Experiment No: 5                                                06/02/2025

title: Expectation Maximization Algorithm for clustering.

Aim: To implement the Expection - Maximization (EM) algorithm for estimating the parameters of a guassian mixture model (GMM) and apply it to cluster of data point.

Theory:

Expectation maximization (EM) Algorithm: The expectation maximization (EM) Algorithm is an iterative optimization technique used to estimate parameters in probobilistic models when some data is missing or hidden. It is widely used in unsupervised learning for clustering, particularly in Gaussion Mixture Model.

A gaussian mixture model assume that data points are generated from a mixture of multiple gaussian distributed, each with its own mean and variance. The EM algorithm find the optimal parameter for each Gaussian components.
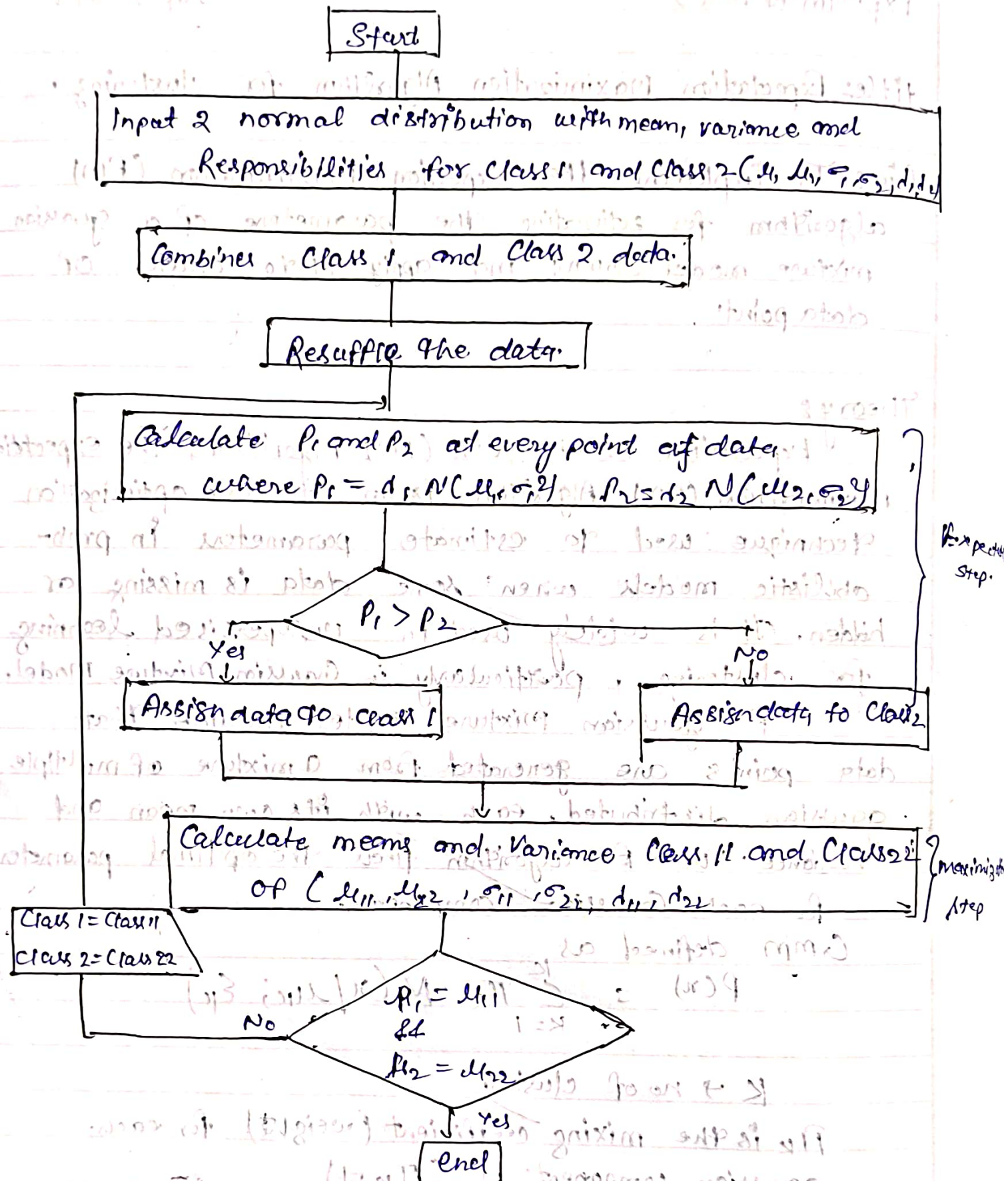
GMM defined as

$$P(x) = \sum_{K=1}^{K} \pi_K \cdot N(x | \mu_K; \Sigma_K)$$

$K \rightarrow$ no of clusters.

$\pi_K$ is the mixing coefficient (weight) for each gaussian component $(\sum \pi_K = 1)$

**Flowchart:**

```
                    ┌─────────┐
                    │  Start  │
                    └────┬────┘
                         │
┌────────────────────────────────────────────────────────────┐
│ Input 2 normal distribution with mean, variance and         │
│    Responsibilities for Class 1 and Class 2 ($\mu_1, \mu_2, \sigma_1, \sigma_2, d_1, d_2$) │
└────────────────────────┬───────────────────────────────────┘
                         │
        ┌────────────────────────────────────┐
        │ Combines Class 1 and Class 2 data. │
        └────────────────┬───────────────────┘
                         │
              ┌─────────────────────────┐
              │  Resuffle the data.     │
              └────────────┬────────────┘
                         │
┌────────────────────────────────────────────────────────┐
│ Calculate $P_1$ and $P_2$ at every point of data        │    ⎫
│   where $P_1 = d_1 \cdot N(\mu_1, \sigma_1^2)$   $P_2 = d_2 \cdot N(\mu_2, \sigma_2^2)$ │    ⎬ Expected
└────────────────────────┬───────────────────────────────┘    ⎭ Step.
                         │
                    ◇ $P_1 > P_2$ ◇──────────── No ──────────┐
                    Yes │                                     │
                        ▼                                     ▼
        ┌──────────────────────────┐         ┌──────────────────────────┐
        │ Assign data to class 1   │         │ Assign data to Class 2   │
        └──────────────┬───────────┘         └──────────────┬───────────┘
                       └────────────┬───────────────────────┘
                                    │
┌────────────────────────────────────────────────────────────┐
│ Calculate means and Variance Class 1 and Class 2            │   ⎫
│   of ($\mu_{11}, \mu_{22}, \sigma_{11}, \sigma_{22}, d_{11}, d_{22}$ )  │   ⎬ maximize
└────────────────────────┬───────────────────────────────────┘   ⎭ step
                         │
┌──────────────────┐     │
│ Class 1 = Class 11│    ◇ $P_1 = \mu_{11}$ ◇
│ Class 2 = Class 22│────── No ──  &&
└──────────────────┘     │     $P_2 = \mu_{22}$ ◇
                         │        │ Yes
                         │        ▼
                         │   ┌─────────┐
                         │   │   end   │
                         │   └─────────┘
```

Gaussian probability density function (Pdf) :

$$N(x|\mu_k, \Sigma_k) = \underset{P(x_i|b)}{\frac{1}{\sqrt{2\pi\sigma_k^2}}} \exp\left(\frac{-(x_i-\mu_{id})^2}{2\sigma_k^2}\right)$$
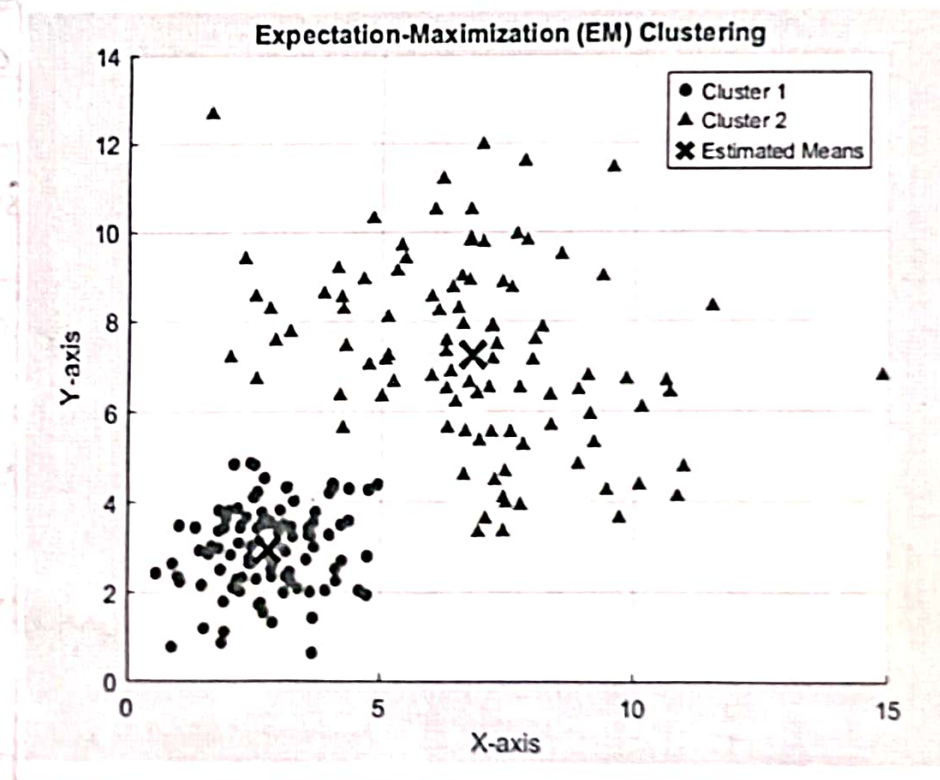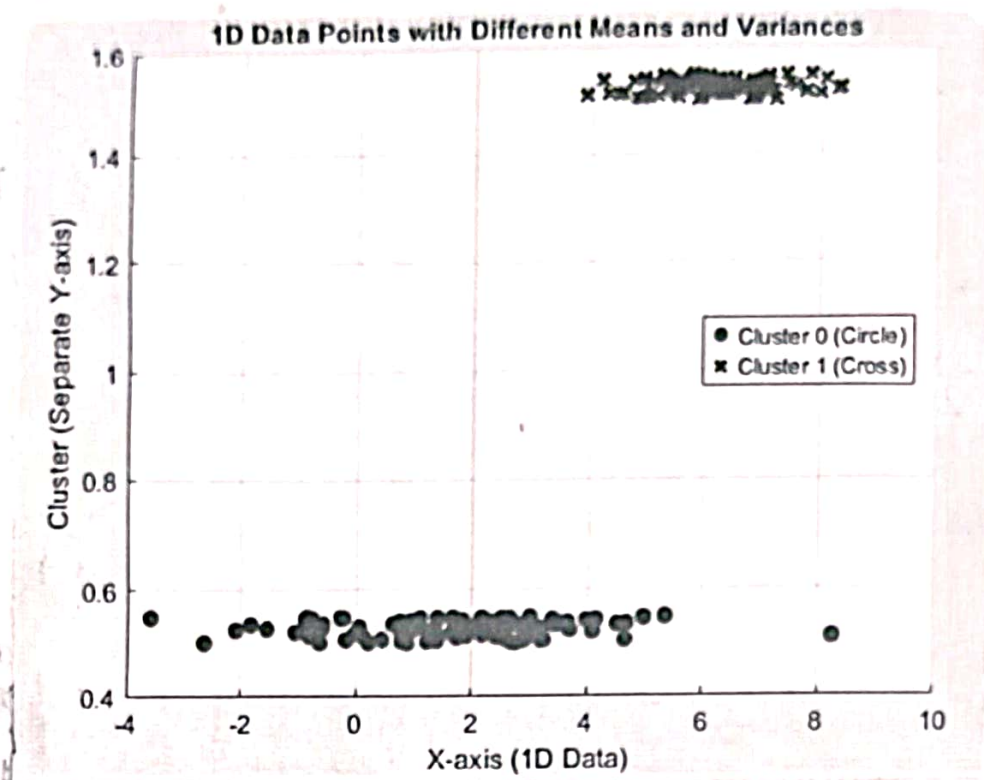
$\sigma_k \rightarrow$ covariance (, matrix)

$\mu_k \rightarrow$ means (vector)

## Algorithm:

$\Rightarrow$ Initilize the Input the normal distribution with means, variance and weights coefficient for 2 clusters.

$\Rightarrow$ Combines the both clusters data.

$\Rightarrow$ Suffle the data.

$\Rightarrow$ Calculate the probabilites for the both cluster. where the $P_1 = \lambda_1 \cdot N(\mu_1, \sigma_1^2)$ and

$P_2 = \lambda_2 N(\mu_2, \sigma_2^2)$

$\Rightarrow$ Assign the data with the highest Probabity Cluster.

$=)$ Calculate the means and variance the new assign clusters.

$\dashv$ if the new means and new variance is same the previous means and variance Respectivily assign that particular cluster. Other wise iterate the step-4.

$\dashv$ end

1D Data Points with Different Means and Variances



Expectation-Maximization (EM) Clustering

Result: Taking random datapoint in 1D & 2D plane. (100 Samples)

for 1D→ initilize mean. [2,6]

variance randomly choose b/w the 0-2

for 2D→ Initilize means [1,5; 7,8];

co-variance matrix 1 = $\begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix}$, Covariance matrix = $\begin{bmatrix} 2 & 0 \\ 0 & 2 \end{bmatrix}$

Discussion!

→ The algorithm starts with arbitrary initial means and variance.

→ A responsibility matrix is created to store the probability that each datapoint belongs to a given cluster.

→ Compute the probability of each data point belongig to each cluster using Gaussian probability density function (PdF)

→ Responsibilities (Y) are computed based on these probabilities.

→ Each datapoints is assigned to the cluster with the highest responsibilities value

→ The Expectation - Maximization Algorithm should be able to estimate these cluster assignments without prior knowledge of them based purly on the distribution of data.

## Conclusion:

The Expectation Maximization Algorithm clusters the given 2D and 1D data into two groups based on their underlying Gaussian distribution. It iteratively refines the cluster parameters. Updating means and covariance matrics until Convergence. The soft clustering approach assign probabilities. rather than strict label making Em effective for overlapping Clusters. Although Em is sensitive to initilization. It perform well with properly chosen starting values. Enhancement like better initilization and regularization can further improve robustness, Overall, Em is a powerful clustering method widely used in the real world application.

13/02/25