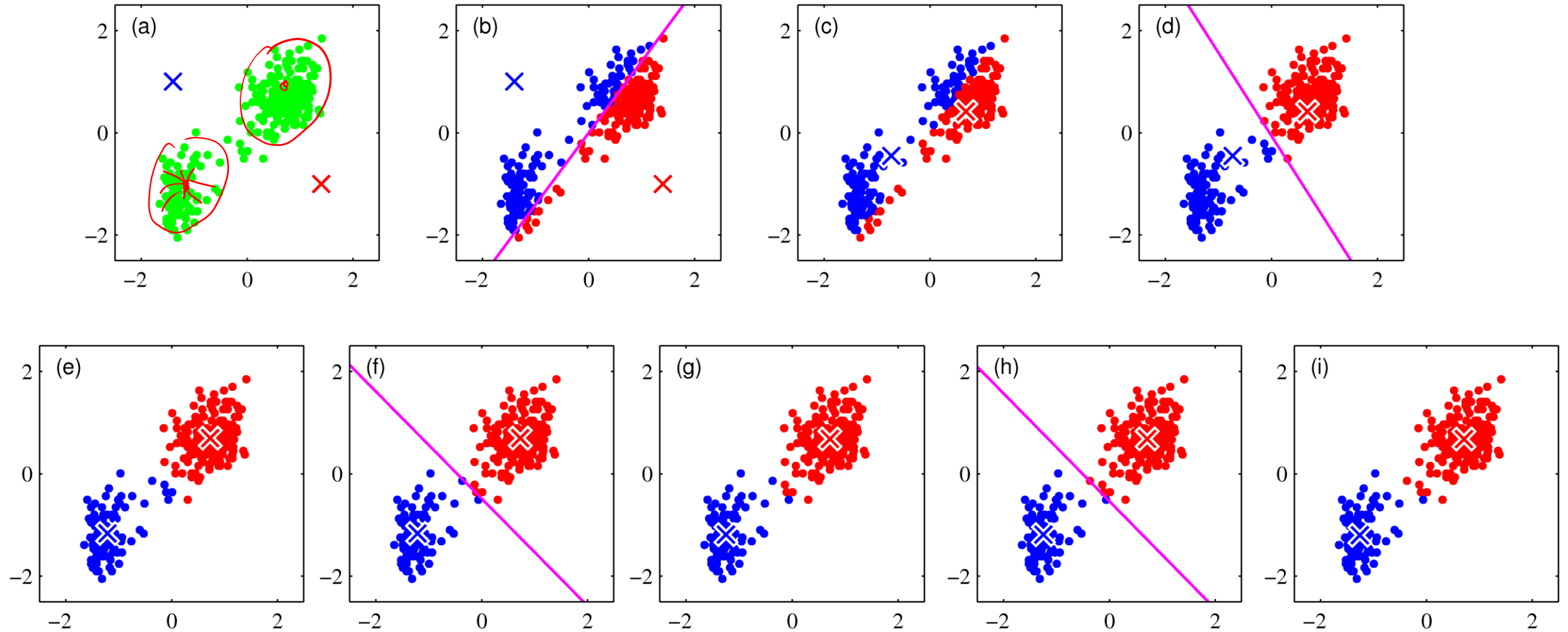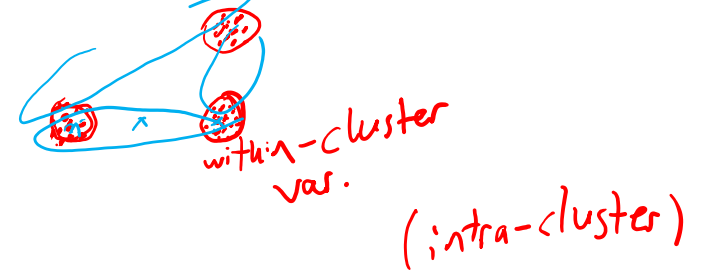# Data Analytics
# EEE 4774 & 6777

Module 3

Clustering

Spring 2022

# Clustering: K-means

# K-means

- Unsupervised method for identifying groups: Clustering

- Data $\{x_1, \ldots, x_N\}$ where $x_n \in \mathbb{R}^D$ — #dim.

  #instances   → #clusters

- $\min E(c_n, m_k) = \sum_{n=1}^{N} \sum_{k=1}^{K} c_{nk} \|x_n - m_k\|^2$ where $c_n = [c_{n1} \ldots c_{nK}]$ and $c_{nk} \in \{0,1\}$

  cluster mean

- Iteratively minimize $E$ over $c_n$ and $m_k$

  cluster assignment var.

Objective: Minimize the within-cluster variances

within-cluster var.

(intra-cluster)

Initialize $m_k$
for i=1:max_iter

Step 1 :   Minimize $E$ with respect to $c_n$ keeping $m_k$ fixed → Update $c_n$

Step 2 :   Minimize $E$ with respect to $m_k$ keeping $c_n$ fixed → Update $m_k$

    if $\dfrac{\left\| c_n^{(i)} - c_n^{(i-1)} \right\|}{\left\| c_n^{(i-1)} \right\|} < \varepsilon$ and $\dfrac{\left\| m_k^{(i)} - m_k^{(i-1)} \right\|}{\left\| m_k^{(i-1)} \right\|} < \varepsilon$

         norm

        *break*

    end
end

# K-means



Step 1:

Assign each data point to the nearest cluster

$$c_{nk} = \begin{cases} 1 & if \quad k = \arg\min_j \|\boldsymbol{x}_n - \boldsymbol{m}_j\|^2 \\ 0 & otherwise \end{cases}$$

n th data instance — cluster center — Euclidean distance

Step 2

$$\boldsymbol{m}_k = \frac{\sum_n c_{nk} \boldsymbol{x}_n}{\sum_n c_{nk}} = mean\ of\ points\ assigned\ to\ cluster\ k$$

# data points in cluster k

- Since $E$ decreases at each iteration, convergence is guaranteed

- However, it may converge to a local minimum

- K-medoids: generalization of K-means to a general distance measure

$$E(\boldsymbol{c}_n, \boldsymbol{m}_k) = \sum_{n=1}^{N} \sum_{k=1}^{K} c_{nk} V(\boldsymbol{x}_n, \boldsymbol{m}_k)$$

non-Euclidean distance

K-means

(a) randomly chosen cluster means

Iteration 1 — Step 1 (b), Step 2 (c)

Iter. 2 — Step 1 (d)

It.2 - Step 2 (e)

I3 - S1 (f)

I3 - S2 (g)

I4 - S1 (h)

I4 - S2 (i)

# Gaussian Mixture Model $(GMM)$



cluster membership var. ~ categorical dist.

Parameters

$C_n$

$\pi^{(K)}$

$\mu_{k,K}$

$\Sigma_{k,K}$

$x_n$ observed data inst.

K Gaussian models

$N$

$$\boldsymbol{C}_n = [C_{nk}]_{k=1,\dots,K} = [0 \cdots 1 \cdots 0] \in [0,1]^K$$

Cluster assign/membership var.    K clusters

$$p(C_{nk} = 1) = \pi_k, \quad \pi_k \in [0,1], \quad \sum_{k=1}^{K} \pi_k = 1$$

Likelihood of $x_n$ under GMM

latent/hidden var.

categorical distr.

cluster mean    covariance

$$p(\boldsymbol{x}_n) = \sum_{\boldsymbol{C}_n \in [0,1]^K} p(\boldsymbol{x}_n, \boldsymbol{C}_n) = \sum_{k=1}^{K} \pi_k \mathcal{N}(\boldsymbol{x}_n | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)$$

joint prob dist.    prior    likelihood of $x_n$ under $k$th Gaus. Model

log-likelihood    # instances

$$\log p(\boldsymbol{X}) = \sum_{n=1}^{N} \log \sum_{k=1}^{K} \pi_k \mathcal{N}(\boldsymbol{x}_n | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)$$

$\log P(X_n)$

$$\log p(X) = \log \prod_{n=1}^{N} P(x_n) = \sum \log p(x_n)$$

$$= \sum_{n=1}^{N} \log \sum_{k=1}^{K} \pi_k \frac{\exp\{-(\boldsymbol{x}_n - \boldsymbol{\mu}_k)^T \boldsymbol{\Sigma}_k^{-1}(\boldsymbol{x}_n - \boldsymbol{\mu}_k)/2\}}{(2\pi)^{d/2} |\boldsymbol{\Sigma}_k|^{1/2}}$$

# ML for GMM

$$\frac{\partial}{\partial \mu_k} \log f(\mu_k) = \frac{\frac{\partial}{\partial \mu_k} f(\mu_k)}{\mu_k}$$

$$\frac{\partial}{\partial \mu_k} e^{f(\mu_k)} = e^{f(\mu_k)} \frac{\partial}{\partial \mu_k} f(\mu_k)$$

$$\max_{\boldsymbol{\mu}_k} \log p(\boldsymbol{X}) \implies \frac{\partial}{\partial \boldsymbol{\mu}_k} \log p(\boldsymbol{X}) = \sum_{n=1}^{N} \frac{\pi_k \mathcal{N}(\boldsymbol{x}_n | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)}{\sum_{j=1}^{K} \pi_j \mathcal{N}(\boldsymbol{x}_n | \boldsymbol{\mu}_j, \boldsymbol{\Sigma}_j)} \boldsymbol{\Sigma}_k^{-1}(\boldsymbol{x}_n - \boldsymbol{\mu}_k) = 0$$

*prior prob. for cluster k*

*derivative of exponent in Gauss. pdf*

*posterior prob. for cluster k*

$$p(C_{nk} = 1 | \boldsymbol{x}_n) = \frac{p(C_{nk} = 1)\, p(\boldsymbol{x}_n | C_{nk} = 1)}{\sum_{j=1}^{K} p(C_{nj} = 1)\, p(\boldsymbol{x}_n | C_{nj} = 1)} = \frac{\pi_k \mathcal{N}(\boldsymbol{x}_n | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)}{\sum_{j=1}^{K} \pi_j \mathcal{N}(\boldsymbol{x}_n | \boldsymbol{\mu}_j, \boldsymbol{\Sigma}_j)} = \gamma(C_{nk})$$

*$\pi_k$*   *$\mathcal{N}(\boldsymbol{x}_n | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)$*

*joint prob. $p(\boldsymbol{x}_n, C_n)$*

*posterior*

**coupled equations no closed-form solution!**

For $\mu_k, \Sigma_k, \pi_k$
Need to find $\gamma(C_{nk})$
which depends
on $\mu_k, \Sigma_k, \pi_k$ !

$$\boldsymbol{\mu}_k = \frac{1}{\sum_{n=1}^{N} \gamma(C_{nk})} \sum_{n=1}^{N} \gamma(C_{nk})\, \boldsymbol{x}_n = \frac{1}{N_k} \sum_{n=1}^{N} \gamma(C_{nk})\, \boldsymbol{x}_n$$

*weight with posterior*

*weighted average (sample mean)*

Similarly,

$$\boldsymbol{\Sigma}_k = \frac{1}{N_k} \sum_{n=1}^{N} \gamma(C_{nk})\, (\boldsymbol{x}_n - \boldsymbol{\mu}_k)(\boldsymbol{x}_n - \boldsymbol{\mu}_k)^T, \qquad N_k = \sum_{n=1}^{N} \gamma(C_{nk}),$$

*weighted sample cov. matrix*

as we have multiple gaussian here and we don't know to which gauss. our data point belongs to

*effective number of points in cluster k*

and

$$\pi_k = \frac{N_k}{N}$$

as we are doing soft assignment

# Iterative Solution: EM for GMM

- Expectation-Maximization for iteratively computing ML in GMM

1. Initialize $\boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k, \pi_k$ and compute the initial value of $\log p(\boldsymbol{X})$

2. **E step:** Compute the posteriors using the current parameter values

$$\gamma(C_{nk}) = \frac{\pi_k \mathcal{N}(\boldsymbol{x}_n | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)}{\sum_{j=1}^{K} \pi_j \mathcal{N}(\boldsymbol{x}_n | \boldsymbol{\mu}_j, \boldsymbol{\Sigma}_j)}$$
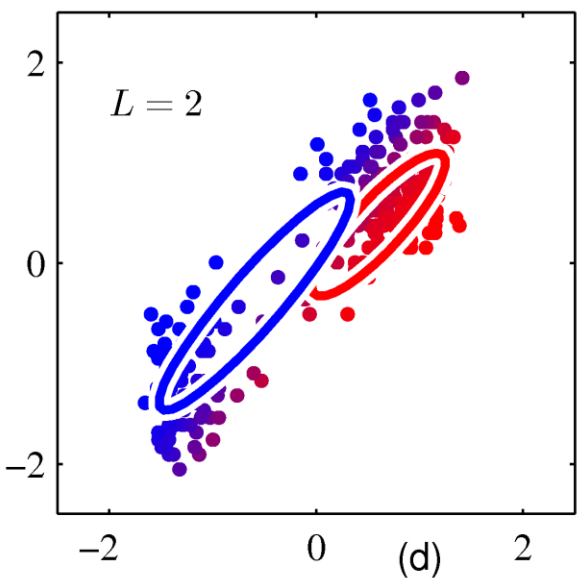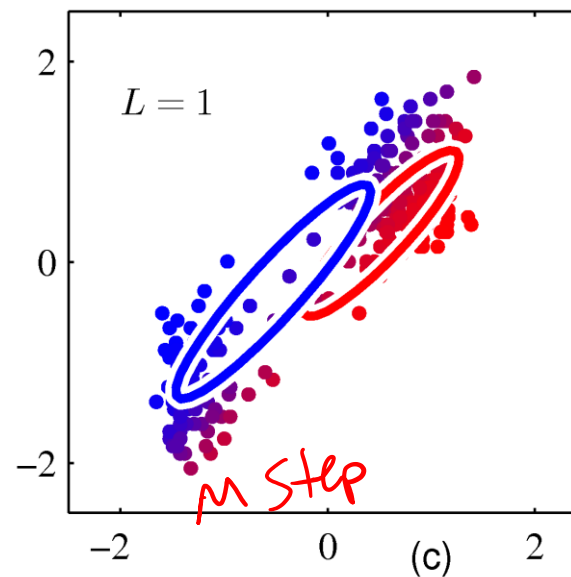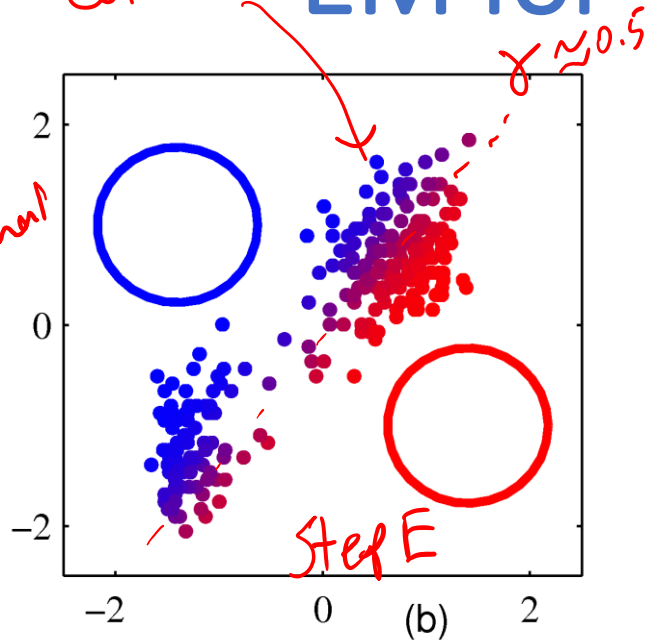
*Iteration*

3. **M step:** Re-estimate the parameters using the current posteriors

$$\boldsymbol{\mu}_k^{new} = \frac{1}{N_k} \sum_{n=1}^{N} \gamma(C_{nk}) \, \boldsymbol{x}_n \, , \quad \boldsymbol{\Sigma}_k^{new} = \frac{1}{N_k} \sum_{n=1}^{N} \gamma(C_{nk}) \, (\boldsymbol{x}_n - \boldsymbol{\mu}_k^{new})(\boldsymbol{x}_n - \boldsymbol{\mu}_k^{new})^T \, , \quad \pi_k^{new} = \frac{N_k}{N}, \quad \text{where} \quad N_k = \sum_{n=1}^{N} \gamma(C_{nk})$$

4. Compute the log-likelihood and check for convergence of either the parameters or the log-likelihood.

   If no convergence, return to step 2.

# EM for GMM



Handwritten annotations on figure (a): $\Sigma_1$: diagonal, $m_1$ X, $\Sigma_2$: diagonal, $m_2$ X

Handwritten annotations on figure (b): posterior prob. $\gamma$ color coded, $\gamma \approx 0.5$, Step E

Figure (c): $L = 1$, M Step

Figure (d): $L = 2$

Figure (e): $L = 5$

Figure (f): $L = 20$

- Many more iterations than K-means, and each iteration much more expensive,

- But provides *probabilistic modeling* with *soft assignments* and *covariance*

- Run K-means to initialize EM for GMM

- Converges to a local maximum

# Expectation-Maximization (EM) Algorithm

- Objective: find ML for models with latent variables $C$ (e.g., missing values in the dataset), observed data $X$, and parameters $\theta$

$$\log p(X|\theta) = \log \sum_C p(X, C|\theta)$$

*used for solving inter-locked equations of latent variable models*

- Assume maximization of the complete-data log-likelihood $\log p(X, C|\theta)$ is easy

1. Initialize $\theta^{old}$

2. E step: Evaluate $p(C|X, \theta^{old})$ and

   *← posterior of latent var. $C$ given data*

   *complete data*

$$Q(\theta, \theta^{old}) = E_{p(C|X, \theta^{old})}[\log p(X, C|\theta)] = \sum_C p(C|X, \theta^{old}) \log p(X, C|\theta)$$

   *posterior*    *complete data log-likelihood*

3. M step: $\theta^{new} = \arg\max_\theta Q(\theta, \theta^{old})$     {maximize $Q(\theta, \theta^{old}) + \log p(\theta)$ for MAP}

   *prior for param. $\theta$*

4. If no convergence, then $\theta^{old} \leftarrow \theta^{new}$ and return to step 2

# GMM by EM vs. K-means

- EM soft assigns data points *softly* to a cluster using posterior $p(C_{nk} = 1|x_n)$,

whereas K-means performs *hard* assignment

- Consider a GMM with covariance $\epsilon I$ for all clusters, where $\epsilon$ is a fixed constant, not a parameter to be re-estimated

$$p(C_{nk} = 1|x_n) = \frac{\pi_k \exp\{-\|x_n - \mu_k\|^2/2\epsilon\}}{\sum_{j=1}^{K} \pi_j \exp\{-\|x_n - \mu_j\|^2/2\epsilon\}}$$

- As $\epsilon \to 0$, in the denominator the smallest $\|x_n - \mu_j\|^2$ will go to 0 most slowly,

hence posterior for that cluster will go to 1 and the others will go to 0 $\implies$ ***Hard assignment to the closest cluster***

- Update for the mean $\mu_k$ also reduces to that of K-means

- K-means does not estimate the covariances of the clusters

# Evaluation of Clustering Results

- Several similarity measures for clusters can be used to evaluate the performance of clustering algorithms

- Can be used to determine the optimum number of clusters

- **Internal Evaluation:** based on the clustered data itself

    - typically assigns good score if high similarity within clusters and low similarity between clusters

    - e.g., Silhouette value (works well with K-means), Dunn index, Davies-Bouldin index

- **External Evaluation:** based on data that was not used for clustering, e.g., ground truth

    - measures how close clustering is to the benchmark classes

    - e.g., Rand index, F-measure, Mutual information, Confusion matrix
    *adjusted Rand index*