#### EM (more formally)

Semi-optional technical material alert!

The formal treatment helps us use EM correctly in more complex situations. However, EM algorithms for many problems can "guessed at" using intuition.

The more formal treatment is not needed for the homework.

#### EM (more formally)

- Assume K clusters. Index over clusters by *k*, over points by *n*.
- New notation for cluster membership:

For each point,  $n, z_n$  is a vector of K values where exactly one

 $z_{n,k} = 1$ , and all others are 0. Note that  $\sum_{k} z_{n,k} = 1$ .

#### EM (more formally)

• Denote cluster priors by:

 $\boldsymbol{\pi}_{k} \equiv p(\boldsymbol{z}_{k}=1)$ 

• Denote the responsibilities that each cluster has for each point by:

$$\gamma(z_{n,k}) \equiv p(z_{n,k} = 1 | x_n, \theta^{(s)}) = \frac{\pi_k p(x_n | z_{n,k} = 1, \theta_k^{(s)})}{\sum_{k'} \pi_{k'} p(x_n | z_{n,k'} = 1, \theta_{k'}^{(s)})}$$

## EM (more formally)

Represent the entire data set of N points,  $\mathbf{x}_n$ , as a matrix X with rows  $\mathbf{x}_n^T$ .

Represent the latent variable assignments with a matrix Z. (For the true assignment, each row is zero except for a single element that is 1.)

We call  $\{Z, X\}$  the *complete* data set (everthing is known). The observed data,  $\{X\}$ , is called the *incomplete* data set.

## EM (more formally)

We assume that computing the MLE of parameters,

 $\arg\max_{\theta}\left\{\log\left\{p\left(Z,X|\theta\right)\right\}\right\}$ 

with complete data is relatively easy.

Recall our intuitive treatment of EM for GMM. If we knew the cluster membership, we would know how to compute the means.

Since we did not know the cluster membership we did a weighted computation.

## EM (more formally)

Notice the complexity of the incomplete log likelihood:

$$\log(p(X|\theta)) = \sum_{n} \log\left(\sum_{k} \pi_{k} p(x_{n}|\theta)\right)$$
nasty sum in log

By constrast, for complete log likelihood we can incorporate the assignment by:

$$p(X,Z|\theta) = \prod_{n} \prod_{k} \pi_{k}^{z_{n,k}} \left\{ p(x_{n}|\theta) \right\}^{z_{n,k}}$$

So

$$\log(p(X,Z|\theta)) = \sum_{n} \sum_{k} \left\{ z_{n,k} \left( \log(\pi_k) + \log(p(x_n|\theta)) \right) \right\}$$

(No nasty sum in log; well suited for the expectation calculation).

# EM (more formally)

For the E step, we compute the responsibilities which is straightforward.

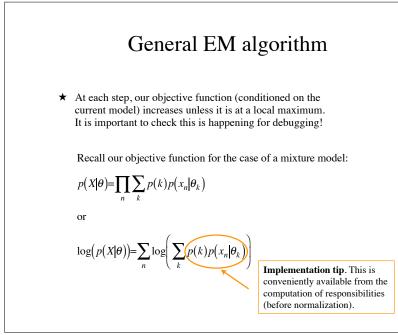
Next, define  $Q(\theta^{(s+1)}, \theta^{(s)}) = \sum_{Z} p(Z|X, \theta^{(s)}) \log(p(X, Z|\theta^{(s+1)}))$ (Expectation of  $\log(p(X, Z|\theta^{(s+1)}))$  over  $p(Z|X, \theta^{(s)})$ ).

The M step then computes  $\theta^{(s+1)} = \arg \max_{a} \left\{ Q(\theta^{(s+1)}, \theta^{(s)}) \right\}$ 

Maximizing Q is generally feasible and corresponds to the intuitive development.

## General EM algorithm

- 1. Choose initial values for  $\theta^{(s=1)}$ (can also do assignments, but then jump to M step). 2. E step: Evalute  $p(Z|X,\theta^{(s)})$ 3. M step: Evalute  $\theta^{(s+1)} = \arg \max_{\theta} \{Q(\theta^{(s+1)},\theta^{(s)})\}$ where  $Q(\theta^{(s+1)},\theta^{(s)}) = \sum_{Z} p(Z|X,\theta^{(s)}) \log(p(X,Z|\theta^{(s+1)}))$ 4. Check for convergence; If not done, goto 2.
  - ★ At each step, our objective function increases unless it is at a local maximum. It is important to check this is happening for debugging!



# Deriving the GMM M-step

Evalute 
$$\theta^{(s+1)} = \arg \max_{\theta} \left\{ Q(\theta^{(s+1)}, \theta^{(s)}) \right\}$$
  
where  $Q(\theta^{(s+1)}, \theta^{(s)}) = \sum_{Z} p(Z|X, \theta^{(s)}) \log(p(X, Z|\theta^{(s)}))$   
Recall that  $\log(p(X, Z|\theta)) = \sum_{n} \sum_{k} \left\{ z_{n,k} \left( \log(\pi_k) + \log(p(x_n|\theta_k)) \right) \right\}$   
So  $Q(\theta^{(s+1)}, \theta^{(s)}) = \sum_{Z} p(Z|X, \theta^{(s)}) \sum_{n} \sum_{k} \left\{ z_{n,k} \left( \log(\pi_k) + \log(p(x_n|\theta_k)) \right) \right\}$ 

Deriving the GMM M-step  $Q(\theta^{(s+1)}, \theta^{(s)}) = \sum_{Z} p(Z|X, \theta^{(s)}) \sum_{n} \sum_{k} \left\{ z_{n,k} \left( \log(\pi_{k}) + \log(p(x_{n}|\theta_{k})) \right) \right\}$   $= \sum_{n} \sum_{k} \sum_{Z} p(Z|X, \theta^{(s)}) \left\{ z_{n,k} \left( \log(\pi_{k}) + \log(p(x_{n}|\theta_{k})) \right) \right\}$ inner sum

This exchanging of summing order says that **instead** of summing over points and clusters for all correspondences Z, we sum over all correspondences for a **given** point and cluster.

We will focus on the inner sum.

Z is all possible correspondences. To generate them all more explicitly, we can consider the first point. For each possible assignment of the first point, we then need all possible combinations of the other points. To get that, we consider all possible assignments of the second point, together with all possible assignments of the remaining points. This shows:

$$\sum_{Z} ( ) \equiv \sum_{\mathbf{z}_{1}} \sum_{\mathbf{z}_{2}} \sum_{\mathbf{z}_{3}} \cdots \sum_{\mathbf{z}_{N}} ( )$$

Note that a sum over  $\mathbf{z}_n$  is short hand for a sum over cluster assignments for point n. Hence each of the sums on the right are over clusters.

$$\sum_{\mathbf{z}_n} (\bullet) \equiv \sum_{k_n=1}^{\kappa} (\bullet)$$

Further, because the points are independent, we have:

$$p(Z|\bullet) = \prod_{\mathbf{z}_i} p(\mathbf{z}_i|\bullet)$$

$$\sum_{k} p(Z|X,\theta^{(i)}) \cdot z_{n,k} \cdot \left(\log(\pi_{k}) + \log(p(x_{n}|\theta_{k}))\right)$$

$$= \sum_{k_{1}=1}^{K} \sum_{k_{2}=1}^{K} \cdots \sum_{k_{n}=1}^{K} \prod_{i} p(z_{i,k_{i}}|x_{i},\theta^{(i)}) \cdot z_{n,k} \cdot \left(\log(\pi_{k}) + \log(p(x_{n}|\theta_{k}))\right)$$

$$= \sum_{k_{n}=1}^{K} p(z_{n,k_{n}}|x_{i},\theta^{(i)}) \cdot z_{n,k} \cdot \left(\log(\pi_{k}) + \log(p(x_{n}|\theta_{k}))\right) \cdot \sum_{k_{n}=1}^{K} \sum_{k_{n}=1}^{K} \cdots \sum_{k_{n}=1}^{K} p(z_{n,k_{n}}|x_{n},\theta^{(i)}) \cdot z_{n,k} \cdot \left(\log(\pi_{k}) + \log(p(x_{n}|\theta_{k}))\right) \cdot \sum_{k_{n}=1}^{K} p(z_{n-1,k_{n-1}}|x_{n-1},\theta^{(i)}) \sum_{k_{n}=1}^{K} p(z_{n,k_{n}}|x_{n},\theta^{(i)}) \cdots \sum_{k_{n}=1}^{K} p(z_{n,k_{n}}|x_{n},\theta^{(i)}) \cdot z_{n,k} \cdot \left(\log(\pi_{k}) + \log(p(x_{n}|\theta_{k}))\right)$$

$$= \sum_{k_{n}=1}^{K} p(z_{n,k_{n}}|x_{n},\theta^{(i)}) \cdot z_{n,k} \cdot \left(\log(\pi_{k}) + \log(p(x_{n}|\theta_{k}))\right)$$

$$\sum_{z} p(Z|X,\theta^{(s)}) \cdot z_{n,k} \cdot (\log(\pi_{k}) + \log(p(x_{n}|\theta_{k}))))$$

$$= \sum_{k_{n}=1}^{K} \sum_{k_{2}=1}^{K} \cdots \sum_{k_{n}=1}^{K} \prod_{i} p(z_{i,k_{i}}|x_{i},\theta^{(s)}) \cdot z_{n,k} \cdot (\log(\pi_{k}) + \log(p(x_{n}|\theta_{k}))))$$

$$= \sum_{k_{n}=1}^{K} p(z_{n,k_{n}}|x_{i},\theta^{(s)}) \cdot z_{n,k} \cdot (\log(\pi_{k}) + \log(p(x_{n}|\theta_{k})))) \cdot \sum_{k_{1}=1}^{K} \sum_{k_{2}=1}^{K} \cdots \sum_{k_{n}=1}^{K} \sum_{k_{n}=1}^{K} \sum_{i=n}^{K} p(z_{i,k_{n}}|x_{i},\theta^{(s)}) \cdot z_{n,k} \cdot (\log(\pi_{k}) + \log(p(x_{n}|\theta_{k}))))$$

$$= p(z_{n,k}=1|X,\theta^{(s)})(\log(\pi_{k}) + \log(p(x_{n}|\theta_{k})))$$

$$= p(z_{n,k})(\log(\pi_{k}) + \log(p(x_{n}|\theta_{k}))) \quad (definition of \gamma(z_{n,k}), the responsibility)$$

Deriving the M-step

$$Q(\theta^{(s+1)}, \theta^{(s)}) = \sum_{Z} p(Z|X, \theta^{(s)}) \sum_{n} \sum_{k} \left\{ z_{n,k} \left( \log(\pi_{k}) + \log(p(x_{n}|\theta_{k})) \right) \right\}$$
$$= \sum_{n} \sum_{k} \left\{ \gamma(z_{n,k}) \left( \log(\pi_{k}) + \log(p(x_{n}|\theta_{k})) \right) \right\}$$

We need to maximize this with respect to the parameters for each cluster, k. Notice that:

$$\frac{\delta}{\delta\theta_{k^*}} Q(\theta^{(s+1)}, \theta^{(s)}) = \sum_{n} \left\{ \gamma(z_{n,k^*}) \frac{\delta}{\delta\theta_{k^*}} \left( \log(\pi_{k^*}) + \log(p(x_n | \theta_{k^*})) \right) \right\}$$

(The values of k not of current interest, i.e., not  $k^*$ , die)

Example—deriving the GMM M-step

$$\frac{\delta}{\delta\mu_{k}} Q(\theta^{(s+1)}, \theta^{(s)}) = \sum_{n} \left\{ \gamma(z_{n,k}) \frac{\delta}{\delta\mu_{k}} \left( \log(\pi_{k}) + \log(p(x_{n}|\theta_{k})) \right) \right\}$$
$$= \sum_{n} \left\{ \gamma(z_{n,k}) \frac{\delta}{\delta\mu_{k}} \left( \log(p(x_{n}|\theta_{k})) \right) \right\}$$
$$\sum_{n} \left\{ \gamma(z_{n,k}) \frac{\delta}{\delta\mu_{k}} \left( \log(N(x_{n}|\mu_{k}, \Sigma_{k})) \right) \right\}$$

Example — deriving the GMM M-step  

$$N(\mathbf{x}_{n}|\mu_{k}, \Sigma_{k}) = \frac{1}{(2\pi)^{D/2} |\Sigma_{k}|^{1/2}} \exp\left(-\frac{1}{2}(\mathbf{x}_{n} - \mu_{k})^{T} \Sigma_{k}^{-1}(\mathbf{x}_{n} - \mu_{k})\right)$$

$$\log(N(\mathbf{x}_{n}|\mu_{k}, \Sigma_{k})) = \log\left(\frac{1}{(2\pi)^{D/2} |\Sigma_{k}|^{1/2}}\right) - \frac{1}{2}(\mathbf{x}_{n} - \mu_{k})^{T} \Sigma_{k}^{-1}(\mathbf{x}_{n} - \mu_{k})$$

$$\frac{\delta}{\delta\mu_{k}} \log(N(\mathbf{x}_{n}|\mu_{k}, \Sigma_{k})) = \Sigma_{k}^{-1}(\mathbf{x}_{n} - \mu_{k})$$
(exercise for the interested)

Example — deriving the GMM M-step  

$$\frac{\delta}{\delta\mu_{k}}Q(\theta^{(s+1)},\theta^{(s)}) = \sum_{n} \left\{ \gamma(z_{n,k}) \frac{\delta}{\delta\mu_{k}} (\log(N(x_{n}|\mu_{k},\Sigma_{k}))) \right\}$$

$$\frac{\delta}{\delta\mu_{k}}Q(\theta^{(s+1)},\theta^{(s)}) = 0 \text{ means that}$$

$$\sum_{n} \left\{ \gamma(z_{n,k}) \Sigma_{k}^{-1}(\mathbf{x}_{n} - \mu_{k}) \right\} = 0 \text{ (Inner sigma is precision matrix, not a sum}$$

$$\sum_{n} \left\{ \gamma(z_{n,k}) (\mathbf{x}_{n} - \mu_{k}) \right\} = 0 \text{ (Multiply by } \Sigma_{k})$$

Example — deriving the GMM M-step So,  $\sum_{n} \{ \gamma(z_{n,k}) (\mathbf{x}_{n} - \mu_{k}) \} = 0$ and  $\mu_{k} \sum_{n} \{ \gamma(z_{n,k}) \} = \sum_{n} \{ \gamma(z_{n,k}) (\mathbf{x}_{n}) \}$ and  $\mu_{k} = \frac{\sum_{n} \{ \gamma(z_{n,k}) (\mathbf{x}_{n}) \}}{\sum_{n} \{ \gamma(z_{n,k}) \}}$  (same as before)

## Example-deriving the GMM M-step

Finding variances/covariances is similar.

Finding the mixing coefficients is also similar, except we also need to enforce that they sum to one.

(Here the equations for the *k*'s are coupled).

So we use Lagrange Multipliers.