happen



for mixed data type we use Gower distance instead of euclidean, otherwise by default it is Euclidean.

→ Inter cluster > Intra cluster

Page

Date

→ between sum of square should be greater than within sum of squares - to make groups different

$$\text{Total SS} = \sum_{i=1}^{100} (\text{age}_i - \bar{\text{age}})^2 + (\text{inc}_i - \bar{\text{inc}})^2 + (\text{wt}_i - \bar{\text{wt}})^2 \text{ for entire sample}$$

Within SS = Same formula as above, however for specific clusters, So if 3 clusters, 3 within SS (of 3 samples)

$$\begin{aligned} \text{Between SS} &= n_1(\bar{\text{age}}_1 - \bar{\text{age}})^2 + n_2(\bar{\text{age}}_2 - \bar{\text{age}})^2 + \\ (\text{for 3 clusters of } 2 \text{ variables}) \quad &n_3(\bar{\text{age}}_3 - \bar{\text{age}})^2 + n_1(\bar{\text{wt}}_1 - \bar{\text{wt}})^2 + \\ &n_2(\bar{\text{wt}}_2 - \bar{\text{wt}})^2 + n_3(\bar{\text{wt}}_3 - \bar{\text{wt}})^2 \end{aligned}$$

Generally $\frac{\text{Btw SS}}{\text{Total SS}} > 70\%$ is a good number

Out of total variability of 100%, I am able to explain/attribute to 3 groups. So I am able to reduce the variability by 57%.

Scaling of variables — In case one variable is significantly big that it would overshadow other variables, you need to scale, normalize, ~~and~~ standardize the data

Height, wt, Income example

Scaled average is on -3 to +3 scale

→ Other normalization technique — Max-Min normalization
— also called Range norm.

→ Log transformation (though it is complicated)

Aurra has to be done on raw/unscaled dataset and not on changed/Scaled/Weighted dataset.

Page Date

→ We can give weight to a particular variable like $d_i \cdot wt \leftarrow 2 * \text{diff} d_i \cdot wt$, say you are promoting a weight loss drink and you have to emphasize on weight by giving specific multiplier weight. It is kind of opposite of scaling.

- Outliers on individual variables can be ignored.
- However, if entire cluster is outlier (say a cluster of 5 elements very different from other clusters), then we need to study it.
- Outlier analysis is not univariate typically, Has to be multivariate so that whole row becomes outlier.
- In K-means, it is possible to give your own center point/centroid instead of allowing K-means algorithm to have random start (random center). However, it is a tedious task.

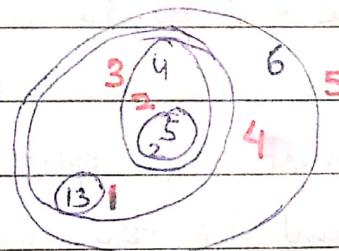
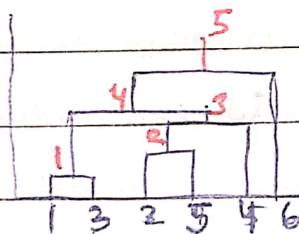
If 3 groups with 20 variables, we will have to key in centroid in the dataset by 60 times (20×3).

Hierarchical Clustering

k-means is used when we want to cluster quickly, we are not bothered about small shift in random center, also we have a huge data set to work with

But many a times you want to know how clusters are forming, at what point they are merging, how story is unfolding, so that your clarity about clustering groups is better, you do hierarchical clustering.

- Produces a set of nested clusters organized as a hierarchical tree
- Can be visualized as a dendrogram (a tree like diagram that records the sequence of merges or splits)
- Each row is treated as 1 cluster, then based on distances b/w each row, rows are merged, and at the top becoming one cluster



Unlike K means, we know here that 1,3 joined, 2,5 joined etc etc and finally reached to one cluster. So we have whole story unfolded in front of us. We can cut the tree based on logical visualization and come up with a quantitative number of clusters.

The higher the join, more distance rows have b/w them -

Strengths of hierarchical clustering

- No assumptions on the number of clusters. Any desired number of clusters can be obtained by cutting the dendrogram at proper level (we can actually see that due to strong commonality, which items are joining)
- Hierarchical clustering may correspond to meaningful taxonomies
In K-means, we start with random # of clusters with no clarity what exactly is going on it.

Two main types of hierarchical clustering

- **Agglomerative** - Bottom up - Mostly used / Popular
 - o Starts with the points as individual clusters
 - o At each step, merge the closest pair of clusters until only one cluster (or K cluster left)
- **Dissimile** - Top down
 - o starts with one, all inclusive cluster
 - o At each step, split a cluster until each cluster contains a point (or there are K clusters)

Traditional hierarchical algorithms use a similarity or distance matrix

- Merge or split one cluster at a time based on distances

Agglomerative clustering algorithm

Suppose there are 10 rows, we can find $^{10}C_2$ combination of distance b/w rows. We first join 2 rows which are minimum distance and form a temp cluster then we either join 2 more rows with min distance or join another row with temp cluster based on distance.

- Compute distance matrix b/w the input data inputs
- Let each data point be a cluster
- Repeat

Merge two closest cluster

Update the distance matrix

Until only a single cluster remains

Imp → Key operation is the computation of the distance b/w two clusters. Eg. how to find distance b/w 1, 2 or 1, 2 or 2, 4 (refer last page)

→ We have 3 ways to compute distances and work in above situation which leads to different. As distance matrix initially ~~consists~~ consists of distances b/w actual rows, we are not sure that after joining 1, 2, we should join 1, 2 first or 1, 2 or 2, 4 etc. (we don't have distances b/w groups/temp clusters). So 3 linkage method.

1) Complete linkage

Simple Single

2) Single linkage

3) Average linkage

Based on the type of linkage, shape of dendrogram changes.

example

Page Date

	1	3	4	5	2	(PD)	max of PD	min of PD	average of PD
Points taken						Complete	Single	Average	
(1,3), 2	7	(1,2) 3,2	8				8	7	7.5
(1,3), 4		(1,4) 3,4							
2, 5			2,5						
4, 5				4,5					
2, 4				2,4					

You will have to use all 3 linkages and see which one is working best for our clustering problem.

~~so~~ → for dynamic data (new data is coming), we need to recluster after a certain frequency (daily for big data or say monthly) so that centroid changes based on new data. For less amount of new data, we may try to fit it in already existing cluster. For large data, we may need to recluster.

Generally for segmenting, grouping, labelling, clustering data is static (say marketing data which won't change overnight), for dynamic data clustering may not be a very good option (daily changing info).

`dist(d2)` would calculate distances (euclidean) b/w all rows. Create distance matrix `dt`.

`hcclust(dt, method = "complete")` → `hcomplete`

`cutree(hcomplete, K = 4)` → `c1`

We can also ~~box~~ plot average of each variable (histogram) cluster wise to analyze our dataset

clustering is done on "Measures" (Numerical Data). Not good for categorical data (Dimensions). For dimensions, "Association Rules" are applied as part of Unsupervised learning.

Page Date

of PD After you cut the tree, you get the cluster # each row belongs to

R-steps

Say if d_1 is large dataset, we create d_2 with 100 rows with age, wt, income variables

→ $d_2 \leftarrow d_1[1:100,]$, check structure `str(d2)`

In hierarchical clustering, we created distance matrix

→ $\text{dist}(d2) \rightarrow dt$ - It will create pairwise euclidean matrix and store in dt (fundamental matrix)

→ $\text{hclust}(dt, \text{method} = \text{"complete"}) \rightarrow hcom$ or `single/average`, it creates a cluster/whole dendrogram using selected method

→ $\text{plot}(hcom)$, $\text{plot}(hcom, hang = -1)$

Visually visually we can see that there are 4 groups. 1,2 and 3,4 join at a very later stage so clearly they are separate groups.

→ For ^{average} $\text{hclust}(dt, \text{method} = \text{"average"}) \rightarrow havg$, $\text{plot}(havg)$

→ for single $\text{hclust}(dt, \text{method} = \text{"single"}) \rightarrow hsin$, $\text{plot}(hsin)$

After analyzing above 3 plots, we can see that complete is giving us 4 clear groups, so we would cut the tree into 4 groups.

→ $\text{cutree}(hcom, K=4) \rightarrow cl$, cl

cl would tell us which row is in which group.

So cluster is reduced to a vector with label names.

→ $d2\$cluster \leftarrow cl$, here you are adding 'cl' vector view(cl) to $d2$ as column cluster,

we can now do ANOVA or summary or boxplot etc to analyze various population parameters "**cluster-wise**"

In dendrogram, you can see which rows joins first so their distance is minimum or they have strong commonality. So here you are able to create complete taxonomy in mind. In K-means you do not know about whole story but just the groups as an output, for the rows which are joining at large distance, you need to study individual variables to determine due to which variable distance is large.

Silhouette coefficient - Once you are done with group membership and you want to verify whether the membership is correct or if there could be some improvement, you calculate silhouette coefficient.

- suppose we consider row i .
- first we find the average distance of row i from all members in its group (say row i belongs to group 2, find distance of row i from all members of group 2, say it is a_i)
- Then find the average distance of row i from other groups/clusters, say 3 more clusters of $b_i = b_1 + b_2 + b_3 / 3$
- Now to conclude that row i is in correct group

$$a_i < b_i$$

distance within own grp should be less than avg distance of other groups.

$$\rightarrow \text{Silhouette coefficient} = \frac{b_i - a_i}{\max(b_i, a_i)}$$

max distance of b_i and a_i

Silhouette coefficient varies b/w -1 to 1

- It is calculated for each row. So, for rows where silo is +ve, it is good, as for these rows $a_i < b_i$.
- However, if silo is -ve, it indicates that $a_i > b_i$ so that row is closer to some other group and its better to recluster.
- In all clustering problems, there would be some -ve silos. Which is OK. However, if more than we need to revisit clusters or need to increase ^{or decrease} number of clusters.

In R

- silhouette (cl, dt) → sl (give input cluster vector and distance matrix)
 - plot(sl)
 - sl - this table shows row's current cluster, its neighbours and siloh coeff
- Ideally if more than 10% rows have -ve siloh, we need to revisit. Lot of subjectivity in clustering. Depends on interpretation
- If there are clusters of very small size, that means those rows are outliers and very different from others.
 - So far we have worked only on 'Real-Valued' variables. Now we work on mixed data (Multiple data types)

Mixed data distance computation - Gower distance

- As we work on categorical, it is a compromise, say distance b/w red and blue car
- Idea is to use distance measure between 0 and 1 for each variable, i.e. standardize it
- Aggregate = average distance over variables
- Binary, nominal - Calculate Jaccard distance

→ Interval-scaled - $\frac{x_{if} - \min(x_{if})}{R_f}$ - Called as Range normalization

x_{if} = value for object i in variable f

R_f = Range of variable f for all objects

say, age = 215

, ~~min range among~~

min age among all ages = 35

max age among all ages = 75

$$\text{After normalization, } \text{age}_i = \frac{45-35}{75-35} = 0.25$$

→ Ordinal - Use normalized ranks, then like interval-scaled based on range, normalized rank

$$\text{Normalized rank} = z_{if} = \frac{r_{if} - 1}{M_f - 1}$$

Say for a grocery, allowed ranks are 1, 2, 3, 4, 5

However no one gave 5, but there are some 4.
So 5 is not present in your data set. Still we would consider 5 as max rank.

So maximum possible rank is 5 and not 6.

If rank given by a person is 3, then normalized

$$\text{rank} = \frac{3-1}{5-1} = 0.5$$

- Once every variable is brought to 0,1 scale, compute Euclidean distance,

In R - Gower Matrix

- Remove cl column from original d2

$$d2 \leftarrow d2[, -4]$$

- adding a column satisfaction score on a scale of 1 to 5

$$d2\$ss \leftarrow sample(1:5, 100, \text{repl} = \text{T}), \text{str}(d2)$$

- Specify right data type for ss as currently showing integer

$$d2\$ss \leftarrow \text{as. ordered}(d2\$ss)$$

- Use this data to calculate Gower matrix.

(library(Cluster))

daisy(d2, metric = "gower") → dt1

It would compute 100×100 matrix, normalize all variables and then calculate euclidean distance.

- **Imp** - If using gower, you need not scale any particular variable, in case one variable is significantly higher, As normalization is inbuilt process in Gower, by normalizing everything to 0-1,

→ hclust(dt1, method = "complete") → hcsm

→ test all is similar as previous.

→ Fuzzy is a technique of classification where membership is not strict but everything is a probability

Page Date

Cluster Sampling - a method / kind of clustering where you are creating a cluster which has all characteristics of population.

Decision Tree

Two kind of problems, one where we already know classifications (y is categorical) and we try to classify various records into classification tree

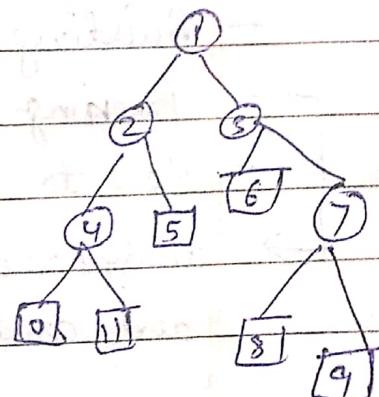
Another is when y is not categorical, it is called decision tree \rightarrow regression tree

Very popular technique (No. 1)

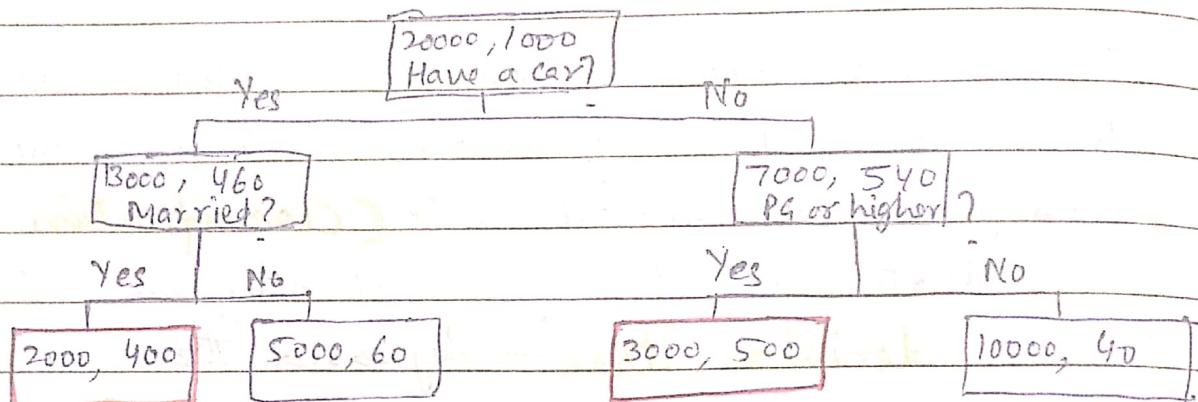
\rightarrow if we are asked to guess a number b/w 1 to 100,
one way is to fall out each number 1 by 1 and reach to the result.

Another is to ask, if more than 50 (No), if more than 25 (Yes), more than 37 — so we should reach the number in 7 hits; DT works on same principle

- DT is an inductive learning task. — it works on multiple iterations and learns various rules on its own, which it uses to make more generalized conclusions.
- A predictive modelling based on a series of Boolean test for branching, for ex number > 50 , so every time left/right node only
 - ① - root node
 - - terminal nodes which do not have further branching (Leaf node)
 - - Non-terminal nodes, have further branches,



2 primary objectives of any model - explanation, Prediction



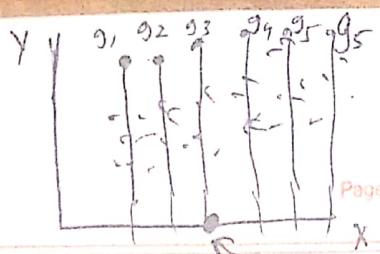
If we are having an exhibition of cars and we invite 20,000 people and only 1000 people turn up for it. So response rate is only 5%

However, ① if we create a decision tree, we found out that out of 2000 people who are having a car and are married, 400 were interested in exhibition.

② Also out of 3000, who do not have a car but have PG/higher degree, 500 turned up for exhibition. So we can focus on these 2 groups for higher response rate instead of complete population.

3 steps of tree construction

- Building the tree
 - Pruning - cross validation
 - How to do it in R
- analyse the tree, accuracy, error, rebuild, again validate and come up with final 'optimal' version.



split partition the data (say x b/w 0 to 100) into multiple partitions (say 10) and then calculate Gini Index $g_1/g_2 \dots g_{10}$. whichever partition gives minimum 'Gini', that x value is picked in tree.

Building the tree -

- build algorithm to choose the best split at every level
- To know what is best, best split is something which reduces the impurity level of the node the most
- Methods used to determine impurity
 - Gini
 - Entropy
 - Classification

To find out impurity

Case is to find which color ball would come if I pick from different bags.

bag 1 [5R5G] - high impurity - unambiguous - not sure which ball would come if I pick one

bag 2 [8R2G] - low impurity - high chances that red would be the outcome

Pure node [10R0G] - 0 impurity - 100% chance of red outcome

For classification - idea is to reduce this impurity at each node

For classification tree, ultimate goal is to find out those boolean rules (Married, Yes/No) which reduces the impurity of each node.

So there are different measures of impurity which ultimately gives you same answer, as they all are based on probability.

Gini Measure

It is the probability of randomly selected elements belonging to the same class

$$\text{Sum}(p(i) \wedge 2) \text{ for all class } i.$$

Higher the gini score, its better.

2 Assumption

- Sample size is large

eg- What is the probability of getting 2 balls of same colour from bag 1

$$\begin{aligned} RR + GC &= \frac{5}{10} \times \frac{5}{10} + \frac{5}{10} \times \frac{5}{10} && \text{(without replacement, probability of 2nd ball would have been } \frac{4}{9}) \\ &= \frac{1}{2} \end{aligned}$$

but because assumption is that the sample size is large, we will consider $\frac{1}{2}$ instead of $\frac{4}{9}$

$$\text{So impurity } = 1 - \text{Purity} = 1 - \frac{1}{2} = \frac{1}{2} = .5$$

$$\text{For Bag 2} = 1 - \left(\frac{8 \times 8}{10 \times 10} + \frac{2 \times 2}{10 \times 10} \right) = \frac{32}{100} = .32$$

So impurity for bag 1 is higher than Bag 2

$$p_1 = P(Q)$$

$$P(i) \wedge 2 = P(Q)P(R) + P(Q)P(R)$$

$$p_2 = P(R)$$

so the formula becomes more complex due to the fact that there are two different probabilities, one for Q and one for R.

Gini is the measure of purity, so higher is better
 $1 - \text{Gini}$ is impurity so lesser is better.

To consider which variable to pick for splitting

100	20
R	G

$$X_1 < 10$$

100	20
R	G

$$X_2 > -7$$

50	R	50	G
----	---	----	---

50	R	15	G
----	---	----	---

$$a = \left(\frac{5}{55}\right)^2 + \left(\frac{50}{55}\right)^2$$

$$b = \left(\frac{15}{65}\right)^2 + \left(\frac{50}{65}\right)^2$$

$$c = \left(\frac{5}{95}\right)^2 + \left(\frac{90}{95}\right)^2$$

$$d = \left(\frac{15}{25}\right)^2 + \left(\frac{10}{25}\right)^2$$

90	R	5	G
----	---	---	---

10	R	15	G
----	---	----	---

a is level of purity

To compare 2 different options, we need to find out the weighted average at each node.

$$\text{pure1} = a \cdot \frac{55}{120} + b \cdot \frac{65}{120}$$

- weighed measure of purity for this split. (no. of observations that land in the node gives weight)

$$\text{pure2} = c \cdot \frac{95}{120} + d \cdot \frac{25}{120}$$

purity should be higher or $1 - \text{purity}$ should be lower. So split 2 is better

for further splits, these a,b,c,d childs would become parent and again we would compute purity (im)

Entropy

Another measures, to find out best split. All gives same results as all are based on probability.

Gini

- Measures reduction in entropy
- Entropy is $\sum [p(i) * \log_2 p(i)]$ where $p(i)$ is probability in class i.
- Information gain

$$\text{Entropy (Parent node)} - \text{Entropy (child nodes)}$$

So instead of $p(i) * p(i)$ we use $p(i) * \log_2 p(i)$

Misclassification

error rate kind of proportion wrongly classified = $\frac{5}{11}$

$$C_1 = 5$$

$$C_2 = 6$$

Say for box 1 [50R 5G] because of the majority, we are classifying this box as red.

However error we are making is $\frac{5}{11}$ as 5 balls are green

So for a, b, c, d nodes calculate misclassification.

Take weighted average

Wherever misclassification is less, choose that split.

Use package "rpart" (Recursive partitioning)
rpart, plot to visualize the tree

Page Date

If too many outliers in the dataset, go for decision tree instead of Regression as it asks only a certain number, does not ask outliers.

If too many variables, use decision tree first to identify important variables and then go for logistic regression.

In DT, you directly talk about variables interaction in Logit you do not directly talk about interaction if you have a car and married as this is a classification tree

`rpart(y ~ age + gender, data = d, method = "class")`

→ dt

`rpart.plot(dt)`

`rnorm(1000, 45, 4) = age`

It creates 1000 records with mean 45 and stdv = 4

`table(d$y)` y is 0 and 1

0 1

523 477

so 0 (blue node) = 523

1 (green node) = 477 = $477/1000 = 47.7\%$

all 100% records

Root node	0 node
	• 48 100 %
	age < 51

blue node (1000) are here as this node is classified as blue node, error is 48% (1 record)

misclassification is 477

blue node

0	0.46 (92%)
---	------------

92% records (920) are here
• 46 is proportion of 1 in this node

Green node

1	0.54 (8%)
---	-----------

had to be 0
we classified them in 1 node

8% records (80) are here
• 66 is proportion of 1 in this so 34 is 0, so 34% is misclassification

So $\text{age} < 51$ is one variable which is helping in the split by reducing the impurity to the maximum.

`dt` would give you detailed picture of the tree

* marked are terminal nodes

`xpart, plot(dt, type = 1, extra = 2)` instead of % provides absolute numbers

By default it is doing Gini measure.

For others you need to specify.

So I can now pass a new record into this tree to do the classification prediction (which node that record would fall into)

`predict(dt, d) → pred`

will give 0 and 1 probability for each record

`View(pred)`

Here it would directly

`predict(dt, d, type = "class") → pred` give the classification

`View(pred)`

each record would fall into.

you can use
test data here with `dt = final_dt`

Confusion Matrix

`table(actual = d$y, predicted = pred)`

To compare how good I predicted as compared to actuals. Comparing of y_i and \hat{y}_i

Any model is developed in training data and after complete satisfaction, it is tested on new/unknown training data which model has never seen and then the accuracy is determined.

Precision - $\frac{TP}{TP+FP}$ - What proportion of the identifications was actually correct

Recall - $\frac{TP}{TP+FN}$ - What proportion of actual +ves was predicted identified correctly

	0	1	
0	408	115	← False positive (was 0, you predicted 1)
1	287	190	
	1		(False Positive Rate - FPR) = $115/(408+115)$ $= 190/(287+190)$

False -ve (was 1, you predicted 0) TPR (True Positive rate)

$$\text{misclassification Rate} = 115 + 287 = 402 = 40.2\%$$

We started with 47.7% misclassification

from 523 initial 0's we are now classifying 408 only as 0

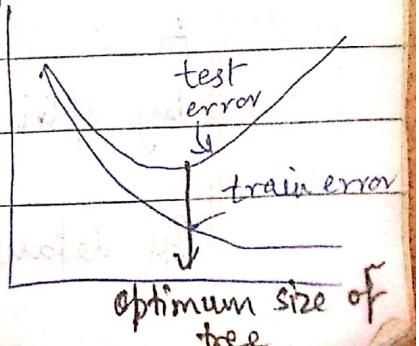
Another algorithm to use for split is chi-sq test using expected and observed frequency for gender, you calculate χ^2 (gender), same for other variables whichever is lower, that variable you choose.

However, these days Gini and Entropy only are used.

So for any regression problem, if we add 1 more variable to already used 5 variables, R^2 would increase (1 more variable would add to explanation) but adjusted R^2 may decrease.

Similarly for DT, if we keep on splitting, accuracy on training data would increase. But accuracy may decrease for an unforeseen dataset of test. ~~so~~ This is called problem of overfitting

So we need to find out optimum size of tree and stop there.



Bigger trees means more complex tree so complexity parameter (like adjusted R^2) is an important parameter for Pruning.

with every increase in node size, what is the reduction in misclassification error. - complexity parameter is an interpretation of that - CP value.

Say 10 nodes - error = .09

11 nodes - error = .085

So reduction in error = .005

Every split leads to a CP. So every leaf will have a CP value.

printcp(dt)

CP	nsplit	no of splits	misclassification error	xerror	sestd
		↓	↓		
0.0545	0		1.0000	1.000	0.0331
0.0167	1		0.9454	1.014	0.0331
0.0157	4		0.8951	1.050	0.0331
0.0104	6		0.8637	0.9937	0.0331
0.0100	8		0.8427	0.9874	0.0330

rel_error = root node error (.477) is considered 1 (baseline)
 0.9454 says misclassification error reduced from 1 to 0.9454 i.e. from .477 to $.477 \times 0.9454 = 451$
 by splitting the tree once.

By default, at rel error = 0.01, tree would stop growing.

By default it is done. Can do 5 for CV also. However 10 fold is most accepted

10 fold cross validation → out of 1000, if tree is performing on 900 records only randomly leaving 100 records each time (which inbuilt it is considering as test dataset), and to cross validate the tree performance, it calculates multiple misclassification errors, mean of these errors is xerror and std dev is xstd.

It is kind of inbuilt process of software to fine tune the tree.

So this error would reduce to a certain point and then again start increasing. So we need to stop building the tree here.

→ To ensure we build big tree

`rpart(y ~ age + gender, data = dt, method = "class", control = rpart.control(cp = .001)) → dt`

it would create much bigger tree as $cp = .001$ so we are allowing even small improvements. By default CP value is 0.01.

(train error)

Here we see that rel-error continuously decreasing till 69 splits, but xerror (test error) started increasing after decrease at 1st split only,

(ie)

So we should stop splitting the tree after 1st level.

std dev rule - Pick that CP value where diff in xerror from higher split to lower split is within $\frac{1}{3}$ std dev i.e. $xerror + std dev$ contains 1 level above xerror.

So pick 1 level above CP instead of the row where CP is lowest

plot cp(dt)

Now you know what level you need to cut the tree based on CP value, cut/prune the tree.

`prune(dt, cp = 0.015) → final dt`

`rpart rpart.plot(finaldt)`

→ Some strategies to make tree more understandable

- * In this example we have used 2 variables, age (continuous) and gender (categorical with only 2 values)
- We can have categorical variables with more levels or
- We can even convert continuous in categorical (say 6-7 ranges of age) for better interpretation

- * DT does not bother about outliers

- * for ex, we added occupation with 4 levels 'academic' 'lawyer', 'doctor', 'entrepreneur' in the existing example.

`rpart(y ~ ., data = d, method = "class") → dt`

`rpart.plot(dt)`

In this decision tree, only age and gender based rules are showing up. Not the occupation. So DT is telling you which variable is important or not.

Now in subsequent logistic Regression you can leave occupation and only important variable age and gender are considered.

→

Here for node 7, it uses condition

occupation != "l", "d", "e", Other condition is

occupation == "a"

So it has again taken the rules as boolean Only

$y \sim \cdot - \text{age} = \text{all except age}$
 $y \sim \cdot = \text{all variables}$

Page Date

Conditions in which we use DT -

- 1) Have very large rows and columns
- 2) Have very large columns (Count), you can use DT to find important ones to be used in subsequent models
- 3) Outliers in the data
- 4) Tells about variable interactions easily

3 performance measures. Higher variation in a variable is good for explanation.

Lesser or no variation does not give much information.

↳ ①

ROC Curve - To determine good cutoff probability (instead of 0.5)

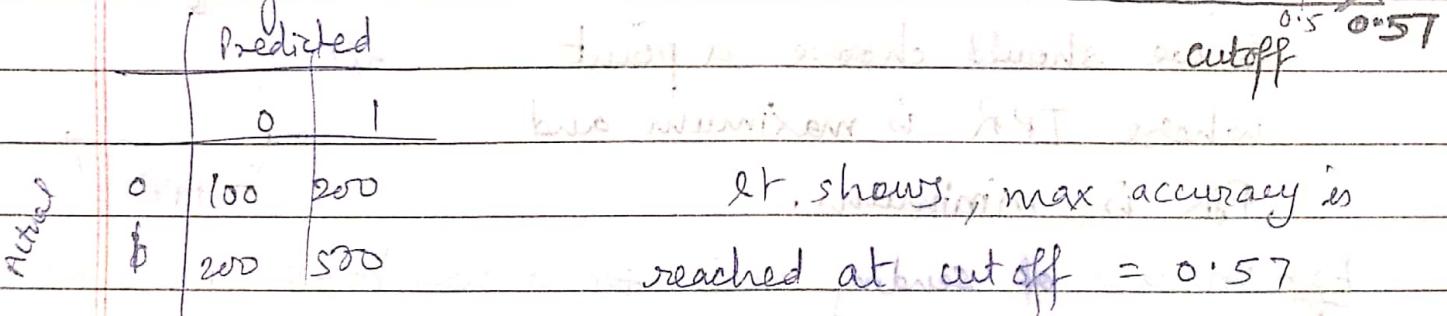
plots probability vs actual 0/1

library (ROCR) takes test prob 1 for each record, instead of d, use test data for fine tuning

```
predObj <- prediction(d$y, pred[,2], d$y)
perf <- performance(predObj, "acc")
plot(perf)
```

0.5 may not necessarily be a right number

So need to find right cutoff to label something as 0 or 1



→ If I need to give loan, I can not take probability = 0.5 everytime for taking the decision.

→ Instead I will double check to find a healthy cutoff (Next level decision)

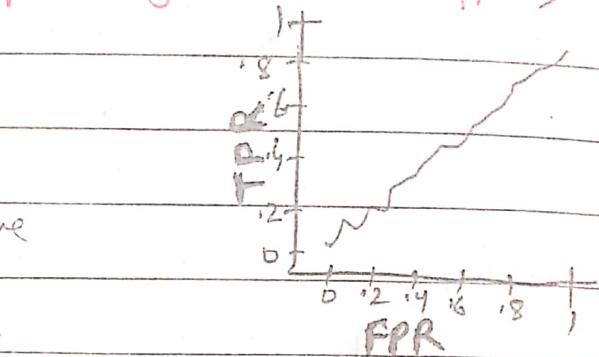
→ predObj has probabilities of all 1000 rows with prediction=1
pred[,2] → we are considering prediction = 1 only.
and d\$y which has actual 0,1 values for y

② Make a Trade off b/w FPR and TPR

`perf1 ← performance(predobj, "tpr", "fpr")
plot(perf1)`

	0	1	
0	a	b	→ false negative
1	c	d	

False -ve



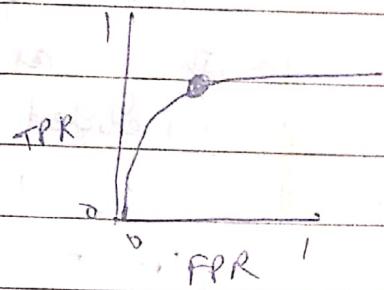
$$\text{accuracy} = (a+d)/N$$

$$\text{TPR} = \frac{d}{c+d}$$

$$\text{FPR} = \frac{b}{a+b}$$

~~Break point~~ High TPR is good as we are predicting 1 as 1
FPR should be less as something which was 0, we predicted as 1

generally, the curve would be like
so we should choose a point
where TPR is maximum and
FPR is minimum.



Eg - In case of sending mail for a new book, False + ve is not bad, because for those ppl (actual 0) who did not buy the book, we are predicting that they will buy the book. So we would send few extra mail.
• But False -ve is not very good. In above example, who actually bought the book (actual 1), we predicted

that they will not buy, so we would miss/lose those customers.

TPR and FPR changes (confusion matrix) when you change cut off, instead of '5, you use '1 or '01.

`plot(perf1, print.cutoffs.at = seq(0, 1, by = '01))`

`plot(perf1, print.cutoffs.at = seq(0, 1, by = '-1))`

③ **Lift of the model** - how much I am gaining from the model. Is the model giving me ROI

Say you have ~~pass~~ responses of 1000 customers buying a product (0/1)

Prior probability is out of 100, 20 would buy, so $P(1) = .2$

$y(1000)$	$b(1000)$	
o] 20 divide into	.27] 35	Help you determine
o] 20 10:10%	.86] 29	probability and
G] 20 deciles	.81] 27	sort them in
o] 20 for each	.74] 21	descending
o] 20 decile $P(1)=20$.71] 17	Now 1st decile
1] 20 customers buying	.17] 1	has say 35
1] 20	.02]	
1] 20	.01]	

Total 200 customers
without using Model (WUM)

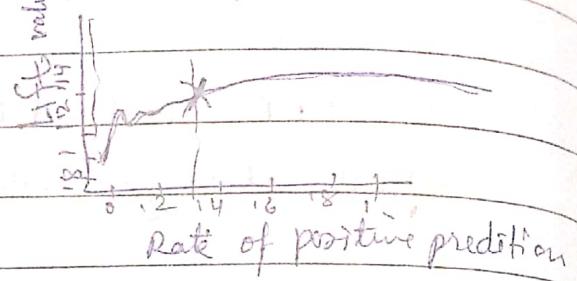
After using Model (AVM)

Now, WUM, if I was spending x amount on to get 200 cust here AVM, I am getting 35 cust by spending same x amount or you can say WUM, I got 200 cust spending y amount.

However AVM, I got those 200 cust way easier by spending less amount. So we have achieved the lift.

`perf2 <- performance(predobj, "lift", "opp")`
`plot(perf2)`

- It says beyond certain point (30% opp) we are not gaining in lift value



Working on real data set G.csv

3 important questions to answer -

- What are the important variables which need to be considered
- What would be the healthy cutoff to decide who is 0 or 1
- What is model's lift. How model is really performing on test dataset (ROI?)

Intermediate question →

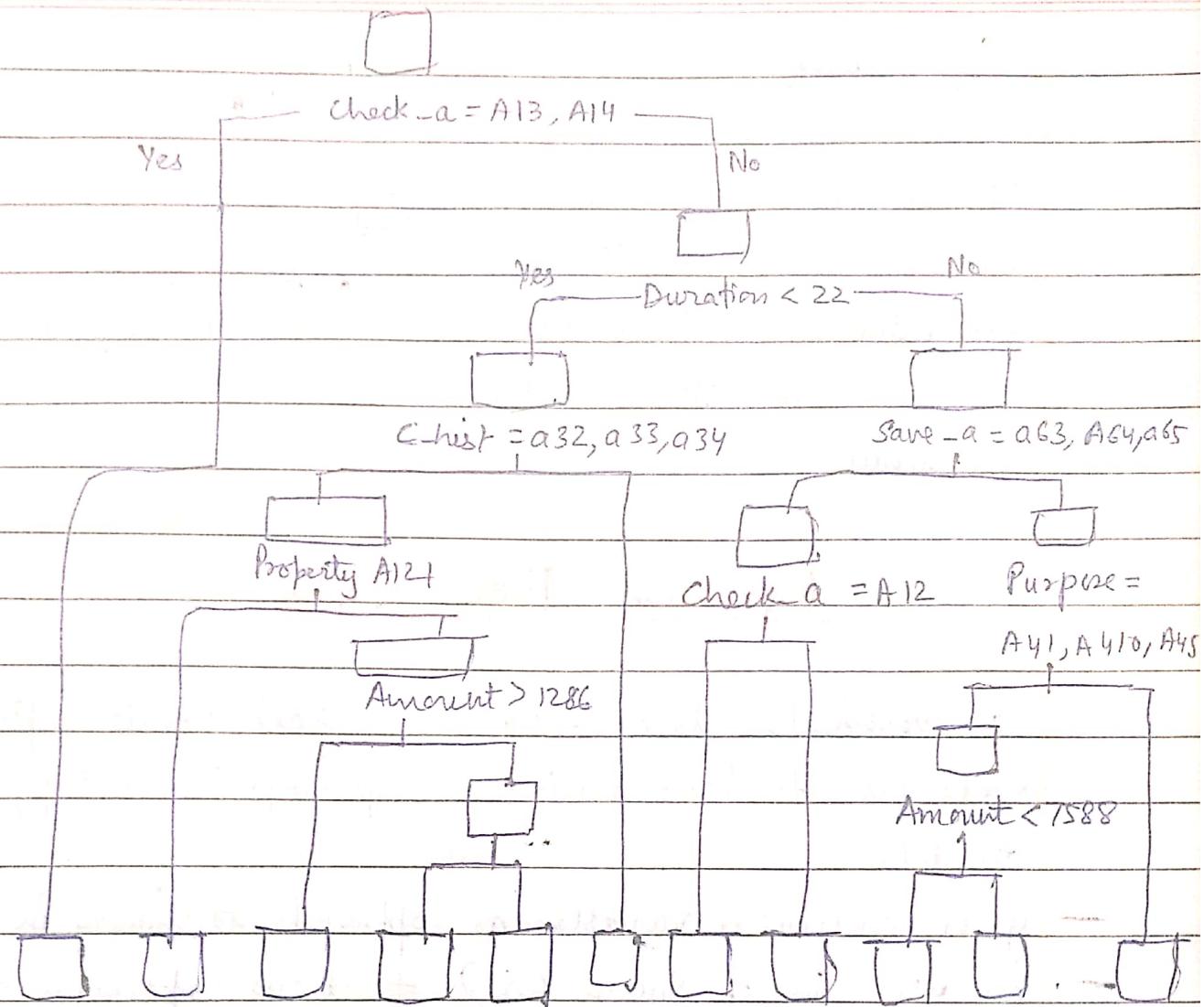
What should be DT size? Do you need to prune it?

At what level you need to pursue?

CP	nsplit	rel - error	error	scaled
0.0822	0	1.000	—	0.05577
0.0400	2	0.8355	0.8355	0.0527
0.0355	3	0.7955	0.8355	0.0527
0.0222	4	0.7600	0.8355	0.0527
0.0133	7	0.6933	0.8222	0.0524
0.0088	10	0.6311	0.7866	0.0516
0.0077	12	0.6044	0.8267	0.0525
0.0059	15	0.5777	0.8577	0.0532

- * Decision Tree serves as an important, quick variable selection tool.
- * tell lot about interaction

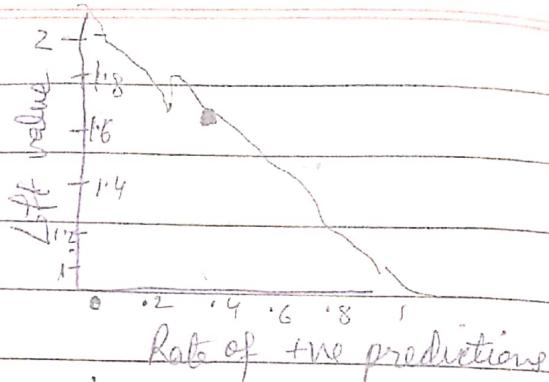
Page _____ Date _____



- (1) This tree shows that, out of 20 variables in the file, check_a, Duration, c_hist, save_a, property, amount check and purpose are important variables. If we have to perform linear reg, we can perform on these variables only. ALSO we can come up with fewer dimensions.
- (2) This tree also shows that there is interaction b/w duration and c_hist or c_hist and property. There is no interaction b/w duration and amount c_hist.
- (3) If duration and c_hist is coming many times in the tree, we can say that they are strongly interacting

Lift curve

It shows with 35 TRR,
we are achieving 1.6 of
lift value i.e. 1.6 times
more value than without using
the model.

Regression Tree (RT)

In classification tree, when our goal while splitting each node was to have minimum impurity (Gini, Entropy - etc.)
In RT

- y is continuous variable as opposed to binary in CT.
- We use Sum of Squares (SS) to have optimum split.
- $Lss + Rss < Tss$, these variables are used in splitting, which makes this equation true (Reduces SS)
where $Tss = \sum (y_i - \bar{y})^2$

$$Lss = \sum (y_{Li} - \bar{y}_L)^2 \quad \text{Yes} \quad | \quad x_i < 10 \quad \text{No}$$

$$Rss = \sum (y_{Ri} - \bar{y}_R)^2 \quad | \quad \bar{y}_{\text{left}}(1-m) \quad | \quad \bar{y}_{\text{right}}(1-n) \quad | \quad Lss \quad | \quad Rss$$

- Not doing any estimation here, we are just calculating ss only.

Most of determining techniques do not have any assumptions.

Page Date

- Linear Reg vs Regression tree - how a linear reg situation can be outsmarted by Regression tree (case)

- R Code

Picked up y variable having 10 values each in ranges

15-25, 0-5, 30-35 and 15-17

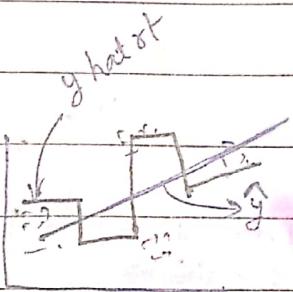
x is variable from 1 to 40

fit lm on it, $\text{lm}(y \sim x) \rightarrow \text{ft}$

predict (ft) $\rightarrow \hat{y}$

lines (x, yhat, col=1)

NOT getting good output from lm.



$\text{rpart}(y \sim x, \text{method} = \text{"anova"}) \rightarrow \text{rt}$

predict (rt) $\rightarrow \hat{y}_{\text{rt}}$

lines (x, yhatrt, col=2)

correctly predicting y, so it says y var is not a linear relationship

^{imb} \rightarrow In regression, you are getting very poor r^2 , it tells that there is no strong linear relationship, So try RT there

To check the extent of agreement (accuracy) in both cases.

$$- \text{LM } \sum (y - \hat{y})^2 = 4401.157$$

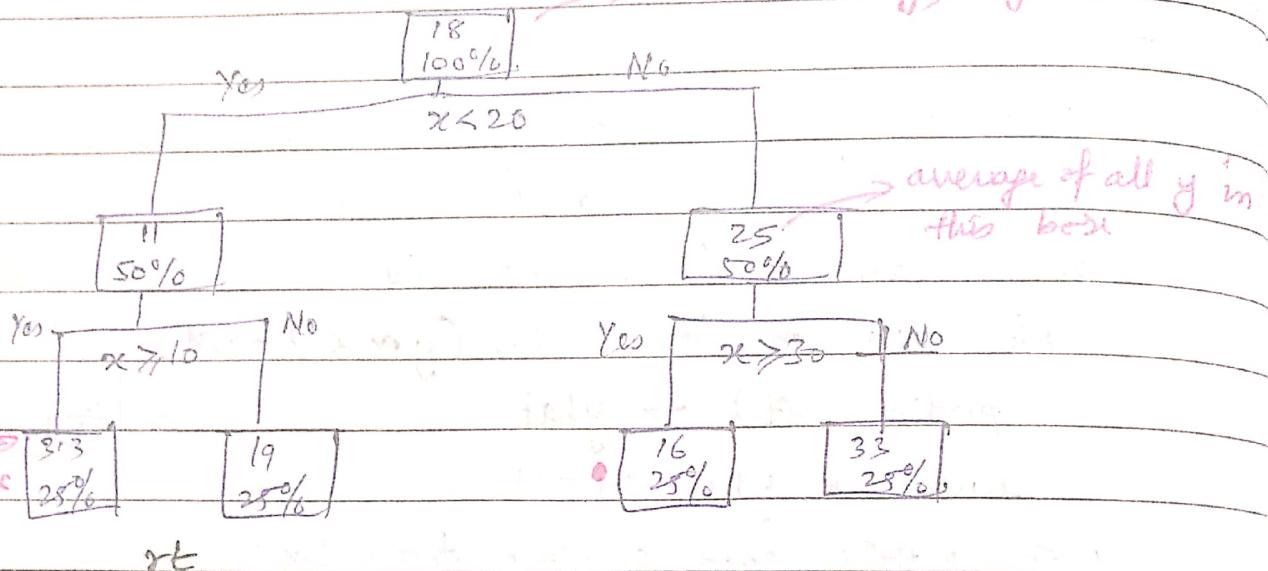
$$- \text{RT } \sum (y - \hat{y}_{\text{rt}})^2 = 92.62$$

LM showing lot of disagreement of \hat{y} values from y as compared to \hat{y}_{rt} /y.

Summary (ft)

Shows $r^2 = 0.0424$ (very poor)

rpart, plot (rt)

 $18 \text{ is } \text{mean}(y) - \bar{y}$ 

rt

which ($x < 20$) \rightarrow idmean($y[\text{id}]$) \rightarrow gives ~~rel error~~ 11print cp(rt) prints complexity parameter (cp)
($= \text{error} + \text{stderror}$)

rel_error of next row = rel_error of previous row - cp

$$\text{Root node error} = \frac{4596.4}{40} = 114.91$$

$$= \text{sum}((y - \text{mean}(y))^2)$$

$$= 4596.4$$

So objective is to reduce this SS at each node
 so that node error is minimum.

predict(rt)

- gives that average y mentioned in each box for corresponding y . So predicted y for $x = 37$ is 16 as x is < 20 and $x > 30$,

Imp -

RT is very helpful in finding missing values

If there are many missing values in data, it's good to have some estimation (mean y in that box) for that value to go further.

- RT does not worry about Outliers. If there is an outlier in y variable, it may change the structure of RT. However outlier in x does not impact much.

Imp -

Influence point - prof added outliers in yx values. Also added z and gender. Then ran lm, which gives high ($.59$) r^2 . It is not that the relationship has become linear, but this one observation is pulling the line from non-linear to linear. So need to be careful about such situation.

Imp -

Cooks distance (ft)

Threshold value of Cooks distance = 1

$$\frac{\text{No of observations} - \text{No of variables}}{\text{No of observations} - 1} = \frac{1}{39} \quad \text{in this case } (40-1) \text{ without extra var}$$

for 41st row, cooks distance = ~~1.23~~ 1.23

for other rows, $\leq \frac{1}{39}$

So we can say 41st row is an influence. Cooks distance is the measure of how each value is influencing the linear line. CD is

fundamental in finding outliers in lm

Random Forest (RF)

Eg - To find out who will win KA election would not settle down on input of 1 person, it instead would ask many, and based on their answers, would find out real answer. Not dependent on 1 predictor, instead dependent on multiple predictors (also called democracy) for an estimate.

mainly 2 Ensemble (collection) Techniques

Extremely popular because of accuracy, not because of explanatory power.

(1) Bagging (2) Boosting

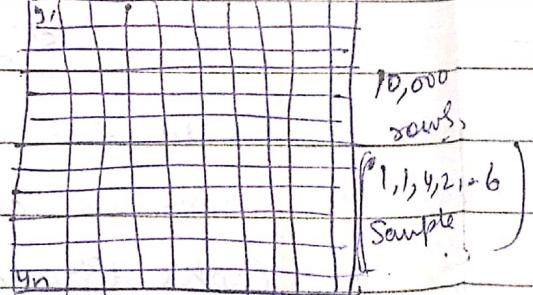
Random Forest	Gradient boosting
(Fits) n trees, $\hat{y}_n = \bar{y}$	(Fits) n trees, $\hat{y}_n = \sum_{i=1}^n \hat{y}_{ni}$

Base of both bagging and boosting is trees (DT)

- Ex. dataset with $x_1 \sim x_n$ var and y var, we create RF of 200 trees

Random Forest (RF) randomizes on following 2 things:

- Random Sample (Sample with replacement)
- Random Columns. \hookrightarrow Bootstrapping.



So if I am creating RF of 200 trees, in each

Sample Some rows would be duplicated and some would be left out (as samples are created with replacement)

These which are left out are called Out of bag (OOB) samples.

Second is we do not use all columns in each sample. In KA ex. input from a family of 4 ppl. would be same (with same mindset). so no point in taking response of complete family.

If we create identical trees with all columns, it would be like same trees. So objective is to create different trees with different no of columns (thumbrule is $\frac{1}{3}$ or \sqrt{n} (no. of columns)) to be considered in a particular tree for splits

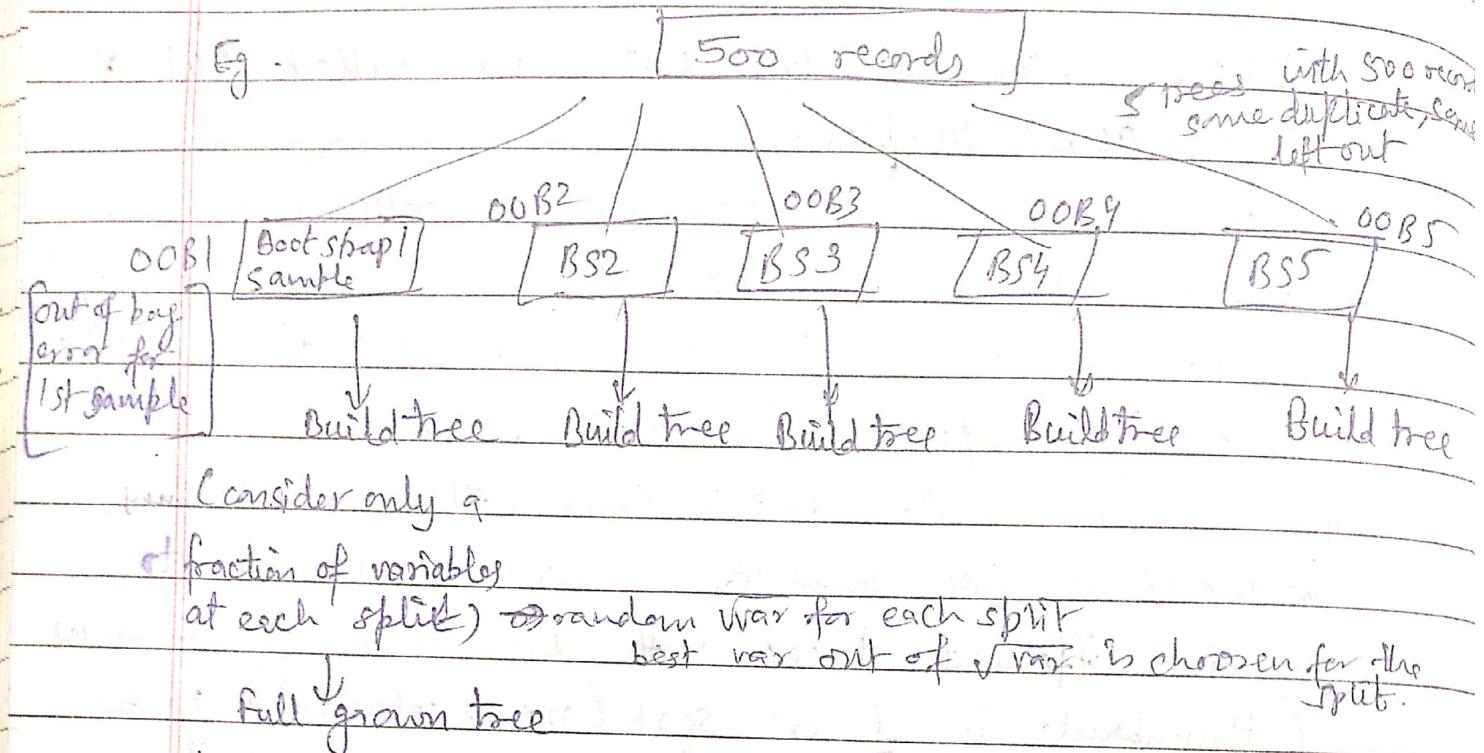
→ Creating Random Samples (10,000 records from the population serves as sample again and again with replacement) is called Bootstrapping/Bootstrap sample.

→ Here we enforce diversity.

→ If I create samples without replacement all 1000 would come each time so all trees would be identical

→ In DT we use all vars of the dataset for splitting and use that var finally for split which reduces the impurity the most. In RF, we use subset of columns each time to create a fully grown tree. So at each node split, \sqrt{n} vars are used to determine impurity. For each node set of \sqrt{n} vars can be different.

Ques → In RF, we create Fully grown trees i.e. no CP value to be used, as it would have cancelling effect.



Now for prediction, classification, we pass a new record to each tree and then take the majority decision. For above ex, if 3 trees predict 1 and 2 trees predict 0, the the outcome should be 1 for prediction, Regression, again we pass the new record to each tree and then take the average / aggregated decision. For ex, we get 5 different numbers for \hat{y}_{new} using 5 trees, now we take the average of all to get final estimate of \hat{y}_{new} .

Result is the ensemble (collection) of outputs in RF. RF can be classification or Regression type based on y variable (if 2 classes then 2 trees and so on for more than 2 classes).

Ques

Rf makes sense when you have large number of rows and columns. For less rows/columns, you will not get enough variation to create a forest.
eg. of classification forest using credit data

library (randomForest)

\times randomforest (CREDIT ~., data = d, ntree = 200, mtry = 5) \rightarrow rf

mtry = 5 means 5 variables to be used
randomforest (as.factor(CREDIT) ~., data = d, ntree = 200, mtry = 5) \rightarrow rf

rf

Output interpretation

OOB estimate of error = 23.2%

Confusion matrix

		0 : 1	class.error
0	634	66	0.0942
1	166	134	0.5533

OOB estimate means those ~~out of bag~~ records which were not the part of the sample are passed through their respective trees, there would be some misclassification. So the count of this misclassification is generated as % error by the output.

Hence, OOB error is basically an estimate of how your RF model will ~~perform with new unknown~~ ~~dataset~~.

In RF, you don't worry about collinearity. So you are able to include anything & everything.

Page Date

rf \$err

would give output for all trees

OOB error rate
with 1 tree

	OOB	O	1
[1, 1]	0.3649	0.2682	0.5752
[2, 1]	0.3299	0.2355	0.5505
[3, 1]	0.3337	0.2390	0.5502
[4, 1]	0.3452	0.2418	0.5826

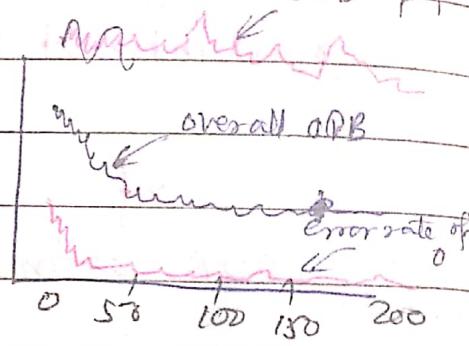
error rate
OOB of 1

OOB error
with 4 trees

plot (rf)

It tells that after 150 trees we are not getting any significant improvement in overall error rate.

So no need to build a forest of 200 trees.



To predict, `predict(rf, test) \rightarrow pred`

`table(test $ CREDIT, pred)`

`pred`: almost all 0's for 0's with 1's for 1's

0 164 11 there error is $164/175 = 0.936$.

1 48 27 and same for 1's

GRIP

To check which variables are important

`sort(Importance(rf))`, `varImpPlot(rf)`

It prints a table of Gini's measure which tells for all possible splits, how much is the contribution of each variable in reducing the impurity.

getTree (rf, 15)

gives complete information about 15th tree. However we cannot visualize RF.

Importance Table

Mean Decrease Gini

check - A

37.33

it checks across the

Duration

32.71

RF, for all possible

C-hist

19.08

splits, contribution of

Purpose

25.13

a particular variable.

1

1

Hence it gives you some

sort of explanatory power

you can use RF misclassification as a benchmark.

(best black box model's accuracy) and use these important variables with further Regression models.

Suppose RF gives error 23% but

DT or any Regression model gives 22%

so you can decide whether, 1% gain in accuracy

by RF is worth the explanation loss which we

are not getting because of using RF

However, if you are not clear about imp variables, have large rows and columns, so you

can go for DT/RF and use that information to

use in further models.

RF gives best case scenario and imp variables.

$$19 - 8 = 11$$

Big Problem

Boosting algorithm can not work with Categorical variables. All variables should be numeric

Page Date

(2) Boosting ensemble techniques - Gradient boosting

Here we build small trees (not fully grown) sequentially. First, we build 1 small tree and calculate the errors. Errors/prediction of this tree is used to build another small tree and so on. This process reduces the gap of errors with each step. You keep on doing it until you don't get further improvement.

Tree 1 (small tree)



now Errors → Tree 2



now Errors → Tree 3

$$\text{prediction}(t) = \text{prediction}(1) + \text{prediction}(2) + \dots$$

$$= \text{actual value} + \text{error}_1 + \text{error}_2 + \dots + \text{error}_{t-1} + \text{prediction}(t-1)$$

Ex -

$y = 40$ values, with 10 values each in range 15-25, 0-5, 30-35 and 15-17.

$x = 1 : 40$

plot(x, y)

create 1st tree $\text{rpart}(y \sim x, \text{method} = \text{"anova"}) \rightarrow \text{rt}$

plot predicted $\text{predict(rt1)} \rightarrow p1$, review

value of \hat{y} with $\text{lines}(x, p1, \text{col} = 1)$

Now to boost the model, we need to check the errors it is making.

$$e1 = y - p1$$

error = actual - predicted

2nd tree: $\text{rpart}(e1 \sim x, \text{method} = \text{"anova"}) \rightarrow \text{st2}$

$$p_2 = p_1 + \text{predict(st2)}$$

lines (x , p_2 , $\text{col} = 2$)

$$\text{sum}((y - p_1)^2) = 92.62$$

$$\text{sum}((y - p_2)^2) = 81.63$$

Here sum of squares error reduced from tree 1 to tree 2.

$$e_2^2 = y - p_2$$

$\text{rpart}(e_2 \sim x, \text{method} = \text{"anova"}) \rightarrow \text{st3}$

$$p_3 = p_2 + \text{predict(st3)}$$

$$\text{sum}((y - p_2)^2) = 73.419$$

Guideline rule: is Error should be related to x

Heteroscedasticity. So it gives strong results

Absence of heteroscedasticity \Rightarrow Homoscedasticity

No relation b/w ~~residuals~~ residuals/errors with x values.

Using R Package xgboost

imp for boosting all x should be numeric, as we use derivatives here

$\frac{\delta}{\delta b} \sum (y - a - bx)^2$ which needs all x to be numeric.

If you are creating credit scoring model, you can't rely on RF as accuracy is not the only intention, You need the explanation there.

Page Date

Imp

Popular method to code factor variables to numeric is One-hot coding - Dummy variable with all levels -

Factor Occupation

Doctor (D)

Engineer (E)

Lawyer (L)

One-hot coding

D E L

1 0 0

0 1 0

0 0 1

Imp

Short cut to do one hot coding for all variables

Say we run `glm` and fetch model matrix from it which we can directly use for one hot coding,

`glm(CREDIT ~ ., data = train) > m
model.matrix(m) = m`

`View(m)` has all variables with one hot coding

Use this now for gradient boosting,

`gboost(data = m, labels = train$CREDIT,`

`nrounds = 3, max_depth = 2, objective = "binary:logistic"`

3 trees sequentially

not fully grown
but only till 2 splits

`rg`

Output - trainerror: 0.264000

train-error: 0.264000

train-error: 0.274667

} misclassification measure

if you increase trees to 20 and depth to 4, error reduces a lot, so play with these parameters

gdp Gradient Boosting works better than RF. On Kaggle these days, boosting is the winner for most competitions.

Gives very good accuracy.

Page _____

Date _____

--	--	--	--	--	--

As test dataset is not in 0,1 column structure which is required for xgb model (in case of factors), so need to convert test dataset to model matrix.

`glm(CREDIT ~ ., data = test) → lt`

`model.matrix(lt) → mtest`

`predict(xgb, mtest) → pred`

`view(pred)` ← predicts probability of 1

`ifelse(pred > 0.5, 1, 0) → pclass` To chk classification

`table(test $ CREDIT, pclass)`

		p class		
		0	1	
		0	152 23	error = 66 / 250
		1	43 32	= 1264

To improve performance of gradient boosting, we use parameter "eta". It tells how much weightage to give on subsequent level to trees generally a value b/w 0 and 1 as low weightage should be given.

As we earlier added previous level p to next level p directly ($p_2 = p_1 + \text{pred}$), using eta we are adding the weighted p to next one.

`xgboost(data = m, label = train $ CREDIT, nrounds = 20, max_depth = 2, objective = "binary:logistic", eta = 0.05)`

Using $\text{eta} = 0.1$ in next step performed better → xgb

Overfitting
Using high eta or max depth, you can get less train error but test error would increase, so be cautious here.

age - 35 \rightarrow buy age - 35 < 0 not buy age \geq 35

age - $x \geq 0$ To keep threshold constant we add that bias
can change anything

Page	[]	Date	[]	[]	[]	[]	[]
------	-----	------	-----	-----	-----	-----	-----

ANN - Artificial Neural Network

1) Massive connectivity and parallel processing

Draws parallels with brain neurons. The way brain neurons learn on its known through learning, experience same way ANN learns on its known. 100 billion neurons in brain, each one connected to 5000 other neurons.

e.g. identifying a, a, A, a

how brain knows that above all comes in 'a' category.

not just by knowledge, brain tries to map slant, shape etc to the learning of a and classify anything in this category as a. So brain parallelly processes many,

2) Robust and Fault Tolerant, Capability to adapt to surroundings, if today you are in one city, move to another city, brain adapts to surroundings

3) Ability to learn & generalizes from known examples - 3 pics of cat, you can identify 4th also as cat

4) Collective behaviors is more imp than individual behavior.

Artificial Neural N/w aim is to mimic Biological Neural N/w

highly complex, non-linear, and parallel processing system - Brain

Eg. No linear answer, why we like a person and why we do not like.

linear. $\text{output} = \text{T}$ when we do not want to apply any activation function on final output. We are good with y_K as linear.

In R

Page Date

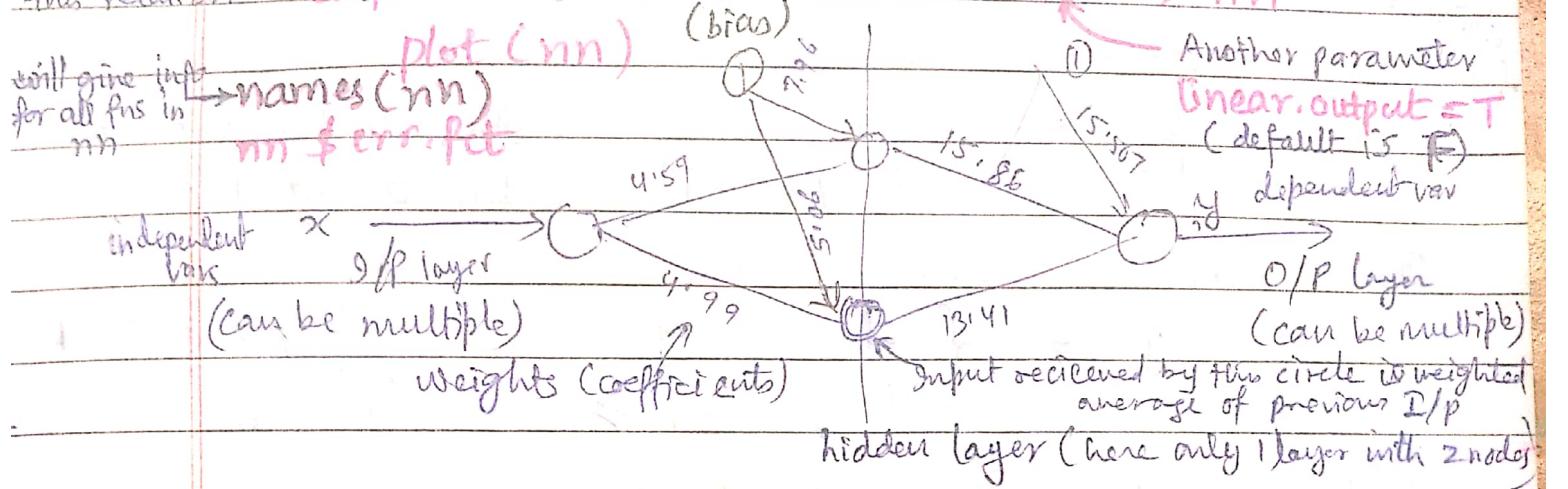
library(neuralnet)

data.frame ($y = \text{runif}(100, 40, 50)$) $\rightarrow d$

$d\$x \leftarrow \text{sqrt}(dfy)$

x, y have sqrt relationship

Building a nn, to make neuralnet (formula = $y \sim x$, data = d, hidden = 2, computer understand this relation. errfn = 'sse', threshold = 0.01) $\rightarrow nn$



typically $\text{hidden} = c(3, 2)$ - would give 2 hidden layers

we start with with 3 and 2 nodes.

more nodes

and then converge

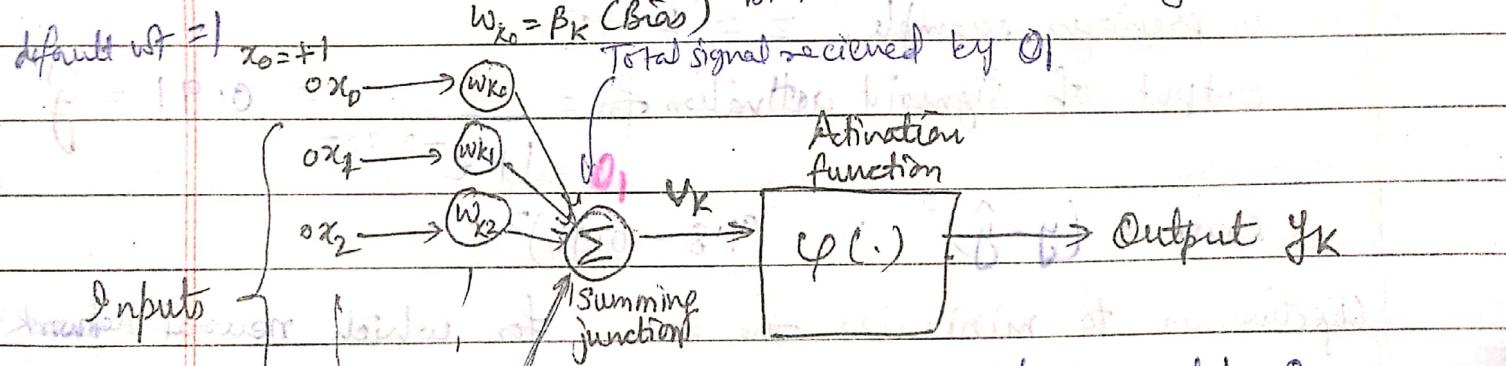
for ex (10, 5, 2) \rightarrow Compute(nn, d\$x) $\rightarrow p$; View(p)

p\$net.result

so when here, it is predicting \hat{y} for all values $y = 44.78$.

No variation so bad nn.

Bias is added at each layer.



where w_{kp} is weight between input x_p and neuron k . b_k is total quantity received by O_1

w_{kp} is weight between input x_p and neuron k . b_k is total quantity received by O_1

$w_{kp} = w_{k0} + w_{k1}x_1 + \dots + w_{kp}x_p$

Synaptic weight neuron (including bias) \rightarrow $w_{kp} = b_k + w_{k0} + w_{k1}x_1 + \dots + w_{kp}x_p$

weight has 2 neurons tells about the connection strength of w_{kp}

say $x_1 = 2, x_2 = -75, y = 1.8$

$b_k = 1, w_{k0} = 1, w_{k1} = 1.5, w_{k2} = -1.5$

$Z = 1x1 + 2x5 + -75x-1.5 = 2.375$

learning rate and structure (no. of hidden layers) bring more variation in prediction. Err-fct and Act.fct do not bring much diff.

Page

Date

Activation functions -

1) Softmax (Generalized logistic function)

$$y_j = \frac{\exp(z_j)}{\sum_{i=1}^n \exp(z_i)} \rightarrow \text{output is b/w 0 and 1}$$

Sum of all output values is 1

default
fn

2) Sigmoid function (most popular) (Logistic function)

$$y_j = \frac{1}{1 + \exp(-z_j)} \rightarrow \text{The output is continuous b/w 0 and 1}$$

3) Heaviside function (also known as Perception)

$$y_j = \begin{cases} 0, & z_j < 0 \\ 1, & z_j \geq 0 \end{cases} \rightarrow \begin{array}{l} \text{Output is 0 or 1} \\ \text{(Binary)} \end{array}$$

4) Hyperbolic (Tanh) function

$$y_j = \frac{\exp(z_j) - \exp(-z_j)}{\exp(z_j) + \exp(-z_j)} \rightarrow \begin{array}{l} \text{Output is continuous b/w} \\ -1 \text{ to } 1 \end{array}$$

5) Rectified Linear Unit (ReLU) - $y = \max(0, z)$ if z is ≥ 0

You have to try with various activation functions, how ever whichever gives best results. However Sigmoid and Tanh are used mostly.

In previous example, $z = 2.375$

$$\text{Output of Sigmoid activation fn} = \frac{1}{1 + e^{-2.375}} = 0.91 = y$$

$$\text{error} = (y - \hat{y})^2 = (1.8 - 0.91)^2$$

Objective is to minimize the error, for which neural network goes back and readjust the weights. nn starts with some numbers for weights (other than 0) and then settle down where weight fn is minimum error. $\rightarrow \text{error} = \text{expt} + \text{act} = S$

9mb

You can't use Categorical variables (Factors) for NN, as we won't be able to calculate Gradient descent. So factor variables need to be converted to 0/1 using One-hot coding

i=1 to n (for all rows)

Date						
------	--	--	--	--	--	--

(quadratic)

$$\rightarrow \text{Sum of square error fn} = \frac{1}{2} \sum_{i=1}^n (y_i - \hat{y}_i)^2$$

In all variables both err. func can be tried

or In continuous variables, use SSE

$$\rightarrow \text{Cross entropy error fn} = \sum (y_i \log p_i + (1-y_i) \log (1-p_i))$$

for classification, better to use Cross Entropy

though SSE can also be used.

As neural network looks at error, across whether large/small, go back to re-calculate weights, this process is called through gradient descent.

If more than one neuron, say 2 neurons, z_1 and z_2 would be calculated. For z_2 , weight would be different. If more layers of neurons, z_1 and z_2 becomes x_1 and x_2 for them. However, weights again would be different. Bias for these layer neurons would also be different.

Gradient descent - fundamental method of estimation

Say $y \approx x$ where b is coefficient of x .

$$\frac{1}{2} (y - bx)^2 = \text{SSE}$$

Take derivative of SSE w.r.t b (which tells about slope)

$$\frac{\partial \text{SSE}}{\partial b} = (y - bx)(-x) = \text{slope} = \text{gradient}$$

It tells with 1 unit change in b , how much SSE changes. Derivative is either an upward or downward slope. To reduce SSE (we will have to descend downwards).

$$\text{Adaptive learning rate} : b_{\text{new}} = b - (y - bx)(-x) \times \eta$$

η (Learning parameter)

However if we descend this fast, it can be a problem

So we want to descend slowly not by complete slope but with a small fraction $\eta = \eta$ (typically 0.1)

~~Matt Mazur step by step back propagation example~~

Page Date

$$SSE_{\text{new}} = \frac{1}{2} (y - b_{\text{new}} X)^2$$

if two outputs

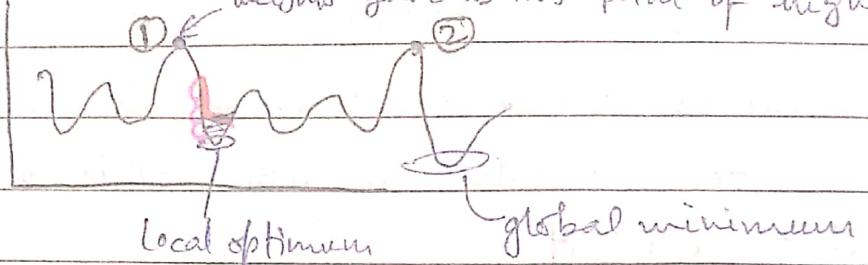
$$SSE_{\text{new}} = SSE_1 + SSE_2$$

Repeat this process again n again for revising weights until there is no more improvement possible

weights gave us this point of high SSE

plot of SSE for multiple wts

values



If we take huge jump in α for descent, we may miss the bottom point (local optimum here) and keep on oscillating as every time we have to go down the slope by α .

So better to descend slowly using small α so that oscillating problem is not there and we reach local optimum.

Secondly, we may be satisfied with local optimum only without looking at "Global minimum" which was further.

To deal with our 'n' multiple times, as starting at point 1, you may start with point 2, and then you reach minimum SSE (global minimum) so that search area is maximized. Plotting SSE with various values of b

Forward Pass — Computing \hat{y} and various wts (1st time)

Back Propagation — As we calculate SSE of last step (\hat{y}) and then propagate back to adjust weights in order to reduce SSE. All weights are recalculated (not just 1 or 2 of previous layer)

No of hidden layers — can't tell beforehand, have to try multiple times.

(Uses chain rule of derivatives in this field)

V Imp.

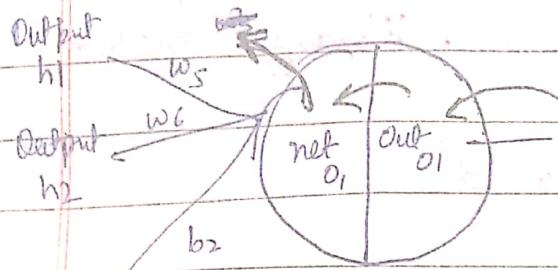
Tip
Range depends
on the problem

Imp tip for NN - Scale down all the variables on a scale of 0 to 1 to have better predictability (or -1 to 1)

$$\text{Range normalization} = \frac{\text{age}_i - \min(\text{age})}{\max(\text{age}) - \min(\text{age})} = \text{age}_{\text{norm}}^{(i)}$$

or score normalization

$$\frac{\partial E_{\text{total}}}{\partial w_5} = \frac{\partial E_{\text{Total}}}{\partial \text{out}_{01}} * \frac{\partial \text{out}_{01}}{\partial \text{net}_{01}} * \frac{\partial \text{net}_{01}}{\partial w_5}$$



$$E_{01} = \frac{1}{2} (\text{target}_{01} - \text{out}_{01})^2$$

$$E_{\text{total}} = E_{01} + E_{02}$$

$$\frac{\partial \left(\frac{1}{1+e^{-x}} \right)}{\partial x} = \frac{\partial \left(\frac{1}{1+e^{-x}} \right)}{\partial (1+e^{-x})} * \frac{\partial (1+e^{-x})}{\partial x}$$

Example

For some variables which we want to give more importance, can be normalized from -3 to +3 (or an expanded range) and others 0 to 1 or -1 to +1. After NN, these variables can be converted back to original using scaling.

Example - Car RAw

Classification exercise for prediction

Dataset is about probability of purchase of 2nd hand cars

V_1 - vhigh, high, medium, low - Price

V_2 - vhigh, high, medium, low - maintenance cost

V_3 - 2, 3, 4, 5 and more - no of passengers it can hold

V_4 - 2, 3, 4, more - no of doors in the car

V_5 - small, medium, big - boot space

V_6 - low, medium, high - safety features

Y - V_7 - unacceptable, acceptable - Dependent variable

This data can't be fit into ANN so need to do One-hot coding
`colnames(d) = c('price', 'mc', 'seats', 'doors', 'bootspace', 'safety', 'buy')`
`install.packages("onehot"); library(onehot)`
`codedd <- onehot(d)`
`dnew <- predict(codedd, d)`
`view(dnew)`

All factor variables have now been converted into one-hot variables. Convert dnew into dataframe

`as.data.frame(dnew) > dnew` (1728 observations, 25 variables)

In case of neural network we can't directly use ~~y > y ~.~~ (~ is not acceptable and it is very cumbersome to type all 25 variables)

`str=replace_all(colnames(dnew), "l", "l-") > colnames(dnew)`
`paste(colnames(dnew)[1:21], collapse = "+") > st`
`paste("buy_acc ~", st) > st1`
`as.formula(st1) > f` to convert string into formula.
 Note: ~~~~~ - ~~~~~

working of softmax with sigmoid

softmax is nothing but probability distribution function

softmax = $\frac{e^{x_i}}{\sum e^{x_i}}$

softmax = $\frac{e^{x_1}}{e^{x_1} + e^{x_2} + \dots + e^{x_n}}$

softmax = $\frac{e^{x_1}}{e^{x_1} + e^{x_2} + \dots + e^{x_n}}$

softmax = $\frac{e^{x_1}}{e^{x_1} + e^{x_2} + \dots + e^{x_n}}$

softmax = $\frac{e^{x_1}}{e^{x_1} + e^{x_2} + \dots + e^{x_n}}$

softmax = $\frac{e^{x_1}}{e^{x_1} + e^{x_2} + \dots + e^{x_n}}$

softmax = $\frac{e^{x_1}}{e^{x_1} + e^{x_2} + \dots + e^{x_n}}$

softmax = $\frac{e^{x_1}}{e^{x_1} + e^{x_2} + \dots + e^{x_n}}$

Time Series Forecasting

Page Date

Time series Analytical techniques

- Moving Average methods
- Stationarity, ARIMA

Objectives

- Understand the importance of forecasting and its impact on the effectiveness of the supply chain and overall performance of the organization
- Learn various components of time series data, such as trend, seasonality, cyclical components and random components.
- Learn different techniques such as moving average, exponential smoothing methods.
- Learn practical challenges associated with forecasting model

We all want to know/predict about the future. However if forecast is inaccurate, it can impact top line & bottom line of the organization. So it becomes necessary to forecast the demand and services for a product as accurately as possible.

e.g - Non availability of products / Too much inventory

- forecasting bench strength based on anticipated projects in Infosys

Forecasting Methodologies

- There are many different TS techniques
- It is usually impossible to know which technique will be best for a particular dataset - Reason is data would be completely different for different techniques.
- It is customary to try out several techniques and select the one that seems to work best
- To be an effective TS modeller, keep several TS techniques in tool box of yours like Simple moving average, Exponential moving average, Holt winter, ARIMA, Auto Regression methods
- Auto Regression is subset of linear regression.
- Eg - In linear, $X_t = \text{ad expenditure}$, $Y_t = \text{Sales in a particular month}$

$Y_t = \beta_0 + \beta_1 X_t$

Now if we change Y_{t-1} as previous month sales and try to predict this month sales, it becomes auto regression problem (extension of simple linear regression)

$$Y_t = a + b Y_{t-1}$$

here instead of using another variable X to predict Y , we use historical data of Y to predict current/future Y . Above is AR model with lag 1. If $Y_t \sim Y_{t-2}$, it is AR model with lag 2

What is time series data?

Time Series data is a sequence of observations collected from a process with equally spaced time periods

Examples — Daily closing stock prices

Libraries used in R - Forecast, date::Table, TTR

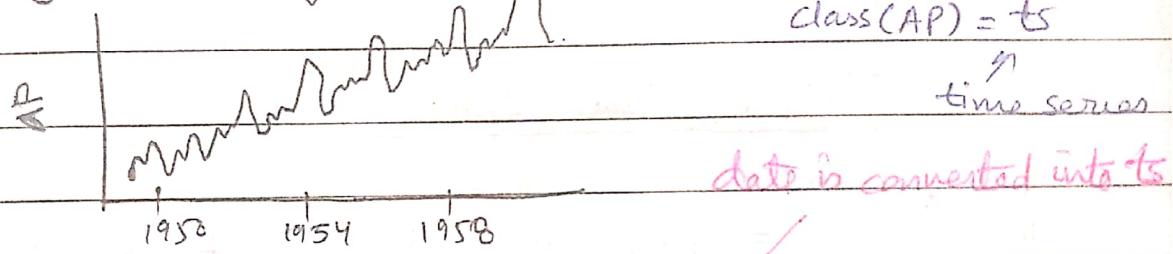
ts data \rightarrow Col 1 (Date/y/m/any timeframe (equidistant)), Col 2 (some measure)

Page Date

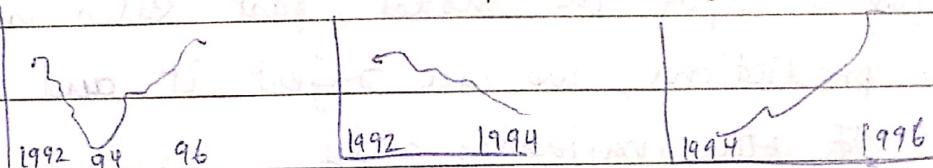
- Daily data on sales
- Monthly inventory
- Daily customers
- Monthly unemployment rates
- GDP

Time series should be equally spaced. Not like you have stock price data of weekly and daily mixed.

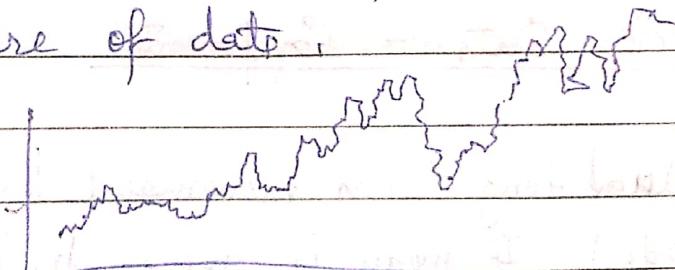
Monthly air passengers data plotted



Same pattern (upward/downward slope) is repeating every year. Also standard deviation and mean is increasing each year. $xrate.ts = ts(xrate, start=1991, frequency=4)$
plot(xrate.ts)



V pattern = breaking into 2 time periods to get clear picture of data.



Microsoft data with
upward trend but
lot of noise (lot
of jumping)

Why do we study TS analysis

- To understand past and to predict future, in order to take informed decisions.
- TS analysis quantifies the main feature in data and the random variation (Random variations are not always bad)
- Widely applicable in govt, industry and commerce,
- Deals with huge datasets,

Typical Time Series

$$\hat{y}_{t+1} = f(y_t, y_{t-1}, y_{t-2}, \dots) + f(x_1, x_2, x_3, \dots)$$

historical y

Other variables

* f can be linear or nonlinear

If after developing the model, past sales is not giving better predictions, we can reject it and consider effect of other variables only,

Time Series Data - Components

Trend - Gradual long term movement (upwards or downwards). If mean is going up, it is upwards. Trend is easiest to detect in the data.

e.g. population growth in India

Cyclical patterns - Results from events which are recurrent but not periodic in nature (no fixed pattern but it repeats). An up/down repetitive movement in demand repeats itself over a long period of time. Hard to detect. Repeating in infrequent interval
Eg. Recession in US economy,

Seasonality - Results from events that are periodic and recurrent in nature. An up/down repetitive movement with a trend, occurring periodically repeating with frequent interval. Seasonality is generally within a year.
Eg. Sales in festive season.

Random Components / Irregular components -

Disturbances or residual that remain after all above behaviours have been accounted for.

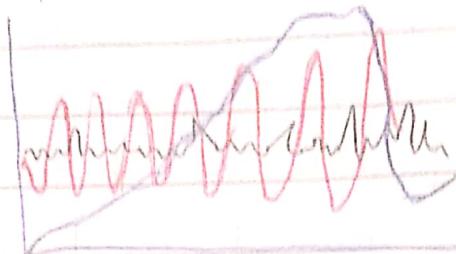
Erratic movements that are non predictable because they do not follow a pattern. Completely noisy.
Eg. Earthquake,

For a given TS data, we can separate out Trend, Seasonality and random components by using R function "Decomposition"

```

Sales <- read.csv("Sales_RegressionTime.csv")
SalesTimeSeries <- ts(sales[,1], frequency=4)
Sales <- data.frame(Sales) # attach sales
Sales <- seq(1:20) # adding 1 to 20 values as time
Saleslm1 <- lm(Sales ~ time, raw = TRUE) # poly means x + x^2
Saleslm2 <- lm(Sales ~ poly(time, 2), raw = TRUE)
for microsoft data, after decomposing, we can get

```



- Seasonal
- Irregular
- Trend

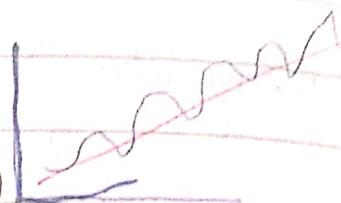
after decomposition, data starts showing patterns and information

Regression against time

predicting sales based on time

Linear Regression \Rightarrow

$$\text{Sales} = \alpha + \beta \text{time} \quad (y_t = 375 + 92.5 x_{1t})$$



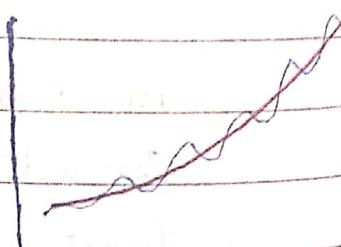
purpose of the model is to minimize RMSE (root mean square error) i.e. error b/w predicted sales and actual sales. Also we need to fit linear line

To improve prediction, we are adding another coefficient

Quadratic equation

$$\text{Sales} = \alpha + \beta \text{time} + \gamma \text{time}^2$$

$$\gamma = \text{gamma}$$



so instead of straight line, we have started getting curvilinear line.

$$y_t = 683.09 + 8.47 x_{1t} + 4.00 x_{2t}$$

Seasonal Regression Model

Looking at date pattern, we can say that it goes up and coming down every quarter

Bring seasonality - Sales & Seasonal \leftarrow as.factor(rep(c(1:4), 5))
 saleslm2s \leftarrow lm(sales ~ poly(time, 2) + seasonal, data = sales)
 points(time, predict(saleslm2s), type = "l", col = "red", lwd = 2)
 Poly(time, 1) = time
 $\text{poly}(\text{time}, 2) = \text{time} + \text{time}^2 = \cancel{\text{time}} + \text{time} + I(\text{time}^2)$
 In general we don't use $\cancel{\text{time}}^2$
 so we can try to bring that quarterly seasonal pattern in our model, by declaring some dummy variables

		value of			
takes care	of any	x_{3t}	x_{4t}	x_{5t}	
Quarter	1	1	0	0	
2	0	1	0		
3	0	0	1		
4	0	0	0		← both 0 for quarter 4

$$y_t = \beta_0 + \beta_1 x_{1t} + \beta_2 x_{2t} + \beta_3 x_{3t} + \beta_4 x_{4t} + \beta_5 x_{5t} + \epsilon_t$$

Above graph shows that forecasted values matched real values very closely by bringing seasonality into the model (was explained B and proved by the professor using R coding)

Forecasting Techniques

- Different forecasting techniques are developed based on different logics.
- Simple techniques like moving average and exponential smoothing predict the future value of a time-series data as a function of past observations.
- Regression based models such as auto-regressive (AR), moving average (MA), auto-regressive

Moving Average method

and moving average (ARMA), auto regressive Integrated moving average (ARIMA), and ARIMA with X (ARMAX) use more sophisticated regression models to forecast the future values of a time series data.

- Important tip is that using complex mathematical models does not guarantee more accurate forecast.
- Simple moving average technique may outperform complex ARIMA models in few cases.

Forecasting accuracy

Some metrics are required to compare and confirm the best model for a particular dataset

- Frequently used forecasting accuracy measures are:
- * Mean absolute error
- * Mean absolute percentage error
- * Mean squared error
- * Root mean squared error

Let y_t is the actual value of Y at time t and f_t is the corresponding forecasted value. Using these two values, we can calculate difference i.e., error. Assume there are $n = \text{ } \cancel{10}$ observations in the dataset

$$\therefore n = 100$$

Mean absolute error (MAE) -

$$MAE = \frac{1}{n} \sum_{t=1}^n |Y_t - F_t|$$

MAE is the average "absolute" error and should be calculated on the test data set (average of all the errors)

Mean absolute Percentage error (MAPE)

$$MAPE = \frac{1}{n} \sum_{t=1}^n \frac{|Y_t - F_t|}{Y_t} \times 100\%$$

- It is the average of all the errors in terms of percentage so it is easy to use and interpret.
- Widely used, since it is dimensionless, it can be used for comparing different models with varying scales.
- Suppose my model has MAPE 5% and Siddharth model has MAPE 6%, we can say, irrespective of different data sets, my model has better accuracy.

Mean Squared error

$$MSE = \frac{1}{n} \sum_{t=1}^n (Y_t - F_t)^2$$

- Lower MSE implies better prediction
- However, it depends on range of TS data!

Root mean squared error (RMSE)

$$\text{RMSE} = \sqrt{\frac{1}{n} \sum_{t=1}^n (Y_t - F_t)^2}$$

$$= \sqrt{\text{MSE}}$$

- RMSE and MAPE are 2 most popular accuracy measures of forecasting.
- RMSE is the standard deviation of errors or residual for sunny or any challenge kind of date, we can set the target for a specific value of RMSE.
- MAPE 5% gives good information. However RMSE 5 units is just a ~~not~~ number for error info.

Moving Average and weighted average method

- MA is one of the simplest forecasting techniques which forecasts the future value of a TS data using average (or weighted average) of past N observations.
- Mathematically, a simple moving average is calculated using the formula

$$F_{t+1} = \frac{1}{N} \sum_{k=t-N}^t Y_k$$

- It is SMA, as N past observations are given equal

weights, i.e. $\frac{1}{N}$

- In weight moving average, past observations are given differential weights (usually latest date gets more weight and oldest date gets least weight)
- Weighted moving average is given by

$$F_{t+1} = \sum_{k=t-N}^t w_k \times Y_k$$

where w_k is the weight given to value of Y at time k (Y_k) and $\sum_{k=t+1-N}^t w_k = 1$

In other words, total weights should be equal to 1.

- Eg, we want to predict sales of 13th month based on 12 month data.

^{SMA formula from TTR} For SMA - $F_{13} = \frac{1}{12} (S_1 + S_2 + \dots + S_{12})$

$$F_{14} = \frac{1}{12} (S_2 + S_3 + S_4 + \dots + S_{13})$$

for WMA - $F_{13} = 0.5 S_{12} + 0.3 S_{11} + 0.1 S_{10} + \dots +$

where all weights sum = 1

max importance is given to 12th month sales

for SMA - $F_{14} = \frac{1}{12} (S_2 + S_3 + \dots + S_{13})$

$$F_{15} = \frac{1}{12} (S_3 + S_4 + \dots + S_{14})$$

If we need to calculate based on last 12 months data, our average would move accordingly. That's why it is called **Moving Average method**

Check wsb Sales Prediction.r

If data is fairly static, MA or WMA gives good forecast (as no random variations in data so averaging works)

Page

Date

--	--	--	--	--	--

Single Exponential smoothing (ES) (Winters, 1960)

- One of the drawbacks of SMA is that it gives equal weight to all previous observations to forecast future value. Also how to conclude that 12 or any other number.
 - Drawback of weighted moving average is you need to assign different weights to past observations.
 - To overcome these problems, we use simple Exponential Smoothing (SES) technique.
 - Like MA, SES assumes fairly steady TS data with no specific trend, seasonal or cyclic components
- Imp: In SES, weights assigned to past data decline exponentially with the most recent observations assigned higher weights than old observations.
- If our data does not have any trend, cyclicity or seasonality and is random in nature, it is best to use MA, WMA or SES.

$$F_{t+1} = \alpha Y_t + (1-\alpha) F_t$$

- The parameter α in the above equation is called smoothing constant and its value lies between 0 and 1.
- Since model uses only one smoothing constant, it is called single exponential smoothing (SES).

- Substituting for F_t recursively in the above equation we get

$$F_{t+1} = \alpha Y_t + \alpha(1-\alpha)Y_{t-1} + \alpha(1-\alpha)^2 Y_{t-2} + \dots + \alpha(1-\alpha)^{t-1} Y_1 + (1-\alpha)^t F_1$$

of observation here we can see that the weights assigned should be taken to older observations decreases exponentially for averaging, i.e. from α to $\alpha(1-\alpha)$ to $\alpha(1-\alpha)^2$ -- to $(1-\alpha)^t$

12 month sales

${}^{12\text{th month}} \uparrow$
12 forecast sales.

Eg.

$$F_{13} = \alpha S_{12} + (1-\alpha) F_{12}$$

$$F_{12} = \alpha (S_{11}) + (1-\alpha) F_{11} \quad \text{say } \alpha = .90$$

$\frac{1}{90\% \text{ of } 11^{\text{th}} \text{ month Sales}}$ $\frac{10\% \text{ of forecasted value of }}{11^{\text{th}} \text{ month}}$

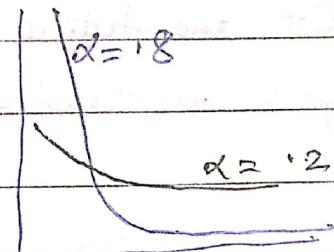
- So this model considers all past observations, though weights may be different.

- If $\alpha = 1$, it would decrease exponentially very fast slowly

- for $\alpha = .2$, we are

giving lot of importance to history. For $\alpha = .8$,

we are giving more importance to latest date.



SMA disadvantages and SES advantages

SMA disadvantages -

- Increasing n makes forecast less sensitive to changes in data
- It always lags behind trend as it is based on past observations, the longer the time period n , the greater the lag as it is slow to recognize the shifts in the level of data points.
- Imp - Not suitable for trend and seasonal pattern data, as it starts exhibiting forecast bias and systematic errors

SES Advantages -

- Uses all historic data, unlike the moving average where only the past few observations are considered to predict the future value (eg. 12 months MA)
- It assigns progressively decreasing weights to older data.

Optimal Smoothing constant in a SES (α)

- Whenever data is smooth (without much fluctuations), we may choose higher value of α .
- Whenever data is highly fluctuating, choose lower value of α .
- Find optimal value of α by solving non-linear

No need to do optimization problem. There is a package in R which can directly be used. Also Holt winter can also be used for getting α optimization problem (subject to constraint $0 < \alpha < 1$)

$$\text{Min}_{\alpha} \left[\sqrt{\frac{1}{n} \sum_t (Y_t - F_t)^2} \right]$$

For WSB Sales data, optimal value of α that minimizes RMSE is 0.1574 and corresponding RMSE is 7393.99 · 76

Double exponential Smoothing (Holt's method)

- SES drawback - it does not do well in the presence of trend
- So we introduce one more equation for capturing trend in the time-series data
- Double exponential smoothing uses 2 equations to forecast future values of T_t , one for forecasting the level (short term average value) and another for capturing trend.

Equation 1: Level (or intercept) or Average [Level is nothing but average(SMA)]

$$L_t = \alpha Y_t + (1-\alpha) F_t$$

Equation 2: Trend

$$T_t = \beta (L_t - L_{t-1}) + (1-\beta) T_{t-1}$$

α, β - smoothing constants. $0 < \alpha, \beta < 1$

Forecast at time $t+1$ is given by

$$F_{t+1} = L_t + T_t$$

$$F_{t+n} = L_t + n T_t$$

to initialize,
 $L_1 = Y_1$

$$T_1 = Y_1$$

where L_t = level which represents the smoothed value upto and including the last date,

T_t = slope of the line (Rate of increase/decrease at period period t)

n = number of periods well into the future

Ex

Initial values : $L_1 = Y_1$

$$T_1 = (Y_2 - Y_1) \quad \text{or}$$

$$\frac{(Y_t - Y_{t-1})}{t-1}$$

For WSB Sales data

$$L_1 = Y_1 = 3002666$$

$$T_1 = \frac{Y_{36} - Y_1}{35} = \frac{4732677 - 3002666}{35} \\ = 49428.88$$

$$F_2 = L_1 + T_1 = 3002666 + 49428.88 \\ = 3052095$$

Calculating forecasting values from 37 to 48 months

by using $\alpha = 0.0328$, $\beta = 0.9486$

2 ways of initializing T_1

$$T_1 = Y_1$$

$$T_1 = \frac{Y_{36} - Y_1}{35}$$

Initialization is very imp.

for the model to perform good.

$$Y = 17$$

$$T + \beta = 17 + 17$$

- SMA, WMA and SES uses past data to average out. So trend is lost. Only average (or level) was captured
- But in DES (Holt), trend is captured along with levels

Triple Exponential Smoothing (TES) Holt Winter

MA and SES / DES handles averaging and Trend but they do not handle seasonality in data.

Above methods indicate systematic error patterns (mainly two errors with -ve errors at fixed intervals, (actual - forecast)). This kind of error pattern implies presence of seasonality.

So TES is used.

BASS model - if a new product is being launched what would be the forecasted sales based on various attributes. BASS diffusion model (by Frank Bass)

Eq.1 Level (or intercept) equation

$$L_t = \alpha Y_t + (1-\alpha) [L_{t-1} + T_{t-1}]$$

S_{t-c}

Eq.2 Trend Equation

$$T_t = \beta (L_t - L_{t-1}) + (1-\beta) T_{t-1}$$

Trend index

Seasonal equation

$$S_t = \gamma Y_t + (1-\gamma) S_{t-c}$$

L_t
Seasonal index

For TES, directly use Holt Winter fn. in R
It would directly give optimum value of α, β, γ and
various coefficients as well as fitted values,
Holt Winter forecast gives Confidence Interval also.

- Off These kind of models are mainly suitable for static
data for high frequency data (where data is coming
every 1 hr, different kind of models are used)

Multiplicative model ↴

$$F_{t+1} = [L_t + T_t] \times S_{t+1} - c$$

where c = number of seasons

if monthly seasonality, $c = 12$

if quarterly seasonality, $c = 4$

if daily data, $c = 7$

To predict seasonality index using method of averages

Step 1 - Calculate average of value of γ for each season

$\bar{Y}_1, \bar{Y}_2, \dots$ for each season

Step 2 - Calculate average of season's average, $\bar{\bar{Y}}$

Step 3 - Seasonality index for season $k = \frac{\bar{Y}_k}{\bar{\bar{Y}}}$

Interpretation of seasonality index - can be understood
as % change from trend line

e.g. - seasonality index $1.088 = 108.8\%$ → In Jan, demand
will be approximately 8.8% more from trend line,
seasonality index for March $= 0.8885 = 88.85\%$ →
means, demand in March would be 11.5% less
than the trend line.

* Exponential Smoothing techniques are very sensitive
to initial values of level, trend & seasonality index

In this case, best model proved out to be Linear Regression model where Sales was regressed against Promotion Expenses and competitor promotion.

LR gave better accuracy than MA, SES, DSS, Holt Winter etc.

Use regression on time when trend is the most pronounced.

Another way of incorporating seasonality is to take the trend prediction and actual value.

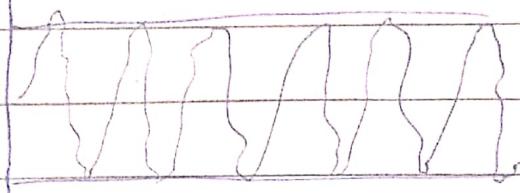
Depending on additive or multiplicative model, compute the deviation and map it as seasonality effect for each prediction.

Multiplicative seasonality means going wider and wider with time (expanding)

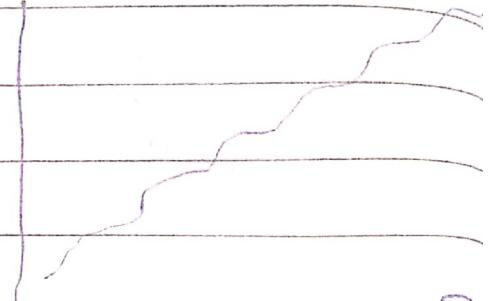
Holt winter method allows to choose seasonality as either additive or multiplicative.
Take, whichever gives best MAPE

Additive seasonality means being constant with time

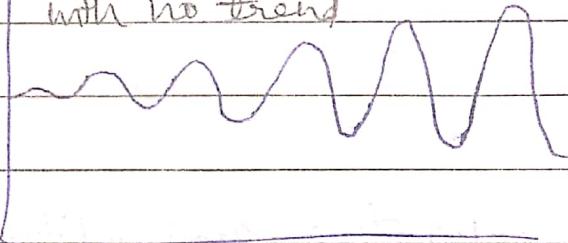
Additive seasonality with
no trend



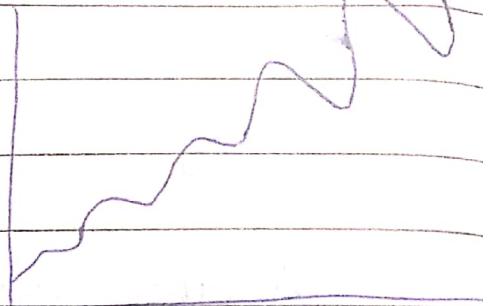
Add. seasonality with tre



Multiplicative seasonality
with no trend



Mult. seas - with pend



If there is no trend or if seasonality and fluctuations are more important than trend, then regression on time's coefficients behave weirdly

right approach: structural time series

structural time series model

additive seasonal component

multiplicative seasonal component

trend component

noise component

structural time series model

additive seasonal component

multiplicative seasonal component

trend component

noise component

ARIMA

Page Date

Box-Jenkins Methodology

- Individual forecasting methods AR and MA are combined to create ARMA and ARIMA models
- ARMA models are regression models - means a variable is regressed on itself measured at different time periods
- Fundamental assumption of AR - Time series is stationary process. (Constant mean/std over pd of time)
 - If TS data is not stationary, we have to convert it to stationary time series data before applying AR.
 - Constant mean/std means if 50 records are picked randomly, multiple times, they should have similar mean/std dev.

Limitations of Holt-Winter

- It did not worry much about the correlation b/w current month sales (y var) and previous month sales.
- Relationship was studied algebraically and not from date relationship (correlation)
 - what is the guarantee that my current month sales would be based on previous month sales.

frequency - 252 for stock data (only weekdays, not weekends)

Page Date / /

Sales

1	10
2	15
3	20
4	18
5	15

Considering $x = \text{Sales}[1:4]$

$y = \text{Sales}[2:5]$

trying to regress current sales on previous sales

Now if $\text{Cor}(x, y)$ is high say > 0.7 , we can say current sales is strongly dependent on previous month sale (Auto correlation)

Auto correlation - Correlation with itself with different lag.

In above case it is autocorrelation with lag 1.

$$\text{Sales}_{\text{Jan 2018}} = \alpha + \beta \text{Sales}_{\text{Jan 2017}}$$

here lag is 12 months

Stationary data - has constant statistical properties

- mean, variance, autocorrelation etc, over time

Non stationary - for stationary data, forecasting is easier

Differencing technique - to convert non stationary to stationary

- Generally 1 or 2 differencing are enough to convert a series to stationary

$$\text{Sales}_{t,d} = \text{Sales}_t - \text{Sales}_{t-1} \quad \text{first difference}$$

$$\text{Sales}_{t,2d} = \text{Sales}_{t,d} - \text{Sales}_{t-1,d} \quad \text{second difference}$$

In 99% cases, it would convert to stationary

non stationary data

Page

Date

miles

→ time

$$\text{miles.d} = \text{diff(miles)}$$

stationary data

Outlier

$$\text{miles.d2} = \text{diff(miles.d)} \leftarrow \text{second diff.}$$

If many outliers in first diff, go to 2nd diff.

AR model

$$\hat{y}_t = f(y_{t-1}, y_{t-2}, y_{t-3}, \dots, e_t)$$

AR(p) - When AR model depends on p past values, where lagged values having significant relationship with most recent value,

$$\hat{y}_t = c + \phi_1 y_{t-1} + \phi_2 y_{t-2} + \dots + \phi_p y_{t-p} + e_t$$

Imp factor - How to decide value of p (how many past values to consider). Can't regress against all values in data.

For this example, $p=12$ and $\beta_{11} = \dots = \beta_{12} = 0$

Imp task in auto Regression model

- identify the values of p (number of lags)
- p only determines right model ($y = f(x)$) equation) to be used for forecasting
- Standard approach used for model identification
 - ACF (Auto correlation function)
 - PACF (Partial auto correlation function)
- Auto correlation is correlation b/w y_t measure at diff time periods
- Say auto correlation with day 1 (γ_1) = $\text{cor}(y_t, y_{t-1})$
 $= 0.9$
- $\gamma_2 = \text{cor}(y_t, y_{t-2}) = 0.8$
- $\gamma_{12} = \text{cor}(y_t, y_{t-12}) = 0.99$
- Here present sales is strongly correlated with 12 months old sales, so we would build the model accordingly
- $\gamma_3 = \dots = \gamma_{11} = 0$

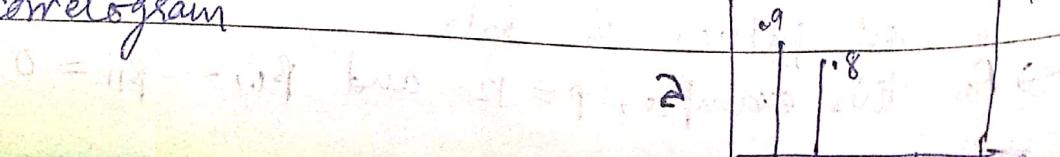
Now diff models like, model 1 = $y_t = \alpha + \beta_1 y_{t-1}$

model 2 = $y_t = \alpha + \beta_1 y_{t-1} + \beta_2 y_{t-2}$

model 3 = model 2 + $\beta_{12} y_{t-12} = y_t$

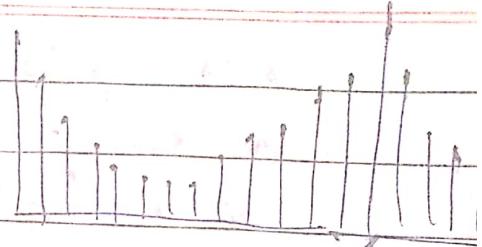
model 4 = $y_t = \alpha + \beta_{12} y_{t-12}$

ACF / correlogram — Plot of autocorrelation for different values of K is called ACF or correlogram



Function in R -

`acf(miles time series)`



Here 12th lag shows significant correlation

PACF - PACF of lag k , ρ_{pk} is the correlation b/w y_t and y_{t-k} when the influence of all intermediate values ($y_{t-1}, y_{t-2}, \dots, y_{t-k+1}$) is removed from both y_t and y_{t-k}

removing influence by other y_s .

$$\text{AR(1)} \quad y_t = \beta_0 + \beta_1 y_{t-1} \quad \text{PACF(lag 1)} = \beta_1$$

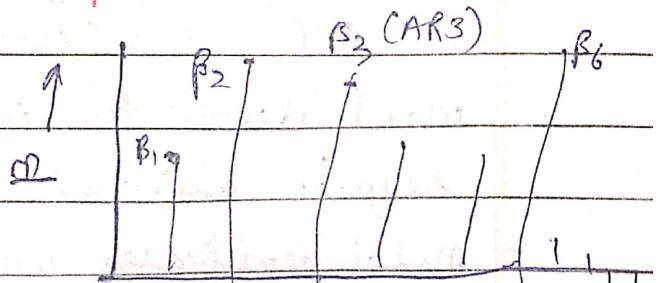
$$\text{AR(2)} \quad y_t = \beta_0 + \beta_1 y_{t-1} + \beta_2 y_{t-2} \quad \text{PACF(lag 2)} = \beta_2$$

AR(p)

if $|\beta_p| < \text{small}$ for all $p > 12$
that means all data points y_{t-13}, \dots are
insignificant. So need to take only from $p=1$
to 12.

So PACF is an iterative process for model selection

PACF graph



here after β_6 , β_s are not significant so consider only lag \rightarrow
only significant ones in your equation

ACF gives the confidence to select based on strong correlation.

PACF gives the order of selection

- Plot of PACF values for diff k is called PACF graph
- Hypothesis tests can be carried out to check whether the auto correlation and partial auto correlation values are different from zero.
- Null and alternate hypothesis
 - $H_0: \rho_k = 0$
 - $H_a: \rho_k \neq 0$
- Null hypothesis is rejected when $|\rho_k| > 1.96$ and $|\rho_k| > \frac{1.96}{\sqrt{n}}$ (95% quantile)
- To select appropriate p in AR model, following thumb rule may be used
 - Number of lags is p when
 - The partial autocorrelation, $|\rho_{pk}| > 1.96$ for first p values (first p lags) and cuts off to 0.
 - ACF, ρ_k decreases exponentially
 - Model identification is an iterative process and may require additional inputs
 - Model identification using ACF and PACF can not be taken as conclusive evidence for number of lags in AR process
 - Domain knowledge plays a key role.