

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
 - Clustering
 - 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.