

# MACHINE LEARNING – UNIT IV

## NOTES

### RGPV Exam-Oriented Notes

 **Easy + Detailed + One-Day Before Exam Preparation**

---

## UNIT IV SYLLABUS

1. Introduction to Clustering
  2. Hierarchical Clustering
  3. AGNES
  4. DIANA
  5. K-Means Clustering
  6. K-Mode Clustering
  7. Self Organizing Map (SOM)
  8. Expectation Maximization (EM)
  9. Gaussian Mixture Models (GMM)
  10. Principal Component Analysis (PCA)
  11. Locally Linear Embedding (LLE)
  12. Factor Analysis
- 

### **1. INTRODUCTION TO CLUSTERING**

#### **Definition**

Clustering is an unsupervised learning technique used to group similar data points into clusters.

---

## Easy Explanation

Clustering groups similar objects together.

Example:

Students with similar marks → One Group

Customers with similar shopping habits → One Group

---

## Important Point

In clustering:

No labeled output is given

Machine finds hidden patterns automatically.

---

## Types of Clustering

Type	Meaning
Hierarchical Clustering	Creates cluster hierarchy
Partitional Clustering	Divides data into groups
Density-Based	Based on data density

---

## Applications

- Customer segmentation
  - Image grouping
  - Market analysis
  - Recommendation systems
- 

## **Advantages**

- ✓ Finds hidden patterns
  - ✓ Works without labeled data
  - ✓ Useful for large datasets
- 

## **Disadvantages**

- ✗ Difficult to choose clusters
  - ✗ Sensitive to noisy data
- 

## **Conclusion**

Clustering helps organize similar data into meaningful groups.

---

# **2. HIERARCHICAL CLUSTERING**

## **Definition**

Hierarchical clustering creates clusters in the form of a hierarchy or tree structure.

---

## **Easy Explanation**

It groups data step by step until all points become part of a hierarchy.

---

## Types

Type	Meaning
AGNES	Bottom-up approach
DIANA	Top-down approach

---

## Diagram

Small Clusters



Bigger Clusters



Final Cluster

---

## Advantages

- ✓ Easy visualization
  - ✓ No need to specify clusters initially
- 

## Disadvantages

- ✗ Slow for large data
  - ✗ Computationally expensive
-

## Applications

- Biological classification
  - Document clustering
  - Social network analysis
- 

## Conclusion

Hierarchical clustering creates nested clusters in tree-like structure.

---

## 3. AGNES

### Full Form

AGNES = Agglomerative Nesting

---

### Definition

AGNES is a bottom-up hierarchical clustering algorithm.

---

### Easy Explanation

Initially every data point is treated as separate cluster.

Nearest clusters are merged repeatedly.

---

## Working

Single Points



Merge Closest Clusters



Larger Clusters

---

## Advantages

- ✓ Simple approach
  - ✓ Good for small datasets
- 

## Disadvantages

- ✗ Slow for large datasets
- 

## Conclusion

AGNES builds clusters by merging nearest clusters step by step.

---

## 4. DIANA

## Full Form

DIANA = Divisive Analysis

---

## Definition

DIANA is a top-down hierarchical clustering algorithm.

---

## Easy Explanation

Initially all data points belong to one cluster.

Then clusters are divided repeatedly.

---

## Working

One Large Cluster



Split into Smaller Clusters



Final Clusters

---

## Advantages

- Better cluster division
  - Useful for structured data
- 

## Disadvantages

- High computational cost

---

## Conclusion

DIANA creates clusters by repeatedly splitting larger clusters.

---

## 5. K-MEANS CLUSTERING

### Definition

K-Means is a partitional clustering algorithm that divides data into K clusters.

---

### Easy Explanation

K-Means groups similar data points around central points called centroids.

---

### Working Steps

1. Choose K clusters.
  2. Select random centroids.
  3. Assign points to nearest centroid.
  4. Recalculate centroid.
  5. Repeat until clusters stabilize.
- 

### Diagram

Data Points → Grouped Around Centers

---

## Advantages

- ✓ Simple and fast
- ✓ Works well for large data
- ✓ Easy implementation

---

## Disadvantages

- ✗ Need to choose K value
- ✗ Sensitive to outliers

---

## Applications

- Customer segmentation
- Image compression
- Pattern recognition

---

## Conclusion

K-Means clusters similar data points around centroids efficiently.

---

## 6. K-MODE CLUSTERING

### Definition

K-Mode clustering is an extension of K-Means used for categorical data.

---

## Easy Explanation

K-Means works with numerical data.

K-Mode works with:

Yes/No

Male/Female

Color Names

---

## Features

K-Means	K-Mode
Numerical data	Categorical data
Uses mean	Uses mode

---

## Advantages

- Handles categorical data
  - Simple implementation
- 

## Disadvantages

- Less accurate for mixed data
-

## **Conclusion**

K-Mode is useful for clustering categorical datasets.

---

## **7. SELF ORGANIZING MAP (SOM)**

### **Definition**

SOM is a neural-network based clustering technique that converts high-dimensional data into low-dimensional maps.

---

### **Easy Explanation**

SOM organizes similar data together automatically.

---

### **Working**

Input Data



Competitive Learning



Organized Map

---

### **Advantages**

- ✓ Good visualization
  - ✓ Reduces dimensionality
- 

## **Disadvantages**

- ✗ Slow training
  - ✗ Complex implementation
- 

## **Applications**

- Pattern recognition
  - Data visualization
  - Image analysis
- 

## **Conclusion**

SOM is useful for clustering and visualizing high-dimensional data.

---

# **8. EXPECTATION MAXIMIZATION (EM)**

## **Definition**

EM is an iterative algorithm used to estimate unknown parameters in statistical models.

---

## **Easy Explanation**

EM repeatedly improves cluster assignments and parameters until best solution is found.

---

## Two Steps of EM

Step	Meaning
Expectation Step	Estimate probabilities
Maximization Step	Update parameters

---

## Working

Initial Parameters  
↓  
Expectation Step  
↓  
Maximization Step  
↓  
Repeat Until Stable

---

## Advantages

- ✓ Works for incomplete data
  - ✓ Produces accurate probability estimates
- 

## Disadvantages

- ✗ Slow convergence
  - ✗ Sensitive to initial values
-

## Applications

- Clustering
  - Image segmentation
  - Pattern recognition
- 

## Conclusion

EM improves parameter estimation iteratively for clustering and probability models.

---

# 9. GAUSSIAN MIXTURE MODEL (GMM)

## Definition

GMM is a probabilistic clustering model that assumes data is generated from multiple Gaussian distributions.

---

## Easy Explanation

Instead of fixed clusters, GMM gives probability of belonging to each cluster.

---

## Example

Point A:

70% Cluster 1

30% Cluster 2

---

## Advantages

- ✓ Flexible clustering
  - ✓ Handles overlapping clusters
- 

## Disadvantages

- ✗ Computationally expensive
  - ✗ Sensitive to initialization
- 

## Applications

- Speech recognition
  - Image processing
  - Pattern analysis
- 

## Conclusion

GMM provides soft clustering using probability distributions.

---

# 10. PRINCIPAL COMPONENT ANALYSIS (PCA)

## Definition

PCA is a dimensionality reduction technique used to reduce large features into smaller important features.

---

## Easy Explanation

PCA removes unnecessary data while preserving important information.

---

## Working

High-Dimensional Data



Important Components Selected



Reduced Data

---

## Advantages

- ✓ Reduces complexity
  - ✓ Faster computation
  - ✓ Removes redundancy
- 

## Disadvantages

✘ Information loss possible

✘ Difficult interpretation

---

## Applications

- Image compression
  - Data visualization
  - Feature extraction
- 

## Conclusion

PCA reduces dimensions while preserving important information.

---

# 11. LOCALLY LINEAR EMBEDDING (LLE)

## Definition

LLE is a non-linear dimensionality reduction technique.

---

## Easy Explanation

LLE preserves local relationships between nearby data points while reducing dimensions.

---

## Advantages

- ✓ Captures non-linear structure
  - ✓ Better visualization
- 

## **Disadvantages**

- ✗ Computationally expensive
- 

## **Applications**

- Image processing
  - Data visualization
- 

## **Conclusion**

LLE reduces dimensions while preserving neighborhood relationships.

---

# **12. FACTOR ANALYSIS**

## **Definition**

Factor Analysis identifies hidden factors responsible for observed data.

---

## **Easy Explanation**

Many variables may depend on a few hidden causes called factors.

---

## **Example**

Student Performance depends on:

- Intelligence
- Study Habit
- Motivation

These hidden causes are factors.

---

## **Advantages**

- ✓ Reduces variables
  - ✓ Identifies hidden relationships
- 

## **Disadvantages**

- ✗ Difficult interpretation
- 

## **Applications**

- Psychology
  - Market research
  - Data analysis
- 

## **Conclusion**

Factor Analysis finds hidden factors affecting observed data.

---



## **MOST IMPORTANT 7-MARK**

### **QUESTIONS**

1. Explain Clustering with types.
  2. Explain Hierarchical Clustering.
  3. Explain AGNES and DIANA.
  4. Explain K-Means Clustering.
  5. Explain PCA.
  6. Explain EM Algorithm.
  7. Explain Gaussian Mixture Model.
  8. Explain SOM.
  9. Explain Factor Analysis.
- 



## **MOST IMPORTANT 14-MARK**

### **QUESTIONS**

★ Very Important

1. Explain K-Means Clustering with steps.
  2. Explain Hierarchical Clustering with AGNES and DIANA.
  3. Explain PCA with applications.
  4. Explain EM Algorithm and GMM.
  5. Explain SOM with applications.
  6. Explain dimensionality reduction techniques.
  7. Explain Factor Analysis in detail.
-

# PYQ-BASED EXPECTED QUESTIONS

## High Probability

- Explain K-Means Clustering.
  - Explain PCA.
  - Explain Hierarchical Clustering.
  - Differentiate AGNES and DIANA.
  - Explain EM Algorithm.
- 

## Medium Probability

- Explain GMM.
  - Explain SOM.
  - Explain LLE.
  - Explain Factor Analysis.
- 

## ONE-NIGHT REVISION

- ✓ Clustering groups similar data
  - ✓ K-Means uses centroids
  - ✓ AGNES = Bottom-up clustering
  - ✓ DIANA = Top-down clustering
  - ✓ PCA reduces dimensions
  - ✓ SOM uses neural networks
  - ✓ EM estimates probabilities iteratively
  - ✓ GMM uses Gaussian distributions
  - ✓ LLE preserves local relationships
  - ✓ Factor Analysis finds hidden factors
-



# SMART STUDY PLAN



## FIRST STUDY

1. K-Means Clustering
  2. PCA
  3. Hierarchical Clustering
  4. AGNES vs DIANA
- 



## SECOND PRIORITY

5. EM Algorithm
  6. GMM
  7. SOM
- 



## LAST REVISION

8. LLE
  9. Factor Analysis
  10. K-Mode Clustering
- 



## FINAL EXAM TIP

For 14-mark answers always write:

1. Definition
2. Working Steps
3. Advantages
4. Disadvantages
5. Applications
6. Conclusion

This presentation style helps score higher marks in RGPV exams.