

Mixtures - E.G. 7 and 7

$$p(\underline{x}) = \sum_{t=1}^k p(\underline{x}|G_t)P(G_t)$$

$G_t$  are the mixture components

$p(\underline{x}|G_t)$  component densities. Different classes  
 $P(G_t)$  mixture proportion.

1) Component densities are Gaussian and  
 have  $p(\underline{x}|G_t) \sim N(\underline{\mu}_t, \Sigma_t)$

$\Phi = \{P(G_t), \underline{\mu}_t, \Sigma_t\}_{t=1}^k$  one parameters that  
 need to be estimated.

In supervised case, we know the classes, and  
 the assignment of instances to classes.

$$r_i^t = 1, \text{ if } \underline{x}^t \in C_i \\ 0, \text{ otherwise}$$

In this case, it is straight-forward to learn  
 the parameters

$$\hat{P}(G_t) = \frac{1}{N} \sum_i r_i^t$$

$$\underline{\mu}_t = \frac{\sum_i r_i^t \underline{x}^t}{\sum_i r_i^t}$$

$$\Sigma_t = \frac{\sum_i r_i^t (\underline{x}^t - \underline{\mu}_t)(\underline{x}^t - \underline{\mu}_t)^T}{\sum_i r_i^t}$$

(2)

## Mixture Densities (cont.)

But, usually we do not know the classes.

This is called unsupervised learning

i.e.  $r^t$  is not known.

We need to estimate  $r^t$  — later, we show that the EM algorithm can be used.

### k-means clustering.

Vector quantization:

Sample  $\underline{X} = \{\underline{x}^t\}_{t=1}^N$

$k$  reference vectors  $\underline{m}_j \quad j=1 \text{ to } k$ .

Represent  $\underline{x}^t$  by the most similar entry  $\underline{m}_i$

$$\|\underline{x}^t - \underline{m}_i\| = \min_j \|\underline{x}^t - \underline{m}_j\|$$

$\underline{m}_i$  — codebook vectors.

Compression — results from quantization.

How to calculate the best  $\{\underline{m}_j\}$ ?



(3)

## K-means (cont)

Reconstruction error

$$E(\sum_{i=1}^k b_i^t \|x^t - \underline{m}_i\|^2 | X) = \sum_t \sum_i b_i^t \|x^t - \underline{m}_i\|^2.$$

where  $b_i^t = \begin{cases} 1 & \text{if } \|x^t - \underline{m}_i\| = \min_j \|x^t - \underline{m}_j\| \\ 0 & \text{otherwise} \end{cases}$

k-means is an iterative procedure that seeks to minimize  $E(\cdot)$ .

Start with random  $\underline{m}_i$ .

Then compute  $b_i^t$ :

re-estimate the means  $\underline{m}_i = \frac{\sum_t b_i^t x^t}{\sum_t b_i^t}$ .

Repeat.

This procedure is guaranteed to converge

(almost always).

But the final results depend on the initial guess.

Alternative Initialization Strategies:

- Take random k instances as initial  $\underline{m}_i$
- Calculate the means of all the data, then add small perturbations
- Calculate the principal components, divide its range into k intervals, etc.

(4)

## Expectation-Maximization (EM) Algorithm.

$$\begin{aligned} \mathcal{L}(\phi | x) &= \log \prod_t p(x^t | \phi) \\ &= \sum_t \log \sum_{i=1}^k P(x^t | G_i) P(G_i) \end{aligned}$$

EM Introduce hidden variable  $Z$  which associate the data to the class.

$$\begin{aligned} Z^t &= \{z_1^t, \dots, z_k^t\} & z_i^t = 1 \text{ if } x^t \\ &\quad \text{below } t \text{ class } G_i \\ P(z^t) &= \prod_{i=1}^k \pi_i^{z_i^t} & 0, \text{ otherwise.} \\ P(x^t | z^t) &= \prod_{i=1}^k P_i(x^t)^{z_i^t} & P_i(x^t) \text{ short} \\ P(x^t, z^t) &= P(z^t) P(x^t | z^t) \end{aligned}$$

$$\begin{aligned} \mathcal{L}_c(\phi | x, z) &= \log \prod_t p(x^t, z^t | \phi) \\ &= \sum_t \log p(x^t, z^t | \phi) \\ &= \sum_t (\log P(z^t | \phi) + \log P(x^t | z^t, \phi)) \\ &= \sum_t \bar{z}_i z_i^t \{\log \pi_i + \log P_i(x^t | \phi)\} \end{aligned}$$

(5) EM contd

E-step

$$\begin{aligned} Q(\phi | \phi^e) &= E[\log P(x, z) | x, \phi^e] \\ &= E[\sum_i L_i(\phi | x, z) | x, \phi^e] \\ &= \sum_t \sum_i E[z_i^t | x, \phi] \quad \left\langle \log \Pi_i + \log p_i(x^t | \phi) \right\rangle \end{aligned}$$

where  $E[z_i^t | x, \phi]$

$$\begin{aligned} &= E[z_i^t | x^t, \phi^e] \\ &= P(z_i^t = 1 | x^t, \phi^e) \\ &= \frac{P(x^t | z_i^t = 1, \phi^e) P(z_i^t = 1 | \phi)}{P(x^t | \phi^e)} \quad \text{Bayes rule} \\ &= \frac{P_i(x^t | \phi^e) \Pi_i}{\sum_j P_j(x^t | \phi^e) \Pi_j} \\ &= \frac{P(x^t | G_i, \phi^e) P(G_i)}{\sum_j P(x^t | G_j, \phi^e) P(G_j)} \\ &= P(G_i | x^t, \phi^e) = h_i^t \end{aligned}$$

Expected value of hidden variable,  $E[z_i^t]$   
 is the posterior probability that  $x^t$  is generated by  $G_i$ .  
 "Soft label"

## (6) EM (Cont)

M-step: maximize to get next of parameter  $\phi^{t+1} = \arg \max_{\phi} Q(\phi | \phi^t)$

$$Q(\phi | \phi^t) = \sum_t \sum_i h_i^t \left\{ \log \pi_i + \log p_i(x^t | \phi^t) \right\}$$

$$= \sum_t \sum_i h_i^t \log \pi_i + \sum_t \sum_i h_i^t \log p_i(x^t | \phi)$$

$$\text{solve to get } \pi_i = \sum_t h_i^t / N$$

If we assume Gaussian components

$$\hat{P}_i(x^t | \phi) \sim N(m_i, S_i)$$

$$M\text{-step is } m_i^{t+1} = \frac{\sum_t h_i^t x^t}{\sum_t h_i^t}$$

$$S_i^{t+1} = \sum_t h_i^t \frac{(x^t - m_i^{t+1})(x^t - m_i^{t+1})^\top}{\sum_t h_i^t}$$

$$\text{where } h_i^t = \frac{\pi_i |S_i|^{-\frac{1}{2}} \exp\{-\frac{1}{2}(x^t - m_i)^\top S_i^{-1}(x^t - m_i)\}}{\sum_j \pi_j |S_j|^{-\frac{1}{2}} \exp\{-\frac{1}{2}(x^t - m_j)^\top S_j^{-1}(x^t - m_j)\}}$$

(7)

## Supervised Learning after Clustering.

Clustering can be used for data exploration by grouping instances.

If such groups are found, they can be named and their attributes defined.

Clustering is frequently used as a pre-processing stage.

## (8) Modern View of EM

Task estimate parameters  $\phi$  from  $p(x, z | \phi)$ .

Problem we can't observe  $z$

Proposed Solution: Estimate  $\phi$  from  $p(x | \phi)$

$$\text{where } p(x | \phi) = \sum_z p(x, z | \phi)$$

use iterative algorithm, EM, to estimate  $\phi$ .

Formulate minimize w.r.t.  $\phi$ .  $-\log p(x | \phi)$

Equivalent to minimize

$$F[\phi, q(z)] = -\log p(x | \phi) + \sum_z q(z) \log \frac{q(z)}{p(z | x, \phi)}$$

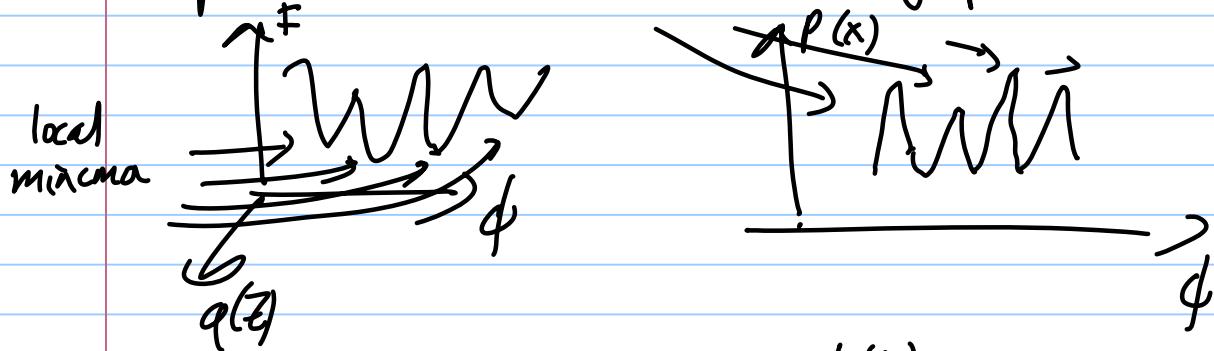
w.r.t.  $\phi$  &  $q(z)$ .

Why  $\rightarrow$  because the second term is  $> 0$  and  $= 0$  only if  $q(z) = p(z | x, \phi)$

Kullback-Leibler divergence.

(9) Algorithm (EM) minimize  $F[\phi, q(z)]$   
w.r.t.  $\phi$  &  $q(z)$  alternately.

This is guaranteed to always decrease  $F$  and converge to a local minimum. This corresponds to a local maximum of  $p(x|\phi)$



Initialize the parameters  $\phi^{(0)}$ ,  
then repeat the two steps : at time  $t$ .

E-step : minimize  $F$  w.r.t.  $q(z)$  gives

$$q^{t+1}(z) = p(z|x, \phi^t)$$

M-step : minimize  $F$  w.r.t.  $\phi$  gives

$$\phi^{t+1} = \arg \min_{\phi} - \sum_z q^{t+1}(z) \log p(x, z | \phi)$$

see next page ...

(10) E-step follows from standard expressions of  $F[\phi, q(z)] = -\log p(x|\phi) + \sum_z q(z) \log \left\{ \frac{q(z)}{p(z|x,\phi)} \right\}$

M-step follows by re-expressing  $q(z) = p(z|x,\phi)$ .  
 minimise by setting

$$F = -\log p(x|\phi) + \sum_z q(z) \log q(z) - \sum_z q(z) \log p(z|x,\phi)$$

$$F = \sum_z q(z) \log (z) - \underbrace{\sum_z q(z) \log p(x|\phi)}_{\text{using } \sum_z q(z) = 1} - \underbrace{\sum_z q(z) \log p(z|x,\phi)}_{-\sum_z q(z) \log \{p(x|\phi)p(z|x,\phi)\}} - \underbrace{\sum_z q(z) \log p(x,z|\phi)}$$

This is the "modern" formulation of EM.  
 It gives a proof of convergence. It shows where the two steps come from.

It also leads to variants of EM (beyond scope of the course).