

Machine Learning

=
DSC 7456

Unsupervised Learning

* Association Rules

* Clustering (Finding similar groups among data points)
(Row similarity)

↳ Hierarchical clustering

↳ DBScan

↳ K-means

* Principal Component Analysis.
(PCA)

Supervised Learning

Regression

* Linear Regress

Classification

* Logistic Reg

* Naive Bayes

* Decision Trees.

* K-NN.

* SVM

* Random

* Forest
Time Series

* Gradient Decent Method → Optimizer
Algorithm

* Singular value Decomposition → Matrix Factorization
(collaborative filtering)

K

Association Rules

Affinity Analysis -

How likely the person would buy 'Y' object if he buys 'X' object?

- Note :-
- Main thing is identifying the important association rules
 - Find the goodness of the A-Rules based on business value.

Basic Definitions :-

(*) Support $\{B\} = \frac{4}{8} = \frac{\text{no. of times bought}}{\text{Total Trans}}$

Transaction	T1	T2	T3	T4	T5	T6	T7	T8
B	✓	✓	✓	✓	✓	✓	✓	✓
D	✓	✓	✓	✓	✓	✓	✓	✓
P	✓	✓	✓	✓	✓	✓	✓	✓
R	✓	✓	✓	✓	✓	✓	✓	✓

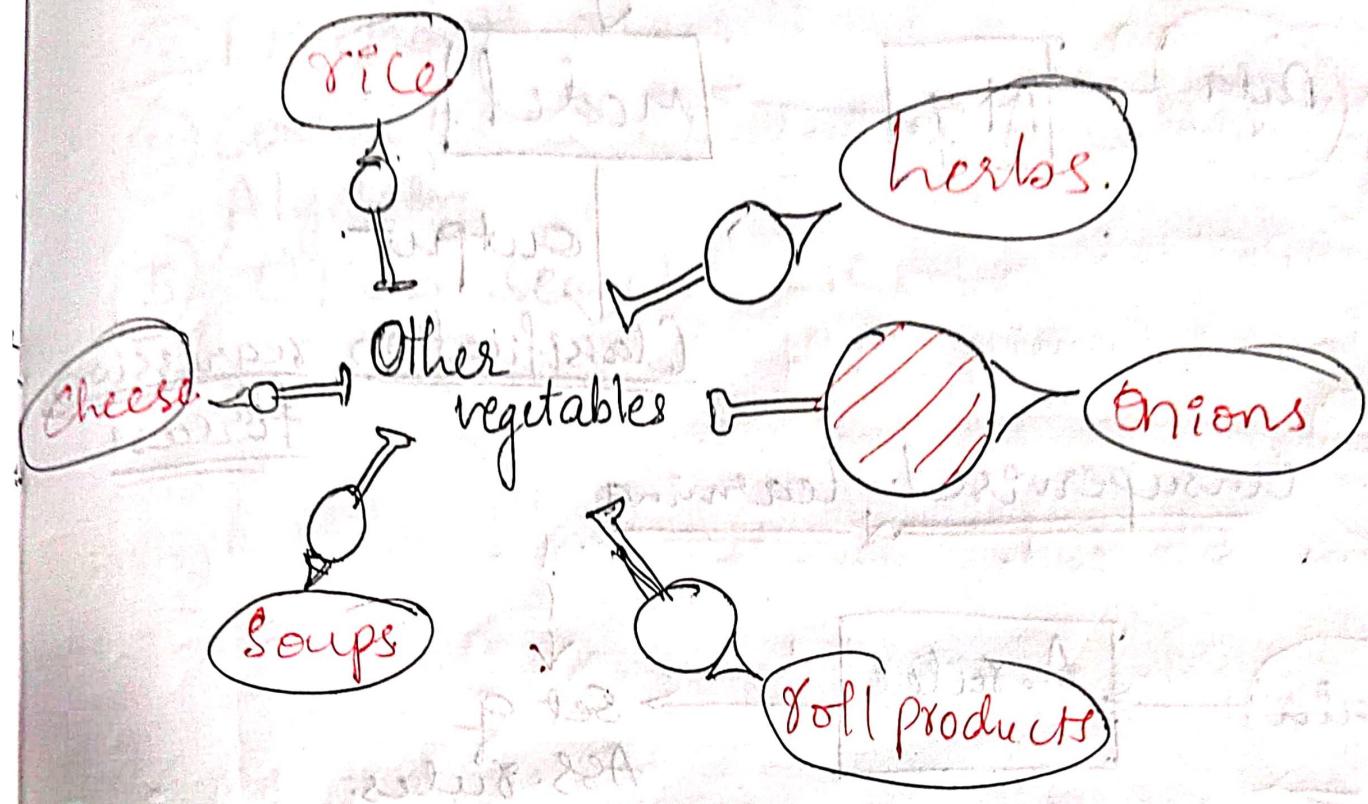
(*) Confidence $\{B \rightarrow D\} = \frac{\text{Support}\{B, D\}}{\text{Support}\{B\}}$

(*) Lift $\{B \rightarrow D\} = \frac{\text{Support}\{B, D\}}{\text{Support}\{B\} \times \text{Support}\{D\}}$

Recall,
A measure of importance of the rule

(*) The above metrics are used to find the goodness of the A-Rules

Simple representation of all the 3 metrics of association rules.



- ① Larger the circle \rightarrow higher Support
- ② Red circles \rightarrow higher lift

Point of Sale Data.

POS data is the data of customers along with prod they bought in any marts, shops, etc.

Steps

- ① Item Sets (sets of 2, 3, etc) which are relevant
- ② Ass. rules for relevant Item Set
- ③ Find Support, lift, confidence of all the relevant item set
- ④ Keep the best one's (with good supp, lift, conf)

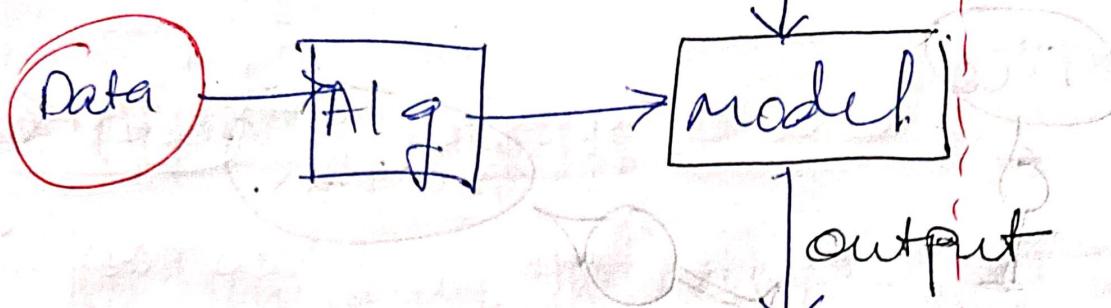
- ⑤ From these Ass. rules we cannot do classification but we can merchandise; bundling, pricing, merchantize insights, etc

~~Distortion S. will be P. minimization for alg. m. 2.~~

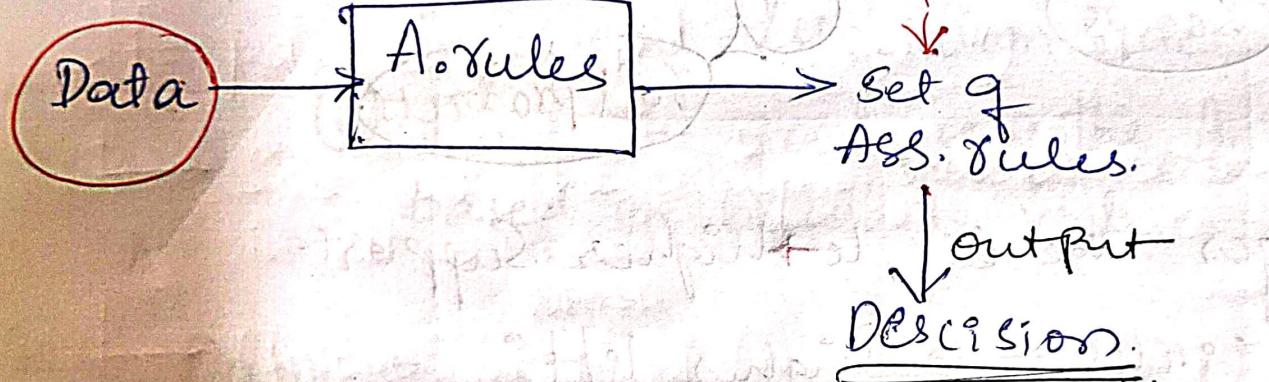


Supervised Learning

New Data

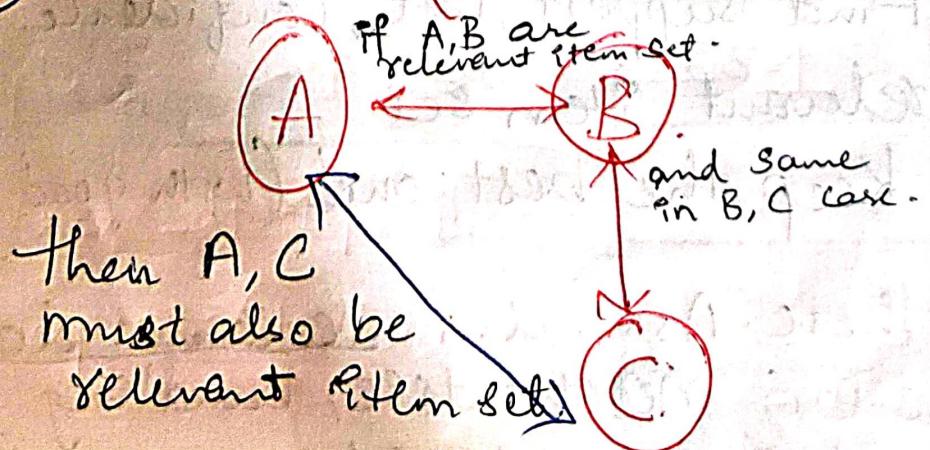


Unsupervised learning



Since hardest part is finding the relevant item set.: there are some algorithms to get them.

① APRIORI (Downward Closure)



Distance

Association Rules in Python

mlxtend → Package is used.

import os

import numpy as np

import pandas as pd.

from mlxtend.frequent_patterns import apriori,

df = pd.read_excel("xxxxx.xls")

association rules

→ See all the unique values using .unique
in all columns

which
→ Choose in level (ex → subcat/category) the ass. rule
must be built

xxxxx_df = pd.crosstab(df['orderID'], df[['sub-cat']])

→ Build crosstab(matrix) of orderID and all
the attributes in chosen level:

→ Replace >1 terms := 1.

xxxxx_df[xxxxx_df > 1] = 1.

itemsets = apriori(transaction_df, min_support=0.01,
use_colnames=True)

→ use apriori to make relevant items

rules = association_rules(itemsets, metric='lift', min_threshold=1.1)

sorted_rules = rules.sort_values(['support', 'confidence', 'lift'], ascending=[True, False, False])

→ Check for the highest confidence value itemsets and make decisions from those rules.

c) sorted_rules.to_csv('rules_extracted.csv')

Distance :- (Euclidean, Manhattan)

(Cosine Similarity), Pearson and Distance

Distan

* The distance b/w two data attributes is used to find the similarity among the data.

↳ distance is very helpful in clustering the datasets

↳ The distance b/w the data attributes is measured and weighted summation of all attributes is taken to know the overall distance (similarity) among the datasets.

* If d_1 is near to d_2 , and d_2 near to d_3 , then d_1 is not far from d_3 .

* Practical data has different types of data

- ① Numeric → Euclidean → Assumption attribut
- ② Interval-scaled. → Manhattan → attributes are L
- ③ Symmetric binary
- ④ Asymmetric binary
- ⑤ Ratio-scaled.
- ⑥ Ordinal and nominal. → asymmetric
- ⑦ Binary → symmetric

Euclidean

$$\sqrt{(x_2 - x_1)^2 + (y_2 - y_1)^2}$$

Manhattan

$$|x_2 - x_1| + |y_2 - y_1|$$

Some other common metrics

① Weighted distances

② Minkowski distance

$$\rightarrow \sqrt[n]{(x_2 - x_1)^n + (y_2 - y_1)^n}$$

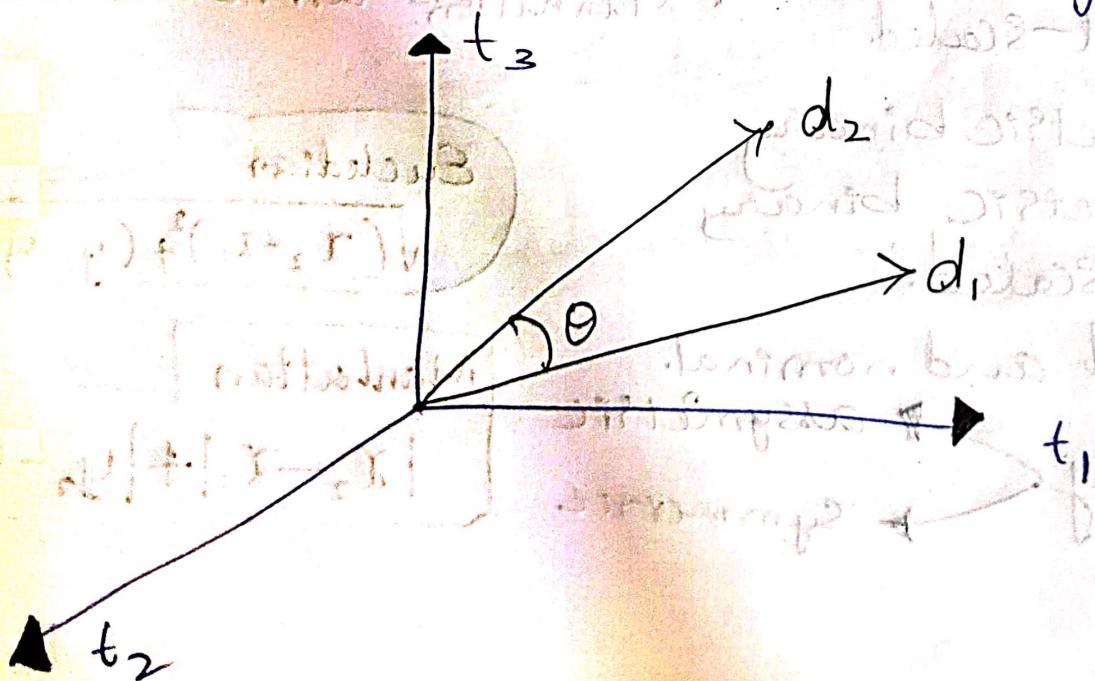
③ The maximum distance amongst all attributes.

④ Co-relation b/w rows.

* Distance in higher dimensions become meaningless.

→ Cosine Similarity

* Distance b/w vectors d_1 & d_2 captured by the cosine of the angle b/w them



$$\text{sim}(d_j, d_k) = \frac{\vec{d}_j \cdot \vec{d}_k}{\|\vec{d}_j\| \|\vec{d}_k\|} = \frac{\sum_{i=1}^n w_{ij} \cdot w_{ik}}{\sqrt{\sum_{i=1}^n w_{ij}^2} \cdot \sqrt{\sum_{i=1}^n w_{ik}^2}}$$

- ④ Cosine of angle b/w two vectors
- ④ The denominator involves the length of the vectors
- ④ The cosine measure is also known as the normalised inner product.

④ For Categorical Attributes in unsupervised settings:

Data :

Data ;	1	0
	a	b
0	c	d

Hamming Distance = $\frac{\text{no. of dissimilar attributes}}{\text{no. of similar attributes}}$

$$\text{Hamming Distance} = \frac{b + c}{b + c + a + d}$$

$$\text{Hamming Distance} = \frac{b + c}{b + c + a + d}$$

Binary

$$\begin{aligned} \text{Asymmetric} &= \frac{b + c}{a + b + c} \\ \text{Symmetric} &= \frac{b + c}{b + c + a + d} \end{aligned}$$

Pearson Linear Correlation

$$P(x, y) = \frac{\sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})}{\sqrt{\sum_{i=1}^n (x_i - \bar{x})^2} \sqrt{\sum_{i=1}^n (y_i - \bar{y})^2}}$$

* Pearson Linear Correlation (PLC) is a measure that is invariant to scaling and shifting of attribute vectors.

* The range is b/w $-1 \leq P \leq +1$

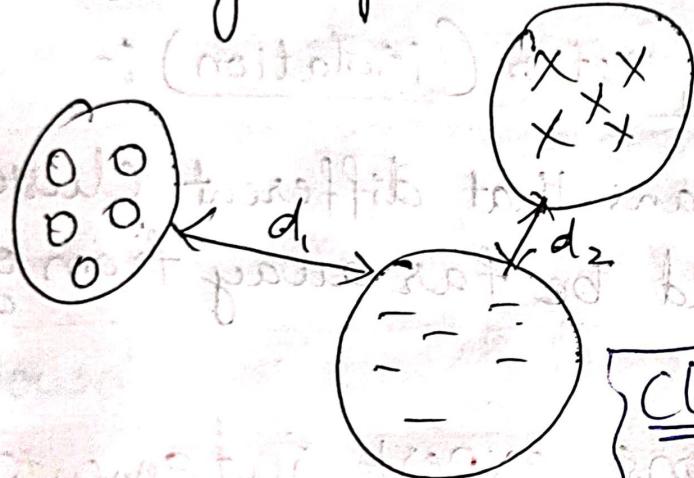
* The Pearson L. Corr is a similarity measure we can make it as dissimilarity measure

$$d_p = \frac{1 - P(x, y)}{2}$$

* If we measure distance by any method, finally we must combine all distance to get a finalized distance by weighted sum.

Clustering

- * In general a grouping of objects such that the objects in a group (cluster) are similar and different from (or unrelated to) the object in other group.



* This is the task (i.e. Clustering) that must be done in Some Space

Cluster → Similar locally, dissimilar globally

Data input for clustering can be in two forms.

① Data matrix x →

$$\begin{bmatrix} x_{11} & x_{12} & \dots & x_{1p} \\ x_{21} & x_{22} & \dots & x_{2p} \\ \vdots & \vdots & \ddots & \vdots \\ x_{n1} & x_{n2} & \dots & x_{np} \end{bmatrix}$$

It consists of the values of all unknown attributes (columns)

② Dissimilarity (Distance matrix) matrix. →

$$\begin{bmatrix} 0 & d(2,1) & \dots & d(n,1) \\ d(2,1) & 0 & \dots & d(n,2) \\ \vdots & \vdots & \ddots & \vdots \\ d(n,1) & d(n,2) & \dots & 0 \end{bmatrix}$$

$\frac{n(n-1)}{2} \rightarrow \text{input}$

Quality of Cluster :- (To find how good is the Cluster)

① Intra-cluster Cohesion (compactness) :-

(*) It measures how near the data points in a cluster are to the cluster centroid.

② Inter-cluster Separation (isolation) :-

(*) Separation means that different cluster centroids should be far away from one another.

(*) In most applications, expert judgments are still the key. (expert judgments involves the type of the cluster, and its compactness)

① Hierarchical Clustering :-

→ There are two main types of hierarchical clustering :-

① Agglomerative Clustering :-

Bottom-up
most preferred

→ Start with the points as individual clusters.

→ At each step, merge the closest pair of clusters until only one cluster is left.

② Divisive :-

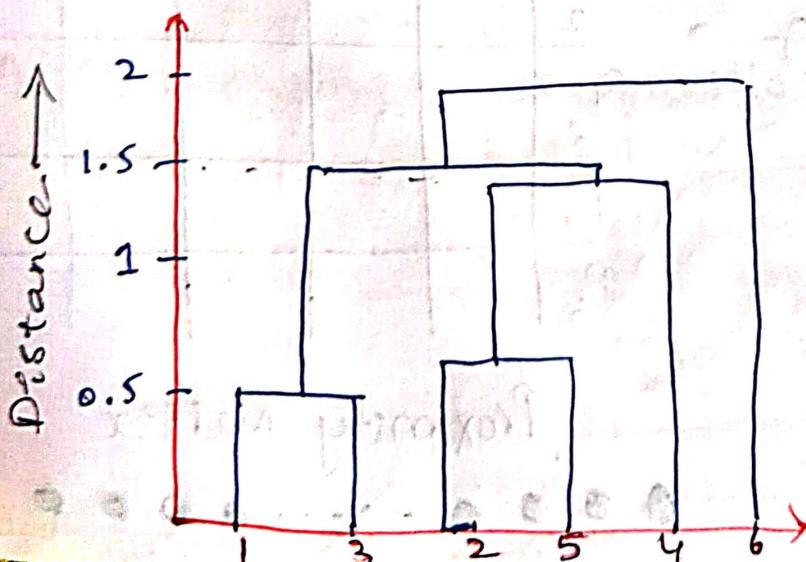
Top down

→ Start with one, all-inclusive cluster.

→ At each step, split a cluster until each cluster contains a point.

→ Produces a set of nested clusters organised as a hierarchical tree.

→ Can be visualized as a dendrogram.



Critical The distance for clustering of the data points into a cluster must be decided according to the business problem.

Agglomerative clustering algorithm.

- ① Compute the proximity matrix (distance matrix)
- ② Let each data point be a cluster
- ③ Repeat
 - Merge the two closest cluster
 - update the proximity matrix
- ④ Until only a single cluster remains.

→ Key operation is the computation of the proximity/distance of two clusters

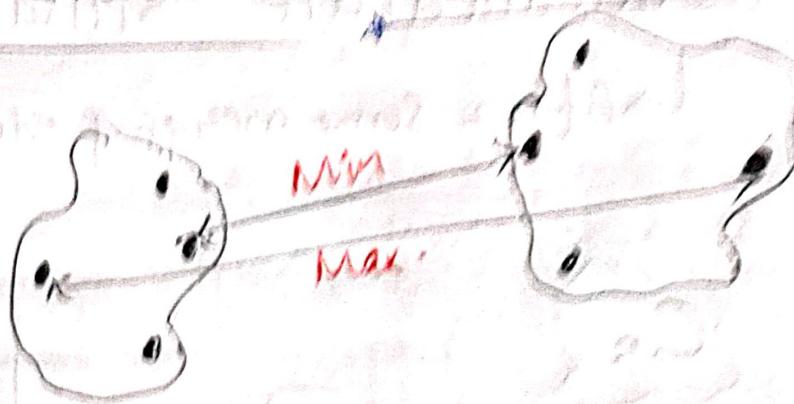
Example:-

→ Starting Situation

→ Start with clusters of individual points and a

How to Define Inter-cluster Similarity

→ MIN



→ MAX.

→ Group Average

→ Distance Between Centroids.

→ Other methods driven by an Objective function

↳ Ward's Method uses squared error

Min

→ Single Link, will merge two clusters when a single pass of elements is linked.

* Min works very well on non-elliptical shapes

* Min is sensitive to noise and outliers.

Max

→ Complete linkage will merge two clusters when all pairs of elements have been linked.

* Less susceptible to noise and outliers

* Tends to break large clusters.

* Biased towards globular clusters.

Group Average

* Compromise b/w Single and Complete link.

* Less susceptible to noise and outliers

* Biased towards globular clusters

Ward's Method

→ Similarity of two clusters is based on the increase in SSE, when two clusters are merged.

* Less susceptible to noise and outliers

* Biased towards globular clusters

K-means
... used to initialize

Problems with Hierarchical clustering

- ① Computational Complexity in time and space
- ② Once a decision is made to combine two clusters, it cannot be undone.
- ③ No objective function is directly minimized.

⑥ DBScan:-

→ In density based clustering we partition points into dense regions separated by not-so-dense regions.

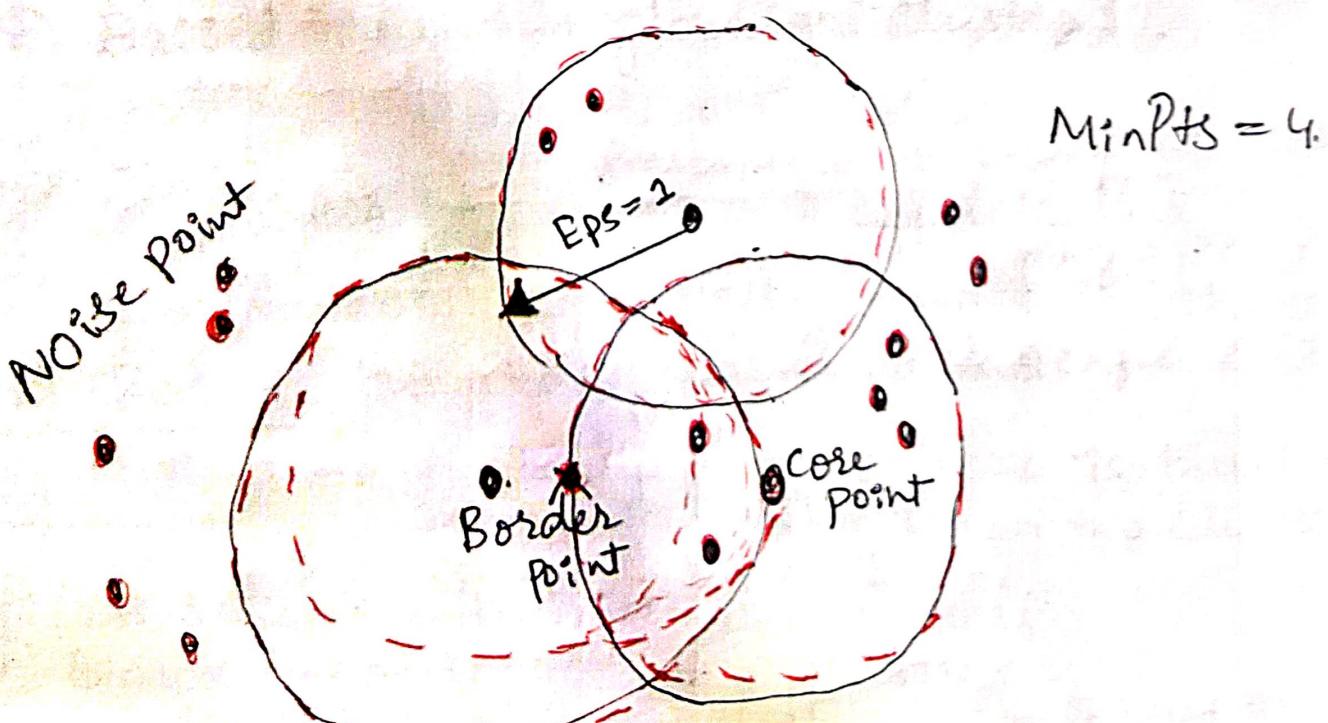
Dense Region → A Circle of radius EPS that contains at least Minpts points

* Characterization of Points :-

→ A point is a **Core point** if it has more than a specified number of points (Minpts) within EPS (epsilon)

→ A **Border point** has fewer than Minpts points within EPS, but is in the neighborhood of a core point

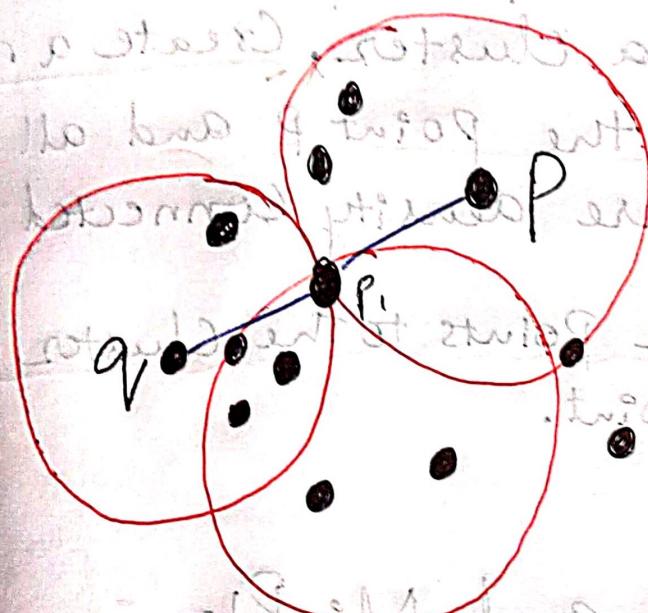
→ A **Noise point** is any point that is not a core point or a border point.



Density - Connected Points :-

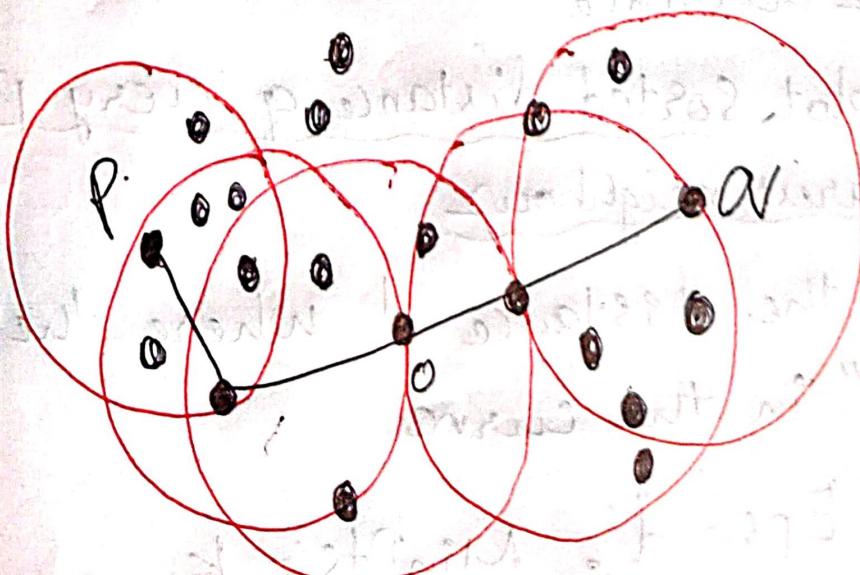
Density Edge :-

→ we place an edge b/w two core points q & P if they are within distance Eps .



Density - Connected :-

→ A point P is density-connected to q if there is a path of edges from P to q .

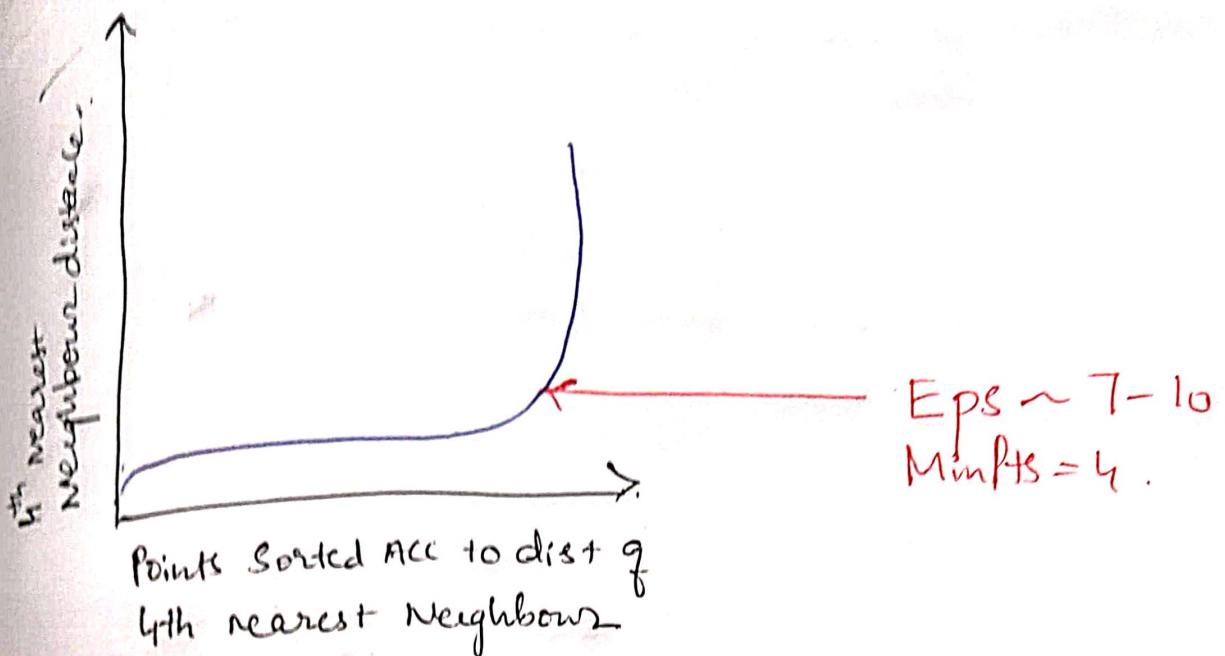


DBSCAN Algorithm

- ① Label points as core, border and noise
- ② Eliminate noise points
- ③ For every ~~as~~ core point p that has not been assigned to a cluster, Create a new cluster with the point p and all the points that are density connected to p .
- ④ Assign border points to the cluster of the closest core point.

Determining the Eps and MinPts :-

- * Idea is that for points in a cluster, their k^{th} nearest neighbours are at roughly the same distance.
- * Noise points have k^{th} nearest neighbour at farther distance
- * So, plot sorted distance of every point to its k^{th} nearest neighbour
- * Find the distance d where there is a "knee" in the curve
- * $\text{Eps} = d$, $\text{MinPts} = k$.



Advantages :-

- ① Resistant to Noise
- ② Can handle clusters of diff shapes & sizes

Disadvantages :-

- ③ does not work well for varying densities
- ④ " High-dimensional data

Partitioning Algorithms \rightarrow K-Means

Takes the whole data & globally cuts into cluster

→ K-means is a ~~partitioning~~ clustering algorithm as it partitions the given data into k-clusters.

→ Each cluster has a center, called centroid

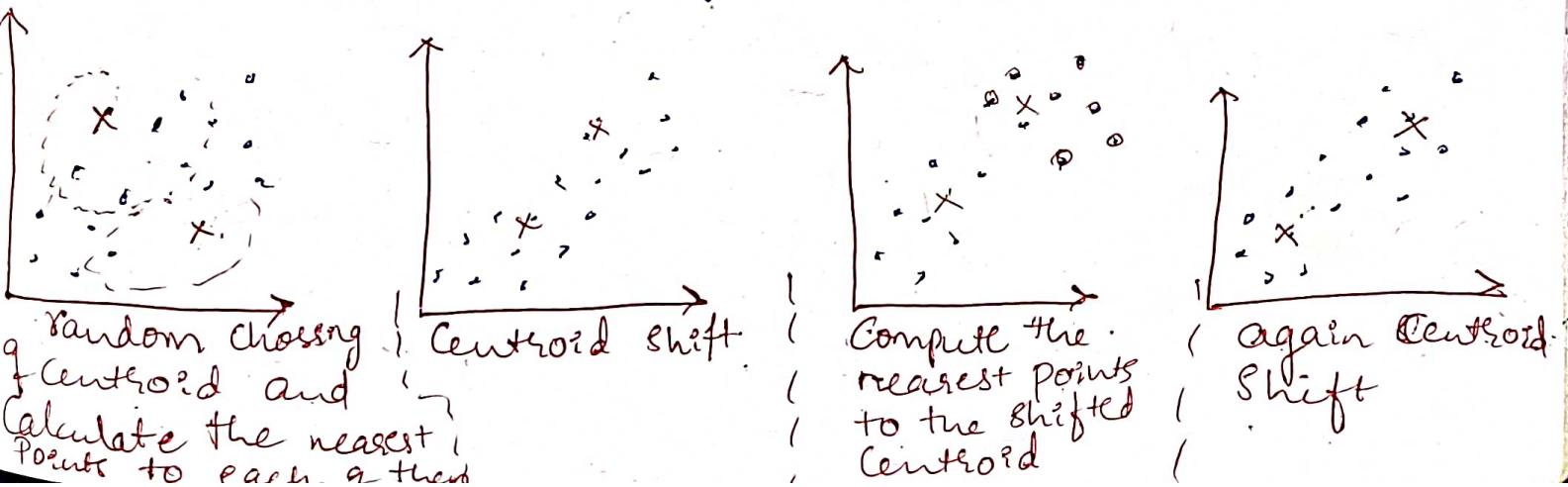
→ k is specified by the user.

K-Means Algorithm:-

→ For a given k , the K-means algo works as follows.

Initial assumption
using expert advise.

- ① Randomly choose k data points (seeds) to be the initial centroids, cluster centers.
- ② Assign each data point to the closest centroid.
- ③ Re-compute the Centroids using the current cluster membership
- ④ If a convergence criterion is not met, or if some clusters dont get any points, go to ②.



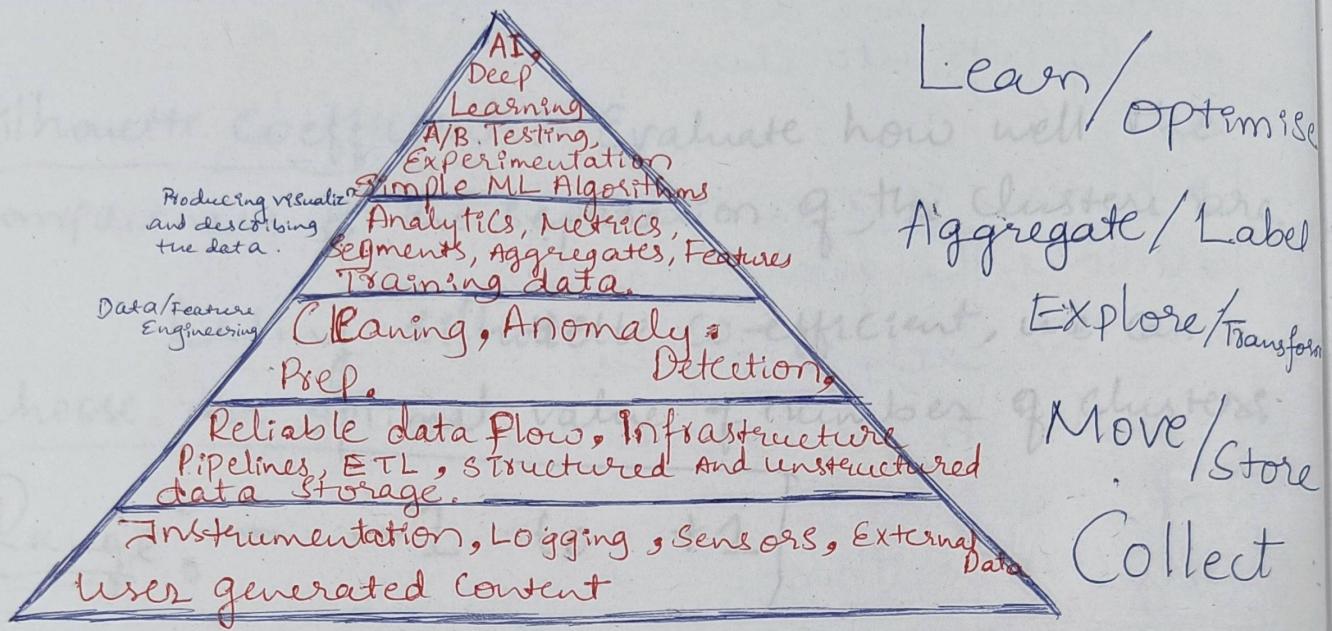
Silhouette Score

Silhouette coefficient :- Evaluate how well the compactness & the separation of the clusters are.

Using silhouette co-efficient, we can choose an optimal value of number of clusters.

$$\text{Range} = -1 \text{ to } +1$$

Supervised learning



Shot on OnePlus

By Krishna

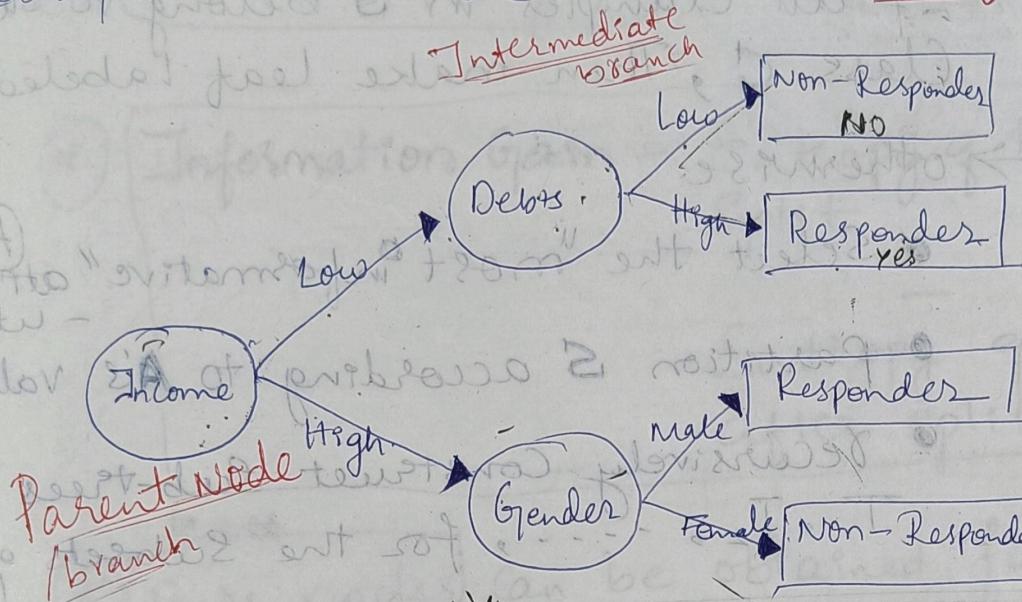
3 Decision Trees

Attribute / Feature → Independent Variable
Label → Dependent variable

Sample Decision Tree

Leaf

Loan



In a D-T every leaf corresponds to a unique rule

* Labels | target variable are always at the leaf (leaf is homogenous subset)

* D-T can be used for Classification & Regression problems.

→ TDIDT → Top-down Induction of D-T.

Recursive partitioning.

→ ID3 (Quinlan) → Iterative Dichotomizer

→ C.A.R.T → Classification And Regression Tree.

→ ASSISTANT

→ C4.5 (Quinlan)

→ See5/C5.0 (Quinlan)

* Small D-Trees are always better.

Shot on OnePlus

By Krishna

T DIDT Algorithm :-

↳ To Construct decision tree T from Learning Set S :-

→ If all examples in S belong to same class C , Then make Leaf Labeled C

→ otherwise

• Select the "most informative" attribute A

• Partition S according to A 's value

• Recursively Construct Sub-trees

T_1, T_2, \dots , for the Subsets of S

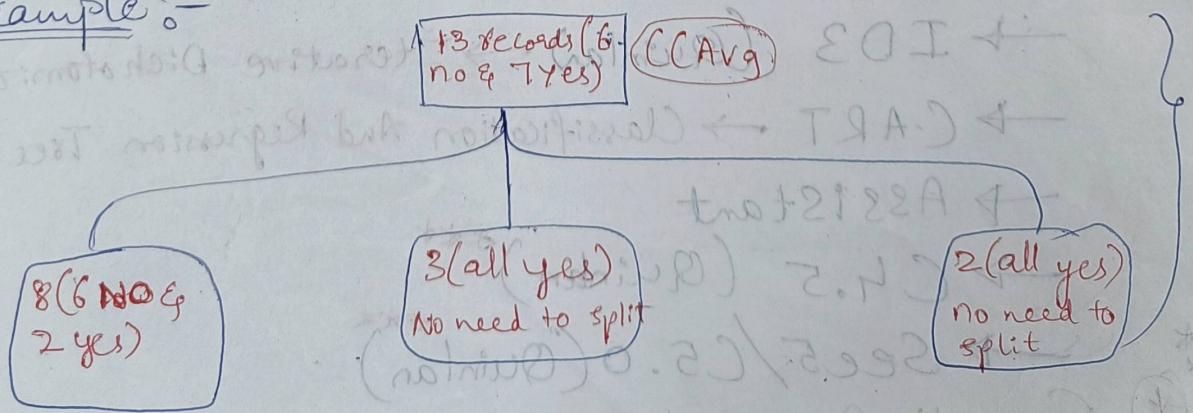
During Constructing a Decision Tree there will be following questions

• What types of decisions

→ Which attribute to choose for split

→ When to stop (to avoid overfitting)

Example :-



(i) Information Gain

Entropy and Information Gain (ID3)

$$H = -\sum_i P_i \log_2 P_i$$

P \rightarrow Class 1
Probability Total

* Information Gain = Entropy of the System before Split - Entropy of the System after split.

* Pure nodes can be obtained faster by minimizing entropy at every split or maximize the Information gain at every split.

* For an impure class, entropy will be zero and entropy will be max(1) when both the classes in the tables are equal.

Entropy before split

$$H = -\frac{6}{13} \times \log_2 \frac{6}{13} - \frac{7}{13} \times \log_2 \frac{7}{13} = 0.9957$$

Entropy after split on CC AVG

$$H = \frac{8}{13} \left(\frac{-6}{8} \times \log_2 \frac{6}{8} - \frac{2}{8} \times \log_2 \frac{2}{8} \right) + \frac{3}{13} \left(\frac{-3}{3} \times \log_2 \frac{3}{3} \right) + \frac{2}{13} \left(\frac{-2}{2} \times \log_2 \frac{2}{2} \right)$$

$$H = 0.4992$$

Shot on OnePlus

By Krishna Information gain = 0.9957 - 0.4992 = 0.4965

Advantages of Trees

- ★ Fast
- ★ Robust
- ★ Explainable (Explainable)
- ★ Requires very little experimentation
- ★ We can easily build some intuitions about the data.

Other measures of purity :-

→ Gini index

→ Information-theoretic

→ χ^2 (Chi-squared)

- ★ Some time Information gain will fail if attribute chosen to split has all unique values and results in leaf which is equal to number of unique value.

To prevent this problem we define a term called Information Content

⑪ Gain Ratio :-

Information Content :-

$$I.C = - \sum f_i \log f_i$$

$f_i \rightarrow$ fraction of the members in a state divided by the total members

I.C for I.D \rightarrow It has 13 states. So, the
 $I.C = \left(-\frac{1}{13} \times \log_2 \left(\frac{1}{13} \right) \right) \times 13 = 3.7$

Gain Ratio :- (used in C4.5)

* Since Information gain is biased towards attributes with many classes (levels)

Gain ratio normalizes information gain by dividing by the information content at the attribute

$$\text{Gain Ratio}(A) = \frac{\text{Gain}(A)}{\text{Information Content}(A)}$$

* Attribute with maximum gain ratio is selected as the splitting attribute.

③ Gini Index :- (Used in CART)

$$1 - \sum_{i=1}^m p_i^2$$

\rightarrow (attribute) most suitable

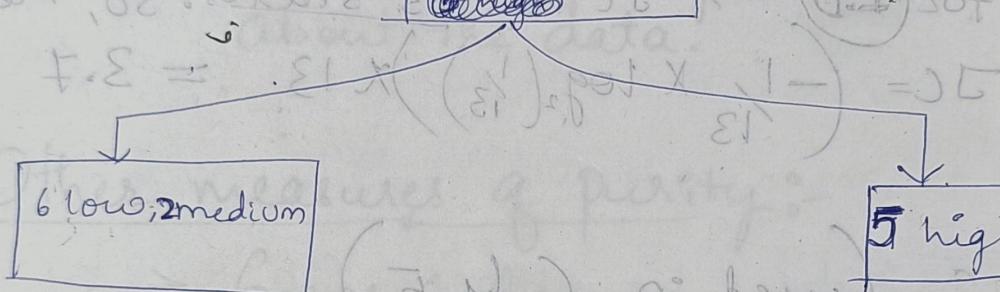
$m \rightarrow$ classes

$P \rightarrow$ Probability

* It is computed on binary splits only.



CC Avg	[6 low, 2 medium]
13 records	5 medium 5 high



{low, medium} $\xleftarrow{\text{Binary split}}$ {High}

* The one with least gini index is picked.

Eg:- Gini index before split = $1 - \left(\frac{6}{13}\right)^2 - \left(\frac{7}{13}\right)^2 = 0.497$

Gini index after split =

$$\frac{8}{13} \left(1 - \left(\frac{6}{8}\right)^2 - \left(\frac{2}{8}\right)^2\right) + \frac{5}{13} \left(1 - \left(\frac{5}{5}\right)^2\right) = 0.231$$

Regression Trees uses Variance for splitting

* Minimizing the variance

\rightarrow Variance before split

\rightarrow weighted average of variance after split

\rightarrow The attribute that reduces the variance most is chosen for the node.

$\uparrow \chi^2$

$\downarrow \text{Gini Index}$

$\uparrow \underline{\text{Info gain}}$

$\downarrow \underline{\text{Variance.}}$

$\uparrow \underline{\text{Gain Ratio}}$

Bayesian optimal classifier → highest accuracy a decision tree can have but it is hard to

④ D-T are unstable at leaf (^{small} change in leaf will ^{change a lot}) but stable at top

Pruning — ④ Avoid overfitting

④ The better method would be growing the entire tree and pruning the unwanted branches later.

④ After chopping, check the errors on training & test data to avoid overfitting. (Cost-Complexity Pruning)

④ The previous method of growing a tree is called Reduced Error Pruning

CART.

→ Binary Split

→ Cost Complexity Pruning

C 4.5

→ multi split

→ Gain ratio

→ Pessimistic Pruning

Shot on OnePlus

By Krishna

Some Notes on Decision Tree Induction.

(*) For decision trees the ~~decision~~ boundary is rectangle.

Tessellation

6 tons medium

process + rental

over no best materials

about 27 + tax

rental
leaving
costs

(the first approach) feel to old stores are T-0
(the second) the one with less cost
got to old + two

Gini index below 0.5
with three boxes

some enterprises produced better better added up
and reduced behaviour at growing bus and
3 products no work with profit: 10% A
Stock best (like Amazon). profit from boxes of
vertical books 3% with a pickup of better answer out
product good

K-Nearest Neighbour :-

(*) It is called lazy learner becoz KNN does not build any model, they just remember the ^{training} data and try to do the classification that we require.

(*) Lazy learner is large group of techniques, KNN is one among them

(*) Can trivially adapt to evolving data, i.e. since it ~~will~~ not built model, change in data can be handled and classify accordingly.

(*) Lot more computation must be done during classifying than that of training.

Eg:- If a customer came to market and wants to buy something, we can compare with the other customers in the data like income, age at what time he came, etc.

By keeping track of these data ~~and~~ and the buying pattern of those similar customer we can recommend the products they have bought earlier.

① In K-NN for every test sample we have to
 ~~to~~ search the nearest neighbours.

→ Another example:-

If in case a new case is filed in the high court and we want to choose the bench members(judge) to take a decision, we can see the similar case in past year and the judges who given judgement and according those judges are preferred for bench members.

② ~~Most~~ Suitable for problems where one cannot assume that training and test samples are drawn from the same distribution.

③ Selecting K is primary challenge.

Lazy learning (Fewer assumption about future is done compared to other algorithms)

Similar Test Sample / Local data

↓
Directly work with data

(Instance

Based Learning)

Case Based Reasoning

↓
Build a local model

Local reg/
local sum etc.

(*) Since we make fewer assumptions, KNN has twice the Bayesian optimal error rate.

(*) KNN Data must not contain NOISE

KNN - Algorithm

(*) Classification → Assign a class to a new data point based on its neighbours (mode)

(*) Regression → Identify a numeric value of a new data point based on its neighbours (mean/median)

(*) weighted mean/mode of entire data.

(*) Store all input data in the training set

For each instance in the test set

Search for the k-nearest instances to the input instance using appropriate distance measure

For classification, compute the confidence for each class as C_i/k

(*) The classification for the input instance is the class with the highest confidence.

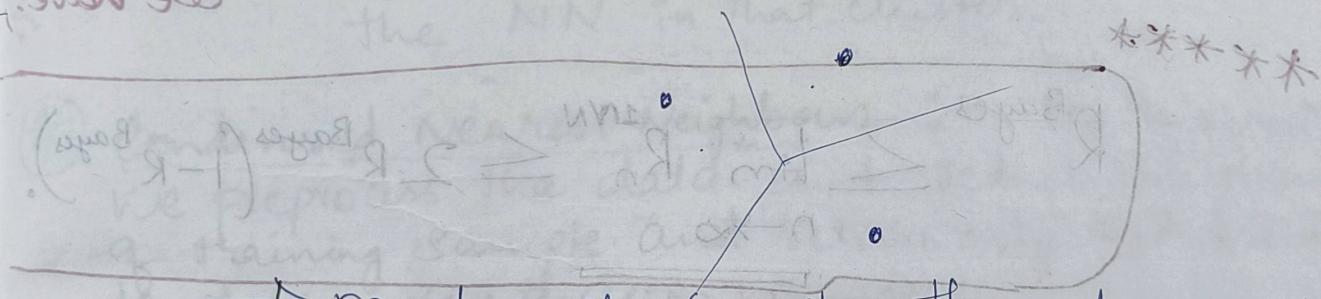
④ Curse of dimensionality \rightarrow k-NN is fully distance based hence heavily impacted by C.g.dimens

In pre-processing, we must reduce the dimension using any feature Selection method.

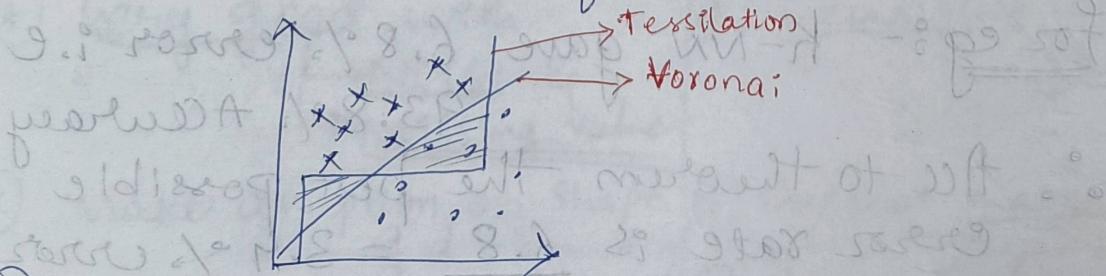
⑤ Decision Surfaces Located by KNN :-

→ Voronoi Diagram

Every point of different class is separated by a perpendicular bisector b/w them.



→ much more complex than decision trees
In Voronoi diagram there is no error less than that of axis parallel boundaries in D-T



* 'k' is the number of the nearest neighbours,
usually a will be odd number.

① Noisy training is an issue

② Can overfit.

③ missing values

Issues of KNN

K should be less than
the square root of the total
number of training patterns

(*) KNN can predict labels of any type.

→ Binary, multiclass, regression
multidimensional regression.

→ Text mining (Part-of-speech Sequence)

(*) Theorem T. → In a two-class classification problem we have.

$$R^{\text{Bayes}} \leq \lim_{n \rightarrow \infty} R_n^{\text{KNN}} \leq 2R^{\text{Bayes}}(1 - R^{\text{Bayes}}).$$

→ we can find how well can machine learning do on the particular data.

for eg:- K-NN gave 6.8% error i.e 93.8% accuracy

∴ acc to theorem the best possible error rate is $\frac{6.8}{2} = 3.4\%$ error rate

∴ we cannot get 96.6% accuracy for that particular data.

* As K-value increases the decision surface become linear and error rate slightly increases

* We can use K-means with KNN to speed up the process

↳ If in KNN, we have given a new test data then we will find the nearest neighbours for that data.

Speeding KNN using K-means

By using K-means, it will give the ~~best~~ clusters that are present in data. Therefore the new data given can be first clustered and only find the NN in that cluster.

* Condensed Nearest Neighbour :- In this method we preprocess the data and reduce the number of training sample and retain only that are needed to define the decision boundary.

Decision Trees (Eager Learning) K-Nearest Neighbours (Lazy learner)

* very good with:

→ Noise &
→ Missing value.

* Make assumption on shape of d-boundary

* DTrees are not good when future looks diff from past.

* Sensitive to noise
→ Missing value

* Does not make any assumption

* Can work when future looks diff from past

* Guaranteed accuracy
 $2 \times R_{Bayes}$

SVD → Singular Value Decomposition *

Eg Recommendation Systems :-

Molecules $\xrightarrow{\text{to understand the molecule}} \text{Composition}$

Number $\xrightarrow[\text{a number}]{\text{to understand}} \text{Factorization}$

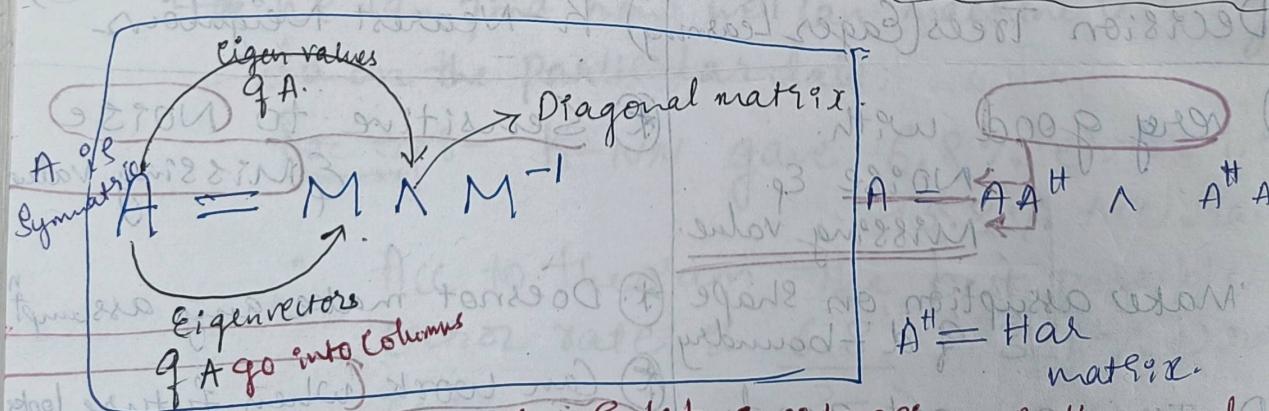
Material $\xrightarrow{\text{to understand a matrix}} \text{Decomposition}$

$$\begin{bmatrix} N_1 & N_2 & N_3 \\ P_1 & P_2 & P_3 \end{bmatrix}$$

In PCA \rightarrow Identifying Principal Component Ex:- 1728

In SVD \rightarrow Factorizing a Matrix $= 64 \times 27$

$$= 2^6 \times 3^3$$

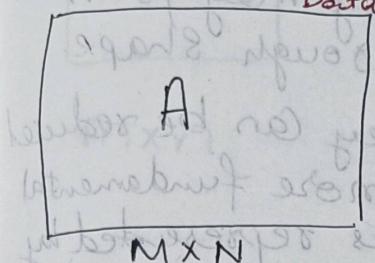


SVD enables us to split overlapping influence data into their own space and remove noise

Rank → must be
(R) decided by user

From these matrix we can find the similar users and similar cells or some relation using Cosine Similarity or Pearson Correlation.

Matrix A Data



=

$$U \Sigma V^T$$

MxR R x N Topic x user

$R \rightarrow$ Rank.

Topic x Topic matrix

→ Diagonal matrix (Every entry other than diagonal are zero)

→ Diagonal values = Eigen values of A & A^T

Each row is Orthogonal (Independent to each other)

④ Σ contains the eigen values in diagonal, as we go down the eigenvalues decrease i.e., the magnitude of stretch decreases. ∴ we can set a threshold

* We use SVD to remove overlapping influences and also removing the noise (least impact on EVal). Since Attributes in rows and columns have overlapping influences.

⑤ These features are called Latent dimensions, in raw data we have to find this.

L This latent content is calculated by matrix ~~correlation~~ Factorization (latent level)

Applications

- ↳ Data Compression.
- ↳ Noise removal
- ↳ Recommender System
- ↳ Text mining (LSI)

SVD intuition of features

- (*) The original space has a cloud (large) of points
- (*) They may form some rough shape.
- (*) They can be reduced into more fundamental factors represented by smaller set of features

Recommender Systems

→ Here Recommender System is studied using an example of user ratings for Netflix.

Context Based → Netflix can be expert to identify qualities of movie/users that tell how similar their tastes are.

User-Behavior Based → Identify similar movies (or users) using distance metric

Recommender through SVD

- Create a user space through SVD
- Remove unwanted columns (Low eigen values)

$$\text{users} \rightarrow \boxed{\begin{matrix} A \\ \text{movies} \end{matrix}} = U \sum V$$

Rank matrix movie space

MXR RXN.

- For every user compute the cosine similarity based distance and identify closest user
- Look up the ratings they gave to the movie

of feature
 space
 (rge) q
 form.
 ape.
 reduced
 amental
 ded by
 feature.

$$A = U \Sigma V^T$$

$$A_k = U_k S_k V_k$$

- ① The reconstructed matrix $A_k = U_k S_k V_k$ is the closest rank- k matrix to the original matrix A .
- ② In rank matrix we will cut down to ' k ' from ' r ' i.e. we remove the dimension which have less eigen values.
- ③ This implies each user is in smaller dimensional space which helps to find user-similarity
- ④ If we make U and V orthogonal and Σ diagonal, this can be the best factorization.

$$A = U \Sigma V^T$$

$$A^T A = V \Sigma^T U^T U \Sigma V^T = V \Sigma^2 V^T$$

V 's \rightarrow Eigen vector of $A^T A$
 U 's \rightarrow Eigen vector of $A A^T$

* Matrix factorization and recommenders are
Quicker than SVD

→ Since Σ is not used.

$$\Rightarrow A_{m \times n} = P_{m \times r} * Q_{r \times n}^T$$

→ There are r latent features that
describe the quality

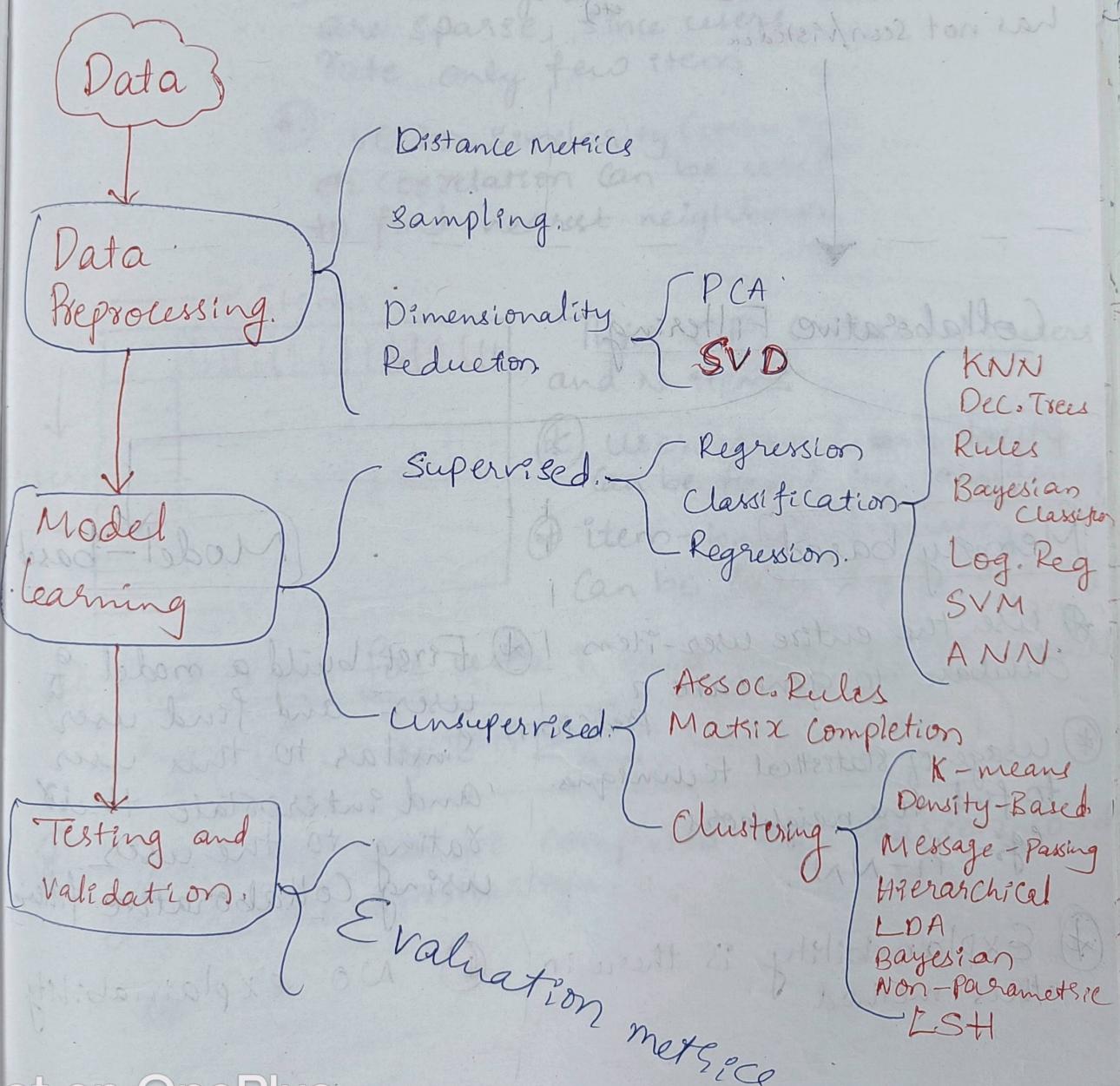
→ Minimize the error by identifying
suitable values of P and Q
(By gradient descent)

In these days Matrix factorization
is preferred more than SVD

Recommender Systems → Estimate a utility function that
automatically predicts how a user
will like an item. Based on:
* Past behaviour
* Relations to other users
* Item Similarity
* Context etc.

Recommendation is a Machine Learning Problem

(*) The core of Recommendation Engine can be assimilated to a general data mining problem



Collaborative Filtering

Content-Based Filtering

- ④ Collaborative filtering.

find the similar rows (eg:- user) and recommends the column (eg:- movies, books, etc) to other user who has not seen/visited etc.

- ④ Content-Based

filtering recommends the column (eg:- movies, books, etc) based.

Collaborative Filtering

Memory-based

- ④ Use the entire user-item database to generate a prediction.

- ④ Usage of statistical techniques to find the neighbors.
eg:- K-NN.

- ④ Explainability is there in this method.

- ④ First build a model of user and find user similar to this user and interpolate their rating to the user using collaborative filtering.

- ④ No explainability

User-based

push() of - has Item-based

"users who liked this item also liked"

"users who are similar to you also liked"

- ④ Construct a vector ~~declaration~~ for each user
 - ④ On average, user vector are sparse, since users rate only few items.
 - ④ Vector similarity (cosine sim) or correlation can be used to find nearest neighbours

A hand-drawn diagram illustrating a matrix structure for user-item ratings. The matrix is enclosed in a blue border. The columns are labeled "items" at the top, with an arrow pointing from the left towards the column headers. The rows are labeled "users" on the far left, with an arrow pointing downwards towards the row labels. The matrix contains numerous red horizontal and vertical lines forming a grid. In the bottom right corner of the matrix area, the word "ratings" is written in red, with a curved arrow pointing from the text towards the matrix.

This CF contains M users
and N items

- (*) User-based similarity
Can be found by considering rows
 - (*) Item-based similarity
Can be found by considering columns

- ④ Prediction for user i and product j is computed

$$\text{by } V_i^* = v_i + k \sum_{v_{kj} \neq ?} y_{jk} (v_{kj} - v_k)$$

- ④ Similarity can be computed by Pearson Co-rel or Cosine Similarity.

For item-based Col-filtering

- ① Look into the items the target user has rated
- ② Compute how similar they are to the target
 - ↳ Similarity only using Past ratings for other users.
- ③ Select k most similar items.
- ④ Compute Prediction by taking weighted average on the target user's ratings on the most similar items.

Recommendation

Content based

Collaborative filtering

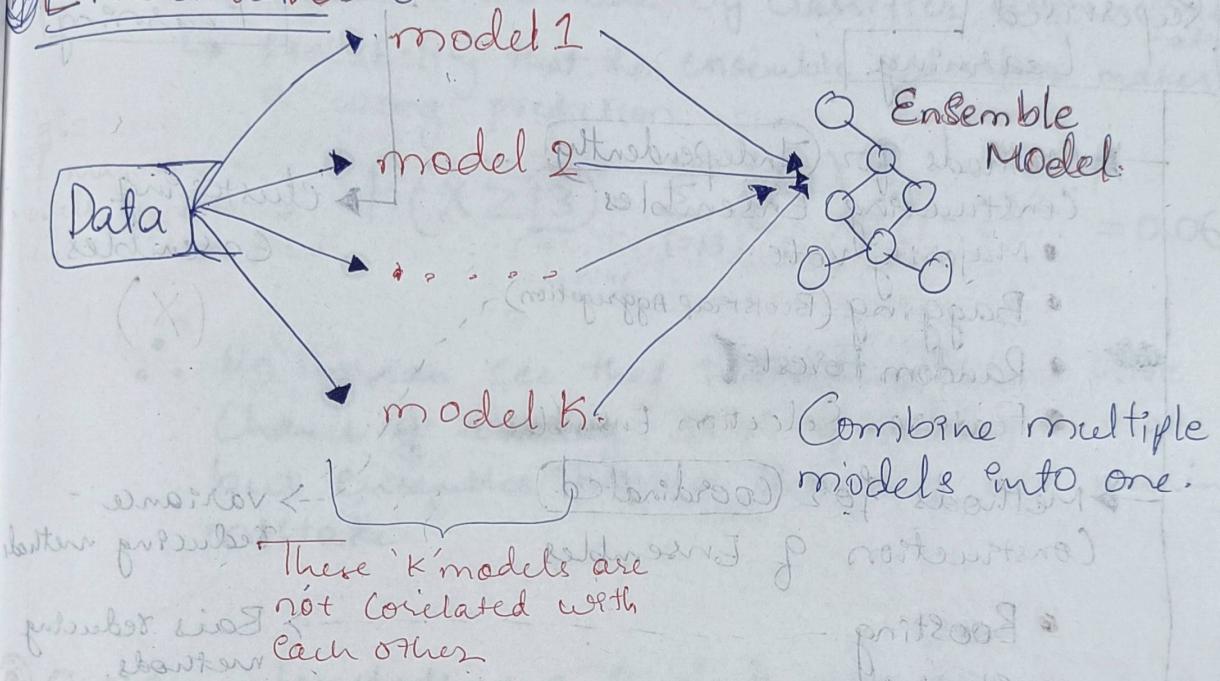
User

Item.

- ① Active users
- ② User-user similarity
- ③ Nearest user
- ④ Interpolate rating
- ⑤ Items user rated
- ⑥ How similar based on ratings
- ⑦ Weighted summations

Ensembles, RF, Stacking & SVM

Ensembles :-



Basic Idea → Build different "experts" and let them Vote.

Advantages :-

- ★ Improve predictive performance.
- ★ other types of classifier can be directly included.
- ★ Easy to implement.
- ★ Not much parameters tuning.

Disadvantages :-

- ★ The combined model is not so transparent (black box).
- ★ Not a compact representation.

Types of Ensembles.

Supervised Learning

→ Methods for Independently Constructing Ensembles

- Majority Vote.
- Bagging (Bootstrap Aggregation)
- Random Forests
- Feature-Selection Ensembles

→ Methods for Coordinated Construction of Ensembles

- Boosting
- Stacking
- Mixture of Experts

Unsupervised Learning

clustering Ensembles (X)

→ variance - Reducing method

Bias reducing methods.

Why do Ensembles work?

Mathematically

① Suppose there are 25 base classifiers

- ↳ Each classifier has error rate, $\epsilon = 0.35$ (Accuracy = 0.65)
- ↳ Assume errors made by classifier are uncorrelated
- ↳ Probability that the ensemble classifier makes a wrong prediction.

If majority in 25 models says fraud, then it is a fraud;

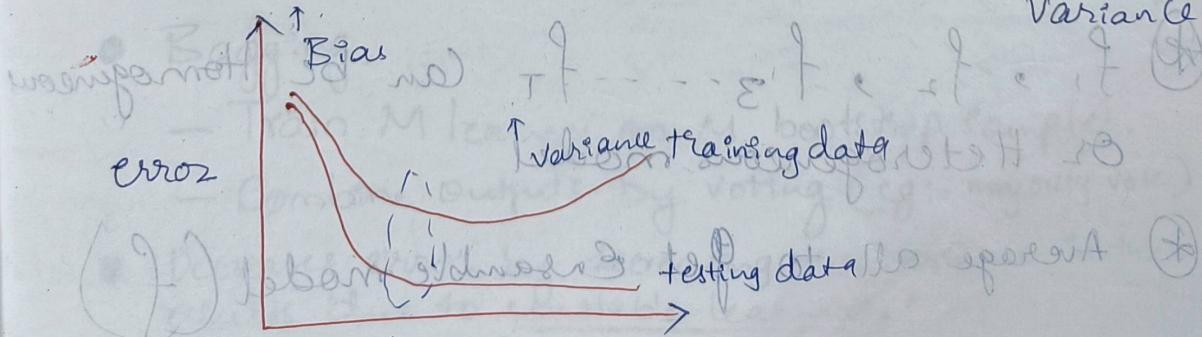
$$P(X \geq 13) = \sum_{i=13}^{25} \binom{25}{i} \epsilon^i (1-\epsilon)^{25-i} = 0.06$$

From above we can see that individual model have chance of ~~making~~ 35% of making mistake but ensembles have a chance of 6% of making mistake.

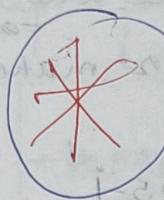
② Overcome limitations of single hypothesis

- ↳ The target function may not be implementable with individual classifiers, but may be approximated by model averaging.

③ Ensemble methods combine learners to reduce bias & variance.



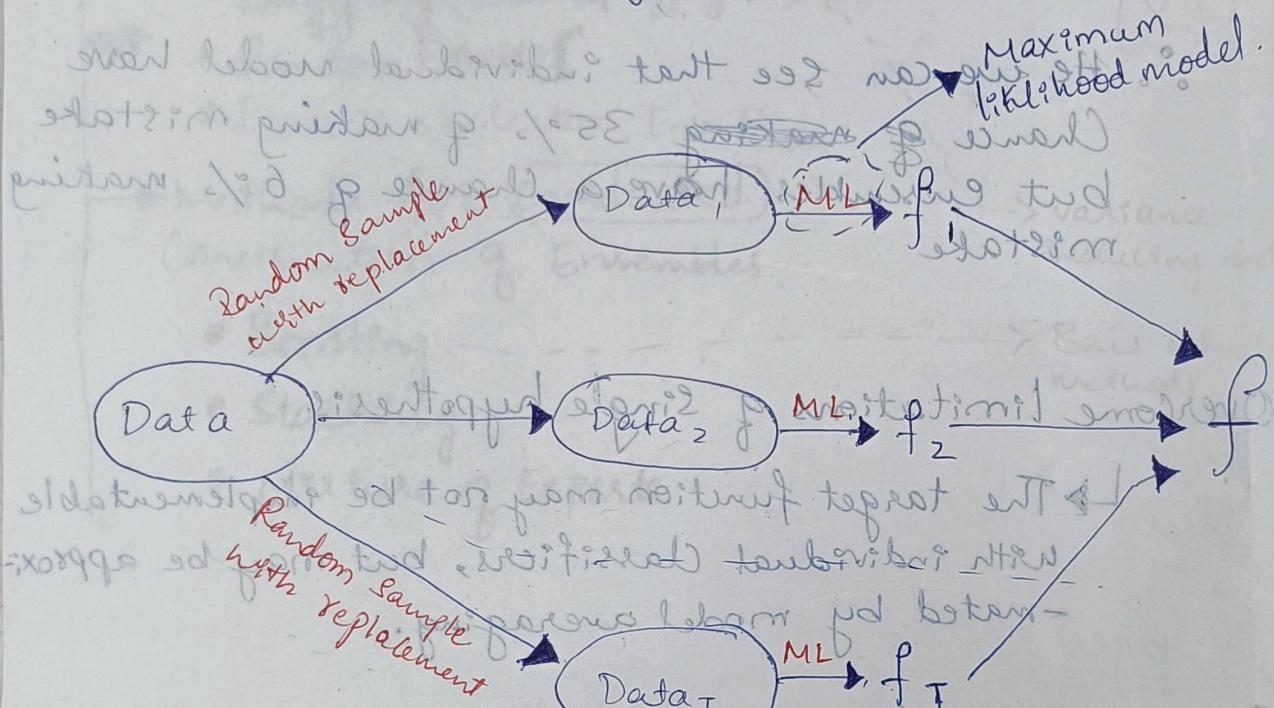
- ★ Ensembles were most popular method before neural network came to existence.



BAGGING

(Ensembles - method 1)

Bootstrap Aggregation



- ★ $f_1, f_2, f_3, \dots, f_T$ can be Homogeneous or Heterogeneous model

- ★ Average all f_i to get Ensemble model (f)

Bootstrap aggregation :- Example. Random Sampling with Replacement
our initial dataset :- $\{1, 2, 3, 4, 5, 6\}$

First draw :

$\{1, 2, 2, 2, 3, 5\}$ Model 1

Second draw :

$\{1, 1, 3, 4, 4, 6\}$ Model 2

Third draw :

$\{2, 3, 3, 4, 5, 5\}$ Model 3

Forth draw :

$\{2, 4, 4, 5, 5, 6\}$ Model 4

uncorrelated datasets

\therefore uncorrelated model.

Average : Final model
mode

Bagging :-

- Create ensembles by "bootstrap aggregation", i.e., repeatedly randomly resampling the training data.

— Bootstrap :- draw N items from \times with replacement

- Bagging

— Train M learners on M bootstrap samples.

— Combine outputs by voting (e.g.: majority vote)

• Decreases error by decreasing the variance in the results due to instable learners.

(*) Bagging is good for unstable learners as it reduces variance and overfitting

→ To generate large number of unstable learners

(Bagging)

Choose records randomly
Choose attributes randomly

(*) In Naive Bayes, Logistic regression, Linear regression

→ Bootstrap of records fails (Because there are global classifiers, Not large variance)

→ but not bootstrap of features

works well

→ Decision trees will do wonders on both.

An example of Bagging

(*) Random Forest

(*) 500-1000 trees commonly used

① Decision trees

↳ Non linear classifiers

↳ Easy to use

↳ Easy to interpret

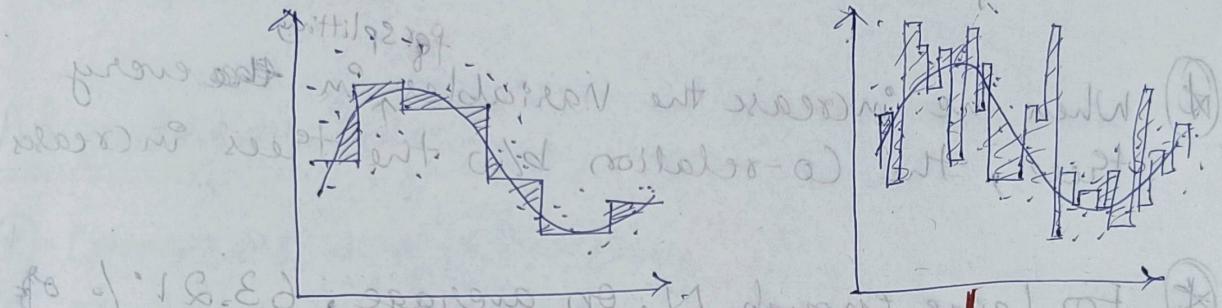
↳ Susceptible to overfitting

$$m = \sqrt{P} \rightarrow \text{regression}$$

$$m = P/3 \rightarrow \text{classification}$$

④ The problem with single decision tree which is not pruned is poor predictor (noisy) on test data. i.e., it is overfit, it does not do well on test data.

We use such unpruned predictors (d-Trees) and average their output in boosting.



Bagging uses this high variance to our advantage.

Random Sampling with Replacement

1000 sample

Choose 10 attributes out of 100 randomly

Best attr to split out of 10

K points

Choose 10 attr out of 100

→ Best attr to split

1000 - K

Choose 10 attr out of 100

→ Best attr out of 100 to split

Berney & Stanford

- ① Bagging introduces randomness into the rows of the data.
- ② Random forest introduces randomness into the rows and columns of the data.

(west-b) Tree
③ Single decision will be used only for understanding the rules in most cases after this random forest is built to solve the prediction problem.

④ When we increase the variables in every tree, the correlation b/w the trees increases for splitting

⑤ For large enough N, on average, 63.21% of the original records end up in any bootstrap sample. Other 36.79% of the observations are not used in the construction of a particular tree.

In these unused 36.79% of attributes can be used as test sample for testing the model.

These observations are called or considered.

Out of Bag (OOB)

⑥ Random forest is extremely useful in feature selection. Since in every tree we can see which attributes are used more and select the important variable.

Parameters for tuning Random forest :-

- ① Number of trees
- ② mtry
- ③ Node size --?
- ④ Sampling Scheme --?
- ⑤ Split rule.

Advantages :-

- ⊕ Competitive performance.
- ⊕ Does not overfit
- ⊕ Robust to outliers
- ⊕ Handles missing data (imputation is not required)
- ⊕ Provide automatic feature Selection.
- ⊕ minimal preprocessing time.

Disadvantages :-

- ⊕ Cannot compete with advance boosting algorithms
- ⊕ Can become slow on large data sets
- ⊕ Less interpretable

← Shot on OnePlus

By Krishna



Support Vector Machine

(SVM)

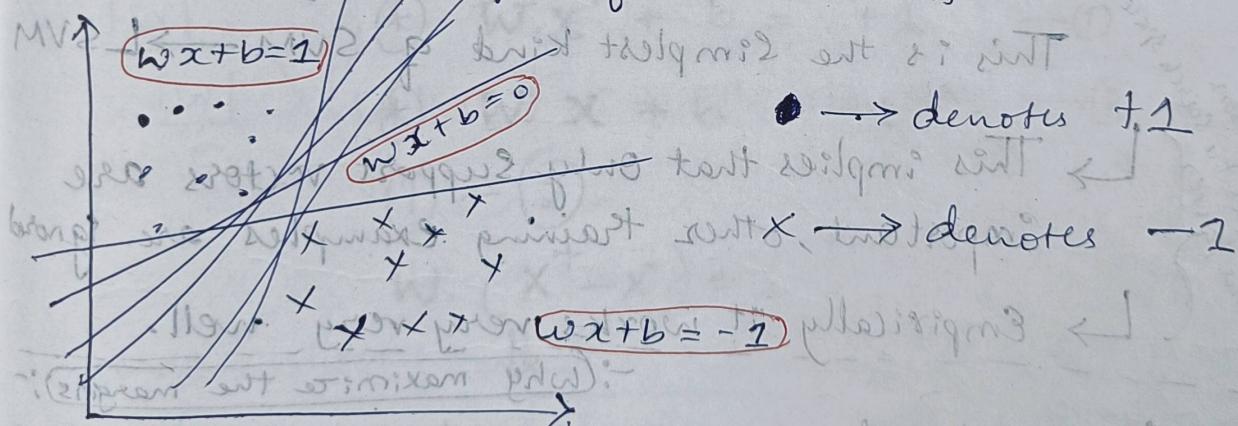
- ⊕ SVMs have clever way to prevent over-fitting.

→ SVM is arguably most imp & interesting recent discovery in ML by Vladimir Nau. Vapnik

WORKS
well in higher
dimensions.

SVMs have a very clever way to use huge number of features without requiring nearly as much computation as seems to be necessary.

- ⊕ Consider an example of linear classifier.



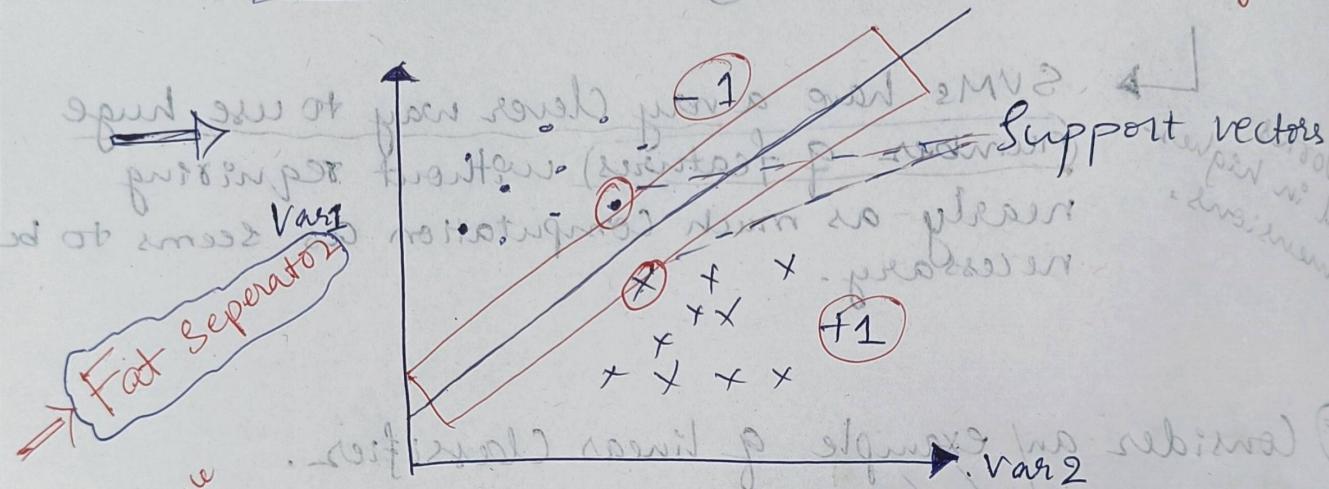
on new unseen data set A

- Here which linear separator to pick, so that it minimizes the error for training set and works well for test set.

To work well on test set the classes must not be nearer to separator (as equally far from both)

i.e. the separator must maximize the margin at the boundary

We need maximum margin linear classifier, which maximizes the margin



This is the simplest kind of SVM \rightarrow LSVM

- \hookrightarrow This implies that only support vectors are important, other training examples are ignored.
- \hookrightarrow Empirically it works very very well.
:(why maximize the margin):

\hookrightarrow A classifier with a large margin makes no low-certainty classification decisions.

Since points near decision surface are uncertain classification decisions (50% either way)

Shot on OnePlus

By Krishna

\hookrightarrow Fewer choices of where it can be put

The Search Problem becomes simpler in higher dimension because in high-dims the margin classifier inflate in all direction and finding the best fit becomes easy.

Linear SVM Mathematically :-

$$\begin{aligned} \text{if } & \Sigma \leq (d+xw), p \quad \text{"Predict class = +1" zone.} \\ \text{if } & \Sigma \geq d+xw \quad \text{if } \Sigma = d+xw \quad \text{if } \Sigma = d+xw \quad \text{if } \Sigma = d+xw \\ \text{if } & \Sigma \leq d+xw \quad \text{"Predict class = -1" zone.} \\ \text{if } & \Sigma \geq d+xw \end{aligned}$$

$$|\Sigma| = M \quad \text{what we know :-}$$

$$\star \quad w \cdot x^+ + b = +1 \quad \text{--- (1)}$$

$$\star \quad w \cdot x^- + b = -1 \quad \text{--- (2)}$$

$$(1) - (2)$$

$$w \cdot (x^+ - x^-) = 2$$

$$\text{For any dimension, } M = \frac{(x^+ - x^-) \cdot w}{\|w\|} \Rightarrow \frac{2}{\|w\|} \quad M \rightarrow \text{Margin width.}$$

④ we must maximize M margin width.

Goal

(1) Correctly Classify all training data

$$w x_i + b \geq 1 \quad \text{if } y_i = +1$$

$$w x_i + b \leq -1 \quad \text{if } y_i = -1$$



Simplifying those two eqn. \Rightarrow

$$y_i (w x_i + b) \geq 1$$

$$\text{if } y_i = +1 \Rightarrow w x_i + b \geq 1$$

$$\text{if } y_i = -1 \Rightarrow -1 (w x_i + b) \geq 1$$

$$\rightarrow \text{multiply } -1 \text{ on both sides}$$

$$w x_i + b \leq -1$$

(2)

Maximize the Margin

$$M = \frac{2}{\|w\|}$$

$$\min_{w} \frac{1}{2} \|w\|^2$$

$$\frac{\|w\|}{2}$$

$$\|w\| \Rightarrow w^T w$$

This is a quadratic optimization problem in linear constraints. Solve for w & b .

Minimize $\Phi(w) = \frac{1}{2} w^T w$ (obj func)

Subject to: $y_i (w x_i + b) \geq 1$ for all i

(Constraints)

for new yet unknown
hyperplane

$$\begin{cases}
 \underline{ax + b = 0} \\
 \underline{a_1x_1 + a_2x_2 + b = 0} \\
 \underline{a_1x_1 + \dots + a_nx_n + b = 0}
 \end{cases} \rightarrow \begin{array}{l} 2-D \\ 3-D (\text{Plane}) \\ \text{Hyperplane} \end{array}$$

\downarrow

$\vec{a} \cdot \vec{x} + b = 0$

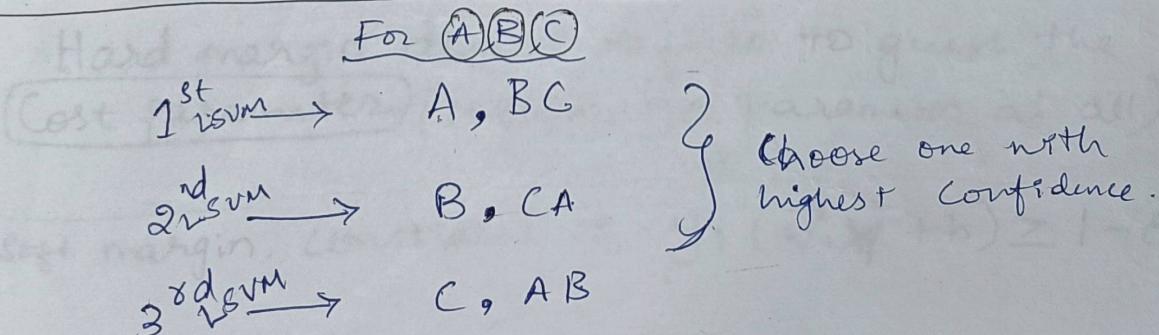
vector representation of hyperplane

If we have several thousand dimensional data points and knowing that there is a linear separator, we can find that linear separator guaranteed by LSVM.

Non-solvable if the data is not linearly separable.

Final Conclusion, For arbitrary no. of dimensions for arbitrary no. of datapoints, if there exists linear separator, we can find the separator by Simple Quadratic Programming.

LSVM is 2-class classifier, for 3-class classifier we build 3 LSVM separator and choose the one with higher confidence.

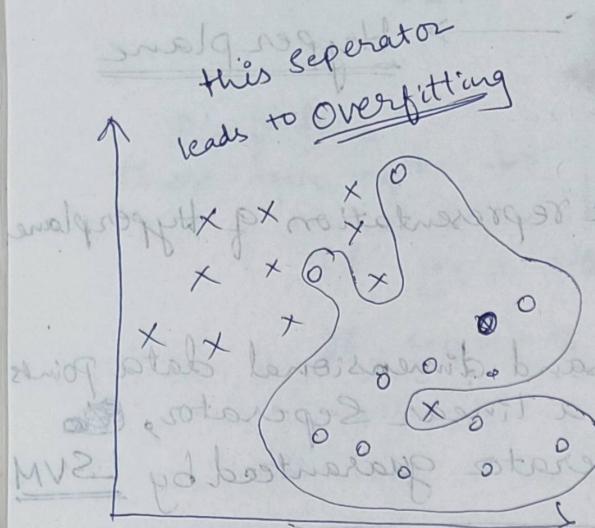




Dataset with noise.

① Initially try with soft margin classifiers with linear separator

② If ① not work try non-linear separator.



x : denotes $\rightarrow +1$

o : denotes $\rightarrow -1$

Hard margin

NO training errors.

So fast we require all data points be classified correctly.

→ What if the training set is noisy.

① Solution 1 :- use very powerful kernels

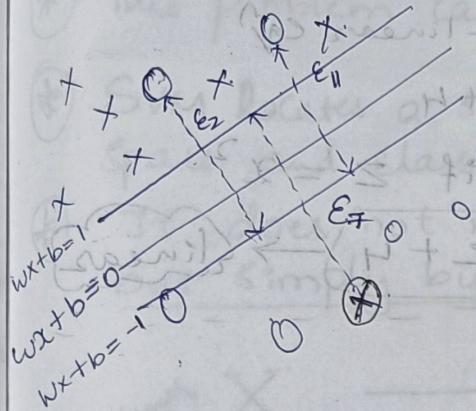
for stronger but not overfitted kernel

e.g. polynomial, sigmoid, gaussian etc.

To avoid overfitting we introduce

Soft Margin Classification.

→ Slack variables can be added to allow misclassification of difficult or noisy examples.



∴ we must minimize

$$\frac{1}{2} \mathbf{w} \cdot \mathbf{w} + C \sum_{k=1}^R \epsilon_k$$

(Weightage Error rate) Cost Parameter

* The distance of noisy points from the separator must be minimal. we add the ϵ_k term of all noisy points to the objective function.

Hard margin vs Soft margin

- Soft margin also, always has a solution
- Soft-margin is more robust to outliers
- Hard margin does not require to guess the Cost parameter (requires no parameter at all)

for soft margin, constraint is $y_i (\mathbf{w} \cdot \mathbf{x}_i + b) \geq 1 - \epsilon_i$

(*) Non-linear SVMs

The main theme of Non-linear SVMs is, a non-linear equation in a particular dimension becomes linear in higher dimensions.

for eg: consider non-linear eq

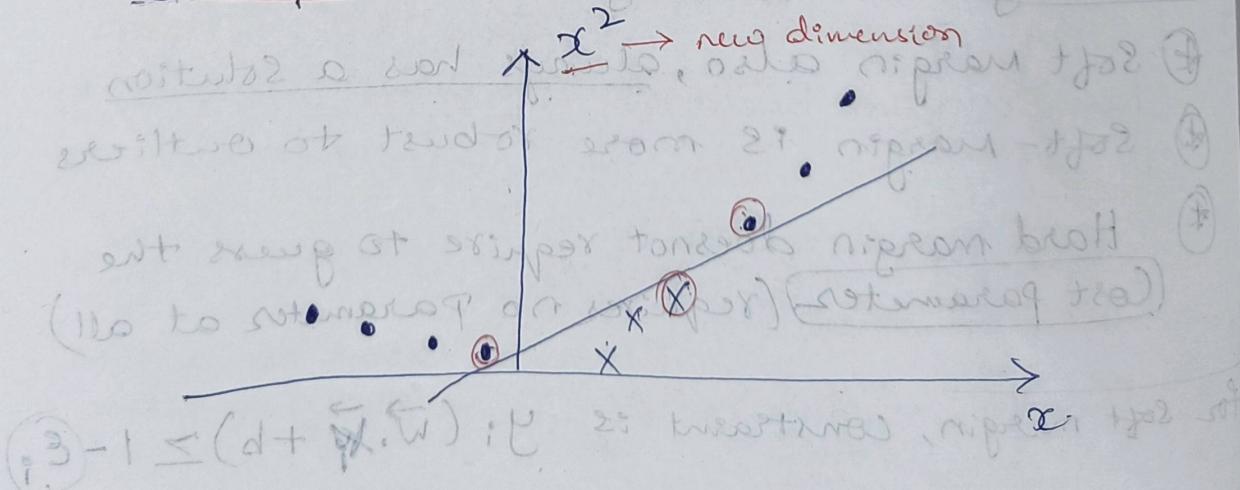
$$y = x^2 + 3x + 4$$

New dimension \rightarrow if $Z = x^2$

$$\Rightarrow [y = 3x + Z + 4] \rightarrow \text{linear}$$

Step 2: Int map 2DQ, piong wntib ent
map (3) int bbs ent to minimum ad fault
map The above points cannot be linearly separable

(*) Map the points to higher dimensions



The original input space can always be mapped to some higher-dimensional feature space where the training set is separable.

$$\phi: X \rightarrow \Psi(x)$$

This problem can be solved using **KERNEL TRICK**

SVM locates a separating hyperplane in the feature space and classify points in that space.

It does not need to represent the space explicitly, but simply by defining a **Kernel function**

$$X_1 \xrightarrow{\text{non-linearly separable}} \phi(x_1) \quad \left. \begin{array}{l} \text{In this space they} \\ \text{are linearly separable} \end{array} \right\}$$
$$X_2 \xrightarrow{\text{separ.}} \phi(x_2)$$

$$\text{Lower dim.} \xrightarrow{\text{higher dim.}} \phi(x_1)^\top \phi(x_2)$$

we find a linear separator in higher dimension by doing the dot product

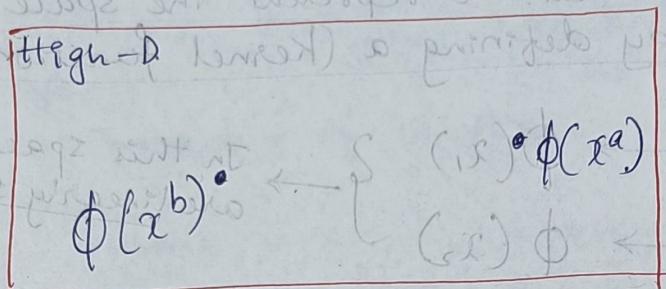
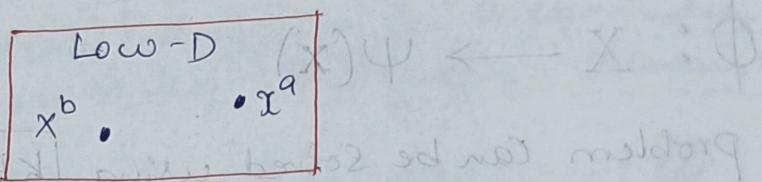
$$\Rightarrow \phi(x_1)^\top \phi(x_2)$$

Therefore kernel function is that function, if it is applied on x_1, x_2 ($f(x_1, x_2)$) such that the output is $\phi(x_1)^\top \phi(x_2)$.

This helps to get these scalar product without moving to higher dimension i.e. by staying in current dimensional space.

Kernel Trick

* The Kernel function plays the role of the dot product in the higher-dimensional feature space.



$$K(x^a, x^b) = \phi(x^a) \cdot \phi(x^b)$$

letting the kernel to do the work

doing the scalar product in the obvious way.

Examples of kernel functions :- (Picking the K-function is our choice)

* Linear : $K(x_i, x_j) = x_i^T x_j$

* Polynomial of power P_p : $K(x_i, x_j) = (1 + x_i^T x_j)^p$

* Gaussian (radial-basis function network) :

$$K(x_i, x_j) = \exp\left(-\frac{\|x_i - x_j\|^2}{2\sigma^2}\right)$$

* Sigmoid : $K(x_i, x_j) = \tanh(\beta_0 x_i^T x_j + \beta_1)$

etc (Many kernels are there)

Properties of SVM :-

* Flexibility in choosing a similarity function - ?

* Ability to handle large feature spaces (Higher dimensions.)

* overfitting can be controlled by soft margin approach

* Nice Math Property :- A simple convex optimization problem which is guaranteed to converge to a single global solution.

* Feature selection - - ?

Examples of kernel functions (Picking the k-function is our choice)

* Linear : $K(x_i, x_j) = x_i^T x_j$

* Polynomial of power p : $K(x_i, x_j) = (1 + x_i^T x_j)^p$

* Gaussian (radial-basis function network) :

$$K(x_i, x_j) = \exp\left(-\frac{\|x_i - x_j\|^2}{2\sigma^2}\right)$$

* Sigmoid : $K(x_i, x_j) = \tanh(\beta_0 x_i^T x_j + \beta_1)$

etc (Many kernels are there)

Properties of SVM :-

* Flexibility in choosing a similarity function - - ?

* Ability to handle large feature spaces (Higher dimensions.)

* overfitting can be controlled by soft margin approach

* Nice Math property :- a simple convex optimization problem which is guaranteed to converge to a single global solution.

* Feature Selection - - ?

SVM Applications :-

- ↳ used in many real world problems
- ↳ $(\xi^T x + \beta) = (\xi, x)^T K + \beta$ → text & hypertext categorizations (Features: word count)
- image classification
- bioinformatics (Protein, cancer classification)
- hand-written character recognition
- DNA array expression data analysis

Disadvantages :-

- ④ Sensitive to more noise.
- ④ It only considers two classes.
- ④ Choice of Kernel.
- ④ Choice of Kernel parameters.

The Classification rule :-

* The final classification rule is quite simple.

$$\text{bias} + \sum_{\text{SSESV}} w_s K(x^{\text{test}}, x^s) > 0. \quad ?$$

\uparrow
The set of support vectors

$> 0 \rightarrow C_1$
 $< 0 \rightarrow C_2$

* All the cleverness goes into Selecting the Support Vectors that maximize the margin and Computing the weight to use each support vector.

Advance ML models :-

* SVM.

* Deep learning

* Ensembles
(RF, XGBoost)

95% Sophisticated ML Models