

Mixture models and the EM algorithm

Qingyun Yu, Ke Zhang

School of Mathematics
Sun Yat-sen University

September 23, 2020

Table of Contents

1 Mixture models

2 The EM algorithm

3 Appendix

- Suppose $z_i \in \{1, \dots, K\}$, representing a discrete latent state, $p(z_i) = \text{Multi}(\pi)$. For the likelihood, we use $p(\mathbf{x}_i | z_i = k) = p_k(\mathbf{x}_i)$, where p_k is the k 'th base distribution for the observations.
- We are mixing together the K base distributions as follows:

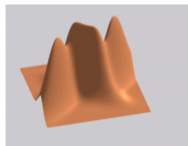
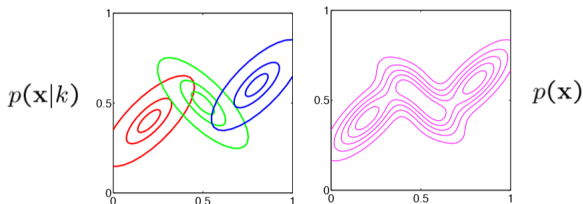
$$p(\mathbf{x}) = \sum_{k=1}^K \pi_k p_k(\mathbf{x}).$$

This is a convex combination of the p_k 's, π_k satisfy $0 \leq \pi_k \leq 1$ and $\sum_{k=1}^K \pi_k = 1$.

Mixtures of Gaussians(GMM)

- Each base distribution in the mixture is a multivariate Gaussian with mean μ_k and covariance matrix Σ_k .

$$p(\mathbf{x}_i | \theta) = \sum_{k=1}^K \pi_k \mathcal{N}(\mathbf{x}_i | \mu_k, \Sigma_k).$$



$p(\mathbf{x})$

Using mixture models for clustering

- Use them for clustering compute $p(z_i = k \mid \mathbf{x}_i, \boldsymbol{\theta})$, which represents the posterior probability that point i belongs to cluster k .

$$r_{ik} \triangleq p(z_i = k \mid \mathbf{x}_i, \boldsymbol{\theta}) \propto p(z_i = k \mid \boldsymbol{\theta}) p(\mathbf{x}_i \mid z_i = k, \boldsymbol{\theta}),$$

Assuming $1 - \max_k r_{ik}$ is small, it may be reasonable to compute a hard clustering using the MAP estimate, given by,

$$z_i^* = \arg \max_k r_{ik} = \arg \max_k \log p(\mathbf{x}_i \mid z_i = k, \boldsymbol{\theta}) + \log p(z_i = k \mid \boldsymbol{\theta}).$$

Note that z_i^* is designed on $\boldsymbol{\theta}$.

Assume that the conditional probability density functions $p(\vec{x} \mid y = k), k = 1, 2, \dots, K$ are both normally distributed with mean and covariance parameters $\theta = (\vec{\mu}_k, \Sigma)$, respectively.

The Bayes optimal solution:

$$p(y = c \mid \mathbf{x}, \theta) \propto \pi_c \exp \left[\mu_c^T -1 \mathbf{x} - \frac{1}{2} \mathbf{x}^T -1 \mathbf{x} - \frac{1}{2} \mu_c^T -1 \mu_c \right].$$

To class a new sample, we need to calculate the parameters θ .

If label missing?

Clustering.

An example for clustering

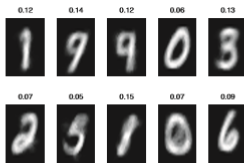


Figure 11.5 We fit a mixture of 10 Bernoullis to the binarized MNIST digit data. We show the MLE for the corresponding cluster means, μ_k . The numbers on top of each image represent the mixing weights π_k . No labels were used when training the model. Figure generated by `mixBerMoistEM`.

suppose our data consist of D -dimensional bit vectors. In this case, an appropriate classconditional density is a product of Bernoullis:

$$p(\mathbf{x}_i \mid z_i = k, \theta) = \prod_{j=1}^D \text{Ber}(x_{ij} \mid \mu_{jk}) = \prod_{j=1}^D \mu_{jk}^{x_{ij}} (1 - \mu_{jk})^{1-x_{ij}}.$$

where μ_{jk} is the probability that bit j turns on in cluster k .

The latent variables do not have to any meaning, we might simply introduce latent variables in order to make the model more powerful. For example, one can show that the mean and covariance of the mixture distribution are given by

$$\begin{aligned}\mathbb{E}[\mathbf{x}] &= \sum_k \pi_k \boldsymbol{\mu}_k, \\ \text{cov}[\mathbf{x}] &= \sum_k \pi_k [\boldsymbol{\Sigma}_k + \boldsymbol{\mu}_k \boldsymbol{\mu}_k^T] - \mathbb{E}[\mathbf{x}] \mathbb{E}[\mathbf{x}]^T,\end{aligned}$$

where $\boldsymbol{\Sigma}_k = \text{diag}(\mu_{jk}(1 - \mu_{jk}))$.

Mixtures of experts

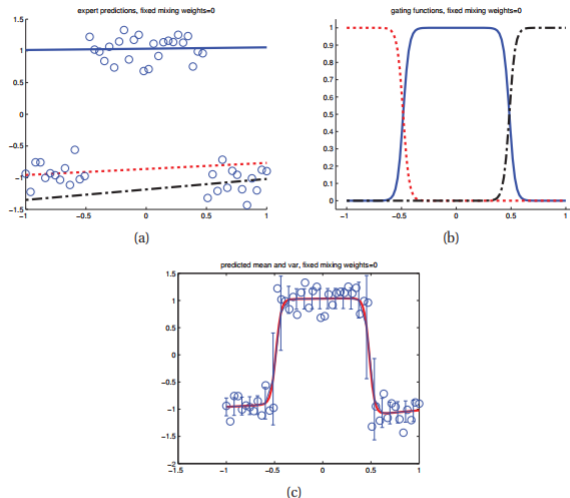


Figure 11.6 (a) Some data fit with three separate regression lines. (b) Gating functions for three different “experts”. (c) The conditionally weighted average of the three expert predictions. Figure generated by mixexpDemo.

Mixtures of experts

We can model this by allowing the mixing weights and the mixture densities to be input-dependent:

$$p(y_i | \mathbf{x}_i, z_i = k, \boldsymbol{\theta}) = \mathcal{N}(y_i | \mathbf{w}_k^T \mathbf{x}_i, \sigma_k^2),$$
$$p(z_i | \mathbf{x}_i, \boldsymbol{\theta}) = \text{Multi}(z_i | \mathcal{S}(\mathbf{V}^T \mathbf{x}_i)).$$

The function $p(z_i = k | \mathbf{x}_i, \boldsymbol{\theta})$ is called a gating function, and decides which expert to use, depending on the input values. The overall prediction of the model, obtained using

$$p(y_i | \mathbf{x}_i, \boldsymbol{\theta}) = \sum_k p(z_i = k | \mathbf{x}_i, \boldsymbol{\theta}) p(y_i | \mathbf{x}_i, z_i = k, \boldsymbol{\theta}).$$

Mixtures of experts

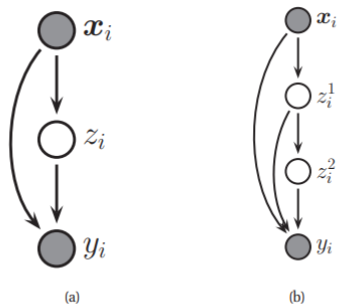


Figure 11.7 (a) A mixture of experts. (b) A hierarchical mixture of experts.

Parameter estimation for mixture models

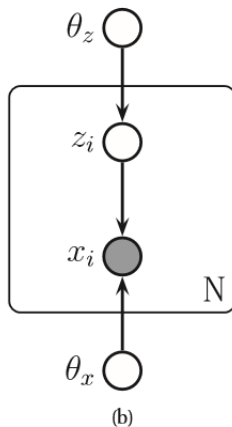
- If the z_i were observed, we see that $\theta_z \perp \theta_x \mid \mathcal{D}$ and hence the posterior will factorize.

$$L(\theta \mid \mathcal{D}) = \prod_{i=1}^N p(x_i \mid \theta_x, z_i).$$

- If the z_i are hidden, the parameters are no longer independent, and the posterior does not factorize, making it much harder to compute.

$$L(\theta \mid \mathcal{D}) = \prod_{i=1}^N \left\{ \sum_{k=1}^K p(z_i = k \mid \theta_z) p(x_i \mid \theta_x, z_i) \right\}.$$

Parameter estimation for mixture models



Unidentifiability

- Definition

Let $\mathcal{P} = \{P_\theta : \theta \in \Theta\}$ be a statistical model. We say that \mathcal{P} is identifiable if the mapping $\theta \mapsto P_\theta$ is one-to-one:

$$P_{\theta_1} = P_{\theta_2} \quad \Rightarrow \quad \theta_1 = \theta_2, \quad \text{for all } \theta_1, \theta_2 \in \Theta.$$

- Consider a GMM. If the z_i were all observed, we would have a unimodal posterior for the parameters:

$$L(\boldsymbol{\theta} \mid \mathcal{D}) = \prod_{i=1}^N p(z_i = k \mid \theta_z) N(x_i \mid \mu_k, \Sigma_k).$$

Consequently we can easily find the globally optimal MLE.

- suppose the z_i 's are hidden, for each of the possible ways of “filling in” the z_i 's, we get a different unimodal likelihood. ($K!$ possible labelings)

Unidentifiability

Consider a GMM,

$$p(\mathbf{x} \mid \boldsymbol{\theta}) = \pi_1 \mathcal{N}(\mathbf{x} \mid \boldsymbol{\mu}_1, \sigma_1^2) + \pi_2 \mathcal{N}(\mathbf{x} \mid \boldsymbol{\mu}_2, \sigma_2^2).$$

$$\pi_1 = \pi_2 = 0.5, \mu_1 = -10, \mu_2 = 10, \sigma_1 = \sigma_2 = 5.$$

$$\pi_1 = \pi_2 = 0.5, \mu_1 = 10, \mu_2 = -10, \sigma_1 = \sigma_2 = 5.$$

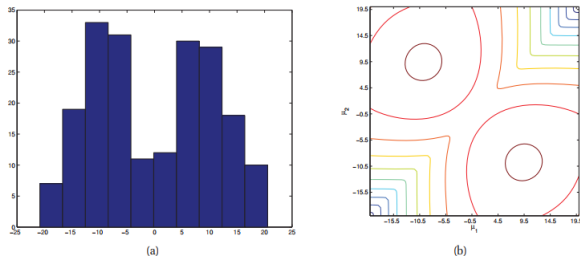


Figure 11.10 Left: $N = 200$ data points sampled from a mixture of 2 Gaussians in 1d, with $\pi_k = 0.5$, $\sigma_k = 5$, $\mu_1 = -10$ and $\mu_2 = 10$. Right: Likelihood surface $p(\mathcal{D} | \mu_1, \mu_2)$, with all other parameters set to their true values. We see the two symmetric modes, reflecting the unidentifiability of the parameters. Figure generated by `mixGaussLikSurfaceDemo`.

Computing a MAP estimate is non-convex

Now suppose the joint probability distribution $p(\mathbf{z}_i, \mathbf{x}_i | \boldsymbol{\theta})$ is in the exponential family, which means it can be written as follows:

$$p(\mathbf{x}, \mathbf{z} | \boldsymbol{\theta}) = \frac{1}{Z(\boldsymbol{\theta})} \exp [\boldsymbol{\theta}^T \boldsymbol{\phi}(\mathbf{x}, \mathbf{z})] ,$$

where $\boldsymbol{\phi}(\mathbf{x}, \mathbf{z})$ are the sufficient statistics, and $Z(\boldsymbol{\theta})$ is the normalization constant. the complete data log likelihood can be written as follows:

$$\ell(\boldsymbol{\theta} | \mathbf{x}, \mathbf{z}) = \sum_{i=1}^N \log p(\mathbf{x}_i, \mathbf{z}_i | \boldsymbol{\theta}) = \boldsymbol{\theta}^T \left(\sum_{i=1}^N \boldsymbol{\phi}(\mathbf{x}_i, \mathbf{z}_i) \right) - N \log Z(\boldsymbol{\theta}).$$

The first term is clearly linear in $\boldsymbol{\theta}$. One can show that $Z(\boldsymbol{\theta})$ is a convex function, so the overall objective is concave.

Computing a MAP estimate is non-convex

when we assumed z_i is missing. The observed data log likelihood is given by

$$\ell(\boldsymbol{\theta}) = \sum_{i=1}^N \log \sum_{\mathbf{z}_i} p(\mathbf{x}_i, \mathbf{z}_i \mid \boldsymbol{\theta}) = \sum_{i=1}^N \log \left[\sum_{\mathbf{z}_i} e^{\boldsymbol{\theta}^T \phi(\mathbf{z}_i, \mathbf{x}_i)} \right] - N \log Z(\boldsymbol{\theta}).$$

One can show that the log-sum-exp function is convex, and we know that $Z(\boldsymbol{\theta})$ is convex. However, the difference of two convex functions may be not convex.

The disadvantage of non-convex functions is that it is usually hard to find their global optimum.

Table of Contents

1 Mixture models

2 The EM algorithm

3 Appendix

The EM algorithm

One approach is to use a generic gradient-based optimizer to find a local minimum of the negative log likelihood or NLL, given by

$$-\ell(\boldsymbol{\theta} \mid \mathbf{x}), \text{ where } \mathbf{x} \text{ is the observed data.}$$

However, we often have to enforce constraints, such as the fact that covariance matrices must be positive definite(for GMM), mixing weights must sum to one, which can be tricky.

A much simpler algorithm — Expectation Maximization.

Theoretical basis for EM

Let \mathbf{x} be the visible or observed variables, and let \mathbf{z} be the hidden or missing variables. Consider an arbitrary distribution $q(\mathbf{z})$ over the hidden variables. The observed data log likelihood can be written as follows:

$$\begin{aligned}\ell(\boldsymbol{\theta} \mid \mathbf{x}) &= \log p(\mathbf{x} \mid \boldsymbol{\theta}) = \log \left[\int p(\mathbf{x}, \mathbf{z} \mid \boldsymbol{\theta}) d\mathbf{z} \right] \\ &= \log \left[\int q(\mathbf{z}) \frac{p(\mathbf{x}, \mathbf{z} \mid \boldsymbol{\theta})}{q(\mathbf{z})} d\mathbf{z} \right]\end{aligned}$$

Now $\log(u)$ is a concave function, so from Jensen's inequality (Equation 2.113) we have the following lower bound:

$$\ell(\boldsymbol{\theta} \mid \mathbf{x}) = \log \mathbb{E}_{\mathbf{z}} \frac{p(\mathbf{x}, \mathbf{z} \mid \boldsymbol{\theta})}{q(\mathbf{z})} \geq \mathbb{E}_{\mathbf{z}} \log \frac{p(\mathbf{x}, \mathbf{z} \mid \boldsymbol{\theta})}{q(\mathbf{z})}, \text{ where } \mathbf{z} \sim q.$$

The above argument holds for any positive distribution q . Which one should we choose? Intuitively we should pick the q that yields the tightest lower bound.

Theoretical basis for EM

Let $\mathcal{F}(\boldsymbol{\theta}, q) := \mathbb{E}_z \log \frac{p(\mathbf{x}, \mathbf{z} | \boldsymbol{\theta})}{q(\mathbf{z})}$, then

$$\begin{aligned}\mathcal{F}(\boldsymbol{\theta}, q) &= \int q(\mathbf{z}) \log \frac{p(\mathbf{x}, \mathbf{z} | \boldsymbol{\theta})}{q(\mathbf{z})} d\mathbf{z} \\&= \int q(\mathbf{z}) \log \frac{p(\mathbf{z} | \mathbf{x}, \boldsymbol{\theta}) p(\mathbf{x} | \boldsymbol{\theta})}{q(\mathbf{z})} d\mathbf{z} \\&= \int q(\mathbf{z}) \log \frac{p(\mathbf{z} | \mathbf{x}, \boldsymbol{\theta})}{q(\mathbf{z})} d\mathbf{z} + \int q(\mathbf{z}) \log p(\mathbf{x} | \boldsymbol{\theta}) d\mathbf{z} \\&= -\mathbb{KL}(q(\mathbf{z}) \| p(\mathbf{z} | \mathbf{x}, \boldsymbol{\theta})) + \log p(\mathbf{x} | \boldsymbol{\theta}),\end{aligned}$$

which means that

$$\ell(\boldsymbol{\theta} | \mathbf{x}) = \mathcal{F}(\boldsymbol{\theta}, q) + \mathbb{KL}(q(\mathbf{z}) \| p(\mathbf{z} | \mathbf{x}, \boldsymbol{\theta}))$$

Lower bound approximation

Given $\theta^{(k-1)}$, define

$$q^{(k)} = \arg \max_q \mathcal{F}(q, \theta^{(k-1)}),$$

where $\mathcal{F}(q^{(k)}, \theta^{(k-1)})$ could be viewed as the "optimal" approximation of $\ell(\theta^{(k-1)} | \mathbf{x})$

After approximation, we could get $\theta^{(k)} = \arg \max_{\theta} \mathcal{F}(\theta, q^{(k)})$.

If the gap between $\mathcal{F}(\theta^{(k-1)}, q^{(k)})$ and $\ell(\theta^{(k-1)} | \mathbf{x})$ is small enough, we could expect that $\ell(\theta^{(k)} | \mathbf{x}) \geq \ell(\theta^{(k-1)} | \mathbf{x})$, since

$$\ell(\theta^{(k)} | \mathbf{x}) \geq \mathcal{F}(\theta^{(k)}, q^{(k)}) \geq \mathcal{F}(\theta^{(k-1)}, q^{(k)}) \approx \ell(\theta^{(k-1)} | \mathbf{x})$$

Lower bound approximation

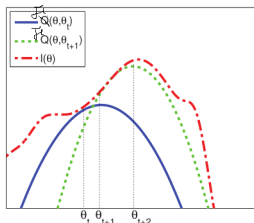


Figure 11.16 Illustration of EM as a bound optimization algorithm. Based on Figure 9.14 of (Bishop 2006a). Figure generated by `emLogLikelihoodMax`.

Lower bound approximation

Now we observed that

$$\begin{aligned} q^{(k)} &= \arg \max_q \mathcal{F} \left(q, \boldsymbol{\theta}^{(k-1)} \right) \\ &\propto \arg \max_q -\mathbb{KL} \left(q(\mathbf{z}) \parallel p \left(\mathbf{z} \mid \mathbf{x}, \boldsymbol{\theta}^{(k-1)} \right) \right) \\ &= p \left(\mathbf{z} \mid \mathbf{x}, \boldsymbol{\theta}^{(k-1)} \right), \end{aligned}$$

which means that

$$\begin{aligned} \boldsymbol{\theta}^k &= \arg \max_{\boldsymbol{\theta}} \mathcal{F} \left(q^{(k)}, \boldsymbol{\theta}^{(k-1)} \right) \\ &= \arg \max_{\boldsymbol{\theta}} \mathbb{E}_{\mathbf{z}} \log \frac{p(\mathbf{x}, \mathbf{z} \mid \boldsymbol{\theta})}{q^{(k)}(\mathbf{z})} \\ &\propto \arg \max_{\boldsymbol{\theta}} \mathbb{E}_{\mathbf{z}} \log p(\mathbf{x}, \mathbf{z} \mid \boldsymbol{\theta}), \text{ where } \mathbf{z} \sim q \left(\mathbf{z} \mid \mathbf{x}, \boldsymbol{\theta}^{(k-1)} \right) \\ &=: \arg \max_{\boldsymbol{\theta}} Q \left(\boldsymbol{\theta}, \boldsymbol{\theta}^{(k-1)} \right) \end{aligned}$$

This is so called EM

Algorithm:

1. Start with initial guesses for the parameter $\theta^{(0)}$.
2. Expectation Step: at the t th iteration, compute $Q(\theta, \theta^{(k-1)})$.
3. Maximization Step: determine the new estimate $\theta^{(k)} = \arg \max_{\theta} Q(\theta, \theta^{(k-1)})$.
4. Repeat steps 2 and 3 until convergence.

The expected complete data log likelihood is given by

$$\begin{aligned} Q(\boldsymbol{\theta}, \boldsymbol{\theta}^{(t-1)}) &\triangleq \mathbb{E} \left[\sum_i \log p(\mathbf{x}_i, z_i | \boldsymbol{\theta}) \right] \\ &= \sum_i \mathbb{E} \left[\log \left[\prod_{k=1}^K (\pi_k p(\mathbf{x}_i | \boldsymbol{\theta}_k))^{I(z_i=k)} \right] \right] \\ &= \sum_i \sum_k \mathbb{E} [\mathbb{I}(z_i = k)] \log [\pi_k p(\mathbf{x}_i | \boldsymbol{\theta}_k)] \\ &= \sum_i \sum_k p(z_i = k | \mathbf{x}_i, \boldsymbol{\theta}^{(t-1)}) \log [\pi_k p(\mathbf{x}_i | \boldsymbol{\theta}_k)] \\ &= \sum_i \sum_k r_{ik} \log \pi_k + \sum_i \sum_k r_{ik} \log p(\mathbf{x}_i | \boldsymbol{\theta}_k) \end{aligned}$$

where $r_{ik} \triangleq p(z_i = k | \mathbf{x}_i, \boldsymbol{\theta}^{(t-1)})$ is the responsibility that cluster k takes for data point i .

- E step: The E step has the following simple form, which is the same for any mixture model:

$$r_{ik} = \frac{\pi_k p(\mathbf{x}_i | \theta_k^{(t-1)})}{\sum_{k'} \pi_{k'} p(\mathbf{x}_i | \theta_{k'}^{(t-1)})}.$$

- M step: In the M step, we optimize π and the θ_k of Q function. For π , we obviously have

$$\pi_k = \frac{1}{N} \sum_i r_{ik} = \frac{r_k}{N},$$

where $r_k \triangleq \sum_i r_{ik}$ is the weighted number of points assigned to cluster k.

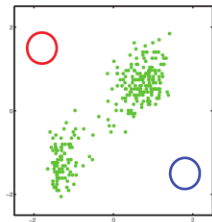
To derive the M step for the μ_k and Σ_k terms, we look at the parts of Q that depend on μ_k and Σ_k ,

$$\begin{aligned}\ell(\mu_k, \Sigma_k) &= \sum_k \sum_i r_{ik} \log p(\mathbf{x}_i | \theta_k) \\ &= -\frac{1}{2} \sum_i r_{ik} \left[\log |\Sigma_k| + (\mathbf{x}_i - \mu_k)^T \Sigma_k^{-1} (\mathbf{x}_i - \mu_k) \right]\end{aligned}$$

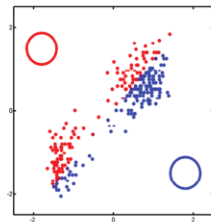
This is just a weighted version of the standard problem of computing the MLEs of an MVN (see Section 4.1.3). the new parameter estimates are given by

$$\begin{aligned}\mu_k &= \frac{\sum_i r_{ik} \mathbf{x}_i}{r_k} \\ \Sigma_k &= \frac{\sum_i r_{ik} (\mathbf{x}_i - \mu_k)(\mathbf{x}_i - \mu_k)^T}{r_k} = \frac{\sum_i r_{ik} \mathbf{x}_i \mathbf{x}_i^T}{r_k} - \mu_k \mu_k^T.\end{aligned}$$

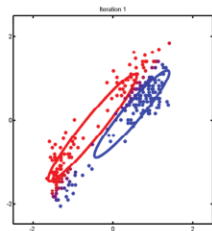
Example: EM for GMMs



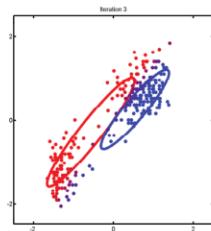
(a)



(b)



(c)



(d)

Example: EM for GMMs

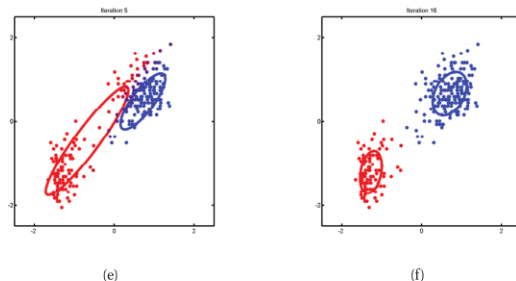


Figure 11.11 Illustration of the EM for a GMM applied to the Old Faithful data. (a) Initial (random) values of the parameters. (b) Posterior responsibility of each point computed in the first E step. The degree of redness indicates the degree to which the point belongs to the red cluster, and similarly for blue; this purple points have a roughly uniform posterior over clusters. (c) We show the updated parameters after the first M step. (d) After 3 iterations. (e) After 5 iterations. (f) After 16 iterations. Based on (Bishop 2006a) Figure 9.8. Figure generated by `mixGaussDemoFaithful`.

K-means algorithm

A popular variant of the EM algorithm for GMMs known as the K-means algorithm.

Consider a GMM in which we assume $\Sigma_k = \sigma^2 \mathbf{I}_D$ and $\pi_k = 1/K$ are fixed, only $\mu_k \in \mathbb{R}^D$, have to be estimated.

During the E step:

$$p(z_i = k \mid \mathbf{x}_i, \theta) \approx \mathbb{I}(k = z_i^*)$$

where $z_i^* = \operatorname{argmax}_k p(z_i = k \mid \mathbf{x}_i, \theta)$. (hard **EM**).

Since we assumed an equal spherical covariance matrix for each cluster,

$$z_i^* = \operatorname{argmin}_k \|\mathbf{x}_i - \mu_k\|_2^2.$$

The M step:

$$\mu_k = \frac{1}{N_k} \sum_{i: z_i=k} \mathbf{x}_i.$$

K-means algorithm

Algorithm 11.1: K-means algorithm

```
1 initialize  $\mathbf{m}_k$ ;  
2 repeat  
3   | Assign each data point to its closest cluster center:  $z_i = \arg \min_k \|\mathbf{x}_i - \boldsymbol{\mu}_k\|_2^2$ ;  
4   | Update each cluster center by computing the mean of all points assigned to it:  
   |  $\boldsymbol{\mu}_k = \frac{1}{N_k} \sum_{i: z_i=k} \mathbf{x}_i$ ;  
5 until converged;
```

Example: K-means algorithm

1. Starts with guesses for the three cluster centers.
2. For each data point, the closest cluster center (in Euclidean distance) is identified;
3. Each cluster center is replaced by the coordinate-wise average of all data points that are closest to it.
4. Repeat 2 and 3.

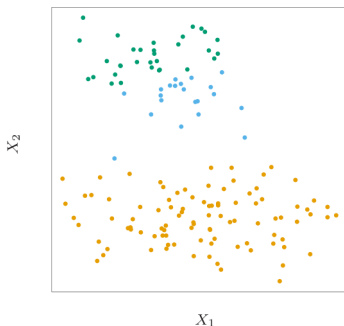


FIGURE 14.4. Simulated data in the plane, clustered into three classes (represented by orange, blue and green) by the K-means clustering algorithm

Example: K-means algorithm

Some of the K-means iterations for the simulated data of Figure 14.4

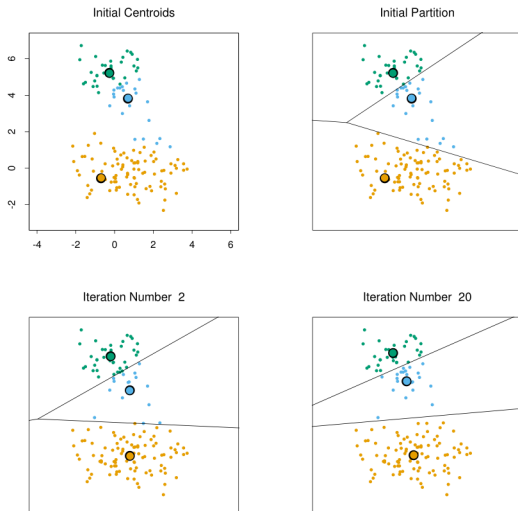


Table of Contents

1 Mixture models

2 The EM algorithm

3 Appendix

EM for mixture of experts

We can fit a mixture of experts model using EM in a straightforward manner. The expected complete data log likelihood is given by

$$Q(\boldsymbol{\theta}, \boldsymbol{\theta}^{old}) = \sum_{i=1}^N \sum_{k=1}^K r_{ik} \log [\pi_{ik} \mathcal{N}(y_i | \mathbf{w}_k^T \mathbf{x}_i, \sigma_k^2)]$$

$$\pi_{i,k} \triangleq \mathcal{S}(\mathbf{V}^T \mathbf{x}_i)_k$$

$$r_{ik} \propto \pi_{ik}^{old} \mathcal{N}\left(y_i | \mathbf{x}_i^T \mathbf{w}_k^{old}, (\sigma_k^{old})^2\right)$$

So the E step is the same as in a standard mixture model, except we have to replace π_k with $\pi_{i,k}$ when computing r_{ik} .

EM for mixture of experts

In the M step, we need to maximize $Q(\boldsymbol{\theta}, \boldsymbol{\theta}^{\text{old}})$ wrt \mathbf{w}_k, σ_k^2 and \mathbf{V} . For the regression parameters for model k, the objective has the form

$$Q(\boldsymbol{\theta}_k, \boldsymbol{\theta}^{\text{old}}) = \sum_{i=1}^N r_{ik} \left\{ -\frac{1}{\sigma_k^2} (y_i - \mathbf{w}_k^T \mathbf{x}_i)^2 \right\}$$

we can immediately write down the MLE as $\mathbf{w}_k = (\mathbf{X}^T \mathbf{R}_k \mathbf{X})^{-1} \mathbf{X}^T \mathbf{R}_k \mathbf{y}$, where $\mathbf{R}_k = \text{diag}(r_{:,k})$. The MLE for the variance is given by

$$\sigma_k^2 = \frac{\sum_{i=1}^N r_{ik} (y_i - \mathbf{w}_k^T \mathbf{x}_i)^2}{\sum_{i=1}^N r_{ik}}$$

We replace the estimate of the unconditional mixing weights π with the estimate of the gating parameters, \mathbf{V} . The objective has the form

$$\ell(\mathbf{V}) = \sum_i \sum_k r_{ik} \log \pi_{i,k}$$

EM for DGMs with hidden variables

For simplicity of presentation, we will assume all CPDs are tabular. Based on Section 10.4.2, let us write each CPT as follows:

$$p(x_{it} \mid \mathbf{x}_{i,\text{pa}(t)}, \boldsymbol{\theta}_t) = \prod_{c=1}^{K_{\text{pa}(t)}} \prod_{k=1}^{K_t} \theta_{tck}^{\mathbb{I}(x_{it}=i, \mathbf{x}_{i,\text{pa}(t)}=c)}$$

The log-likelihood of the complete data is given by

$$\log p(\mathcal{D} \mid \boldsymbol{\theta}) = \sum_{t=1}^V \sum_{c=1}^{K_{\text{pa}(t)}} \sum_{k=1}^{K_t} N_{tck} \log \theta_{tck}$$

where $N_{tck} = \sum_{i=1}^N \mathbb{I}(x_{it} = i, \mathbf{x}_{i,\text{pa}(t)} = c)$ are the empirical counts. Hence the expected complete data log-likelihood has the form,

$$\mathbb{E}[\log p(\mathcal{D} \mid \boldsymbol{\theta})] = \sum_t \sum_c \sum_k \bar{N}_{tck} \log \theta_{tck}$$

EM for DGMs with hidden variables

where

$$\bar{N}_{tck} = \sum_{i=1}^N \mathbb{E} [\mathbb{I}(x_{it} = i, \mathbf{x}_{i,\text{pa}(t)} = c)] = \sum_i p(x_{it} = k, \mathbf{x}_{i,\text{pa}(t)} = c \mid \mathcal{D}_i)$$

where \mathcal{D}_i are all the visible variables in case i . The \bar{N}_{tck} are the expected sufficient statistics, and constitute the output of the E step. Given these ESS, the M step has the simple form

$$\hat{\theta}_{tck} = \frac{\bar{N}_{tck}}{\sum_{k'} \bar{N}_{tjk'}}$$

This can be proved by adding Lagrange multipliers (to enforce the constraint $\sum_k \theta_{tjk} = 1$) to the expected complete data log likelihood, and then optimizing each parameter vector θ_{tc} separately.

Batch EM review

Before explaining online EM, we review batch EM in a more abstract setting. Let $\phi(\mathbf{x}, \mathbf{z})$ be a vector of sufficient statistics for a single data case. Let $\mathbf{s}_i = \sum_{\mathbf{z}} p(\mathbf{z} | \mathbf{x}_i, \boldsymbol{\theta}) \phi(\mathbf{x}_i, \mathbf{z})$ be the expected sufficient statistics for case i , and $\boldsymbol{\mu} = \sum_{i=1}^N \mathbf{s}_i$ be the sum of the ESS. Given $\boldsymbol{\mu}$, we can derive an ML or MAP estimate of the parameters in the M step; we will denote this operation by $\boldsymbol{\theta}(\boldsymbol{\mu})$.

Algorithm 11.2: Batch EM algorithm

```
1 initialize  $\boldsymbol{\mu}$ ;  
2 repeat  
3    $\boldsymbol{\mu}^{new} = \mathbf{0}$  ;  
4   for each example  $i = 1 : N$  do  
5      $\mathbf{s}_i := \sum_{\mathbf{z}} p(\mathbf{z} | \mathbf{x}_i, \boldsymbol{\theta}(\boldsymbol{\mu})) \phi(\mathbf{x}_i, \mathbf{z})$  ;  
6      $\boldsymbol{\mu}^{new} := \boldsymbol{\mu}^{new} + \mathbf{s}_i$  ;  
7    $\boldsymbol{\mu} := \boldsymbol{\mu}^{new}$  ;  
8 until converged;
```

Incremental EM

In incremental EM (Neal and Hinton 1998), we keep track of μ as well as the s_i . When we come to a data case, we swap out the old s_i and replace it with the new s_i^{new} .

Algorithm 11.3: Incremental EM algorithm

```
1 initialize  $s_i$  for  $i = 1 : N$ ;  
2  $\mu = \sum_i s_i$ ;  
3 repeat  
4   for each example  $i = 1 : N$  in a random order do  
5      $s_i^{new} := \sum_{\mathbf{z}} p(\mathbf{z} | \mathbf{x}_i, \theta(\mu)) \phi(\mathbf{x}_i, \mathbf{z})$  ;  
6      $\mu := \mu + s_i^{new} - s_i$ ;  
7      $s_i := s_i^{new}$ ;  
8 until converged;
```

This can be viewed as maximizing the lower bound $Q(\theta, q_1, \dots, q_N)$ by optimizing q_1 , then θ , then q_2 , then θ , etc. As such, this method is guaranteed to monotonically converge to a local maximum of the lower bound and to the log likelihood itself.

Stepwise EM

In stepwise EM, whenever we compute a new s_i , we move μ towards it, as shown in Algorithm 7. At iteration k , the stepsize has value η_k , which must satisfy the Robbins-Monro conditions in Equation 8.82. We can get somewhat better behavior by using a minibatch of size m before each update. It is possible to optimize m and k to maximize the training set likelihood, by trying different values in parallel for an initial trial period; this can significantly speed up the algorithm.

Algorithm 11.4: Stepwise EM algorithm

```
1 initialize  $\mu$ ;  $k = 0$  ;
2 repeat
3   for each example  $i = 1 : N$  in a random order do
4      $s_i := \sum_{\mathbf{z}} p(\mathbf{z}|\mathbf{x}_i, \theta(\mu)) \phi(\mathbf{x}_i, \mathbf{z})$  ;
5      $\mu := (1 - \eta_k)\mu + \eta_k s_i$ ;
6      $k := k + 1$ 
7 until converged;
```

Fitting models with missing data

Suppose we want to fit a joint density model by maximum likelihood, but we have “holes” in our data matrix, due to missing data (usually represented by NaNs). More formally, let $O_{ij} = 1$ if component j of data case i is observed, and let $O_{ij} = 0$ otherwise. Let $\mathbf{X}_v = \{x_{ij} : O_{ij} = 1\}$ be the visible data, and $O_{ij} = 0$ otherwise. Let $\mathbf{X}_h = \{x_{ij} : O_{ij} = 0\}$ be the missing or hidden data.

Our goal is to compute $\hat{\theta} = \underset{\theta}{\operatorname{argmax}} p(\mathbf{X}_v | \theta, \mathbf{O})$. Under the missing at random assumption (see Section 8.6.2), we have

$$p(\mathbf{X}_v | \theta, \mathbf{O}) = \prod_{i=1}^N p(\mathbf{x}_{iv} | \theta)$$

where \mathbf{x}_{iv} is a vector created from row i and the columns indexed by the set $\{j : O_{ij} = 1\}$.

Fitting models with missing data

Hence the log-likelihood has the form

$$\log p(\mathbf{X}_v | \boldsymbol{\theta}) = \sum_i \log p(\mathbf{x}_{iv} | \boldsymbol{\theta})$$

where

$$p(\mathbf{x}_{iv} | \boldsymbol{\theta}) = \sum_{\mathbf{x}_{ih}} p(\mathbf{x}_{iv}, \mathbf{x}_{ih} | \boldsymbol{\theta})$$

where \mathbf{x}_{iv} is the vector of hidden variables for case i (assumed discrete for notational simplicity). Substituting in, we get

$$\log p(\mathbf{X}_v | \boldsymbol{\theta}) = \sum_i \log \left[\sum_{\mathbf{x}_{ih}} p(\mathbf{x}_{iv}, \mathbf{x}_{ih} | \boldsymbol{\theta}) \right]$$

Unfortunately, this objective is hard to maximize. since we cannot push the log inside the sum. However, we can use the EM algorithm to compute a local optimum.

EM for the MLE of an MVN with missing data

Suppose we want to fit an MVN by maximum likelihood, but we have missing data. We can use EM to find a local maximum of the objective, as we explain below.

To get the algorithm started, we can compute the MLE based on those rows of the data matrix that are fully observed. If there are no such rows, we can use some ad-hoc imputation procedures, and then compute an initial MLE.

E step

Once we have θ^t , we can compute the expected complete data log likelihood at iteration t as follows:

$$\begin{aligned} Q(\theta, \theta^{t-1}) &= \mathbb{E} \left[\sum_{i=1}^N \log \mathcal{N}(\mathbf{x}_i \mid \boldsymbol{\mu}, \boldsymbol{\Sigma}) \mid \mathcal{D}, \theta^{t-1} \right] \\ &= -\frac{N}{2} \log |2\pi \boldsymbol{\Sigma}| - \frac{1}{2} \sum_i \mathbb{E} \left[(\mathbf{x}_i - \boldsymbol{\mu})^T \boldsymbol{\Sigma}^{-1} (\mathbf{x}_i - \boldsymbol{\mu}) \right] \\ &= -\frac{N}{2} \log |2\pi \boldsymbol{\Sigma}| - \frac{1}{2} \text{tr} \left(\boldsymbol{\Sigma}^{-1} \sum_i \mathbb{E} \left[(\mathbf{x}_i - \boldsymbol{\mu}) (\mathbf{x}_i - \boldsymbol{\mu})^T \right] \right) . \\ &= -\frac{N}{2} \log |\boldsymbol{\Sigma}| - \frac{N}{2} \log(2\pi) - \frac{1}{2} \text{tr} (\boldsymbol{\Sigma}^{-1} \mathbb{E}[\mathbf{S}(\boldsymbol{\mu})]) \end{aligned}$$

where

$$\mathbb{E}[\mathbf{S}(\boldsymbol{\mu})] \triangleq \sum_i \left(\mathbb{E} [\mathbf{x}_i \mathbf{x}_i^T] + \boldsymbol{\mu} \boldsymbol{\mu}^T - 2\boldsymbol{\mu} \mathbb{E} [\mathbf{x}_i]^T \right)$$

E step

We see that we need to compute $\sum_i \mathbb{E}[\mathbf{x}_i]$ and $\sum_i \mathbb{E}[\mathbf{x}_i \mathbf{x}_i^T]$ these are the expected sufficient statistics. Specifically, consider case i, where components v are observed and components h are unobserved. We have

$$\begin{aligned}\mathbf{x}_{ih} \mid \mathbf{x}_{iv}, \boldsymbol{\theta} &\sim \mathcal{N}(\mathbf{m}_i, \mathbf{V}_i) \\ \mathbf{m}_i &\triangleq \boldsymbol{\mu}_h + \boldsymbol{\Sigma}_{hv} \boldsymbol{\Sigma}_{vv}^{-1} (\mathbf{x}_{iv} - \boldsymbol{\mu}_v) \\ \mathbf{V}_i &\triangleq \boldsymbol{\Sigma}_{hh} - \boldsymbol{\Sigma}_{hv} \boldsymbol{\Sigma}_{vv}^{-1} \boldsymbol{\Sigma}_{vh}\end{aligned}$$

Hence the expected sufficient statistics are

$$\mathbb{E}[\mathbf{x}_i] = (\mathbb{E}[\mathbf{x}_{ih}]; \mathbf{x}_{iv}) = (\mathbf{m}_i; \mathbf{x}_{iv})$$

where we have assumed (without loss of generality) that the unobserved variables come before the observed variables in the node ordering.

E step

To compute $\sum_i \mathbb{E} [\mathbf{x}_i \mathbf{x}_i^T]$, we use the result that $\text{cov}[\mathbf{x}] = \mathbb{E} [\mathbf{x} \mathbf{x}^T] - \mathbb{E}[\mathbf{x}] \mathbb{E} [\mathbf{x}^T]$. Hence

$$\mathbb{E} [\mathbf{x}_i \mathbf{x}_i^T] = \mathbb{E} \left[\begin{pmatrix} \mathbf{x}_{ih} \\ \mathbf{x}_{iv} \end{pmatrix} \begin{pmatrix} \mathbf{x}_{ih}^T & \mathbf{x}_{iv}^T \end{pmatrix} \right] = \begin{pmatrix} \mathbb{E} [\mathbf{x}_{ih} \mathbf{x}_{ih}^T] & \mathbb{E} [\mathbf{x}_{ih}] \mathbf{x}_{iv}^T \\ \mathbf{x}_{iv} \mathbb{E} [\mathbf{x}_{ih}]^T & \mathbf{x}_{iv} \mathbf{x}_{iv}^T \end{pmatrix}$$

$$\mathbb{E} [\mathbf{x}_{ih} \mathbf{x}_{ih}^T] = \mathbb{E} [\mathbf{x}_{ih}] \mathbb{E} [\mathbf{x}_{ih}]^T + \mathbf{V}_i$$

M step

By solving $\nabla Q(\boldsymbol{\theta}, \boldsymbol{\theta}^{(t-1)}) = \mathbf{0}$, we can show that the M step is equivalent to plugging these ESS into the usual MLE equations to get

$$\begin{aligned}\boldsymbol{\mu}^t &= \frac{1}{N} \sum_i \mathbb{E}[\mathbf{x}_i] \\ \boldsymbol{\Sigma}^t &= \frac{1}{N} \sum_i \mathbb{E}[\mathbf{x}_i \mathbf{x}_i^T] - \boldsymbol{\mu}^t (\boldsymbol{\mu}^t)^T\end{aligned}$$

Thus we see that EM is not equivalent to simply replacing variables by their expectations and applying the standard MLE formula; that would ignore the posterior variance and would result in an incorrect estimate. Instead we must compute the expectation of the sufficient statistics, and plug that into the usual equation for the MLE.

MAP estimation

As usual, the MLE may overfit.

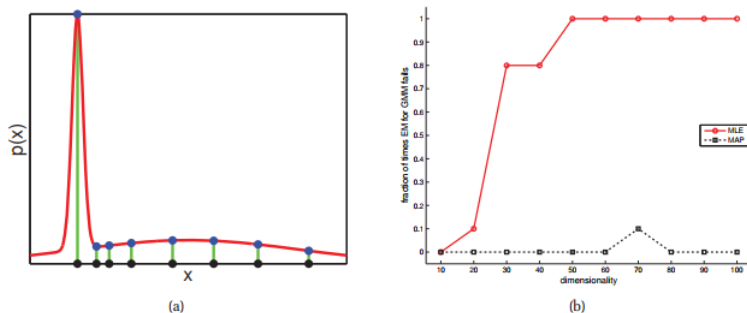


Figure 11.13 (a) Illustration of how singularities can arise in the likelihood function of GMMs. Based on (Bishop 2006a) Figure 9.7. Figure generated by `mixGaussSingularity`. (b) Illustration of the benefit of MAP estimation vs ML estimation when fitting a Gaussian mixture model. We plot the fraction of times (out of 5 random trials) each method encounters numerical problems vs the dimensionality of the problem, for $N = 100$ samples. Solid red (upper curve): MLE. Dotted black (lower curve): MAP. Figure generated by `mixGaussMLvsMAP`.

MAP estimation

An easy solution to this is to perform MAP estimation. The new auxiliary function:

$$Q(\boldsymbol{\theta}, \boldsymbol{\theta}^{\text{old}}) = \left[\sum_i \sum_k r_{ik} \log \pi_{ik} + \sum_i \sum_k r_{ik} \log p(\mathbf{x}_i | \boldsymbol{\theta}_k) \right] \\ + \log p(\boldsymbol{\pi}) + \sum_k \log p(\boldsymbol{\theta}_k)$$

The E step remains unchanged, it is natural to use a Dirichlet prior, $\boldsymbol{\pi} \sim \text{Dir}(\boldsymbol{\alpha})$, since this is conjugate to the categorical distribution. The MAP estimate is given by

$$\pi_k = \frac{r_k + \alpha_k - 1}{N + \sum_k \alpha_k - K}.$$

For simplicity, we consider a conjugate prior of the form

$$p(\boldsymbol{\mu}_k, k) = \text{NIW}(\boldsymbol{\mu}_k, k | \mathbf{m}_0, \kappa_0, \nu_0, \mathbf{S}_0).$$

MAP estimation

$$\begin{aligned}\hat{\boldsymbol{\mu}}_k &= \frac{r_k \bar{\mathbf{x}}_k + \kappa_0 \mathbf{m}_0}{r_k + \kappa_0} \\ \bar{\mathbf{x}}_k &\triangleq \frac{\sum_i r_{ik} \mathbf{x}_i}{r_k} \\ \hat{\mathbf{S}}_k &= \frac{\mathbf{S}_0 + \mathbf{S}_k + \frac{\kappa_0 r_k}{\kappa_0 + r_k} (\bar{\mathbf{x}}_k - \mathbf{m}_0)(\bar{\mathbf{x}}_k - \mathbf{m}_0)^T}{\nu_0 + r_k + D + 2} \\ \mathbf{S}_k &\triangleq \sum_i r_{ik} (\mathbf{x}_i - \bar{\mathbf{x}}_k)(\mathbf{x}_i - \bar{\mathbf{x}}_k)^T.\end{aligned}$$

how to set \mathbf{S}_0 ,

$$\mathbf{S}_0 = \frac{1}{K^{1/D}} \text{diag}(s_1^2, \dots, s_D^2)$$

where $s_j = (1/N) \sum_{i=1}^N (x_{ij} - \bar{x}_j)^2$, The parameter ν_0 controls how strongly we believe this prior. The weakest prior we can use, while still being proper, is to set $\nu_0 = D + 2$, so this is a common choice.