

Introduction to Pattern Recognition

→ Patterns are recognised by the help of algorithm used in machine learning. Recognising patterns is the process of classifying the data based on the model that is created by training data, which then detects patterns & characteristics from the patterns.

→ Pattern Recognition is the process which can detect different categories and get information about particular data.

→ Features of Pattern Recognition:

- 1) Pattern Recognition learns from the data.
- 2) Automatically recognizes the patterns even when partially visible.
- 3) Should be able to recognise pattern which are familiar.
- 4) The pattern should be recognized from different angles & shapes.

→ Applications of Pattern Recognition:

- | | |
|-----------------------|-------------------------|
| 1) Computer Vision | 5) Biometrics. |
| 2) Engineering | 6) Geology. |
| 3) Speech Recognition | 7) Intrusion Detection. |
| 4) Image Recognition | 8) Civil Administration |

* Design Principles of Pattern Recognition:

In pattern recognition system, for recognising patterns or structures 2 basic approaches are used which can be implemented in different techniques. These are

- 1) Statistical Approach
- 2) Structural Approach

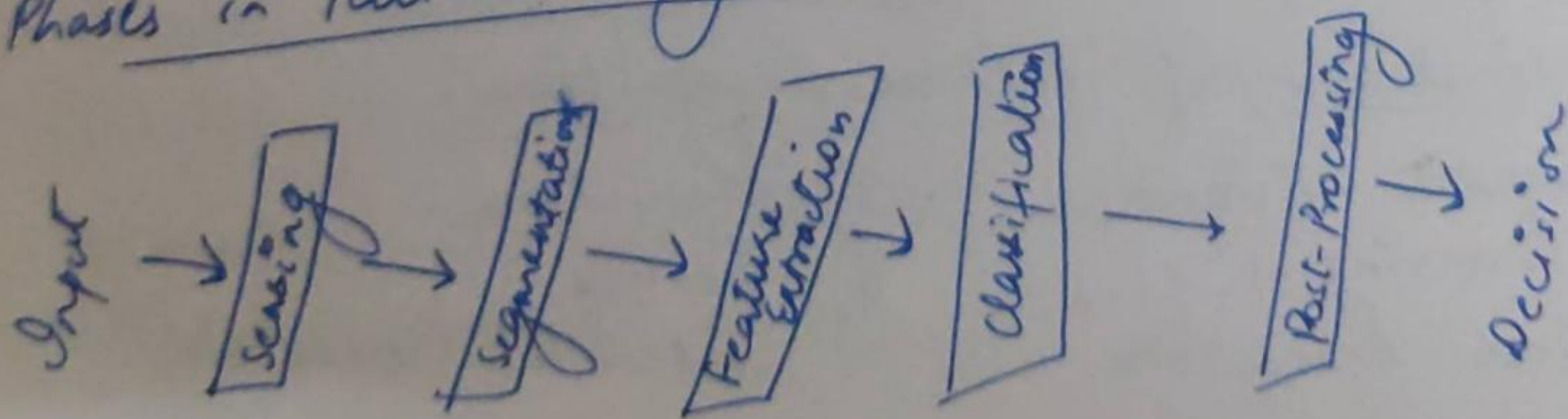
1) Statistical methods are mathematical formulas, models & techniques that are used in statistical analysis of raw research data. The application of statistical method extract information from research data and provides different ways to assess the robustness of research output.

The 2 main statistical methods are:

- a) Descriptive - summarises data from sample ^{using indices.}
- b) Inferential - draws conclusion from data that are subject to random variation.

2) Structural Approach is a technique wherein the learner masters the pattern of sentence structures. are different arrangements of words in one accepted style or other. Some of them are sentence or phrase pattern, formulas and idioms.

* Phases in Pattern Recognition:



Statistical Pattern Recognition

- It refers to the use of statistics to learn from examples.
- It means to collect observation, study & digest them in order to infer general rules / concepts that can be applied to the new, unseen observation.
- It relates to the use of statistical techniques for analysing data measurement in order to extract info. and justified decision.
- Statistical approach in PR is represented in D features in d dimensional space as a point.
- Objective to establish decision boundaries in feature space which separate pattern of different classes.
- Discriminate Analysis based approach for classification and uses mean squared error criteria.
- Construct the decision boundaries of the specified form.
- Statistical Approach follows up:
 - 1) Linear Discriminant Analysis (LDA)
 - 2) Principal Component Analysis (PCA)
 - 3) Posterior Probability Estimation
 - 4) Window-Density based classifier.
 - 5) Edited K-NN Rule.

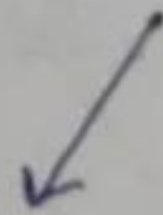
(PEM)

Parameter
Estimation
Method① Principal Component Analysis

There are collection of points in a real p -space are a sequence of p -direction vectors.

It is a process of computing the principal components & using them to perform a change of basis on the data.

- used in exploratory data analysis
- used in making predictive models
- commonly used for dimensionality reduction.



By projecting each data point onto only the first few principal components to obtain lower-dimensional data while preserving as much of the data's variation as possible.

- PCA is defined as orthogonal linear transformation that transform the data to a new co-ordinate system such that the greatest variances by some scaled projection of the data comes to lie on the first co-ordinate; similarly the second variance on the second co-ordinate.

Numerical Example: PCA

* Consider the dataset with 5 observation of Math, English & Arts marks; arranged in matrix A

$$A = \begin{bmatrix} 90 & 60 & 90 \\ 90 & 90 & 30 \\ 60 & 60 & 60 \\ 60 & 60 & 90 \\ 30 & 30 & 30 \end{bmatrix}$$

Matrix A

$$\bar{A} = [66 \quad 60 \quad 60]$$

mean of matrix A

- Compute the covariance matrix of the whole dataset

$$\text{cov}(X, Y) = \frac{1}{n-1} \sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})$$

find the covariance matrix of A, the result would be square matrix of $d \times d$ dimension.

→ Covariance Matrix

$$\begin{array}{c} M \quad E \quad A \\ M \begin{bmatrix} 504 & 360 & 180 \\ 360 & 360 & 0 \\ 180 & 0 & 720 \end{bmatrix} \\ E \\ A \end{array}$$

~~~~~  
least more  
variance variance

→ Compute Eigen values & Eigen vectors.

$$\det(A - \lambda I) = 0$$

$$\det \left( \begin{bmatrix} 504 & 360 & 180 \\ 360 & 360 & 0 \\ 180 & 0 & 720 \end{bmatrix} - \lambda \begin{bmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \end{bmatrix} \right) = \det \begin{bmatrix} 504 - \lambda & 360 & 180 \\ 360 & 360 - \lambda & 0 \\ 180 & 0 & 720 - \lambda \end{bmatrix} = \det \begin{bmatrix} \lambda & 0 & 0 \\ 0 & \lambda & 0 \\ 0 & 0 & \lambda \end{bmatrix}$$



Determinant:  $-\lambda^3 + 1584\lambda^2 - 641520\lambda + 25660800 = 0$

on solving we get -

Eigen values  $\lambda \approx 44.82; 629.11; 910.069$  (approx.)

- Solving eigenvector we get -

$$\begin{pmatrix} -3.75100 \\ 4.28441 \\ 1 \end{pmatrix} \begin{pmatrix} -0.50494 \\ -0.67548 \\ 1 \end{pmatrix} \begin{pmatrix} 1.05594 \\ 0.69108 \\ 1 \end{pmatrix}$$

- Sort the eigenvectors by decreasing eigenvalues  
& choose  $k$  eigenvectors with the largest eigenvalue  
to form a  $d \times k$  dimensional matrix  $W$

- After sorting the eigenvalues in decreasing order,

$$\begin{pmatrix} 910.06995 \\ 629.11039 \\ 44.81966 \end{pmatrix}$$

- Eigen vectors corresponding to 2 max. eigenvalues are:

$$W = \begin{bmatrix} 1.05594 & -0.50494 \\ 0.69108 & -0.67548 \\ 1 & 1 \end{bmatrix}$$

- Transform the samples onto the new subspace.

$$y = W' \times x$$

↓  
transpose of the matrix  $W$



PEM ② Linear Discriminant Analysis

→ The objective of LDA is to perform dimensionality reduction while preserving as much of the class discriminatory info. as possible.

→ most famous example of dimensionality reduction is - "PCA"; but question is how to utilize the label info. in finding info. projections? Thus, for that purpose -

Fisher-LDA consider maximizing the following obj.

$$J(w) = \frac{w^T S_B w}{w^T S_W w}$$

B/w classes scatter matrix.

↪ within classes scatter matrix.

→ LDA approximates the Bayes classifier rule by modeling conditional class densities as multivariate normal. For classes  $c \in 1 \dots K$  let a feature vector  $x \in \mathbb{R}^p$ . This can be expressed:

$$P(X=x | C=j) = N(\mu_j, \Sigma)$$

, each class  $j$  has its own mean  $\mu_j \in \mathbb{R}^p$

but classes together share a covariance matrix  $\Sigma \in \mathbb{R}^{p \times p}$

→ Various steps of LDA process are:

- ① Centre Data
- ② Calculate Estimators
- ③ Sphere variables
- ④ Project using  $P_m$  - (covariance matrix  $\Sigma$ )



## Example: Fisher Linear Discriminant Analysis

Data: 1) class 1 has 5 samples  $C_1 = [(1,2), (2,3), (3,3), (4,5), (5,5)]$   
2) " 2 has 6 "  $C_2 = [(1,0), (2,1), (3,1), (3,2), (5,3), (6,5)]$

Step

1: Arrange data in 2 separate matrices:

$$C_1 = \begin{bmatrix} 1 & 2 \\ 2 & 3 \\ 3 & 3 \\ 4 & 5 \\ 5 & 5 \end{bmatrix}$$

$$C_2 = \begin{bmatrix} 1 & 0 \\ 2 & 1 \\ 3 & 1 \\ 3 & 2 \\ 5 & 3 \\ 6 & 5 \end{bmatrix}$$

\* PCA performs poorly on this due to direction of largest variance is not helpful for classification.

Step 2: Compute the mean for each class

$$\mu_1 = \text{mean}(C_1) = [3 \quad 3.6]$$

$$\mu_2 = \text{mean}(C_2) = [3.3 \quad 2]$$

Step 3: Compute the scatter matrices  $S_1$  &  $S_2$  for each class

$$S_1 = 4 * \text{cov}(C_1) = \begin{bmatrix} 10 & 8.0 \\ 8.0 & 7.2 \end{bmatrix}$$

$$S_2 = 5 * \text{cov}(C_2) = \begin{bmatrix} 17.3 & 16 \\ 16 & 16 \end{bmatrix}$$

Step 4: Within the class scatter,

$$S_w = S_1 + S_2 = \begin{bmatrix} 27.3 & 24 \\ 24 & 23.2 \end{bmatrix}$$

No need to solve for eigenvalues, as it has full rank

Step 5: Inverse  $(S_w^{-1}) = \begin{bmatrix} 0.39 & -0.41 \\ -0.41 & 0.47 \end{bmatrix}$

$$v = S_w^{-1}(\mu_1 - \mu_2) = \begin{bmatrix} -0.79 \\ 0.89 \end{bmatrix}$$



→ As long as the line has the right direction its exact position doesn't matter

Step 6: Compute the actual 1D vector  $y$  for each class.

$$y_1 = v^t c_1^t = [-0.65 \quad 0.73] \begin{bmatrix} 1 & 2 & 3 & 4 & 5 \\ 2 & 3 & 3 & 5 & 5 \end{bmatrix} = [0.81 \quad -0.4]$$

$$y_2 = v^t c_2^t = [-0.65 \quad 0.73] \begin{bmatrix} 1 & 2 & 3 & 3 & 5 & 6 \\ 0 & 1 & 1 & 2 & 3 & 5 \end{bmatrix} = [-0.65 \quad -0.25]$$

— X —  
 —————  
 —————



# Classification Technique :

## Nearest Neighbour Rule

→ It is the simplest decision procedures - NN Rule.

③ Easy to get and interpret output. ← ② These classifiers essentially involve finding the similarity b/w the test pattern & every pattern in training set. ← ① It classifies a sample based on the category of its nearest neighbours.

↓  
④ calculation time is less  
↳ more predictive power

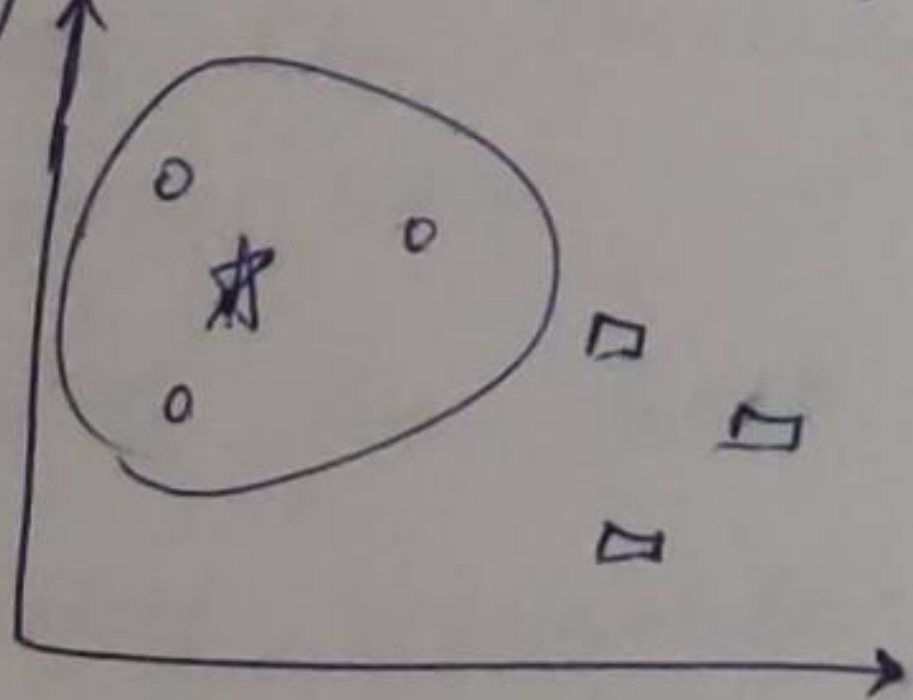
→ It can be used for both the classification & regression predictive problems.

→ for eg: Find the class of star.

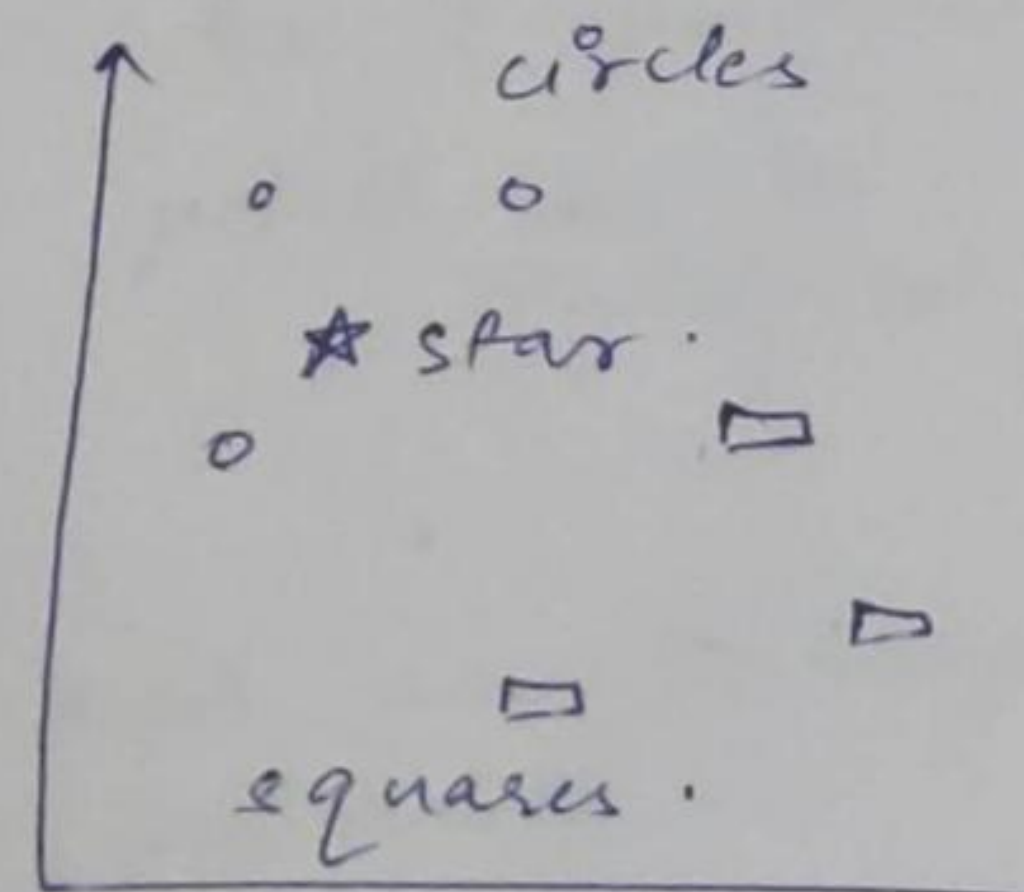
can be circle square

say  $k=3$  for first time.

①



• As all the circles are nearest to the star, so, the confidence level is high.



• so, the choice becomes very obvious.

• As the observations become constant, we can make boundaries of each class.

• Gradually, after each iteration the boundaries gets specific and smoother and clear; sorted form.

• Training depends on no of observation.

• Get the predicted class at end.



## Lecture: 41.

### Bayes Classifier.

- Important for statistical learning.
- A classifier is a machine learning model that is used to discriminate different object based on certain features.
- A Naive Bayes classifier is a probabilistic mlc learning model that's used for classification task.
- Base is on Bayes theorem:

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

- Using this Bayes theorem, we can find the prob. of A happening given that B has occurred. Here B, is the evidence & A is the hypothesis.
- Here, the predictors are independent. It means presence of one feature doesn't affect the other.



## Support Vector Machine

→ It is a supervised learning technique that can be used for both classification & regression problem.

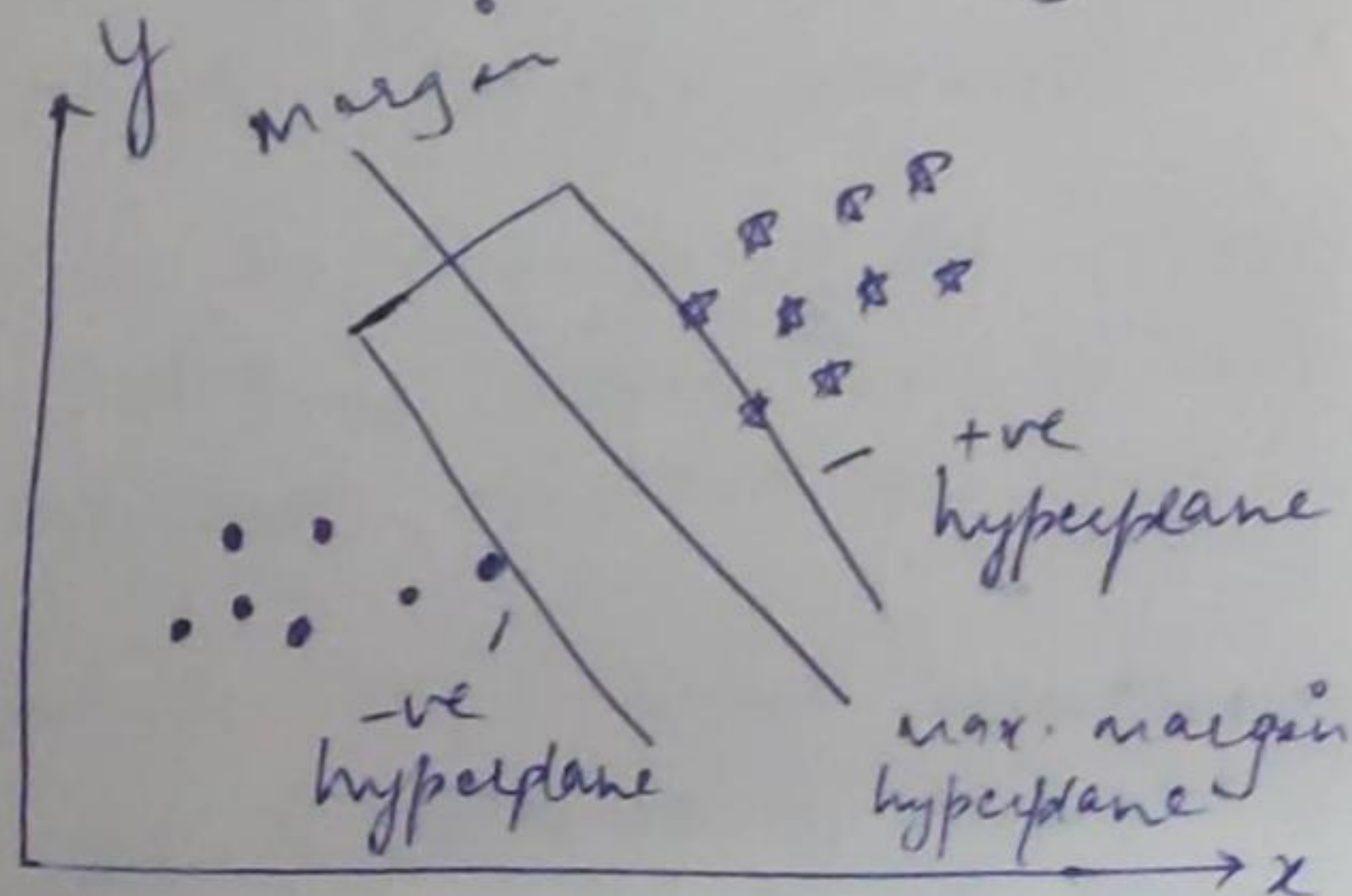
→ Basic Concepts:

- a) Support vectors - These are simply the co-ordinates of individual observation. These are near to hyperplane and influence the position & orientation of hyperplane.
- b) Hyperplane - These are the decision boundaries that help classify the data points.

c) Margin - This is the  $\perp^{\text{th}}$  from both hyperplane & more the margin more better machine. As it provides some re-inforcement so that future data points can be classified.

→ Working of SVM -

- 1) Take the linearly separable data and non-linearly separable.
- 2) Apply SVM on both
- 3) With the help of SVM kernel tricks to make it linearly separable. and then solve.
- 4) It takes data points in consideration & gives out a hyperplane, which divides the both classes.
- 5) Many hyperplanes can be generated but the best is the one which divided the classes from maximum distance.





## Lecture 43

### k-means clustering

- It is an iterative algorithm that tries to partition the dataset into  $k$ -pre-defined distinct non-overlapping sub-groups where each data points belong to only one group.
- It tries to make the intra-cluster data points as similar as possible while keeping the cluster as different as possible.
- It assign data points to a cluster such that the sum of squared distance b/w the data points and cluster's centroid is at the minimum.
- The less variation we've within cluster the more homogeneous the data points are within the same cluster.

#### Working of k-means:

- 1) Specify the no of clusters  $k$ .
- 2) Initialise centroids by shuffling the dataset & then randomly selecting  $k$  data points for the centroids.
- 3) Keep iterating until there is no change to the centroid.  
i.e. the assignment of data point to cluster isn't changing.
  - compute the sum of squared distance b/w data points & centroid.
  - Assign each data point to the closest cluster.
  - compute the centroid for the cluster by taking the avg. of the all data points that belong to each cluster.