

Statistical clustering

Basic, very general, generative model

$$p(\mathbf{x}) = \sum_k p(k) p(\mathbf{x}|k)$$

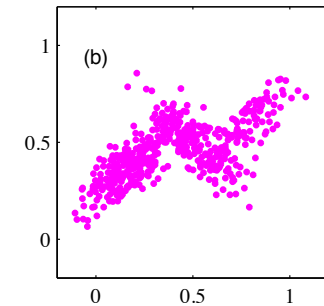
Data density, given the above model.

$$p(\{\mathbf{x}_i\}) = \prod_i \left\{ \sum_k p(k) p(\mathbf{x}_i|k) \right\}$$

Example: Three bivariate Gaussians

Points sampled according to:

$$p(\mathbf{x}) = \sum_k p(k) p(\mathbf{x}|k)$$



Even if we know the exact model, we cannot be sure from the data which point comes from which cluster. We only have the distribution for this.

EM for statistical clustering

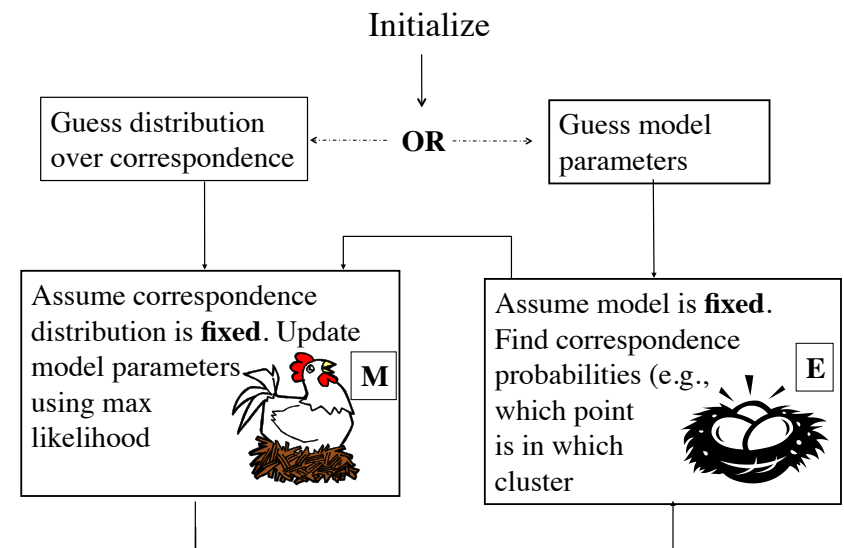
Given the number of clusters and some data, we fit the cluster parameters by maximizing the objective function:

$$\log(p(\{\mathbf{x}_i\})) = \sum_i \log \left\{ \sum_k p(k) p(\mathbf{x}_i|k) \right\}$$

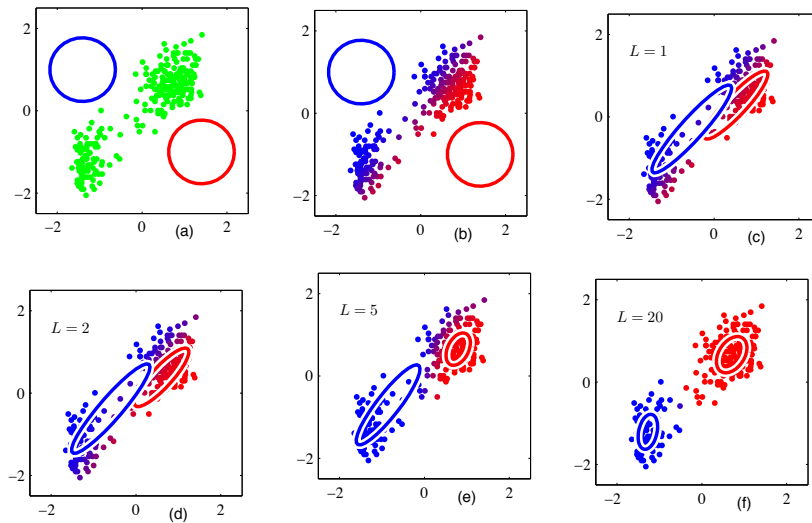
This is generally intractable. (We can only find a local maximum).

We do this with Expectation Maximization (EM).

EM flow chart



EM illustrated



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

EM (more formally)

- Denote cluster priors by:

$$\pi_k \equiv p(z_k = 1)$$

- 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)

- Responsibilities For GMM:

$$\gamma(z_{n,k}) = \frac{\pi_k N(x_n | u_k^{(s)}, \Sigma_k^{(s)})}{\sum_{k'} \pi_{k'} N(x_n | u_{k'}^{(s)}, \Sigma_{k'}^{(s)})}$$

EM (more formally)

Represent the entire data set of N points, \mathbf{x}_n ,
with 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 (everything is known).
The observed data, $\{X\}$, is called the *incomplete* data set.

EM (more formally)

We assume that computing the MLE of parameters,

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

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)

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 over $p(Z | X, \theta^{(s)})$).

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

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

EM (more formally)

Notice the complexity of the incomplete log likelihood:

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

The 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

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

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

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 $\theta^{(s+1)} = \arg \max_{\theta} \{Q(\theta^{(s+1)}, \theta^{(s)})\}$
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 is increases unless it is at a local maximum. It is important to check this is

GMM M-step

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

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

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

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

GMM M-step

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

(Intuitive?)

GMM M-step

$$\begin{aligned} &\sum_Z p(Z|X, \theta^{(s)}) f(n, k) \\ &= \sum_Z \prod_{n'} p(z_{n'}|X, \theta^{(s)}) f(n, k) \quad (\text{independence!}) \\ &= \sum_{z_n} p(z_n|X, \theta^{(s)}) f(n, k) \underbrace{\sum_{z_1} \dots \sum_{z_{n-1}} \sum_{z_{n+1}} \dots \sum_N \prod_{n' \neq n} p(z_{n'}|X, \theta^{(s)})}_{\text{all possibilities without point n.}} \\ &= \sum_{z_n} p(z_n|X, \theta^{(s)}) f(n, k) \end{aligned}$$

GMM M-step

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

(Here we use the special nature of the indicator variables.

All possible values of z_n are the K cluster assignments for point n .)

$$= \sum_n \sum_k \left\{ \gamma(n, k) \left(\log(\pi_k) + \log(p(x_n | \theta_k)) \right) \right\}$$

(We know that $z_{n,k}$ is zero for $k' \neq k$)

GMM M-step

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

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