

DECISION TREE

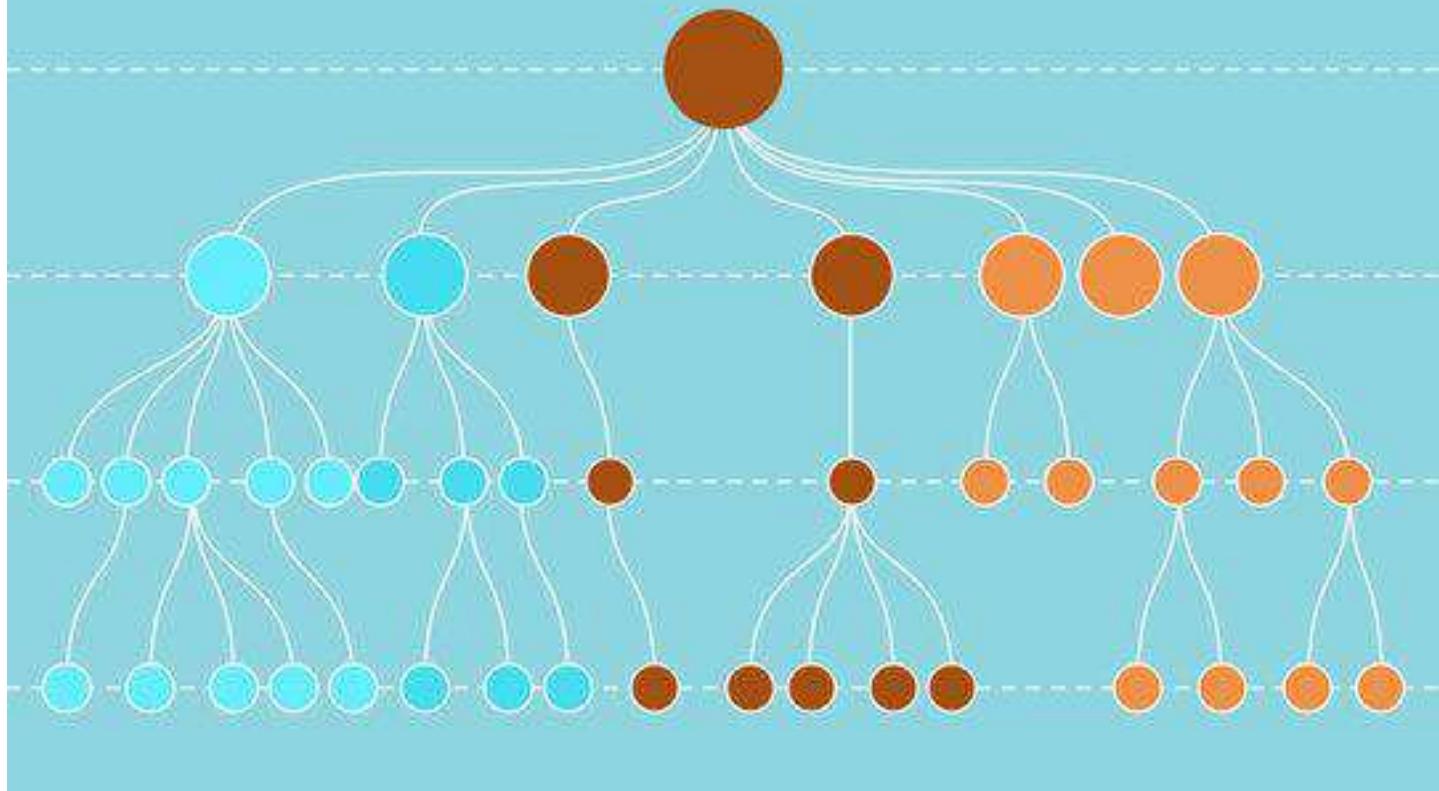


Table Of Contents

1. Why do we need Decision Trees
2. How it works
3. How do we select a root node
4. Understanding Entropy, Information Gain
5. Solving an Example on Entropy
6. Understanding Gini Impurity
7. Solving an Example on Gini Impurity
8. Decision tree for Regression
9. Why Decision Trees are Greedy Approach
10. Understanding Pruning

scaling is not needed for DT

③ Decision tree

for both Regression & classification
(supervised learning)

→ As of now we have Regression models,
which makes a best fit line.
or a best fit polynomial.

→ But In case of classification model,
we have this logistic regression
which can only classify by using line.

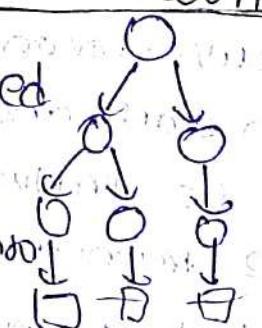
* what if our data is classified in polynomial like

 } Logistic Regression
can't be used here

So, this is where "decision tree" came to picture

→ As the name suggests, this algorithm works
by dividing the whole dataset into a
"tree-like structure" based on some rules &
conditions & then give predictions based
on those conditions. Let's see how it works.

- 1 → First it selects a root node based on given cond.
- 2 → then the root node will split into child nodes based on cond.
- 3 → If any node doesn't full fill the cond., then it again, splits into new cond.
- 4 → continues till all the cond. are met or if you have pre-defined the depth of your tree



Q) But, how do we select a root node?

Ans: This is where mathematical induction comes in.

→ usually Regression tree are used for quantitative data.

For scale or qualitative/categorical data we use classification trees.

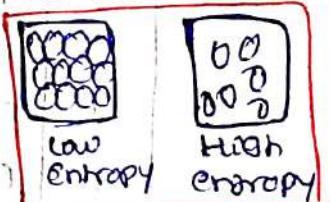
→ In regression tree, we split the nodes based on RSS (Residual sum of squares) criteria.

In classification, it is done using classification error rate, Gini impurity & Entropy.

We need a pure element (root node) to build tree.
So to find that we can use entropy.

* Entropy is the measure of randomness of data.
It gives the impurity present in the dataset.

→ When we split our node into two regions and put different observations in both the regions, the main goal is to reduce the randomness in the region & divide our data cleanly than it was in the previous node. If splitting the node doesn't lead to entropy reduction, we try to split based on different cond, or we stop.



→ A region is clean (low entropy) when it contains data with the same labels & random if there is a mixture of labels present (high entropy).

We need to calculate Entropy of each column u.

$$E = -P \cdot \log(P)$$

P = Probability,

$$E = -P \cdot \log_2 P = (P \cdot \log_2 1/P)$$

* Information gain (G_n) calculates the decrease in entropy after splitting a node. It is the difference b/w entropies before & after the split. The more the information gain, the more entropy is removed, & more clean the node.

$$G(n) = 1 - \sum \left(\frac{s_v}{S} \cdot E_i \right)$$

$$G(n) = \text{Entropy}(T) - \text{Entropy}(T^n)$$

Q) Let's take three data,

m	x	y	class
①	1	1	I
①	1	0	I
0	0	1	II
①	0	0	II

At we introduced new data,

new data
x: 0, y: 1, what is class?

First we need individual column entropy.

(E_x, E_y, E_z)

In each column we have (1,0) & total 2 class (I, II).

$$E_x(I) = \left(\frac{-2}{3} \cdot \log \frac{2}{3} \right) + \left(\frac{-1}{3} \cdot \log \frac{1}{3} \right) = -0.276$$

we have '3' ①s in x column,
out of which '2' are I class.

∴ ① out of 3 is I class.

I

$$E_I(0) = \left(\frac{1}{2} \cdot \log \left(\frac{1}{2} \right) \right) + \left(\frac{1}{2} \cdot \log \left(\frac{1}{2} \right) \right) = 0$$

$$E_I(1) = \left(\frac{1}{2} \cdot \log \left(\frac{1}{2} \right) \right) + \left(\frac{1}{2} \cdot \log \left(\frac{1}{2} \right) \right) = 0$$

$$E_I(0) = \left(\frac{1}{2} \cdot \log \left(\frac{1}{2} \right) \right) + \left(\frac{1}{2} \cdot \log \left(\frac{1}{2} \right) \right) = 0$$

$$E_I(1) = \left(\frac{1}{2} \cdot \log \left(\frac{1}{2} \right) \right) + \left(-\frac{1}{2} \cdot \log \left(\frac{1}{2} \right) \right) = 1$$

$$E_I(0) = \left(-\frac{1}{2} \cdot \log \left(\frac{1}{2} \right) \right) + \left(-\frac{1}{2} \cdot \log \left(\frac{1}{2} \right) \right) = 1$$

Information Gain

↑ No. of 1's
↑ No. of 0's.

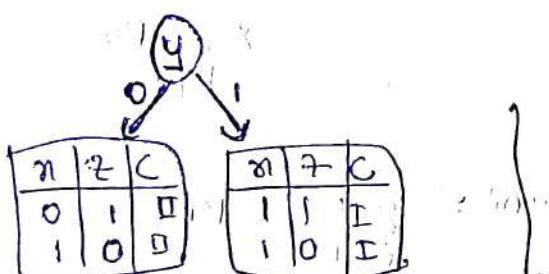
$$G_I = 1 - \left[\frac{S_0}{S} \cdot E_I(1) + \frac{S_1}{S} \cdot E_I(0) \right] \quad \text{Total no. of evaluation}$$

$$G_I = 1 - \left[\frac{3}{4} \times 0.276 + \frac{1}{4} \times 0 \right] \approx 0.27$$

$$G_I = 1 - \left[\frac{2}{4} \times 0 + \frac{2}{4} \times 0 \right] = 0$$

$$G_I = 1 - \left[\frac{2}{4} \times 1 + \frac{2}{4} \times 1 \right] = 0$$

$G_I = 1$, is the highest info gain, so loc entropy.
Hence it is Root Node: (y)

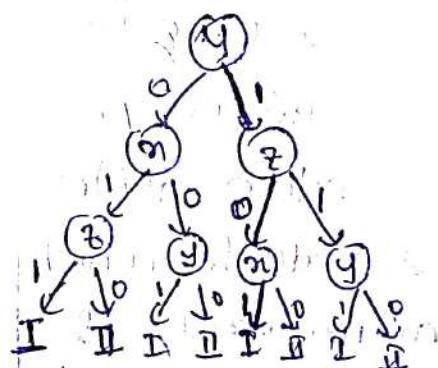


Again we have to find $E_{I0}, E_{I1}, G_{I0}, G_{I1}$.

Find largest if parent node.

until leaf node gives a class.

Finally it can be written:



So, $1 - 0 = 1 \rightarrow \text{I class}$

Gini Impurity

Gini impurity is a measure of how often a randomly chosen element from the set would be incorrectly labelled if it was randomly labelled according to the distribution of labels in the subset.

→ It is calculated by multiplying the probability that a given observation is classified into the correct class & sum of all the probabilities when that particular observation is classified into the wrong class.

→ Let's suppose there are k number of classes and an observation belongs to the class C_i , & its probability is given by P_i .

then probability that observation belongs to any other class other than C_i is,

$$\sum_{k \neq i} P_k = 1 - P_i$$

then gini impurity = $\sum_{i=1}^k P_i \times \sum_{k \neq i} P_k$

gini impurity value

will be 0 & 1

0 - no impurity

1 - Random distn.
the node for which

gini impurity is least is selected
as Root node
to split,

$$Gini = \sum_{i=1}^k P_i (1 - P_i)$$

$$Gini = \sum_{i=1}^k P_i - P_i^2$$

$$Gini = \sum_{i=1}^k P_i - \sum_{i=1}^k P_i^2$$

$$Gini = 1 - \sum_{i=1}^k P_i^2$$

b) cover data

<u>class</u>	<u>Gender</u>	<u>Stay in hotel</u>
9	m	Yes
10	F	No
8	F	Yes
8	F	No
9	m	Yes
10	m	No
11	F	Yes
11	m	Yes
8	F	Yes
9	M	No
11	M	No
10	M	Yes
10	M	No

Let's understand how root node is selected by calculating impurity.

We have 2 features which can use for hotel class & gender. We will calculate gini for each of the features & then select that feature which has least gini impurity.

Step 1

<u>class</u>	<u>stay</u>	<u>Total</u>	<u>P(class)</u>	<u>P(No)</u>
8	Yes = 2, No = 1	3	2/3	1/3
9	Yes = 2, No = 1	3	2/3	1/3
10	Yes = 1, No = 3	4	1/4	3/4
11	Yes = 3, No = 1	4	3/4	1/4
		14		

Gini for class feature

$$G(class=8) = 1 - (P(class))^2 - (P(No))^2 = 1 - (2/3)^2 - (1/3)^2 = 4/9$$

$$G(class=9) = 1 - (2/3)^2 - (1/3)^2 = 4/9$$

$$G(class=10) = 1 - (1/4)^2 - (3/4)^2 = 6/16$$

$$G(class=11) = 1 - (3/4)^2 - (1/4)^2 = 6/16$$

Weighted sum of gini for class

$$G(class) = \frac{\text{no of class 8}}{\text{Total}} \times G(8) + \frac{\text{no of class 9}}{\text{Total}} \times G(9) + \dots$$

$$G(class) = \frac{3}{14} \times \frac{4}{9} + \frac{3}{14} \times \frac{4}{9} + \frac{4}{14} \times \frac{6}{16} + \frac{4}{14} \times \frac{6}{16}$$

$$G(class) = 0.404$$

ii) Gini for gender

Gender	stain	n	D(G)	P(N)
m	4el>5, No>3	8	5/8	3/8
F	4el>3, No>3	6	4/6	1/2

$$a(m) = 1 - (5/8)^2 - (3/8)^2 = 0.167$$

$$a(F) = 1 - (1/2)^2 - (1/2)^2 = 0.5$$

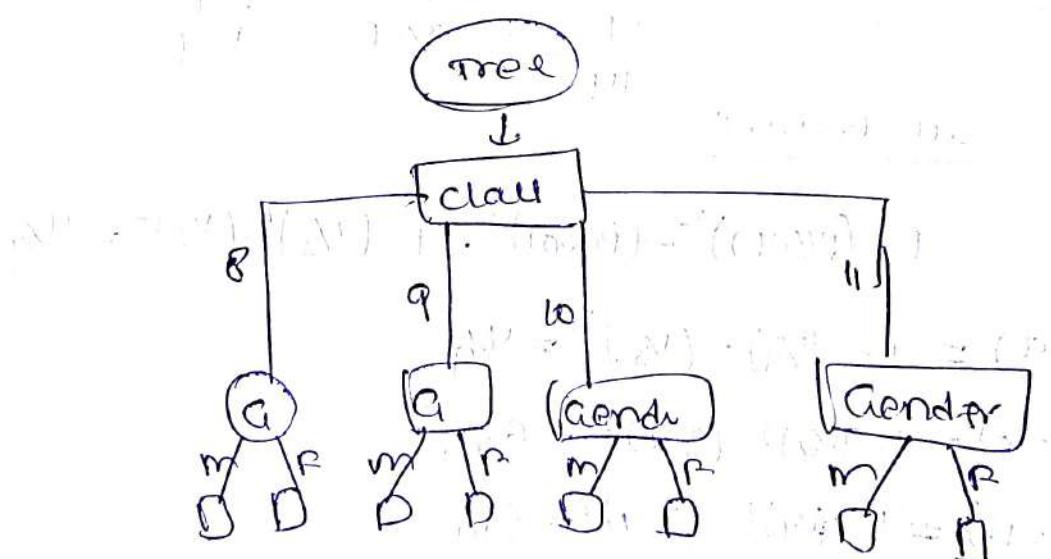
$$a(\text{Gender}) = \frac{8}{14} \times 0.167 + \frac{6}{14} \times 0.5$$

$a(\text{Gender}) = 0.16822$

we can see

$a(\text{class}) < a(\text{Gender})$

thus, root node
is class



this is how decision tree node is selected by calculating gini impurity for each node individually. If other features, then we need to repeat same process after selecting each node.

Note F

- ① we should control growth of tree by "min-sample-split" for which we will not, because a infinitely grown tree with depth (~100) also has a high chance at some point to overfitting.

Decision Tree for Regression

→ when performing regression with a decision tree, we try to divide these given values or x into distinct & non-overlapping regions.

Ex: For a set of possible values, x_1, x_2, \dots, x_p , we will try to divide them into J distinct & non-overlapping regions R_1, R_2, \dots, R_J . For a given observation falling into the region " R_j ", the prediction is equal to the mean of the response value for each training observation in the Region R_j .

- ② Regions are selected in a way to reduce residual sum,

$$\sum_{j=1}^J E[(y_i - \hat{y}_{R_j})^2]$$

y_i = mean draw the response var. in the region R_j !

Drawback of DT

- | | | |
|---|--|---------------------------------------|
| ① Small change in data can cause instability in the model because of Greedy approach. | ② Proba of overfitting is very high for DT | ③ Take more time to train a DT model. |
|---|--|---------------------------------------|

Recursive Binary splitting (Greedy Approach)

- As mentioned, we try to divide the x values into S regions, but is very expensive in terms of computational time to try to fit every set of x values into S regions.
- Thus decision tree opt for a greedy approach in which nodes are divided into regions based on the given condition is called greedy, but it does the best split at a given step at that point of time rather than looking for splitting a step for a better tree in upcoming steps.
- It decides threshold to divide others into different regions.

$$\sum_{i \in R_1(S)} (y_i - \hat{y}_{R_1})^2 + \sum_{i \in R_2(S)} (y_i - \hat{y}_{R_2})^2$$

These process have a high chance of overfitting the training data as it is very complex.

Tree Pruning

→ It is the method of trimming down a full tree to reduce the complexity & variance in the data. Just as we regularized linear regression, we can also regularize the decision tree model by adding a new term.

$$\sum_{m=1}^{|T|} \sum_{i: x_i \in R_m} (y_i - \hat{y}_{\text{fit}})^2 + \alpha |T|$$

where,

'T' is the subtree which is a subset of full tree T_0 . ' α ' is non-negative tuning parameter, which penalizes the MSE with an increase in tree len.

Post-Pruning (Backward pruning)

→ It is the process where DT is generated first & then the non-significant branches are removed. Cross-validation set or data is used to check the effect pruning & test whether expanding a node will make an improvement or not. If node is having improvement, we continue else convert to leaf node.

Pre-Pruning (Forward Pruning)

→ It stops the non-significant branch from generating, or uses a condition to decide when should it terminate splitting of some of the branch prematurely at tree is generated.

* Different algorithms for Decision tree

① ID3 (Iterative Dichotomiser) :-
It is one of the algorithms used to construct decision trees for classification. It uses information gain as the criteria for finding the root node and splitting them.

It only accepts categorical attributes

② C4.5 :-

It is an extension of ID3 algorithm, better than ID3 as it deals with both continuous & discrete values. Also used for classification purposes.

③ Classification & Regression Algo (CART) :-

It is the most popular algorithm for constructing decision tree. It uses gini impurity as the default calculation for selecting root node, however one can use "entropy" for criteria as well. Both for regression & classification problem.

Qn:- If asked to calculate, Entropy complicated.

Ans:- Which or why CART uses 'Gini'

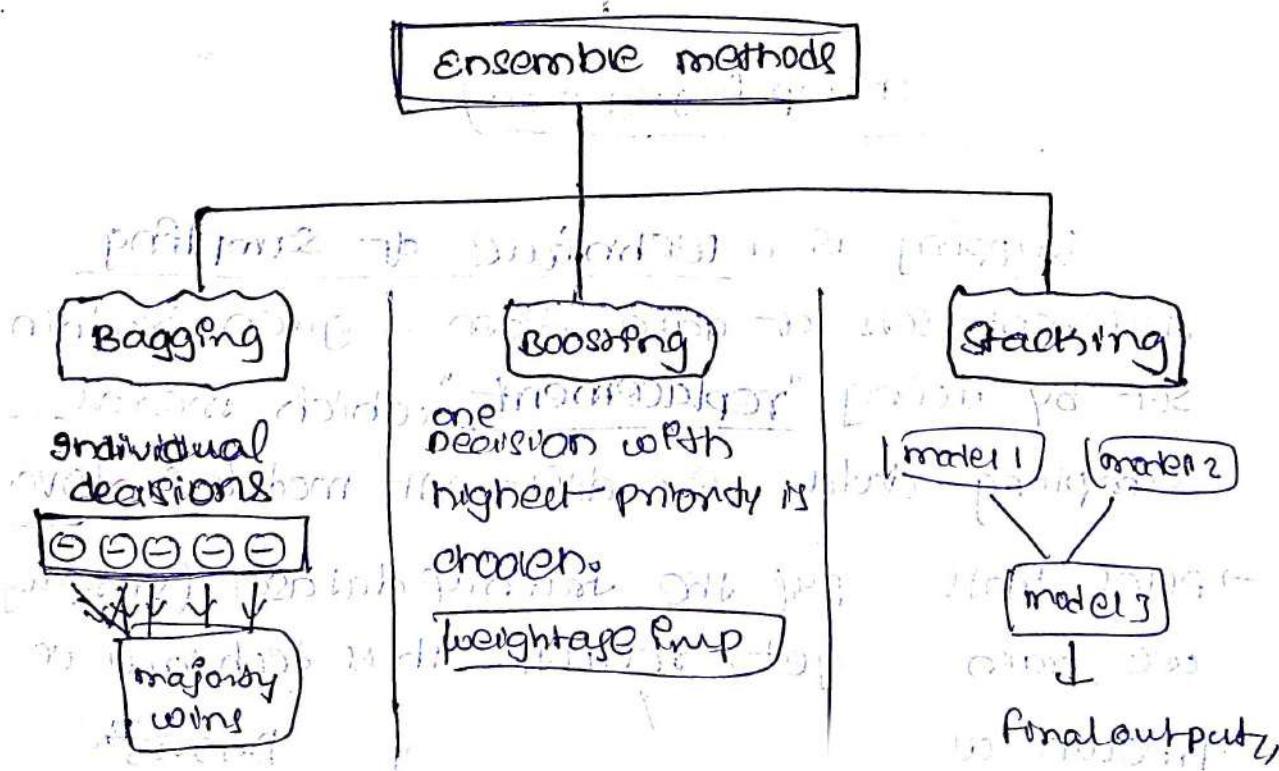
multiple
decision maker

* Ensemble Techniques & Random Forest

→ until we are predicting using one model, but what if we can take models and predict the output, won't be better than one. Of course it will be. This is the idea behind Ensemble.

→ we regularly come across the option of "Audience poll". And constant, mostly win when goel for option which had highest vote from audience. So, we can say, taking opinions from a majority or people is much more preferred than option of single person.

→ Ensemble techniques has a similar underlying idea where we aggregate group of predictions from a group of predictors. Such algorithms are called "Ensemble methods" and such predictors are called "Ensembles".



④ Let's suppose we have n predictors (models).

→ z_1, z_2, \dots, z_n with a standard deviation of σ

$$\text{var}(z) = \sigma^2 \quad \text{var of single predictor}$$

$$\text{Avg}(u) = \frac{(z_1 + z_2 + \dots + z_n)}{n} \quad \begin{array}{l} \text{Average of Predictors,} \\ (\text{Expected val}) \end{array}$$

If we use u as predictor then expected val still remains the same, but variance is reduced so much.

$$\text{var}'(u) = \frac{\sigma^2}{n}$$

This is why taking mean is preferred over single predictor.

* So Ensemble methods use multiple small models and combine their predictions to obtain a more powerful predictive power.

Bagging (Bootstrap Aggregation)

→ Bootstrapping is a technique of sampling different sets of data from a given training set by using "replacement", which means Sampling data for different models can overlap.

→ After bootstrapping the training dataset (sampling), we train & get results. This technique is known as Bagging / bootstrapping aggregation.

→ Bagging is a type of ensemble technique in which a single training algorithm is used on different subsets of training data where sampling is done with replacement (bootstrap).

→ In case of Regression, Bagging prediction is simply the mean of all the predictions & in case of分类, Bagging prediction is the most frequent prediction (majority vote).

④ Bagging is also known as "parallel mode", since we run all models parallelly and combine results at the end.

Advantages

- ① Bagging significantly decreases variance without bias.
- ② It works as well bootstrapping training data.
- ③ Works well with small dataset.
- ④ If training set is very huge, it can save computational time by training model on relatively smaller dataset & still can increase the accuracy of model.

Disadvantages

- It improves the accuracy of model on the expense of interpretability.
we can't interpret all those will be so many models.

Sampling → Similar to bagging but no replacement is done. This causes less diversity in the sampled datasets and data ends up being correlated. That's why Bagging is more preferable than sampling.

Boosting

→ Boosting is an ensemble technique (involves several trees) that starts from a weaker decision & keeps on building the model such that the final prediction is the weighted sum of all the weaker decision models.

The weights are assigned based on performance of individual tree.

→ In boosting, ensemble parameters are calculated in "stagewise way" which means that while calculating

the subsequent weight, the learning from previous tree is considered as well.

* i.e. the learning of previous tree boosts the learning of next tree and goes on -

→ we mostly used decision tree as weak classifier. Any other algorithm can be used as a base; but reasons for choosing tree are:-

Pros	Cons
1) Robust to outliers	1) Inability to extract a useful combination of features
2) Feature scaling not required	2) High variance leading to small computational power
3) Handles missing values	
4) can deal with irrelevant inputs	
5) computational scalability	

→ Boosting minimizes the variance by taking into consideration the result from small trees.

\$ General understanding eg. of boost

① you wanna travel to different place.

→ And you know a friend who is a traveller, so you ask him about the place (give your info) he will tell what he know & he will ask his friends.

→ He asks his friend who visited there to give some phone no or any agency.

→ And that friend ask another friend who live there.

First, we have to boost other person with info we have, and he will boost next person with info he have & finally get result.

② Dad wants to buy car.

→ He asks you & without any knowledge you say Yes. (Bad weightage)

→ Grandpa thinks and says what model should be bought consider grand children also wants car. (Good weightage)

→ consider son & grandpa, mom says what to buy, first new or second hand, based on financial status (Better weightage)

Stacking (stacked generalization)

→ stacking is a type of ensemble technique which combines the predictions of 2 or more models, also called base models and use the combination of the input for a new model

(re)

(meta model)

meta model - A new model is trained on the predictions of base models

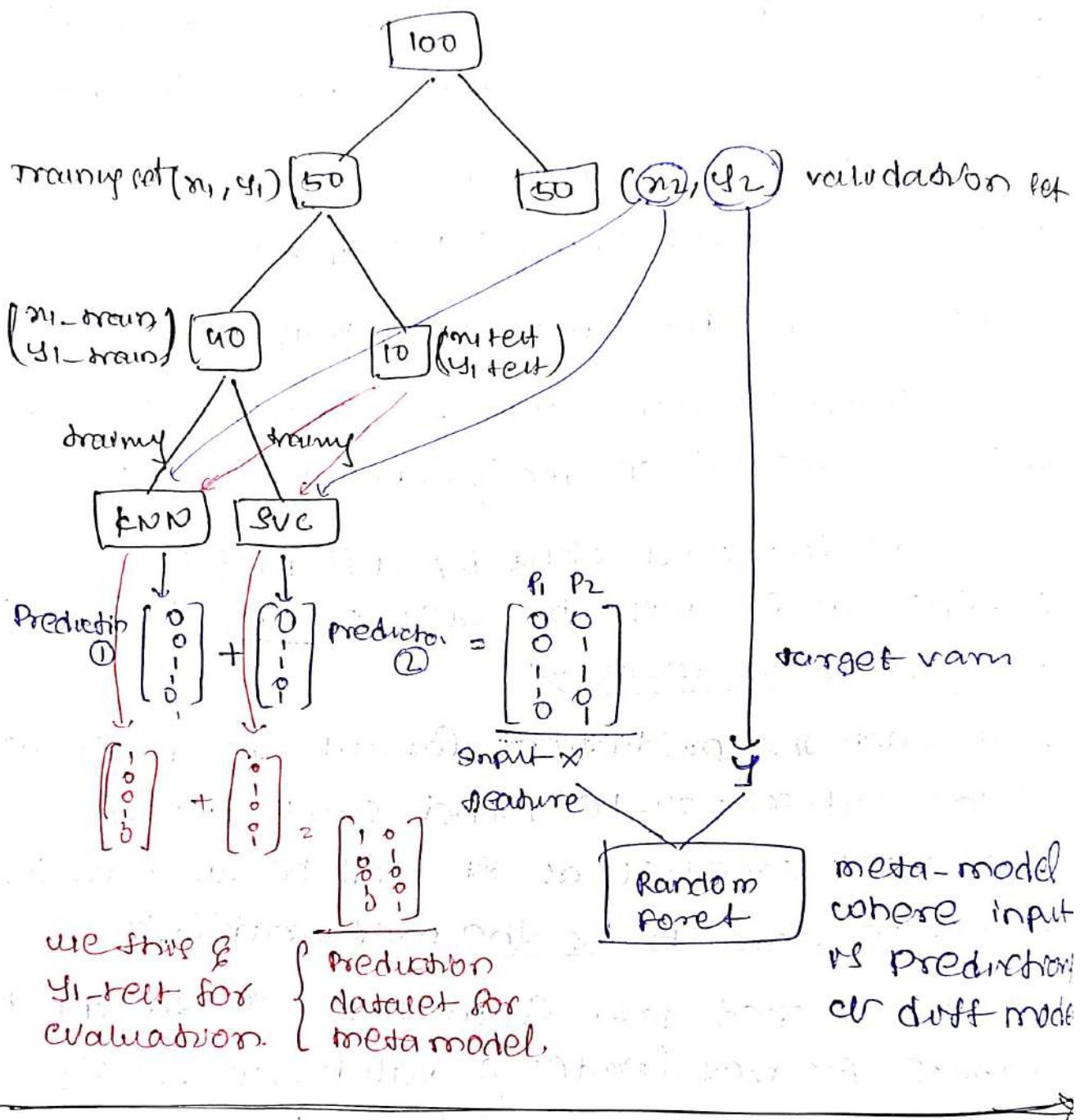
- suppose you have a classification problem & you can use several models like logitReg, SVM, KNN, Random Forest etc. the idea is to use few models like KNN, SVM as the base model & make prediction using that.
- Now predictions made by these models are used as an input feature for Random Forest to train on & give prediction.
- Stacking can be multi-level; using base model at level 1 the pass predictions onto another sub-base models at level 2 & so on

then at last any meta-model which take predictions of last sub-base models as input & do predictions,

working

- ① split the dataset into a training set & a holdout set.
generally we do a 50-50 split of training & holdout.
training set = x_1, y_1 ; Holdoutset = x_2, y_2 [validation set]
- ② split the training set again into training & test sets.
like, $x_1\text{-train}, y_1\text{-train}, x_1\text{-test}, y_1\text{-test}$
- ③ train all the base models on training set $x_1\text{-train}$, $y_1\text{-train}$
- ④ after training is done, get the predictions of all
the base models on the validation set x_2
- ⑤ stack all the predictions by diff models
together as it will be used as input feature
for the meta-model.
- ⑥ again, get the predictions for all the base models
on the test set $x_1\text{-test}$, and stack all their
predictions together as it will be used as the
- ⑦ prediction dataset for the metamodel.
- ⑧ use the stacked data from step 5 as the input
feature for metamodel & validation set x_2
as the target variable & train the model on
these data.
- ⑨ once the training is done check the accuracy
of metamodel by using data from step 7
for prediction & $y_1\text{-test}$ for evaluation
- ⑩ end

visual interpretation



out-of-Bag evaluation

- In Bagging, there might be some data which are never sampled at all. The remaining data which are not sampled are called out-of-bag instances.
- Since model never trains over these data, they can be used for evaluating the accuracy of model by using these data for prediction.

④ Random Forest

- Decision tree have low bias, but high variance. And we know bagging technique is a very good solution for decreasing the variance for a decision tree.
- Instead of using a bagging model with underlying model as a decision tree, we can also use Random Forest which is more convenient & well optimized for decision tree.
- The main issue of bagging is that there is not much independence among sampled datasets (there is correlation).
- The adv. of RF over bagging is that RF makes a tweak to working algorithm of bagging model to decrease the correlation in trees. The idea is to introduce more randomness while creating tree which will help in reducing correlation.

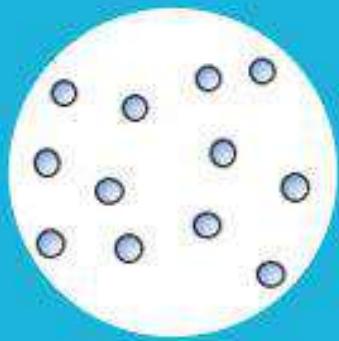
working of RF model

- ① different samples are collected from training dataset using bootstrapping.
- ② on each sample we train our tree model & we allow the tree to grow depth.
- ③ once the trees are formed, prediction is made by the random forest by aggregating the predictions of all the model. for regression model, the mean of all predictions is the final & for classification mode, the mode of all the predictions is called the final.

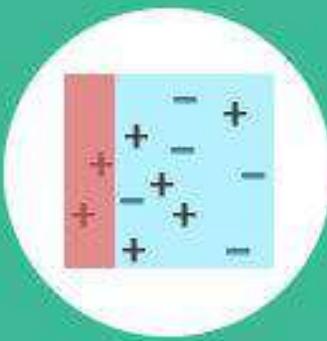
adv & disadv of Random Forest

- ① it can be used for both regres & classif problems.
- ② since base model is a tree, handling of missing values is easy.
- ③ it gives very accurate result with very low variance.
- ④ results of a random forest are hard to interpret in comparison with DTs.
- ⑤ higher computational time than other recursive models.
- ⑥ should be used when accuracy is upmost priority.

Monday, Jan 15, 2018



**Gradient
boosting**



AdaBoost

www.educba.com

Table of Contents

1. Understanding Boosting
2. Understanding AdaBoost
3. Solving and Example on AdaBoost
4. Understanding Gradient Boosting
5. Solving an Example on Gradient Boosting
6. AdaBoost Vs Gradient Boosting

Boosting

→ Boosting is an ensemble technique (involves several trees) that starts from a weaker decision & keeps on building the model such that the final prediction is the weighted sum of all the weaker decisions made.

the weights are assigned based on performance of individual tree

→ in boosting, ensemble parameters are calculated in "stagewise way" which means that while calculating the subsequent weight, the learning from previous tree is considered as well.

* i.e. the learning of previous tree boosts the learning of next tree and goes on -

→ we mostly used decision tree as weak classifier.
Any other algorithm can be used as base;

but reason for choosing tree are -

Pros	Cons
1) Robust to outliers	1) Inability to extract a linear combination of features
2) Feature scaling not required	2) High variance leading to small computational power
3) Handles missing values	
4) can deal with irrelevant inputs	
5) Computational scalability	

→ Boosting minimize the variance by taking into consideration the result from random tree.

⑧ general understanding eg dr boost

① you wanna travel to different place.

→ And you know a friend who is a traveller, so you ask him about the place (give your info, he will tell what he know & he will ask his friends).

→ He asks his friend who visited there to give some phone no or any agency.

→ And that friend ask another friend who live there.

First, we have to boost other person's worth into we have, and he will boost next person with into he have & finally get result.

② dad wants to buy car.

→ He asks you & without any knowledge you say yes. (Bad weightage)

→ Grandpa thinks and says what model should be bought consider grand children also wants car. (Good weightage)

→ Consider son & grandpa, mom says what to buy, first new or second hand, belief on financial status (Better weightage)

(6) AdaBoost (Adaptive Boosting)

$$\begin{matrix} w_1 \\ w_2 \\ \vdots \\ w_n \end{matrix}$$

→ Assume that the no. of training samples are "N", and the no. of iterations (tree created) is "M". Note that possible class outputs are $Y = \{-1, 1\}$.

① Initially, total weightage should be equal "1".

so. Initially consider observation weight as $w_i = \frac{1}{N}$

② for $m = 1$ to M ,

→ fit a classifier $G_m(m)$ to the training data with w_i

→ compute $\text{err}_m = \sum_{i=1}^N w_i \cdot I(y_i \neq G_m(m))$ → Prediction

→ compute $\alpha_m = \frac{1}{2} \log \frac{(1 - \text{err}_m)}{\text{err}_m}$ → This is the contribution of that tree to final result

→ calculate new weights using the formula,

$$w_i \leftarrow w_i \cdot \exp[\alpha_m \cdot I(y_i \neq G_m(m))]$$

③ normalize the new sample weights, so that their sum is (1) .

④ construct the new tree using the new weights,

i.e; consider weights obtained in 1st tree as the weights for the second tree and, weights of 2nd tree to the 3rd and so on.

⑤ at end, compare summation of results from all trees and final result of either the one with highest sum (for Regr) or the one with most weightage (for classification),

Eg

- ① understand ~~read~~ the Adaboost by simple dataset for heart patient prediction.

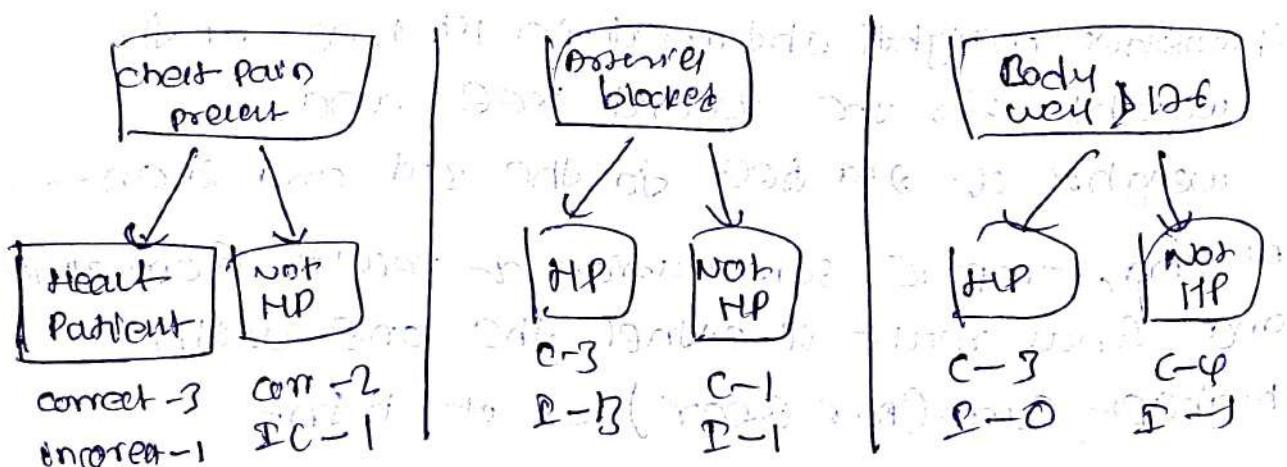
Ex	is chest pain	Arteries blocked	weight or person	is Heart Patient	new weight
0	Yes	Yes	205	Yes	0.05
1	No	Yes	180	Yes	0.05
2	Yes	No	210	Yes	0.05
3	Yes	Yes	162	Yes	0.33
4	No	Yes	156	No	0.05
5	No	Yes	125	No	0.05
6	Yes	No	160	No	0.05
7	Yes	Yes	172	No	0.05

Step ① - Generalize weights

→ there are 8 rows in our dataset.

$$\therefore w = \frac{1}{N} = \frac{1}{8} \quad (\text{samples are equally important})$$

Step ② - consider columns to create a weak decision model and then try to figure out what all correct & incorrect predictions based on that column



→ we will calculate "Gini Index" of each column.

And we select tree with the lowest GI.

& this will be the first decision maker for our model.

over 100

Step ③ Now, we calculate the contribution of tree to our final decision (column)

→ we got weight of first decision maker & we have one incorrectly predicted by it.

$$\text{error} = \frac{1}{8}$$

$$\text{contribution} = \frac{1}{8} \log\left(\frac{1-\text{err}}{\text{err}}\right) \approx 0.97$$

→ Step ④ → modify the weights

→ increase the sample weight for

$$\text{incorrectly predicted, new weight} = \text{old weight} \cdot e^{0.97}$$

$$w_{10} = \frac{1}{8} \cdot e^{0.97} = 0.33$$

→ decrease the sample weight for

$$\text{correctly predicted, new weight} = \text{old weight} \cdot e^{-0.97}$$

$$w_7 = \frac{1}{8} \cdot e^{-0.97} = 0.08$$

And add weights & divide by summation to

normalize it.

there, new normalized weights will act as sample weights for the next iteration,

Step ①

- And we create new tree which consider dataset prepared with new sample weight.
- suppose, m tree classify a person as HP &
n tree classify a person as not HP.

then contribution of all m tree & all n-tree
are added to

Higher is taken

②

* Gradient Boosting Tree

- Gradient Boosted tree use decision tree as estimator.
It can work with different loss function,
evaluate its gradient & approximate with simple tree

Adaboost is a special case of Gradient boosting where it use exponential loss function.

Algorithm

- ① calculate the avg. of the label column at initially this avg shall minimize the total error
- ② calculate pseudo residual.

$$\text{pseudo residual} = \text{actual label} - \text{predicted (avg) result}$$

mathematically,

$$\textcircled{4} \text{ derivative of pseudo residual} = \frac{\delta L(y_i, p_{mi})}{\delta (f_m))}$$

→ here, gradient or error term is getting calculated at the goal is to minimize the error. Hence, name is gradient boosted tree.

\textcircled{3} create a tree to predict the pseudo residuals instead of a tree to predict for actual val. of target variable.

$$\text{new result} = (\text{previous result}) + (\text{learning rate} \times \text{residual})$$

$$f_t(n) = f_0(n) + v \cdot \delta$$

v - learning rate
δ - residual.

\textcircled{4} Repeat these steps until residual stop decreasing.

EBM to understand GBM.

Q) we need to develop a model to predict an apartment rent based on sq. footage data.

Sq feet	Rent
750	1160
800	1200
850	1200
900	1400
950	1600

STEP 0 - set initial model

lets set avg rent price in our model: $f_0(n) \geq 1160$

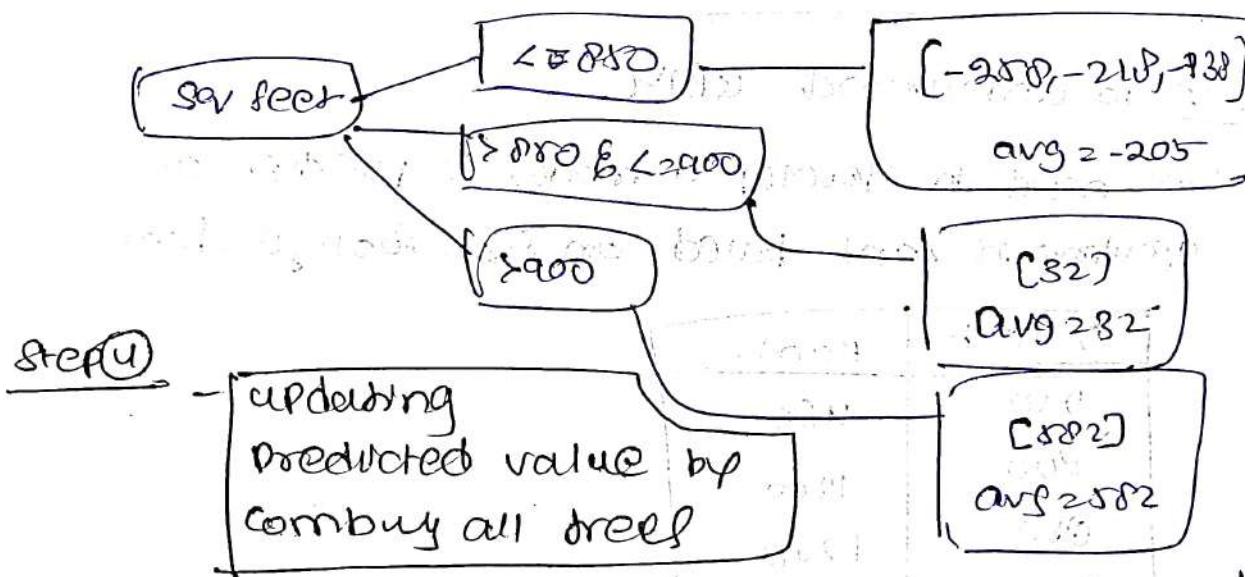
Step ② - calculate the residuals

→ For every instance, we compute the residual as diff b/w actual & predicted val (avg)

<u>sq feet</u>	<u>actual rent</u>	<u>Predicted rent (avg)</u>	<u>Residual</u>
750	1160	1018	-288
800	1200	1018	-218
850	1280	1018	-138
900	1400	1018	32
950	2000	1018	582

Step ③ - Build a model on residuals, instead on actual col. values

we build a decision tree to predict residual



Step ④

Updating Predicted value by combining all trees

$$\text{from updated predicted val} = P_0(n) + (\text{learning rate} \times \text{Residual})$$

<u>$P_0(n)$</u>	<u>Residual</u>	<u>Learn rate</u>	<u>$P_1(n)$</u>
1018	-205	1	1213
1018	-205	1	1213
1018	-205	1	1213
1018	32	1	1050
1018	582	1	2000

step 6 — repeat step 2-4, for the predefined no. of iterations,

step 6 — finally, the composite model sums together all of the weak models into one strong prediction.

GBM vs AdaBoost

→ essentially both AdaBoost & GBM have similar intuition, i.e., to convert a set of weak learners into a single strong learner.

⊕ however, the only point of difference is,

"How they create a weak learner during process?"

→ AdaBoost, changes the weights assigned to each of the instances, such that every subsequent learner focused more on the difficult ones and therefore, contributed to creating a strong learner.

→ Gradient boosting, trains the weak learner on errors. Each iteration computes the residuals & fits the next learner. Finally, using an optimization process, the algorithm computes the contribution of each weak learner which minimized the overall error of the model.

⑧ XG Boost (eXtreme Gradient Boosting)

- XG Boost regularizes data better than normal gradient boosted tree.
- XG Boost's objective function is the sum of loss function evaluated over all the predictions & regularization func for all predictions (of tree).

$$\text{obj}(\theta) = \sum_{i=1}^n l(y_i - \hat{y}_i) + \sum_{j=1}^J \lambda L(f_j)$$

- loss function depends on the task being performed and a regularization term is demanded,

$$L(f) = YT + \frac{\lambda}{2} \sum_{j=1}^J w_j^2$$

watcher over score.

for controlling overall no. of new created leaves.

- unlike other tree-building algorithms, XG Boost doesn't use entropy or Gini index. Instead, it utilizes gradient/error term and a "heuristic" for measuring the tree.

- Heuristic for a Regr. Problem is residual and for classification problem, Heuristic is a second order derivative of the loss at the current estimate.

$$h(m) = \frac{\partial^2 L(y, f(m))}{\partial f(m)^2}$$

$m = f^{(m-1)}(m)$

Tree Buildup in XG Boost

- ① Initialize the tree with one leaf.
- ② compute the "similarity" using,

$$\text{similarity} = \frac{\text{gradient}^2}{\text{hessian} + \lambda}$$

where, λ - Regularization term.

- ③ Now, for splitting data into a tree form, calculate,

$$\text{gain} = (\text{left similarity}) + (\text{right similarity}) - (\text{similarity for Root})$$

- ④ For tree pruning, the parameter γ is used. The algorithm starts from the lowest level of the tree & then starts pruning based on the value of γ .

If $(\text{Gain} - \gamma) \leq 0$, then remove that branch.
Else, keep the branch.

- ⑤ learning or done using new eq.

$$\text{new value} = \text{old value} + \eta \times \text{Prediction}$$

K Nearest Neighbors

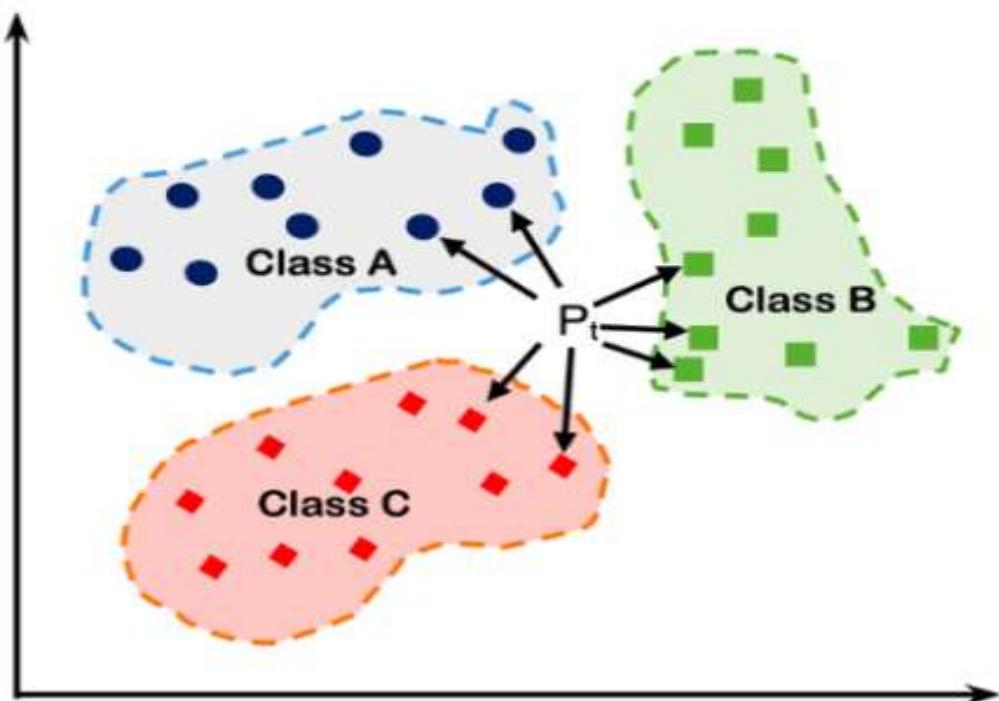


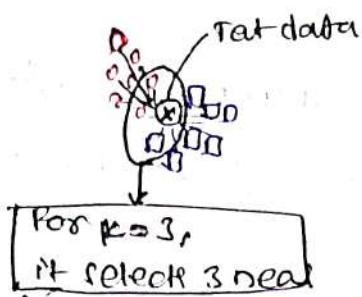
Table Of Contents

1. How does K-Nearest Neighbours work
2. How is Distance Calculated
 - Eculidean Distance
 - Hamming Distance
 - Manhattan Distance
3. Why is KNN a Lazy Learner
4. Effects of Choosing the value of K
5. Different ways to perform KNN
6. Understanding KD-Tree
7. Solving an Example of KD Tree
8. Understanding Ball Tree

(a) K-Nearest neighbour (KNN)

- KNN is a type of supervised learning algorithm which is used for both Regression & classification, mostly for classification.
- Given a dataset with different classes, KNN tries to predict the correct class of test data by calculating the distance between the test data and ALL THE TRAINING POINTS!! All
Each element in test is measured with all the training points.
- And then select the k points which are closest to the test data.
Hence, "k" in KNN is,
"How many closest point are you selecting to predict data?"
- Once the points are selected, the algorithm calculates the probability (for classification) of the test point belonging to the class of the k-training points and the class with the highest probability is selected.
- In case of Regression Problem, the predicted value is the mean of the k-selected training points.

Algorithm / visual illustration



KNN algorithm calculates distance b/w the test data and the given training data. And then it will select k points which are nearest.

→ And there are red, 1 blue. Hence, it is Red.

→ If it is Regr prob, the predicted value will be mean of those three.

How distance calculated?

* Euclidean distance:-

- It is most commonly used method to calculate the distance b/w two points.
- The Euclidean distance b/w two points, $P(p_1, p_2)$ & $Q(q_1, q_2)$ is calculated as

$$d(P, Q) = \sqrt{(q_1 - p_1)^2 + (q_2 - p_2)^2}$$

* Hamming distance:-

- Hamming distance is the distance metric that measures the no. of mismatched b/w two vectors. It is most widely used on categorical data.

- Generally, if we have features all categorical data then we consider the difference to be 0 if both values are same, difference is 1 if both are different.

e.g) $(A, B, C, D) - (A, D, C, B) = (0, 1, 0, 1) = 2$,

* Manhattan distance

- also known as L₁ norm, taxicab form, rectilinear dist, city block distance.

→ this distance represents the sum of absolute differences b/w the opposite values in vector.

$$MD(x, y) = \sum_{i=1}^n |x_i - y_i|$$

→ manhattan distance is less influenced by outliers than the Euclidean distance, with very high dimensional datasets more preferred.

* Lazy learners

- KNN algorithm is often termed as lazy learner.
- And most of the algorithms like Bayesian classifier, logistic regression & SVM etc. are called lazy learner. These algorithms generalize over the training set before receiving the test data.

i.e; they create the model based on the entire training data before receiving the test data, and then do the prediction for the data.

- But this is not the case with KNN algorithm, it does not create a generalized model (equation) for the training set but waits for the test data. Once the test data is provided then only it starts generalizing the training data to classify the test data. Here, entire training data is the model, waits for test set.

→ So, a lazy learner just stores training data &

* Weighted nearest neighbour

- As name suggest, we assign weight to the K-nearest neighbour. The weights are typically assigned on the basis of distance.
- Sometimes rest of data are assigned '0' also.
- The main intuition is that the point in neighbor should have more weight than farther points.

* Choosing the value of K

- The value of K affects the KNN classifier drastically.
- The flexibility of model decreases with the increase of K.
- If 'K' value is low, then model has high variance & low bias.
As 'K' value is increased, then variance decreased & bias increased.
- With very low value of K, there is a chance of algorithm overfitting the data, where as with very high value of K, there is a chance of 'underfitting'.

* Different ways to Perform K-NN

→ the way K-NN classifies the data by calculating the distance of test data from each of the observations and selecting K-value, this approach is known as "Brute Force KNN" [which is computationally very expensive (time taking)]

→ so, the idea behind using other algorithm for KNN classifier is to reduce the time during test period by preprocessing the training data in such a way that the test data can be easily classified in the "appropriate clusters".

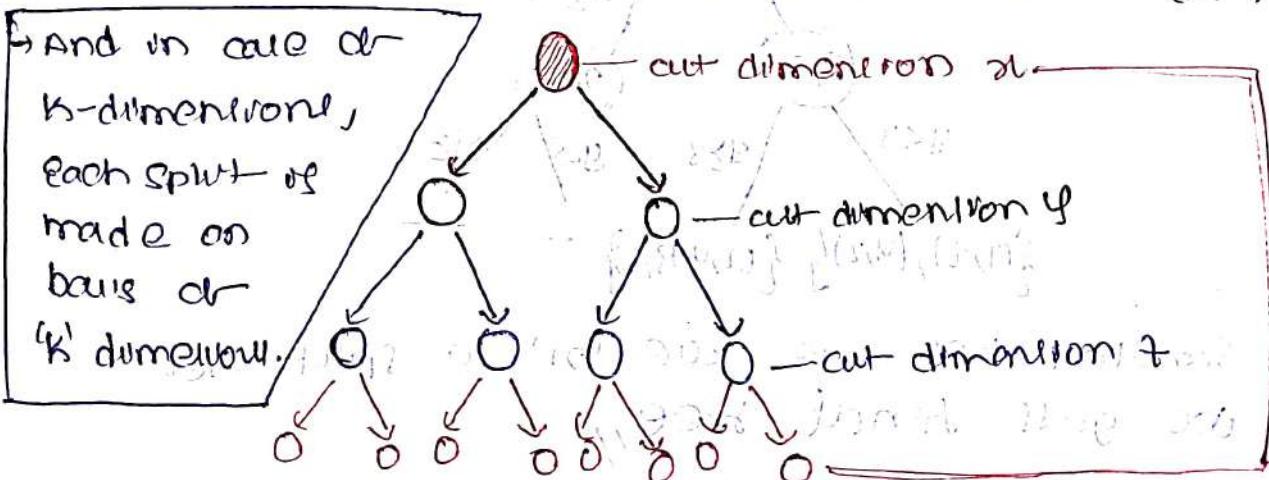
* K-dimensional tree (Kd Tree)

→ K-d Tree is a hierarchical binary tree.

→ It rearranges the whole dataset in a binary tree structure, so that when test data is provided, it could give out the result by traversing through the tree (every split eliminates half data), which takes less time than brute search.

Ex

let's say we have 3 dimensional data (2, 4, 8).



Eg)

Q1 training data $\Rightarrow \{ (1,2), (2,3), (2,4), (3,6), (4,2), (5,7), (6,8), (2,5), (8,5), (9,1), (9,3) \}$

H

Here, $k=2$

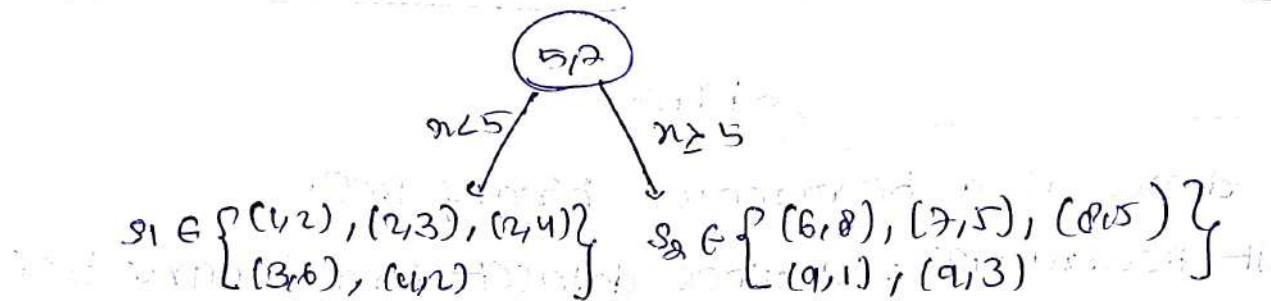
lets build our ad-tree

(i) Sort data & choose median as split point.

ie, $\{ 1, 2, 2, 3, 4, 5, 6, 7, 8, 9, 9 \}$

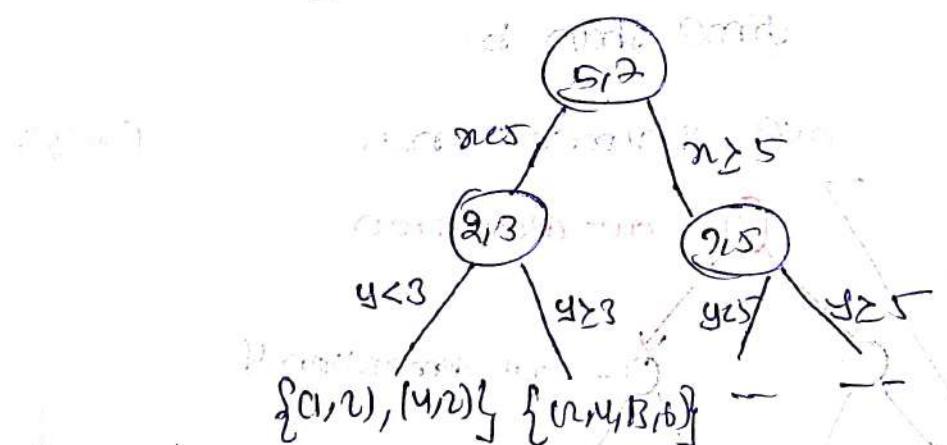
at index 5 \swarrow median.

our first node will be $(5,7)$; $x \geq 5$ (SPLIT cond)



(ii) splitting of right side nodes into further

ie, $y_1, y_2 \in \{ 2, 3, 4, 6 \}$ & $y_{S2} \in \{ 1, 3, 5, 7, 8 \}$



Similarly, here we can go to split and we get final tree.

Note (kd-tree)

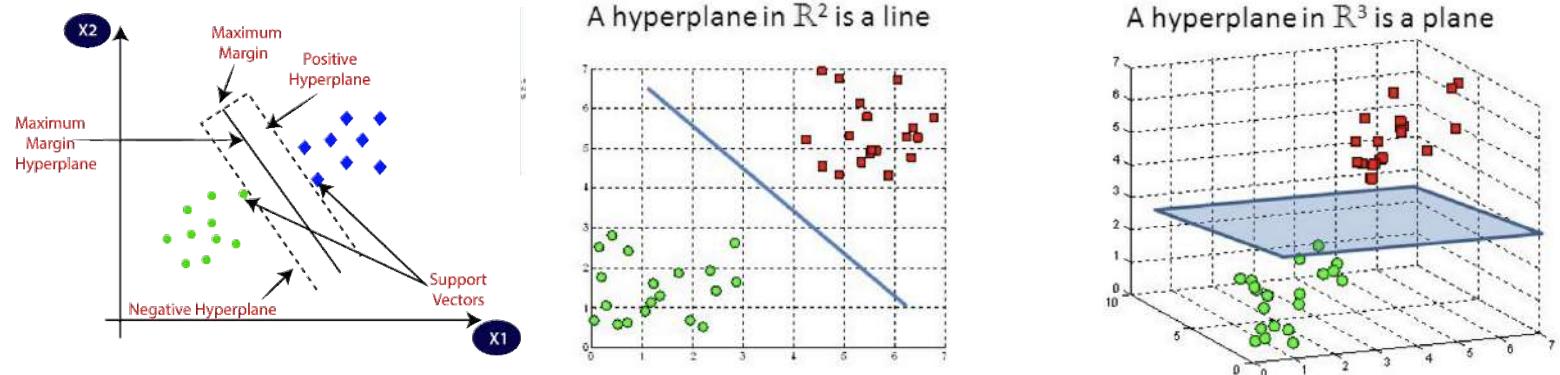
- once the tree is formed, it is easy for algorithm to search for the probable nearest neighbor just by traversing tree.
- the main problem of kd tree is that it gives probable nearest neighbors but can miss out actual nearest neighbors.

ii) Ball-tree

- similar to kd tree, ball tree are also hierarchical data structures. they are very efficient specially in case of higher dimensions.

Formed by following steps:

- two clusters are created initially.
- all data points must belong to atleast one cluster.
- one point cannot be in both clusters.
- distance of point is calculated from centroid of each cluster. the point closer to centroid goes into that particular cluster.
- each cluster is then divided into subclusters again, and then the points are classified into each cluster on the basis of distance from centroid.
- thus, this is how clusters are kept divided till a certain depth.



Understanding Support Vector Machines

Table Of Contents

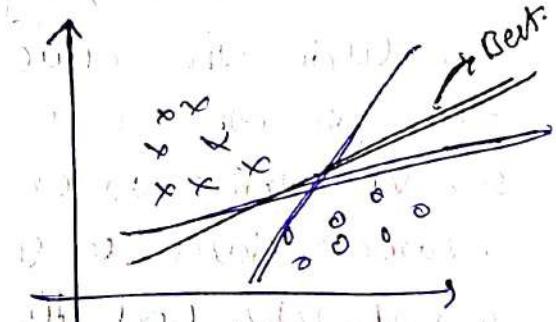
1. Understanding Concept of SVC
2. What are Support Vectors
3. What is Margin
4. Hard Margin and Soft Margin
5. Kernelized SVC
6. Types of Kernels
7. Understanding SVR

* Support Vectors Machine (SVM)

- support vector machine is a supervised machine learning algorithm, which can be a classifier as well as regressor. This is a linear model, and if consider linear classification models we have
- Logistic Regression and support vector classifier.
- In case of logistic regression, we ultimately build a line/plane/hyperplane, and classify the points to the side of the line, and given a new point based on the probability, we can say which part or plane it belongs. Now, what about SVM.

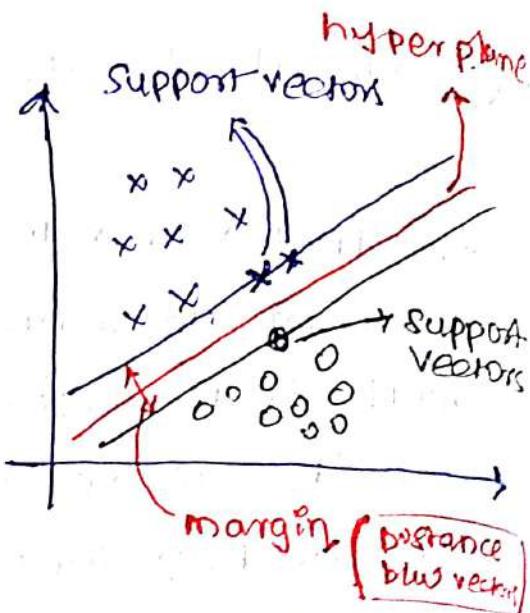
Concept of SVM

- Similar to logistic regression, we get a plane/hyperplane in SVM, but how do we decide on the best hyperplane. There can be several planes passing through the data points.
- So, the key principle behind SVM is to find the hyperplane / decision boundary, that maximally separates data points belonging to distinct classes, that is, maximizing the margin between datapoints.



① what is a support vector

→ let's suppose we have a decision boundary and to both sides of line we take the closest points and build a parallel plane to our decision boundary.

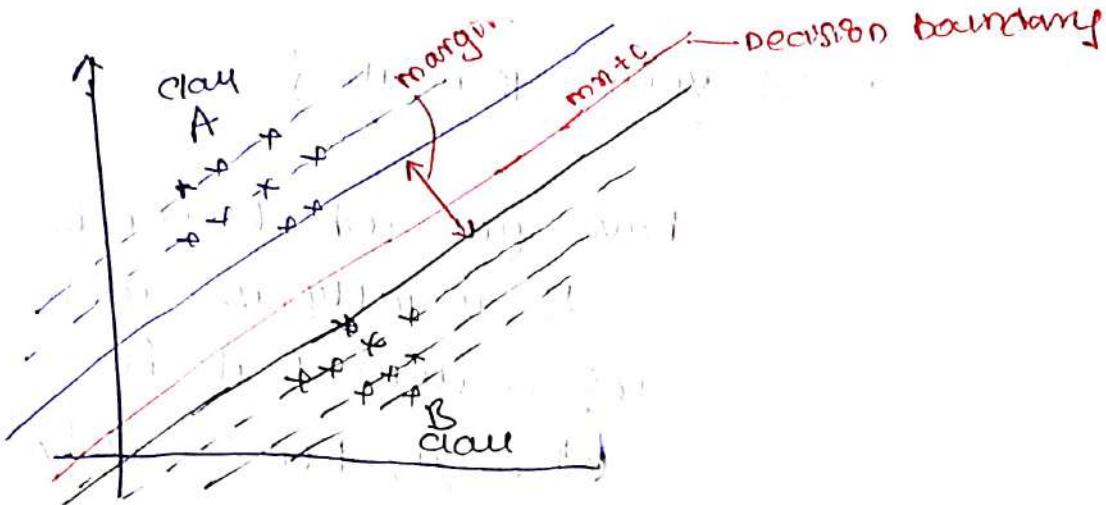


→ the points that are close to decision boundary are called "support vectors". And these support vectors are crucial for defining hyperplane so these support vectors help us find optimal plane.

② what is Margin

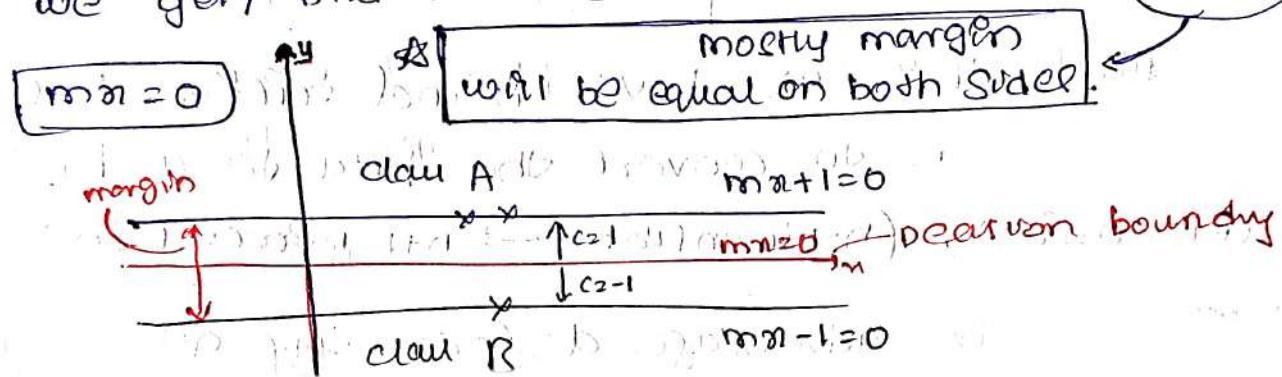
→ If we consider logistic Regression, based on the probability we will decide the class for a new point, but there are chances for things to go wrong here, say probability is 0.9 & so we conclude the class, but it can be wrong!

so, SVM tries to create maximum separation between data to get rid off the confusion or probability. And the maximum separation can be achieved through the distance between two closest vector lines from opp. sides.



→ Decision boundary can be represented as, $y = mx + c$,

→ consider decision boundary is at origin, we get, and assume support vectors have $c=1$.



+ Now, if we want to predict a new datapoint, then we first checks the y-value of $f(x)$. If it falls in (0 to 1), it comes under class A.

or if it falls in (0 to -1), it comes under class B and we don't need probabilities. Just by the range of y , we are predicting the new datapoint.

① soft margin and Hard Margin

→ the concept of SM & HM, refer to different approaches in handling the margin violations & the presence of noise or outliers in the data.

Hard Margin

- HM aims to find the optimal hyperplane that completely separates the classes without any margin violations. OR
- it works well when the data is linearly separable & free from noise.

$$y_i^*(w^T n_i + b) \geq 1$$

→ signifies all data points are correctly classified.

Soft Margin

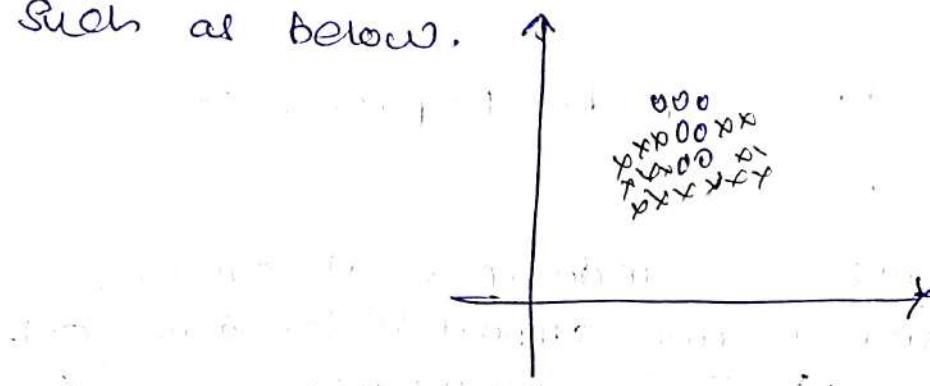
- SM allows for some margin violations & misclassifications, thus providing flexibility when dealing with noisy / overlapping data.
- OR introduced a slack variable (ϵ_i) to handle margin violations, allowing some points to fall within the margin / on wrong side of boundary.

$$y_i^*(w^T n_i + b) \geq 1 - \epsilon_i$$

→ where $\epsilon_i \geq 0$ for all points,

⑥ Kernelized support vector machines

+ It's not always we have get which can be easily separable by lines!! usually, we will be having data in such a way that it cannot be separable through lines easily. Such as below.



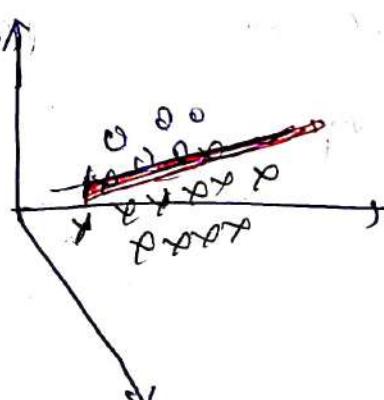
→ To deal this, we have kernel tricks. Basically, the idea is to convert the data to a higher dimension, (n -dimension \rightarrow $n+1$ dimension).

→ What is the advantage of increasing dimension?

If we take above example, there is no way to apply SVM, we can't get a decision boundary line (as the points are clustered together).

(The data is linearly separable in the next slide)

Now, imagine it is converted to 3D dimension, then for sure we can get a layer or separation between the data points.



① How do we increase the dimension?

→ say we have a 2D dataset as (n, y) .

Now if we want to add another dimension, we can simply use some mathematical function to find z .

Eg $z = n + y$, $z = n - y$, $z = n^2 + y^2$, $z = n^2 - y^2$ etc.

→ so we can project the points in a 3D space and then derive a plane which divides the data into two parts. In theory, that's what a Kernel function does without computing additional co-ordinates for the higher dimension.

② Types of Kernels In sum

→ there are many kernels, and here are a few popular ones.

1) Linear kernel → compute dot product between two feature vectors.

$$K(n, y) = n^T y$$

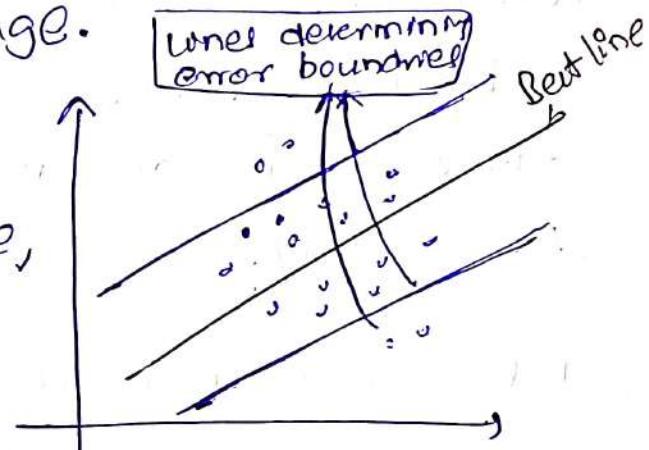
2) Polynomial kernel → $K(n, y) = (n^T y + c)^d$

3) Radial Basis Function (RBF) kernel →
RBF kernel measures the similarity between two samples using a Gaussian radial function.

$$K(n, y) = \exp(-\gamma \|n - y\|^2)$$

⑥ support vector Regressor

- we know how Linear Regression works, where we determine the best fit line. And in LR, the idea is to create a line which minimizes the total residual error.
- In SVR, the approach is a bit different. Here, instead of trying to minimize the error, SVR focuses on keeping the error in a fixed range.
- Here, the middle line is the best fit regressor line, and the other two lines are the bounding ones which denote the range of error.
- The Best-fit line/Hyperplane, will be the one which goes through the maximum number of data points and the error boundaries are chosen to ensure maximum inclusion.
- So, we simply bring majority of data points in between the margin lines.



Naive Bayes Classifier

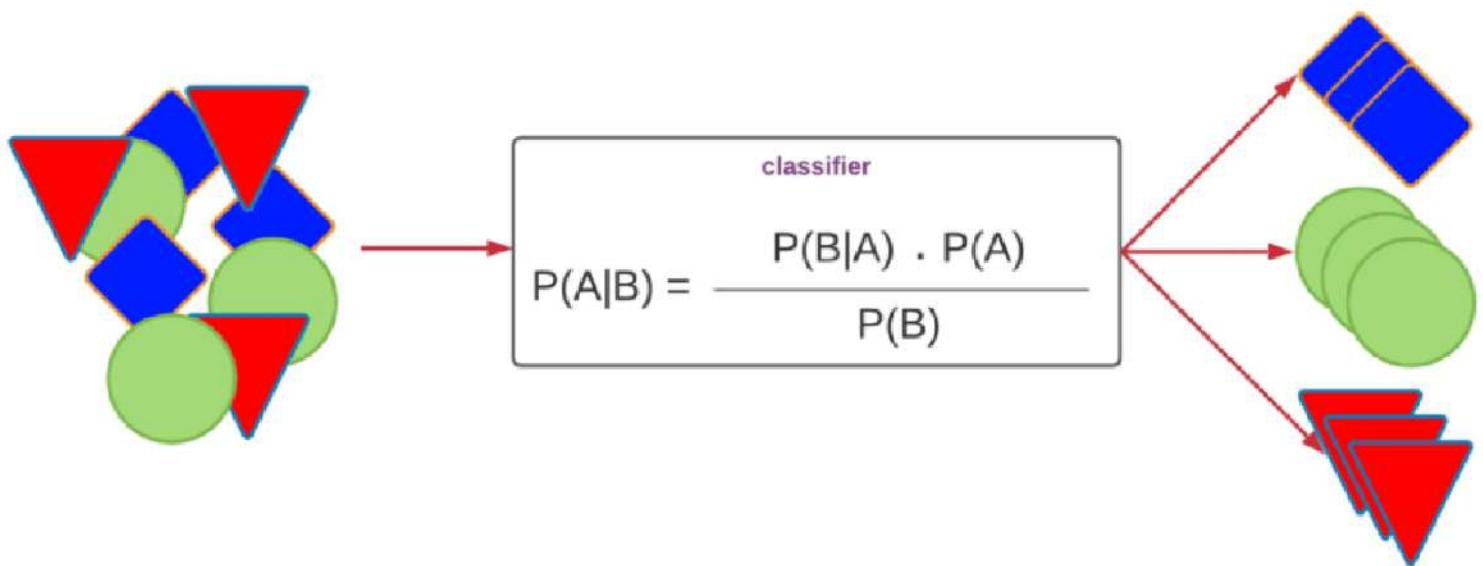


Table Of Contents

1. Why do we need Naive Bayes
2. Concept of how it works
3. Mathematical Intuition of Naive Bayes
4. Solving an Example on Naive Bayes
5. Other Bayes Classifiers
 - Gaussian Naive Bayes Classifier
 - Multinomial Naive Bayes Classifier
 - Bernoulli Naive Bayes Classifier

① Naive Bayes classifier

→ If we have a dataset, where relationship with one single feature may effect the output label, then we cannot generalize about all the features at a time!!

Ex:-

- 1) The food I ordered is good.
- 2) The food I ordered is bad.

Here only the last word signifies about the sentiment or statement.

② This is where Naive Bayes helps!

→ Naive Bayes will not try to generalize complete features at a time (and instead tries to build a independent relationship each & every input feature with output label, so that we get a better idea.)

③ How it works?

→ Naive Bayes is not a single algorithm, but a family of algorithms where all of them share a common principle, i.e; they are based on Bayes.

→ Bayes' theorem finds the probability of an event occurring given the probability of another event that has already occurred.

mathematically,

$$P(A|B) = \frac{P(B|A) \cdot P(A)}{P(B)}$$

where A & B are events & $P(B) \neq 0$

→ Basically we are trying to find the Probability of A, given the event B is true. Event B is also termed as, evidence.

→ $P(A)$ is the priori of A, that is the probability of event before evidence is seen (prior Probability).

→ $P(A|B)$ is the posterior probability, that is the probability of event after evidence is seen. Here, event B occurrence

Ex: Let's take a dataset of Playing golf based on climate.

	outlook	temperature	humidity	windy	play golf
0	Rainy	Hot	High	False	No
1	Rainy	Hot	High	True	No
2	overcast	Hot	High	False	Yes
3	Sunny	Mild	High	False	Yes
4	Sunny	Cool	Normal	False	Yes
5	Sunny	Cool	Normal	True	No
6	overcast	Cool	Normal	True	Yes
7	Rainy	Mild	High	False	No
8	Rainy	Cool	Normal	False	Yes
9	Sunny	Mild	Normal	False	Yes
10	Rainy	Mild	Normal	True	Yes
11	overcast	Mild	High	True	Yes
12	overcast	Hot	Normal	False	Yes
13	Sunny	Mild	High	True	No

- Here we can say that if it's sunny, mostly likely person wouldn't go for game. Similar case if it's too windy / too rainy.
- so we can clearly see that a single feature will affect the label directly. Hence, we will build model using independent relationship of each feature with output label by leveraging Naive Bayes.

- But before we apply Naive Bayes, we should make sure it follows the assumptions.

★ ① The fundamental Naive Bayes Assumption is that each feature makes an independent, and equal contribution to the outcome.

- And we can clearly say no pair of features are dependent, eg - the temperature being hot has nothing to do with humidity / being rainy. Hence, features are independent.

- knowing only temperature & humidity alone can't predict the outcome, hence all features contribute equally to outcome. we have, Bayes' theorem as

$$P(y/x) = \frac{P(x/y) \cdot P(y)}{P(x)}$$

y - all variable
x - dependent feature.

→ And by taking naive assumption do bayes theorem, which is independence among the features, we get.

$$P(A, B) = P(A) \cdot P(B)$$

$$P(y | n_1, n_2, \dots, n_n) = \frac{P(n_1, n_2, \dots, n_n | y) \cdot P(y)}{P(n_1, n_2, \dots, n_n)}$$

Based on naive assumption, or independence,

$$= \frac{[P(n_1 | y) \cdot P(n_2 | y) \cdots P(n_n | y)] \cdot P(y)}{P(n_1) \cdot P(n_2) \cdots P(n_n)}$$

$$= \frac{P(y) \cdot \prod_{i=1}^n P(n_i | y)}{\prod_{i=1}^n P(n_i)}$$

constant for any class variable in a class.

Now, to make a classifier model, we find the probability of a given set of inputs for all possible values of the class variable y and pick up the output with maximum probability.

$$\therefore y = \operatorname{argmax}_y P(y) \prod_{i=1}^n P(n_i | y)$$

so, to solve this problem we need to calculate $P(y) \& P(n_i | y)$ \rightarrow class probability, conditional Probability.

so now we calculate $P(\text{sunny} | \text{yes})$, $P(\text{overcast} | \text{yes})$, $P(\text{rainy} | \text{yes})$, $P(\text{sunny} | \text{no})$, $P(\text{overcast} | \text{no})$, $P(\text{rainy} | \text{no})$, similarly for each feature. $[P(n_i | y_j)]$

we get,

outlook features

	Y	N	$P(Y)$	$P(N)$
Sunny	3	2	$3/9$	$2/5$
overcast	4	0	$4/9$	$0/5$
Rainy	2	3	$2/9$	$3/5$
Total	9	5	100%	100%

temperature feature

	Y	N	$P(Y)$	$P(N)$
Hot	2	2	$2/9$	$2/5$
mild	4	2	$4/9$	$2/5$
cool	3	1	$3/9$	$1/5$
Total	9	5	100%	100%

humidity feature

	Y	N	$P(Y)$	$P(N)$
High	3	4	$3/9$	$4/5$
Normal	6	1	$6/9$	$1/5$
Total	9	5	100%	100%

wind feature

	Y	N	$P(Y)$	$P(N)$
false	6	2	$6/9$	$2/5$
true	3	3	$3/9$	$3/5$
Total	9	5	100%	100%

play $P(Y|yes) / P(Y|no)$

yes	9	$9/14$
No	5	$5/14$
Total	14	100%

→ Here we calculated $P(n^i | y_j)$ for each $n^i \in Y_j$, and also $P(y_j)$ for each class variable y_j .

Prediction

→ now, if we have a new datapoint as

(Sunny, Hot, Normal, False) = today climate.

$$P(Y | \text{today climate}) = P(\text{Sunny} | \text{Yes}) \cdot P(\text{Hot} | \text{Yes}) \cdot$$

$$\frac{P(\text{Normal} | \text{Yes}) \cdot P(\text{False} | \text{Yes}) \cdot P(\text{Yes})}{P(\text{today})}$$

$$P(\text{Yes} | \text{today})$$

$$\propto \frac{3}{9} \cdot \frac{2}{9} \cdot \frac{6}{9} \cdot \frac{6}{9} \cdot \frac{9}{14} \approx 0.0211$$

$$P(N \mid \text{today}) \propto \frac{2}{5} \cdot \frac{2}{5} \cdot \frac{1}{5} \cdot \frac{2}{5} \cdot \frac{5}{11} \approx 0.1142$$

However, we know,

$$P(Y \mid \text{today}) + P(N \mid \text{today}) = 1$$

APPLY reciprocal,

$$1 = \frac{1}{P(Y \mid \text{today}) + P(N \mid \text{today})} \quad \text{multiply } P(Y \mid \text{today})$$

$$P(Y \mid \text{today}) = \frac{P(Y \mid \text{today})}{P(Y \mid \text{today}) + P(N \mid \text{today})}$$

$$P(Y \mid \text{today}) = \frac{0.024}{0.0211 + 0.1142} = \frac{0.024}{0.1353} \approx 0.155$$

$$P(N \mid \text{today}) = \frac{0.1142}{0.0211 + 0.1142} = \frac{0.1142}{0.1353} \approx 0.844$$

So, clearly

$$\boxed{P(N \mid \text{today}) > P(Y \mid \text{today})}$$

Hence, "play golf \rightarrow No" for today's climate.

→ However, this above method is for the categorical data.

→ And in case of continuous data we need to make assumptions regarding the distribution of values of each feature.

★ The different Naive Bayes classifier differ mainly by the assumptions we make regarding the distribution of $P(x \mid y)$

① Gaussian Naïve Bayes classifier

- If the feature variables are continuous variables, then we cannot use the Naïve Bayes, and instead we should be using Gaussian Naïve Bayes Classifier.
- In Gaussian Naïve Bayes, continuous values associated with each feature are assumed to be distributed according to a Gaussian distribution, also called "Normal distribution".
- The likelihood of the features is assumed to be Gaussian, hence conditional probability is,

$$P(x_i | y) = \frac{1}{\sqrt{2\pi\sigma_y^2}} \exp\left(-\frac{(x_i - \mu_y)^2}{2\sigma_y^2}\right)$$

Other popular Naïve Bayes classifier are -

- ① multinomial Naïve Bayes: feature vector represent the frequencies with which certain events have been generated by multinomial distribution. This is typically used for document classification.
- ② Bernoulli Naïve Bayes: Assumes that the features are binary or categorical. It is particularly useful for document classification like mnb, and where binary term occurrence features are used rather than term features.

Types of Clustering Algorithms

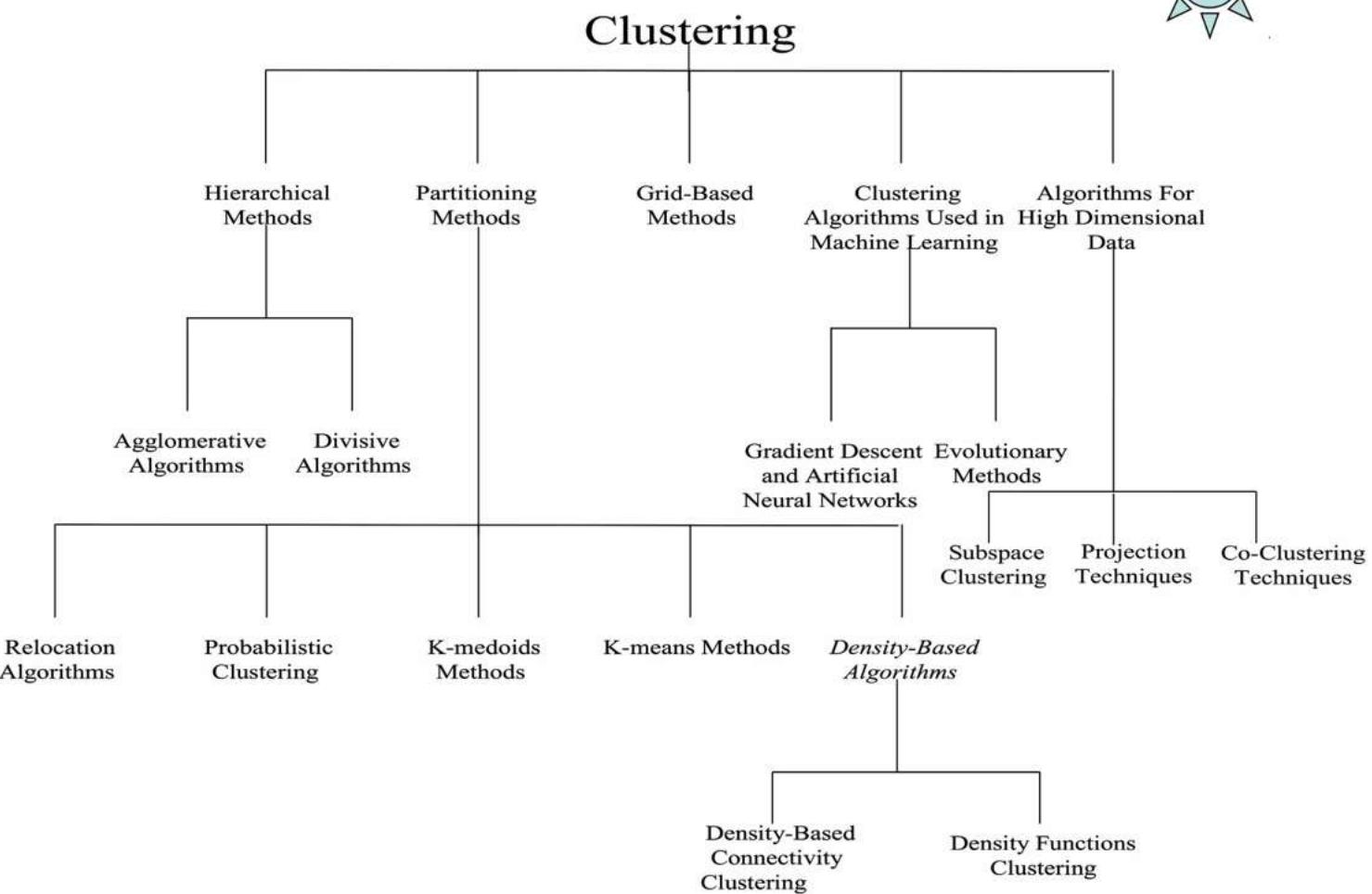
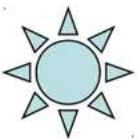


Table of Contents

1. How clustering is different from classification
2. Applications of Clustering
3. What are density based methods
4. What are Hierarchical based methods
5. What are partitioning methods
6. What are Grid Based methods
7. Main Requirements for Clustering Algorithms

Unsupervised Learning Algorithms. (pg 100 & 5 References)

- i) k-Means
- ii) Hierarchical clustering
- iii) DBSCAN
- iv) Performance measurement
- v) Principal component Analysis
- vi) Dimensionality Reduction.

} clustering -

Evaluation
metrics
for
clustering
are
Important



clustering:-

→ First we need to understand clustering is different from classification.

Parameter	classification	clustering
1) Type	used for supervised learning.	used for unsupervised learning.
2) Basic	Process of classifying the input instances based on their corresponding "labels".	Grouping the instances based on their <u>similarity</u> <u>without</u> the help of <u>class labels</u> .
3) complexity	more complex	less complex
4) Example Algo	Logistic Regression, naive Bayes, SVM etc	K-means clustering, Hierarchical clustering, DBSCAN, BIRCH Algorithm, etc

Application of clustering

→ for customer segmentation

You can cluster your customer based on their purchase, their activity on website and so on. And this can be used in recommendation system to suggest content that other user in same cluster enjoyed (Chirag)

Eg:- marketing, insurance, librairie etc.

→ Earthquake prediction

By learning old data it can make cluster and determine dangerous zone.

→ for search engines

As you search, similar image would end up on same cluster.

→ to segment an image

By clustering pixels according to their color, then replacing each pixel's color with the mean color of its cluster, it is possible to reduce no. of different colors in image.

Eg:- used in object detection & tracking system.

→ for Anomaly detection (outlier)

→ As a dimensionality Reduction Technique.

→ For Data Analysis

clustering methods

① density-Based methods

these methods consider the clusters as dense region having some similarity and different from the lower dense region or Spurious.

These methods have good accuracy & ability to merge two clusters.

ESL

i) DBSCAN [Density Based Spatial Clustering of Applications with Noise]

ii) OPTICS [Ordering Points to Identify Clustering Structure.]

② Hierarchical-Based methods

The clusters formed in this method forms a tree-type structure based on the hierarchy. New clusters are formed using the previously formed one. It has two categories

i) Agglomerative

Bottom up approach -

form a single cluster & expand

ii) divisive

Top down approach -

form a big cluster & divide

iii) CURE

[Clustering Using Representatives]

iv) BIRCH

[Balanced Iterative Reducing

Clustering and using Hierarchical]

③ Partitioning methods-

. these methods partition the objects into k -clusters and each partition forms one cluster. this method is used to optimize an objective criterion similarity function such as when the distance is a major parameter.

Eg:-

i) K-Means

ii) CLARANS

[clustering Large Applications based upon Randomized Search]

④ Grid-Based methods

In this method the data space is formulated into a finite no. of cells that form a grid like structure. All the clustering operation done on these grids are fast and independent of the no. of data objects.

Eg:-

i) STING

[Statistical Information Grid]

ii) CLIQUE

[Clustering In QUEST]

iii) Wave, cluster, etc - all grid methods

Main Requirements for clustering Algorithms

- ① It should be Scalable.
- ② It should be able to deal with attributes of different types.
- ③ It should be able to discover arbitrary shape clusters.
- ④ It should have an inbuilt ability to deal with noise & outliers.
- ⑤ The clusters should not vary with the order of input records.
- ⑥ It should be able to handle data of high dimension.
- ⑦ It should be easy to interpret & use.