

# Model Based Clustering

## Fundamentals of Probabilistic Data Mining

---

Fei Zheng

October, 2019

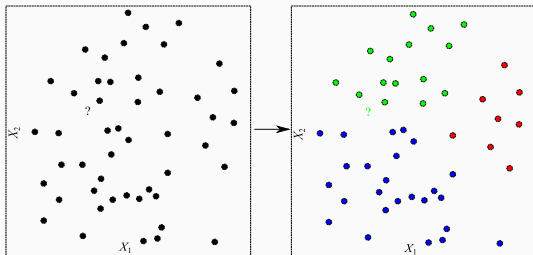


# Clustering

**Data:** points  $\{\mathbf{x}_n\}_{n=1}^N$  in  $\mathbb{R}^d$

**Aim:** find  $K$  clusters ( $K$  fixed here)

- Distance-based approaches: close points tend to be in the same cluster. No explicit assumption required.



Given a data set  $\{\mathbf{x}_n\}_{n=1}^N$ . Let  $z_n$  denote the cluster label of  $\mathbf{x}_n$ .

## K-means (MacQueen 1967)

**Objective:** The sum of the squared distance of each  $\mathbf{x}_n$  to its closest cluster centroid  $\mu_k$  is a minimum.

**Repeat until converge:**

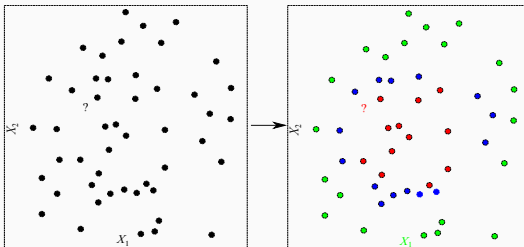
- Assign  $\mathbf{x}_n$  to the cluster indexed by the closest  $\mu_k$ .
- Update  $\mu_k$  using the newly grouped data.

## Clustering vs Classification

	Clustering	Classification
Applied case	suggest groups based on patterns in data	classify new sample into known classes
Prior knowledge	No prior knowledge	A training set
Data needs	Unlabeled samples	Labeled samples from known classes

# Clustering

- Model-based approaches: If  $z_i = z_j = k$ , then  $\mathbf{x}_i$  and  $\mathbf{x}_j$  should have the same (conditional) distribution  $p_k$ .



## Definition (McLachlan & Peel, 2000)

Let  $\{p_\theta\}_{\theta \in \Theta}$  be a parametric family of distributions.  $\mathbf{x} \rightarrow p(\mathbf{x})$  is a mixture of distributions iff there exists  $K, \pi_1, \dots, \pi_K, \theta_1, \dots, \theta_K$ , such that

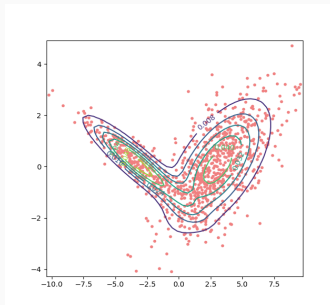
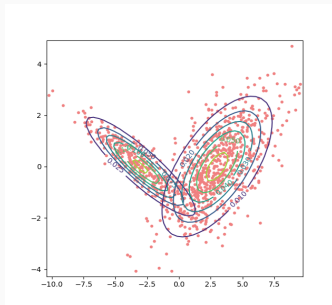
$$p = \sum_{k=1}^K \pi_k p_{\theta_k}, \quad 0 \leq \pi_k \leq 1, \quad \sum_{k=1}^K \pi_k = 1$$

- Mixture model is convex combination of distributions .
- Defines new parametric families of distributions.  
Parameters:  $\Theta = (\pi_1, \dots, \pi_K, \theta_1, \dots, \theta_K)$  (for given  $K$ ).

# Mixture Models

## Example: 2-component Gaussian Mixture Model (GMM)

- $K = 2$ ;  $\pi_1 = 0.3, \pi_2 = 0.7$ ;  $p_{\theta_1} = \mathcal{N}(\mu_1, \Sigma_1), p_{\theta_2} = \mathcal{N}(\mu_2, \Sigma_2)$ .
- $p(\mathbf{x}) = \pi_1 p_{\theta_1}(\mathbf{x}) + \pi_2 p_{\theta_2}(\mathbf{x})$ .



# Mixture Model and Clustering

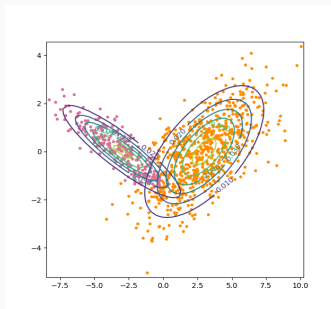
Equivalence of mixture representation and existence of some discrete hidden random variable  $Z$  (label of cluster):

$$\forall k \in \{1, \dots, K\}, \pi_k = p(Z = k)$$

$$\forall \mathbf{x}, p_{\theta_k}(\mathbf{x}) = p(\mathbf{x}|Z = k)$$

$$p(\mathbf{x}) = \sum_{k=1}^K p(Z = k)p(\mathbf{x}|Z = k)$$

- **Clustering:** to infer  $Z$  interpreted as the cluster label (heterogeneous sources) given  $\mathbf{x}$ .





- $X$ : weight of some rodent.
- Proportion of females  $1/3$ , males  $2/3$ .
- Females generally lighter than males, the heaviest females potentially heavier than the lightest males.
- The weights ( $X$ ) are (conditionally) Gaussian distributed for each gender, depend on mean weight and variance.
- Unknown genders in population – both genders are mixed (“mixture of 2 Gaussians”).

# Interpretation

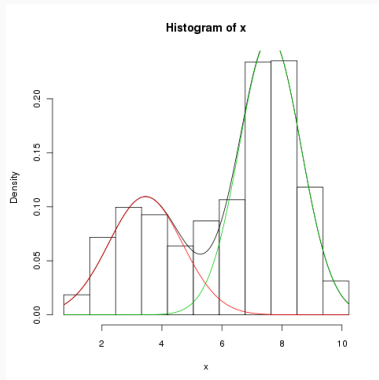
- Parameters:

$$\pi_1 = 1/3, \pi_2 = 2/3;$$

$$\mu_1 = 3 \text{ (kg)}, \mu_2 = 7 \text{ (kg)};$$

$$\sigma_1 = \sigma_2 = 2.$$

- If some rodent weighs 3 (kg), what is its probability to be a female ( $Z = 1$ ) ?



## Remark

- Possible extensions to  $p(\cdot|z = k)$  being in different parametric families.
- Example: mixtures of Weibull and Gamma distributions.

for  $x \geq 0$ ,  $p(x) =$

$$0.2 \frac{a}{b} \left(\frac{x}{b}\right)^{a-1} \exp\left[-\left(\frac{x}{b}\right)^a\right] + 0.8 \frac{c^k}{(k-1)!} x^{k-1} \exp(-cx).$$

# Identifiability Issues

- Generally, a parametric family of models  $\{p_{\Theta}\}_{\Theta \in \Theta}$  has identifiable parameter iff

$$\forall (\Theta, \Theta') \in \Theta^2, p_{\Theta} = p_{\Theta'} \Rightarrow \Theta = \Theta'.$$

- Ensures uniqueness of parameter (necessary condition for unique estimation).
- Identifiability cannot be achieved for mixtures even with fixed  $K$ , since for true  $\Theta = (\pi_1, \dots, \pi_K, \theta_1, \dots, \theta_K)$ , for any permutation  $\kappa \in \mathcal{S}_K$  of the labels (categorical value  $Z$ ), if set  $\Theta' = (\pi_{\kappa(1)}, \dots, \pi_{\kappa(K)}, \theta_{\kappa(1)}, \dots, \theta_{\kappa(K)})$ , we have

$$p_{\Theta} = \sum_{k=1}^K \pi_k p_{\theta_k} = \sum_{k=1}^K \pi_{\kappa(k)} p_{\theta_{\kappa(k)}} = p_{\Theta'}.$$

# Identifiability for Mixtures

- For mixture models: we define identifiability up to a permutation of the labels (equivalence classes), requiring that

$$\sum_{k=1}^K \pi_k p_{\theta_k} = \sum_{k=1}^{K'} \pi'_k p'_{\theta_k}$$
$$\Rightarrow K = K' \text{ and } \exists \kappa \in \mathcal{S}_K \forall k, \pi_k = \pi'_{\kappa(k)} \text{ and } \theta_k = \theta'_{\kappa(k)}.$$

- To ensure this, we constraint the mixture models to satisfy  $\forall k, \pi_k > 0$  (negative case: Zhang & Zhang).
- Additional sufficient condition for mixture identifiability:  $\{p_{\theta}\}_{\theta \in \Theta}$  are linearly independent PDFs.

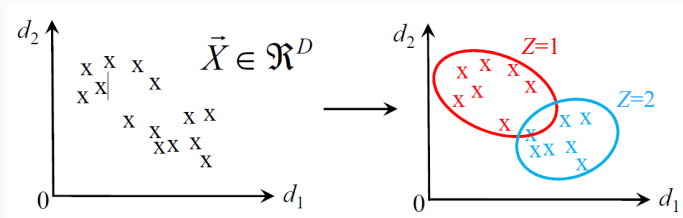
## Exercise 1

Proof that the mixtures of uniform distributions are not identifiable.

# Clustering with Mixtures in 3 Steps

**Assumption:**  $\{Z_n, X_n\}_{n=1}^N$  are independent. The families of  $p_{\theta_k}$  are known (good candidates).

1. Parameter estimation through Maximum Likelihood: learning  $\hat{\Theta}$  from an unlabelled sample set of size  $N$ .
2.  $\forall (n, k)$ , compute  $p_{\hat{\Theta}}(Z_n = k | X_n = x_n)$ .
3. MAP:  $\forall n, \hat{Z}_n = \arg \max_k p_{\hat{\Theta}}(Z_n = k | X_n = x_n)$ .



Maximum Likelihood Estimation (MLE):

$$\hat{\Theta} = \arg \max_{\Theta \in \mathcal{C}} \ell_{\mathbf{x}}(\Theta) = \arg \max_{\Theta \in \mathcal{C}} \sum_{n=1}^N \ln \underbrace{\left[ \sum_{k=1}^K \pi_k p_{\theta_k}(x_n) \right]}_{z \text{ is hidden inside}}$$

with  $\mathcal{C} = \{ (\pi_1, \dots, \pi_K, \theta_1, \dots, \theta_K) \mid \forall k, \pi_k \geq 0 \text{ and } \sum_k \pi_k = 1 \}$ .

Set derivatives to 0:

$$\frac{\partial \ell_{\mathbf{x}}}{\partial \pi_k}(\Theta) = 0, \quad \nabla_{\theta_k} \ell_{\mathbf{x}}(\Theta) = 0$$

No close form solution...

$$\hat{\Theta} = \arg \max_{\Theta \in \mathcal{C}} \sum_{n=1}^N \ln \left[ \sum_{k=1}^K \pi_k p_{\theta_k}(x_n) \right]$$

If the hidden states  $(z_1, \dots, z_N)$  are known:

$$\begin{aligned} \sum_{n=1}^N \sum_{k=1}^K \mathbb{1}_{\{z_n=k\}} \ln [\pi_k p_{\theta_k}(x_n)] &= \sum_{n=1}^N \sum_{k=1}^K \mathbb{1}_{\{z_n=k\}} \ln(\pi_k) + \\ &\quad \sum_{n=1}^N \sum_{k=1}^K \mathbb{1}_{\{z_n=k\}} \ln p_{\theta_k}(x_n). \end{aligned}$$

with the same complexity as to MLE estimation K different distributions within family  $\{p_{\theta}\}_{\theta \in \Theta}$  on K independent sample sets.



# EM Algorithm: Principle

$$\hat{\Theta} = \arg \max_{\Theta} \ln [p_{\Theta}(\mathbf{x})] = \arg \max_{\Theta} \ln \left[ \sum_{\mathbf{z}} p_{\Theta}(\mathbf{x}, \mathbf{z}) \right]$$

where:

- $\mathbf{x} = \{x_1, \dots, x_N\}$  observed,  $\mathbf{z} = \{z_1, \dots, z_N\}$  hidden discrete finite values labels.
- $p_{\Theta}(\mathbf{x}, \mathbf{z})$  easy to maximize,  $p_{\Theta}(\mathbf{x})$  difficult to maximize.

Principle:

- Maximize  $\ln p_{\Theta}(\mathbf{x}, \mathbf{z})$  instead of  $\ln p_{\Theta}(\mathbf{x})$ .
- Since  $\mathbf{z}$  is hidden, consider  $E_{(\mathbf{z}|\mathbf{x}, \Theta)}[\ln p_{\Theta}(\mathbf{X}, \mathbf{Z}) | \mathbf{X} = \mathbf{x}]$ .
- Maximize  $\sum_{\mathbf{z}} p(\mathbf{z}|\mathbf{x}, \Theta) \ln p_{\Theta}(\mathbf{x}, \mathbf{z})$  iteratively.

# EM Algorithm: Formulation

1. Initialize  $\Theta^{(i=0)}$ .
2. **E-step** (expectation): compute  $\mathcal{Q}(\Theta, \Theta^{(i)})$  that

$$\mathcal{Q}(\Theta, \Theta^{(i)}) = \sum_{\mathbf{z}} p(\mathbf{z}|\mathbf{x}, \Theta^{(i)}) \ln p_{\Theta}(\mathbf{x}, \mathbf{z})$$

or at least any relevant quantity.

3. **M-step** (maximization): update  $\Theta$  by

$$\Theta^{(i+1)} = \arg \max_{\Theta} \mathcal{Q}(\Theta, \Theta^{(i)})$$

4. Check the convergence of the log-likelihood  $\ln p_{\Theta}(\mathbf{x})$ .  
Back to step 2 if convergence criterion is not satisfied.

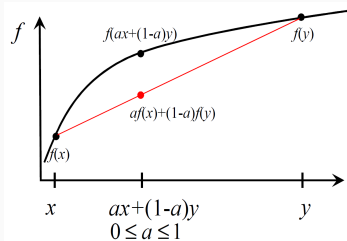
# EM Algorithm: Main Property

## Theorem (Dempster *et al.*, 1977)

$(\ln p_{\Theta^{(i)}}(x))_{i \geq 0}$  is a non-decreasing sequence.

$$\begin{aligned} \ln \left[ \sum_z p_{\Theta}(x, z) \right] &= \ln \left[ \sum_z q(z) \frac{p_{\Theta}(x, z)}{q(z)} \right] \\ &\geq \underbrace{\sum_z q(z) \ln \left[ \frac{p_{\Theta}(x, z)}{q(z)} \right]}_{\mathcal{L}(q, \Theta)} \end{aligned}$$

with equality when  $\frac{p_{\Theta}(x, z)}{q(z)} = C$

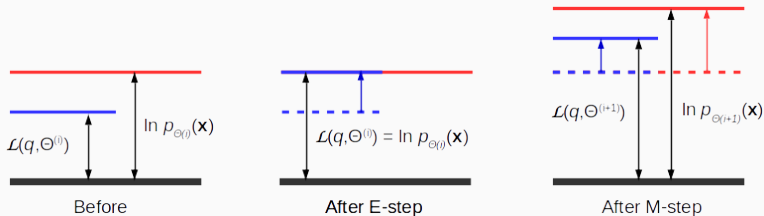


Some concave function  $f(\cdot)$

# EM Algorithm: Main Property

**E-step:** adjust  $q(z) = p(z|x, \Theta^{(i)})$  to make  $\ln p_{\Theta^{(i)}}(\mathbf{x}) = \mathcal{L}(q, \Theta^{(i)})$

**M-step:** adjust  $\Theta$  (derivative = 0) to maximize  $\mathcal{L}(q, \Theta^{(i)})$



- **Remark:**  $\ln p_{\Theta^{(i)}}(\mathbf{x})$  may converge to saddle point, local maximum or even not converge.

## Mixtures of Bernoulli distributions

Consider a set of  $M$  binary variables  $\mathbf{x} = (x_1, \dots, x_M)^\top$  follows

$$p(\mathbf{x}|\boldsymbol{\pi}, \boldsymbol{\mu}) = \sum_{k=1}^K \pi_k p(\mathbf{x}|\boldsymbol{\mu}_k) \text{ where } p(\mathbf{x}|\boldsymbol{\mu}_k) = \prod_{m=1}^M \mu_{km}^{x_m} (1 - \mu_{km})^{1-x_m}$$

- To cluster a given data set  $\mathbf{X} = \{\mathbf{x}_1, \dots, \mathbf{x}_N\}$ :

1. Estimate the parameters of the clusters  $\hat{\Theta} = \{\hat{\boldsymbol{\pi}}, \hat{\boldsymbol{\mu}}\}$

$$\text{E-step: } \gamma^{(i)}(z_{nk}) = p(Z_n = k | \mathbf{x}_n, \Theta^{(i)}) = \frac{\pi_k^{(i)} p(\mathbf{x}_n | \boldsymbol{\mu}_k^{(i)})}{\sum_{l=1}^K \pi_l^{(i)} p(\mathbf{x}_n | \boldsymbol{\mu}_l^{(i)})}.$$

$$\text{M-step: } \pi_k^{(i+1)} = \frac{N_k^{(i)}}{N}, \boldsymbol{\mu}_k^{(i+1)} = \frac{1}{N_k^{(i)}} \sum_{n=1}^N \gamma^{(i)}(z_{nk}) \mathbf{x}_n \text{ with } N_k^{(i)} = \sum_{n=1}^N \gamma^{(i)}(z_{nk}).$$

2. Estimate label  $\hat{Z}_n = \arg \max p(Z_n = k | \mathbf{x}_n, \hat{\Theta})$

## Exercise 2

Compute the so-called complete-data log-likelihood:

$$\ell_{\mathbf{x}, \mathbf{z}}(\Theta) = \sum_{n=1}^N \sum_{k=1}^K \mathbb{1}_{\{z_n=k\}} \ln [\pi_k p_{\theta_k}(x_n)]$$

and its maximizer  $\hat{\Theta}_{\mathbf{x}, \mathbf{z}}$  if  $X \in \mathbb{R}^d$  has conditional multivariate Gaussian distribution that  $p_{\theta_k}(X_n | Z_n = k) = \mathcal{N}(\mu_k, \Sigma_k)$ .

We have:

$$\nabla_{\mu} \left[ (x - \mu)^T \Sigma^{-1} (x - \mu) \right] = -2 \Sigma^{-1} (x - \mu)$$

$$\nabla_{\Sigma} \left[ (x - \mu)^T \Sigma^{-1} (x - \mu) \right] = -\Sigma^{-2} (x - \mu)(x - \mu)^T, \quad \nabla_{\Sigma} [\ln(\det(\Sigma))] = \Sigma^{-1}$$

## Exercise 3

- Read and answer the preparatory questions for next lab session (Mixture models).
- Give the reestimation formulas of the EM algorithm for mixtures with multivariate Gaussian distributions.

- [1] Arthur P Dempster, Nan M Laird, and Donald B Rubin.  
**Maximum likelihood from incomplete data via the em algorithm.**  
*Journal of the Royal Statistical Society: Series B (Methodological)*, 39(1):1–22, 1977.
  
- [2] James MacQueen et al.  
**Some methods for classification and analysis of multivariate observations.**  
In *Proceedings of the fifth Berkeley symposium on mathematical statistics and probability*, volume 1, pages 281–297. Oakland, CA, USA, 1967.
  
- [3] Geoffrey McLachlan and David Peel.  
**Finite mixture models.**  
John Wiley & Sons, 2004.



[4] Baibo Zhang and Changshui Zhang.

**Finite mixture models with negative components.**

*In International Workshop on Machine Learning and Data Mining in Pattern Recognition*, pages 31–41. Springer, 2005.