# Mixture models (1)

- Consider the old problem of modeling a pdf given a dataset of examples X={x<sup>(1</sup>, x<sup>(2</sup>, ..., x<sup>(N</sup>})
  - If the form of the underlying pdf was known (e.g. Gaussian), the problem could be solved using Maximum Likelihood (Lecture 6)
  - If the form of the pdf was unknown, the problem had to be solved with nonparametric density estimation methods such as Parzen windows (Lectures 7-8)

### We will now consider an alternative density estimation method: modeling the pdf with a <u>mixture</u> of parametric densities

- These methods are sometimes known as semi-parametric
  - Think of the individual components in the mixture as kernels, except for there is only a few of them, as opposed to one per data point as in Lecture 7
- In particular, we will focus on mixture models of Gaussian densities (surprised?)

$$P(\mathbf{x} \mid \boldsymbol{\theta}) = \sum_{c=1}^{C} P(\mathbf{x} \mid \boldsymbol{\theta}_{c}) P(\boldsymbol{\omega}_{c})$$



# Mixture models (2)

#### The mixture model problem can be posed in terms of the ML criterion

• Given a dataset of examples X={x<sup>(1</sup>, x<sup>(2</sup>, ..., x<sup>(N</sup>}, find the parameters of the model that maximize the log likelihood of the data

$$\hat{\theta} = \operatorname{argmax}[p(X \mid \theta)] = \operatorname{argmax}\left[\sum_{n=1}^{N} \log p(x^{(n} \mid \theta))\right] = \operatorname{argmax}\left[\sum_{n=1}^{N} \log \sum_{c=1}^{C} p(x^{(n} \mid \theta_{c}) P(\omega_{c}))\right]$$

- where  $\theta_c = {\mu_c, \Sigma_c}$  and  $P(\omega_c)$  are the parameters and mixing coefficient of the c-th mixture component, respectively
  - The mixing coefficients may also be interpreted as priors

#### • We could try to find the maximum of this function by differentiation

• For  $\Sigma_i = \sigma_i I$ , it can be shown [Bishop, 1995] that the solution becomes

$$\begin{split} & \frac{\partial}{\partial \mu_{c}} [\cdot] = 0 \qquad \Rightarrow \quad \hat{\mu}_{c} = \frac{\sum_{n} P(\omega_{c} \mid x^{(n)}) x^{(n)}}{\sum_{n} P(\omega_{c} \mid x^{(n)})} \\ & \frac{\partial}{\partial \sigma_{c}} [\cdot] = 0 \qquad \Rightarrow \quad \hat{\sigma}_{c}^{2} = \frac{1}{d} \frac{\sum_{n} P(\omega_{c} \mid x^{(n)}) \left\| x^{(n)} - \hat{\mu}_{c} \right\|^{2}}{\sum_{n} P(\omega_{c} \mid x^{(n)})} \\ & \frac{\partial}{\partial P(\omega_{c})} [\cdot] = 0 \quad \Rightarrow \quad \hat{P}(\omega_{c}) = \frac{1}{N} \sum_{n} P(\omega_{c} \mid x^{(n)}) \end{split}$$



# Mixture models (3)

- Notice that the previous equations are not a closed form solution
  - The model parameters  $\mu_c$ ,  $\Sigma_c$ , and  $P(\omega_c)$  also appear on the RHS as a result of Bayes rule!
  - Therefore, these expressions represent a highly non-linear coupled system of equations
- However, these expressions suggest that we may be able to use a fixed-point algorithm to find the maxima
  - **1.** Begin with some value of the model parameters. Call these the "old" values
  - 2. Evaluate the RHS of the equations to obtain "new" values for the parameters
  - 3. Let these "new" values become the "old" ones and repeat the process
- Surprisingly, an algorithm of this simple form can be found which is guaranteed to increase the log-likelihood with every iteration!
  - This example represents a particular case of a more general procedure known as the <u>Expectation-Maximization</u> algorithm



## The Expectation-Maximization algorithm (1)

#### The EM is a general method for finding the ML estimate of the parameters of a pdf when the data has missing values

- There are two main applications of the EM algorithm
  - When the data indeed has incomplete, missing or corrupted values as a result of a faulty observation process
  - When assuming the existence of missing or hidden parameters can simplify the likelihood function, which would otherwise lead to an analytically intractable optimization problem. This is the case that occupies our discussion

### Assume a dataset containing two types of features

- A set of features X whose value is known. We call these the incomplete data
- A set of features Z whose value is unknown. We call these the missing data

### • We now define a joint pdf $p(X,Z|\theta)$ called the complete-data likelihood

- This function is a random variable since the features Z are unknown
- You can think of p(X,Z|θ)=h<sub>X,θ</sub>(Z), for some function h<sub>X,θ</sub>(·), where X and θ are constant and Z is a random variable

#### As suggested by its name, the EM algorithm operates by performing two basic operations over and over:

- An Expectation step
- A Maximization step



## The Expectation-Maximization algorithm (2)

### EXPECTATION

• Find the <u>expected</u> value of the log-likelihood log[p(X,Z| $\theta$ )] with respect to the unknown data Z, <u>given</u> the data X and the current parameter estimates  $\theta^{(i-1)}$ 

$$\mathbf{Q}\left(\boldsymbol{\theta} \mid \boldsymbol{\theta}^{(i-1)}\right) = \mathbf{E}_{z}\left[\mathbf{logp}(\mathbf{X}, \mathbf{Z} \mid \boldsymbol{\theta}) \mid \mathbf{X}, \boldsymbol{\theta}^{(i-1)}\right]$$

- where  $\theta$  are the new parameters that we will have to optimize to increase Q
- Note that X and  $\theta^{(i-1)}$  are constants,  $\theta$  is the variable that we wish to adjust, and Z is a random variable defined by a pdf p(Z|X, $\theta^{(i-1)}$ ). Therefore Q( $\theta|\theta^{(i-1)}$ ) is just a function of  $\theta$

 $\mathsf{E}_{z}\left[\mathsf{logp}(X, Z \mid \theta) \mid X, \theta^{(i-1)}\right] = \int_{z \in Z} \mathsf{logp}(X, z \mid \theta) p(z \mid X, \theta^{(i-1)}) dz$ 

#### MAXIMIZATION

• Find the argument  $\theta$  that <u>maximizes</u> the expected value defined by Q( $\theta|\theta^{(i-1)}$ )

$$\theta^{(i)} = \operatorname{argmax} Q(\theta \,|\, \theta^{(i-1)})$$

#### Convergence properties

 It can be shown that (1) each iteration (E+M) is guaranteed to increase the loglikelihood and (2) the EM algorithm is guaranteed to converge to a local maximum of the likelihood function



### The Expectation-Maximization algorithm (3)

#### The two steps of the EM algorithm are illustrated in the figure below

- During the E step, the unknown features Z are integrated out assuming the current values of the parameters  $\theta^{(i-1}$
- During the M step, the values of the parameters that maximize the expected value of the log likelihood are obtained



IN A NUTSHELL: since Z are unknown, the best we can do is maximize the average log-likelihood across all possible values of Z



Introduction to Pattern Analysis Ricardo Gutierrez-Osuna Texas A&M University

# The EM algorithm and mixture models (1)

- Having formalized the expectation maximization algorithm, we are now ready to find the solution to the mixture model problem
  - To keep things simple, we will assume a *univariate* mixture model where all the components have the same known standard deviation  $\sigma$

### Problem formulation

- As usual, we are given a dataset X={ $x^{(1)}, x^{(2)}, ..., x^{(N)}$ }, and we are asked to estimate the model parameters  $\theta$ ={ $\mu_1, \mu_2, ..., \mu_C$ }
- The following process is assumed to have generated each random variable  $\boldsymbol{x}^{(n)}$ 
  - First, a Gaussian component is selected according to the mixture coefficients  $P(\omega_c)$
  - Then,  $x^{(n)}$  is generated according to the likelihood  $p(x|\mu_c)$  of that particular component
- In a mixture model problem, the hidden variables Z={z<sub>1</sub><sup>(n</sup>,z<sub>2</sub><sup>(n</sup>,...z<sub>C</sub><sup>(n</sup>} are used to indicate which of the C Gaussian components generated data point x<sup>(n</sup>

### Solution

• The probability  $p(x,z|\theta)$  for a specific example is

$$p(x^{(n)}, z^{(n)}_{_{1}}, z^{(n)}_{_{2}}, ..., z^{(n)}_{_{c}} \mid \theta) = \frac{1}{\sqrt{2\pi\sigma^{2}}} exp\left(-\frac{1}{2\sigma^{2}} \sum_{c=1}^{C} z^{(n)}_{_{c}} (x^{(n)} - \mu_{_{c}})^{2}\right)$$

• Notice than only one of the  $z_c^{(n)}$  can have a value of 1,and all others are zero



### The EM algorithm and mixture models (2)

• The log-likelihood of the entire dataset is then:

$$\log p(X, Z \mid \theta) = \log \prod_{n=1}^{N} p(x^{(n}, z^{(n} \mid \theta)) = \sum_{n=1}^{N} \left( \log \frac{1}{\sqrt{2\pi\sigma^{2}}} - \frac{1}{2\sigma^{2}} \sum_{c=1}^{C} z_{c}^{(n} (x^{(n} - \mu_{c})^{2}) \right)$$

• To obtain  $Q(\theta|\theta^{(i-1)})$  we must then take the expectation over Z:

$$E_{Z}[\log p(X, Z | \theta)] = \sum_{n=1}^{N} \left( \log \frac{1}{\sqrt{2\pi\sigma^{2}}} - \frac{1}{2\sigma^{2}} \sum_{c=1}^{C} E[z_{c}^{(n)}] (x^{(n)} - \mu_{c})^{2} \right)$$

- where we have used the fact that E[f(z)]= f(E[z]) for a linear function f(z)
- $E[z_c^{(n)}]$  is simply the probability that example  $x^{(n)}$  was generated by the c-th Gaussian component *given the current model parameters*  $\theta^{(i-1)}$

$$E[z_{c}^{(n)}] = \frac{p(x = x^{(n)} | \mu = \mu_{c}^{(i-1)})}{\sum_{q=1}^{C} p(x = x^{(n)} | \mu = \mu_{q}^{(i-1)})} = \frac{exp\left(-\frac{1}{2\sigma^{2}}(x^{(n)} - \mu_{c}^{(i-1)})^{2}\right)}{\sum_{q=1}^{C} exp\left(-\frac{1}{2\sigma^{2}}(x^{(n)} - \mu_{q}^{(i-1)})^{2}\right)}$$
(1)

• These two expressions define the Q function

