

EM (more formally)

We assume that computing the MLE of parameters,

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

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, which happens to be an expectation of the complete log likelihood, over the assignment (responsibility) distribution.

EM (more formally)

Notice the complexity of the incomplete log likelihood:

$$\log(p(X|\theta)) = \sum_n \log \left(\underbrace{\sum_k \pi_k p(x_n|\theta)}_{\text{nasty sum in log}} \right)$$

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

$$p(X, Z | \theta) = \prod_n \prod_k \pi_k^{z_{n,k}} \{ p(x_n | \theta) \}^{z_{n,k}}$$

So

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

(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.

$$\text{Define } Q(\theta^{(s+1)}, \theta^{(s)}) = \sum_Z p(Z | X, \theta^{(s)}) \log(p(X, Z | \theta^{(s+1)}))$$

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

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

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: Evaluate $p(Z | X, \theta^{(s)})$

3. M step: Evaluate $\theta^{(s+1)} = \arg \max_{\theta} \{ Q(\theta^{(s+1)}, \theta^{(s)}) \}$

$$\text{where } Q(\theta^{(s+1)}, \theta^{(s)}) = \sum_Z p(Z | X, \theta^{(s)}) \log(p(X, Z | \theta^{(s+1)}))$$

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!

General EM algorithm

- ★ At each step, our objective function (conditioned on the current model) increases unless it is at a local maximum. It is important to check this is happening for debugging!

Recall our objective function:

$$p(X) = \prod_n \sum_k p(k) p(x_n | k)$$

or

$$\log(p(X)) = \sum_n \log \left(\sum_k p(k) p(x_n | k) \right)$$

Implementation tip. This is conveniently available from the computation of responsibilities (before normalization).

Deriving the GMM M-step

$$\text{Evaluate } \theta^{(s+1)} = \arg \max_{\theta} \{Q(\theta^{(s+1)}, \theta^{(s)})\}$$

$$\text{where } Q(\theta^{(s+1)}, \theta^{(s)}) = \sum_Z p(Z | X, \theta^{(s)}) \log(p(X, Z | \theta^{(s)}))$$

$$\text{Recall that } \log(p(X, Z | \theta)) = \sum_n \sum_k \{z_{n,k} (\log(\pi_k) + \log(p(x_n | \theta_k)))\}$$

$$\text{So } Q(\theta^{(s+1)}, \theta^{(s)}) = \sum_Z p(Z | X, \theta^{(s)}) \sum_n \sum_k \{z_{n,k} (\log(\pi_k) + \log(p(x_n | \theta_k)))\}$$

Deriving the GMM M-step

$$\begin{aligned} Q(\theta^{(s+1)}, \theta^{(s)}) &= \sum_Z p(Z | X, \theta^{(s)}) \sum_n \sum_k \{z_{n,k} (\log(\pi_k) + \log(p(x_n | \theta_k)))\} \\ &= \sum_n \sum_k \sum_Z p(Z | X, \theta^{(s)}) \{z_{n,k} (\log(\pi_k) + \log(p(x_n | \theta_k)))\} \end{aligned}$$

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.

The quantity in parentheses only interacts with the given point and cluster from the outside. So intuitively, most of the sum ver

Deriving the GMM M-step

$$\begin{aligned} Q(\theta^{(s+1)}, \theta^{(s)}) &= \sum_Z p(Z | X, \theta^{(s)}) \sum_n \sum_k \{z_{n,k} (\log(\pi_k) + \log(p(x_n | \theta_k)))\} \\ &= \sum_n \sum_k \underbrace{\sum_Z p(Z | X, \theta^{(s)}) \{z_{n,k} (\log(\pi_k) + \log(p(x_n | \theta_k)))\}}_{\text{inner sum}} \end{aligned}$$

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:

$$\sum_Z () \equiv \sum_{z_1} \sum_{z_2} \sum_{z_3} \cdots \sum_{z_N} ()$$

Note that a sum over 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.

Further, because the points are independent, we have:

$$p(Z|\bullet) = \prod_{z_i} p(z_i|\bullet)$$

$$\underbrace{\sum_Z p(Z|X, \theta^{(s)}) \cdot z_{n,k} \cdot (\log(\pi_k) + \log(p(x_n|\theta_k)))}_{\text{inner sum from previous}}$$

$$= \sum_{z_1} \sum_{z_2} \cdots \sum_{z_N} \prod_{n'} p(z_{n'}|X, \theta^{(s)}) \cdot z_{n,k} \cdot (\log(\pi_k) + \log(p(x_n|\theta_k)))$$

(Using the two formulas from the previous page)

$$\underbrace{\sum_Z p(Z|X, \theta^{(s)}) \cdot z_{n,k} \cdot (\log(\pi_k) + \log(p(x_n|\theta_k)))}_{\text{inner sum from previous}}$$

$$= \sum_{z_1} \sum_{z_2} \cdots \sum_{z_N} \prod_{n'} p(z_{n'}|X, \theta^{(s)}) \cdot z_{n,k} \cdot (\log(\pi_k) + \log(p(x_n|\theta_k)))$$

$$= \sum_{z_n} p(z_n|X, \theta^{(s)}) \cdot z_{n,k} \cdot (\log(\pi_k) + \log(p(x_n|\theta_k))) \underbrace{\sum_{z_1} \cdots \sum_{z_{n-1}} \sum_{z_{n+1}} \cdots \sum_{z_N} \prod_{n' \neq n} p(z_{n'}|X, \theta^{(s)})}_{\substack{\text{All possibilities without point } n. \\ \text{This entire mess evaluates to unity!}}}$$

(Here we move in all the sums that do not interact with the function we are taking the expectation over that do not interact with the outer n variable.)

$$\underbrace{\sum_Z p(Z|X, \theta^{(s)}) \cdot z_{n,k} \cdot (\log(\pi_k) + \log(p(x_n|\theta_k)))}_{\text{inner sum from previous}}$$

$$= \sum_{z_1} \sum_{z_2} \cdots \sum_{z_N} \prod_{n'} p(z_{n'}|X, \theta^{(s)}) \cdot z_{n,k} \cdot (\log(\pi_k) + \log(p(x_n|\theta_k)))$$

$$= \sum_{z_n} p(z_n|X, \theta^{(s)}) \cdot z_{n,k} \cdot (\log(\pi_k) + \log(p(x_n|\theta_k))) \underbrace{\sum_{z_1} \cdots \sum_{z_{n-1}} \sum_{z_{n+1}} \cdots \sum_{z_N} \prod_{n' \neq n} p(z_{n'}|X, \theta^{(s)})}_{\substack{\text{All possibilities without point } n. \\ \text{This entire mess evaluates to unity!}}}$$

$$= \sum_{z_n} p(z_n|X, \theta^{(s)}) \cdot z_{n,k} \cdot (\log(\pi_k) + \log(p(x_n|\theta_k)))$$

(As noted, the big sum on the right is unity. It is the probability of all possible configurations that do not involve point n . Since this covers all cases, it is one.)

$$\underbrace{\sum_Z p(Z|X, \theta^{(s)}) \cdot z_{n,k} \cdot (\log(\pi_k) + \log(p(x_n|\theta_k)))}_{\text{inner sum from previous}}$$

$$\begin{aligned} &= \sum_{z_1} \sum_{z_2} \cdots \sum_{z_N} \prod_{n'} p(z_{n'}|X, \theta^{(s)}) \cdot z_{n,k} \cdot (\log(\pi_k) + \log(p(x_n|\theta_k))) \\ &= \sum_{z_n} p(z_n|X, \theta^{(s)}) \cdot z_{n,k} \cdot (\log(\pi_k) + \log(p(x_n|\theta_k))) \underbrace{\sum_{z_1} \cdots \sum_{z_{n-1}} \sum_{z_{n+1}} \cdots \sum_N \prod_{n' \neq n} p(z_{n'}|X, \theta^{(s)})}_{\substack{\text{All possibilities without point } n. \\ \text{This entire mess evaluates to unity!}}} \\ &= \sum_{z_n} p(z_n|X, \theta^{(s)}) \cdot z_{n,k} \cdot (\log(\pi_k) + \log(p(x_n|\theta_k))) \\ &= p(z_{n,k} = 1|X, \theta^{(s)}) (\log(\pi_k) + \log(p(x_n|\theta_k))) \end{aligned}$$

(Here, remember that k is sitting outside the sum. The indicator variable $z_{n,k}$, selects the probability for k from the sum over z_n , which is a sum over clusters.)

$$\underbrace{\sum_Z p(Z|X, \theta^{(s)}) \cdot z_{n,k} \cdot (\log(\pi_k) + \log(p(x_n|\theta_k)))}_{\text{inner sum from previous}}$$

$$\begin{aligned} &= \sum_{z_1} \sum_{z_2} \cdots \sum_{z_N} \prod_{n'} p(z_{n'}|X, \theta^{(s)}) \cdot z_{n,k} \cdot (\log(\pi_k) + \log(p(x_n|\theta_k))) \\ &= \sum_{z_n} p(z_n|X, \theta^{(s)}) \cdot z_{n,k} \cdot (\log(\pi_k) + \log(p(x_n|\theta_k))) \underbrace{\sum_{z_1} \cdots \sum_{z_{n-1}} \sum_{z_{n+1}} \cdots \sum_N \prod_{n' \neq n} p(z_{n'}|X, \theta^{(s)})}_{\substack{\text{All possibilities without point } n. \\ \text{This entire mess evaluates to unity!}}} \\ &= \sum_{z_n} p(z_n|X, \theta^{(s)}) \cdot z_{n,k} \cdot (\log(\pi_k) + \log(p(x_n|\theta_k))) \\ &= p(z_{n,k} = 1|X, \theta^{(s)}) (\log(\pi_k) + \log(p(x_n|\theta_k))) \\ &= \gamma(z_{n,k}) (\log(\pi_k) + \log(p(x_n|\theta_k))) \quad (\text{definition of } \gamma(z_{n,k}), \text{ the responsibility}) \end{aligned}$$

Deriving the M-step

$$\begin{aligned} Q(\theta^{(s+1)}, \theta^{(s)}) &= \sum_Z p(Z|X, \theta^{(s)}) \sum_n \sum_k \left\{ z_{n,k} (\log(\pi_k) + \log(p(x_n|\theta_k))) \right\} \\ &= \sum_n \sum_k \left\{ \gamma(z_{n,k}) (\log(\pi_k) + \log(p(x_n|\theta_k))) \right\} \end{aligned}$$

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

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

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

Example—deriving the GMM M-step

$$\begin{aligned} \frac{\delta}{\delta \mu_k} Q(\theta^{(s+1)}, \theta^{(s)}) &= \sum_n \left\{ \gamma(z_{n,k}) \frac{\delta}{\delta \mu_k} (\log(\pi_k) + \log(p(x_n|\theta_k))) \right\} \\ &= \sum_n \left\{ \gamma(z_{n,k}) \frac{\delta}{\delta \mu_k} (\log(p(x_n|\theta_k))) \right\} \\ &\quad \sum_n \left\{ \gamma(z_{n,k}) \frac{\delta}{\delta \mu_k} (\log(N(x_n|\mu_k, \Sigma_k))) \right\} \end{aligned}$$

Example—deriving the GMM M-step

$$N(\mathbf{x}_n | \mu_k, \Sigma_k) = \frac{1}{(2\pi)^{D/2} |\Sigma_k|^{1/2}} \exp\left(-\frac{1}{2}(\mathbf{x}_n - \mu_k)^T \Sigma_k^{-1} (\mathbf{x}_n - \mu_k)\right)$$

$$\log(N(\mathbf{x}_n | \mu_k, \Sigma_k)) = \log\left(\frac{1}{(2\pi)^{D/2} |\Sigma_k|^{1/2}}\right) - \frac{1}{2}(\mathbf{x}_n - \mu_k)^T \Sigma_k^{-1} (\mathbf{x}_n - \mu_k)$$

$$\frac{\delta}{\delta \mu_k} \log(N(\mathbf{x}_n | \mu_k, \Sigma_k)) = \Sigma_k^{-1} (\mathbf{x}_n - \mu_k)$$

Example—deriving the GMM M-step

$$\frac{\delta}{\delta \mu_k} Q(\theta^{(s+1)}, \theta^{(s)}) = \sum_n \left\{ \gamma(z_{n,k}) \frac{\delta}{\delta \mu_k} \left(\log(N(x_n | \mu_k, \Sigma_k)) \right) \right\}$$

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

$$\sum_n \left\{ \gamma(z_{n,k}) \Sigma_k^{-1} (\mathbf{x}_n - \mu_k) \right\} = 0 \quad (\text{Inner sigma is a matrix, not a sum}).$$

$$\sum_n \left\{ \gamma(z_{n,k}) (\mathbf{x}_n - \mu_k) \right\} = 0 \quad (\text{Multiply by } \Sigma_k^{-1})$$

Example—deriving the GMM M-step

$$\text{So, } \sum_n \left\{ \gamma(z_{n,k}) (\mathbf{x}_n - \mu_k) \right\} = 0$$

$$\text{and } \mu_k \sum_n \left\{ \gamma(z_{n,k}) \right\} = \sum_n \left\{ \gamma(z_{n,k}) (\mathbf{x}_n) \right\}$$

$$\text{and } \mu_k = \frac{\sum_n \left\{ \gamma(z_{n,k}) (\mathbf{x}_n) \right\}}{\sum_n \left\{ \gamma(z_{n,k}) \right\}} \quad (\text{same as before})$$