

Basics of Pattern Recognition:

What is Pattern Recognition?

- Science of identifying patterns & regularization in data
- Involves classifying IP data into identifiable categories based on key features/patterns.

Pattern →

Set of attributes/features - help in identifying/categorizing data

Ex: A handwritten digit like "5" is a pattern.

The strokes & curves are features.

Recognition →

Labeling / classifying a given pattern into a known class.

Ex - Recognizing whether a shape is circle/pattern

Classifier →

An algorithm/model that assigns labels to IP patterns.

Ex - Decision Tree, K-NN, SVM, Neural Networks

Types of Pattern Recognition →

Supervised learning → Training with labeled data.
The system learns from ex.

Unsupervised learning → No labels provided. The system groups data based on similarity.

Reinforcement learning → learning from feedback (reward/punishment) in an environment.

Steps in PR process:

1. Data collection → Gathering IP data
2. Preprocessing → Cleaning & transforming data
(eg. noise removal, normalization)
3. feature Extraction → Selecting imp. info. from raw data.
4. Classification → Using algo. to assign categories to the data.
5. Post-processing → Improving / interpreting the output
(e.g. minority voting?)

Application of PR →

- Handwriting recognition (OCR)
- Face Recognition
- Speech Recognition

common Algorithms used →

K-Nearest Neighbor (K-NN)

Support Vector Mle (SVM)

Naive Bayes

Decision Trees

Artificial Neural Networks
Clustering (eg. K-Means - Unsupervised Learning)

features selection vs. feature extraction

- | | |
|--|--|
| <ul style="list-style-type: none"> - selects subset of original features - reduces dimensionality & remove irrelevant features. - easy to interpret - techniques:- filter (chi-sq.) - less intensive (comput. cost) | <ul style="list-style-type: none"> - creates new features from original data - create new transforming original features into a new set of features. - Harder to interpret - techniques - Autoencoders - more intensive |
|--|--|

Ex:- Reducing 100 features to 10 using PCA

Ex:- Selecting top 10 features based on correlation scores.

MOD-II

- ① Bayesian Decision Theory
- fundamental statistical approach to the problem of pattern classification.
 - uses probability & Bayes' Theorem
 - to minimize risk (error) of making wrong decisions.
 - Key components:-

Classes (w_1, w_2, \dots, w_n): The categories / labels.
Prior Probability $P(w)$: The prob. that a random sample belongs to class w before seeing the data.

- Likelihood $P(x|w)$: The probability of observing a feature vector x given class w .
- Posterior Probability $P(w|x)$: Prob. that the pattern belongs to class w given feature x (calculated using Bayes' theorem).

Bayes' Theorem

$$P(w_i|x) = \frac{P(x|w_i) \cdot P(w_i)}{P(x)}$$

where:-

$$P(x) = \sum_j P(x|w_j) \cdot P(w_j)$$

Bayesian Decision Rule (Min. Error)

- choose the class w_i for which $P(w_i|x)$ is max
- or, choose w_i if $P(w_i|x) > P(w_j|x), \forall j \neq i$

② Classifiers:-

- used to assign IP patterns (like images, text) to predefined categories / classes.
- most classifiers operate in supervised learning setting, where the model is trained using labeled data.
- include K-NN, SVM, Decision Trees, Naive Bayes Neural Networks.

- classifier creates decision boundary, which separates diff. classes in the feature space.
- shape of this boundary depends on type of classifier.
- classifier effectiveness - measured using tools like confusion matrix, f1-score, recall, etc.

Discriminant Functions

- score fun. used to make decisions.
- is a formula / rule helps us decide which class a data point belongs to.
- helps in classifying objects / patterns by finding which class has highest score / value for a given IP.
- ex:- if you're classifying emails as spam or not spam, a discriminant fun. will give a score for each class. whichever score is higher determines the final label.
- Decision Rule :- Assign x to class w_i if $g_i(x) > g_j(x)$, for all $j \neq i$.
- for bayesian classifier:- $g_i(x) = P(w_i|x)$ or $\ln P(x|w_i) + \ln P(w_i)$

- ### Decision Surface
- boundary separates regions in the feature space where diff. decision rules apply.
 - ex:- for 2D feature space, it could be a line (linear classifier) or curve (non-linear).
 - created by discriminant functions.
 - can be diff. shapes:- Line (2D), plane (3D) or complex curve depending on classifier (linear / non-linear).
 - ex:- if you're classifying fruits as apples or oranges based on color & size, the decision surface would be a line / curve that separates apples from oranges in a color-size graph.

③ Normal Density & Discriminant Functions

Normal Density

- Gaussian Distribution.
- Bell-shaped curve.
- Shows how data is spread out - Most values are around avg. (mean) & fewer are far away.
- It is used when we assume that the data for each class follows a natural pattern (ex:- cats vs dogs)
- N.D. formula - helps us calculate probability that a data point belongs to a class based on its features.

$$P(x|\mu_i) = \frac{1}{(2\pi)^{d/2} |\Sigma_i|^{1/2}} \exp \left[-\frac{1}{2} (x-\mu_i)^T \Sigma_i^{-1} (x-\mu_i) \right]$$

where:-

μ_i : mean vector

Σ_i : covariance matrix

d : dimensionality of x

- ex:- Height of stu. in class follows normal curve. most are near avg., few are very tall (short)

Discriminant function: for Gaussian

If covariance matrices are equal for all classes then discriminant becomes linear -

$$g_i(x) = w_i^T x + w_{i0} \quad \text{where: } (w_i = \Sigma_i^{-1} \mu_i)$$

$$\text{constant term: } w_{i0} = \frac{1}{2} \mu_i^T \Sigma_i^{-1} \mu_i + \ln P(w_i)$$

If covariance are different, it becomes a quadratic function of x .

④ Discrete features in Detail:
when features are discrete (categorical)
rather than continuous.

- ✓ use Probability Tables:
- Each feature has a finite no. of values.
- compute conditional probabilities directly.

$$P(x|w_i) = P(x_1, x_2, \dots, x_n | w_i)$$

If assuming independence of features (Naive Bayes)

$$P(x|w_i) = \prod_{j=1}^n P(x_j|w_i)$$

Example:

for Spam detection using words (discrete feature)

x_1 : "free" appears / not

x_2 : "Offer" appears / not

learn probabilities like $P(x_1 = 1 | w_{\text{spam}})$

use Naive Bayes classifier:

$$g_i(x) = \ln P(w_i) + \sum_j \ln P(x_j | w_i)$$

MOD-III

- ① Parameter Estimation Methods:
- Maximum Likelihood Estimation (MLE)
 - Gaussian Mixture Models (GMM)
 - Expectation-Maximization (EM) Algorithm
 - Bayesian Estimation

Parameter Estimation:

- process of estimating parameters (like mean, variance, etc.) of a probability distribution from observed data.
For ex:- estimating mean (μ) and variance (σ^2) of a normal distribution.

Types of Estimation:

1. Point estimation → gives a single value estimate (e.g. Sample mean).

2. Interval estimation → gives a range of values (confidence interval).

Common estimators

- Sample mean (\bar{x}) estimates pop. mean (μ)
- Sample variance (s^2) estimates pop. variance (σ^2)
- Sample proportion (\hat{p}) estimates pop. proportion (p)

→ It helps us make decisions & predictions when we don't know the true values for the entire population.

a) Maximum Likelihood Estimation

- to estimate unknown parameters of a probability distribution by maximizing the likelihood that the observed data came from that distribution.

→ Likelihood - probability of data given certain parameters.

In MLE, we try to maximize this likelihood.

→ Steps → write the likelihood function.

→ Take the log (called log-likelihood)

→ Differentiate the log-likelihood

→ Solve the eqn. by setting derivative = 0

→ gives MLE estimate.

Common ex:- coin Toss Ex. (Bernoulli Trial)

- Suppose we toss a coin 10 times, & get 7 heads and 3 tails.
- Let p = probability of head
- Likelihood: $L(p) = p^7 \cdot (1-p)^3$
- Take log, differentiate, and solve \rightarrow get MLE
MLE \rightarrow widely used in

ML, Statistics, Data Science, PR.

Properties of MLE:

consistent, efficient,

Limitations: - can be difficult for complex distribution

Real-life Analogy: -

Imagine guessing someone's weight by trying diff. values & seeing which guess best explains how their clothes fit \rightarrow MLE is similar.

(b) Gaussian Mixture Models

- way to group or cluster data that assumes data comes from a mixture of several normal (Gaussian) distributions.
- Each grp. (or cluster) is represented by bell-shaped curve (Gaussian distribution).
- Gaussian distribution is bell-shaped curve often seen in statistics.

Defined by Mean (μ): center of curve.

Standard Deviation (σ): width of the curve.

- Why mixture? bcz. in real life, data is often a mix of diff. groups. Group stays \rightarrow Let's assume data comes from several Gaussian distributions combined together.

- Ex. - Imagine you're looking at the heights of people in a room that has:

Adults, Teenagers, Kids.

Each grp. will have a diff. avg. height so, GMM would try to separate the data into 3 groups each with its mean & std. deviation.

Application \rightarrow Image Segmentation
Speech Recognition.

(c) Expectation-Maximization method (EM)

- to find missing / hidden info. in data by repeating 2 steps:-
 - Guessing
 - Improving the guess.
- It is mainly used in problems where some data is hidden / incomplete, like in GMM.
- Ice-Cream Analogy:
Imagine 2 ice-cream mics are making cones. You see only the cones but you don't know which mic made which cone. You want to guess!
 - The avg. size of cones made by each mic.
 - Which mic likely made which cone?But you don't know which cone came from which mic - that's the hidden data. EM helps you figure out that.
- 2 main steps in EM.
 - i) Expectation Step (E-step) - "Guess who did what"
 - calculate probability
 - ii) Maximization Step (M-step) - "Update your guess"
 - recalculate best possible parameters.

- Repeat doing E-step & M-step again & again until the results stop changing much.

- Applications \rightarrow GMM
 - \rightarrow Clustering
 - \rightarrow Image Processing

(d) Bayesian Estimation

→ Method of estimating unknown values (called parameters) using :- what we already know (prior knowledge)
New data we observe (evidence)
It uses Bayes' Theorem to update our belief.

→ Weather forecast (ex)

- You believe there's 70% chance it will rain tomorrow (prior belief)
- You now check the weather app, & it shows dark clouds (new evidence).

Using Bayesian estim., you update your belief based on this new evidence.

Now, you might say:- now I believe there's 90% chance it will rain.

→ Formula (Bayes Th.)

$$\text{Posterior} = \frac{\text{Likelihood} \times \text{Prior}}{\text{Evidence}}$$

Prior → what you believe before seeing the data.

Likelihood → how likely its the new data, assuming your belief is true.

Posterior → updated belief after seeing the new data.

Evidence → A normalizing factor

→ Application: - Spam filters,
ML, Robotics & AI.

MOD-IV

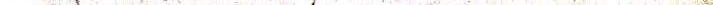
Module-IV
① Hidden Markov Models (HMMs) for sequential classification:-

- 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 examples:
 - Speech recognition
 - Bioinformatics (DNA Sequences)
 - Handwriting recognition

- Components of HMM : (with diagram) 

$N \rightarrow$ No. of hidden States.

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

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

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

$\Pi = \{ \Pi_i \}_{i=1}^n \rightarrow$ initial state distribution.

Together, HMM is represented as:

$$A \in \mathcal{A} \text{ defined by } (A, B, \pi)$$

3 main problems solved by this :-

- 1) Evaluation problem \rightarrow forward Algo.
 - 2) Decoding Problem \rightarrow viterbi Algo.
 - 3) Learning Problem \rightarrow Baum-welch

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 → choose a subset of original features

feature extraction → 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 are 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

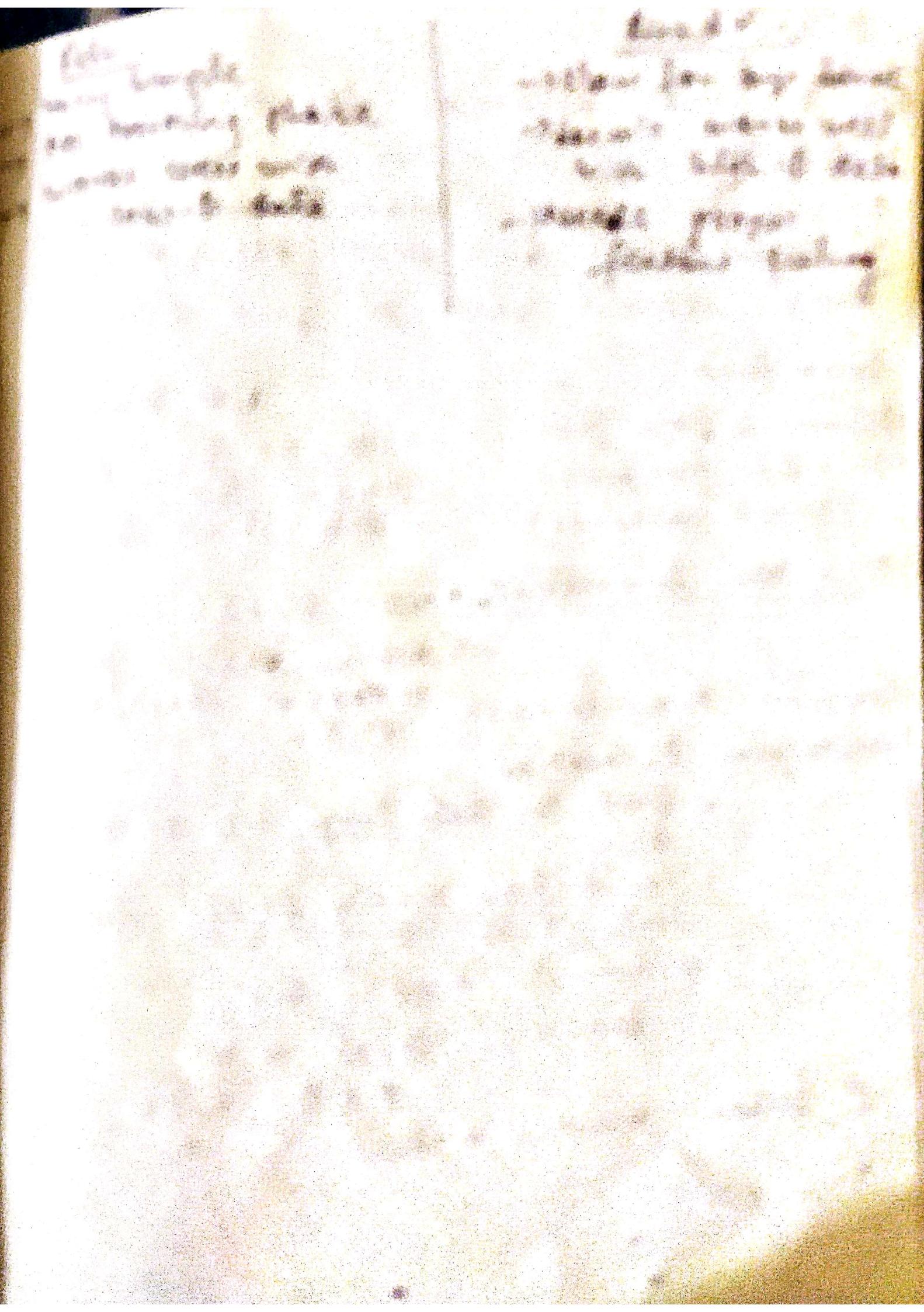
3) Minkowski Distance

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

① Linear Discriminant Function Based Classifier
 → to classify data points by drawing a str. decision boundary (hyperplane) that separates classes.

Linear Discriminant function (LDF)

to assign a class label to an input vector x using a linear function.

function form:- $g(x) = w^T x + w_0$

w : weight vector

x : input vector

w_0 : bias (threshold)

Decision rule (for 2-class problem):

If $g(x) > 0$, assign to class 1

If $g(x) \leq 0$, assign to class 2

Ex:- If it has 2 features (x_1, x_2) , the decision boundary is:-

$$w_1 x_1 + w_2 x_2 + w_0 = 0$$

This is the eqn. of a line.

Properties: → simple & fast

→ can't model non-linear reln.

→ works best when data is linearly separable.

Simple mathematical tool, used to classify data into diff. categories by drawing a flat line (2D) or a flat plane that separates the data.

Ex:- You have 2 types of fruits - Apples, Bananas

Each fruit has 2 features - weight, color.

We can plot them on a graph. If a straight line can separate apples from bananas, then you can use linear discriminant function based classifier.

② Perceptron

- Simplest type of neural network.
- One of the earliest algo. used for classification.
- tries to mimic how a human neuron works & is mainly used to classify data into 2 grp. (binary classification)
- A perceptron takes I/P, multiplies them with weights, adds them up, and then decides whether the O/P should be 1 (yes) or 0 (no).
- formula: - Output = $\begin{cases} 1 & \text{if } (w \cdot x + b) > 0 \\ 0 & \text{otherwise} \end{cases}$
where:- x = I/P features (ht., wt.)
 w = weights (imp. of each I/P)
 b = bias.
 $w \cdot x$ = weighted sum of I/P.
- Perceptron draws a st. line (2D) to separate 2 classes.
- ex:- Want to separate cat(1) from Dog(0).
Based on their height & weight.
Perceptron will try to draw a line on a graph so that
 - cats are on one side
 - Dogs are on other side.

Adv.

Simple & fast
works well for linearly separable data.

Limitations

- Only works if data can be separated by a st. line
- Doesn't handle multi-class problems well.

③ Support Vector Machines (SVM)

- powerful ML algo.
- used for classif. & regression.
- mostly used to classify data into 2 classes.
(Spam or not spam, cat or dog)
- tries to find the best boundary that separates 2 classes with largest possible gap b/w them.
- Ex:- Imagine you have 2 types of shoes
 - . Sports shoes
 - . Party shoes.You want to separate them on a graph based on 2 features: color & heel size.
SVM draws a line b/w the 2 types but not just any line.
It picks the best possible line that is far away from both classes.
That way, new shoes can be placed on the correct side of the line.

Adv.

- effective when clear margin of separation exists.
- can be used for both linear & non-linear classification.
- works well with high-D data.

Disadv.

- slow with large datasets.
- hard to choose the right kernel sometimes
- doesn't perform well if classes overlap a lot.

Real-life uses of SVM:

- Handwriting recog.
- face detection.
- Spam email filtering

• Bioinformatics -

and rounding error
point arithmetic, Preparation
of error

MOD-VIII

① Non-metric methods for pattern recognition classification

Non-metric methods

- Classifⁿ methods that don't rely on distⁿ or geometric measurements
- Not based on metrics like Euclidean dist
- Useful when data is non-numeric
- Nominal data - data that represents categorical numbers
- Ex: - colors: {R, G, B} ; weather: {Sunny, Rainy, Cloudy} ; Gender: {Male, Female} ; Blood Types: {A, B, AB, O}
- Why metric methods don't work.
You can't say "Red" is closer to "Blue" than "Green".
Thus non-metric method is needed.

② Decision Trees

- A tree-like str. where:
- Internal nodes represent decision rules on attributes.
 - Branches rep. outcomes of decisions
 - Leaves rep. class labels (final decision)

Ex:- For a weather Dataset:-

<u>Outlook</u>	<u>Temp.</u>	<u>Humidity</u>	<u>windy</u>	<u>Play Tennis</u>
Sunny	HOT	High	False	No
Overcast	Cool	Normal	True	Yes
Rainy	Mild	High	False	Yes

		[Outlook?]	
Sunny	Overcast	Rainy	
[Humidity?]	Yes	[Windy?]	
High	Normal	false	True
No	Yes	Yes	No

Adv.

- easy to visualize & interpret
- NO need for feature scaling
- can handle non-numeric & categorical data directly.

Disadv.

- . can become complex if not pruned.
- , sensitive to small changes in data
- . prone to overfitting,

① Unsupervised Learning & Clustering

- Type of ML where the comp. learns pattern from data without any labels.
- No one tells the mlc what the correct answers are.
- The mlc figure things out on its own.
 - Ex:- Imagine you have been given a bunch of photos of animals, but no one tells you which one is a dog, cat, or rabbit. Your task is to group similar-looking animals together. That's exactly what unsupervised learning does!
- Unsupervised Algorithms
 - . K-Means Clustering
 - . Principal Component Analysis (PCA)

Clustering :-

- an unsupervised learning method where we group similar data points into clusters.
- Ex:- Given customer data (age, income), clustering can grp. them into segments like:-
Low income, Young age
High income, middle age
Retired, fixed income.

Application → Market Segmentation

- Doc / Image classification
- Anomaly detection

- ② Criterion functions for clustering
- technique in unsupervised learning where we grp. similar data points together into clusters.
 - ex:- Grouping students by similar marks.
 - Grouping customers based on shopping habits.
 - Criterion fun. → like rule / formula
→ helps comp. measure how good a clustering is.
 - goals → points in the same cluster should be close together.
 - points in diff. clusters should be far apart.

- common criterion functions \rightarrow intra-cluster distance
- Within-Cluster Distance (WCD)
 - measures how close the points are inside each cluster.
make this as small as possible
smaller value - better grouping;
ex - group friends by hobbies.
 - Between-Cluster Distance (BCD) \rightarrow inter-cluster distance
 - measures how far apart the clusters are from each other.
goal - make this as large as possible.
Bigger dist. \rightarrow better separation b/w clusters.
 - Ex - you want grp. of Stu. with very diff. Study habits in separate clusters.

Sum of Squared Errors (SSE)

- Adds up the sq. dist. of all points from their cluster center.
- minimize SSE
- used in K-Means clustering

(3) Algorithms for clustering

A) K-Means clustering

1. Choose a no. of clusters k
2. Initialize k random centroids
3. Assign each data point to nearest centroid
4. Update centroids by averaging points in each cluster
5. Repeat steps 3-4 until convergence.

- unsupervised learning algo.
- used to grp. data points into ' k ' clusters based on how similar they are.
- like putting similar items into same box.
- Imagine you have a bag of colorful balls & you want to group them by color, but the balls are not labeled.
- K-means will automatically find groups based on color without you telling it.
- $K =$ no. of clusters. You want ~~to~~ choose this no. before the algo. runs.
for ex: $K = 3$ → it will make 3 groups from the data.

→ works?

1. Choose k points from data
2. Place k random points on data
3. Assign each point to nearest centroid
4. Move the centroids to the centre of their assigned points
5. Repeat 3 & 4 until centroids don't move much.

→ Real-life use → Image compression.

→ Market Segmentation

④ Hierarchical clustering

→ method of unsupervised learning

→ grp. data into tree-like str. called dendrogram.

→ doesn't need you to choose the no. of clusters at the start like k-Means.

Instead, it builds a hierarchy of clusters -
from small to big or big to small.

→ Ex - Organizing a family tree.

You first grp. close relatives

Then you merge them into bigger family groups

Eventually you get 1 big family tree.

Other methods

1) Density Based spatial Clustering of Applications with Noise.

→ Groups points are closely packed

2) Mean Shift

3) Gaussian Mixture Models (GMM).