

Mixture Models and Expectation-Maximization

Yaqiang Yao

School of Computer Science and Technology
University of Science and Technology of China
Hefei China

April 9, 2018

Outline

- 1 K -Means Clustering
- 2 Mixtures of Gaussians
- 3 An Alternative View of EM
- 4 The EM algorithm in General
- 5 Hidden Markov Model

Outline

- 1 *K*-Means Clustering
- 2 Mixtures of Gaussians
- 3 An Alternative View of EM
- 4 The EM algorithm in General
- 5 Hidden Markov Model

Motivation



Motivation



K -Means Clustering: Distortion Measure

- Dataset $\{\mathbf{x}_1, \dots, \mathbf{x}_N\}$
- Partition in K clusters
- Introduce a set of vectors (Cluster prototype): $\boldsymbol{\mu}_k$
- Introduce a corresponding set of binary indicator variable $r_{nk} \in \{0, 1\}$ (1-of- K coding scheme), such that

$$r_{nk} = \begin{cases} 1 & \text{if } \mathbf{x}_n \text{ is assigned to cluster } k \\ 0 & \text{otherwise} \end{cases}$$

- Distortion measure

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

K-Means Clustering: Expectation Maximization

- Find values for $\{r_{nk}\}$ and $\{\mu_k\}$ to minimize

$$J = \sum_{n=1}^N \sum_{k=1}^K r_{nk} \|\mathbf{x}_n - \mu_k\|^2$$

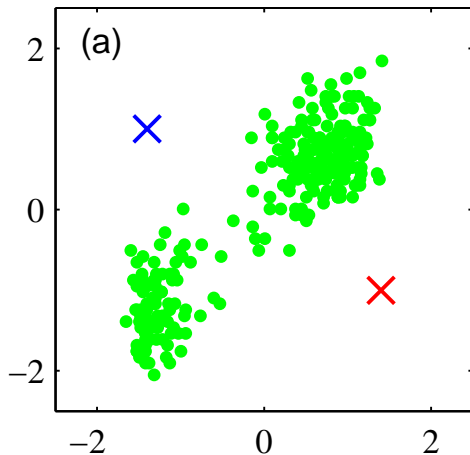
- Interactive procedure
 1. Minimize J w.r.t r_{nk} , keep μ_k fixed (Expectation)

$$r_{nk} = \begin{cases} 1 & \text{if } k = \arg \min_j \|\mathbf{x}_n - \mu_j\|^2 \\ 0 & \text{otherwise} \end{cases}$$

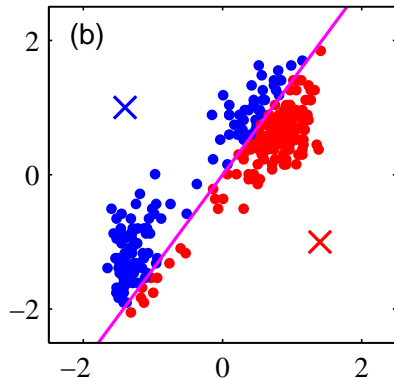
2. Minimize J w.r.t μ_k , keep r_{nk} fixed (Maximization)

$$\begin{aligned} 2 \sum_{n=1}^N r_{nk} (\mathbf{x}_n - \mu_k) &= 0 \\ \implies \mu_k &= \frac{\sum_n r_{nk} \mathbf{x}_n}{\sum_n r_{nk}} \end{aligned}$$

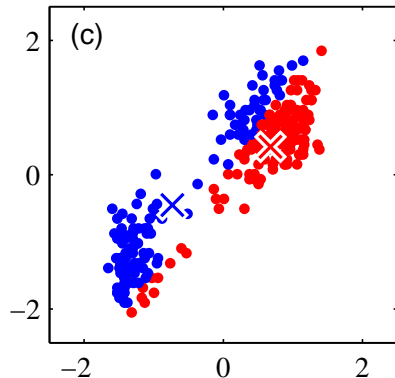
K -Means Clustering: Example



K -Means Clustering: Example cont.

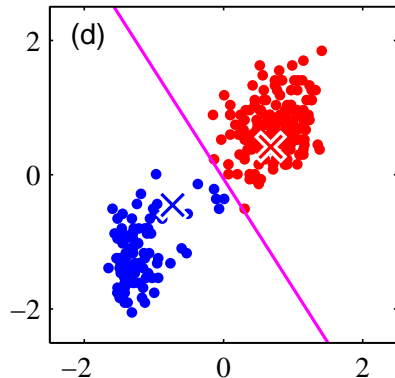


E Step

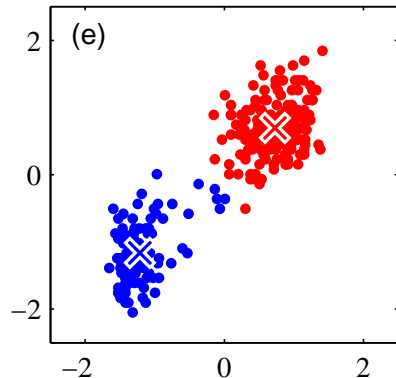


M Step

K -Means Clustering: Example cont.

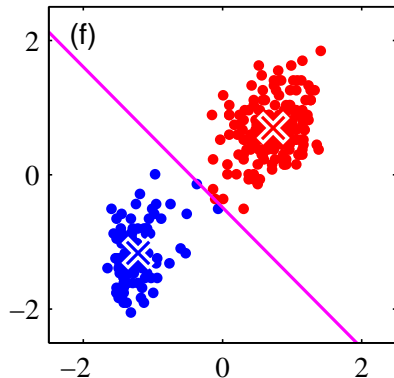


E Step

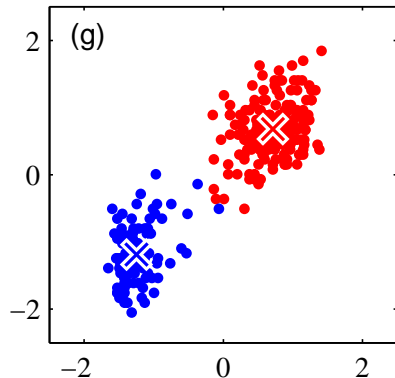


M Step

K -Means Clustering: Example cont.

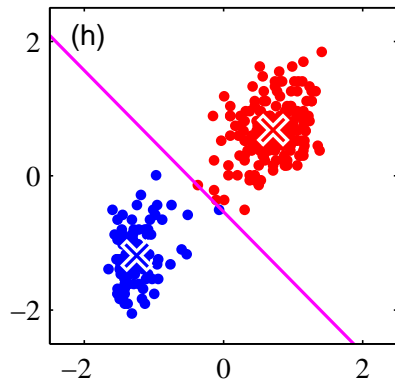


E Step

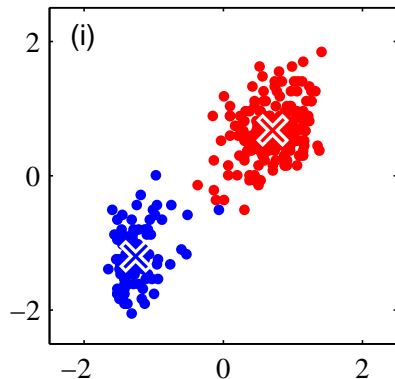


M Step

K -Means Clustering: Example cont.

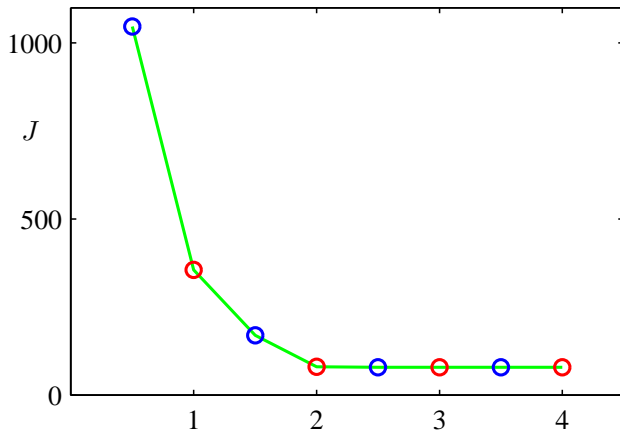


E Step



M Step

K -Means Clustering: Example cont.



- Each E or M step reduces the value of the objective function J
- Convergence to a global or local maximum

K-Means Clustering: Concluding Remarks

1. Direct implementation of *K*-Means can be slow
2. Online version:

$$\boldsymbol{\mu}_k^{\text{new}} = \boldsymbol{\mu}_k^{\text{old}} + \eta_n(\mathbf{x}_n - \boldsymbol{\mu}_k^{\text{old}})$$

3. *K*-medoids, general distortion measure

$$\tilde{J} = \sum_{n=1}^N \sum_{k=1}^K r_{nk} \mathcal{V}(\mathbf{x}_n, \boldsymbol{\mu}_k)$$

where $\mathcal{V}(\cdot)$ is any kind of dissimilarity measure.

4. Hard assignment: $r_{nk} = 1$ and $r_{nj} = 0$ for $j \neq k$.

Image Segmentation and Compression Example

Original image (96,615 colors)



100%

Image Segmentation and Compression Example

Quantized image (10colors, K-Means)



23.2%

Image Segmentation and Compression Example

Quantized image (3colors, K-Means)



8.7%

Image Segmentation and Compression Example

Quantized image (2colors, K-Means)



4.1%

Outline

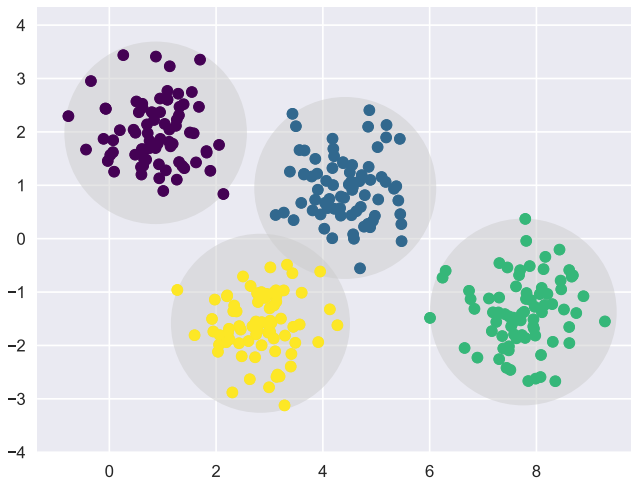
- 1 K -Means Clustering
- 2 Mixtures of Gaussians
- 3 An Alternative View of EM
- 4 The EM algorithm in General
- 5 Hidden Markov Model

Weakness of K -Means

- Non-probabilistic nature
- Use simple distance-from-cluster-center to assign cluster membership

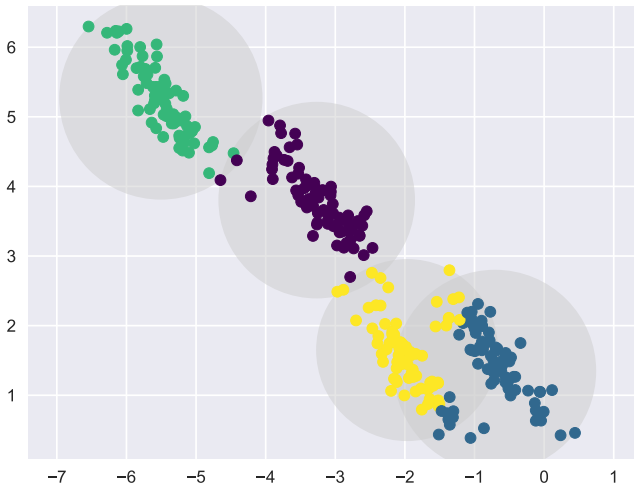
Weakness of K -Means

- Non-probabilistic nature
- Use simple distance-from-cluster-center to assign cluster membership



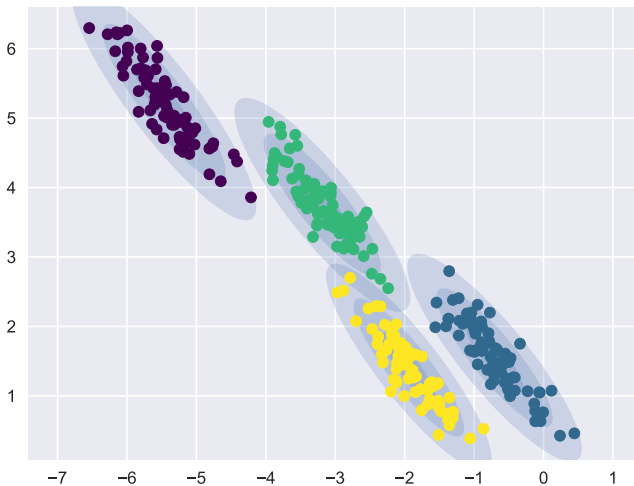
Weakness of K -Means

- Non-probabilistic nature
- Use simple distance-from-cluster-center to assign cluster membership



Weakness of K -Means

- Non-probabilistic nature
- Use simple distance-from-cluster-center to assign cluster membership



Mixture of Gaussians: Latent Variables

- Gaussian Mixture Distribution:

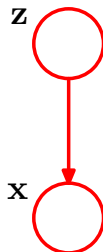
$$p(\mathbf{x}) = \sum_{k=1}^K \pi_k \mathcal{N}(\mathbf{x} | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)$$

where

$$\mathcal{N}(\mathbf{x} | \boldsymbol{\mu}, \boldsymbol{\Sigma}) = \frac{1}{(2\pi)^{D/2}} \frac{1}{|\boldsymbol{\Sigma}|^{1/2}} \exp \left\{ -\frac{1}{2} (\mathbf{x} - \boldsymbol{\mu})^T \boldsymbol{\Sigma}^{-1} (\mathbf{x} - \boldsymbol{\mu}) \right\}$$

- Introduce latent variable z

1. z is binary 1-of- K coding variable
2. $p(\mathbf{x}, z) = p(z)p(\mathbf{x}|z)$



Mixture of Gaussians: Conditional Distribution

- The marginal distribution of z (prior distribution)

$$p(z_k = 1) = \pi_k$$

where $0 \leq \pi_k \leq 1$ and $\sum_k \pi_k = 1$

- Since z uses a 1-of- K representation

$$p(\mathbf{z}) = \prod_k \pi_k^{z_k}$$

Mixture of Gaussians: Conditional Distribution

- The marginal distribution of \mathbf{z} (prior distribution)

$$p(z_k = 1) = \pi_k$$

where $0 \leq \pi_k \leq 1$ and $\sum_k \pi_k = 1$

- Since \mathbf{z} uses a 1-of- K representation

$$p(\mathbf{z}) = \prod_k \pi_k^{z_k}$$

- Given a particular value of \mathbf{z} , the conditional distribution of \mathbf{x} is a Gaussian

$$p(\mathbf{x}|z_k = 1) = \mathcal{N}(\mathbf{x}|\boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)$$

- Another form of the above conditional distribution

$$p(\mathbf{x}|\mathbf{z}) = \prod_{k=1}^K \mathcal{N}(\mathbf{x}|\boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)^{z_k}$$

Mixture of Gaussians: Marginal Distribution

- Joint distribution

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

- Marginal distribution

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

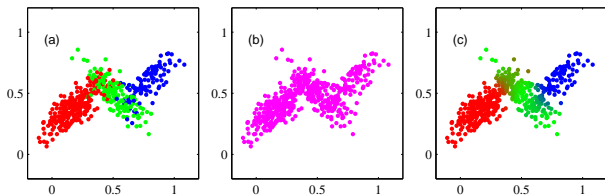
- The use of the joint probability $p(\mathbf{x}, \mathbf{z})$, leads to significant simplifications.

Mixture of Gaussians: Posterior Distribution

- The **responsibility** of component k to generate observation \mathbf{x}

$$\begin{aligned}\gamma(z_k) \equiv p(z_k = 1 | \mathbf{x}) &= \frac{p(z_k = 1)p(\mathbf{x} | z_k = 1)}{\sum_k p(z_k = 1)p(\mathbf{x} | z_k = 1)} \\ &= \frac{\pi_k p(\mathbf{x} | z_k = 1)}{\sum_k \pi_k p(\mathbf{x} | z_k = 1)}\end{aligned}$$

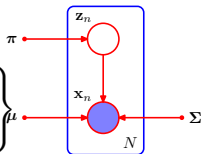
- Generate random samples with ancestral sampling
 - $\hat{z} \sim p(\mathbf{z})$
 - $\mathbf{x} \sim p(\mathbf{x} | \hat{z})$



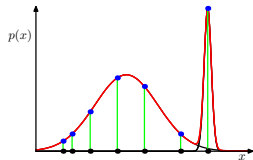
Mixture of Gaussians: Posterior Distribution

- Log Likelihood

$$\ln p(\mathbf{X}|\boldsymbol{\pi}, \boldsymbol{\mu}, \boldsymbol{\Sigma}) = \sum_{n=1}^N \ln \left\{ \sum_{k=1}^K \pi_k \mathcal{N}(\mathbf{x}|\boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k) \right\}$$



- Singularity**: when a mixture component collapses on a datapoint
- Identifiability**: for a ML solution in a K -component mixture there are $K!$ equivalent solutions.



EM for Gaussian Mixtures

- Maximum of log likelihood: derivatives of $\ln p(\mathbf{X}|\boldsymbol{\pi}, \boldsymbol{\mu}, \boldsymbol{\Sigma})$ w.r.t parameters to 0.
- For the $\boldsymbol{\mu}_k$:

$$0 = - \sum_{n=1}^N \underbrace{\frac{\pi_k \mathcal{N}(\mathbf{x}_n | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)}{\sum_k \pi_k \mathcal{N}(\mathbf{x}_n | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)}}_{\gamma(z_{nk})} \boldsymbol{\Sigma}_k^{-1} (\mathbf{x}_n - \boldsymbol{\mu}_k)$$
$$\Rightarrow \boldsymbol{\mu}_k = \frac{1}{\sum_n \gamma(z_{nk})} \sum_{n=1}^N \gamma(z_{nk}) \mathbf{x}_n$$

- For the $\boldsymbol{\Sigma}_k$:

$$\boldsymbol{\Sigma}_k = \frac{1}{\sum_n \gamma(z_{nk})} \sum_{n=1}^N \gamma(z_{nk}) (\mathbf{x}_n - \boldsymbol{\mu}_k)(\mathbf{x}_n - \boldsymbol{\mu}_k)^T$$

EM for Gaussian Mixtures cont.

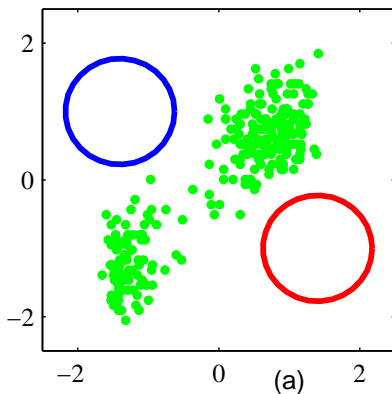
- For the π_k :
 1. Take account of the constraint $\sum_k \pi_k = 1$
 2. Lagrange multiplier

$$\begin{aligned} & \ln p(\mathbf{X}|\boldsymbol{\pi}, \boldsymbol{\mu}, \boldsymbol{\Sigma}) + \lambda \left(\sum_{k=1}^K \pi_k - 1 \right) \\ \implies 0 &= \sum_{n=1}^N \frac{\mathcal{N}(\mathbf{x}_n | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)}{\sum_{j=1}^K \pi_j \mathcal{N}(\mathbf{x}_n | \boldsymbol{\mu}_j, \boldsymbol{\Sigma}_j)} + \lambda \\ \implies \lambda \sum_{k=1}^K \pi_k &= - \sum_{n=1}^N \frac{\sum_{k=1}^K \pi_k \mathcal{N}(\mathbf{x}_n | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)}{\sum_{j=1}^K \pi_j \mathcal{N}(\mathbf{x}_n | \boldsymbol{\mu}_j, \boldsymbol{\Sigma}_j)} = -N \\ \implies \pi_k &= \frac{\sum_{n=1}^N \gamma(z_{nk})}{N} \equiv N_k \end{aligned}$$

- N_k can be interpreted as the effective number of points assigned to cluster k .

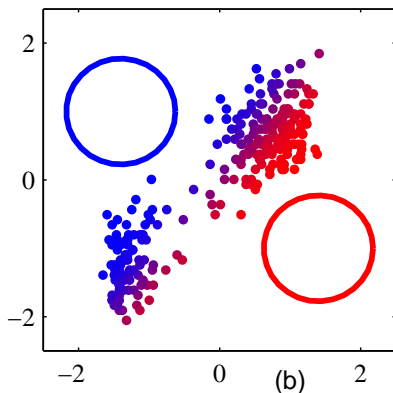
EM for Gaussian Mixtures: Example

- No closed form solutions: $\gamma(z_{nk})$ depends on parameters
- But these equations suggest simple iterative scheme for finding maximum likelihood
- Alternate between estimating the current $\gamma(z_{nk})$ and updating the parameters $\{\boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k, \pi_k\}$



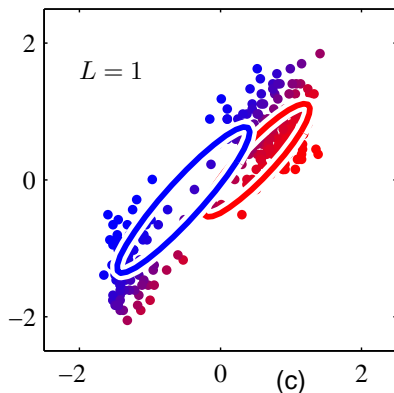
EM for Gaussian Mixtures: Example

- No closed form solutions: $\gamma(z_{nk})$ depends on parameters
- But these equations suggest simple iterative scheme for finding maximum likelihood
- Alternate between estimating the current $\gamma(z_{nk})$ and updating the parameters $\{\mu_k, \Sigma_k, \pi_k\}$



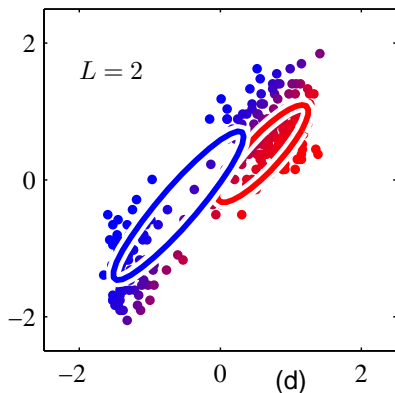
EM for Gaussian Mixtures: Example

- No closed form solutions: $\gamma(z_{nk})$ depends on parameters
- But these equations suggest simple iterative scheme for finding maximum likelihood
- Alternate between estimating the current $\gamma(z_{nk})$ and updating the parameters $\{\boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k, \pi_k\}$



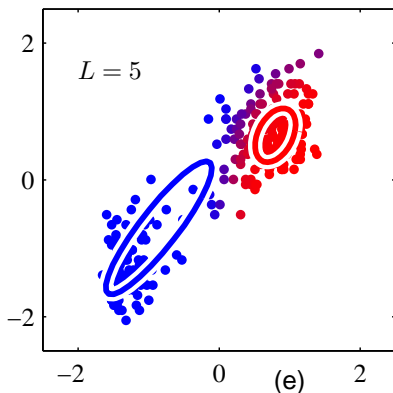
EM for Gaussian Mixtures: Example

- No closed form solutions: $\gamma(z_{nk})$ depends on parameters
- But these equations suggest simple iterative scheme for finding maximum likelihood
- Alternate between estimating the current $\gamma(z_{nk})$ and updating the parameters $\{\boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k, \pi_k\}$



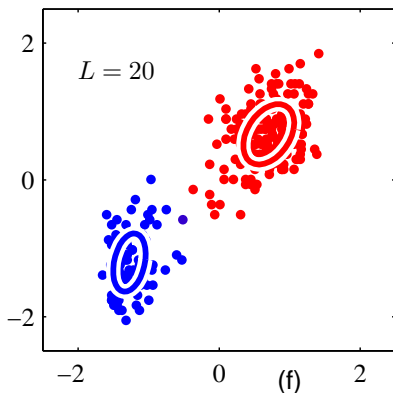
EM for Gaussian Mixtures: Example

- No closed form solutions: $\gamma(z_{nk})$ depends on parameters
- But these equations suggest simple iterative scheme for finding maximum likelihood
- Alternate between estimating the current $\gamma(z_{nk})$ and updating the parameters $\{\boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k, \pi_k\}$



EM for Gaussian Mixtures: Example

- No closed form solutions: $\gamma(z_{nk})$ depends on parameters
- But these equations suggest simple iterative scheme for finding maximum likelihood
- Alternate between estimating the current $\gamma(z_{nk})$ and updating the parameters $\{\boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k, \pi_k\}$



EM for Gaussian Mixtures: Summary

1. Initialize $\{\boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k, \pi_k\}$ and evaluate log-likelihood
2. **E-Step** Evaluate responsibilities

$$\gamma(z_{nk}) = \frac{\pi_k \mathcal{N}(\mathbf{x}_n | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)}{\sum_k \pi_k \mathcal{N}(\mathbf{x}_n | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)}$$

3. **M-Step** Re-estimate parameters, using current responsibilities:

$$\begin{aligned}\boldsymbol{\mu}_k &= \frac{1}{N_k} \sum_{n=1}^N \gamma(z_{nk}) \mathbf{x}_n \\ \boldsymbol{\Sigma}_k &= \frac{1}{N_k} \sum_{n=1}^N \gamma(z_{nk}) (\mathbf{x}_n - \boldsymbol{\mu}_k)(\mathbf{x}_n - \boldsymbol{\mu}_k)^T \\ \pi_k &= \frac{N_k}{N}\end{aligned}$$

4. Evaluate log-likelihood $\ln p(\mathbf{X} | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k, \pi_k)$ and check for convergence (go to step 2).

Outline

- 1 K -Means Clustering
- 2 Mixtures of Gaussians
- 3 An Alternative View of EM**
- 4 The EM algorithm in General
- 5 Hidden Markov Model

An Alternative View of EM: Latent Variables

- Let \mathbf{X} observed data, \mathbf{Z} latent variables, $\boldsymbol{\theta}$ parameters
- Goal: maximize marginal log-likelihood of observed data

$$\ln p(\mathbf{X}|\boldsymbol{\theta}) = \ln \left\{ \sum_{\mathbf{Z}} p(\mathbf{X}, \mathbf{Z}|\boldsymbol{\theta}) \right\}$$

- Optimization problematic due to log-sum

An Alternative View of EM: Latent Variables

- Let \mathbf{X} observed data, \mathbf{Z} latent variables, θ parameters
- Goal: maximize marginal log-likelihood of observed data

$$\ln p(\mathbf{X}|\theta) = \ln \left\{ \sum_{\mathbf{Z}} p(\mathbf{X}, \mathbf{Z}|\theta) \right\}$$

- Optimization problematic due to log-sum
- Assume straightforward maximization for complete data

$$\ln p(\mathbf{X}, \mathbf{Z}|\theta)$$

- Latent \mathbf{Z} is known only through $p(\mathbf{Z}|\mathbf{X}, \theta)$
- Consider expectation of complete data log-likelihood

An Alternative View of EM: Algorithm

1. **Initialization:** Choose initial set of parameters θ^{old}
2. **E-step:** use current parameters θ^{old} to compute $p(\mathbf{Z}|\mathbf{X}, \theta^{old})$ to find expected complete-data log-likelihood for general θ

$$\mathcal{Q}(\theta, \theta^{old}) = \sum_{\mathbf{Z}} p(\mathbf{Z}|\mathbf{X}, \theta^{old}) \ln p(\mathbf{X}, \mathbf{Z}|\theta)$$

3. **M-step:** determine θ^{new} by maximizing above formula

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

4. **Check convergence:** stop, or $\theta^{old} \leftarrow \theta^{new}$ and go to **E-step**

An Alternative View of EM: Gaussian Mixtures Revisited

- Complete-data (log-) likelihood

$$p(\mathbf{X}, \mathbf{Z} | \boldsymbol{\mu}, \boldsymbol{\Sigma}, \boldsymbol{\pi}) = \prod_{n=1}^N \prod_{k=1}^K \pi_k^{z_{nk}} \mathcal{N}(\mathbf{x}_n | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)^{z_{nk}}$$

$$\ln p(\mathbf{X}, \mathbf{Z} | \boldsymbol{\mu}, \boldsymbol{\Sigma}, \boldsymbol{\pi}) = \sum_{n=1}^N \sum_{k=1}^K z_{nk} (\ln \pi_k + \ln \mathcal{N}(\mathbf{x}_n | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k))$$

- Expectation

$$\mathbb{E}[z_{nk}] = \gamma(z_{nk})$$

$$\mathbb{E}[\ln p(\mathbf{X}, \mathbf{Z} | \boldsymbol{\mu}, \boldsymbol{\Sigma}, \boldsymbol{\pi})] = \sum_{n=1}^N \sum_{k=1}^K \gamma(z_{nk}) \{ \ln \pi_k + \ln \mathcal{N}(\mathbf{x}_n | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k) \}$$

Outline

- 1 K -Means Clustering
- 2 Mixtures of Gaussians
- 3 An Alternative View of EM
- 4 The EM algorithm in General**
- 5 Hidden Markov Model

The EM algorithm in General

- Let \mathbf{X} observed data, \mathbf{Z} latent variables, θ parameters
- Goal: maximize marginal log-likelihood of observed data

$$\ln p(\mathbf{X}|\theta) = \ln \left\{ \sum_{\mathbf{Z}} p(\mathbf{X}, \mathbf{Z}|\theta) \right\}$$

- Maximization of $p(\mathbf{X}, \mathbf{Z}|\theta)$ is simple, but difficult for $p(\mathbf{X}|\theta)$
- Given any $q(\mathbf{Z})$, we decompose the data log-likelihood

$$\ln p(\mathbf{X}|\theta) = \mathcal{L}(q, \theta) + KL(q(\mathbf{Z})||p(\mathbf{Z}|\mathbf{X}, \theta))$$

$$\mathcal{L}(q, \theta) = \sum_{\mathbf{Z}} q(\mathbf{Z}) \ln \frac{p(\mathbf{X}, \mathbf{Z}|\theta)}{q(\mathbf{Z})}$$

$$KL(q(\mathbf{Z})||p(\mathbf{Z}|\mathbf{X}, \theta)) = - \sum_{\mathbf{Z}} q(\mathbf{Z}) \ln \frac{p(\mathbf{Z}|\mathbf{X}, \theta)}{q(\mathbf{Z})} \geq 0^1$$

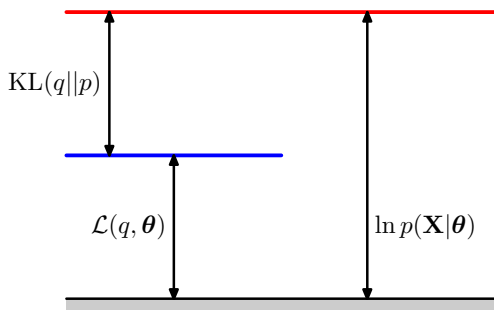
¹Need Verification

The EM algorithm in General: The EM Bound

- $\mathcal{L}(q, \theta)$ is a **lower bound on the data log-likelihood**

$$\mathcal{L}(q, \theta) = \ln p(\mathbf{X}|\theta) - KL(q(\mathbf{Z})||p(\mathbf{Z}|\mathbf{X}, \theta)) \leq \ln p(\mathbf{X}|\theta)$$

- **The EM algorithm performs coordinate ascent on \mathcal{L}**
 1. E-step: maximizes \mathcal{L} w.r.t. q for fixed θ
 2. M-step: maximizes \mathcal{L} w.r.t. θ for fixed q

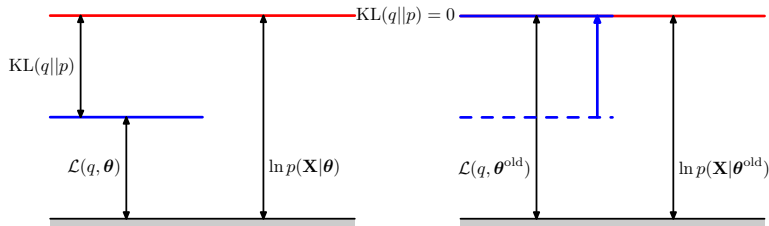


The EM algorithm in General: The E-step

- E-step: maximizes \mathcal{L} w.r.t. q for fixed θ

$$\mathcal{L}(q, \theta) = \ln p(\mathbf{X}|\theta) - KL(q(\mathbf{Z})||p(\mathbf{Z}|\mathbf{X}, \theta))$$

- \mathcal{L} maximized for $q(\mathbf{Z}) \leftarrow p(\mathbf{Z}|\mathbf{X}, \theta)$

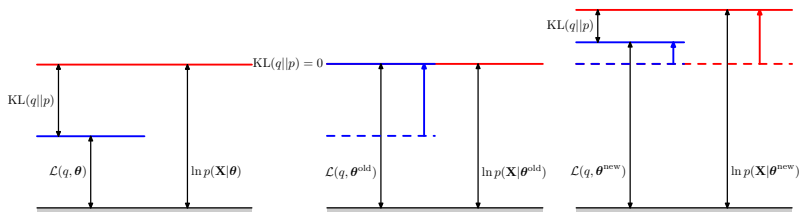


The EM algorithm in General: The M-step

- M-step: maximizes $\mathcal{L}(q, \boldsymbol{\theta})$ w.r.t. $\boldsymbol{\theta}$ for fixed q

$$\mathcal{L}(q, \boldsymbol{\theta}) = \sum_{\mathbf{Z}} q(\mathbf{Z}) \ln p(\mathbf{X}, \mathbf{Z} | \boldsymbol{\theta}) - \sum_{\mathbf{Z}} q(\mathbf{Z}) \ln q(\mathbf{Z})$$

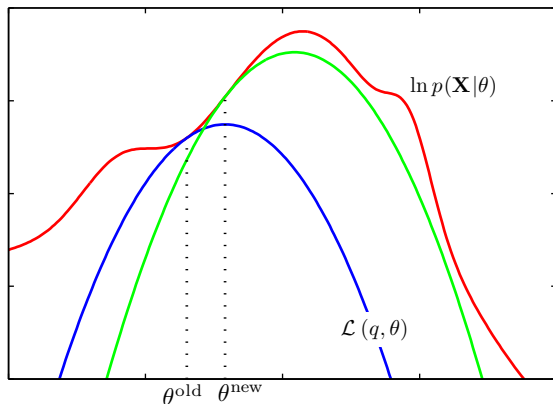
- \mathcal{L} maximized for $\boldsymbol{\theta} = \arg \max_{\boldsymbol{\theta}} \sum_{\mathbf{Z}} q(\mathbf{Z}) \ln p(\mathbf{X}, \mathbf{Z} | \boldsymbol{\theta})$



The EM algorithm in General: Picture in Parameter Space

E-step resets bound $\mathcal{L}(q, \theta)$ on $\ln p(\mathbf{X}|\theta)$ at $\theta = \theta^{old}$, it is

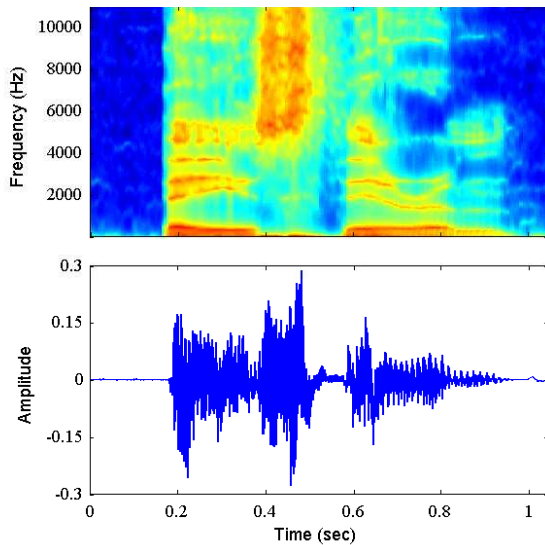
- tight as $\theta = \theta^{old}$
- tangential at $\theta = \theta^{old}$
- convex in θ for exponential family mixture components



Outline

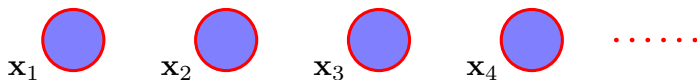
- 1 K -Means Clustering
- 2 Mixtures of Gaussians
- 3 An Alternative View of EM
- 4 The EM algorithm in General
- 5 Hidden Markov Model**

Motivation



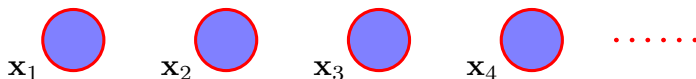
Markov Model

- Simplest way: model as independent:

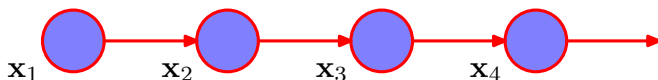


Markov Model

- Simplest way: model as independent:

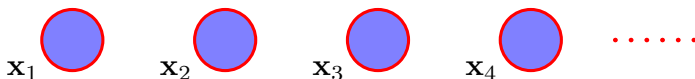


- Better to link observations, e.g. first-order Markov model, condition on previous observation

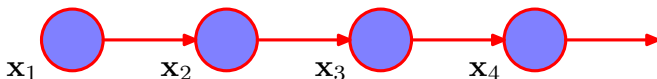


Markov Model

- Simplest way: model as independent:



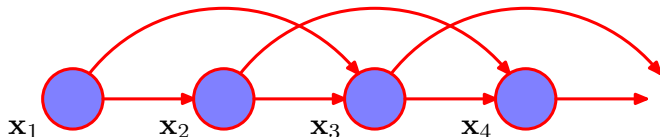
- Better to link observations, e.g. first-order Markov model, condition on previous observation



- Joint distribution

$$\begin{aligned} p(\mathbf{x}_1, \dots, \mathbf{x}_N) &= \prod_{n=1}^N p(\mathbf{x}_n | \mathbf{x}_1, \dots, \mathbf{x}_{n-1}) \\ &= p(\mathbf{x}_1) \prod_{n=2}^N p(\mathbf{x}_n | \mathbf{x}_{n-1}) \end{aligned}$$

Exercise 1



Q: Show that second-order Markov chain described by the joint distribution

$$p(\mathbf{x}_1, \dots, \mathbf{x}_N) = p(\mathbf{x}_1)p(\mathbf{x}_2|\mathbf{x}_1) \prod_{n=3}^N p(\mathbf{x}_n|\mathbf{x}_{n-1}, \mathbf{x}_{n-2})$$

satisfies the conditional independence property

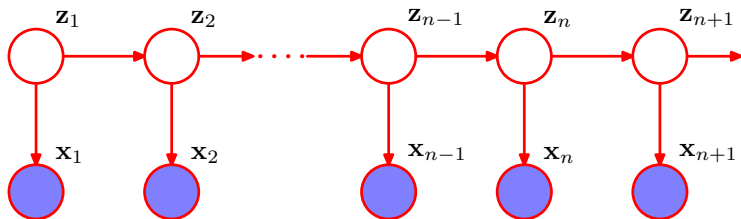
$$p(\mathbf{x}_n|\mathbf{x}_1, \dots, \mathbf{x}_{n-1}) = p(\mathbf{x}_n|\mathbf{x}_{n-1}, \mathbf{x}_{n-2})$$

Hidden Markov Model

- Introduce additional latent variables to permit a rich class of models to be constructed out of simple components

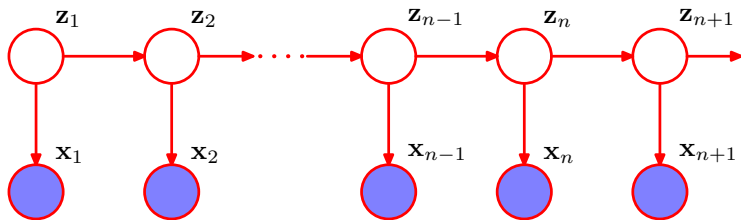
Hidden Markov Model

- Introduce additional latent variables to permit a rich class of models to be constructed out of simple components
- For each observation x_n , introduce a corresponding latent variable z_n



Hidden Markov Model

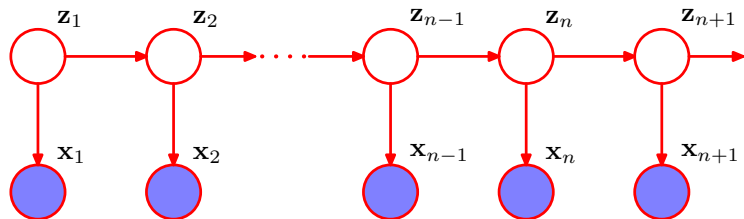
- Introduce additional latent variables to permit a rich class of models to be constructed out of simple components
- For each observation \mathbf{x}_n , introduce a corresponding latent variable \mathbf{z}_n



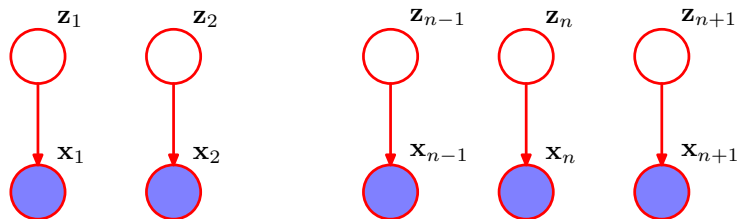
- Joint distribution

$$p(\mathbf{x}_1, \dots, \mathbf{x}_N, \mathbf{z}_1, \dots, \mathbf{z}_N) = p(\mathbf{z}_1) \left[\prod_{n=2}^N p(\mathbf{z}_n | \mathbf{z}_{n-1}) \right] \prod_{n=1}^N p(\mathbf{x}_n | \mathbf{z}_n)$$

Relationship with Mixture Model

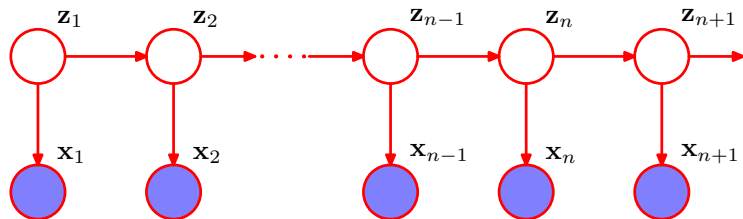


Relationship with Mixture Model



- Examining a single time slice of the model, it corresponds to a mixture distribution with component densities given by $p(x|z)$

Relationship with Mixture Model



- Examining a single time slice of the model, it corresponds to a mixture distribution with component densities given by $p(x|z)$
- HMM can be interpreted as an extension of a mixture model
- The choice of mixture component for each observation is not selected independently but **depends on the choice of component for the previous observation**

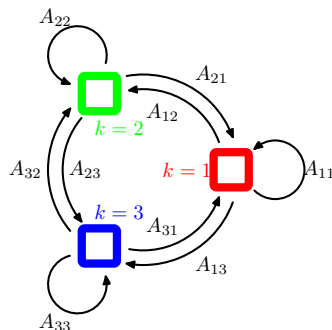
Transition Probabilities

- Elements of transition matrix

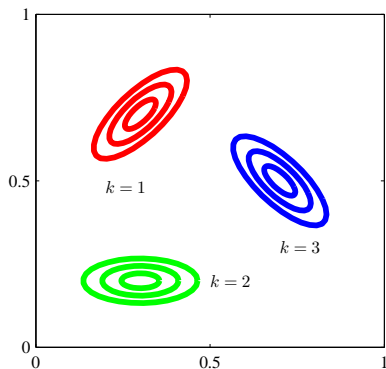
$$A_{jk} \equiv p(z_{nk} = 1 | z_{n-1,j} = 1)$$

- Constraints: $0 \leq A_{jk} \leq 1$ with $\sum_k A_{jk} = 1$
- Conditional distribution

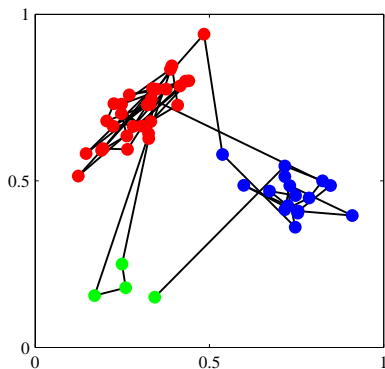
$$p(\mathbf{z}_n | \mathbf{z}_{n-1}, \mathbf{A}) = \prod_{k=1}^K \prod_{j=1}^K A_{jk}^{z_{n-1,j}, z_{nk}}$$



Sampling from HMMs



Emission Distribution



Successive Observations

Maximum Likelihood for HMM

- **E-Step:** forward-backward algorithm

$$\gamma(z_{tk}) = \mathbb{E}[z_{tk}] = \sum_{\mathbf{z}} \gamma(\mathbf{z}) z_{tk}$$

$$\xi(z_{t-1,i}, z_{tj}) = \mathbb{E}[z_{t-1,i}, z_{tj}] = \sum_{\mathbf{z}} \gamma(\mathbf{z}) z_{t-1,i} z_{tj}$$

- **M-Step:** expected complete-data log likelihood function

$$\begin{aligned} Q(\boldsymbol{\theta}, \boldsymbol{\theta}^{old}) &= \sum_{\mathbf{Z}} p(\mathbf{Z}|\mathbf{X}, \boldsymbol{\theta}^{old}) \ln p(\mathbf{X}, \mathbf{Z}|\boldsymbol{\theta}) \\ &= \sum_{k=1}^K \gamma(z_{1k}) \ln \pi_k + \sum_{t=2}^T \sum_{i=1}^K \sum_{j=1}^K \xi(z_{t-1,i}, z_{tj}) \ln A_{ij} \\ &\quad + \sum_{t=1}^T \sum_{k=1}^K \gamma(z_{tk}) \ln p(\mathbf{x}_t|\phi_k) \end{aligned}$$

Exercise 2

Q: Verify the M-step equations

$$\pi_k = \frac{\gamma(z_{1k})}{\sum_{i=1}^K \gamma(z_{1i})}$$
$$A_{ij} = \frac{\sum_{t=2}^T \xi(z_{t-1,i}, z_{tj})}{\sum_{k=1}^K \sum_{t=2}^T \xi(z_{t-1,i}, z_{tk})}$$

by maximization of the expected complete-data log likelihood function, using appropriate **Lagrange multipliers** to enforce the summation constraints on the components of $\boldsymbol{\pi}$ and \boldsymbol{A} .