

# CSCE 633: Machine Learning

## Lecture 24: Unsupervised Learning: EM

Texas A&M University

10-18-19

# Last Time

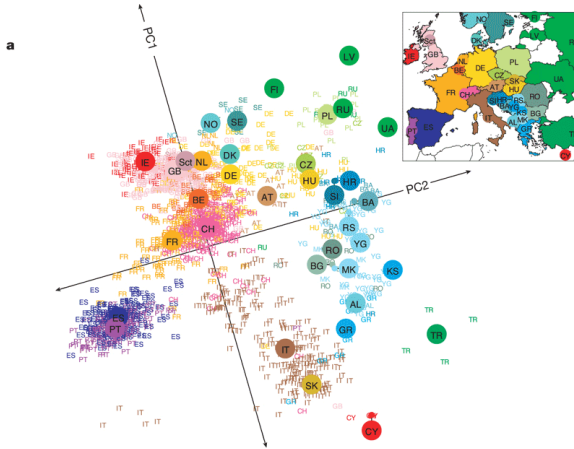
- PCA
- Clustering, K-Means

# Goals of this lecture

- Expectation Maximization

# Clustering

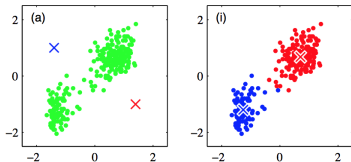
Finding patterns/structure/sub-populations in data



# K-means Clustering

## Representation

- **Input:** Data  $\mathcal{D} = \{\mathbf{x}_1, \dots, \mathbf{x}_N\}$
- **Output:** Clusters  $\mu_1, \dots, \mu_K$
- **Decision:** Cluster membership, the cluster id assigned to sample  $\mathbf{x}_n$ , i.e.  $A(\mathbf{x}_n) \in \{1, \dots, K\}$
- **Evaluation metric:** Distortion measure
$$J = \sum_{n=1}^N \sum_{k=1}^K r_{nk} \|\mathbf{x}_n - \mu_k\|_2^2, \text{ where } r_{nk} = 1 \text{ if } A(\mathbf{x}_n) = k, 0 \text{ otherwise}$$
- **Intuition:** Data points assigned to cluster  $k$  should be close to centroid  $\mu_k$



## K-means Clustering

Evaluation metric:  $\min_{r_{nk}} J = \min_{r_{nk}} \sum_{n=1}^N \sum_{k=1}^K r_{nk} \|\mathbf{x}_n - \boldsymbol{\mu}_k\|_2^2$

Optimization:

- **Step 0:** Initialize  $\boldsymbol{\mu}_k$  to some values
- **Step 1:** Assume the current value of  $\boldsymbol{\mu}_k$  fixed, minimize  $J$  over  $r_{nk}$ , which leads to the following cluster assignment rule
$$r_{nk} = \begin{cases} 1, & \text{if } k = \arg \min_j \|\mathbf{x}_n - \boldsymbol{\mu}_j\|_2^2 \\ 0, & \text{otherwise} \end{cases}$$
- **Step 2:** Assume the current value of  $r_{nk}$  fixed, minimize  $J$  over  $\boldsymbol{\mu}_k$ , which leads to the following rule to update the prototypes of the clusters 
$$\boldsymbol{\mu}_k = \frac{\sum_n r_{nk} \mathbf{x}_n}{\sum_n r_{nk}}$$
- **Step 3:** Determine whether to stop or return to Step 1

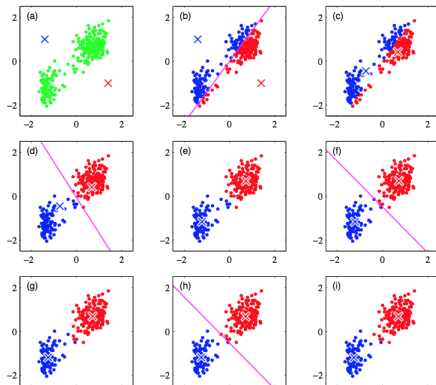
# K-means Clustering

## Remarks

- The centroid  $\mu_k$  is the means of data points assigned to the cluster  $k$ , hence the name K-means clustering.
- The procedure terminates after a finite number of steps, as the procedure reduces  $J$  in both Step 1 and Step 2
- There is no guarantee the procedure terminates at the global optimum of  $J$ . In most cases, the algorithm stops at a **local optimum**, which depends on the initial values in Step 0  $\rightarrow$  **random restarts** to improve chances of getting closer to global optima

# K-means Clustering

## Example





# K-means Clustering

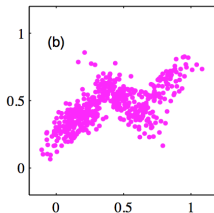
## Application: vector quantization

- We can replace our data points with the centroids  $\mu_k$  from the clusters they are assigned to → **vector quantization**
- We have compressed the data points into
  - a codebook of all the centroids  $\{\mu_1, \dots, \mu_K\}$
  - a list of indices to the codebook for the data points (created based on  $r_{nk}$ )
- This compression is obviously lossy as certain information will be lost if we use a very small  $K$



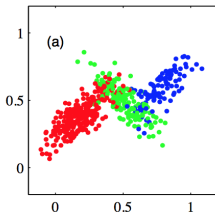
## Probabilistic interpretation of clustering

- We want to find  $p(\mathbf{x})$  that best describes our data
- The data points seem to form 3 clusters
- However, we cannot model  $p(\mathbf{x})$  with simple and known distributions, e.g. one Gaussian



## Probabilistic interpretation of clustering

- Instead, we will model each region with a Gaussian distribution → Gaussian mixture models (GMMs)
- Question 1: How do we know which (color) region a data point comes from?
- Question 2: What are the parameters of Gaussian distributions in each region?
- We will answer both in an unsupervised way from data  $\mathcal{D} = \{\mathbf{x}_1, \dots, \mathbf{x}_n\}$



## Gaussian mixture models: formal definition

A Gaussian mixture model has the following density function for  $\mathbf{x}$

$$p(\mathbf{x}) = \sum_{k=1}^K \omega_k \mathcal{N}(\mathbf{x}; \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)$$

- $K$ : number of Gaussians
- $\boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k$ : mean & covariance of  $k^{th}$  component
- $\omega_k$ : component weights, how much component  $k$  contributes to the final distribution

$$\omega_k > 0, \quad \forall k \quad \text{and} \quad \sum_{k=1}^K \omega_k = 1$$

$\omega_k$  can be represented by the **prior** distribution:  $\omega_k = p(z = k)$ , which decides which mixture to use

## GMM as the marginal distribution of a joint distribution

- Consider the following joint distribution  $p(\mathbf{x}, z) = p(z)p(\mathbf{x}|z)$
- $z$  is a discrete random variable taking values between 1 and  $K$ , “selects” a Gaussian component
- We denote  $\omega_k = p(z = k)$
- Assume Gaussian conditional distributions  $p(\mathbf{x}|z = k) = \mathcal{N}(\mathbf{x}; \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)$
- Then the marginal distribution of  $\mathbf{x}$  is

$$p(\mathbf{x}) = \sum_{k=1}^K \omega_k \mathcal{N}(\mathbf{x}; \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)$$

(which is the Gaussian mixture model)

## GMM as the marginal distribution of a joint distribution

- The joint distribution between  $\mathbf{x}$  and  $z$  (representing color) are

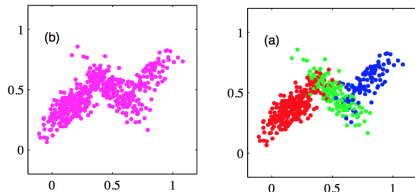
$$p(\mathbf{x}|z = \text{red}) = \mathcal{N}(\mathbf{x}; \mu_1, \Sigma_1)$$

$$p(\mathbf{x}|z = \text{blue}) = \mathcal{N}(\mathbf{x}; \mu_2, \Sigma_2)$$

$$p(\mathbf{x}|z = \text{green}) = \mathcal{N}(\mathbf{x}; \mu_3, \Sigma_3)$$

- The marginal distribution is thus

$$p(\mathbf{x}) = p(\text{red})\mathcal{N}(\mathbf{x}; \mu_1, \Sigma_1) + p(\text{blue})\mathcal{N}(\mathbf{x}; \mu_2, \Sigma_2) \\ + p(\text{green})\mathcal{N}(\mathbf{x}; \mu_3, \Sigma_3)$$



## Parameter estimation for GMMs: the easy case with complete data

We know the component in which each sample belongs to

- Data  $\mathcal{D} = \{(\mathbf{x}_1, z_1), \dots, (\mathbf{x}_N, z_N)\}$
- We want to find  $\boldsymbol{\theta} = \{\boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k, \omega_k\}$
- Maximum log-likelihood:  $\boldsymbol{\theta}^* = \arg \max \log(p(\mathcal{D})|\boldsymbol{\theta})$

Solution

$$\omega_k = \frac{\sum_n \gamma_{nk}}{\sum_k \sum_n \gamma_{nk}} \quad \boldsymbol{\mu}_k = \frac{1}{\sum_n \gamma_{nk}} \sum_n \gamma_{nk} \mathbf{x}_n$$

$$\boldsymbol{\Sigma}_k = \frac{1}{\sum_n \gamma_{nk}} \sum_n \gamma_{nk} (\mathbf{x}_n - \boldsymbol{\mu}_k)(\mathbf{x}_n - \boldsymbol{\mu}_k)^T$$

where  $\gamma_{nk} = 1$  if  $z_n = k$

## Parameter estimation for GMMs: the easy case with complete data

### Understanding the intuition

$$\omega_k = \frac{\sum_n \gamma_{nk}}{\sum_k \sum_n \gamma_{nk}} \quad \mu_k = \frac{1}{\sum_n \gamma_{nk}} \sum_n \gamma_{nk} \mathbf{x}_n$$

$$\Sigma_k = \frac{1}{\sum_n \gamma_{nk}} \sum_n \gamma_{nk} (\mathbf{x}_n - \mu_k)(\mathbf{x}_n - \mu_k)^T$$

- For  $\omega_k$ : count the number of data points whose  $z_n$  is  $k$  and divide by the total number of data points
- For  $\mu_k$ : get all the data points whose  $z_n$  is  $k$ , compute their mean
- For  $\Sigma_k$ : get all the data points whose  $z_n$  is  $k$ , compute their covariance



## Parameter estimation for GMMs: incomplete data

Trick: estimation with soft  $\gamma_{nk}$

$$\gamma_{nk} = p(z_n = k | \mathbf{x}_n)$$

$$\omega_k = \frac{\sum_n \gamma_{nk}}{\sum_k \sum_n \gamma_{nk}}$$

$$\boldsymbol{\mu}_k = \frac{1}{\sum_n \gamma_{nk}} \sum_n \gamma_{nk} \mathbf{x}_n$$

$$\boldsymbol{\Sigma}_k = \frac{1}{\sum_n \gamma_{nk}} \sum_n \gamma_{nk} (\mathbf{x}_n - \boldsymbol{\mu}_k)(\mathbf{x}_n - \boldsymbol{\mu}_k)^T$$

Every data point  $\mathbf{x}_n$  is assigned to a component fractionally according to  $p(z_n = k | \mathbf{x}_n)$ , also called **responsibility**

## Parameter estimation for GMMs: incomplete data

Trick: estimation with soft  $\gamma_{nk}$

Since we do not know  $\theta$  to begin with, we cannot compute the soft  $\gamma_{nk}$ . But we can invoke an iterative procedure and alternate between estimating  $\gamma_{nk}$  and using the estimated  $\gamma_{nk}$  to compute  $\theta = \{\mu_k, \Sigma_k, \omega_k\}$

- **Step 0:** guess  $\theta$  with initial values
- **Step 1:** compute  $\gamma_{nk}$  using current  $\theta$
- **Step 2:** update  $\theta$  using computed  $\gamma_{nk}$
- **Step 3:** go back to Step 1

Questions: i) is this procedure correct, for example, optimizing a sensible criterion? ii) practically, will this procedure ever stop instead of iterating forever? The answer lies in the **Expectation Maximization (EM)** algorithm, a powerful procedure for model estimation with unknown data

## Expectation Maximization: motivation and setup

- EM is used to estimate parameters for probabilistic models with hidden/latent variables

$$p(\mathbf{x}|\theta) = \sum_z p(\mathbf{x}, \mathbf{z}|\theta)$$

where  $\mathbf{x}$  is the observed random variable and  $\theta$  is the hidden

- We are given data  $\mathcal{D} = \{\mathbf{x}_1, \dots, \mathbf{x}_N\}$ , where the corresponding hidden variable values  $z$  are not included
- Our goal is to obtain the maximum likelihood estimate of  $\theta$

$$\theta^* = \arg \max_{\theta} \sum_{n=1}^N \log p(\mathbf{x}_n|\theta) = \arg \max_{\theta} \sum_{n=1}^N \log \sum_{\mathbf{z}_n} p(\mathbf{x}_n, \mathbf{z}_n|\theta)$$

The above objective function is called **incomplete** log-likelihood and it is **computationally intractable** (since it needs to sum over all possible values of  $\mathbf{z}$  and then take the logarithm)

## Expectation Maximization: Complete log-likelihood

- Incomplete log-likelihood

$$l(\theta) = \sum_{n=1}^N \log \sum_{\mathbf{z}_n} p(\mathbf{x}_n, \mathbf{z}_n | \theta)$$

- EM uses a clever trick to change this into sum-log form by defining:

$$\begin{aligned} Q_q(\theta) &= \sum_{n=1}^N \log \mathbb{E}_{\mathbf{z}_n \sim q(\mathbf{z}_n)} p(\mathbf{x}_n, \mathbf{z}_n | \theta) \\ &= \sum_{n=1}^N \sum_{\mathbf{z}_n} q(\mathbf{z}_n) \log p(\mathbf{x}_n, \mathbf{z}_n | \theta) \end{aligned}$$

where  $q(\mathbf{z})$  is a distribution over  $\mathbf{z}$

The above is called **expected (complete)** log-likelihood, since it takes the form of log-sum which is **computationally tractable**

## Expectation Maximization: Choice of $q(\mathbf{z})$

- We will choose a special  $q(\mathbf{z}) = p(\mathbf{z}|\mathbf{x}, \theta)$ , i.e., the posterior probability of  $\mathbf{z}$

$$Q(\theta) = Q_{\mathbf{z} \sim p(\mathbf{z}|\mathbf{x}, \theta)}(\theta)$$

- We will show that

$$l(\theta) = Q(\theta) + \sum_{n=1}^N \mathbb{H}[p(\mathbf{z}|\mathbf{x}_n, \theta)]$$

where  $\mathbb{H}$  is the entropy of the probabilistic distribution  $p(\mathbf{z}|\mathbf{x}, \theta)$

$$\mathbb{H}[p(\mathbf{z}|\mathbf{x}, \theta)] = - \sum_{\mathbf{z}_n} p(\mathbf{z}_n|\mathbf{x}, \theta) \log p(\mathbf{z}_n|\mathbf{x}, \theta)$$

## Expectation Maximization: Choice of $q(\mathbf{z})$

### Proof

$$Q_q(\theta) = \sum_n \sum_{\mathbf{z}_n} q(\mathbf{z}_n) \log p(\mathbf{x}_n, \mathbf{z}_n | \theta)$$

$$\begin{aligned} Q(\theta) &= Q_{\mathbf{z} \sim p(\mathbf{z}|\mathbf{x}, \theta)}(\theta) = \sum_n \sum_{\mathbf{z}_n} p(\mathbf{z}_n | \mathbf{x}_n, \theta) \log p(\mathbf{x}_n, \mathbf{z}_n | \theta) \\ &= \sum_n \sum_{\mathbf{z}_n} p(\mathbf{z}_n | \mathbf{x}_n, \theta) [\log p(\mathbf{x}_n | \theta) + \log p(\mathbf{z}_n | \mathbf{x}_n, \theta)] \\ &= \sum_n \sum_{\mathbf{z}_n} p(\mathbf{z}_n | \mathbf{x}_n, \theta) \log p(\mathbf{x}_n | \theta) + \sum_n \sum_{\mathbf{z}_n} p(\mathbf{z}_n | \mathbf{x}_n, \theta) \log p(\mathbf{z}_n | \mathbf{x}_n, \theta) \\ &= \sum_n \log p(\mathbf{x}_n | \theta) \sum_{\mathbf{z}_n} p(\mathbf{z}_n | \mathbf{x}_n, \theta) - \mathbb{H}[p(\mathbf{z} | \mathbf{x}_n, \theta)] \\ &= l(\theta) - \mathbb{H}[p(\mathbf{z} | \mathbf{x}_n, \theta)] \end{aligned}$$

## Expectation Maximization

- As before,  $Q(\theta)$  cannot be computed, as  $p(\mathbf{z}|\mathbf{x}, \theta)$  depends on the unknown parameter values  $\theta$
- Instead, we will use a known value  $\theta^{old}$  to compute the expected likelihood

$$Q(\theta, \theta^{old}) = \sum_n \sum_{\mathbf{z}_n} p(\mathbf{z}_n | \mathbf{x}_n, \theta^{old}) \log p(\mathbf{x}_n, \mathbf{z}_n | \theta)$$

- In the above, the variable is  $\theta$ , while  $\theta^{old}$  is known
- By its definition  $Q(\theta) = Q(\theta, \theta^{old})$
- But how does  $Q(\theta, \theta^{old})$  relates to  $I(\theta)$ ?
- We will show that

$$I(\theta) \geq Q(\theta, \theta^{old}) + \sum_n \mathbb{H} [p(\mathbf{z} | \mathbf{x}_n, \theta^{old})]$$

- Thus,  $Q(\theta)$  is better than  $Q(\theta, \theta^{old})$ , except that we cannot compute the former

# Expectation Maximization

## Proof

$$\begin{aligned}l(\theta) &= \sum_n \log p(\mathbf{x}_n|\theta) = \sum_n \log \sum_{\mathbf{z}_n} p(\mathbf{x}_n, \mathbf{z}_n|\theta) \\&= \sum_n \log \sum_{\mathbf{z}_n} p(\mathbf{z}|\mathbf{x}_n, \theta^{old}) \frac{p(\mathbf{x}_n, \mathbf{z}_n|\theta)}{p(\mathbf{z}|\mathbf{x}_n, \theta^{old})} \\&\geq \sum_n \sum_{\mathbf{z}_n} p(\mathbf{z}|\mathbf{x}_n, \theta^{old}) \log p(\mathbf{x}_n, \mathbf{z}_n|\theta) \\&\quad - \sum_n \sum_{\mathbf{z}_n} p(\mathbf{z}|\mathbf{x}_n, \theta^{old}) \log p(\mathbf{z}|\mathbf{x}_n, \theta^{old}) \\&= \underbrace{Q(\theta, \theta^{old}) + \sum_n \mathbb{H} [p(\mathbf{z}|\mathbf{x}_n, \theta^{old})]}_{A(\theta, \theta^{old}) : \text{auxiliary function}}\end{aligned}$$

In the above we have used that  $\log \sum_i w_i x_i \geq \sum_i w_i \log x_i$ ,  $\forall w_i \geq 0$  s.t.  $\sum_i w_i = 1$



## Expectation Maximization

- Suppose we have an initial guess  $\theta^{old}$  and we maximize the auxiliary function

$$\theta^{new} = \arg \max_{\theta} A(\theta, \theta^{old})$$

- With the new guess, we have

$$l(\theta^{new}) \geq A(\theta^{new}, \theta^{old}) \geq A(\theta^{new}, \theta^{old}) = l(\theta^{old})$$

- By maximizing the auxiliary function, we will keep increasing the likelihood (core of EM)

# Expectation Maximization

## Algorithm outline

- **Step 0:** Initialize  $\theta$  with  $\theta^{(0)}$
- **Step 1 (E-step):** Compute the auxiliary function using the current value of  $\theta$

$$A(\theta, \theta^{(t)})$$

- **Step 2 (M-step):** Maximize the auxiliary function

$$\theta^{(t+1)} \leftarrow \arg \max A(\theta, \theta^{(t)})$$

- **Step 3:** Increase  $t$  to  $t + 1$  and back to Step 1; or stop if  $l(\theta^{(t+1)})$  does not improve much

# Expectation Maximization

## Remarks

- EM converges but only to a local optimum: global optimum is not guaranteed
- The E-step depends on computing the posterior probability  $p(\mathbf{z}_n | \mathbf{x}_n, \theta^{(t)})$
- The M-step does not depend on the entropy term, so we need only to do the following

$$\theta^{(t+1)} \leftarrow \arg \max_{\theta} A(\theta, \theta^{(t)}) = \arg \max_{\theta} Q(\theta, \theta^{(t)})$$

We often call the last term Q-function

## Expectation Maximization for GMMs

- Hidden variable follows Multinomial distribution  
 $\mathbf{z}_n \sim \text{Multinomial}(\omega_1, \dots, \omega_K)$
- The likelihood of observation  $\mathbf{x}_n$  equals the probability of corresponding component  $p(\mathbf{x}_n | \mathbf{z}_n, \theta) = \prod_{k=1}^K \mathcal{N}(\mathbf{x}_n | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)^{\gamma_{nk}}$
- Complete data log-likelihood

$$\begin{aligned} Q(\theta) &= \mathbb{E} \left[ \sum_n \log p(\mathbf{x}_n, \mathbf{z}_n | \theta) \right] = \mathbb{E} \left[ \sum_n \log (p(\mathbf{z}_n | \theta) p(\mathbf{x}_n | \mathbf{z}_n, \theta)) \right] \\ &= \mathbb{E} \left[ \sum_n \log \left( \prod_k \omega_k^{\gamma_{nk}} \prod_k \mathcal{N}(\mathbf{x}_n | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)^{\gamma_{nk}} \right) \right] \\ &= \mathbb{E} \left[ \sum_n \sum_k \gamma_{nk} (\log \omega_k + \log \mathcal{N}(\mathbf{x}_n | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)) \right] \\ &= \sum_n \sum_k \mathbb{E}[\gamma_{nk} | \mathbf{x}_n, \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k] \gamma_{nk} (\log \omega_k + \log \mathcal{N}(\mathbf{x}_n | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)) \end{aligned}$$

## Expectation Maximization for GMMs

What is the E-step in GMM?

$$\mathbb{E}[\gamma_{nk} | \mathbf{x}_n, \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k] = p(\gamma_{nk} = 1 | \mathbf{x}_n, \boldsymbol{\theta}_k) = \frac{p(\mathbf{x}_n | \gamma_{nk} = 1, \boldsymbol{\theta}) p(\gamma_{nk} = 1 | \boldsymbol{\theta}_k)}{\sum_k p(\mathbf{x}_n | \gamma_{nk} = 1, \boldsymbol{\theta}) p(\gamma_{nk} = 1 | \boldsymbol{\theta}_k)}$$

We compute the probability

$$\gamma_{nk} = p(z = k | \mathbf{x}_n, \boldsymbol{\theta}^{(t)})$$
$$\gamma_{nk} = \frac{\omega_k p(\mathbf{x}_n | \boldsymbol{\mu}_k^{(t-1)}, \boldsymbol{\Sigma}_k^{(t-1)})}{\sum_k \omega_k p(\mathbf{x}_n | \boldsymbol{\mu}_k^{(t-1)}, \boldsymbol{\Sigma}_k^{(t-1)})}$$

## Expectation Maximization for GMMs

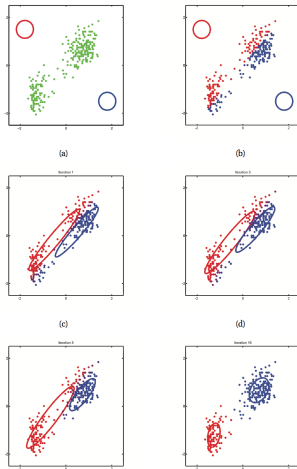
What is the M-step in GMM?

Maximize the auxiliary function

$$\begin{aligned}Q(\theta, \theta^{(t)}) &= \sum_n \sum_k p(z = k | \mathbf{x}_n, \theta^{(t)}) \log p(\mathbf{x}_n, z = k | \theta) \\&= \sum_n \sum_k \gamma_{nk} \log p(\mathbf{x}_n, z = k | \theta) \\&= \sum_n \sum_k \gamma_{nk} \log(p(z = k) p(\mathbf{x}_n | z = k)) \\&= \sum_n \sum_k \gamma_{nk} [\log \omega_k + \log \mathcal{N}(\mathbf{x}_n | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)] \\&\Rightarrow \omega_k = \frac{\sum_n \gamma_{nk}}{N} \\&\Rightarrow \boldsymbol{\mu}_k = \frac{1}{\sum_n \gamma_{nk}} \sum_n \gamma_{nk} \mathbf{x}_n \\&\Rightarrow \boldsymbol{\Sigma}_k = \frac{1}{\sum_n \gamma_{nk}} \sum_n \gamma_{nk} (\mathbf{x}_n - \boldsymbol{\mu}_k)(\mathbf{x}_n - \boldsymbol{\mu}_k)^T\end{aligned}$$

# Expectation Maximization for GMMs

## Example of EM implementation



## GMMs and K-means

GMMs provide probabilistic interpretation for K-means

- Assume all Gaussian components have  $\sigma^2 \mathbf{I}$  as their covariance matrices
- Further assume  $\sigma \rightarrow 0$
- Thus, we only need to estimate  $\mu_k$ , i.e., means
- Then, the EM for GMM parameter estimation simplifies to K-means

For this reason, K-means is often called **hard GMM** or GMMs is called **soft K-means**. The soft posterior  $\gamma_{nk}$  provides a probabilistic assignment for  $\mathbf{x}_n$  to cluster  $k$  represented by the corresponding Gaussian distribution



# Takeaways and Next Time

- Expectation Maximization and GMM
- Next Time: More Unsupervised Learning