

[Spring, 2017]

Mixture Models and EM

Pattern Recognition (BRI623)



Heung-II Suk

hisuk@korea.ac.kr

<http://www.ku-milab.org>



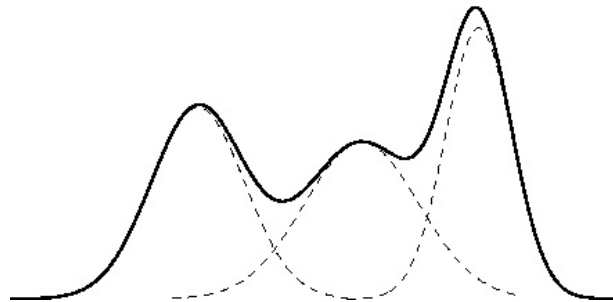
Department of Brain and Cognitive Engineering,
Korea University

Contents

- 1 Introduction
- 2 K -Means Clustering
- 3 Mixtures of Gaussians
- 4 An Alternative View of EM
- 5 EM Algorithm in General

Introduction

- Introduction of **latent variables** allows complicated distributions to be formed from simpler components
 - ▶ Mixture distributions (e.g., Gaussian mixture): discrete latent variables
 - ▶ Continuous latent variables (Chapter 12)

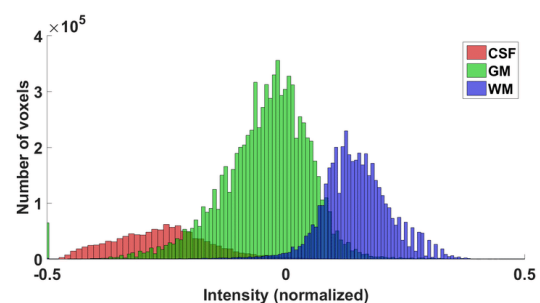


- Mixture models
 - ▶ provide a framework for building more complex probability distributions
 - ▶ used to cluster data (*clustering*)



2/95

Brain tissue segmentation



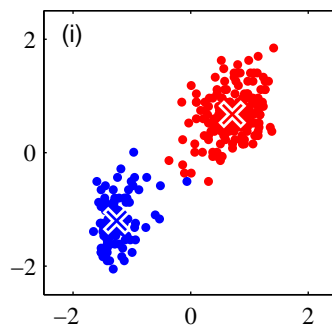
3/95

- K -means algorithm: non-probabilistic technique
 - ▶ Identifying groups, or clusters, of data points in a multidimensional space
- Latent variable view of mixture distributions
 - ▶ Discrete latent variables: interpreted as defining assignments of data points to specific components of the mixture
- Expectation-Maximization (EM) algorithm

K -Means Clustering

K-means Clustering

- Given a data set $\{\mathbf{x}_1, \dots, \mathbf{x}_N\}$ in a D -dimensional Euclidean space
- Goal: to partition the data set into some number K of clusters
 - ▶ Intuitively, comprising a group of data points whose inter-point distances are small compared with the distances to points outside of the cluster



6/95

[Problem definition]

- Prototype μ_k associated with the k -th cluster; center of the cluster
- Binary indicator variable $r_{nk} \in \{0, 1\}$
 - ▶ which of the K clusters the data point \mathbf{x}_n is assigned to
 - ▶ known as the '1-of- K coding scheme'
- Objective: to find (1) an assignment of data points to clusters $\{r_{nk}\}$ as well as (2) a set of vectors $\{\mu_k\}$, such that the sum of the squares of the distances of each data point to its closest vector μ_k , is a minimum.

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



7/95

K -means algorithm

Iterative and successive optimizations with respect to the $\{r_{nk}\}$ and the $\{\mu_k\}$

- ① Choose some initial values for the $\{\mu_k\}$
- ② Minimize J w.r.t. the $\{r_{nk}\}$, keeping the $\{\mu_k\}$ fixed (expectation)
 - ▶ Estimating the expected cluster
- ③ Minimize J w.r.t. the $\{\mu_k\}$, keeping the $\{r_{nk}\}$ fixed (maximization)
 - ▶ Maximizing the likelihood
- ④ Repeat this two-stage optimization until convergence



8/95

Determination of the $\{r_{nk}\}$ (expectation)

- J : a linear function of $\{r_{nk}\}$
 - ▶ Optimization: a closed form solution
- Independence among terms involving different n ,
 - ▶ Optimization for each n separately
 - ▶ By choosing r_{nk} to be 1 for whichever value of k gives the minimum value of $\|\mathbf{x}_n - \mu_k\|^2$

$$r_{nk} = \begin{cases} 1 & \text{if } k = \underset{j}{\operatorname{argmin}} \|\mathbf{x}_n - \mu_k\|^2 \\ 0 & \text{otherwise} \end{cases}$$



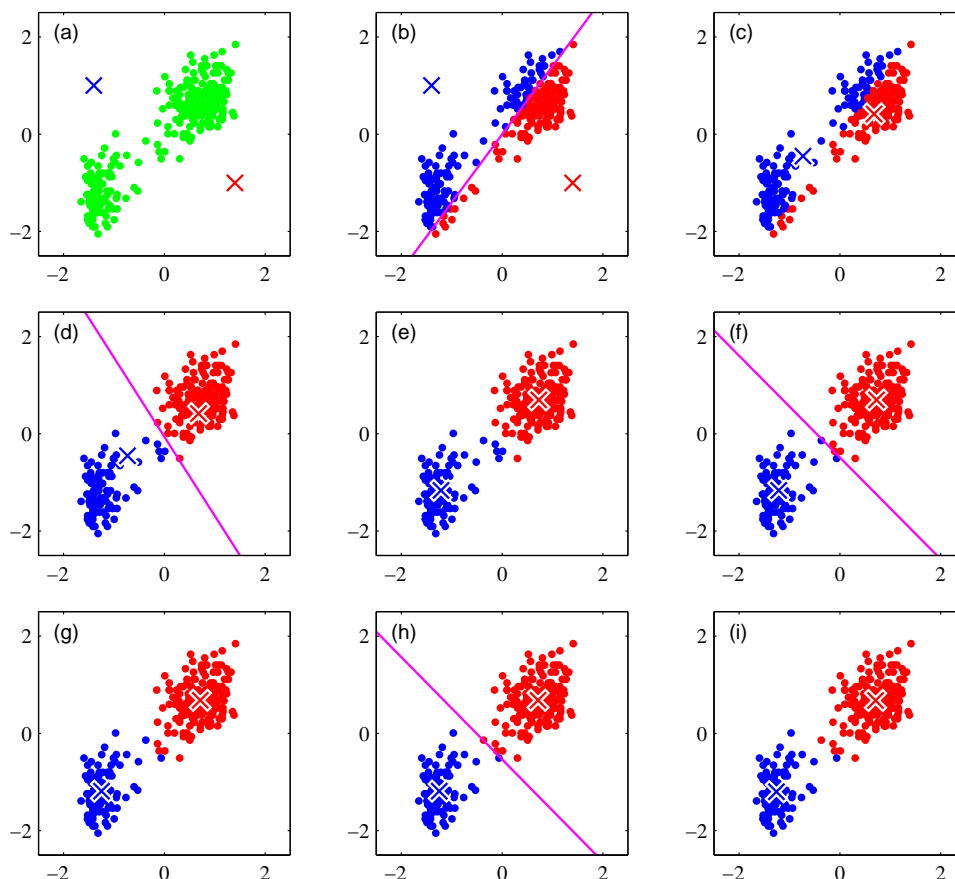
9/95

Optimization of the $\{\mu_k\}$ (maximization)

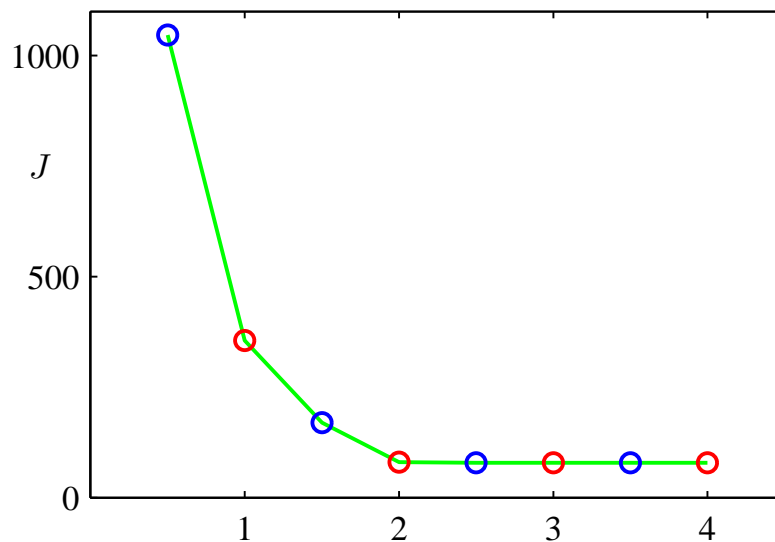
- J : a quadratic function of μ_k .
- Optimization: by setting its derivative w.r.t. μ_k to zero

$$\mu_k = \frac{\sum_n r_{nk} \mathbf{x}_n}{\sum_n r_{nk}}$$

- μ_k : mean of all of the data points \mathbf{x}_n assigned to cluster k



Plot of the cost function J after each step



assignment (expectation) step, updating (maximization) step

Non-decreasing after each stage in iterations

- Poor initial values for cluster centers
 - ▶ Several steps are involved for convergence
- Better initialization is to choose μ_k to be a random subset of K data points
- K -means algorithm itself is often used to initialize the parameters in a Gaussian mixture model before applying the EM algorithm

Implementation of the K -means algorithm

- Direct implementation can be relatively slow
 - ▶ in each expectation step, need to compute the Euclidean distance between every prototype vector and every data point

$$\|\mathbf{x}_n - \boldsymbol{\mu}_k\|^2$$

- Speeding up
 - ▶ Precomputing a data structure (e.g., tree) such that nearby points are in the same subtree (Ramasubramanian and Paliwal, 1990; Moore, 2000)
 - ▶ Making use of the triangle inequality for distances, thereby avoiding unnecessary distance calculations (Hodgson, 1998; Elkan, 2003)



14/95

Online stochastic algorithm (MacQueen, 1967)

- By applying the Robbins-Monro procedure (Chapter 2.3.4) to the problem of finding the roots of the regression function

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

given by the derivatives of J w.r.t. $\boldsymbol{\mu}_k$

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

- ▶ η_n : learning rate parameter, typically made to decrease monotonically as more data points are considered



15/95

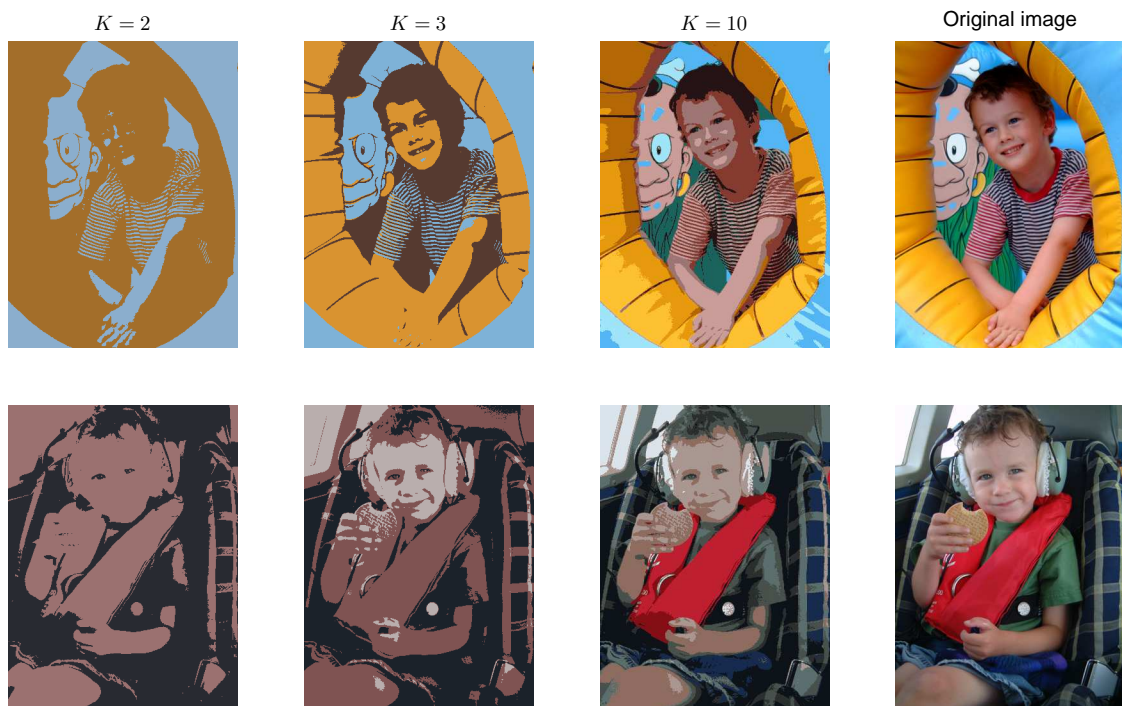
Application of K -Means

Image segmentation

- Partition an image into regions
 - ▶ each of which has homogeneous visual appearance
 - ▶ or corresponds to objects
 - ▶ or parts of objects
- Each pixel is a point in $[R, G, B]$ space
- K -means clustering is used with a palette of K colors
- Methods does not take into account proximity of different pixels.



18/95



19/95

Lossy data compression

- Accept some errors in the reconstruction in return for higher levels of compression than can be achieved in the lossless case
- For each of the N data points, we store only the identity k of the cluster to which it is assigned.
- Also store the values of the K cluster centers μ_k , where $K \ll N$
- known as *vector quantization*; $\{\mu_k\}$ called *code-book vectors*



20/95

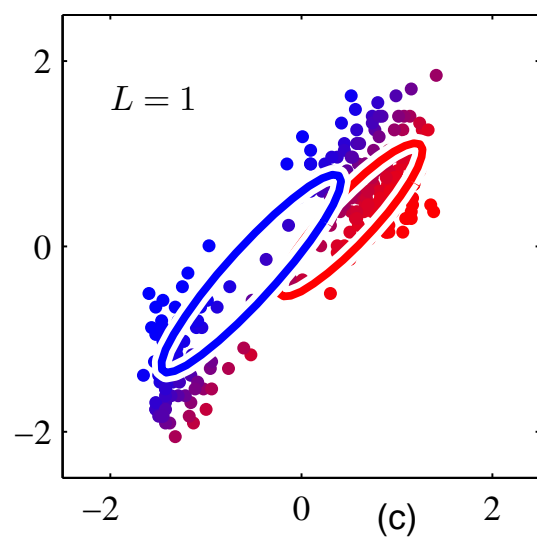
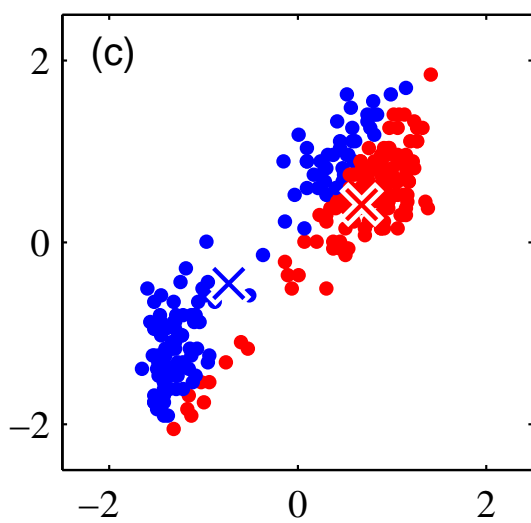
Limitation of K -means

- Hard assignment: every data point is assigned uniquely to one and only one cluster
- A point may lie roughly midway between cluster centers.
- A *probabilistic approach* will have a 'soft' assignment of data points to clusters in a way that reflects the **level of uncertainty** over the most appropriate assignment



21/95

Mixtures of Gaussians

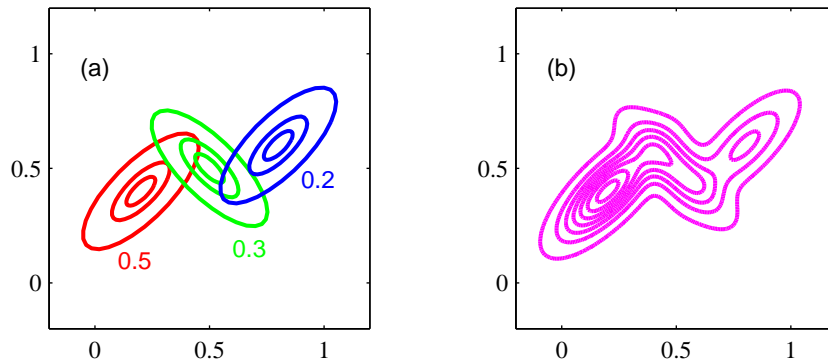


Gaussian Mixture Model (GMM)

A simple linear superposition of Gaussian components

- Providing a richer class of density models than the single Gaussian

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



Formulation of Gaussian mixtures in terms of **discrete latent variables**.

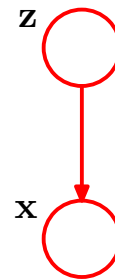
- Introduce a K -dimensional binary random variable \mathbf{z} having a 1-of- K coding scheme
 - ▶ a particular element z_k equal to 1 and all the other elements equal to 0

$$\sum_{k=1}^K z_k = 1$$

- K possible states for the vector \mathbf{z} according to which element is nonzero

- Define the joint distribution

$$p(\mathbf{x}, \mathbf{z}) = \underbrace{p(\mathbf{x}|\mathbf{z})}_{\text{conditional}} \underbrace{p(\mathbf{z})}_{\text{marginal}}$$



Graphical representation of a mixture model

► Probabilistic Graphical Models

- Marginal distribution over \mathbf{z} : $p(\mathbf{z})$
 - Specified in terms of the mixing coefficients $\{\pi_k\}$

$$p(z_k = 1) = \pi_k \quad 0 \leq \pi_k \leq 1, \sum_{k=1}^K \pi_k = 1$$

- Due to a 1-of- K representation

$$p(\mathbf{z}) = \prod_{k=1}^K \pi_k^{z_k}$$

- Conditional distribution of \mathbf{x} given a particular value for \mathbf{z}

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

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

- Marginal distribution of \mathbf{x} : $p(\mathbf{x})$

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

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

- Given several observations $\mathbf{x}_1, \dots, \mathbf{x}_N$, for every observed data point \mathbf{x}_n , a corresponding latent variable \mathbf{z}_n
- Instead of the marginal distribution $p(\mathbf{x})$, **we work with the joint distribution $p(\mathbf{x}, \mathbf{z})$**
 - ▶ Leading to significant simplifications
 - ▶ Most notably through the introduction of the Expectation-Maximization (EM) algorithm

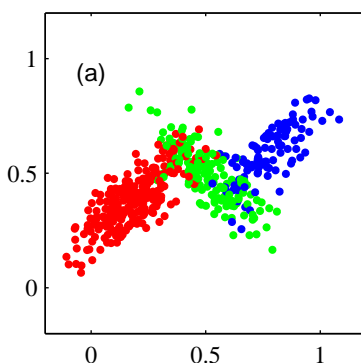
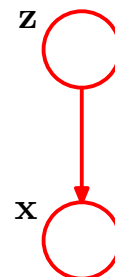
- Conditional probability of \mathbf{z} given \mathbf{x}

$$\begin{aligned}
 p(z_k = 1|\mathbf{x}) &= \frac{p(z_k = 1)p(\mathbf{x}|z_k = 1)}{\sum_{j=1}^K p(z_j = 1)p(\mathbf{x}|z_j = 1)} \\
 &= \frac{\pi_k \mathcal{N}(\mathbf{x}|\boldsymbol{\mu}_k, \Sigma_k)}{\sum_{j=1}^K \pi_j \mathcal{N}(\mathbf{x}|\boldsymbol{\mu}_j, \Sigma_j)} \equiv \gamma(z_k)
 \end{aligned}$$

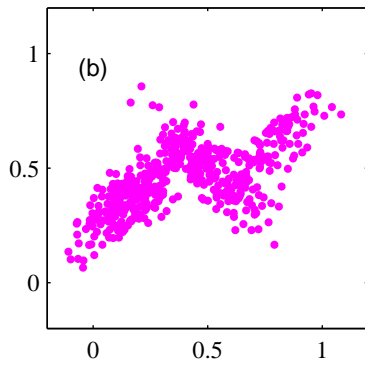
- View π_k as the prior probability of $z_k = 1$
- $\gamma(z_k)$ as the posterior probability once we have observed \mathbf{x}
 - ▶ also viewed as the *responsibility* that component k takes for 'explaining' the observation \mathbf{x}

Ancestral sampling to generate random samples distributed according to the Gaussian mixture model

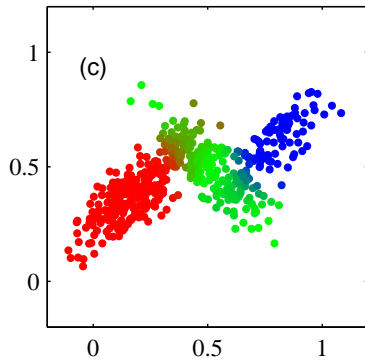
- 1 Generate a value for \mathbf{z} , denoted as $\hat{\mathbf{z}}$, from the marginal distribution $p(\mathbf{z})$
- 2 Generate a value for \mathbf{x} from the conditional distribution $p(\mathbf{x}|\hat{\mathbf{z}})$



(a) Samples from $p(\mathbf{x}, \mathbf{z})$ are plotted according to value of \mathbf{x} and colored with value of \mathbf{z}



(b) Samples from marginal distribution $p(\mathbf{x})$ obtained by ignoring values of \mathbf{z}



(c) Representing the value of the responsibilities $\gamma(z_{nk})$ associated with data point \mathbf{x}_n by plotting the corresponding point using proportions of red, blue, and green ink given by $\gamma(z_{nk})$ for $k = 1, 2, 3$, respectively

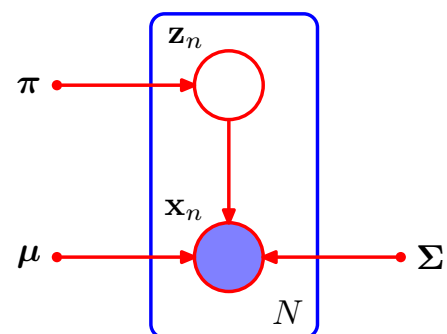
Maximum Likelihood for GMM

Given a set of N i.i.d. observations $\{\mathbf{x}_n\}_{n=1}^N$, model this data using a mixture of Gaussians

$$\mathbf{X} = \begin{bmatrix} \mathbf{x}_1^\top; \cdots; \mathbf{x}_N^\top \end{bmatrix} \in \mathbb{R}^{N \times D}$$

$$\mathbf{Z} = \begin{bmatrix} \mathbf{z}_1^\top; \cdots; \mathbf{z}_N^\top \end{bmatrix} \in \mathbb{R}^{N \times K}$$

$$\mathbf{z}_n \in \{0, 1\}^K, \quad \sum_{k=1}^K z_{nk} = 1$$



Goal: to estimate the three sets of parameters

$$\{\pi_k\}_{k=1}^K, \quad \{\boldsymbol{\mu}_k\}_{k=1}^K, \quad \{\Sigma_k\}_{k=1}^K$$

Likelihood function

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

Log of the likelihood function

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



34/95

- A more complex problem than for the case of a single Gaussian
 - ▶ Presence of the summation over k that appears inside the logarithm
 - ▶ No longer obtain a closed form solution
- Iterative approaches
 - ▶ Gradient-based optimization techniques (Fletcher, 1987; Nocedal and Wright, 1999; Bishop and Nabney, 2008)
 - ▶ **EM algorithm**
 - Broad applicability
 - Foundations for a discussion of [variational inference](#) techniques



35/95

Two significant problems associated with the maximum likelihood framework applied to Gaussian mixture models

- Presence of singularities
- Problem of identifiability (Casellan and Berger, 2002)

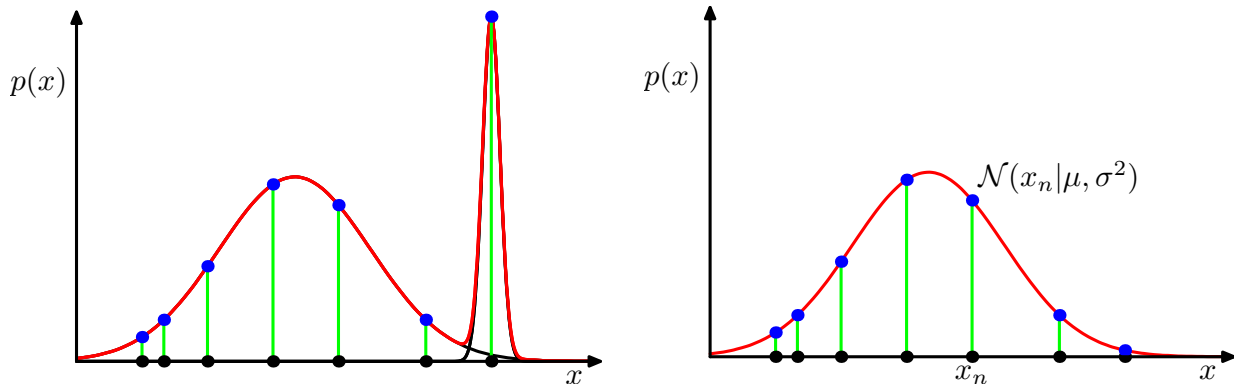
Singularities with Gaussian mixtures

- Consider a Gaussian mixture
 - ▶ Components have covariance matrices $\Sigma_k = \sigma_k^2 \mathbf{I}$
- Suppose that the j -th component has its mean μ_j exactly equal to one of the data points \mathbf{x}_n

$$\mathcal{N}(\mathbf{x}_n | \mathbf{x}_n, \sigma_j^2 \mathbf{I}) = \frac{1}{(2\pi)^{1/2}} \frac{1}{\sigma_j}$$

- As $\sigma_j \rightarrow 0$, $\mathcal{N}(\mathbf{x}_n | \mathbf{x}_n, \sigma_j^2 \mathbf{I})$ goes to infinity

- Thus, the maximization of the log likelihood function is not well-posed.
 - ▶ Such singularities will occur whenever one of the Gaussian components 'collapses' onto a specific data point.
 - ▶ Not arise in the case of a single Gaussian distribution



- If a single Gaussian collapses onto a data point, it will contribute multiplicative factors to the likelihood function arising from the other data points and these factors will go to zero exponentially fast, giving an overall likelihood that goes to zero rather than infinity.
- However, once we have (at least) two components in the mixture, one of the components can have a finite variance and therefore assign finite probability to all of the data points while the other component can shrink onto one specific data point and thereby contribute an ever increasing additive value to the log likelihood. (*overfitting*)
- This difficulty does not occur if we adopt a *Bayesian* approach.

- For the moment, however, we simply note that in applying maximum likelihood to Gaussian mixture models we must take steps to avoid finding such pathological solutions and instead seek local maxima of the likelihood function that are well behaved.
- We can hope to avoid the singularities by using suitable heuristics, for instance by detecting when a Gaussian component is collapsing and **resetting its mean to a randomly chosen value while also resetting its covariance to some large value**, and then continuing with the optimization.

Problem of identifiability

- For any given maximum likelihood solution, a K -component mixture will have a total of $K!$ equivalent solutions
 - ▶ corresponding to the $K!$ ways of assigning K sets of parameters to K components
- In other words, for any given (nondegenerate) point in the space of parameter values there will be a further $K! - 1$ additional points all of which give rise to exactly the same distribution.
- **An important issue when interpreting the parameter values discovered by a model**
- However, for the purposes of finding a good density model, it is irrelevant because any of the equivalent solutions is as good as any other.

EM for Gaussian Mixtures

- Expectation-Maximization (EM) algorithm
 - ▶ a method for finding *maximum likelihood* solutions for models with latent variables (Dempster *et al.*, 1977; McLachlan and Krishnan, 1997)

- In the context of the Gaussian mixture model

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

- Taking derivatives w.r.t. $\boldsymbol{\mu}_k$, $\boldsymbol{\Sigma}_k$, π_k and setting to zero



42/95

Maximizing $\ln p(\mathbf{X}|\boldsymbol{\pi}, \boldsymbol{\mu}, \boldsymbol{\Sigma})$ w.r.t. $\boldsymbol{\mu}_k$

- Taking derivatives w.r.t. $\boldsymbol{\mu}_k$ and setting to zero

$$\begin{aligned} 0 &= - \sum_{n=1}^N \underbrace{\frac{\pi_k \mathcal{N}(\mathbf{x}_n | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)}{\sum_{j=1}^K \pi_j \mathcal{N}(\mathbf{x}_n | \boldsymbol{\mu}_j, \boldsymbol{\Sigma}_j)}}_{p(z_k=1|\mathbf{x}_n)} \boldsymbol{\Sigma}_k^{-1} (\mathbf{x}_n - \boldsymbol{\mu}_k) \\ &= - \sum_{n=1}^N \gamma(z_k) \boldsymbol{\Sigma}_k^{-1} (\mathbf{x}_n - \boldsymbol{\mu}_k) \end{aligned}$$

- Multiplying both sides by $\boldsymbol{\Sigma}_k$ (assuming non-singular)

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

where $N_k = \sum_{n=1}^N \gamma(z_k)$: effective number of points assigned to cluster k



43/95

Maximizing $\ln p(\mathbf{X}|\boldsymbol{\pi}, \boldsymbol{\mu}, \boldsymbol{\Sigma})$ w.r.t. Σ_k

- Taking derivatives w.r.t. Σ_k and setting to zero
 - Making use of the result for the maximum likelihood solution for the covariance matrix of a single Gaussian

$$\Sigma_k = \frac{1}{N_k} \sum_{n=1}^N \gamma(z_k) (\mathbf{x}_n - \boldsymbol{\mu}_k) (\mathbf{x}_n - \boldsymbol{\mu}_k)^\top$$

where $N_k = \sum_{n=1}^N \gamma(z_k)$: effective number of points assigned to cluster k



44/95

Maximizing $\ln p(\mathbf{X}|\boldsymbol{\pi}, \boldsymbol{\mu}, \boldsymbol{\Sigma})$ w.r.t. π_k with a constraint $\sum_{k=1}^K \pi_k = 1$

- Constrained optimization: Lagrangian multiplier

$$\ln p(\mathbf{X}|\boldsymbol{\pi}, \boldsymbol{\mu}, \boldsymbol{\Sigma}) + \lambda \left(\sum_{k=1}^K \pi_k - 1 \right)$$

- Taking derivatives w.r.t. π_k and setting to zero

$$0 = \sum_{n=1}^N \frac{\mathcal{N}(\mathbf{x}_n|\boldsymbol{\mu}_k, \Sigma_k)}{\sum_{j=1}^K \pi_j \mathcal{N}(\mathbf{x}_n|\boldsymbol{\mu}_j, \Sigma_j)} + \lambda$$

- Multiplying both sides by π_k and sum over k : $\lambda = -N$

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

- Mixing coefficient for the k -th component is given by the average responsibility which that component takes for explaining the data points



45/95

- The results for μ_k , Σ_k , and π_k are **not a closed-form solution**
 - ∴ responsibilities $\gamma(z_k)$ depend on those parameters in a complex way

$$\frac{\pi_k \mathcal{N}(\mathbf{x}_n | \mu_k, \Sigma_k)}{\sum_{j=1}^K \pi_j \mathcal{N}(\mathbf{x}_n | \mu_j, \Sigma_j)}$$

- However, the results suggest a simple **iterative scheme** for finding a solution to the maximum likelihood problem.

EM algorithm for GMM

- 1 Choose some initial values for the means $\{\mu_k\}$, covariances $\{\Sigma_k\}$, and mixing coefficients $\{\pi_k\}$
- 2 Alternate between the following two updates until convergence
 - (**E-step**) Use the current values for the parameters to **evaluate the posterior probabilities, or responsibilities**

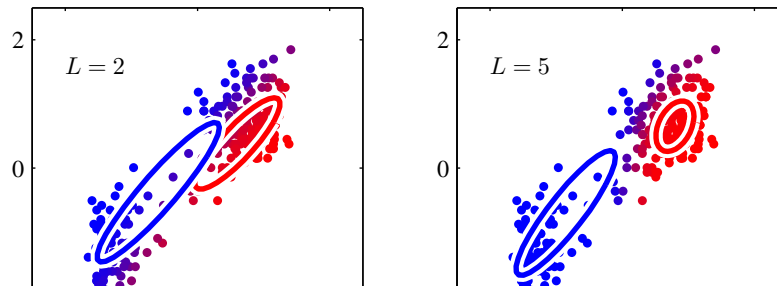
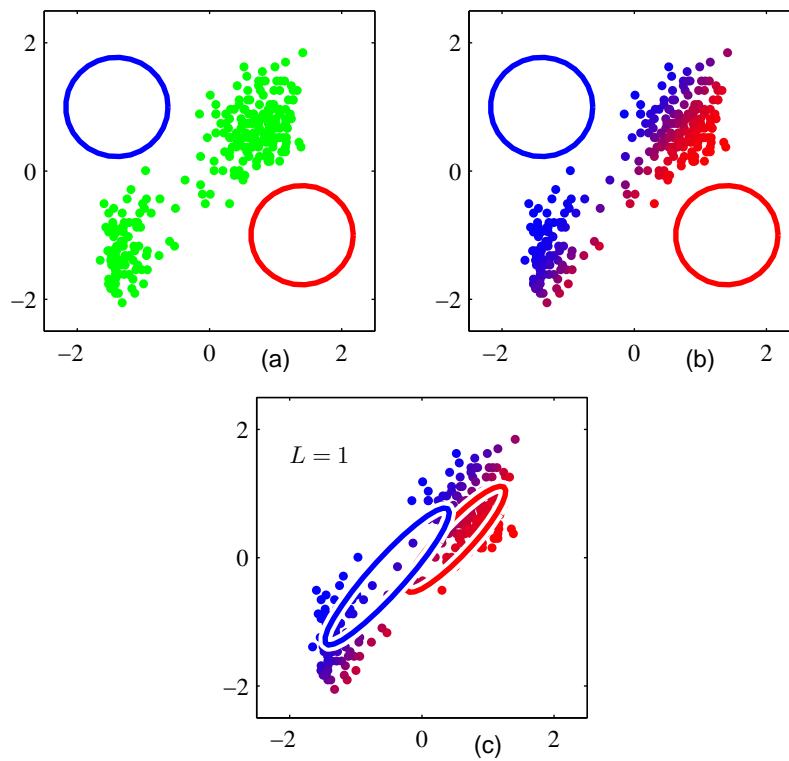
$$p(z_k = 1 | \mathbf{x}) = \gamma(z_k) = \frac{\pi_k \mathcal{N}(\mathbf{x}_n | \mu_k, \Sigma_k)}{\sum_{j=1}^K \pi_j \mathcal{N}(\mathbf{x}_n | \mu_j, \Sigma_j)}$$

- (**M-step**) Use these probabilities to **re-estimate the means, covariances, and mixing coefficients**

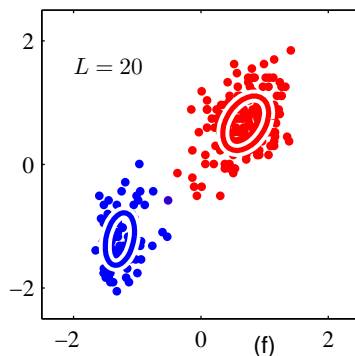
$$\mu_k^{\text{new}} = \frac{1}{N_k} \sum_{n=1}^N \gamma(z_k) \mathbf{x}_n$$

$$\Sigma_k^{\text{new}} = \frac{1}{N_k} \sum_{n=1}^N \gamma(z_k) \left(\mathbf{x}_n - \mu_k^{\text{new}} \right) \left(\mathbf{x}_n - \mu_k^{\text{new}} \right)^{\top}$$

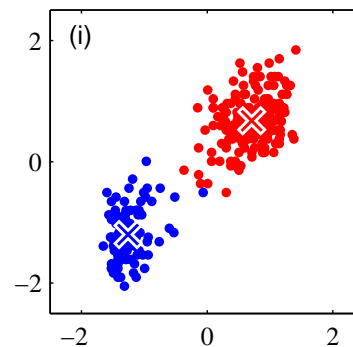
$$\pi_k^{\text{new}} = \frac{N_k}{N} \quad \text{where } N_k = \sum_{n=1}^N \gamma(z_k)$$



EM result



K -means result



- EM takes many more iterations to reach convergence.
- Each cycle requires significantly more computation.
- Common to run K -means first in order to find a suitable initialization
 - ▶ Covariance matrices: sample covariances of the clusters
 - ▶ Mixing coefficients: fractions of data points assigned to the respective clusters
- EM is not guaranteed to find the global maximum of the log likelihood function.

An Alternative View of EM

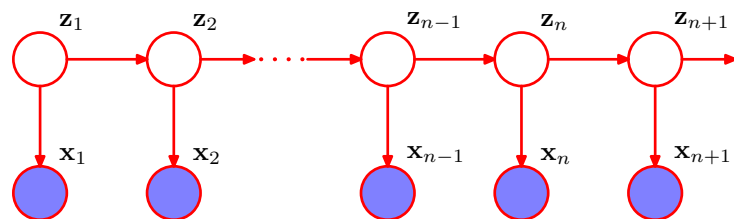


50/95

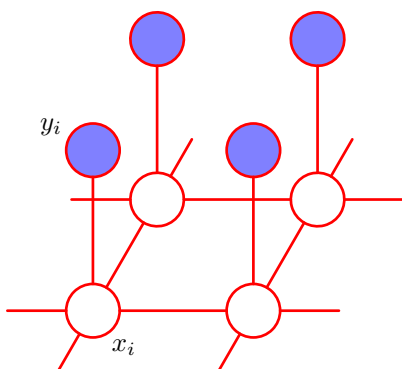
Latent Variables



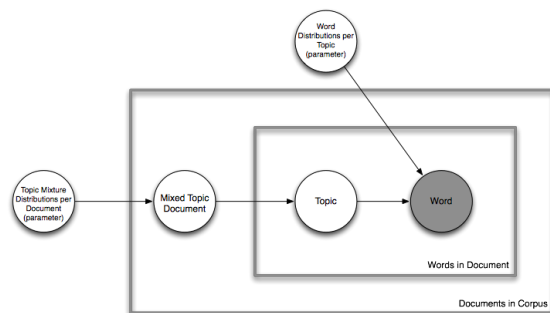
Mixture model



Hidden Markov Model (HMM) [Rabiner, 1989]



Markov Random Field (MRF)



Latent Dirichlet Allocation (LDA) [Blei *et al.*, 2003]



51/95

Likelihood Function with Latent Variables

Goal of EM: to find maximum likelihood solutions for models with latent variables

- Observed data: $\mathbf{X} = [\mathbf{x}_1^T; \dots; \mathbf{x}_n^T; \dots; \mathbf{x}_N^T]$
- Latent variables: $\mathbf{Z} = [\mathbf{z}_1^T; \dots; \mathbf{z}_n^T; \dots; \mathbf{z}_N^T]$
- Set of all model parameters: θ
- Log-likelihood function

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

$$\ln p(\mathbf{X}|\theta) = \ln \left\{ \int p(\mathbf{X}, \mathbf{Z}|\theta) d\mathbf{Z} \right\} \quad (\mathbf{Z} : \text{continuous})$$



52/95

Complete data $\{\mathbf{X}, \mathbf{Z}\}$

- For each observation in \mathbf{X} , we know corresponding value of latent variable \mathbf{Z}
- Log-likelihood
$$\ln p(\mathbf{X}, \mathbf{Z}|\theta)$$

- Straightforward for maximization of the log-likelihood function

Incomplete data $\{\mathbf{X}\}$

- Unobservable variable \mathbf{Z}
- Log-likelihood

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

- Challenging for maximization of the log-likelihood function



53/95

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

- Summation over the latent variables inside the logarithm
- Preventing the logarithm from acting directly on the joint distribution
- Resulting in complicated expressions for the ML solution

Latent Variables in EM

- Our state of knowledge of the values of the latent variables in \mathbf{Z} is given only by the posterior distribution $p(\mathbf{Z}|\mathbf{X}, \theta)$: E-step of EM algorithm
- Since we cannot use the complete-data log likelihood, we consider instead its expected value under the posterior distribution of the latent variable

$$\mathbb{E} [\ln p(\mathbf{X}, \mathbf{Z}|\theta)] = \sum_{\mathbf{Z}} p(\mathbf{Z}|\mathbf{X}, \theta) \ln p(\mathbf{X}, \mathbf{Z}|\theta)$$

- Maximize this expectation: M-step of EM algorithm

- **E-step**: use the current parameter values θ^{old} to find the posterior distribution of the latent variables given by $p(\mathbf{Z}|\mathbf{X}, \theta^{\text{old}})$
- Use this posterior to find expectation of the complete-data log-likelihood evaluated for some general parameter value θ

$$\mathcal{Q}(\theta, \theta^{\text{old}}) \equiv \mathbb{E} \left[\ln p(\mathbf{X}, \mathbf{Z}|\theta) | \theta^{\text{old}} \right] = \sum_{\mathbf{Z}} p(\mathbf{Z}|\mathbf{X}, \theta^{\text{old}}) \ln p(\mathbf{X}, \mathbf{Z}|\theta)$$

- **M-step**: determine the revised parameter estimate θ^{new} by maximizing

$$\theta^{\text{new}} = \underset{\theta}{\operatorname{argmax}} \mathcal{Q}(\theta, \theta^{\text{old}})$$

- ▶ Since the logarithm acts directly on the joint distribution $p(\mathbf{X}, \mathbf{Z}|\theta)$, maximization will be tractable now.



56/95

General EM Algorithm

Given a joint distribution $p(\mathbf{X}, \mathbf{Z}|\theta)$ over observed variables \mathbf{X} and latent variables \mathbf{Z} , governed by parameters θ , the goal is to maximize the likelihood function $p(\mathbf{X}|\theta)$ with respect to θ .

- 1 Choose an initial setting for the parameters θ^{old}
- 2 **(E-step)** Evaluate $p(\mathbf{Z}|\mathbf{X}, \theta^{\text{old}})$
- 3 **(M-step)** Evaluate θ^{new} given by

$$\theta^{\text{new}} = \underset{\theta}{\operatorname{argmax}} \mathcal{Q}(\theta, \theta^{\text{old}})$$

$$\mathcal{Q}(\theta, \theta^{\text{old}}) \equiv \mathbb{E} \left[\ln p(\mathbf{X}, \mathbf{Z}|\theta) | \theta^{\text{old}} \right] = \sum_{\mathbf{Z}} p(\mathbf{Z}|\mathbf{X}, \theta^{\text{old}}) \ln p(\mathbf{X}, \mathbf{Z}|\theta)$$

- 4 If not converged, then let $\theta^{\text{old}} \leftarrow \theta^{\text{new}}$ and return to step 2.



57/95

EM for MAP Solutions

To find MAP solutions for models in which a prior $p(\theta)$ is defined over the parameters

- **(E-step)** Remains the same as in the ML case
- **(M-step)** Quantity to be maximized is given by

$$\mathcal{Q}(\theta, \theta^{\text{old}}) + \ln p(\theta)$$

- ▶ Suitable choices for the prior will remove the singularities

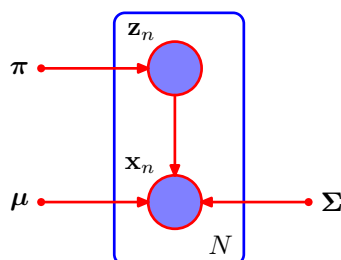


58/95

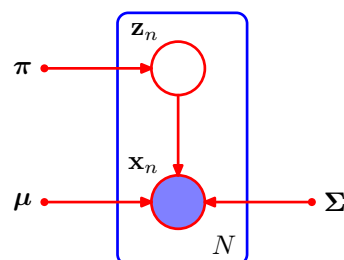
Gaussian Mixtures Revisited

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

Complete-data



Incomplete-data



60/95

For complete-data

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

- A sum of K independent contributions, one for each mixture component
- Maximization w.r.t. $\boldsymbol{\mu}_k$ or $\boldsymbol{\Sigma}_k$ is exactly as for a single Gaussian, but involving only the subset of data points 'assigned' to that component
- Maximization w.r.t. π_k : fractions of data points assigned to the corresponding components

For incomplete-data

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

- Due to unknown latent variables, we obtain its **expectation** w.r.t. the posterior distribution of latent variables



61/95

Latent variable view of EM to a Gaussian mixture model

- Log of the likelihood function

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

- Posterior distribution of latent variables \mathbf{Z}

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

- ▶ Factorizes over n so that under the posterior distribution the $\{\mathbf{z}_n\}$ are independent



62/95

- Expected value of the indicator variable z_{nk} under this posterior probability

$$\begin{aligned}\mathbb{E}[z_{nk}] &= \frac{\sum_{z_{nk}} z_{nk} [\pi_k \mathcal{N}(\mathbf{x}_n | \boldsymbol{\mu}_k, \Sigma_k)]^{z_{nk}}}{\sum_{z_{nj}} [\pi_j \mathcal{N}(\mathbf{x}_n | \boldsymbol{\mu}_j, \Sigma_j)]^{z_{nj}}} \\ &= \frac{\pi_k \mathcal{N}(\mathbf{x}_n | \boldsymbol{\mu}_k, \Sigma_k)}{\sum_{j=1}^K \pi_j \mathcal{N}(\mathbf{x}_n | \boldsymbol{\mu}_j, \Sigma_j)} = \gamma(z_{nk})\end{aligned}$$

- ▶ '*Responsibility*' of component k for data point \mathbf{x}_n

- Expected value of the complete-data log likelihood function

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

EM algorithm for a Gaussian mixture model

- 1 Choose some initial values for the parameters $\boldsymbol{\mu}^{\text{old}}$, $\boldsymbol{\Sigma}^{\text{old}}$, and $\boldsymbol{\pi}^{\text{old}}$
- 2 **(E-Step)** Evaluate the responsibilities $\gamma(z_{nk})$
- 3 **(M-Step)** Maximize $\mathbb{E}_{\mathbf{Z}} [\ln p(\mathbf{X}, \mathbf{Z} | \boldsymbol{\pi}, \boldsymbol{\mu}, \boldsymbol{\Sigma})]$ by keeping $\gamma(z_{nk})$ fixed to obtain $\boldsymbol{\mu}^{\text{new}}$, $\boldsymbol{\Sigma}^{\text{new}}$, and $\boldsymbol{\pi}^{\text{new}}$
- 4 Iterate **(E-Step)** and **(M-Step)** until convergence

Relation to K -Means

K -Means

- Finding code vectors that minimize the difference between \mathbf{x}_n and $\boldsymbol{\mu}_k$
- (E-step) Indication of belonging to a cluster: '*hard*' assignment
- (M-step) Updating means values

EM

- Look for density parameters that maximize the likelihood of samples
- (E-step) Probability of belonging to clusters: '*soft*' assignment
- (M-step) Updating parameters of densities

- K means algorithm does not estimate the covariances of the clusters but only the cluster means.
- *Elliptical K-means algorithm* [Sung and Poggio, 1994]: a hard-assignment version of the Gaussian mixture model with general covariance matrices



65/95

EM Algorithm in General



75/95

EM Algorithm in General

- Observed variables $\mathbf{X} = \{\mathbf{x}_n\}_{n=1}^N$ and hidden variables $\mathbf{Z} = \{\mathbf{z}_n\}_{n=1}^N$

- Goal: to maximize the likelihood function

$$p(\mathbf{X}|\theta) = \sum_{\mathbf{Z}} p(\mathbf{X}, \mathbf{Z}|\theta) \quad \text{or} \quad p(\mathbf{X}|\theta) = \int_{\mathbf{Z}} p(\mathbf{X}, \mathbf{Z}|\theta)$$

- Suppose that direct optimization of $p(\mathbf{X}|\theta)$ is difficult, but that optimization of the complete-data likelihood function $p(\mathbf{X}, \mathbf{Z}|\theta)$ is significantly easier.



76/95

- Introduce a distribution $q(\mathbf{Z})$ defined over the latent variables
- Then for any choice of $q(\mathbf{Z})$, the following decomposition holds

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

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

$$\mathcal{KL}(q||p) = - \int_{\mathbf{Z}} q(\mathbf{Z}) \ln \left\{ \frac{p(\mathbf{Z}|\mathbf{X}, \theta)}{q(\mathbf{Z})} \right\}$$

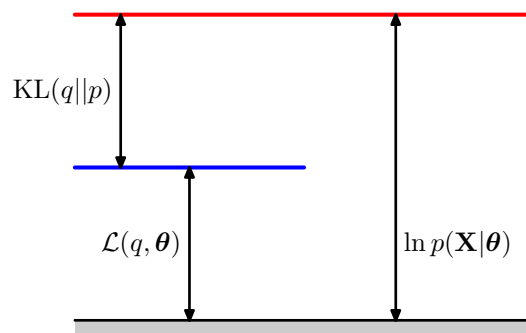
- ▶ $\mathcal{L}(q, \theta)$: a functional of the distribution $q(\mathbf{Z})$ and a function of the parameter θ



80/95

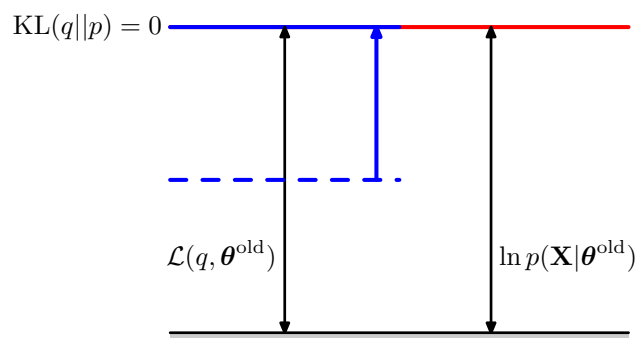
$$\ln p(\mathbf{X}|\theta) = \mathcal{L}(q, \theta) + \mathcal{KL}(q||p)$$

- Since $\mathcal{KL}(q||p) \geq 0$, $\mathcal{L}(q, \theta) \leq \ln p(\mathbf{X}|\theta)$
- $\mathcal{L}(q, \theta)$: a lower bound on $\ln p(\mathbf{X}|\theta)$

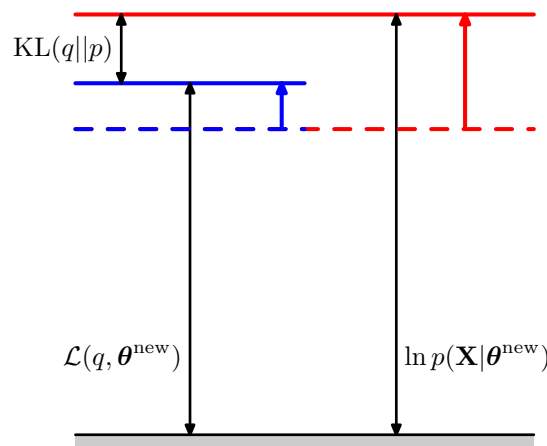


Suppose that the current value of the parameter vector is θ^{old} .

- E-step: the lower bound $\mathcal{L}(q, \theta)$ is maximized with respect to $q(\mathbf{Z})$, while holding θ^{old} fixed
 - ▶ The solution to this maximization problem is easily seen by noting that the value of $\ln p(\mathbf{X}|\theta^{\text{old}})$ does not depend on $q(\mathbf{Z})$
 - ▶ The largest value of $\mathcal{L}(q, \theta)$ will occur when the KL-divergence vanishes.
 - ▶ When $q(\mathbf{Z})$ is equal to the posterior distribution $p(\mathbf{Z}|\mathbf{X}, \theta^{\text{old}})$
 - ▶ The lower bound will equal the log likelihood.

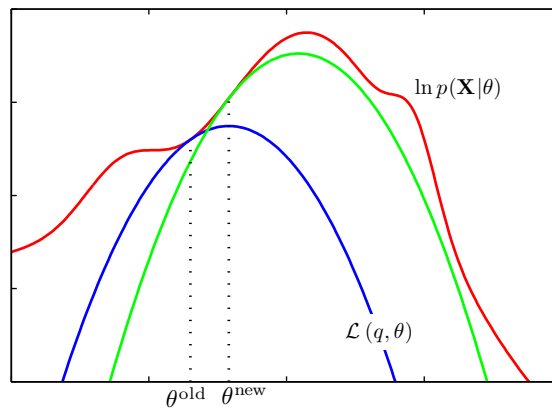


- M-step: the distribution $q(\mathbf{Z})$ is held fixed and the lower bound $\mathcal{L}(q, \theta)$ is maximized with respect to θ to give some new value θ^{new} .
 - ▶ Cause the lower bound \mathcal{L} to increase and the corresponding log likelihood function
 - ▶ Because the distribution q is determined using the old parameter values rather than the new values and is held fixed during the M-step, it will not equal the new posterior distribution $p(\mathbf{Z}|\mathbf{X}, \theta^{\text{new}})$, and hence there will be a nonzero KL divergence.



$$\begin{aligned}
 \mathcal{L}(q, \theta) &= \sum_{\mathbf{Z}} p(\mathbf{Z}|\mathbf{X}, \theta^{\text{old}}) \ln p(\mathbf{X}, \mathbf{Z}|\theta) - \sum_{\mathbf{Z}} p(\mathbf{Z}|\mathbf{X}, \theta^{\text{old}}) \ln p(\mathbf{X}, \mathbf{Z}|\theta^{\text{old}}) \\
 &= \underbrace{Q(\theta, \theta^{\text{old}})}_{\text{Entropy} \geq 0, \text{ const.}} + \underbrace{\mathcal{H}\{p(\mathbf{Z}|\mathbf{X}, \theta^{\text{old}})\}}_{\text{Entropy} \geq 0, \text{ const.}}
 \end{aligned}$$

- By substituting $q(\mathbf{Z}) = p(\mathbf{Z}|\mathbf{X}, \theta^{\text{new}})$, after the E-step, the lower bound takes the form
 - ▶ Const is simply the negative entropy of the q distribution and is therefore independent of θ
- In the M-step, the quantity that is being maximized is the expectation of the complete-data log likelihood



- **Red**: (incomplete data) log likelihood function to maximize
- In the E-step, we evaluate the posterior distribution over latent variables, which gives rise to a lower bound $\mathcal{L}(q, \theta)$ whose value equals the log likelihood at θ^{old} , so that both curves have the same gradient (**Blue**).
- In the M-step, the bound is maximized giving the value θ^{new} , which gives a larger value of log likelihood than θ^{old}
- The subsequent E-step then constructs a bound that is tangential at θ^{new} (**Green**).

EM for MAP

Using the EM algorithm to **maximize the posterior distribution** $p(\theta|\mathbf{X})$ for models in which we have introduced **a prior** $p(\theta)$ over the parameters

$$\begin{aligned}
 \ln p(\theta|\mathbf{X}) &= \ln p(\theta, \mathbf{X}) - \ln p(\mathbf{X}) \\
 &= \ln p(\mathbf{X}|\theta) + \ln p(\theta) - \ln p(\mathbf{X}) \\
 &= \underbrace{\mathcal{L}(q, \theta) + \mathcal{KL}(q||p)}_{=\ln p(\mathbf{X}|\theta)} + \ln p(\theta) - \ln p(\mathbf{X}) \\
 &\geq \mathcal{L}(q, \theta) + \ln p(\theta) - \underbrace{\ln p(\mathbf{X})}_{\text{const.}}
 \end{aligned}$$

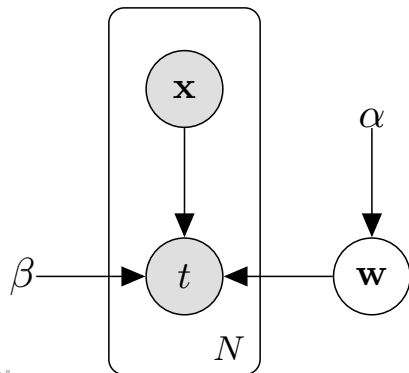
- E-step: same as for the MLE-EM (q only appears in $\mathcal{L}(q, \theta)$)
- M-step: modified through the introduction of the prior term $\ln p(\theta)$

EM for Bayesian Linear Regression

Given a data set $\mathcal{D} = \{\mathbf{x}_n, t_n\}_{n=1}^N$

$$t = \mathbf{w}^\top \mathbf{x} + \epsilon; \quad \epsilon \sim \mathcal{N}(0, \beta^{-1})$$

$$p(\mathbf{t}|\mathbf{X}, \mathbf{w}; \beta) = \prod_{n=1}^N p(t_n|\mathbf{x}_n, \mathbf{w}; \beta)$$



$$p(t|\mathbf{x}, \mathbf{w}; \beta) = \mathcal{N}(t; \mathbf{w}^\top \mathbf{x}, \beta^{-1})$$

$$p(\mathbf{w}; \alpha) = \mathcal{N}(\mathbf{w}; \mathbf{0}, \alpha^{-1} \mathbf{I})$$



87/95

E-Step: $(\theta = [\alpha, \beta])$

$$\begin{aligned} \mathcal{Q}(\theta, \theta^{\text{old}}) &= \int p(\mathbf{w}|\mathbf{t}; \alpha^{\text{old}}, \beta^{\text{old}}) \ln p(\mathbf{t}, \mathbf{w}; \alpha^{\text{old}}, \beta^{\text{old}}) d\mathbf{w} \\ &= \int p(\mathbf{w}|\mathbf{t}; \alpha^{\text{old}}, \beta^{\text{old}}) \{ \ln p(\mathbf{t}|\mathbf{w}; \alpha, \beta) \ln p(\mathbf{w}; \alpha, \beta) \} d\mathbf{w} \\ &= \frac{N}{2} \ln \beta - \frac{\beta}{2} \left(\|\mathbf{t} - \mathbf{X} \boldsymbol{\mu}^{\text{old}}\|^2 + \text{Tr}[\mathbf{X}^\top \boldsymbol{\Sigma}^{\text{old}} \mathbf{X}] \right) \\ &\quad + \frac{D}{2} \ln \alpha - \frac{\alpha}{2} \left(\|\boldsymbol{\mu}^{\text{old}}\|^2 + \text{Tr}[\boldsymbol{\Sigma}^{\text{old}}] \right) + \text{const.} \end{aligned}$$

$$\text{where } \begin{cases} \boldsymbol{\mu}^{\text{old}} = \beta^{\text{old}} \boldsymbol{\Sigma}^{\text{old}} \mathbf{X}^\top \mathbf{t} \\ \boldsymbol{\Sigma}^{\text{old}} = (\beta^{\text{old}} \mathbf{X}^\top \mathbf{X} + \alpha^{\text{old}} \mathbf{I})^{-1} \end{cases}$$



88/95

M-Step:

$$(\alpha^{\text{new}}, \beta^{\text{new}}) = \underset{\alpha, \beta}{\operatorname{argmin}} \mathcal{Q} \left(\underbrace{\alpha, \beta}_{\theta}, \underbrace{\alpha^{\text{old}}, \beta^{\text{old}}}_{\theta^{\text{old}}} \right)$$

$$\begin{aligned} \frac{\partial \mathcal{Q}(\theta, \theta^{\text{old}})}{\partial \alpha} &= \frac{D}{2\alpha} - \frac{1}{2} \left(\|\boldsymbol{\mu}^{\text{old}}\|^2 + \operatorname{Tr}[\boldsymbol{\Sigma}^{\text{old}}] \right) \\ \frac{\partial \mathcal{Q}(\theta, \theta^{\text{old}})}{\partial \beta} &= \frac{N}{2\beta} - \frac{1}{2} \left(\|\mathbf{t} - \mathbf{X}\boldsymbol{\mu}^{\text{old}}\|^2 + \operatorname{Tr}[\mathbf{X}^{\top} \boldsymbol{\Sigma}^{\text{old}} \mathbf{X}] \right) \end{aligned}$$

$$\begin{aligned} \alpha^{\text{new}} &= \frac{D}{\|\boldsymbol{\mu}^{\text{old}}\|^2 + \operatorname{Tr}[\boldsymbol{\Sigma}^{\text{old}}]} \\ \beta^{\text{new}} &= \frac{N}{\|\mathbf{t} - \mathbf{X}\boldsymbol{\mu}^{\text{old}}\|^2 + \operatorname{Tr}[\mathbf{X}^{\top} \boldsymbol{\Sigma}^{\text{old}} \mathbf{X}]} \end{aligned}$$



89/95

Partially Hidden Data

- We can learn when there are missing (hidden) variables on some cases and not on others.

$$\ln p(\mathcal{D}|\theta) = \sum_{\text{complete}} \ln p(\mathbf{X}^c, \mathbf{Z}^c|\theta) + \sum_{\text{incomplete}} \ln \int p(\mathbf{X}, \mathbf{Z}^m|\theta) d\mathbf{Z}$$

(E-Step) estimate the hidden variables on the incomplete case only

(M-Step) optimize the log-likelihood on the complete data plus the expected log-likelihood on the incomplete data



91/95

Supplementary



Probabilistic Graphical Models (PGM)

Framework for representing dependencies
among the random variables

Probability Theory

Graph Theory



Probability Theory

- Sum rule

$$p(A) = \int p(A, B) dB$$

- Product (chain) rule

$$p(A, B) = p(B|A) p(A)$$

- Bayes rule

$$p(B|A) = \frac{p(A|B) p(B)}{p(A)}$$



(Statistical) Independence

- Marginal independence

$$A \perp\!\!\!\perp B \equiv p(A, B) = p(A) p(B)$$

- Conditional independence

$$\begin{aligned} A \perp\!\!\!\perp B|C &\equiv p(A, B|C) = p(A|C) p(B|C) \\ &\equiv p(A|B, C) = p(A|C) \end{aligned}$$

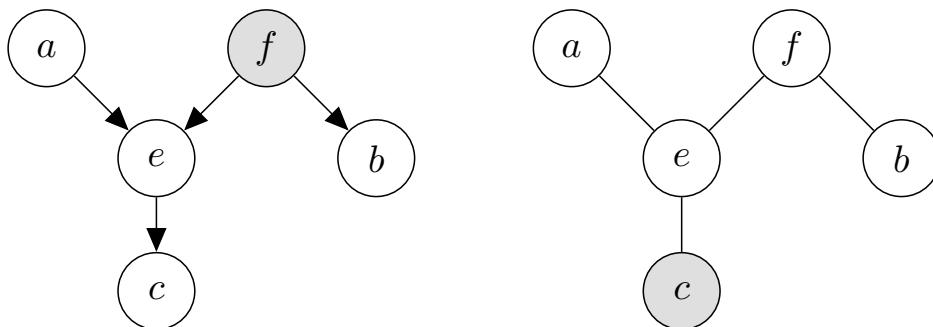
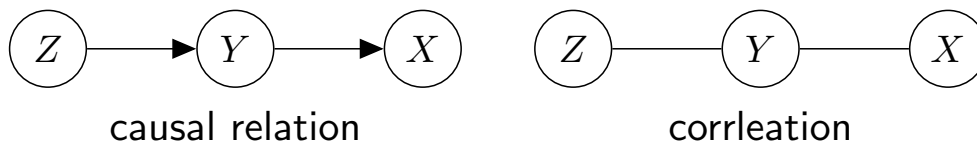


Graph theory:

$$\mathbb{G}(V, E)$$

► Graph Terminology

- V (vertices or nodes): represents random variables
 - Observed measurements, parameters, latent variables, hypothesis
- E (edges or links): represents *probabilistic* relationships between variables



- The graph captures the way in which *the joint distribution over all of the random variables can be decomposed into a product of factors* each depending only on a subset of the variables.



Useful properties of PGM

- A simple way to visualize the structure of a probabilistic model
- Used to design and motivate new models
- Insights into the properties of the model, including conditional independence properties
- Complex computations can be expressed in terms of graphical manipulation, in which underlying mathematical expressions are carried along implicitly.



- The pattern of edges in the graph represents the **qualitative dependencies** between the variables; the absence of an edge between two nodes means that any statistical dependency between these two variables is mediated via some other variable or set of variables.
- The **quantitative dependencies** between variables which are connected via edges are specified via parameterized conditional distributions, or more generally non-negative potential functions. The pattern of edges and the potential functions together specify a joint probability distribution over all the variables in the graph.



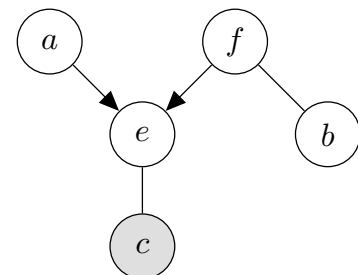
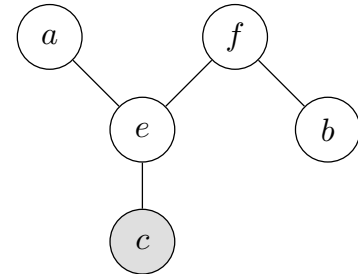
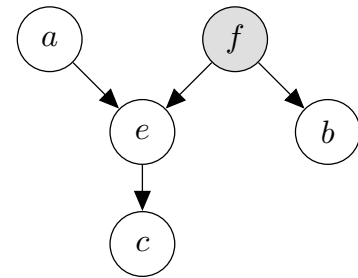
- **Bayesian Networks (BN)**

- ▶ *Directed* Acyclic Graph (DAG)
- ▶ Useful for expressing **causal relationships** between random variables
- ▶ A graphical way to represent a particular factorization of a joint distribution

- **Markov Random Fields (MRF)**

- ▶ *Undirected* Graphical Models
- ▶ The links do not carry arrows and have no directional significance
- ▶ **Better suited to expressing soft constraints** between random variables

- **Chain Graphs:** including both directed and undirected links



▶ Go Back