

MOD-IV

- (1) Hidden Markov Models (HMMs) for sequential classification :-
- Discrete Hidden Markov Models (DHMM)
 - Continuous Density Hidden Markov Models (CDHMM)

(2) Hidden Markov Models (HMMs)

- Statistical model used to describe sequential data (like speech, handwriting, DNA, etc.) where the system is modeled as a Markov process with hidden (unobservable) states.

- Real world Ex :-

→ Speech recognition.

→ Bioinformatics (DNA Sequences).

→ Handwriting recognition.

- Components of HMM:

Symbol.

$N \rightarrow$ No. of hidden states.

$M \rightarrow$ No. of possible observⁿ. (for discrete models)

$A = \{a_{ij}\} \rightarrow$ State transition probability matrix. (State i to j)

$B = \{b_i(k)\} \rightarrow$ Observ. probability matrix.

$\pi = \{\pi_i\} \rightarrow$ initial state distriⁿ.

Together, HMM is represented as :-

$$1. A \text{ HMM} = (A, B, \pi)$$

3 main problems solved by this :-

1) Evaluation problem → forward Algo.

2) Decoding Problem → Viterbi Algo.

3) Learning Problem → Baum-welch Algo.

Hidden Markov Model (DHMM)

③ Discrete

- observⁿ comes from a finite set:
(like words, phonemes, symbols).
- each state emits one of M discrete obs based
on a probability distribution.

Ex:- Let's say M = 3 (Symbols: A, B, C).

If you're in state S_i :-

Probability of A = 0.6

B = 0.3

C = 0.1

Then you emit one of them randomly
based on this distribution.

→ Modelling Steps :-

1. Define states & obs. symbols.
2. Initialize A, B and π.
3. use forward - Backward ~~or~~ to compute
req. probabilities.
4. Train using EM if model is unknown.

→ Applications → Text seq. classificⁿ.

→ speech recognition.

→ POS tagging.

④ Continuous Density Hidden Markov Models (CDHMM)

Observⁿ are continuous-valued vectors

each state emits observⁿ acc. to a continuous
prob. distrib., usually a Gaussian / mixture
of Gaussians (GMM).

Ex:- If you're in State S, your obs. is a
feature vector like :-

$$x = [2.3, -1.5, 0.7]$$

This is drawn from:

$$P(x|s_i) \sim N(\mu_i, \Sigma_i)$$

Or more generally:

$$P(x|s_j) = \sum_{m=1}^M c_{jm} N(x|\mu_{jm}, \Sigma_{jm})$$

c_{jm} → mixture weights

μ_{jm}, Σ_{jm} :- mean & covariance of m-th Gaussian in state

Real-world data like speech, images or sensor signals are not discrete - they are continuous.

CDHMMs provides a more realistic & powerful modelling tool for such data.

MOD-V

Dimension Reduction Methods

(1) Dimension Reduction Methods
Process of reducing the no. of I/P variables in a dataset while preserving as much information as possible.

- Why??
• removes irrelevant features.
• reduces computational cost.
• helps in visualizing high-Dimensional data.
• improves classification.

2 main Approaches:-

feature selection \rightarrow choose a subset of original features

feature extraction \rightarrow create new features from combination of old ones.

1) Fisher Discriminant Analysis

2) Principal component Analysis

(1) Fisher Discriminant Analysis

Goal:- find a new axis that maximizes class separability.

Key Idea:- Maximize b/w-class Variance
minimize within-class Variance.

How it works:-

1. compute mean of each class & overall mean.
2. compute b/w-class scatter & within-class scatter.
3. find projection w that maximizes the ratio:

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

$S_B \rightarrow$ b/w-class scatter matrix

$S_w \rightarrow$ within-class scatter matrix

Pros

- Supervised
- Great for classification task
- Better class separation than PCA
- Assumes normal distributed feature
- works best when class covariance is similar.

cons

② Principal component Analysis

Goal: Transform original data into a new coordinate system such that:-

- 1st new axis captures max. Variance
- Next axes capture remaining variance in decreasing order

How it works

i. Standardize data (mean = 0)

ii. Compute covariance matrix.

iii. Compute eigenvectors & eigenvalues

iv. Select top k eigenvectors \rightarrow new feature space

v. Transform data into new space

Pros:

- Reduces overfitting
- Captures max. Variance
- useful for visualizing data in 2D/3D.
- Principal comp. are.
- Not easily interpretable
- Assumes linear relationship.

cons

PCA vs LDA

• unsupervised

• Supervised

• Max. variance

• Maximize class separation

• No class labels needed

• Needs class labels

• Principal components

• Discriminant vectors

Parzen - window methods

Non-parametric technique

To estimate the probability density funct.
(PDF) of a random variable.

It helps us guess the shape of the data distribution, without assuming any specific shape.

- used in - pattern recognition,
- classification
- density estimation.

Imagine

You are trying to find how people are spread across a park.

You don't know the exact distribution.

So, you place a small window around every person & count how many people fall within that window.

If many people fall into a window, you say "density is high here".

If few people fall in, you say "density is low here".

This is Parzen - window idea!

Adv.

- Doesn't assume any fixed shape (flexible)
- Works for any distribution.

Disadv.

• Can be slow for large data sets.

• Choice of window size (h) is critical.

Too small \rightarrow noisy

Too big \rightarrow oversmooths.

⑤ K-Nearest Neighbour method

→ K-NN - supervised learning algo:
used for:- classifier
regression.

→ Simplest algo. in ML.

→ Ex:- Imagine you have a dataset of fruits.

<u>color</u>	<u>size</u>	<u>fruit</u>
Red	Big	Apple
Yellow	Big	Banana
Red	Small	Cherry

Now you get a new fruit:-

Color: Red.

Size: Medium

Using K-NN:

- You find its k-closest fruits in terms of color & size.
- Suppose $k=3$, and 2 of the 3 closest fruits are Apples \rightarrow the new fruit is classified as Apple.

Dist. measures

1) Euclidean dist.:

$$d(P, Q) = \sqrt{(P_1 - Q_1)^2 + (P_2 - Q_2)^2 + \dots + (P_n - Q_n)^2}$$

2) Manhattan Distance

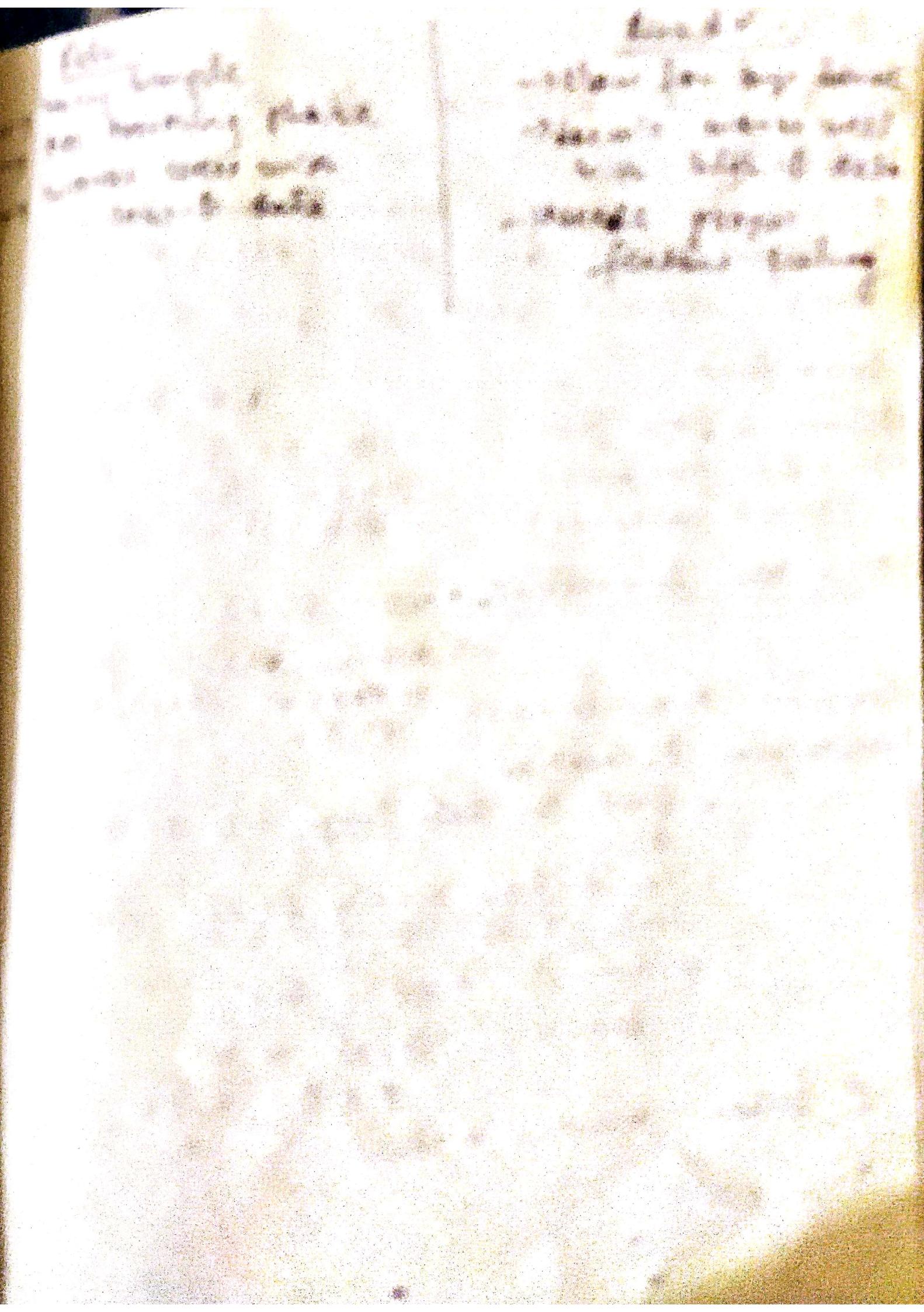
3) Minkowski Distance

Choosing K: - Small K (like 1 or 3)

Sensitive to noise & outliers.

Large K (like 10. or 15.):

Smoothen decision boundaries
but may ignore local patterns.



Non Parametric techniques for density estimation:-

- Density Estimation \rightarrow Process of estimating the probability distribution (PDF) of a random variable based on given data.

Parametric

- Assumes a fixed form
- estimates parameters (mean, variance)

Non-Parametric

Makes no assumption about the distribution shape

Why Non-Parametric \rightarrow more flexible

\rightarrow works better for unknown distribution

Common Non-Parametric Methods :-

1. Histogram Estimation:-

- . divide the data range into eq-sized bins
- . count how many points fall into each bin
- . Estimate density by:-

$$\text{Density} = \frac{\text{count in bin}}{n \times \text{Bin width}}$$

2. Parzen window method } module - 5
3. k-NN