

EXPECTATION-MAXIMISATION ALGORITHM

Derivation and a cookbook for further recepies

Sven Laur
University of Tartu
email `swen@math.ut.ee`

May 8, 2019

1 Mixture model

To simplify the treatment, we consider here only how to derive and apply the expectation-maximisation algorithm in the context of soft clustering. The treatment of other applications is analogous. The soft-clustering is based on a probabilistic model, which combines k different data sources. More precisely, we assume that individual observations $\mathbf{x}_1, \dots, \mathbf{x}_n$ are generated by k different data sources $\mathcal{D}_1, \dots, \mathcal{D}_k$. Each data source \mathcal{D}_j is modelled by a probability distribution with parameters Θ_j and thus we can always compute the likelihood $p[\mathbf{x}_i|\Theta_j]$. Furthermore, we assume that all data points are independently sampled and the probability that \mathbf{x}_i is generated by the data source \mathcal{D}_j is λ_j . These parameters are also known as mixture proportions. Now let

$$\Theta = (\lambda_1, \dots, \lambda_k, \Theta_1, \dots, \Theta_k)$$

denote the set of all parameters. Then it is straightforward to estimate the likelihood of an individual observation

$$p[\mathbf{x}_i|\Theta] = \sum_{j=1}^k \lambda_j \cdot p[\mathbf{x}_i|\Theta_j] \ .$$

As a concrete example, consider a two-dimensional mixture model where both components are a two-dimensional normal distributions with parameters:

$$\begin{aligned} \Sigma_1 &= \begin{pmatrix} 4.125 & -3.875 \\ -3.875 & 4.125 \end{pmatrix} & \mu_1 &= \begin{pmatrix} 4 \\ -4 \end{pmatrix} \\ \Sigma_2 &= \begin{pmatrix} 4.125 & 3.875 \\ 3.875 & 4.125 \end{pmatrix} & \mu_2 &= \begin{pmatrix} -4 \\ 4 \end{pmatrix} \end{aligned}$$

and equal mixture proportions:

$$\lambda_1 = 0.5 \qquad \lambda_2 = 0.5 \ .$$

The corresponding model is depicted in Figure 1.

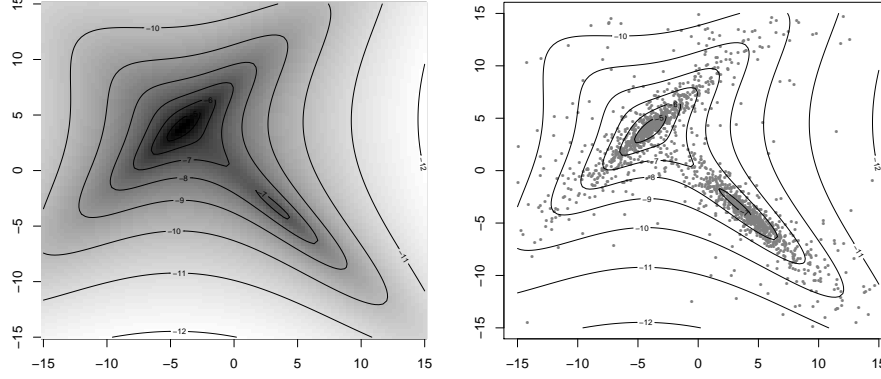


Figure 1: Two-dimesnional Gaussian mixture model. Probability density function is depicted on the left and 2000 element sample is depicted on the right.

2 Approximation of log-likelihood

Let $\mathbf{z} = (z_1, \dots, z_n)$ denote a potential labelling of points where $z_i \in \{1, \dots, k\}$ shows from which data source \mathcal{D}_j the datapoint \mathbf{x}_i is generated. For brevity, let

$$\mathcal{Z} = \{1, \dots, k\}^n$$

denote the set of all possible labelings. Then the posterior distribution of model parameters

$$p[\Theta | \mathbf{x}_1, \dots, \mathbf{x}_n] = \sum_{\mathbf{z} \in \mathcal{Z}} p[\mathbf{z}, \Theta | \mathbf{x}_1, \dots, \mathbf{x}_n] \quad (1)$$

can be further decomposed by applying product rule to individual terms

$$p[\mathbf{z}, \Theta | \mathbf{x}_1, \dots, \mathbf{x}_n] = \Pr[\mathbf{z} | \Theta, \mathbf{x}_1, \dots, \mathbf{x}_n] \cdot p[\Theta | \mathbf{x}_1, \dots, \mathbf{x}_n]$$

where the first term assigns probabilities to individual labelings and the second term is the term we want to express¹. Let us assume that we have somehow gotten a rather good estimate $q : \mathcal{Z} \rightarrow \mathbb{R}$ on the labelling probabilities, i.e.,

$$\forall \mathbf{z} \in \mathcal{Z} : \quad q(\mathbf{z}) \approx \Pr[\mathbf{z} | \Theta, \mathbf{x}_1, \dots, \mathbf{x}_n]$$

and we want to use and improve this estimation further. Without loss of generality, we can assume that the approximation is a proper probability distribution:

$$\sum_{\mathbf{z} \in \mathcal{Z}} q(\mathbf{z}) = 1 \quad . \quad (2)$$

¹Although it is formally correct we have gotten nowhere, yet

If $q(\mathbf{z})$ is not a probability distribution, we can always renormalise it. The latter does not change the values much, since

$$\sum_{\mathbf{z} \in \mathcal{Z}} \Pr[\mathbf{z} | \Theta, \mathbf{x}_1, \dots, \mathbf{x}_n] = 1 \quad .$$

Under this assumptions, we can approximate the sought posterior probability using individual labellings \mathbf{z} :

$$p[\Theta | \mathbf{x}_1, \dots, \mathbf{x}_n] = \frac{p[\mathbf{z}, \Theta | \mathbf{x}_1, \dots, \mathbf{x}_n]}{p[\mathbf{z} | \Theta, \mathbf{x}_1, \dots, \mathbf{x}_n]} \approx \frac{p[\mathbf{z}, \Theta | \mathbf{x}_1, \dots, \mathbf{x}_n]}{q(\mathbf{z})} \quad (3)$$

By substituting both approximations into the equation (1), we get another circular statement which is obvious:

$$p[\Theta | \mathbf{x}_1, \dots, \mathbf{x}_n] = \sum_{\mathbf{z} \in \mathcal{Z}} q(\mathbf{z}) \cdot \frac{p[\mathbf{z}, \Theta | \mathbf{x}_1, \dots, \mathbf{x}_n]}{q(\mathbf{z})} \quad .$$

To maximise this, we can take logarithms from both sides and maximise the right hand side of the following statement

$$\log p[\Theta | \mathbf{x}_1, \dots, \mathbf{x}_n] = \log \left(\sum_{\mathbf{z} \in \mathcal{Z}} q(\mathbf{z}) \cdot \frac{p[\mathbf{z}, \Theta | \mathbf{x}_1, \dots, \mathbf{x}_n]}{q(\mathbf{z})} \right) \quad .$$

Since logarithm is a convex-cap function, the Jensen inequality allows us to push the logarithm through the sum to get a lower bound

$$F(q, \Theta) := \sum_{\mathbf{z} \in \mathcal{Z}} q(\mathbf{z}) \cdot \log \left(\frac{p[\mathbf{z}, \Theta | \mathbf{x}_1, \dots, \mathbf{x}_n]}{q(\mathbf{z})} \right) \leq \log p[\Theta | \mathbf{x}_1, \dots, \mathbf{x}_n] \quad . \quad (4)$$

As $q(\mathbf{z})$ is a good approximation of $p[\mathbf{z} | \Theta, \mathbf{x}_1, \dots, \mathbf{x}_n]$, we can use the approximation (3) to show that the lower bound $F(q, \Theta)$ is close to to the actual value. Indeed, note that

$$\sum_{\mathbf{z} \in \mathcal{Z}} q(\mathbf{z}) \cdot \log \left(\frac{p[\mathbf{z}, \Theta | \mathbf{x}_1, \dots, \mathbf{x}_n]}{q(\mathbf{z})} \right) \approx \sum_{\mathbf{z} \in \mathcal{Z}} q(\mathbf{z}) \cdot \log(p[\Theta | \mathbf{x}_1, \dots, \mathbf{x}_n])$$

which together with the equation (2) itself implies

$$\sum_{\mathbf{z} \in \mathcal{Z}} q(\mathbf{z}) \cdot \log \left(\frac{p[\mathbf{z}, \Theta | \mathbf{x}_1, \dots, \mathbf{x}_n]}{q(\mathbf{z})} \right) \approx \log(p[\Theta | \mathbf{x}_1, \dots, \mathbf{x}_n])$$

To summarise, we have shown the following important claim.

Claim 1. *Assume that $q(\mathbf{z})$ is a good approximation of a posterior probability $\Pr[\mathbf{z} | \Theta, \mathbf{x}_1, \dots, \mathbf{x}_n]$. Then the lower bound $F(q, \Theta) \approx \log(p[\Theta | \mathbf{x}_1, \dots, \mathbf{x}_n])$.*

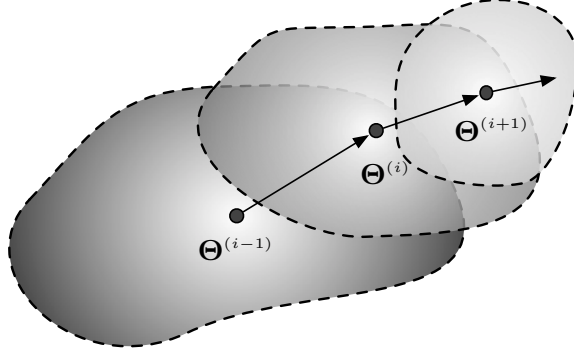


Figure 2: Iterative approximate search for maximal a posteriori estimate. Given a current solution $\Theta^{(i)}$, approximate the log-likelihood and find the next optimum $\Theta^{(i+1)}$. The algorithm is guaranteed to work if the new solution $\Theta^{(i+1)}$ is inside the grey region where the approximation is good. In each step, we have to recalibrate the approximation by setting $q(z) = \Pr[z|\Theta^{(i)}, \mathbf{x}_1, \dots, \mathbf{x}_n]$.

3 Iterative minimisation algorithm

In most cases, exact minimisation of the log-likelihood $p[\Theta|\mathbf{x}_1, \dots, \mathbf{x}_n]$ is intractable and even the exact computation of partial derivatives for the iterative minimisation is a challenging task. Hence, the iterative minimisation of $F(q, \Theta)$ together with the recalibration of the approximation $q(z)$ on the posterior probability $\Pr[z|\Theta, \mathbf{x}_1, \dots, \mathbf{x}_n]$ is a compelling alternative, see Figure 2.

Let us study this iterative maximisation procedure in more detail. In particular, let $q_i(z) = \Pr[z|\Theta^{(i)}, \mathbf{x}_1, \dots, \mathbf{x}_n]$ and $q_{i+1}(z) = \Pr[z|\Theta^{(i+1)}, \mathbf{x}_1, \dots, \mathbf{x}_n]$ be the approximations of the labelling distributions used in steps i and $i+1$. Then the next value of parameters is computed by maximising the lower bound:

$$\Theta^{(i+1)} = \underset{\Theta}{\operatorname{argmax}} F(q_i, \Theta) . \quad (5)$$

Let us assume further that $\Theta^{(i+1)}$ is so close to $\Theta^{(i)}$ that

$$q_i(z) \approx q_{i+1}(z) .$$

Then Claim 1 assures that the lower bound is rather close to the actual log-likelihood:

$$F(q_i, \Theta^{(i+1)}) \approx \log(p[\Theta^{(i+1)}|\mathbf{x}_1, \dots, \mathbf{x}_n]) .$$

Moreover, the approximation becomes more and more accurate if the difference between individual steps converges to zero:

$$\|\Theta^{(i)} - \Theta^{(i+1)}\| \rightarrow 0 .$$

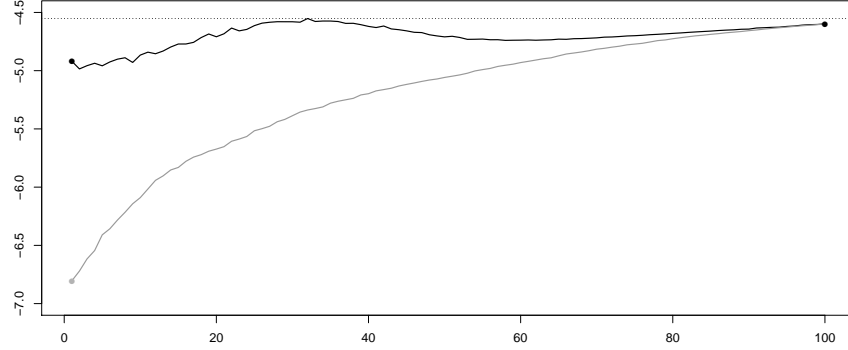


Figure 3: Convergence of the EM-algorithm. The x -axis depicts iterations. The y -axis depicts the values of the log-likelihood and its lower bound. The grey line shows how the lower bound increases during iterative minimisation. The black line shows the maximal value of the log-likelihood in a fixed size surrounding of $\Theta^{(i)}$. The iterative maximisation quickly converges to a local maxima. However, it also misses a another maximum, since this maximum is rather far away from the search path and thus the approximation is not good enough to detect it.

As a direct consequence note that if we manage to show that the iterative maximisation procedure converges

$$\Theta^{(i)} \rightarrow \Theta^*$$

then also the lower bound converges to the log-likelihood:

$$F(q_i, \Theta^{(i)}) \rightarrow \log(p[\Theta^* | \mathbf{x}_1, \dots, \mathbf{x}_n])$$

Moreover, as Θ^* is the local maximum of the lower bound $F(q_*, \Theta)$ then it must be also a local maximum of $\log(p[\Theta | \mathbf{x}_1, \dots, \mathbf{x}_n])$. Informal reasoning is the following. In a small surrounding of Θ^* , the approximation

$$F(q_*, \Theta) \approx \log(p[\Theta | \mathbf{x}_1, \dots, \mathbf{x}_n])$$

is so good that the maximum condition

$$F(q_*, \Theta) \leq F(q_*, \Theta^*)$$

also implies the same for log-likelihoods

$$\log(p[\Theta | \mathbf{x}_1, \dots, \mathbf{x}_n]) \leq \log(p[\Theta^* | \mathbf{x}_1, \dots, \mathbf{x}_n]) .$$

These basic concepts behind the algorithm are illustrated in Figure 3.

1. Choose a suitable initial estimate for model parameters $\Theta^{(0)}$.
2. E-STEP. Given an estimate $\Theta^{(t)}$ update labelling probability:

$$\forall \mathbf{z} \in \mathcal{Z} : \quad q_t(\mathbf{z}) = \Pr[\mathbf{z} | \Theta^{(t)}, \mathbf{x}_1, \dots, \mathbf{x}_n] \quad .$$

3. M-STEP. Find the next estimate $\Theta^{(t+1)}$ by minimising the lower bound:

$$\Theta^{(t+1)} = \underset{\Theta}{\operatorname{argmax}} F(q_t, \Theta) \quad .$$

4. Repeat from the step 2 until convergence.

Figure 4: Expectation-maximisation algorithm

4 Expectation-Maximisation algorithm

Our treatment of the iterative minimisation algorithm described above was rather informal. In particular, we did not prove that the sequence of iterative updates leads to the sequence of estimates

$$\Theta^{(0)}, \Theta^{(1)}, \dots, \Theta^{(i)}, \Theta^{(i+1)} \dots$$

that converges to a limiting value Θ^* . It is straightforward to see that

$$F(q_t, \Theta^{(i)}) \leq F(q_t, \Theta^{(i+1)}) \quad .$$

since the update step (5) maximises $F(q_t, \Theta)$. However, the recalibration of the labelling probabilities:

$$q_{t+1}(\mathbf{z}) = \Pr[\mathbf{z} | \Theta^{(t+1)}, \mathbf{x}_1, \dots, \mathbf{x}_n]$$

is completely detached from $F(q, \Theta)$ and thus we must do further analysis to show

$$F(q_t, \Theta^{(t+1)}) \leq F(q_{t+1}, \Theta^{(t+1)}) \quad ,$$

which would assure the convergence since the log-likelihood of the model must be bounded and thus also its lower bound $F(q, \Theta)$. The corresponding proof can be obtained through straightforward manipulations

$$\begin{aligned} F(q_{t+1}, \Theta^{(t+1)}) &= \sum_{\mathbf{z} \in \mathcal{Z}} q_{t+1}(\mathbf{z}) \cdot \log \left(\frac{p[\mathbf{z}, \Theta^{(t+1)} | \mathbf{x}_1, \dots, \mathbf{x}_n]}{\Pr[\mathbf{z} | \Theta^{(t+1)}, \mathbf{x}_1, \dots, \mathbf{x}_n]} \right) \\ &= \sum_{\mathbf{z} \in \mathcal{Z}} q(\mathbf{z}) \cdot \log(p[\Theta^{(t+1)} | \mathbf{x}_1, \dots, \mathbf{x}_n]) \\ &= \log(p[\Theta^{(t+1)} | \mathbf{x}_1, \dots, \mathbf{x}_n]) \end{aligned}$$

where the second equation from the Bayes equation

$$p[\mathbf{z}, \Theta^{(t+1)} | \mathbf{x}_1, \dots, \mathbf{x}_n] = \Pr [\mathbf{z} | \Theta^{(t+1)}, \mathbf{x}_1, \dots, \mathbf{x}_n] \cdot p[\Theta^{(t+1)} | \mathbf{x}_1, \dots, \mathbf{x}_n] .$$

Hence we have shown that after the recalibration step the lower bound matches the value it bounds:

$$F(q_{i+1}, \Theta^{(t+1)}) = \log(\Pr [\Theta^{(t+1)} | \mathbf{x}_1, \dots, \mathbf{x}_n])$$

and thus also

$$F(q_t, \Theta^{(t+1)}) \leq F(q_{t+1}, \Theta^{(t+1)}) .$$

Hence, we have shown the following claim.

Theorem 1. *The expectation-maximisation algorithm is guaranteed to converge if the log-likelihood function $\log(p[\Theta | \mathbf{x}_1, \dots, \mathbf{x}_n])$ is bounded.*

Proof sketch. Since both steps of the algorithm increase the lower bound and the log-likelihood is bounded the sequence of values

$$F(q_1, \Theta^{(1)}), \dots, F(q_i, \Theta^{(i)}), \dots$$

must converge to a limiting value. Hence the algorithm will stop. It is much more difficult to formally show that the convergence also takes place in terms of parameters $\Theta^{(0)}, \Theta^{(1)}, \dots, \Theta^{(i)}, \dots$, as large change in arguments might lead to small change in the output. To guarantee convergence in parameters, we must assume regularity from the function, i.e., there are no plateaus of constant values, where the algorithm could get stuck. \square

5 Semantics behind the maximisation step

The general description of the M-step is rather void of any meaning. In the following, we show how one can give a fractional semantics to this step. The main difference between the EM-algorithm and the hard clustering algorithm is in the E-step. The hard clustering algorithm chooses the most probable label z_i for the datapoint \mathbf{x}_i and then proceeds to find optimal parameters for each cluster. The soft-clustering algorithm assigns weights

$$w_{ij} = \Pr [z_i = j | \mathbf{x}_i, \Theta^{(t)}]$$

to each data point and uses all points with appropriate weights to re-estimate cluster parameters Θ . Next, we will show that the steps in the soft-clustering algorithm correspond to the E-step and the M-step.

First of all note that the recalibration rule in the E-step can be decomposed further as all data points are assumed to be independently drawn and labelled:

$$q_t(\mathbf{z}) = \frac{p[\mathbf{z}, \mathbf{x}_1, \dots, \mathbf{x}_n | \Theta^{(t)}]}{p[\mathbf{x}_1, \dots, \mathbf{x}_n | \Theta^{(t)}]} = \prod_{i=1}^n \frac{p[z_i, \mathbf{x}_i | \Theta^{(t)}]}{p[\mathbf{x}_i | \Theta^{(t)}]} = \prod_{i=1}^n \Pr [z_i | \mathbf{x}_i, \Theta^{(t)}] ,$$

which itself can be restated in terms of weights

$$q_t(\mathbf{z}) = \prod_{i=1}^n w_{iz_i} \ .$$

Secondly, we will decompose also the lower bound minimised in the M-step

$$\begin{aligned} F(q_t, \Theta) &= \sum_{\mathbf{z} \in \mathcal{Z}} q_t(\mathbf{z}) \cdot \log \left(\frac{p[\mathbf{z}, \Theta | \mathbf{x}_1, \dots, \mathbf{x}_n]}{q_t(\mathbf{z})} \right) \\ &= \sum_{\mathbf{z} \in \mathcal{Z}} q_t(\mathbf{z}) \cdot \log (p[\mathbf{z}, \Theta | \mathbf{x}_1, \dots, \mathbf{x}_n]) - \sum_{\mathbf{z} \in \mathcal{Z}} q_t(\mathbf{z}) \cdot \log(q_t(\mathbf{z})) \ . \end{aligned}$$

As the right term is independent of Θ , we can maximise only the first term

$$F_1(\Theta) = \sum_{\mathbf{z} \in \mathcal{Z}} q_t(\mathbf{z}) \cdot \log (p[\Theta, \mathbf{z} | \mathbf{x}_1, \dots, \mathbf{x}_n]) \ .$$

By the Bayes rule

$$p[\mathbf{z}, \Theta | \mathbf{x}_1, \dots, \mathbf{x}_n] = p[\mathbf{z}, \mathbf{x}_1, \dots, \mathbf{x}_n | \Theta] \cdot \frac{p[\Theta]}{p[\mathbf{x}_1, \dots, \mathbf{x}_n]}$$

where by our assumption the prior on the model parameters is constant. Thus, the second sum in the following expression

$$F_1(\Theta) = \sum_{\mathbf{z} \in \mathcal{Z}} q_t(\mathbf{z}) \cdot \log (p[\mathbf{z}, \mathbf{x}_1, \dots, \mathbf{x}_n | \Theta]) + \sum_{\mathbf{z} \in \mathcal{Z}} q_t(\mathbf{z}) \cdot \log \left(\frac{p[\Theta]}{p[\mathbf{x}_1, \dots, \mathbf{x}_n]} \right)$$

is constant during the maximisation and we ignore it. Hence, we can maximise the following function

$$F_2(\Theta) = \sum_{\mathbf{z} \in \mathcal{Z}} q_t(\mathbf{z}) \cdot \log (p[\mathbf{z}, \mathbf{x}_1, \dots, \mathbf{x}_n | \Theta])$$

instead. Now we can apply the independence of \mathbf{x}_i, z_i pairs and further simplify

$$\begin{aligned} F_2(\Theta) &= \sum_{\mathbf{z} \in \mathcal{Z}} q_t(\mathbf{z}) \cdot \log \left(\prod_{i=1}^n p[\mathbf{x}_i, z_i | \Theta] \right) \\ &= \sum_{i=1}^n \sum_{\mathbf{z} \in \mathcal{Z}} q_t(\mathbf{z}) \cdot \log (p[\mathbf{x}_i, z_i | \Theta]) \\ &= \sum_{i=1}^n \sum_{\mathbf{z} \in \mathcal{Z}} q_t(\mathbf{z}) \cdot \log (p[\mathbf{x}_i, z_i | \Theta]) \\ &= \sum_{i=1}^n \sum_{\mathbf{z} \in \mathcal{Z}} \prod_{\ell=1}^n w_{\ell z_\ell} \cdot \log (p[\mathbf{x}_i, z_i | \Theta]) \ . \end{aligned}$$

For any fixed i value, the inner sum

$$\sum_{\mathbf{z} \in \mathcal{Z}} \prod_{\ell=1}^n w_{\ell z_\ell} \cdot \log(p[\mathbf{x}_i, z_i | \Theta])$$

can be viewed as a standard sum over products

$$\sum_{\mathbf{z} \in \mathcal{Z}} \prod_{\ell=1}^n a_{\ell z_\ell} = \prod_{\ell=1}^n \sum_{j=1}^k a_{\ell j}$$

where

$$a_{\ell z_\ell} = \begin{cases} w_{\ell z_\ell} & \text{if } i \neq \ell, \\ w_{i z_i} \cdot \log(p[\mathbf{x}_i, z_i | \Theta]) & \text{if } i = \ell, \end{cases}$$

we get that

$$\sum_{\mathbf{z} \in \mathcal{Z}} \prod_{\ell=1}^n w_{\ell z_\ell} \cdot \log(p[\mathbf{x}_i, z_i | \Theta]) = \prod_{\substack{\ell=1 \\ \ell \neq i}}^n \left(\sum_{j=1}^k w_{\ell j} \right) \cdot \left(\sum_{j=1}^k w_{i j} \cdot \log(p[\mathbf{x}_i, z_i = j | \Theta]) \right).$$

Since all weights $w_{\ell j}$ connected to the data point \mathbf{x}_ℓ sum up to one we get

$$\sum_{\mathbf{z} \in \mathcal{Z}} \prod_{\ell=1}^n w_{\ell z_\ell} \cdot \log(p[\mathbf{x}_i, z_i = j | \Theta]) = \sum_{j=1}^k w_{i j} \cdot \log(p[\mathbf{x}_i, z_i = j | \Theta])$$

and thus the triple sum simplifies further

$$F_2(\Theta) = \sum_{i=1}^n \sum_{\mathbf{z} \in \mathcal{Z}} \prod_{\ell=1}^n w_{\ell z_\ell} \cdot \log(p[\mathbf{x}_i, z_i | \Theta]) = \sum_{i=1}^n \sum_{j=1}^k w_{i j} \cdot \log(p[\mathbf{x}_i, z_i = j | \Theta]).$$

Hence, we have obtained the following claim.

Claim 2. *The maximisation task in the M-step is equivalent to the simplified maximisation task*

$$F_2(\Theta) = \sum_{i=1}^n \sum_{j=1}^k w_{i j} \cdot \log(p[\mathbf{x}_i, z_i = j | \Theta]) \rightarrow \max.$$

General recipe: The simplified minimisation task provides you a general way to derive the update formulae for the M-step.

In the clustering context, where the likelihood of \mathbf{x}_i is determined completely by the cluster parameters Θ_j , which generate the point, the task can be further

simplified:

$$\begin{aligned} F_2(\Theta) &= \sum_{i=1}^n \sum_{j=1}^k w_{ij} \cdot \log (p[\mathbf{x}_i | z_i = j, \Theta] \cdot \Pr [z_i = j | \Theta]) \\ &= \sum_{i=1}^n \sum_{j=1}^k w_{ij} \cdot \log (p[\mathbf{x}_i | \Theta_j]) + \sum_{i=1}^n \sum_{j=1}^k w_{ij} \cdot \log (\Pr [z_i = j | \lambda_j]) \quad . \end{aligned}$$

Recipe for mixture proportions: As both terms can be independently maximised, we can directly find optimal update for the mixture proportions:

$$\lambda_j = \frac{\sum_{i=1}^n w_{ij}}{\sum_{i=1}^n \sum_{j=1}^k w_{ij}} \quad .$$

If mixture proportions are fixed by the model, then the update step is not possible and is ignored.

Recipe for cluster parameters: To fix the cluster parameters, we must minimise the following terms

$$G_j(\Theta_j) = \sum_{i=1}^n w_{ij} \cdot \log (p[\mathbf{x}_i | \Theta_j]) \quad .$$

Connection with factional counts: In many cases, direct derivation of optimal parameters is not analytical and there is only iterative minimisation algorithm for hard clustering problem:

$$H(\Theta_j) = \sum_{i \in \mathcal{I}_j} \log p[\mathbf{x}_i | \Theta_j] \rightarrow \min$$

where \mathcal{I}_j is the set of points assigned to the cluster. Note that we can rewrite this as a double sum using indicator functions

$$H(\Theta_j) = \sum_{i=1}^n [i \in \mathcal{I}_j] \cdot \log p[\mathbf{x}_i | \Theta_j] \rightarrow \min \quad .$$

Since some points can coincide, i.e. $\mathbf{x}_i = \mathbf{x}_\ell$, then the procedure that can find a solution for $H(\Theta)$ can find a solution for

$$H_*(\Theta_j) = \sum_{i=1}^n c_{ij} \cdot \log p[\mathbf{x}_i | \Theta_j] \rightarrow \min \quad .$$

where c_{ij} are integer counts. We just have to duplicate \mathbf{x}_i c_{ij} times to generate modified data set for which the original procedure knows how to find the maximising assignment. Obviously, we can approximate $G_j(\Theta_j)$ by the following

sum with integer coefficients

$$G_j(\Theta_j) \approx \frac{1}{\ell} \cdot \sum_{i=1}^n \sum_{j=1}^k \text{round}(w_{ij} \cdot \ell) \cdot p[\mathbf{x}_i | \Theta_j]$$

where the size of ℓ determines the approximation precision. As a result, we can use the update step from the hard-clustering algorithm to maximise

$$H_*(\Theta_j) = \sum_{i=1}^n \sum_{j=1}^k \text{round}(w_{ij} \cdot \ell) \cdot p[\mathbf{x}_i | \Theta_j] \ .$$

6 Examples

6.1 Gaussian mixture model

Recall that the closed form solution for updating the parameters of the Gaussian cluster are

$$\begin{aligned} \boldsymbol{\mu}_j &= \frac{1}{|\mathcal{I}_j|} \cdot \sum_{i \in \mathcal{I}_j} \mathbf{x}_i \\ \Sigma_j &= \frac{1}{|\mathcal{I}_j|} \cdot X_c^T X_c \end{aligned}$$

where X_c is matrix consisting of centred row vectors $\mathbf{x}_i - \boldsymbol{\mu}_j$ belonging to the cluster. Now it is straightforward to verify that with integral multiplicity c_{ij} the formula for the cluster centres becomes

$$\begin{aligned} n_j &= \sum_{i=1}^n c_{ij} \\ \boldsymbol{\mu}_j &= \frac{1}{n_j} \cdot \sum_{i=1}^n c_{ij} \mathbf{x}_i \end{aligned}$$

and thus the optimal solution for soft-clustering task is

$$\begin{aligned} n_j &= \sum_{i=1}^n w_{ij} \\ \boldsymbol{\mu}_j &= \frac{1}{n_j} \cdot \sum_{i=1}^n w_{ij} \mathbf{x}_i \end{aligned}$$

The analysis of the shape parameter is a bit more involved, because we must understand how replicated rows in the centred matrix X_c affect Σ_j . Let us do this by the example. Assume that the original centred matrix is 2×5 matrix

$$\begin{pmatrix} a_1 & a_2 & a_3 & a_4 & a_5 \\ b_1 & b_2 & b_3 & b_4 & b_5 \end{pmatrix}$$

Then the corresponding covariance matrix is computed as

$$\Sigma = \frac{1}{2} \cdot \begin{pmatrix} a_1 & b_1 \\ a_2 & b_2 \\ a_3 & b_3 \\ a_4 & b_4 \\ a_5 & b_5 \end{pmatrix} \begin{pmatrix} a_1 & a_2 & a_3 & a_4 & a_5 \\ b_1 & b_2 & b_3 & b_4 & b_5 \end{pmatrix}$$

If the first row (data point) has has multiplicity 2 and the second multiplicity 3 then the corresponding covariance matrix is computed

$$\Sigma = \frac{1}{5} \cdot \begin{pmatrix} a_1 & a_1 & b_1 & b_1 & b_1 \\ a_2 & a_2 & b_2 & b_2 & b_2 \\ a_3 & a_3 & b_3 & b_3 & b_3 \\ a_4 & a_4 & b_4 & b_4 & b_4 \\ a_5 & a_5 & b_5 & b_5 & b_5 \end{pmatrix} \begin{pmatrix} a_1 & a_2 & a_3 & a_4 & a_5 \\ a_1 & a_2 & a_3 & a_4 & a_5 \\ b_1 & b_2 & b_3 & b_4 & b_5 \\ b_1 & b_2 & b_3 & b_4 & b_5 \\ b_1 & b_2 & b_3 & b_4 & b_5 \end{pmatrix}$$

and thus the element

$$\Sigma_{ij} = \frac{1}{5} \cdot (a_i a_j + a_i a_j + b_i b_j + b_i b_j + b_i b_j) = \frac{1}{5} \cdot (2a_i a_j + 3b_i b_j)$$

The same result can be obtained if we multiply

$$\Sigma = \frac{1}{5} \cdot \begin{pmatrix} a_1 & b_1 \\ a_2 & b_2 \\ a_3 & b_3 \\ a_4 & b_4 \\ a_5 & b_5 \end{pmatrix} \begin{pmatrix} 2 & 0 \\ 0 & 3 \end{pmatrix} \begin{pmatrix} a_1 & a_2 & a_3 & a_4 & a_5 \\ b_1 & b_2 & b_3 & b_4 & b_5 \end{pmatrix}$$

Hence we have empirically derived the following update rule

$$\Sigma_j = \frac{1}{n_j} \cdot X_c^T \text{diag}(\mathbf{w}_{*,j}) X_c$$

where $\text{diag}(\mathbf{w}_{*,j})$ is a diagonal matrix with elements w_{1j}, \dots, w_{nj} .

6.2 Laplacian distribution

For this cluster model we do not have closed form solution for the shape parameters but we do have a closed form solution for the cluster centres. Namely,

$$\mu_{j\ell} = \text{median}(x_{1k}, \dots, x_{nk})$$

Hence, we can quickly derive what is the corresponding update step if the data points have multiple copies. In this case, we have the following scheme

$$\underbrace{a_1, \dots, a_1}_{c_{1j}}, \underbrace{a_2, \dots, a_2}_{c_{2j}}, \dots, \underbrace{a_n, \dots, a_n}_{c_{nj}}$$

where a_1, \dots, a_n are ordered increasingly and we have to choose the middle element. This lead to the following curious update step

1. Order the elements increasingly a_1, \dots, a_n
2. Let w_1, \dots, w_n be the corresponding weights. Compute cumulative sums $d_i = w_1 + \dots + w_i$ and choose the element a_k for which $d_{k-1} < 0.5$ and $d_k > 0.5$. If there is a tie $d_k = 0.5$ then return $\frac{1}{2} \cdot (a_{k-1} + a_k)$.