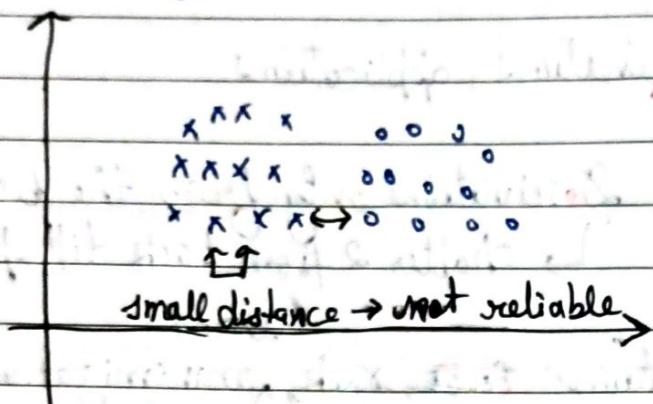


Clustering

- It is a type of
- Unsupervised learning.



- We take a central value to determine if a data point belongs to one cluster.
- Clustering is an unsupervised machine learning technique used to group similar data points together based on their characteristics.
- It is useful in Exploratory Data Analysis to uncover patterns & relationships in data.
- The fundamental concept of similarity & dissimilarity are crucial in clustering as they define how points are grouped.

Similarity in clustering :-

Similarity measures how close or alike two data points are. It is often quantified using various distance or similarity metrics.

Common similarity measures are :-

- (i) Euclidean distance
- (ii) Cosine similarity
- (iii) Jaccard similarity
- (iv) Pearson correlation coefficient

Dissimilarity in clustering :-

It measures how different two data points are from each other. It is often referred as distance in data space.

Common dissimilarity measures :-

- (i) Manhattan distance
- (ii) Mahalanobis distance
- (iii) Hamming distance.

Minimum within clustered distance criteria

→ The minimum within ~~clustered~~ distance aims to ensure that data points within a cluster are as similar as possible, which means, minimizing the distance between all data points in the same cluster.

→ This leads to tighter, more cohesive clusters.

Within cluster distance :-

It is measure of how far apart the data points within a single cluster are from each other.

The goal is to minimize this distance to ensure that all data points in a cluster are closely related or similar.

→ Within cluster sum of squares :-

For a cluster C_i , the within cluster distance is often calculated as the sum of distances between each point in the cluster and the centroid μ_i of that cluster.

$$WCSS = \sum_{x_j \in C_i} \|x_j - \mu_i\|^2$$

i) Minimizing within cluster distance improves the homogeneity of the cluster.

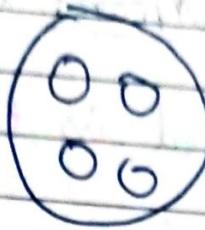
(ii) Lower within cluster distance leads to compact clusters.

There is often a trade-off between this and creating well separated clusters.

Major clustering approaches :-

1. Partitioning based algo : k-means & K-medoid DBScan
2. Hierarchical algos : agglomerative or divisive
3. Density based algos : sparse or dense area (non linear clustering)
4. Grid based algos : multidimensional
5. Model based algos : density function
6. Graph based algos : minimum spanning tree.

(i) Partition based :



we define no. of clusters, 'K' prior.

→ this is also a limitation of this method, that we need to decide 'K' in prior

(ii) Hierarchical



Agglomerative



top-down

Divisive

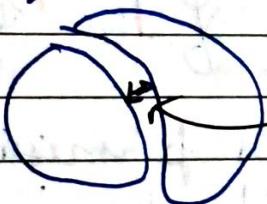


bottom-up

→ assume each data point is a class & keep joining similar points

→ assume all points are in same class and then split acc. to dissimilarity

(iii) Density instead of distance is used.



even though there is small distance but density is used to group points

→ Nowadays, fuzzy means clustering is used:-

→ Fuzzy mean clustering → same point can belong to more than one class.

Symbols

Partitioning : Kmean & K medoid

Density based : DB-Scan

(data point)
↓

K-Means Clustering

- (i) Objects are defined in terms of set of attributes
each $A = (A_1, A_2 \dots, A_n)$
where, A_i is a continuous data type.
- (ii) Distance computation : Any distance such as L1, L2, L3 or cosine similarity.
- (iii) Minimum distance is the measure of closeness between an object and centroid
- (iv) Centroid calculation : It is the mean value of each attribute values of all objects.
- (v) Convergence criteria : Any one of the following are termination condition of the algorithm.

- ① No. of max. iteration permissible
- ② No change of centroid values in any cluster.
- ③ Zero or no significant movement of object from one cluster to another.

Q There are 15 data points, $k=3$. Use Euclidean distance

- Step-1 Plot the points and make guess on potential centroids
- Step-2 Calculate distance of each point from centroids. Point belongs to the group whose centroid is nearest
- Step-3 Calculate new centroid by taking into account all points belonging to a group.
→ Now begin next iteration.

S.No	X	Y
1	2	10
2	2	6
3	11	11
4	6	9
5	6	4
6	1	2
7	5	10
8	4	9
9	10	12
10	7	5
11	9	11
12	4	6
13	3	10
14	3	8
15	6	11

Ans :-



way

Ans Set C₁ (4, 9), C₂ (6, 4) and C₃ (10, 12)

As distance

	(4, 9)	(6, 4)	(10, 12)	class
1	2.236	7.211	8.24	C ₁
2	3.606	4.47	10	C ₁
3	7.28	8.602	1.41	C ₃
4	2	5	5	C ₁
5	5.38	8.602	8.94	C ₂
6	7.616	8.602	5.38	C ₂
7	1.414	5.38	6.08	C ₁
8	0	8.602	5.38	C ₁
9	6.7	8.94	0	C ₃
10	5	1.41	7.616	C ₂
11	5.38	7.615	1.41	C ₃
12	3	2.82	8.48	C ₂
13	1.41	6.71	7.28	C ₁
14	1.41	5	8.06	C ₁
15	2.83	7	4.12	C ₁

1

2

Date / /

$$C_1 \text{ mean} = (1.51, 3.875, 9.125)$$

$$C_2 \text{ mean} = (4.5, 4.25)$$

$$C_3 \text{ mean} = (10, 11.33)$$

distance

$$(3.875, 9.125)$$

$$(4.5, 4.25)$$

$$(10, 11.33)$$

1

2

3

4

Advantages

- (i) k-means clustering is simple and can be used for a variety of objects
- (ii) It is efficient from storage point of view as well as from time point of view.

Disadvantages :

- (i) It does not work on categorical data.
- (ii) If there is some outlier, the performance of k-means clustering degrades
- (iii) It finds a local optima and may actually minimize global optima.

→ k-mean gets affected in presence of outliers
 ∵ k-medoid is used.

Partition around medoid

* k-medoid (PAM)

(i) Initialize : Select k random points out of the N data points as medoid.

(ii) Associate each data point to the closest medoid by using any common distance measure.

(iii) Repeat until the cost decreases

- ↳ For each medoid 'm', for each data point 'o' which is not a medoid
- ↳ swap 'm' and 'o', associate each data point to the closest medoid and recompute the cost
- ↳ if ~~the~~ the total cost is more than in the previous step, then, undo the swap.

Area :-

~~1 + 3 + 16 = 20~~ don't take bcz. cost

e.g.	S.No	x_1	x_2	$C_1(3,4)$	$C_2(7,4)$	Cost	<u>$C_1(3,4)$</u>	<u>$C_2(7,3)$</u>	high
	1	7	6	6	2	Cost	6	3	C_2
	2	2	6	3	7	Cost	3	8	C_1
	3	3	8	4	8	Cost	4	9	C_1
	4	8	5	6	2	Cost	6	3	C_2
	5	7	4	4	6	Cost	4	1	C_2
	6	4	7	4	6	Cost	4	7	C_1
	7	6	2	5	3	Cost	5	2	C_2
	8	7	3	5	1	Cost	5	6	C_2
	9	6	4	3	1	Cost	3	2	C_2
	10	3	7	0	4	Cost	0	5	C_1
						<u>Cost = 20</u>			<u>Cost = 22</u>

Q For a given $k=2$, cluster the above dataset using PAM. Choose $(3,4)$ and $(7,4)$ as medoid and suppose you are using Manhattan distance measure.

Ans e.g. ~~selected~~,

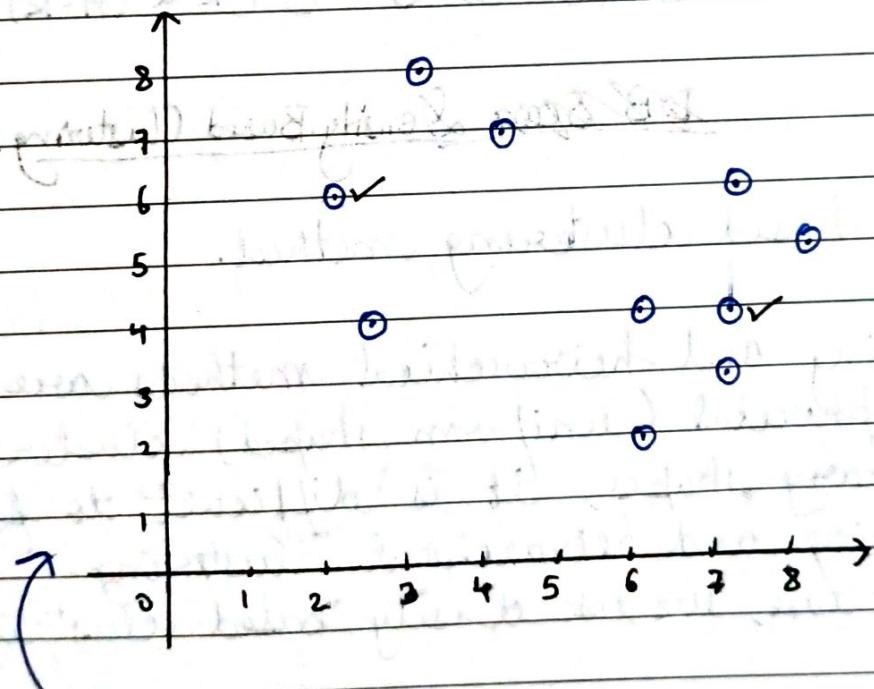
$$\text{Cost} = \text{Cost}((7,4), (7,6)) + \text{Cost}((7,4), (8,5)) + \dots$$

*

2

$$\therefore \text{cost} = 20$$

→ Now, replace $(7,4)$ with some other point and then determine cost.



Draw graph to get better
idea about desired medoids

k-mean v/s k-medoid

- (i) → Both algorithms ~~must~~ need predefined value for k, i.e. number of cluster prior to the training of algorithm.
- (ii) Both algorithm arbitrarily choose the initial cluster centred or medoid.
- (iii) The k-medoid method is more robust than k-mean in the presence of outliers.
- (iv) The time complexity of k-mean is $O(n * k * d)$ while in k-medoids it is $O(k * (n-k)^2)$

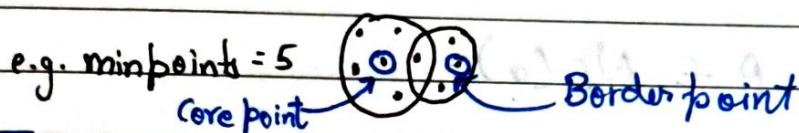
~~DB Scan~~ Density Based Clustering :-

- Density based clustering method.
- Partitioning and hierarchical methods are designed to find spherical (uniform shaped) clusters. In case of arbitrary shapes, it is difficult to find with partitioning and hierarchical clustering.
In that case, we use density based clustering methods.
- Basic idea : We make clusters where data points are dense.
A cluster is defined as a maximal set of density connected points.
- It is based on ~~two~~ two parameters
 - ① ϵ (maximum radius of the neighbourhood)
 - ② min points (minimum no. of points in an ϵ -neighbourhood of that point)

$\therefore \epsilon$ is max radius & min points is min. no. of points in a cluster.

- Core point : 'minpoints' points in ϵ boundary
- Border point : less than 'minpoints' points in ϵ boundary but some of these points are part of core point.
- Outlier point : none of above two hold

e.g. minpoints = 5



- ϵ -neighbourhood objects within a radius of ϵ of an object

$$N_{\epsilon}(P) = \{q \mid d(p,q) \leq \epsilon\}$$

↑ neighbourhood
of P

↑ points with distance less than ' ϵ '

- High density : ϵ -neighbourhood of an object contains at least minpoints of objects.
- Low density : otherwise

- Given ϵ and minpoints, categorize the objects into 3 exclusive groups

- ① A point is a core point if it has more than or equal minpoints objects within ϵ . These are the point that are at interior of the cluster.

② Border point has data objects fewer than minpoints within ϵ , but it is in neighbourhood of a core point.

③ Outlier (noise point) is any point that is neither a core point nor a boundary point.

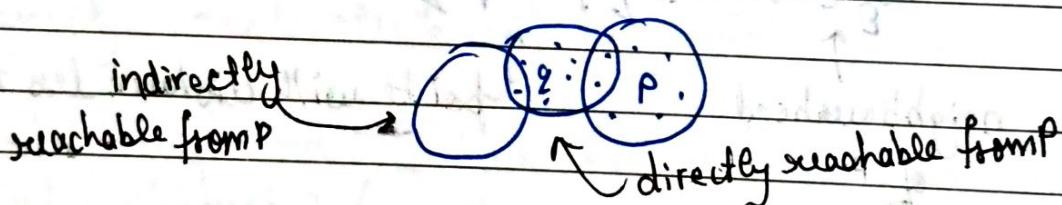
→ Directly density reachable:-

A point P is directly density reachable from a point q w.r.t. ϵ , minpoints, if P belongs to neighbourhood of q

$$\text{i.e. } P \in N_\epsilon(q)$$

otherwise indirectly density reachable

→ Density reachable or indirectly density reachable



DB Scan

→ Density Based Spatial Clustering of Application with Noise

Algorithm :-

1. Arbitrarily select point p
2. Retrieve all points, ~~stars~~ which are density reachable from P w.r.t. ϵ and ~~for~~ minpoints.

3. If P is a core point, a cluster is formed
- ↳ If P is a border point, no points are density reachable from P then DBSCAN visits the next point of the dataset
4. Continue the process until all the points are accessed.

→ Determining ϵ and minpoints

→ Idea is that for points in a cluster, their k^{th} nearest neighbour are at roughly same distance..

→ Noise points have the k^{th} nearest neighbour at farther distance.

Unit-5 : Feature extraction

Dimensionality Reduction

→ Why? to avoid curse of dimensionality.

→ Curse of dimensionality :- no. of features i.e. dimension of feature vector

As no. of features increases, i.e., we fine tune to dataset, then it nicely follows the dataset.

∴ It will learn & classify very nicely.

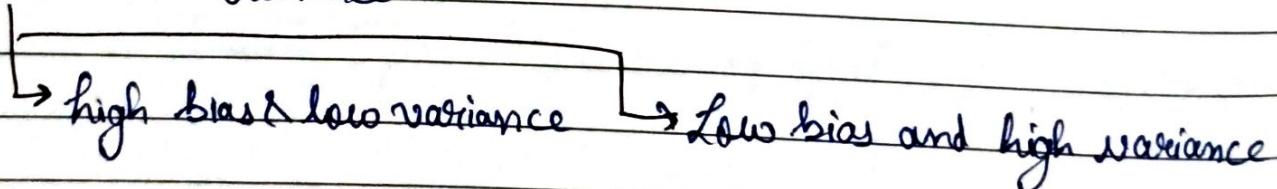
But, when test dataset is used or the same algo is applied to other dataset then it is difficult for our model to classify properly

Therefore, No. of features should not be too much.

∴ if no. of features is too much, it shows too good accuracy at training data and low accuracy at test data.

Also, if features are too less, model will miss minor details and misclassification rate increases.

→ Bias and Variance

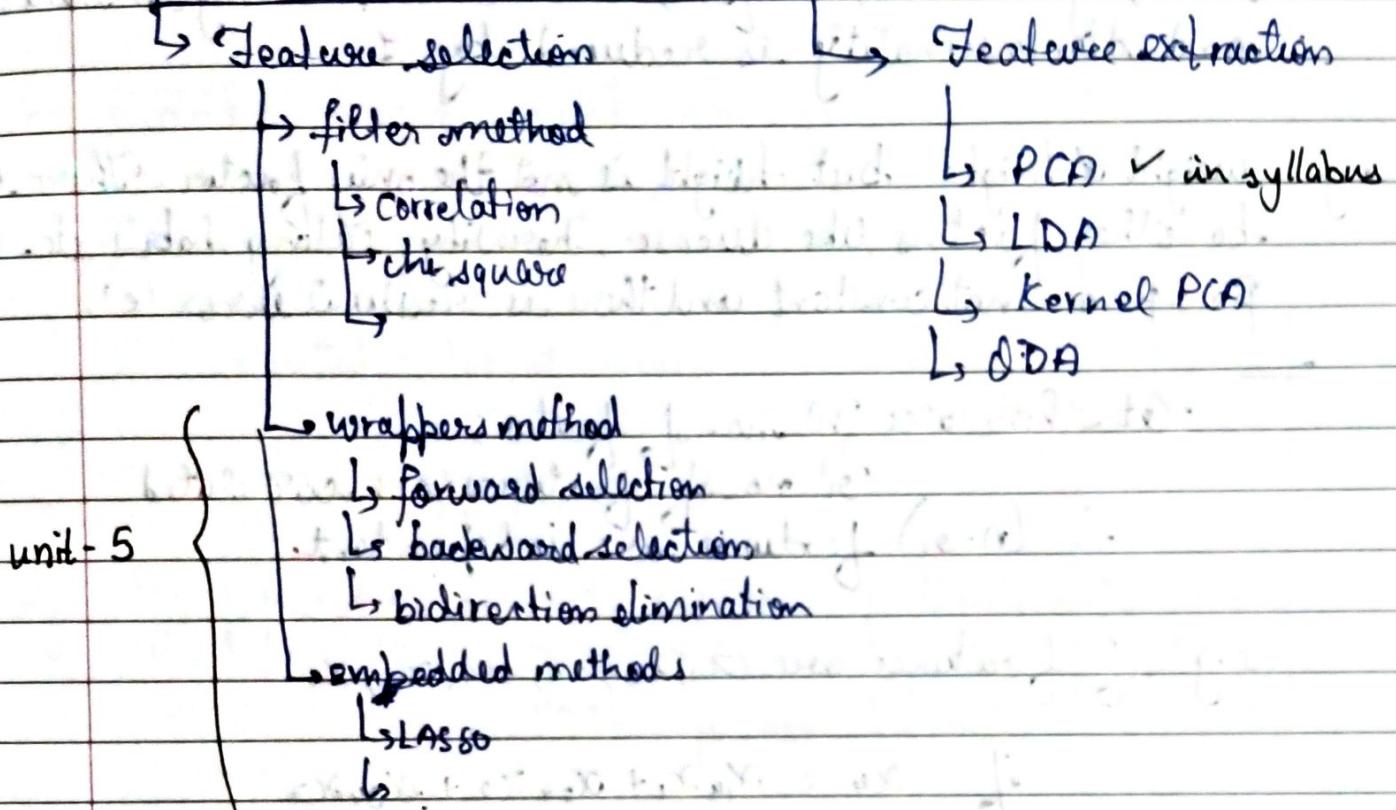


→ Dimensionality reduction : If two or more features give similar information then they can be reduced to one.

- Benefits :
- space required by system reduces
 - less computation
 - easy visualization
 - remove redundancy

Disadvantages : (i) some data might be lost
 (ii) in PCA, principal components are unknown.

→ Approaches



Principal Component Analysis (PCA)

- Part of unsupervised machine learning
- dimensionality reduction + feature extraction
- Increases interpretability and minimizes data loss
- most important features are those which are independent, otherwise one feature can be used to derive the other.

e.g. $y = \beta_0 + \beta_1 x + e$ in linear regression
 ↑ ↑ ↑
 intercept slope residual error

then, $y \propto x$ and, therefore, we need not store y in data-table and dimensionality is reduced by 1.

e.g. weight & height but height is not the only factor. There can be other factors like disease, heredity, eating habits etc. $\therefore \beta_0$ & β_1 are not constant and there is residual error ' e '.

Let there are 'n' no. of features

'r' no. of features are correlated
 $\therefore (n-r)$ features are independent.

e.g. if features are $x_1, x_2, x_3, x_4, x_5, x_6$

$$\text{if } x_4 = \alpha_{14}x_1 + \alpha_{24}x_2 + \alpha_{34}x_3$$

$$x_5 = \alpha_{15}x_1 + \alpha_{25}x_2 + \alpha_{35}x_3$$

$$x_6 = \alpha_{16}x_1 + \alpha_{26}x_2 + \alpha_{36}x_3$$

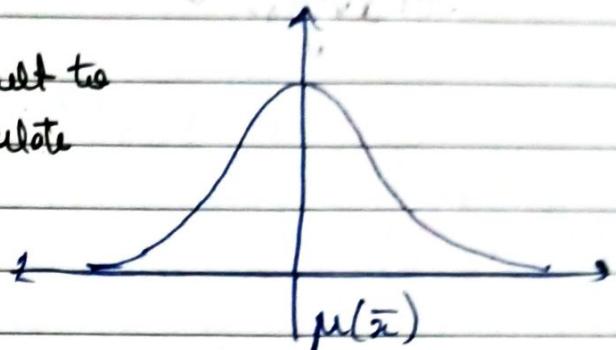
then 6 features can be reduced to 3: x_1, x_2 and x_3

→ Basic terminologies of PCA

- ↳ Variance → spread of data
- ↳ Covariance → dependencies & relationship
- ↳ Standardization → scaling each data within a range
i.e. normalization

$$\text{Covariance matrix, } \text{cov}(x, y) = \frac{\sum (x - \bar{x})(y - \bar{y})}{n}$$

- As we know, most of real life functions follows parametric distribution.
- Model is based on prediction which use central features such as mean.
- μ and σ are different difficult to calculate
 $\mu = \text{mean of population}$
 $\bar{x} = \text{expected value of mean}$ ↳ sample which is a representative of population

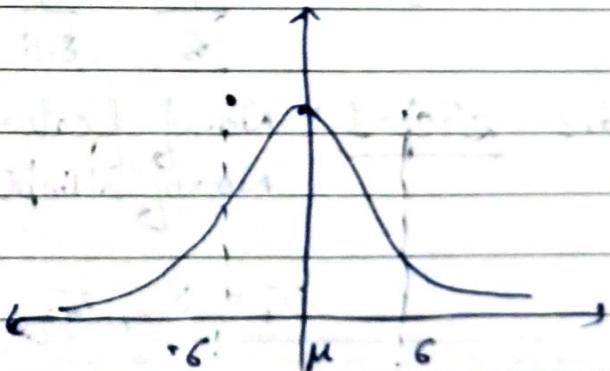


e.g. if population is 20% children
 $30\% > 60$
 $50\% \text{ younger}$ } then sample is a small part of population with similar distribution.

∴ use ground truth to find distribution & then use it to generate a small sample.

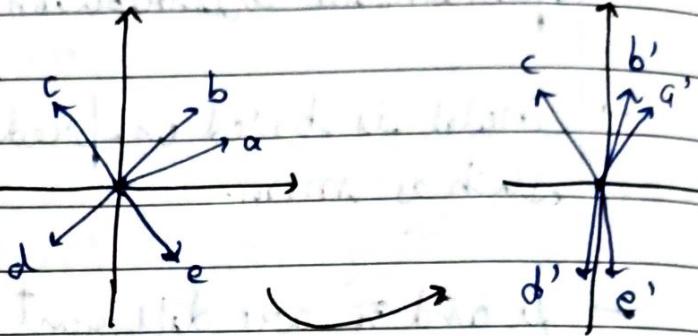
$$\therefore \mu \approx \bar{x} \text{ but } \mu \neq \bar{x}$$

- Around 60-70% data is from $-\sigma$ to σ .
- ∴ Covariance is used to find deviation from mean.



→ Eigenvalues & eigenvectors

even after the transformation, its orientation remains same even if its magnitude changes.



Linear transformation

→ direction of all vectors changes except 'c' ∵ 'c' is eigenvector.

$$\rightarrow x - \mu \quad , \quad \frac{z - \mu}{\sqrt{6}}$$

Normal distribution

standard form of normalization

Sample	x_1	x_2
S_1	4	11
S_2	8	4
S_3	13	5
S_4	7	14
	8	8.5

Ans Step-1 → No. of features = 2

No. of samples = 4

$$\bar{x}_1 = 8, \bar{x}_2 = 8.5$$

Step-2: Find cov. matrix, $\Sigma = \frac{1}{(N-1)} \begin{bmatrix} \text{cov}(a,a) & \text{cov}(a,b) \\ \text{cov}(b,a) & \text{cov}(b,b) \end{bmatrix}$

$$\therefore \Sigma = \begin{bmatrix} 14 & -11 \\ -11 & 23 \end{bmatrix}$$

Step-3 : Calculate eigenvalue & eigenvector

$$|\Sigma - \lambda I| = 0$$

$$\therefore \begin{vmatrix} 14-\lambda & -11 \\ -11 & 23-\lambda \end{vmatrix} = 0 \quad \therefore \lambda^2 - 37\lambda + 201 = 0$$

on solving quad. eqn

$$\therefore \lambda_1 = 30.38, \lambda_2 = 6.62$$

important

* Always calculate eigenvector according to highest eigenvalue first.

: arrange eigenvalues in descending order.

↳ Note that eigenvalues are always unique but eigenvectors are not. ∵ after this step, students might get different answers i.e. sign might be different in answers.

Step-4 : $(\Sigma - \lambda I) U_1 = 0$

↑ ↑ → eigenvector
cov. matrix eigenvalue

$$\therefore U_1 = \begin{bmatrix} u_1 \\ u_2 \end{bmatrix} \quad \text{apply any method to find } u_1 \text{ & } u_2$$

in this example, $\frac{u_1}{11} = \frac{11.2}{14-30.38}$

$$\therefore \text{one possible } U_1 = \begin{bmatrix} 11 \\ -16.38 \end{bmatrix}$$

Normalization :-

$$\bar{u}_1 = \frac{11}{\sqrt{(11)^2 + (-16.38)^2}} = 0.5575$$

$$\bar{u}_2 = \frac{-16.38}{\sqrt{(11)^2 + (-16.38)^2}} = -0.8862$$

similarly

$$\therefore U_2 = \begin{bmatrix} 0.88308 \\ 0.5574 \end{bmatrix}$$

Step-5

first
principal
component

$$P_{11} = U_1^T \begin{bmatrix} 4.8 \\ 11.85 \end{bmatrix} = -4.305$$

↓
 x_1 of S_1 ↓
 \bar{x}_1

↓
 x_2 of S_1 ↑
 \bar{x}_2

$$P_{12} = U_1^T \begin{bmatrix} x_1 - \bar{x}_1 \\ x_2 - \bar{x}_2 \end{bmatrix} = 3.7859$$

↑
 x_1 of S_2
↓
 x_2 of S_2

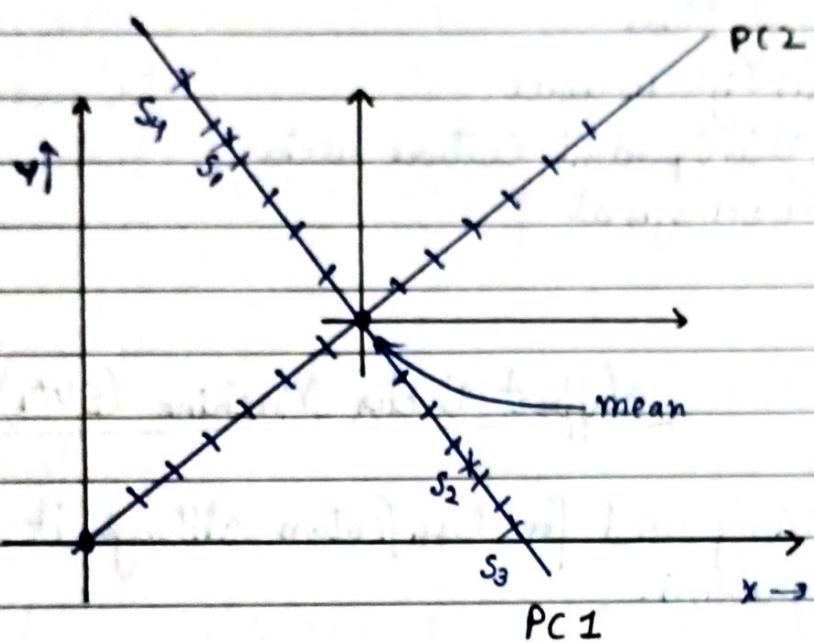
$$P_{13} = 5.692$$

$$P_{14} = -5.123$$

$$\therefore PC1 = -4.305 \quad 3.7859 \quad 5.692 \quad -5.123$$

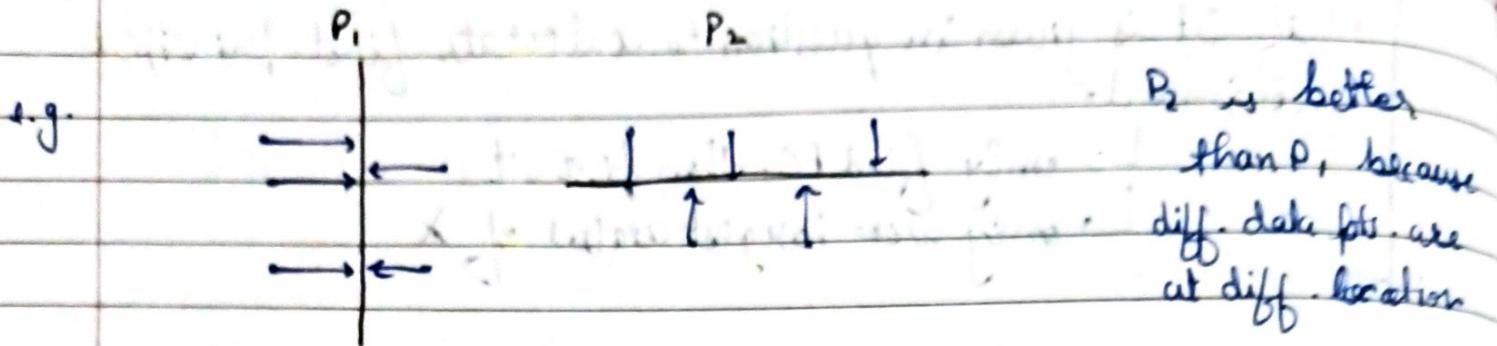
- if it is given in question to calculate first principal component
 - ↳ only find PC₁, then plot it
 - ↳ only use largest value of λ

Plotting ↗



→ we can take any of the two new perpendicular lines as PC₁ and PC₂, main importance is that they pass through mean & are perpendicular

∴ in this example, we reduced 2D into 1D bcz max info (90%) is coming from PC₁.



Eligibility of plane

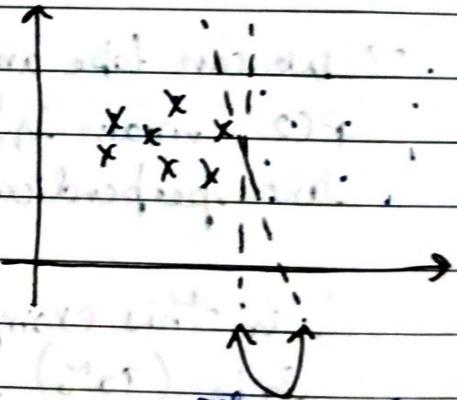
- ↳ (i) high variance
- ↳ (ii) independent feature vector
- ↳ (iii) orthogonal

Support Vector Machine (SVM)

→ It is mainly used for classification although it can also be used for regression.

→ Our goal is to create a decision boundary

→ We need a criteron to determine best decision boundary

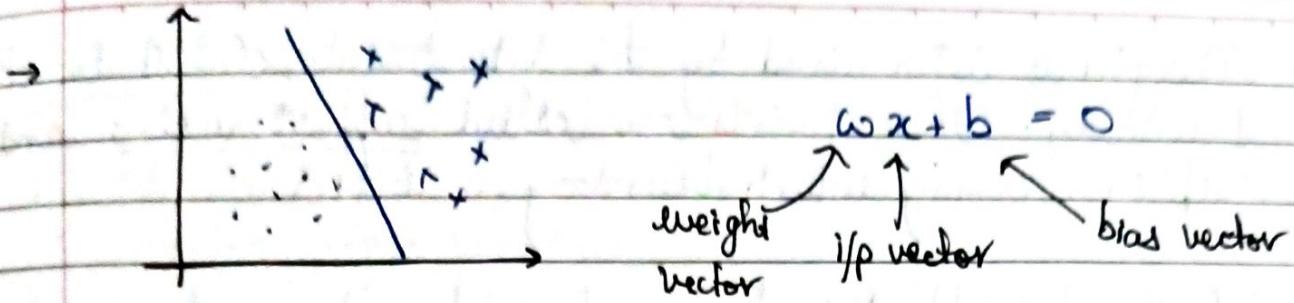


There can be infinite types of decision boundaries

→ Let us relate SVM with perceptron.

In perceptron

- ↳ (i) choose random decision boundary
- ↳ (ii) calculate error
- ↳ (iii) update decisions boundary
- ↳ repeat till min. error.



The distance of boundary from itself is zero.

LHS : $w \cdot x + b > 0 \therefore$ +ve class

RHS : $w \cdot x + b < 0 \therefore$ -ve class

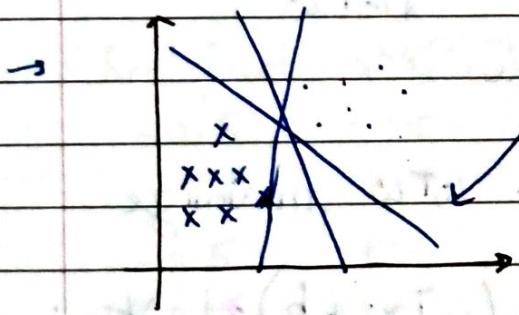
Margin : Best decision boundary has the largest boundary

→ 2D : line

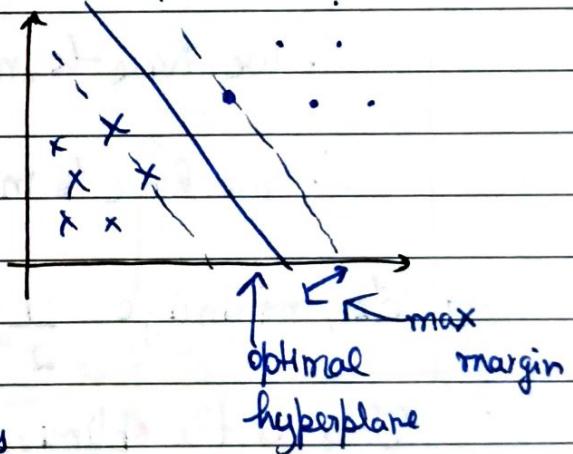
3D : plane

> 3D : hyperplane

In general we call it hyperplane in n-dimensional space.



There can be multiple possible hyperplanes



→ There is chance that some data comes whose features are close to boundary.

∴ smaller the margin, more difficult to classify unique features

greater margin ensures that in order to change nature there is larger distance to travel to cross threshold.

Linear classifier : $f(x, w, b) = \text{sign}(w \cdot x + b)$

+ve or -ve class

- Margin is determined by the data points closest to the boundary. These points are called support vectors and help in determining which boundary is better.
- ∴ farther the support vectors are from boundary, bigger the margin and better boundary.
- Margin can be defined as max width the
- since we only consider support vectors, it requires less calculation.

linear SVM:

$$\text{Distance} \rightarrow \|w\| = \frac{1}{2}$$

important: derivation for distance in linear SVM.

$$\rightarrow \text{linear SVM distance} = \frac{1}{\|w\|} = \frac{1}{\sqrt{w^T w}} = \lambda$$

∴ we have to maximise distance

∴ we have to minimize $\frac{1}{2} w^T w$ i.e. $w^T w$ minimize

$$\rightarrow \text{do, minimize } \frac{1}{2} w^T w \text{ given } y_i(w^T x_i + b) \geq 1 \quad \forall x_i, y_i$$

This is the optimization problem.

⇒ In case of noise, we use more powerful kernels to avoid overfitting.

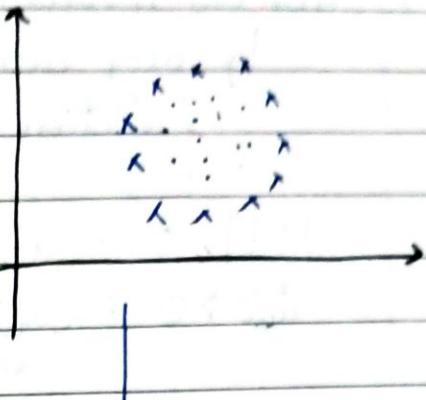
$$\text{Soft margin classification: minimize } \frac{1}{2} w^T w + C \sum_{k=1}^n \epsilon_k$$

constant to
avoid overfitting

error for
each misclassification

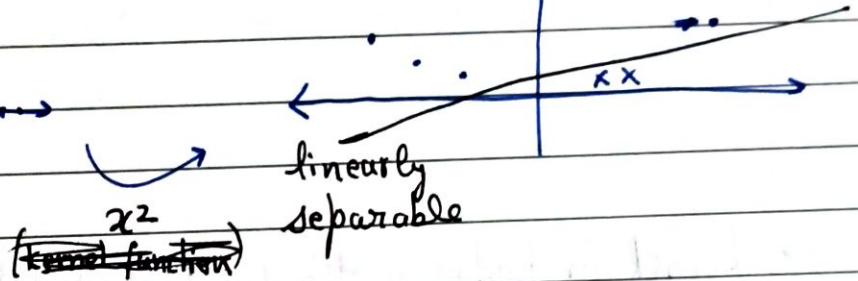
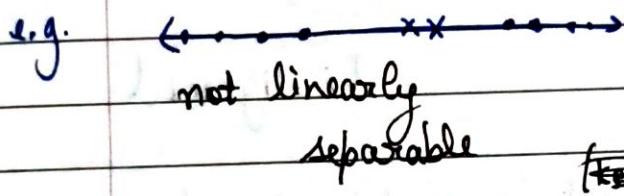
Non-linear SVM

→ Sometimes data is not linearly separable & we need non-linear hyperplane.



→ We increase dimension i.e. make higher dimension space to separate points easily.

transform using $z = x^2 + y^2$
then ~~use~~ they are linearly separable.



e.g. Red points x^2 ~~Blue points~~

-0.5 0.25

-1 1

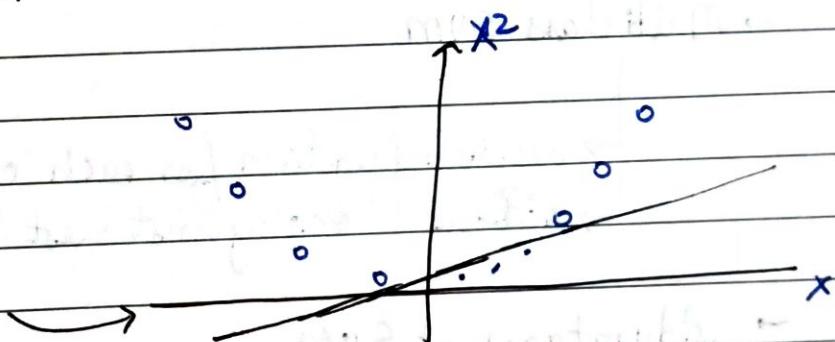
-2.5 6.25

-3.5 12.25

2.5 6.25

3.5 12.25

4.0 16



Blue points

~~0.25~~ 0.5 0.25

1 1

1.5 2.25

→ Non-linear boundary increases complexity : we map to higher dimension where it is linearly separable.

→ Mercer's theorem

↳ every semi-positive definite symmetric function is a kernel function.

$|\det| \geq 0$: semidefinite

$\det > 0$: definite positive.

We can use gram matrix to show.

→ Examples of kernel fn $k(x_i, x_j) = x_i^T x_j \hookrightarrow$ linear

↳ polynomial
↳ gaussian

↳ sigmoid

→ Kernel fn performs the function of dot product in non linear SVM.

→ Multi class SVM

↳ decision function for each class.

↳ ~~one~~ "one against rest", "one against one".

→ Advantages of SVM

↳ effective in high dimension

↳ good space complexity

↳ custom kernel fn

↳ good decision boundary with large margin

→ Weakness of SVM

↳ sensitive to noise

↳ SVM underperform in small train than test data

↳ no probability explanation

↳ not effective in large dataset

Q Consider hyperplane $y = x_1 - 2x_2$
predict class for

- (i) (1, 0)
- (ii) (1, 1)

Ans distance $w^T x + b$

$$\text{here, } x_1 - 2x_2 = w^T x \therefore w = \begin{bmatrix} 1 \\ -2 \end{bmatrix}$$

$$\therefore \text{(i)} \underbrace{\begin{bmatrix} 1 \\ -2 \end{bmatrix}}_{w^T x_1} \begin{bmatrix} 1 \\ 0 \end{bmatrix} = 1 > 0 \therefore \text{+ve class}$$

$$\text{(ii)} \underbrace{\begin{bmatrix} 1 \\ -2 \end{bmatrix}}_{w^T x_2} \begin{bmatrix} 1 \\ 1 \end{bmatrix} = -1 < 0 \therefore \text{-ve class.}$$

Step-1 : find support vectors

Step-2 : augment bias term

e.g. Support vectors : (2, 3) for +ve class = (S₁)
(4, 5) for -ve class = (S₂)

let bias is 1

$$\therefore \begin{bmatrix} 2 \\ 3 \\ 1 \end{bmatrix}, \begin{bmatrix} 4 \\ 5 \\ 1 \end{bmatrix}$$

Now, we have to find Lagrangian multipliers

$$\alpha_1 S_1 S_1 + \alpha_2 S_2 S_2 + \cancel{\alpha_1 S_1 S_2 + \alpha_2 S_2 S_1} = +1 \quad (\cancel{\alpha_1 S_1 S_2 + \alpha_2 S_2 S_1} \text{ is common l+ve})$$

$$\alpha_1 S_2 S_1 + \alpha_2 S_1 S_2 = -1 \quad (S_2 \text{ is common l-ve})$$

$$\alpha_1 \begin{bmatrix} 2 \\ 3 \\ 1 \end{bmatrix} \cdot \begin{bmatrix} 2 \\ 3 \\ 1 \end{bmatrix} + \alpha_2 \begin{bmatrix} 2 \\ 3 \\ 1 \end{bmatrix} \cdot \begin{bmatrix} 4 \\ 5 \\ 1 \end{bmatrix} = 1 \therefore \alpha_1 (4 + 9 + 1) + \alpha_2 (8 + 15 + 1) = 1$$

$$\therefore 14\alpha_1 + 24\alpha_2 = 1$$

$$24\alpha_1 + 46\alpha_2 = -1$$

Now we can calculate Lagrangian multipliers α_1 and α_2 .

$$w = \sum_i \alpha_i \cdot \bar{s}_i = \alpha_1 \begin{bmatrix} 2 \\ 3 \\ 1 \end{bmatrix} + \alpha_2 \begin{bmatrix} 4 \\ 5 \\ 1 \end{bmatrix} = \text{a vector on putting } \alpha_1 \text{ and } \alpha_2$$

Let $w = \begin{bmatrix} w_1 \\ w_2 \\ b \end{bmatrix}_{3 \times 1}$ \therefore Here we will get actual value of b

\therefore eqn: $y = w^T x + b$ \leftarrow substitute w & b here
 This is the eqn of hyperplane.

If A is a set then $|P(A)| = 2^{|A|}$

DOMS Page No.
Date 22 / 10 / 24

Clustering

Hard Soft

each data point belongs only to single class

a membership function is defined such that point belongs to multiple classes

→ Crisp set :-

~~Set~~ $\Rightarrow A = \{a_1, a_2, a_3, \dots, a_n\}$ is a set of events

$\Rightarrow A = \{x | p(x), p: \text{property}\}$

\Rightarrow using char. function,

$$\mu_A(x) = \begin{cases} 1, & \text{if } x \in A \\ 0, & \text{otherwise} \end{cases}$$

e.g. ~~Set~~ $A = \{1, 2, 3, 4, 5, 6, 7, 8\}$ and property 'p' says even, then $A = \{2, 4, 6, 8\}$

→ Fuzzy set : boundary is not crisp.

→ Membership Function

⇒ Fuzzy sets : sets with imprecise or vague boundaries

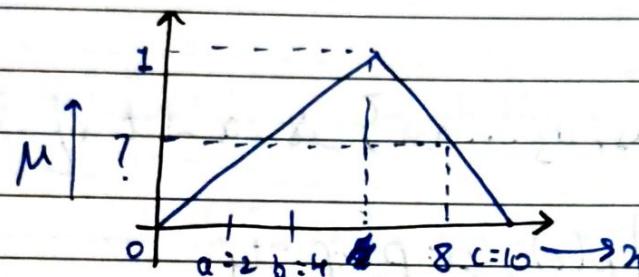
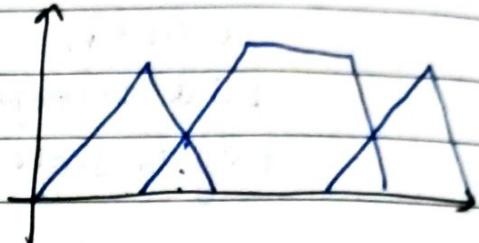
↳ It is a potential tool for handling imprecision and uncertainties

crisp set

Fuzzy set is more general concept of classical set.

How to find the membership function

e.g. fuzzy set will have overlapping areas



Q Calculate membership function μ at $x = 8$

Ans General formula

$$\mu_{\text{triangle}} = \max \left[\min \left(\frac{x-a}{b-a}, \frac{c-x}{c-b} \right), 0 \right]$$

$$\therefore \mu = \max \left[\min \left(\frac{8-2}{6}, \frac{10-8}{6} \right), 0 \right]$$

$$= \max \left(\min (3, \frac{1}{3}), 0 \right) = \frac{1}{3}$$

For trapezium, compare slope with 1 as well

and $d = 10$

$$(1) \quad a=2, b=4, c=8, d=10, \text{ at } x = 3.5$$

$$\text{Ans} \quad \mu = \max \left[\min \left[\frac{(x-a)}{b-a}, \frac{(d-x)}{d-c} \right], 0 \right]$$

$$= \max \left(\min \left(\underbrace{\frac{1.5}{2}, \frac{6.5}{2}}_{1.5}, 0 \right) \right)$$

$$= \max \left(\frac{1.5}{2}, 0 \right) = \cancel{0.75} \times \frac{3}{4} = 0.75$$

Fuzzy C Mean Clustering (FCM)

↳ It is example of soft clustering

↳ (i) membership fn

Steps : (i) plot to get an idea

(ii) assign random prob. for each point

(iii) choose m = fuzziness parameter

↳ usually taken as 2.

Centroid :

$$V_{ij} = \frac{\sum_{k=1}^n r_{ik}^m \cdot x_k}{\sum_{k=1}^m r_{ik}^m}$$

n = no. of cluster
 r = fuzziness membership value
 x = data point
 m = fuzziness parameter

(1,3) (2,5) (4,8) (7,9)

	C_1	0.9	0.7	0.2	0.1
C_2	0.1	0.3	0.8	0.9	

Ans $V_{11} = \frac{r_{11}^2 x_1 + r_{12}^2 x_2 + r_{13}^2 x_3 + r_{14}^2 x_4}{r_{11}^2 + r_{12}^2 + r_{13}^2 + r_{14}^2} = \frac{3.8}{1.9} = 2$

$$= \frac{2.02}{1.35} = 1.496$$

$$Y_{12} = \frac{r_{11}^2 y_1 + r_{12}^2 y_2 + r_{13}^2 y_3 + r_{14}^2 y_4}{r_{11}^2 + r_{12}^2 + r_{13}^2 + r_{14}^2} = \frac{5.29}{1.35} = 3.918$$

$$Y_{21} = \frac{8.42}{1.55} = 5.43$$

$$Y_{22} = \frac{12.89}{1.55} = 8.316$$

\therefore Centroids, $C_1: (1.496, 3.918)$
 $C_2: (5.43, 8.316)$

(iv) Calculate the distance of each data point from
 centroids of C_1 and C_2

$$D_{11} = 1.043$$

$$D_{31} = 4.789$$

$$D_{12} = \cancel{1.94} 6.92$$

$$D_{32} = 1.464$$

$$D_{21} = \cancel{4.789} 1.194$$

$$D_{41} = 7.491$$

$$D_{22} = \cancel{7.491} 4.77$$

$$D_{42} = 1.712$$

(v) Update membership values

$$Y_{ki} = \left(\sum_{j=1}^n \left\{ \frac{d_{kj}}{d_{kj}} \right\}^{\frac{1}{m-1}} \right)^{-1}$$

\Rightarrow Repeat until difference b/w new and old Y is $\neq 0.01$.

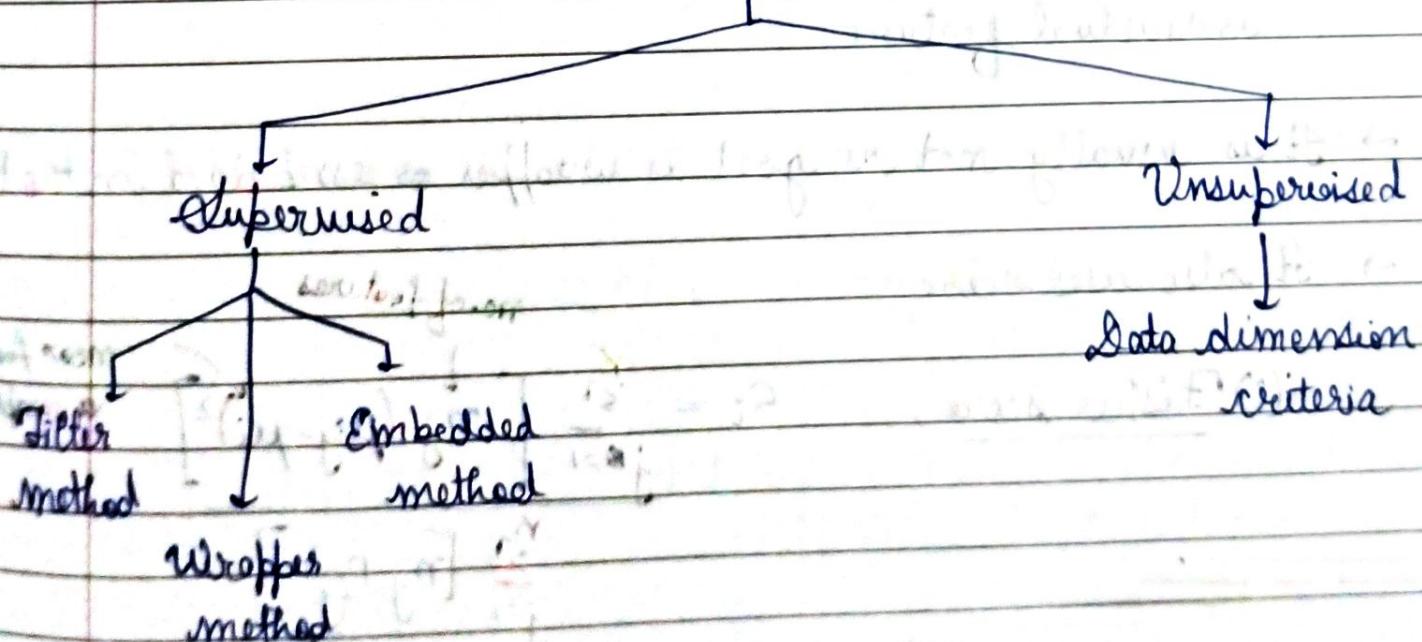
Feature Selection

- Raw data → Sensors → Features → preprocessing
- Feature selection
- After collecting data, we don't know which features are more important and which are irrelevant or redundant.
- ↳ Irrelevant features increase complexity & are not wanted
 - ↳ Redundant features are only duplicates and dont provide any value

e.g. Price | Address | No. of rooms | House age | Owner | Area of rooms

These, here, while determining price, 'Owner' is irrelevant.
Also, 'No. of rooms' is redundant due to already provided information 'Area of rooms'.

Feature Selection Models



Filter method

→ Some considerations for feature selection :-

1. Classification frequency must be same as earlier.
2. Classification of most data points will remain same as earlier.

↳ We are reducing features to :-

- ↳ To reduce time & space complexity.
- ↳ To avoid overfitting.

Filter method :-

→ Univariate

Multivariate

↳ Each feature is independent

↳ we take a batch of features.

↳ so, we consider individual feature

→ It is usually not as good as wrapper or embedded method

→ It also uses score.

(i) Fisher score,

$$S_i = \sum_{j=1}^k [n_j (\bar{x}_{ij} - \mu_i)^2]$$

↑
no. of features
mean for a sample

$$\sum_{j=1}^k [n_j p_{ij}^2]$$

↑
variance

$i \rightarrow$ feature $\mu \rightarrow$ mean
 $j \rightarrow$ class $p \rightarrow$ variance

Data table :-

Instances / samples	x_1	x_2	x_3	\dots	x_k
$\rightarrow s_1$					
$\rightarrow s_2$					
$\rightarrow s_n$					

Entropy : gives mutual information

↳ Randomness gives hidden patterns : more the randomness, higher the information

∴ Information gain, $IG_i(F_i, c) = H(F_i) - H(F_i/c)$

$$\text{where, } H(F_i) = \left[- \sum_j P(x_j) \log(P(x_j)) \right]$$

$$H(F_i/c) = \left[- \sum_k P(c_k) \sum_j P(x_j/c_k) \log_2(P(x_j/c_k)) \right]$$

In such cases, we determine a threshold value and determine

$$SU_i(F_i/c) = \frac{\alpha IG_i(F_i/c)}{H(F_i) + H(c)}$$

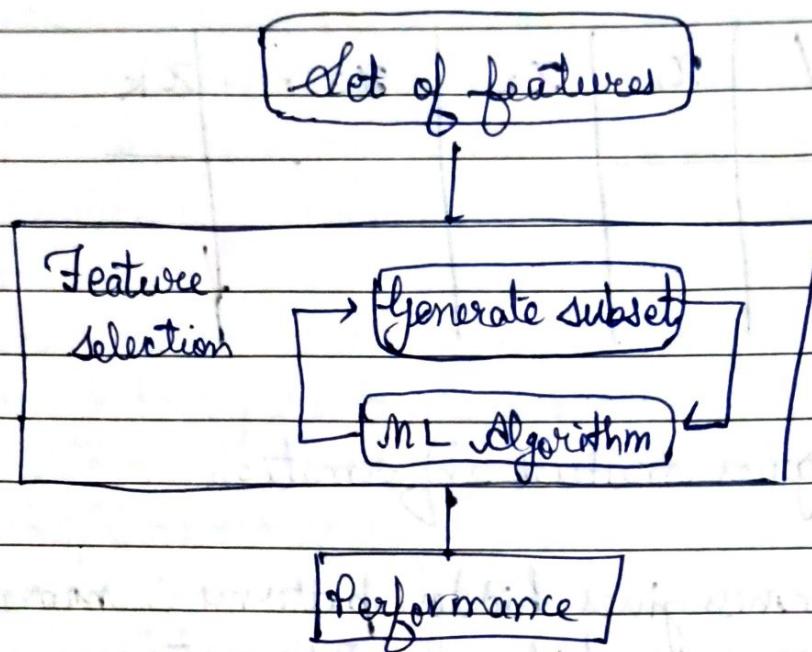
if $SU_i(F_i/c) > T$ where T is threshold, we call it ^{pre}dominant value.

∴ Such features are selected.

⇒ Filter methods are independent of specific classifier.

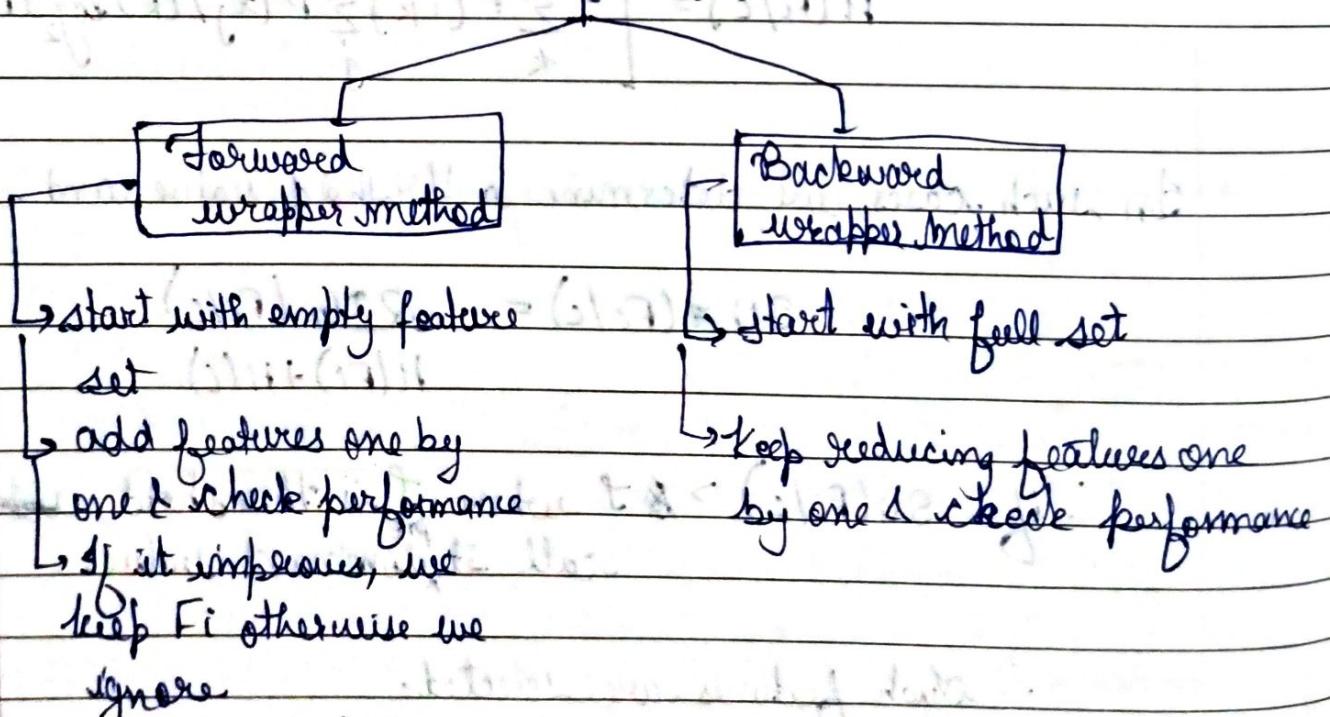
→ also called greedy approach

Wrapper methods



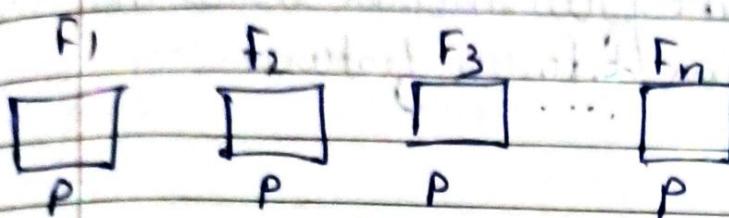
If performance improves : we take that feature
 remains same : redundant feature : ignore
 reduces : irrelevant feature : ignore

→ wrapper methods are of two types



∴ for n features, in first iter
 we take n models.

Forward wrapper



① n models for each feature

② Let us say F_i is best performance
∴ now we make $(n-1)$ models

by combining $F_i F_j$ where $j \neq i$

③ Again take the model with max value of P . Let us say $F_i F_j$

take $F_i F_j$ as fixed :- make $(n-2)$ models

$F_i F_j F_1, F_i F_j F_2, \dots, F_i F_j F_{n-2}$

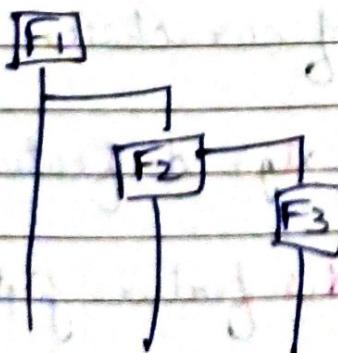
④ Again take model with max performance

Repeat until no more combination is left.

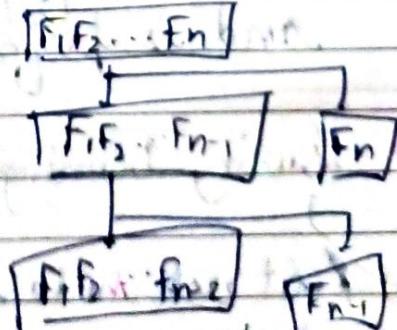
⇒ Hence we receive a subset of best performance model.

→ We take a threshold value to stop when sufficient features have been obtained.

Forward method



Backward method

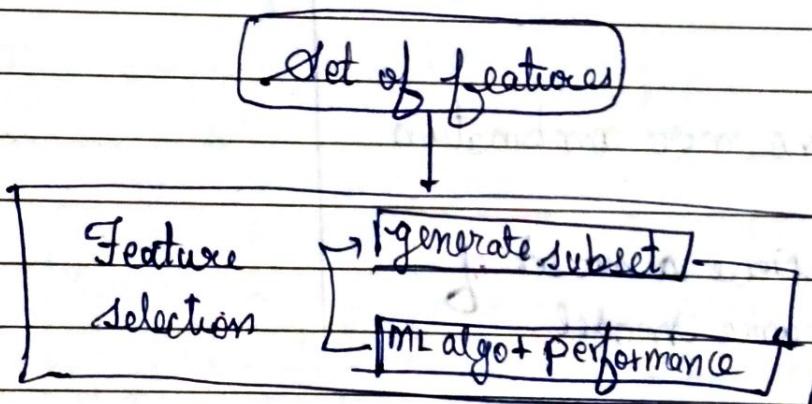


→ Wrapper method uses relation b/w various features using statistical methods, unlike filter method which does not take into account relation between features.

Statistical methods

- | (i) → correlation: If x_1 and x_2 are highly correlated, coefficient one of them is redundant.
- | (ii) chi-square test
- | (iii) ANOVA
- | (iv) Variance thresholding

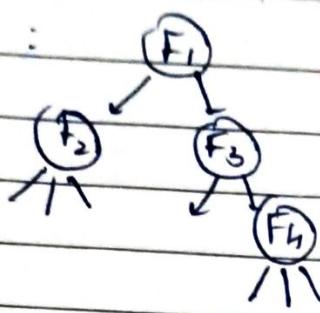
Embedded method - (intrinsic method)



→ We are calculating performance as well ~~as well~~ within model

- ① Selection of feature happens simultaneously ~~with~~ with and is performed by ML algs of our choice.
- ② During training, some steps of ML algo do feature selection
- ③ At each node split they choose better feature just like decision tree.

Decision tree :



} → initially, features are random
→ matrices are calculated
which determine which
feature is better at a
particular level.

- ⇒ Filter methods are less used due to time & space complexity
- ⇒ wrapper & embedded methods are widely used.
- ⇒ embedded is considered best method due to its efficiency