

# Clustering and Expectation Maximization

Shan-Hung Wu  
*shwu@cs.nthu.edu.tw*

Department of Computer Science,  
National Tsing Hua University, Taiwan

NetDB-ML, Spring 2015

## 1 Clustering

- Why Clustering?
- $k$ -Means Clustering
- Semiparametric Density Estimation
- Hierarchical Clustering
- Spectral Clustering
- Practical Considerations

## 2 Expectation Maximization

- Latent Variables and Complete Likelihood
- EM Steps
- EM for Mixture Models
- EM for Mixtures of Gaussians

## 1 Clustering

- Why Clustering?
- $k$ -Means Clustering
- Semiparametric Density Estimation
- Hierarchical Clustering
- Spectral Clustering
- Practical Considerations

## 2 Expectation Maximization

- Latent Variables and Complete Likelihood
- EM Steps
- EM for Mixture Models
- EM for Mixtures of Gaussians

# Clustering

- We now consider the unsupervised datasets  $\mathcal{X} = \{\mathbf{x}^{(t)}\}_{t=1}^N$  where labels  $\mathbf{r}^{(t)}$  are missing
  - Learning the a posteriori knowledge from unlabeled data is called the **unsupervised learning**
- **Clustering** is one unsupervised learning technique used to identify the groups  $G_1, \dots, G_K$  in each which instances are similar (or close) to each other
  - $K$  could be either predefined (a hyperparameter) or not (a parameter)
- Output:  $\mathcal{Z} := \{\mathbf{z}^{(t)}\}_t$ , where
  - Hard labeling:  $\mathbf{z}^{(t)} \in \{0, 1\}^K$  and  $z_i^{(t)} = 1$  iff the instance  $t$  belongs to group  $i$
  - Soft labeling:  $\mathbf{z}^{(t)} \in \mathbb{R}^K$  and  $z_i^{(t)}$  denotes the degree (e.g., probability) the instance  $t$  belongs to group  $i$

# Applications

- Pattern recognition: groups may be meaningful
  - E.g., product/user cluster in market analysis
- Compression: instances in the same group can be represented by a prototype
- Data labeling: groups are good hints for labels
- Data reprocessing for classification/regression: attributes of instances can be augmented by group information; or we can identify groups in each class to estimate  $P[\mathbf{x}|C_i]$  and  $P[C_i]$  more precisely
- And so on...

# Clustering vs. Dimensionality Reduction

- In dimensionality reduction, we find correlations between *attributes* and “group” (i.e., select/extract) attributes
- In clustering, we find similarities between *instances* and group instances

## 1 Clustering

- Why Clustering?
- $k$ -Means Clustering
- Semiparametric Density Estimation
- Hierarchical Clustering
- Spectral Clustering
- Practical Considerations

## 2 Expectation Maximization

- Latent Variables and Complete Likelihood
- EM Steps
- EM for Mixture Models
- EM for Mixtures of Gaussians

# K-Means Clustering (1/2)

- Suppose each group  $G_i$  is parametrized by a prototype  $\mathbf{m}_i$ , the mean of all instances in this group
- Hard labeling:  $z_i^{(t)} = 1$  iff  $\mathbf{x}^{(t)}$  is the closest to  $\mathbf{m}_i$ ; i.e.,  
$$\|\mathbf{x}^{(t)} - \mathbf{m}_i\| = \min_j \|\mathbf{x}^{(t)} - \mathbf{m}_j\|$$
- The objective of **K-means clustering** is to find  $\mathbf{m}_i$  such that the total **reconstruction error**  $rec(\{\mathbf{m}_i\}_{i=1}^K; \mathcal{X}) = \sum_{t=1}^N \sum_{i=1}^K z_i^{(t)} \|\mathbf{x}^{(t)} - \mathbf{m}_i\|^2$  is minimized



# K-Means Clustering (2/2)

**Input:**  $\mathcal{X} \leftarrow \{\mathbf{x}^{(t)}\}_{t=1}^N$ ,  $K$

**Output:** The prototypes  $\mathbf{m}_i$ ,  $1 \leq i \leq K$

Initialize each  $\mathbf{m}_i$  to a random example  $\mathbf{x}^{(t)}$ ;

**repeat**

**foreach**  $\mathbf{x}^{(t)} \in \mathcal{X}$  **do**

$z_i^{(t)} \leftarrow \begin{cases} 1 & \text{if } \|\mathbf{x}^{(t)} - \mathbf{m}_i\| = \min_j \|\mathbf{x}^{(t)} - \mathbf{m}_j\| \\ 0 & \text{otherwise} \end{cases}$  ;

**end**

**foreach**  $\mathbf{m}_i$  **do**

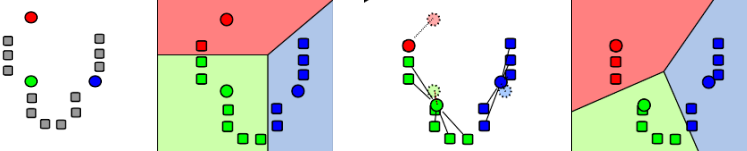
$\mathbf{m}_i \leftarrow \sum_{t=1}^N z_i^{(t)} \mathbf{x}^{(t)} / \sum_{t=1}^N z_i^{(t)}$ ;

**end**

**until** all  $\mathbf{m}_i$  converge;

**Algorithm 1:** The  $K$ -means algorithm.

# Example

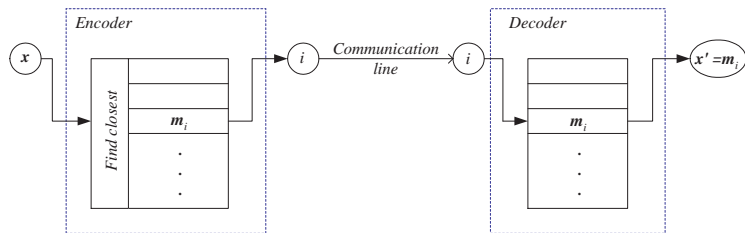


# Applications (1/2)

- One famous application of the  $K$ -means clustering is **vector quantization**, which aims to find a discrete set of vectors  $\{\mathbf{m}_i\}_{i=1}^K$  representative of the whole, possibly continuous, set of data points
  - E.g., in **color quantization**, we seek the best 256 colors of a 24 bits/pixel (16 million) color image
  - Once we get these 256 colors, for each pixel we only need to store the 8 bits color index
- We can quantize the 16 million colors uniformly into 256, but some of these 256 colors may be wasted when there is no nearby color appears in the image
  - We want nonuniform quantization where  $\mathbf{m}_i$  sit at the most dense areas of the whole dataset
- The  $K$ -means clustering minimizes  $rec(\{\mathbf{m}_i\}_{i=1}^K; \mathcal{X}) = \sum_{t=1}^N \sum_{i=1}^K z_i^{(t)} \|\mathbf{x}^{(t)} - \mathbf{m}_i\|^2$  and finds prototypes at the center of the dense regions

## Applications (2/2)

- Another example is the use of **codebooks** in telecommunication systems
  - Each point in the dataset is a vector storing the sample of a voice signal
  - We want to quantize samples into  $K$  representative vectors
  - If we store these  $K$  vectors in each device, the signal can be sent by indexes (of  $\lg K$  bits each) only



**Figure :** Given  $x$ , the encoder sends the index  $i$  of the nearest **codeword**  $m_i$  and the decoder receives  $x' = m_i$ . The error is  $\|x - x'\|^2$ .

# Limitations (1/2)

- The main disadvantage of the  $K$ -means clustering is that it is a local search procedure
  - The final prototypes  $\mathbf{m}_i$  may not be the optimal ones, and highly depend on the initial  $\mathbf{m}_i$
  - Can you give an example dataset based on which the  $K$ -means returns bad clusters? [Homework]
- Generally, the initial  $\mathbf{m}_i$  should a) locate at regions where instances occur; b) be far away from each other
- The  $K$ -means++ proposes one possible initialization step:
  - 1 Choose an instance uniformly at random to be  $\mathbf{m}_1$
  - 2 For each  $\mathbf{x}^{(t)}$ , compute  $d(\mathbf{x}^{(t)})$ , the distance between  $\mathbf{x}^{(t)}$  and the nearest  $\mathbf{m}_i$  that has already been determined
  - 3 Assign another instance to  $\mathbf{m}_{i+1}$ , but this time an instance  $\mathbf{x}$  is chosen with probability  $\frac{d(\mathbf{x})^2}{\sum_{t=1}^N d(\mathbf{x}^{(t)})^2}$
  - 4 Repeat Steps 2 and 3 until  $K$  initial prototypes are determined

## Limitations (2/2)

- Another shortcoming of the  $K$ -means is that clusters are assumed to be spherical and with equal size
  - Due to that the Euclidean distance is used when updating the cluster assignment  $z_i^{(t)}$  for each instance
- In practice, clusters may have different sizes
  - Next, we see how the above assumption can be relaxed using the probability framework we are already familiar

## 1 Clustering

- Why Clustering?
- $k$ -Means Clustering
- Semiparametric Density Estimation
- Hierarchical Clustering
- Spectral Clustering
- Practical Considerations

## 2 Expectation Maximization

- Latent Variables and Complete Likelihood
- EM Steps
- EM for Mixture Models
- EM for Mixtures of Gaussians

# Mixture Models

- Basic assumption: the dataset  $\mathcal{X}$  is a mixture of groups  $G_1, \dots, G_K$ 
  - E.g., in the hand-written digit recognition,  $\mathcal{X}$  consists of images of “0,” “1,” “2,” and so forth
  - Even if  $\mathcal{X}$  are images of the same digit (say “1”) there are still typical different ways to write the digit (with or without head)
- Soft labeling:  $\mathcal{Z} = \{\mathbf{z}^{(t)} \in \mathbb{R}^K\}_t$
- The **mixture density** of an instance  $\mathbf{x}$  can be expressed as
$$p(\mathbf{x}) = \sum_{i=1}^K p(\mathbf{x}|G_i)P[G_i]$$
- Model: a collection of groups, i.e.,  $\{G_i\}_{i=1}^K$
- Parameters:  $G_1, \dots$ , and  $G_K$
- Objective:  $\arg_{G_1, \dots, G_K} \max \prod_{t=1}^N \sum_{i=1}^K p(\mathbf{x}^{(t)}|G_i)P[G_i]d\mathbf{x}$
- A generative model this case



# Parametric vs. Nonparametric vs. Semiparametric

- Parametric models: models that can be completely described by (a small number of) parameters
- Nonparametric models: those that cannot be described by parameters
- Semiparametric models: those that can be partially described by parameters
  - Each cluster is parametric
  - But the mixture of clusters,  $\mathcal{Z} = \{\mathbf{z}^{(t)}\}_t$ , is not (i.e., we do not assume the mixture to follow some distribution)

# Semiparametric Clustering vs. Parametric Classification

- Parametric classification is a special case of mixture model where the groups (i.e., classes) are known in advance:

$$P[\mathbf{x}^{(t)}] = \sum_{i=1}^K p(\mathbf{x}^{(t)}|C_i)P[C_i]d\mathbf{x}$$

- Assume  $p(\mathbf{x}^{(t)}|C_i)d\mathbf{x}$  and  $P[C_i]$  follow Gaussian and Bernoulli distributions parametrized by  $\theta_i = (\boldsymbol{\mu}_i, \boldsymbol{\Sigma}_i)$  and  $\theta'_i = p_i$  respectively
- Since we know which instance belongs to which class by  $\mathbf{r}^{(t)}$ , we can estimate  $\theta_i$  and  $\theta'_i$  analytically by maximizing  $P[\mathcal{X}|\theta_i]$  and  $P[\mathcal{X}|\theta'_i]$ :
  - $\hat{p}_i = \frac{N_i}{N}$ , where  $N_i = \sum_{t=1}^N r_i^{(t)}$
  - $\mathbf{m}_i = \frac{1}{N_i} \sum_{t=1}^N \mathbf{x}^{(t)} r_i^{(t)}$  and  $\mathbf{S}_i = \frac{1}{N_i - 1} \sum_{t=1}^N r_i^{(t)} (\mathbf{x}^{(t)} - \mathbf{m}_i)(\mathbf{x}^{(t)} - \mathbf{m}_i)^\top$
- Unfortunately, in semiparametric clustering we don't know  $\mathbf{z}^{(t)}/\mathbf{r}^{(t)}$  so we cannot solve  $p(\mathbf{x}^{(t)}|C_i)d\mathbf{x}$  and  $P[C_i]$  analytically

# If No $\{\theta_i, \theta'_i\}_{i=1}^K$ , Make Them Up

- How?

# If No $\{\theta_i, \theta'_i\}_{i=1}^K$ , Make Them Up

- How? Borrowing the iterations from  $K$ -means
- Start from a random guess of  $\{\theta_i, \theta'_i\}_{i=1}^K$  and then perform the following two steps iteratively:
  - 1 For each instance  $\mathbf{x}^{(t)}$ , update its  $\mathbf{z}^{(t)}$  based on the current  $G_1, \dots, G_K$  parametrized by  $\theta_1, \dots, \theta_K$
  - 2 Update  $\theta_1, \dots, \theta_K$  based on the current  $\mathbf{z}^{(t)}$
- Stop until the groups do not change in Step 2 (or the changes of groups are smaller than a threshold  $\epsilon$ )

# Semiparametric Density Estimation (1/2)

- Suppose in the mixture density  $p(\mathbf{x}) = \sum_{i=1}^K p(\mathbf{x}|G_i)P[G_i]$ , each  $p(\mathbf{x}^{(t)}|G_i)$  and  $P[G_i]$  are Gaussian and Bernoulli distributions parametrized by  $\theta_i = (\boldsymbol{\mu}_i, \boldsymbol{\Sigma}_i)$  and  $\theta'_i = \pi_i$  respectively
- Denote the collection of estimators by  $\Theta = (\mathbf{m}_i, \mathbf{S}_i, \pi_i)_{i=1}^K$
- We guess initial  $\Theta$ , and then:

$$\textcircled{1} \text{ Update mixture: } z_i^{(t)} = P[z_i^{(t)}|\mathbf{x}^{(t)}; \Theta] = \frac{p(\mathbf{x}^{(t)}|z_i^{(t)}; \Theta)P[z_i^{(t)}; \Theta]}{p(\mathbf{x}^{(t)}; \Theta)} = \frac{p(\mathbf{x}^{(t)}|z_i^{(t)}; \Theta)\pi_i}{\sum_{j=1}^K p(\mathbf{x}^{(t)}|z_j^{(t)}; \Theta)\pi_j} = \frac{\det(\mathbf{S}_i)^{-1/2} \exp[-(1/2)(\mathbf{x}^{(t)} - \mathbf{m}_i)^\top \mathbf{S}_i^{-1}(\mathbf{x}^{(t)} - \mathbf{m}_i)]\pi_i}{\sum_{j=1}^K \det(\mathbf{S}_j)^{-1/2} \exp[-(1/2)(\mathbf{x}^{(t)} - \mathbf{m}_j)^\top \mathbf{S}_j^{-1}(\mathbf{x}^{(t)} - \mathbf{m}_j)]\pi_j}$$

- Unlike in  $K$ -means, we assign **soft labels** to  $z_i^{(t)}$
- Update  $\Theta$ : knowing  $z_i^{(t)}$ , we can update  $\pi_i$ ,  $\mathbf{m}_i$ , and  $\mathbf{S}_i$  by, e.g., maximizing the likelihood  $P[\mathcal{X}|\Theta]$

# Semiparametric Density Estimation (2/2)

**Input:**  $\mathcal{X} \leftarrow \{\mathbf{x}^{(t)}\}_{t=1}^N$ ,  $K$

**Output:**  $\Theta = (\mathbf{m}_i, \mathbf{S}_i, \pi_i)_{i=1}^K$

Initialize  $\Theta$  by performing several iterations of  $K$ -means;

**repeat**

**foreach**  $\mathbf{x}^{(t)} \in \mathcal{X}$  **do**

$$z_i^{(t)} \leftarrow \frac{\det(\mathbf{S}_i)^{-1/2} \exp[-(1/2)(\mathbf{x}^{(t)} - \mathbf{m}_i)^\top \mathbf{S}_i^{-1} (\mathbf{x}^{(t)} - \mathbf{m}_i)] \pi_i}{\sum_{j=1}^K \det(\mathbf{S}_j)^{-1/2} \exp[-(1/2)(\mathbf{x}^{(t)} - \mathbf{m}_j)^\top \mathbf{S}_j^{-1} (\mathbf{x}^{(t)} - \mathbf{m}_j)] \pi_j};$$

**end**

**foreach**  $\pi_i$ ,  $\mathbf{m}_i$ , **and**  $\mathbf{S}_i$  **do**

$$\pi_i \leftarrow \frac{\sum_{t=1}^N z_i^{(t)}}{N};$$

$$\mathbf{m}_i \leftarrow \frac{\sum_{t=1}^N \mathbf{x}^{(t)} z_i^{(t)}}{\sum_{t=1}^N z_i^{(t)}} \text{ and } \mathbf{S}_i \leftarrow \frac{\sum_{t=1}^N z_i^{(t)} (\mathbf{x}^{(t)} - \mathbf{m}_i) (\mathbf{x}^{(t)} - \mathbf{m}_i)^\top}{\sum_{t=1}^N z_i^{(t)}};$$

**end**

**until**  $\Theta$  *converges*;

**Algorithm 2:** Semiparametric density estimation for Gaussian mixtures.

# Simplifications

- As in parametric classification, with small training dataset and large dimensionality we can regularize our model by simplifying assumptions
- When the priors  $P[G_i] = \pi_i$  are all equal and  $S_i = s^2\mathbf{I}$ , we have
$$z_i^{(t)} = \frac{\exp[-(1/2s^2)\|\mathbf{x}^{(t)} - \mathbf{m}_i\|^2]}{\sum_{j=1}^K \exp[-(1/2s^2)\|\mathbf{x}^{(t)} - \mathbf{m}_j\|^2]}$$
- We thus see that the  $K$ -means clustering is just a special case of the semiparametric density estimation applied to Gaussian mixtures, where
  - Attributes of instances are independent and with equal variance
  - All groups have equal priors
  - Labels are hardened

## 1 Clustering

- Why Clustering?
- $k$ -Means Clustering
- Semiparametric Density Estimation
- **Hierarchical Clustering**
- Spectral Clustering
- Practical Considerations

## 2 Expectation Maximization

- Latent Variables and Complete Likelihood
- EM Steps
- EM for Mixture Models
- EM for Mixtures of Gaussians



# Hierarchical Clustering (1/2)

- So far, we assume that clusters are independent groups (although they may overlap)
- In some applications, we may want to find the hierarchy of clusters
- Two common types of algorithms:
  - **Agglomerative**: Starting from  $N$  groups, each with single instance, iteratively merging two most similar groups to form a larger one, until there remains a single group
  - **Divisive**: Starting one group containing all instances, dividing large groups into smaller ones, until there are  $N$  groups

## Hierarchical Clustering (2/2)

- When deciding which groups should be merged (or split), a measure of similarity, or equivalently distance  $d$ , is required
  - One common choice is the Minkowski distance:
$$d(\mathbf{x}^{(r)}, \mathbf{x}^{(s)}) = \left( \sum_{i=1}^d |x_i^{(r)} - x_i^{(s)}|^p \right)^{1/p} \text{ for some } p$$
- But how to calculate the distance between two groups?

# Hierarchical Clustering (2/2)

- When deciding which groups should be merged (or split), a measure of similarity, or equivalently distance  $d$ , is required

- One common choice is the Minkowski distance:

$$d(\mathbf{x}^{(r)}, \mathbf{x}^{(s)}) = \left( \sum_{i=1}^d |x_i^{(r)} - x_i^{(s)}|^p \right)^{1/p} \text{ for some } p$$

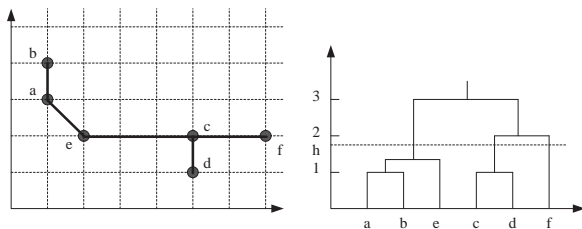
- But how to calculate the distance between two groups?

- **Single-link** metric:  $d(G_i, G_j) = \min_{\mathbf{x}^{(r)} \in G_i, \mathbf{x}^{(s)} \in G_j} d(\mathbf{x}^{(r)}, \mathbf{x}^{(s)})$

- **Complete-link** metric:  $d(G_i, G_j) = \max_{\mathbf{x}^{(r)} \in G_i, \mathbf{x}^{(s)} \in G_j} d(\mathbf{x}^{(r)}, \mathbf{x}^{(s)})$

# Dendrograms

- The result of hierarchical clustering can be shown as the *dendrogram*:



- Each internal node corresponds to a group
- The height of the internal node denote the distance between groups
- The dendrogram can be intersected at a user-specific level  $h$  to get the clusters
- In each cluster, instances in the input space are connected as a tree

# Single- or Complete-Link?

# Single- or Complete-Link?

- With the complete-link metric, all instance in a group have distance less than  $h$ 
  - Assumes that each cluster is spherical
  - Similar to k-means and semiparametric density estimation
  - Used only when this assumption is likely to be true
- Single-link clusters may have diameter (i.e., the greatest length of the shortest paths between instances) much larger than  $h$ 
  - With the single-link metric, two instance are grouped together at level  $h$  if
    - The distance between them is less than  $h$ ; or
    - There exists a path between them such that any two consecutive instances along the path have mutual distance less than  $h$
  - Each final cluster may have an arbitrary shape
  - Suitable for clusters backed by respective underlying manifolds

## 1 Clustering

- Why Clustering?
- $k$ -Means Clustering
- Semiparametric Density Estimation
- Hierarchical Clustering
- Spectral Clustering
- Practical Considerations

## 2 Expectation Maximization

- Latent Variables and Complete Likelihood
- EM Steps
- EM for Mixture Models
- EM for Mixtures of Gaussians

# Global vs. Local Models

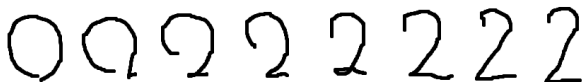
- We have seen models that find clusters by assuming some structure for each cluster
- Global structure: each cluster represents a dense region of a known shape
  - E.g.,  $k$ -means, semiparametric density estimation, hierarchical clustering with complete-link metric
- Local similarity: each instance in a cluster is similar to its nearby instances
  - E.g., hierarchical clustering with single-link metric
  - Local models can produce clusters of arbitrary shapes
  - Suitable to datasets where clusters are backed by respective underlying manifolds

A series of five handwritten digits '2' in a cursive style, arranged horizontally. The digits are slightly different in shape and orientation, illustrating local similarity where each instance is similar to its neighbors.



# More Local Models

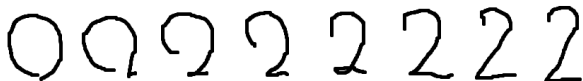
- In local models, any two instances in the same cluster are *not* necessarily similar
  - This is both an advantage and disadvantage
- Cons: they tends to find clusters of unbalanced sizes



- Outliers form singleton clusters
- How to make clusters balanced?

# More Local Models

- In local models, any two instances in the same cluster are *not* necessarily similar
  - This is both an advantage and disadvantage
- Cons: they tends to find clusters of unbalanced sizes



- Outliers form singleton clusters
- How to make clusters balanced?
- We consider only the the flat clustering in the next

# Balanced Cut of Local Similarity Graph (1)

- Given a set of data points  $\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(N)}$ . Let  $\mathbf{S} \in \mathbb{R}^{N \times N}$  be the local similarity matrix where  $s_{ij} \geq 0$  is the similarity between instance between instances  $i$  and  $j$  *if they are neighbors*
  - Euclidean distance is clearly not a good choice
- Local similarity measure?

# Balanced Cut of Local Similarity Graph (1)

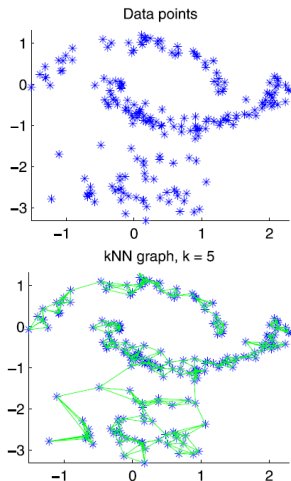
- Given a set of data points  $\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(N)}$ . Let  $\mathbf{S} \in \mathbb{R}^{N \times N}$  be the local similarity matrix where  $s_{ij} \geq 0$  is the similarity between instance between instances  $i$  and  $j$  *if they are neighbors*
  - Euclidean distance is clearly not a good choice
- Local similarity measure?
  - $\epsilon$ -NN similarity:  $s_{ij}$  inverse proportional to the Euclidean distance between  $i$  and  $j$  if  $i$  is a  $\epsilon$ -nearest neighbor of  $j$  or vice versa; otherwise 0
  - Gaussian similarity (soft  $\epsilon$ -NN):  $s_{ij} = \exp(-\frac{\|\mathbf{x}^{(i)} - \mathbf{x}^{(j)}\|^2}{\sigma^2})$  for some hyperparameter  $\sigma$

# Balanced Cut of Local Similarity Graph (2)

- Consider the graph  $G = (V, E)$  where  $V$  denotes the set of instances and  $E$  denotes the set of non-zero local similarity scores
- Given a set of nodes  $A \subset V$ , define  $Cut(A) := \sum_{i \in A, j \notin A} S_{ij}$
- We want to find a  $k$ -partition  $A_1, \dots, A_K$  of  $V$  that solve the problem:

$$\arg \min_{A_1, \dots, A_K \subset V} RatioCut(A_1, \dots, A_K) := \frac{1}{2} \sum_{i=1}^K \frac{Cut(A_i)}{|A_i|}$$

- Cross-partition links are edges are minimized
- $|A_1|, \dots, |A_K|$  are balanced



- Unfortunately, although the min-cut problems can be solved efficiently, the balanced min-cut problems are NP-hard
- **Spectral clustering** solves a relaxation of the above problem
  - Finds the eigenvectors of a **graph Laplacian matrix** induced from the local similarity graph
  - Efficient

# Graph Laplacian

- Given a (local) similarity matrix  $S$ , the **graph Laplacian matrix** is defined as

$$\mathbf{L} = \mathbf{D} - \mathbf{S},$$

where  $\mathbf{D}$  is an  $N \times N$  diagonal matrix with  $d_i = \sum_{j=1}^N s_{ij}$  on the diagonal

- For any vector  $\mathbf{f} \in \mathbb{R}^N$ , we have  $\mathbf{f}^\top \mathbf{L} \mathbf{f} = \frac{1}{2} \sum_{i,j=1}^N s_{ij} (f_i - f_j)^2$   
[Homework]
- $\mathbf{L}$  is symmetric and positive semi-definite
  - The smallest eigenvalue of  $\mathbf{L}$  is 0, and the constant one vector  $\mathbf{1} \in \mathbb{R}^N$  must be (one of) the corresponding eigenvector
  - $\mathbf{L}$  has  $N$  non-negative eigenvalues  $0 = \lambda_1 \leq \dots \leq \lambda_N$ .

# Spectral Clustering

- Idea: map each  $\mathbf{x}^{(t)} \in \mathbb{R}^N$  to  $\mathbf{z}^{(t)} \in \mathbb{R}^m$  in some low dimensional space such that  $\mathbf{z}^{(i)}$  and  $\mathbf{z}^{(j)}$  are similar if they belong to the same cluster
  - Then apply a traditional clustering algorithm (e.g.,  $k$ -means) to obtain the final cluster
- Based on  $\mathbf{f}^\top \mathbf{L} \mathbf{f} = \frac{1}{2} \sum_{i,j=1}^N s_{i,j} (f_i - f_j)^2$ , we can first solve

$$\arg \min_{\mathbf{F} = [\mathbf{f}^{(1)}, \dots, \mathbf{f}^{(m)}] \in \mathbb{R}^{N \times m}} \text{tr}(\mathbf{F}^\top \mathbf{L} \mathbf{F}) = \sum_{i=1}^m \mathbf{f}^{(i)\top} \mathbf{L} \mathbf{f}^{(i)},$$

subject to  $\mathbf{F}^\top \mathbf{F} = \mathbf{I}$

and then let  $\mathbf{z}^{(t)}$  be the  $t$ -th row of  $\mathbf{F}$

- $\mathbf{f}_i$  and  $\mathbf{f}_j$  are orthogonal so that they provide complementary perspectives
- Each  $\mathbf{f}_i$  is normalized so that the clusters are balanced (to be explained later)



# Spectrum of $L$ (1)

- From the Rayleigh-Ritz theorem,  $\mathbf{f}^{(1)}, \dots, \mathbf{f}^{(m)}$  are the eigenvectors corresponding to the smallest eigenvalues of  $L$

## Theorem

*Let  $G = (V, E)$  be an undirected graph with non-negative weights. Then the multiplicity  $K$  of the eigenvalue 0 of  $L$  equals the number of connected components  $A_1, \dots, A_K \subset V$  in the graph. The eigenspace of eigenvalue 0 is spanned by the indicator vectors  $\mathbf{1}_{A_1}, \dots, \mathbf{1}_{A_K} \in \mathbb{R}^N$  of those components.*

## Spectrum of $L$ (2)

### Proof.

Assume that  $\mathbf{f}$  is an eigenvector with eigenvalue 0. We know that  $0 = \mathbf{f}^\top \mathbf{L} \mathbf{f} = \sum_{i,j=1}^N s_{ij} (f_i - f_j)^2$ . As  $s_{ij}$  is non-negative, the sum can only vanish if all terms vanish. Thus, if two vertices  $v_i$  and  $v_j$  are connected (i.e.,  $s_{ij} > 0$ ), then  $f_i = f_j$ . When  $K = 1$ ,  $\mathbf{f}$  needs to be constant one vector and  $\mathbf{L}$  has eigenvalue 0 with multiplicity 1. When  $K > 1$ , without loss of generality we assume that the vertices are ordered according to the connected components they belong to. Then  $\mathbf{S}$  has a block diagonal form, and the same is true for  $\mathbf{L}$ :

$$\mathbf{L} = \begin{pmatrix} \mathbf{L}_1 & & \\ & \ddots & \\ & & \mathbf{L}_K \end{pmatrix}.$$

Since the spectrum of  $\mathbf{L}$  is given by the union of the spectra of  $\mathbf{L}_i$ , and the corresponding eigenvectors of  $\mathbf{L}$  are the eigenvectors of  $\mathbf{L}_i$ , filled with 0 at the positions of the other blocks. □

## Spectrum of $L$ (3)

- Based on the above theorem, we should make sure that #connected components  $< m$  when constructing the local similarity graph
  - Otherwise, some cluster may contain one connected component, and some may contain multiple
- In practice, we usually construct a fully-connected graph
  - The eigenvector of 0 is  $\mathbf{1}$
- Other than  $\mathbf{1}$ , what  $\mathbf{f}$  makes  $\mathbf{f}^\top \mathbf{L} \mathbf{f} = \frac{1}{2} \sum_{i,j=1}^N s_{ij} (f_i - f_j)^2$  small?
  - Those  $\mathbf{f}$ 's with value levels
  - Coordinates corresponding to the same  $G_i$  have the same value (forming a level)
  - The gap between different levels corresponds to the min-cuts

## Spectrum of $L$ (4)

- Besides,  $\|f\| = 1$  makes gap correspond to the balanced min-cuts
- For example, suppose  $K = 2$ . Let

$$f_i = \begin{cases} \sqrt{\frac{|G|}{|V||G|}}, & \text{if } \mathbf{x}^{(i)} \in G \\ -\sqrt{\frac{|G|}{|V||G|}}, & \text{otherwise} \end{cases}.$$

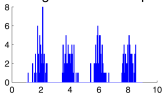
We have

$$\text{tr}(\mathbf{f}^\top \mathbf{L} \mathbf{f}) = \text{RatioCut}(G, \bar{G}),$$

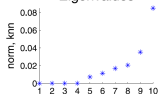
furthermore,  $\mathbf{f}^\top \mathbf{1} = 0$  and  $\|\mathbf{f}\| = 1$  [Homework]

# $m = K$ is enough due to orthogonality

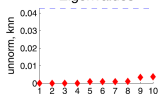
Histogram of the sample



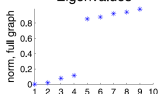
Eigenvalues



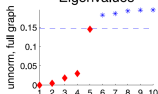
Eigenvalues



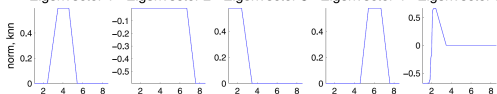
Eigenvalues



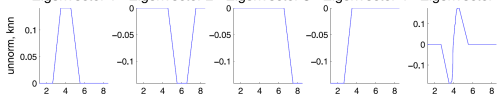
Eigenvalues



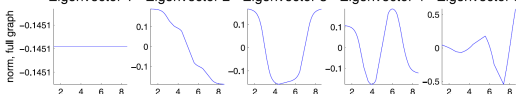
Eigenvector 1 Eigenvector 2 Eigenvector 3 Eigenvector 4 Eigenvector 5



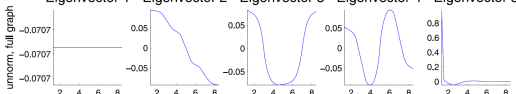
Eigenvector 1 Eigenvector 2 Eigenvector 3 Eigenvector 4 Eigenvector 5



Eigenvector 1 Eigenvector 2 Eigenvector 3 Eigenvector 4 Eigenvector 5



Eigenvector 1 Eigenvector 2 Eigenvector 3 Eigenvector 4 Eigenvector 5



# Spectral Clustering Algorithms

**Input:** Similarity matrix  $\mathbf{S}$ , number of clusters  $K$

**Output:** Clusters  $A_1, \dots, A_K$

Compute the Laplacian  $\mathbf{L}$ ;

Compute the first  $K$  eigenvectors  $\mathbf{u}_1, \dots, \mathbf{u}_K$  of  $\mathbf{L}$ . Let  $\mathbf{U} \in \mathbb{R}^{N \times K}$  be the matrix containing the vectors  $\mathbf{u}_1, \dots, \mathbf{u}_K$  as columns.;

For  $i = 1, \dots, N$ , let  $\mathbf{y}_i \in \mathbb{R}^K$  be the vector corresponding to the  $i$ -th row of  $\mathbf{U}$ ;

Cluster the points  $(\mathbf{y}_i)_{i=1, \dots, N}$  with the  $K$ -means algorithm;

**Algorithm 3:** The spectral clustering algorithm.

# Pros and Cons

- Pros:
  - Local model, balanced
  - Efficient event for large datasets (as  $S$  is sparse)
  - No issue of getting stuck in local minimum (e.g., as in  $k$ -means due to bad initializations)
- Cons:
  - Performance sensitive to the quality of the local similarity graph
  - Relaxation is loose: no guarantee that the final clusters correspond to the balanced min-cuts
- Which local similarity is better?
  - Empirically,  $\epsilon$ -NN graph is less vulnerable to the imperfect choice of parameters ( $\epsilon$ ,  $\sigma$ )
  - Graph is sparse

## 1 Clustering

- Why Clustering?
- $k$ -Means Clustering
- Semiparametric Density Estimation
- Hierarchical Clustering
- Spectral Clustering
- Practical Considerations

## 2 Expectation Maximization

- Latent Variables and Complete Likelihood
- EM Steps
- EM for Mixture Models
- EM for Mixtures of Gaussians



# Evaluating the Clusters

- How to evaluate the clusters we found?

# Evaluating the Clusters

- How to evaluate the clusters we found?
- If labels are not available during evaluation:
  - $\frac{\text{intercluster separation}}{\text{intracluster cohesion}} = \frac{\sum_{i,j} (m_i - m_j)^2}{\sum_i \frac{1}{|G_i|} \sum_{x \in G_i} (x - m_i)^2}$  (the higher the better)
- If labels (i.e.,  $\{\mathbf{r}^{(t)}\}_{t=1}^N$ ,  $\mathbf{r}^{(t)} \in \mathbb{R}^K$ ) are available during the evaluation:
  - $\text{entropy}(G_i) = -\sum_{j=1}^K P_i[r_j^{(t)} = 1] \lg P_i[r_j^{(t)} = 1]$ , where  $P_i[r_j^{(t)} = 1]$  denotes the portions of instances in  $G_i$  which belong to class  $j$ 
    - Here we define  $\lg 0 = 0$
  - $\text{entropy}_{total}(\mathcal{X}) = \sum_{i=1}^K \frac{|G_i|}{N} \text{entropy}(G_i)$  (the lower the better)
- Indirect evaluation: if clustering is used to help perform another task, then we can measure the performance of that task instead
  - E.g., click-through rate of the recommended item in a website (where clustering is used to group similar items/users)

# Deciding the Number of Clusters $K$

- In the previous semiparametric methods,  $K$  is determined in advance
  - We can decide  $K$  using the cross validation technique
  - Plot the reconstruction error against  $K$  and pick the “elbow”
- In hierarchical clustering,  $K$  is decided along with  $h$ 
  - $h$  should be set to cut the “big jump”
- $K$  can be either a parameter or a hyperparameter
- There are extensions for semiparametric methods that adapt  $K$  during the iteration
  - E.g.?

# Deciding the Number of Clusters $K$

- In the previous semiparametric methods,  $K$  is determined in advance
  - We can decide  $K$  using the cross validation technique
  - Plot the reconstruction error against  $K$  and pick the “elbow”
- In hierarchical clustering,  $K$  is decided along with  $h$ 
  - $h$  should be set to cut the “big jump”
- $K$  can be either a parameter or a hyperparameter
- There are extensions for semiparametric methods that adapt  $K$  during the iteration
  - E.g.? at each iteration, we can drop groups that are too small and/or split groups that are too large

## 1 Clustering

- Why Clustering?
- $k$ -Means Clustering
- Semiparametric Density Estimation
- Hierarchical Clustering
- Spectral Clustering
- Practical Considerations

## 2 Expectation Maximization

- Latent Variables and Complete Likelihood
- EM Steps
- EM for Mixture Models
- EM for Mixtures of Gaussians

# Why Iterative Methods Work?

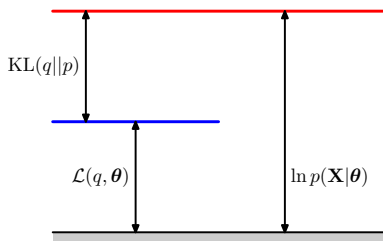
- We have seen the iterative methods for clustering
  - $K$ -means
  - Semiparametric density estimation
- But we haven't answered the following questions:
  - Why does the iteration end?
  - Why is the clusters found in Step 2 better than the ones found in the previous iteration?

# Latent Variables and Complete Likelihood (1/2)

- Problem definition: given a dataset  $\mathcal{X} = \{\mathbf{x}^{(t)}\}_{t=1}^N$ , suppose  $p(\mathbf{x}) = \sum_{i=1}^K p(\mathbf{x}|G_i)P[G_i]$  and denote  $\Theta = (\theta_i, \pi_i)_{i=1}^K$  where  $\theta_i$  parametrizes  $p(\mathbf{x}|G_i)$  and  $\pi_i = P[G_i]$ , we want to find  $\Theta$  such that the log likelihood  $\ln P[\mathcal{X}|\Theta]$  is maximized
  - $\ln P[\mathcal{X}|\Theta] = \sum_{t=1}^N \ln \sum_{i=1}^K p(\mathbf{x}^{(t)}|\theta_i)\pi_i d\mathbf{x}$
  - Unfortunately, since we don't know which instance belongs to which group, we cannot solve this objective analytically
- Now suppose there is a set  $\mathcal{Z} = \{\mathbf{z}^{(t)}\}_{t=1}^N$  of **latent variables**, the **complete likelihood** can be written as:  $\ln P[\mathcal{X}, \mathcal{Z}|\Theta]$ 
  - $z_i^{(t)} = 1$  if  $\mathbf{x}^{(t)}$  belongs to group  $i$ ; 0 otherwise
  - If we have  $\mathcal{Z}$ , we can solve this objective as we did in the parametric classification
  - Unfortunately, we don't know  $\mathcal{Z}$
  - So let's create it and maximize  $E_{\mathcal{Z}}[\ln P[\mathcal{X}|\Theta]]$

# Latent Variables and Complete Likelihood (2/2)

- Observe that
$$\ln P[\mathcal{X}, \mathcal{Z} | \Theta] = \ln P[\mathcal{Z} | \mathcal{X}, \Theta] + \ln P[\mathcal{X} | \Theta]$$
- We have
$$E_{\mathcal{Z}}[\ln P[\mathcal{X} | \Theta]] = L(q, \Theta) + KL(q || P) \text{ for any distribution } q \text{ of } \mathcal{Z}$$
  - $L(q, \Theta) = \sum_{\mathcal{Z}} q(\mathcal{Z}) \ln \left( \frac{P[\mathcal{X}, \mathcal{Z} | \Theta]}{q(\mathcal{Z})} \right)$
  - $KL(q || P) = - \sum_{\mathcal{Z}} q(\mathcal{Z}) \ln \left( \frac{P[\mathcal{Z} | \mathcal{X}, \Theta]}{q(\mathcal{Z})} \right)$
- Both  $L(q, \Theta)$  and  $KL(q || P)$  are **functional** of  $q$
- Since  $KL(q || P)$  is the **relative entropy** (or **Kullback-Leibler divergence**) and is always greater than 0 [Proof: by Jensen's inequality or  $\ln x \leq x - 1$ ], we have the figure at right:



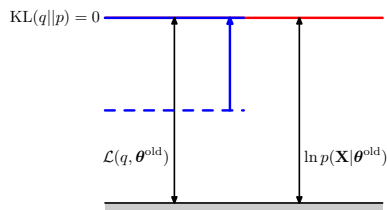


# Iterative Methods: A Functional Perspective (1/3)

- To maximize  $E_{\mathcal{Z}}[\ln P[\mathcal{X}|\Theta]]$ , we can employ an iterative method based on  $E_{\mathcal{Z}}[\ln P[\mathcal{X}|\Theta]] = L(q, \Theta) + KL(q||P)$ 
  - Since  $q$  is unknown, we make up  $q$
  - We don't have to make up  $\mathcal{Z}$  this time because we try out all possible  $\mathcal{Z}$  in  $L(q, \Theta)$  and  $KL(q||P)$
- Start from a random guess about  $\Theta$ , iterate the following steps:
  - 1 Update  $q$  based on current  $\Theta$  such that the blue line is up-aligned with the red
  - 2 Update  $\Theta$  based on current  $q$  to raises the red line
- Stop until  $\Theta$  converges
- Why another version?
  - We are sure that  $E_{\mathcal{Z}}[\ln P[\mathcal{X}|\Theta]]$  (i.e., red line) can be raised at each iteration (although up to a local optimal)

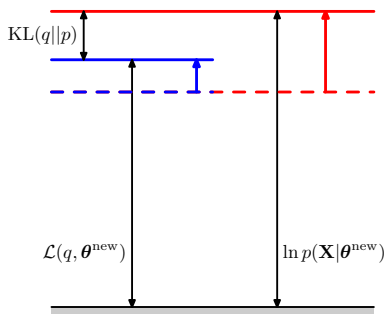
# Iterative Methods: A Functional Perspective (2/3)

- Denote by  $\Theta^{old}$  the parameters found in the previous iteration
- In Step 1, if we update  $q$  such that  $q(\mathcal{Z}) = P[\mathcal{Z}|\mathcal{X}, \Theta^{old}]$ 
  - $KL(q||P) = -\sum_{\mathcal{Z}} q(\mathcal{Z}) \ln 1 = 0$
  - $\ln P[\mathcal{X}|\Theta^{old}] = L(q, \Theta^{old}) + 0$
- Note the value of  $E_{\mathcal{Z}}[\ln P[\mathcal{X}|\Theta^{old}]]$  won't change as we vary  $q$
- So this step basically raises  $L(q, \Theta^{old})$  such that the blue line is up-aligned with the red



# Iterative Methods: A Functional Perspective (3/3)

- Fixing  $q$ , for any  $\Theta$  we have
$$E_{\mathcal{Z}}[\ln P[\mathcal{X}|\Theta]] = L(q, \Theta) + KL(q||P) = \sum_{\mathcal{Z}} P[\mathcal{Z}|\mathcal{X}, \Theta^{old}] \ln \left( \frac{P[\mathcal{X}, \mathcal{Z}|\Theta]}{P[\mathcal{Z}|\mathcal{X}, \Theta^{old}]} \right) - \sum_{\mathcal{Z}} P[\mathcal{Z}|\mathcal{X}, \Theta^{old}] \ln \left( \frac{P[\mathcal{Z}|\mathcal{X}, \Theta]}{P[\mathcal{Z}|\mathcal{X}, \Theta^{old}]} \right)$$
- In Step 2, we find  $\Theta^{new}$  maximizing  $L(q, \Theta)$  (blue line)
  - $KL(q||P) = -\sum_{\mathcal{Z}} P[\mathcal{Z}|\mathcal{X}, \Theta^{old}] \ln \left( \frac{P[\mathcal{Z}|\mathcal{X}, \Theta^{new}]}{P[\mathcal{Z}|\mathcal{X}, \Theta^{old}]} \right) \geq 0$
  - $E_{\mathcal{Z}}[\ln P[\mathcal{X}|\Theta^{new}]] \geq E_{\mathcal{Z}}[\ln P[\mathcal{X}|\Theta^{old}]]$
- So this step basically raises the red line, meanwhile leaving the blue behind
- Repeating Steps 1 and 2 lifts  $E_{\mathcal{Z}}[\ln P[\mathcal{X}|\Theta]]$  till some local optimum



## 1 Clustering

- Why Clustering?
- $k$ -Means Clustering
- Semiparametric Density Estimation
- Hierarchical Clustering
- Spectral Clustering
- Practical Considerations

## 2 Expectation Maximization

- Latent Variables and Complete Likelihood
- EM Steps
- EM for Mixture Models
- EM for Mixtures of Gaussians

# Expectation Maximization (1/2)

- Note that in step 1, we don't need to write down  $q$  explicitly
  - We just need to evaluate terms in  $E_{\mathcal{Z}}[\ln P[\mathcal{X}|\Theta]]$  (a function of  $\Theta$  to be maximized in step 2) that are related to  $q$

- Fixing  $q(\mathcal{Z}) = P[\mathcal{Z}|\mathcal{X}, \Theta^{old}]$ , we have

$$\begin{aligned}L(q, \Theta) &= \sum_{\mathcal{Z}} P[\mathcal{Z}|\mathcal{X}, \Theta^{old}] \ln \left( \frac{P[\mathcal{X}, \mathcal{Z}|\Theta]}{P[\mathcal{Z}|\mathcal{X}, \Theta^{old}]} \right) \\&= \sum_{\mathcal{Z}} P[\mathcal{Z}|\mathcal{X}, \Theta^{old}] \ln (P[\mathcal{X}, \mathcal{Z}|\Theta]) - \sum_{\mathcal{Z}} P[\mathcal{Z}|\mathcal{X}, \Theta^{old}] \ln (P[\mathcal{Z}|\mathcal{X}, \Theta^{old}]) \\&= E_{\mathcal{Z}}[\ln (P[\mathcal{X}, \mathcal{Z}|\Theta]) | \mathcal{X}, \Theta^{old}] + \text{constant}\end{aligned}$$

- The second term is an entropy of  $\mathcal{Z}$  and is independent of  $\Theta$
- Steps 1 and 2 can be refined to be:
  - **Expectation step (E-step)**: Formulate  $Q(\Theta; \Theta^{old}) = E_{\mathcal{Z}}[\ln (P[\mathcal{X}, \mathcal{Z}|\Theta]) | \mathcal{X}, \Theta^{old}]$  and evaluate the terms related to  $P[\mathcal{Z}|\mathcal{X}, \Theta^{old}]$
  - **Maximization step (M-step)**: Solve  $\Theta^{new} = \arg_{\Theta} \max Q(\Theta; \Theta^{old})$

# Expectation Maximization (2/2)

- The **Expectation Maximization (EM)** algorithm is a general technique to find the maximum likelihood solutions for probabilistic models having latent variables
  - Typically, latent variables are discrete, and there is one latent variable per observed instance

**Input:**  $\mathcal{X} \leftarrow \{\mathbf{x}^{(t)}\}_{t=1}^N$

**Output:**  $\Theta^{new}$ , a local optimizer of  $E_{\mathcal{Z}}[P[\mathcal{X}|\Theta]]$

Choose an initial  $\Theta^{new}$ ;

**repeat**

$\Theta^{old} \leftarrow \Theta^{new}$ ;

    Formulate  $\mathcal{Q}(\Theta; \Theta^{old}) = E_{\mathcal{Z}}[\ln(P[\mathcal{X}, \mathcal{Z}|\Theta])|\mathcal{X}, \Theta^{old}]$  and evaluate the terms related to  $P[\mathcal{Z}|\mathcal{X}, \Theta^{old}]$ ; // E-step

$\Theta^{new} \leftarrow \arg_{\Theta} \max \mathcal{Q}(\Theta; \Theta^{old})^{new}$ ; // M-step

**until**  $\Theta^{new}$  converges;

**Algorithm 4:** The general EM algorithm.

## 1 Clustering

- Why Clustering?
- $k$ -Means Clustering
- Semiparametric Density Estimation
- Hierarchical Clustering
- Spectral Clustering
- Practical Considerations

## 2 Expectation Maximization

- Latent Variables and Complete Likelihood
- EM Steps
- EM for Mixture Models
- EM for Mixtures of Gaussians

# I.I.D. Instances

- Assume that  $(\mathbf{x}^{(t)}, \mathbf{z}^{(t)})$  are i.i.d. samples drawn from some distribution
- By definition,  $P[\mathcal{X}, \mathcal{Z}] = \prod_{t=1}^N P[\mathbf{x}^{(t)}, \mathbf{z}^{(t)}]$
- Denote  $\mathbf{e}_1 = [1, 0 \cdots, 0]^\top$ ,  $\mathbf{e}_2 = [0, 1 \cdots, 0]^\top, \dots, \mathbf{e}_K = [0, 0 \cdots, 1]^\top \in \mathbb{R}^K$
- We have  $P[\mathcal{X}] = \sum_{\mathcal{Z}} P[\mathcal{X}, \mathcal{Z}] = \sum_{\mathcal{Z}} \prod_{t=1}^N P[\mathbf{x}^{(t)}, \mathbf{z}^{(t)}] = \sum_{\mathbf{z}^{(1)}=\mathbf{e}_1}^{\mathbf{e}_K} \cdots \sum_{\mathbf{z}^{(N)}=\mathbf{e}_1}^{\mathbf{e}_K} \prod_{t=1}^N P[\mathbf{x}^{(t)}, \mathbf{z}^{(t)}] = \prod_{t=1}^N \sum_{\mathbf{z}^{(t)}=\mathbf{e}_1}^{\mathbf{e}_K} P[\mathbf{x}^{(t)}, \mathbf{z}^{(t)}] = \prod_{t=1}^N P[\mathbf{x}^{(t)}]$
- So,

$$P[\mathcal{Z}|\mathcal{X}, \Theta^{old}] = \frac{P[\mathcal{X}, \mathcal{Z}|\Theta^{old}]}{P[\mathcal{X}|\Theta^{old}]} = \frac{\prod_{t=1}^N P[\mathbf{x}^{(t)}, \mathbf{z}^{(t)}|\Theta^{old}]}{\prod_{t=1}^N P[\mathbf{x}^{(t)}|\Theta^{old}]} = \prod_{t=1}^N P[\mathbf{z}^{(t)}|\mathbf{x}^{(t)}, \Theta^{old}]$$



# Formulating $\mathcal{Q}(\Theta; \Theta^{old})$ (1/3)

- Next, we formulate  $\mathcal{Q}(\Theta; \Theta^{old})$  given  $P[\mathcal{Z}|\mathcal{X}, \Theta^{old}] = \prod_{t=1}^N P[\mathbf{z}^{(t)}|\mathbf{x}^{(t)}, \Theta^{old}]$  (due to the i.i.d. instances), and the assumption of mixture density:
  - $\Theta = (\theta_i, \pi_i)_{i=1}^K$  where  $\theta_i$  parametrizes  $p(\mathbf{x}|G_i)$  and  $\pi_i = P[G_i]$
- Denote by  $d(\mathbf{z}^{(t)})$  the index of attribute of  $\mathbf{z}^{(t)}$  equal to 1
  - $P[\mathbf{x}^{(t)}, \mathbf{z}^{(t)}|\Theta] = P[\mathbf{x}^{(t)}|\mathbf{z}^{(t)}, \Theta]P[\mathbf{z}^{(t)}|\Theta] = P[\mathbf{x}^{(t)}|\mathbf{z}^{(t)}, \theta_{d(\mathbf{z}^{(t)})}]\pi_{d(\mathbf{z}^{(t)})}$
- For brevity, we use the shorthand  $P[z_i^{(t)}]$  for  $P[\mathbf{z}^{(t)} = \mathbf{e}_i]$  (or equivalently  $P[z_i^{(t)} = 1]$ )

## Formulating $Q(\Theta; \Theta^{old})$ (2/3)

$$Q(\Theta; \Theta^{old}) = E_Z[\ln(P[\mathcal{X}, Z|\Theta]) | \mathcal{X}, \Theta^{old}] = \sum_Z \ln(P[\mathcal{X}, Z|\Theta]) P[Z|\mathcal{X}, \Theta^{old}]$$

## Formulating $Q(\Theta; \Theta^{old})$ (2/3)

$$\begin{aligned} Q(\Theta; \Theta^{old}) &= E_{\mathcal{Z}}[\ln(P[\mathcal{X}, \mathcal{Z}|\Theta]) | \mathcal{X}, \Theta^{old}] = \sum_{\mathcal{Z}} \ln(P[\mathcal{X}, \mathcal{Z}|\Theta]) P[\mathcal{Z}|\mathcal{X}, \Theta^{old}] \\ &= \sum_{\mathcal{Z}} \sum_{t=1}^N \ln(P[\mathbf{x}^{(t)}, \mathbf{z}^{(t)}|\Theta]) \prod_{j=1}^N P[\mathbf{z}^{(j)}|\mathbf{x}^{(j)}, \Theta^{old}] \end{aligned}$$

## Formulating $\mathcal{Q}(\Theta; \Theta^{old})$ (2/3)

$$\begin{aligned}\mathcal{Q}(\Theta; \Theta^{old}) &= E_{\mathcal{Z}}[\ln(P[\mathcal{X}, \mathcal{Z}|\Theta]) | \mathcal{X}, \Theta^{old}] = \sum_{\mathcal{Z}} \ln(P[\mathcal{X}, \mathcal{Z}|\Theta]) P[\mathcal{Z}|\mathcal{X}, \Theta^{old}] \\ &= \sum_{\mathcal{Z}} \sum_{t=1}^N \ln(P[\mathbf{x}^{(t)}, \mathbf{z}^{(t)}|\Theta]) \prod_{j=1}^N P[\mathbf{z}^{(j)}|\mathbf{x}^{(j)}, \Theta^{old}] \\ &= \sum_{\mathcal{Z}} \sum_{t=1}^N \ln\left(P[\mathbf{x}^{(t)}|\mathbf{z}^{(t)}, \theta_{d(\mathbf{z}^{(t)})}] \pi_{d(\mathbf{z}^{(t)})}\right) \prod_{j=1}^N P[\mathbf{z}^{(j)}|\mathbf{x}^{(j)}, \Theta^{old}]\end{aligned}$$

## Formulating $Q(\Theta; \Theta^{old})$ (2/3)

$$\begin{aligned}Q(\Theta; \Theta^{old}) &= E_{\mathcal{Z}}[\ln(P[\mathcal{X}, \mathcal{Z}|\Theta]) | \mathcal{X}, \Theta^{old}] = \sum_{\mathcal{Z}} \ln(P[\mathcal{X}, \mathcal{Z}|\Theta]) P[\mathcal{Z}|\mathcal{X}, \Theta^{old}] \\&= \sum_{\mathcal{Z}} \sum_{t=1}^N \ln(P[\mathbf{x}^{(t)}, \mathbf{z}^{(t)}|\Theta]) \prod_{j=1}^N P[\mathbf{z}^{(j)}|\mathbf{x}^{(j)}, \Theta^{old}] \\&= \sum_{\mathcal{Z}} \sum_{t=1}^N \ln\left(P[\mathbf{x}^{(t)}|\mathbf{z}^{(t)}, \theta_{d(\mathbf{z}^{(t)})}] \pi_{d(\mathbf{z}^{(t)})}\right) \prod_{j=1}^N P[\mathbf{z}^{(j)}|\mathbf{x}^{(j)}, \Theta^{old}] \\&= \sum_{\mathcal{Z}} \sum_{t=1}^N \sum_{e=e_1}^{e_K} \delta_{\mathbf{z}^{(t)}, e} \ln(P[\mathbf{x}^{(t)}|e, \theta_{d(e)}] \pi_{d(e)}) \\&\quad \prod_{j=1}^N P[\mathbf{z}^{(j)}|\mathbf{x}^{(j)}, \Theta^{old}] \quad // \delta_{a,b} = 1 \text{ if } a = b; 0 \text{ otherwise}\end{aligned}$$

## Formulating $\mathcal{Q}(\Theta; \Theta^{old})$ (2/3)

$$\begin{aligned}\mathcal{Q}(\Theta; \Theta^{old}) &= E_{\mathcal{Z}}[\ln(P[\mathcal{X}, \mathcal{Z}|\Theta]) | \mathcal{X}, \Theta^{old}] = \sum_{\mathcal{Z}} \ln(P[\mathcal{X}, \mathcal{Z}|\Theta]) P[\mathcal{Z}|\mathcal{X}, \Theta^{old}] \\ &= \sum_{\mathcal{Z}} \sum_{t=1}^N \ln(P[\mathbf{x}^{(t)}, \mathbf{z}^{(t)}|\Theta]) \prod_{j=1}^N P[\mathbf{z}^{(j)}|\mathbf{x}^{(j)}, \Theta^{old}] \\ &= \sum_{\mathcal{Z}} \sum_{t=1}^N \ln\left(P[\mathbf{x}^{(t)}|\mathbf{z}^{(t)}, \theta_{d(\mathbf{z}^{(t)})}] \pi_{d(\mathbf{z}^{(t)})}\right) \prod_{j=1}^N P[\mathbf{z}^{(j)}|\mathbf{x}^{(j)}, \Theta^{old}] \\ &= \sum_{\mathcal{Z}} \sum_{t=1}^N \sum_{e=e_1}^{e_K} \delta_{\mathbf{z}^{(t)}, e} \ln(P[\mathbf{x}^{(t)}|e, \theta_{d(e)}] \pi_{d(e)}) \\ &\quad \prod_{j=1}^N P[\mathbf{z}^{(j)}|\mathbf{x}^{(j)}, \Theta^{old}] \quad // \delta_{a,b} = 1 \text{ if } a = b; 0 \text{ otherwise} \\ &= \sum_{t=1}^N \sum_{e=e_1}^{e_K} \ln(P[\mathbf{x}^{(t)}|e, \theta_{d(e)}] \pi_{d(e)}) \sum_{\mathcal{Z}} \delta_{\mathbf{z}^{(t)}, e} \prod_{j=1}^N P[\mathbf{z}^{(j)}|\mathbf{x}^{(j)}, \Theta^{old}]\end{aligned}$$

## Formulating $\mathcal{Q}(\Theta; \Theta^{old})$ (2/3)

$$\begin{aligned}\mathcal{Q}(\Theta; \Theta^{old}) &= E_{\mathcal{Z}}[\ln(P[\mathcal{X}, \mathcal{Z}|\Theta]) | \mathcal{X}, \Theta^{old}] = \sum_{\mathcal{Z}} \ln(P[\mathcal{X}, \mathcal{Z}|\Theta]) P[\mathcal{Z}|\mathcal{X}, \Theta^{old}] \\ &= \sum_{\mathcal{Z}} \sum_{t=1}^N \ln(P[\mathbf{x}^{(t)}, \mathbf{z}^{(t)}|\Theta]) \prod_{j=1}^N P[\mathbf{z}^{(j)}|\mathbf{x}^{(j)}, \Theta^{old}] \\ &= \sum_{\mathcal{Z}} \sum_{t=1}^N \ln\left(P[\mathbf{x}^{(t)}|\mathbf{z}^{(t)}, \theta_{d(\mathbf{z}^{(t)})}] \pi_{d(\mathbf{z}^{(t)})}\right) \prod_{j=1}^N P[\mathbf{z}^{(j)}|\mathbf{x}^{(j)}, \Theta^{old}] \\ &= \sum_{\mathcal{Z}} \sum_{t=1}^N \sum_{e=e_1}^{e_K} \delta_{\mathbf{z}^{(t)}, e} \ln(P[\mathbf{x}^{(t)}|e, \theta_{d(e)}] \pi_{d(e)}) \\ &\quad \prod_{j=1}^N P[\mathbf{z}^{(j)}|\mathbf{x}^{(j)}, \Theta^{old}] \quad // \delta_{a,b} = 1 \text{ if } a = b; 0 \text{ otherwise} \\ &= \sum_{t=1}^N \sum_{e=e_1}^{e_K} \ln(P[\mathbf{x}^{(t)}|e, \theta_{d(e)}] \pi_{d(e)}) \sum_{\mathcal{Z}} \delta_{\mathbf{z}^{(t)}, e} \prod_{j=1}^N P[\mathbf{z}^{(j)}|\mathbf{x}^{(j)}, \Theta^{old}] \\ &= \sum_{t=1}^N \sum_{e=e_1}^{e_K} \ln(P[\mathbf{x}^{(t)}|e, \theta_{d(e)}] \pi_{d(e)}) \\ &\quad \sum_{\mathbf{z}^{(1)}=e_1}^{e_K} \cdots \sum_{\mathbf{z}^{(N)}=e_1}^{e_K} \delta_{\mathbf{z}^{(t)}, e} \prod_{j=1}^N P[\mathbf{z}^{(j)}|\mathbf{x}^{(j)}, \Theta^{old}]\end{aligned}$$

# Formulating $\mathcal{Q}(\Theta; \Theta^{old})$ (2/3)

$$\begin{aligned}\mathcal{Q}(\Theta; \Theta^{old}) &= E_{\mathcal{Z}}[\ln(P[\mathcal{X}, \mathcal{Z}|\Theta]) | \mathcal{X}, \Theta^{old}] = \sum_{\mathcal{Z}} \ln(P[\mathcal{X}, \mathcal{Z}|\Theta]) P[\mathcal{Z}|\mathcal{X}, \Theta^{old}] \\ &= \sum_{\mathcal{Z}} \sum_{t=1}^N \ln(P[\mathbf{x}^{(t)}, \mathbf{z}^{(t)}|\Theta]) \prod_{j=1}^N P[\mathbf{z}^{(j)}|\mathbf{x}^{(j)}, \Theta^{old}] \\ &= \sum_{\mathcal{Z}} \sum_{t=1}^N \ln\left(P[\mathbf{x}^{(t)}|\mathbf{z}^{(t)}, \theta_{d(\mathbf{z}^{(t)})}] \pi_{d(\mathbf{z}^{(t)})}\right) \prod_{j=1}^N P[\mathbf{z}^{(j)}|\mathbf{x}^{(j)}, \Theta^{old}] \\ &= \sum_{\mathcal{Z}} \sum_{t=1}^N \sum_{e=e_1}^{e_K} \delta_{\mathbf{z}^{(t)}, e} \ln(P[\mathbf{x}^{(t)}|e, \theta_{d(e)}] \pi_{d(e)}) \\ &\quad \prod_{j=1}^N P[\mathbf{z}^{(j)}|\mathbf{x}^{(j)}, \Theta^{old}] \quad // \delta_{a,b} = 1 \text{ if } a = b; 0 \text{ otherwise} \\ &= \sum_{t=1}^N \sum_{e=e_1}^{e_K} \ln(P[\mathbf{x}^{(t)}|e, \theta_{d(e)}] \pi_{d(e)}) \sum_{\mathcal{Z}} \delta_{\mathbf{z}^{(t)}, e} \prod_{j=1}^N P[\mathbf{z}^{(j)}|\mathbf{x}^{(j)}, \Theta^{old}] \\ &= \sum_{t=1}^N \sum_{e=e_1}^{e_K} \ln(P[\mathbf{x}^{(t)}|e, \theta_{d(e)}] \pi_{d(e)}) \\ &\quad \sum_{\mathbf{z}^{(1)}=e_1}^{e_K} \cdots \sum_{\mathbf{z}^{(N)}=e_1}^{e_K} \delta_{\mathbf{z}^{(t)}, e} \prod_{j=1}^N P[\mathbf{z}^{(j)}|\mathbf{x}^{(j)}, \Theta^{old}] \\ &= \sum_{t=1}^N \sum_{e=e_1}^{e_K} \ln(P[\mathbf{x}^{(t)}|e, \theta_{d(e)}] \pi_{d(e)}) \sum_{\mathbf{z}^{(1)}=e_1}^{e_K} \cdots \sum_{\mathbf{z}^{(t-1)}=e_1}^{e_K} \\ &\quad \sum_{\mathbf{z}^{(t+1)}=e_1}^{e_K} \cdots \sum_{\mathbf{z}^{(N)}=e_1}^{e_K} \prod_{j=1, j \neq t}^N P[\mathbf{z}^{(j)}|\mathbf{x}^{(j)}, \Theta^{old}] P[e|\mathbf{x}^{(t)}, \Theta^{old}]\end{aligned}$$



# Formulating $\mathcal{Q}(\Theta; \Theta^{old})$ (2/3)

$$\begin{aligned}
 \mathcal{Q}(\Theta; \Theta^{old}) &= E_{\mathcal{Z}}[\ln(P[\mathcal{X}, \mathcal{Z}|\Theta]) | \mathcal{X}, \Theta^{old}] = \sum_{\mathcal{Z}} \ln(P[\mathcal{X}, \mathcal{Z}|\Theta]) P[\mathcal{Z}|\mathcal{X}, \Theta^{old}] \\
 &= \sum_{\mathcal{Z}} \sum_{t=1}^N \ln(P[\mathbf{x}^{(t)}, \mathbf{z}^{(t)}|\Theta]) \prod_{j=1}^N P[\mathbf{z}^{(j)}|\mathbf{x}^{(j)}, \Theta^{old}] \\
 &= \sum_{\mathcal{Z}} \sum_{t=1}^N \ln\left(P[\mathbf{x}^{(t)}|\mathbf{z}^{(t)}, \theta_{d(\mathbf{z}^{(t)})}] \pi_{d(\mathbf{z}^{(t)})}\right) \prod_{j=1}^N P[\mathbf{z}^{(j)}|\mathbf{x}^{(j)}, \Theta^{old}] \\
 &= \sum_{\mathcal{Z}} \sum_{t=1}^N \sum_{e=e_1}^{e_K} \delta_{\mathbf{z}^{(t)}, e} \ln(P[\mathbf{x}^{(t)}|e, \theta_{d(e)}] \pi_{d(e)}) \\
 &\quad \prod_{j=1}^N P[\mathbf{z}^{(j)}|\mathbf{x}^{(j)}, \Theta^{old}] \quad // \delta_{a,b} = 1 \text{ if } a = b; 0 \text{ otherwise} \\
 &= \sum_{t=1}^N \sum_{e=e_1}^{e_K} \ln(P[\mathbf{x}^{(t)}|e, \theta_{d(e)}] \pi_{d(e)}) \sum_{\mathcal{Z}} \delta_{\mathbf{z}^{(t)}, e} \prod_{j=1}^N P[\mathbf{z}^{(j)}|\mathbf{x}^{(j)}, \Theta^{old}] \\
 &= \sum_{t=1}^N \sum_{e=e_1}^{e_K} \ln(P[\mathbf{x}^{(t)}|e, \theta_{d(e)}] \pi_{d(e)}) \\
 &\quad \sum_{\mathbf{z}^{(1)}=e_1}^{e_K} \cdots \sum_{\mathbf{z}^{(N)}=e_1}^{e_K} \delta_{\mathbf{z}^{(t)}, e} \prod_{j=1}^N P[\mathbf{z}^{(j)}|\mathbf{x}^{(j)}, \Theta^{old}] \\
 &= \sum_{t=1}^N \sum_{e=e_1}^{e_K} \ln(P[\mathbf{x}^{(t)}|e, \theta_{d(e)}] \pi_{d(e)}) \sum_{\mathbf{z}^{(1)}=e_1}^{e_K} \cdots \sum_{\mathbf{z}^{(t-1)}=e_1}^{e_K} \\
 &\quad \sum_{\mathbf{z}^{(t+1)}=e_1}^{e_K} \cdots \sum_{\mathbf{z}^{(N)}=e_1}^{e_K} \prod_{j=1, j \neq t}^N P[\mathbf{z}^{(j)}|\mathbf{x}^{(j)}, \Theta^{old}] P[e|\mathbf{x}^{(t)}, \Theta^{old}] \\
 &= \sum_{t=1}^N \sum_{e=e_1}^{e_K} \ln(P[\mathbf{x}^{(t)}|e, \theta_{d(e)}] \pi_{d(e)}) \left( \sum_{\mathbf{z}^{(1)}=e_1}^{e_K} \cdots \sum_{\mathbf{z}^{(t-1)}=e_1}^{e_K} \right. \\
 &\quad \left. \sum_{\mathbf{z}^{(t+1)}=e_1}^{e_K} \cdots \sum_{\mathbf{z}^{(N)}=e_1}^{e_K} \prod_{j=1, j \neq t}^N P[\mathbf{z}^{(j)}|\mathbf{x}^{(j)}, \Theta^{old}] \right) P[e|\mathbf{x}^{(t)}, \Theta^{old}]
 \end{aligned}$$

## Formulating $\mathcal{Q}(\Theta; \Theta^{old})$ (3/3)

$$= \sum_{t=1}^N \sum_{\mathbf{e}=e_1}^{e_K} \ln (P[\mathbf{x}^{(t)}|\mathbf{e}, \theta_{d(\mathbf{e})}]\pi_{d(\mathbf{e})}) \left( \sum_{z^{(1)}=e_1}^{e_K} \cdots \sum_{z^{(t-1)}=e_1}^{e_K} \sum_{z^{(t+1)}=e_1}^{e_K} \cdots \sum_{z^{(N)}=e_1}^{e_K} \prod_{j=1, j \neq t}^N P[z^{(j)}|\mathbf{x}^{(j)}, \Theta^{old}] \right) P[\mathbf{e}|\mathbf{x}^{(t)}, \Theta^{old}]$$

# Formulating $\mathcal{Q}(\Theta; \Theta^{old})$ (3/3)

$$\begin{aligned} &= \sum_{t=1}^N \sum_{\mathbf{e}=\mathbf{e}_1}^{\mathbf{e}_K} \ln (P[\mathbf{x}^{(t)}|\mathbf{e}, \theta_{d(\mathbf{e})}]\pi_{d(\mathbf{e})}) \left( \sum_{z^{(1)}=\mathbf{e}_1}^{\mathbf{e}_K} \cdots \sum_{z^{(t-1)}=\mathbf{e}_1}^{\mathbf{e}_K} \right. \\ &\quad \left. \sum_{z^{(t+1)}=\mathbf{e}_1}^{\mathbf{e}_K} \cdots \sum_{z^{(N)}=\mathbf{e}_1}^{\mathbf{e}_K} \prod_{j=1, j \neq t}^N P[z^{(j)}|\mathbf{x}^{(j)}, \Theta^{old}] \right) P[\mathbf{e}|\mathbf{x}^{(t)}, \Theta^{old}] \\ &= \sum_{t=1}^N \sum_{\mathbf{e}=\mathbf{e}_1}^{\mathbf{e}_K} \ln (P[\mathbf{x}^{(t)}|\mathbf{e}, \theta_{d(\mathbf{e})}]\pi_{d(\mathbf{e})}) \\ &\quad \left( \prod_{j=1, j \neq t}^N \sum_{z^{(j)}=\mathbf{e}_1}^{\mathbf{e}_K} P[z^{(j)}|\mathbf{x}^{(j)}, \Theta^{old}] \right) P[\mathbf{e}|\mathbf{x}^{(t)}, \Theta^{old}] \end{aligned}$$

# Formulating $\mathcal{Q}(\Theta; \Theta^{old})$ (3/3)

$$\begin{aligned} &= \sum_{t=1}^N \sum_{\mathbf{e}=\mathbf{e}_1}^{\mathbf{e}_K} \ln (P[\mathbf{x}^{(t)}|\mathbf{e}, \theta_{d(\mathbf{e})}]\pi_{d(\mathbf{e})}) \left( \sum_{z^{(1)}=\mathbf{e}_1}^{\mathbf{e}_K} \cdots \sum_{z^{(t-1)}=\mathbf{e}_1}^{\mathbf{e}_K} \right. \\ &\quad \left. \sum_{z^{(t+1)}=\mathbf{e}_1}^{\mathbf{e}_K} \cdots \sum_{z^{(N)}=\mathbf{e}_1}^{\mathbf{e}_K} \prod_{j=1, j \neq t}^N P[z^{(j)}|\mathbf{x}^{(j)}, \Theta^{old}] \right) P[\mathbf{e}|\mathbf{x}^{(t)}, \Theta^{old}] \\ &= \sum_{t=1}^N \sum_{\mathbf{e}=\mathbf{e}_1}^{\mathbf{e}_K} \ln (P[\mathbf{x}^{(t)}|\mathbf{e}, \theta_{d(\mathbf{e})}]\pi_{d(\mathbf{e})}) \\ &\quad \left( \prod_{j=1, j \neq t}^N \sum_{z^{(j)}=\mathbf{e}_1}^{\mathbf{e}_K} P[z^{(j)}|\mathbf{x}^{(j)}, \Theta^{old}] \right) P[\mathbf{e}|\mathbf{x}^{(t)}, \Theta^{old}] \\ &= \sum_{t=1}^N \sum_{\mathbf{e}=\mathbf{e}_1}^{\mathbf{e}_K} \ln (P[\mathbf{x}^{(t)}|\mathbf{e}, \theta_{d(\mathbf{e})}]\pi_{d(\mathbf{e})}) P[\mathbf{e}|\mathbf{x}^{(t)}, \Theta^{old}] \end{aligned}$$

# Formulating $\mathcal{Q}(\Theta; \Theta^{old})$ (3/3)

$$\begin{aligned} &= \sum_{t=1}^N \sum_{\mathbf{e}=e_1}^{e_K} \ln (P[\mathbf{x}^{(t)}|\mathbf{e}, \theta_{d(\mathbf{e})}]\pi_{d(\mathbf{e})}) \left( \sum_{z^{(1)}=e_1}^{e_K} \cdots \sum_{z^{(t-1)}=e_1}^{e_K} \right. \\ &\quad \left. \sum_{z^{(t+1)}=e_1}^{e_K} \cdots \sum_{z^{(N)}=e_1}^{e_K} \prod_{j=1, j \neq t}^N P[z^{(j)}|\mathbf{x}^{(j)}, \Theta^{old}] \right) P[\mathbf{e}|\mathbf{x}^{(t)}, \Theta^{old}] \\ &= \sum_{t=1}^N \sum_{\mathbf{e}=e_1}^{e_K} \ln (P[\mathbf{x}^{(t)}|\mathbf{e}, \theta_{d(\mathbf{e})}]\pi_{d(\mathbf{e})}) \\ &\quad \left( \prod_{j=1, j \neq t}^N \sum_{z^{(j)}=e_1}^{e_K} P[z^{(j)}|\mathbf{x}^{(j)}, \Theta^{old}] \right) P[\mathbf{e}|\mathbf{x}^{(t)}, \Theta^{old}] \\ &= \sum_{t=1}^N \sum_{\mathbf{e}=e_1}^{e_K} \ln (P[\mathbf{x}^{(t)}|\mathbf{e}, \theta_{d(\mathbf{e})}]\pi_{d(\mathbf{e})}) P[\mathbf{e}|\mathbf{x}^{(t)}, \Theta^{old}] \\ &= \sum_{t=1}^N \sum_{i=1}^K \ln (\pi_i) P[z_i^{(t)}|\mathbf{x}^{(t)}, \Theta^{old}] + \\ &\quad \sum_{t=1}^N \sum_{i=1}^K \ln (P[\mathbf{x}^{(t)}|z_i^{(t)}, \theta_i]) P[z_i^{(t)}|\mathbf{x}^{(t)}, \Theta^{old}] \end{aligned}$$

# Evaluating $P[\mathcal{Z}|\mathcal{X}, \Theta^{old}]$

- Given mixtures of i.i.d. samples, we have
$$Q(\Theta; \Theta^{old}) = \sum_{t=1}^N \sum_{i=1}^K \ln(\pi_i) P[z_i^{(t)} | \mathbf{x}^{(t)}, \Theta^{old}] + \sum_{t=1}^N \sum_{i=1}^K \ln \left( P[\mathbf{x}^{(t)} | z_i^{(t)}, \theta_i] P[z_i^{(t)} | \mathbf{x}^{(t)}, \Theta^{old}] \right)$$
- The problem evaluating  $P[\mathcal{Z}|\mathcal{X}, \Theta^{old}]$  is thus reduced to evaluating  $P[z_i^{(t)} | \mathbf{x}^{(t)}, \Theta^{old}]$  for all  $1 \leq i \leq K$  and  $1 \leq t \leq N$

## 1 Clustering

- Why Clustering?
- $k$ -Means Clustering
- Semiparametric Density Estimation
- Hierarchical Clustering
- Spectral Clustering
- Practical Considerations

## 2 Expectation Maximization

- Latent Variables and Complete Likelihood
- EM Steps
- EM for Mixture Models
- EM for Mixtures of Gaussians

# Evaluating $P[z_i^{(t)} | \mathbf{x}^{(t)}, \Theta^{old}]$

- Problem: given  $\Theta^{old}$ , evaluate  $P[z_i^{(t)} | \mathbf{x}^{(t)}, \Theta^{old}]$  for all  $1 \leq i \leq K$  and  $1 \leq t \leq N$

- From Bayes' theorem,

$$P[z_i^{(t)} | \mathbf{x}^{(t)}, \Theta^{old}] = \frac{P[\mathbf{x}^{(t)} | z_i^{(t)}, \Theta^{old}] P[z_i^{(t)} | \Theta^{old}]}{\sum_{j=1}^K P[\mathbf{x}^{(t)} | z_j^{(t)}, \Theta^{old}] P[z_j^{(t)} | \Theta^{old}]} = \frac{P[\mathbf{x}^{(t)} | z_i^{(t)}, \theta_i^{old}] \pi_i^{old}}{\sum_{j=1}^K P[\mathbf{x}^{(t)} | z_j^{(t)}, \theta_j^{old}] \pi_j^{old}}$$

- If we further assume that instances in each group are normally distributed, then  $\theta_i^{old} = (\boldsymbol{\mu}_i^{old}, \boldsymbol{\Sigma}_i^{old})$  and we can easily obtain

$P[\mathbf{x}^{(t)} | z_i^{(t)}, \theta_i^{old}]$  based on the normal distribution

- For brevity, we denote the evaluated  $P[z_i^{(t)} | \mathbf{x}^{(t)}, \Theta^{old}]$  by  $h_i^{(t)}$ 
  - $h_i^{(t)}$  aligns with the soft label  $z_i^{(t)}$  in semiparametric density estimation



# Solving $\arg_{\Theta} \max \mathcal{Q}(\Theta; \Theta^{old})$

- $\mathcal{Q}(\Theta; \Theta^{old}) = \sum_{t=1}^N \sum_{i=1}^K \ln(\pi_i) h_i^{(t)} + \sum_{t=1}^N \sum_{i=1}^K \ln(P[\mathbf{x}^{(t)} | z_i^{(t)}, \theta_i]) h_i^{(t)}$
- Observe that the first term of  $\mathcal{Q}(\Theta; \Theta^{old})$  depends only on  $\{\pi_i\}_{i=1}^K$ ; and the second depends only on  $\{\theta_i\}_{i=1}^K$
- We can obtain  $\Theta^{new}$  by solving the two problems individually:
  - $\arg_{\pi_1, \dots, \pi_K} \max \sum_{t=1}^N \sum_{i=1}^K \ln(\pi_i) h_i^{(t)}$  subject to  $\sum_{i=1}^K \pi_i = 1$
  - $\arg_{\theta_1, \dots, \theta_K} \max \sum_{t=1}^N \sum_{i=1}^K \ln(P[\mathbf{x}^{(t)} | z_i^{(t)}, \theta_i]) h_i^{(t)}$

# Solving $\{\pi\}_{i=1}^K$

- Lagrangian:

$$L(\{\mathbf{x}^{(t)}\}_{t=1}^N, \{\pi_i\}_{i=1}^K, \alpha) = \sum_{t=1}^N \sum_{i=1}^K \ln(\pi_i) h_i^{(t)} - \alpha \left( \sum_{i=1}^K \pi_i - 1 \right)$$

- Taking the partial derivatives of  $L$  with respect to  $\alpha, \pi_1, \dots, \pi_K$  and setting them to zero we have  $\sum_{i=1}^K \pi_i = 1$  and

$$\sum_{t=1}^N \frac{1}{\pi_i} h_i^{(t)} - \alpha = 0 \Rightarrow \sum_{t=1}^N h_i^{(t)} = \pi_i \alpha \text{ for } i = 1, \dots, K$$

- Summing the equations with  $\alpha$  above we have  $\sum_{i=1}^K \sum_{t=1}^N h_i^{(t)} = \sum_{i=1}^K \pi_i \alpha \Rightarrow \alpha = \frac{\sum_{t=1}^N \sum_{i=1}^K h_i^{(t)}}{\sum_{i=1}^K \pi_i} = \sum_{t=1}^N \sum_{i=1}^K P[z_i^{(t)} | \mathbf{x}^{(t)}, \Theta^{old}] = N$

- Substituting  $N$  for  $\alpha$  in each of the above equation we have

$$\pi_i = \frac{\sum_{t=1}^N h_i^{(t)}}{N}$$

- This aligns with the  $\pi_i$  in semiparametric density estimation

# Solving $\{\theta_i\}_{i=1}^K$ (1/2)

- Objective:  $\arg_{\theta_1, \dots, \theta_K} \max \sum_{t=1}^N \sum_{i=1}^K \ln \left( P[\mathbf{x}^{(t)} | z_i^{(t)}, \theta_i] \right) h_i^{(t)}$
- Since the groups in the mixtures are independent with each other, we can solve  $\theta_i = (\boldsymbol{\mu}_i, \boldsymbol{\Sigma}_i)$  one by one
  - $\arg_{\theta_i} \max \sum_{t=1}^N \ln \left( P[\mathbf{x}^{(t)} | z_i^{(t)}, \theta_i] \right) h_i^{(t)}$
- With the Gaussian mixture, we have
$$\sum_{t=1}^N \ln \left( P[\mathbf{x}^{(t)} | z_i^{(t)}, \theta_i] \right) h_i^{(t)} = -\frac{N_i d}{2} \log(2\pi) - \frac{N_i}{2} \log(\det(\boldsymbol{\Sigma}_i)) - \frac{1}{2} \sum_{t=1}^N h_i^{(t)} (\mathbf{x}^{(t)} - \boldsymbol{\mu}_i)^\top \boldsymbol{\Sigma}_i^{-1} (\mathbf{x}^{(t)} - \boldsymbol{\mu}_i) = -\frac{N_i d}{2} \log(2\pi) + \frac{N_i}{2} \log(\det(\boldsymbol{\Sigma}_i^{-1})) - \frac{1}{2} \sum_{t=1}^N h_i^{(t)} \text{tr}(\boldsymbol{\Sigma}_i^{-1} (\mathbf{x}^{(t)} - \boldsymbol{\mu}_i)(\mathbf{x}^{(t)} - \boldsymbol{\mu}_i)^\top),$$
 where  $N_i = \sum_{t=1}^N h_i^{(t)}$

## Solving $\{\theta_i\}_{i=1}^K$ (2/2)

- Taking the partial derivatives of the above objective with respect to  $\mu_i$  and  $\Sigma_i^{-1}$  and setting them to zero we have

$$\begin{cases} \sum_{t=1}^N h_i^{(t)} (\mathbf{x}^{(t)} - \mu_i)^\top \Sigma_i^{-1} = \mathbf{0}^\top \\ \frac{N_i}{2} \Sigma_i - \frac{1}{2} \sum_{t=1}^N h_i^{(t)} (\mathbf{x}^{(t)} - \mu_i)(\mathbf{x}^{(t)} - \mu_i)^\top = \mathbf{O} \end{cases}$$

- $$\mathbf{m}_i = \frac{\sum_{t=1}^N \mathbf{x}^{(t)} h_i^{(t)}}{\sum_{t=1}^N h_i^{(t)}}$$
- $$\mathbf{S}_i = \frac{\sum_{t=1}^N (\mathbf{x}^{(t)} - \mathbf{m}_i)(\mathbf{x}^{(t)} - \mathbf{m}_i)^\top h_i^{(t)}}{\sum_{t=1}^N h_i^{(t)}}$$

- Again, these results align with the  $\mathbf{m}_i$  and  $\mathbf{S}_i$  in semiparametric density estimation

- Both the  $K$ -means and semiparametric density estimation are EM algorithms

- The iteration ends and  $\Theta$  converges to a local optimum

- In particular, when assuming that the priors  $\pi_i$  are all equal and  $\Sigma_i = \sigma^2 \mathbf{I}$ , we have

- $$h_i^{(t)} = \frac{\exp[-(1/2)(s^{old})^2 \|\mathbf{x}^{(t)} - \mathbf{m}_i^{old}\|^2]}{\sum_{j=1}^K \exp[-(1/2)(s^{old})^2 \|\mathbf{x}^{(t)} - \mathbf{m}_j^{old}\|^2]}$$

- The objective  $\arg_{\theta_1, \dots, \theta_K} \max \sum_{t=1}^N \sum_{i=1}^K \ln \left( P[\mathbf{x}^{(t)} | z_i^{(t)}, \theta_i] \right) h_i^{(t)}$  can be

rewritten as  $\arg_{\mathbf{m}_1, \dots, \mathbf{m}_K, s} \min \sum_{t=1}^N \sum_{i=1}^K \frac{\|\mathbf{x}^{(t)} - \mathbf{m}_i\|^2}{s^2} h_i^{(t)}$

- This is equivalent to minimizing the reconstruction error in the  $K$ -means