

The Basic Idea of EM

Jianxin Wu

LAMDA Group

National Key Lab for Novel Software Technology
Nanjing University, China
wujx2001@gmail.com

June 7, 2017

Contents

1 Introduction	1
2 GMM: A working example	2
2.1 Gaussian mixture model	2
2.2 The hidden variable interpretation	3
2.3 What if we can observe the hidden variable?	5
2.4 Can we imitate an oracle?	6
3 An informal description of the EM algorithm	6
4 The Expectation-Maximization algorithm	7
4.1 Jointly-non-concave incomplete log-likelihood	7
4.2 (Possibly) Concave complete data log-likelihood	8
4.3 The general EM derivation	9
4.4 The E- & M-steps	11
4.5 The EM algorithm	12
4.6 Will EM converge?	12
5 EM for GMM	13

1 Introduction

Statistical learning models are very important in many areas inside computer science, including but not confined to machine learning, computer vision, pattern recognition and data mining. It is also important in some deep learning models, such as the Restricted Boltzmann machine (RBM).

Statistical learning models have parameters, and estimating such parameters from data is one of the key problems in the study of such models. Expectation-Maximization (EM) is arguably the most widely used parameter estimation technique. Hence, it is worthwhile to know some basics of EM.

However, although EM is a must-have knowledge in studying statistical learning models, it is not easy for beginners. This note introduces the basic idea behind EM.

I want to emphasize that the main purpose of this note is to introduce the *basic idea* (or, emphasizing the *intuition*) behind EM, not for covering all details of EM or presenting rigorous mathematical derivations.¹

2 GMM: A working example

Let us start from a simple working example: the Gaussian Mixture Model (GMM).

2.1 Gaussian mixture model

In Figure 1, we show three curves corresponding to three different probability density functions (p.d.f.). The blue curve is the p.d.f. of a normal distribution $N(10, 16)$, i.e., a Gaussian distribution with the mean $\mu = 10$ and the standard deviation $\sigma = 4$ (and $\sigma^2 = 16$). We denote this p.d.f. as $p_1(x) = N(x; 10, 16)$. The red curve is another normal distribution $N(30, 49)$ with $\mu = 30$ and $\sigma = 7$. Similarly, we denote it as $p_2(x) = N(x; 30, 49)$.

We are interested in the black curve, whose first half is similar to the blue one, while the second half is similar to the red one. This curve is also the p.d.f. of a distribution, denoted by p_3 . Since the black curve is similar to parts of the blue and red curves, it is reasonable to conjecture that p_3 are related to both p_1 and p_2 .

Indeed, p_3 is a *weighted* combination of p_1 and p_2 . In this example,

$$p_3(x) = 0.2p_1(x) + 0.8p_2(x). \quad (1)$$

Because $0.2 + 0.8 = 1$, it is easy to verify that $p_3(z) \geq 0$ always holds and $\int_{-\infty}^{\infty} p_3(x) dx = 1$. Hence, p_3 is a valid p.d.f.

p_3 is a *mixture* of two Gaussians (p_1 and p_2), hence a *Gaussian mixture model* (GMM). The definition of a GMM is in fact more general: it can have more than two components, and the Gaussians can be multivariate.

¹The first version of this note was written in Chinese, and was started as a note-taking in a course in the Georgia Institute of Technology while I was a graduate student there. That version was typeset in Microsoft Word. Unfortunately, that version contained a lot of errors and I did not have a chance to check it again. This version (written in 2016) is started while I am preparing materials for the *Pattern Recognition* course I will teach in the Spring Semester in Nanjing University. It is greatly expanded, and the errors that I found are corrected.

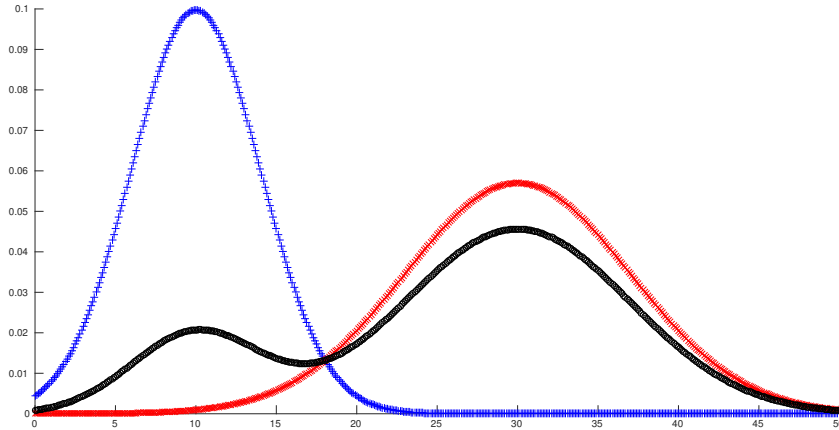


Figure 1: A simple GMM illustration.

A GMM is a distribution whose p.d.f. has the following form:

$$p(\mathbf{x}) = \sum_{i=1}^N \alpha_i N(\mathbf{x}; \boldsymbol{\mu}_i, \Sigma_i) \quad (2)$$

$$= \sum_{i=1}^N \frac{\alpha_i}{(2\pi)^{d/2} |\Sigma_i|^{1/2}} \exp\left(-\frac{1}{2}(\mathbf{x} - \boldsymbol{\mu}_i)^T \Sigma_i^{-1} (\mathbf{x} - \boldsymbol{\mu}_i)\right), \quad (3)$$

in which \mathbf{x} is a d -dimensional random vector.

In this GMM, there are N Gaussian components, with the i -th Gaussian has the mean vector $\boldsymbol{\mu}_i \in \mathbb{R}^d$ and the covariance matrix $\Sigma_i \in \mathbb{R}^{d \times d}$.² These Gaussian components are combined together using a linear combination, where the weight for the i -th component is α_i (called the *mixing coefficients*). The mixing coefficients must satisfy the following conditions

$$\sum_{i=1}^N \alpha_i = 1, \quad (4)$$

$$\alpha_i \geq 0, \forall i. \quad (5)$$

It is easy to verify that under these conditions, $p(\mathbf{x})$ is a valid multivariate probability density function.

2.2 The hidden variable interpretation

We can have a different interpretation of the Gaussian mixture model, using the hidden variable concept, as illustrated in Figure 2.

²We will use boldface letters to denote a vector.

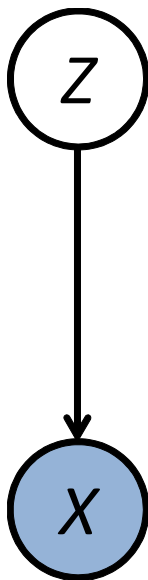


Figure 2: GMM as a graphical model.

In Figure 2, the random variable X follows a Gaussian mixture model (cf. Equation 3). Its parameter is

$$\boldsymbol{\theta} = \{\alpha_i, \boldsymbol{\mu}_i, \boldsymbol{\Sigma}_i\}_{i=1}^N . \quad (6)$$

If we want to sample an instance from this GMM, we could directly sample from the p.d.f. in Equation 3. However, there is another two-step way to perform the sampling.

Let us define a random variable Z . Z is a multinomial discrete distribution, taking values from the set $\{1, 2, \dots, N\}$. The probability that Z takes the value $Z = i$ is α_i , i.e., $\Pr(Z = i) = \alpha_i$, for $1 \leq i \leq N$. Then, the two-step sampling procedure is:

Step 1 Sample from Z , and get a value i ($1 \leq i \leq N$);

Step 2 Sample \boldsymbol{x} from the i -th Gaussian component $N(\boldsymbol{\mu}_i, \boldsymbol{\Sigma}_i)$.

It is easy to verify that the sample \boldsymbol{x} achieved from this two-step sampling procedure follows the underlying GMM distribution in Equation 3.

In learning GMM parameters, we are given a sample set $\{\boldsymbol{x}_1, \boldsymbol{x}_2, \dots, \boldsymbol{x}_M\}$, where \boldsymbol{x}_i are i.i.d. (independently and independently distributed) instances sampled from the p.d.f. in Equation 3. From this set of samples, we want to estimate or learn the GMM parameters $\boldsymbol{\theta} = \{\alpha_i, \boldsymbol{\mu}_i, \boldsymbol{\Sigma}_i\}_{i=1}^N$.

Because we are given the samples \boldsymbol{x}_i , the random variable X (cf. Figure 2) are called *observed* (or observable) random variables. As shown in Figure 2, observed random variables are usually shown as a filled circle.

The random variable Z , however, is not observable, and is called a *hidden variable* (or a latent variable). Hidden variables are shown as a circled, as the Z node in Figure 2.

2.3 What if we can observe the hidden variable?

In real applications, we do not know the value (or instantiation) of Z , because it is hidden (not observable). This fact makes estimating GMM parameters rather difficult, and techniques such as EM (the focus of this note) have to be employed.

However, for the sample set $\mathcal{X} = \{\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_M\}$, let us consider the scenario in which we can further suppose that some oracle has given us the value of Z : $\mathcal{Z} = \{z_1, z_2, \dots, z_M\}$. In other words, we know that \mathbf{x}_i is sampled from the z_i -th Gaussian component.

In this case, it is easy to estimate the parameters θ . First, we can find all those samples that are generated from the i -th component, and use \mathcal{X}_i to denote this subset of samples. In precise mathematical languages,

$$\mathcal{X}_i = \{\mathbf{x}_j | z_j = i, 1 \leq j \leq M\}. \quad (7)$$

The mixing coefficient estimation is a simple counting. We can count the number of examples which are generated from the i -th Gaussian component as $m_i = |\mathcal{X}_i|$, where $|\cdot|$ is the size (number of elements) of a set. Then, the maximum likelihood estimation for α_i is

$$\hat{\alpha}_i = \frac{m_i}{\sum_{j=1}^M m_j} = \frac{m_i}{M}. \quad (8)$$

Second, it is also easy to estimate the μ_i and Σ_i parameters for any $1 \leq i \leq N$. The maximum likelihood estimation solutions are the same as those single Gaussian equations:³

$$\hat{\mu}_i = \frac{1}{m_i} \sum_{\mathbf{x} \in \mathcal{X}_i} \mathbf{x}, \quad (9)$$

$$\hat{\Sigma}_i = \frac{1}{|m_i|} \sum_{\mathbf{x} \in \mathcal{X}_i} (\mathbf{x} - \hat{\mu}_i)(\mathbf{x} - \hat{\mu}_i)^T. \quad (10)$$

In short, if we know the hidden variable's instantiations, the estimation is straightforward. Unfortunately, we are only given the observed sample set \mathcal{X} . The hidden variable instantiations \mathcal{Z} is unknown to us. This fact complicates the entire parameter estimation process.

³Please refer to my note on properties of normal distributions for derivation of these equations.

2.4 Can we imitate an oracle?

A natural question to ask ourselves is: if we do not have an oracle to teach us, can we imitate the oracle's teaching? In other words, can we guess the value of z_j for \mathbf{x}_j ?

A natural choice is to use the posterior $p(z_j|\mathbf{x}_j, \boldsymbol{\theta}^{(t)})$ as a replacement for z_j . This term is the posterior probability given the sample \mathbf{x}_j and the current parameter value $\boldsymbol{\theta}^{(t)}$.⁴ The posterior probability is the best educated guess we can have given the information that is handy.

In this guessing game, we have at least two issues in our way. First, an oracle is supposed to know everything, and will be able to tell us that \mathbf{x}_7 comes from the third Gaussian component, with 100% confidence. If an oracle exists, we can simply say $z_7 = 3$ in this example. However, our guess will never be deterministic—it can at best be a probability distribution about the *random variable* z_j .

Hence, we will assume that for every observed sample \mathbf{x} , there is a corresponding hidden vector \mathbf{z} , whose values can be guessed but cannot be observed. We still use Z to denote the underlying random variable, and use \mathcal{Z} to denote the set of hidden vectors. In the GMM example, a vector \mathbf{z}_j will have N dimensions, but one and only one of these dimensions will be 1, and all others will be 0.

Second, the guess we have about \mathbf{z}_j is a distribution determined by the posterior $p(\mathbf{z}_j|\mathbf{x}_j, \boldsymbol{\theta}^{(t)})$. However, what we really want are values instead of a distribution. How are we going to use this guess? A common trick in statistical learning is to use its expectation. We will leave the details about how the expectation is used to later sections.

3 An informal description of the EM algorithm

Now we are ready to give an informal description of the EM algorithm.

- We first initialize the values of $\boldsymbol{\theta}$ in any reasonable way;
- Then, we can estimate the best possible \mathcal{Z} (expectation of its posterior distribution) using \mathcal{X} and the current $\boldsymbol{\theta}$ estimation;
- With this \mathcal{Z} estimation, we can find a better estimate of $\boldsymbol{\theta}$ using \mathcal{X} ;
- A better $\boldsymbol{\theta}$ (combined with \mathcal{X}) will lead to a better guess of \mathcal{Z} ;
- This process (estimating $\boldsymbol{\theta}$ and \mathcal{Z} in alternating order) can proceed until the change in $\boldsymbol{\theta}$ is small (i.e., the procedure converges).

In still more informal languages, after proper initialization of the parameters, we can:

⁴As we will see, EM is an iterative process, in which the variable t is the iteration index. We will update the parameter $\boldsymbol{\theta}$ is every iteration, and use $\boldsymbol{\theta}^{(t)}$ to denote its value in the t -th iteration.

E-Step Find a better guess of the non-observable hidden variables, by using the data and current parameter values;

M-Step Find a better parameter estimation, by using the current guess for the hidden variables and the data;

Repeat Repeat the above two steps until convergence.

In the EM algorithm, the first step is usually called the Expectation step, abbreviated as the E-step; while the second step is usually called the Maximization step, abbreviated as the M-step. The EM algorithm repeats E- and M-steps in alternating order. When the algorithm converges, we get the desired parameter estimations.

4 The Expectation-Maximization algorithm

Now we will show more details of the EM algorithm. Suppose we are dealing with two sets of random variables: the observed variables X and the hidden variables Z . The joint p.d.f. is $p(X, Z; \theta)$, where θ are the parameters. We are given a set of instances of X to learn the parameters, as $\mathcal{X} = \{\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_M\}$. The task is to estimate θ from \mathcal{X} .

For every \mathbf{x}_j , there is a corresponding \mathbf{z}_j . And we want to clarify that θ now include the parameters that are associated with Z . In the GMM example, z_{ij} are estimates for Z , $\{\alpha_i, \boldsymbol{\mu}_i, \Sigma_i\}_{i=1}^N$ are parameters specifying X , and θ include both sets of parameters.

4.1 Jointly-non-concave incomplete log-likelihood

If we use the maximum likelihood (ML) estimation technique, the ML estimate for θ is

$$\hat{\theta} = \arg \max_{\theta} p(\mathcal{X}|\theta). \quad (11)$$

Or equivalently, we can maximize the log-likelihood

$$\hat{\theta} = \arg \max_{\theta} \ln p(\mathcal{X}|\theta), \quad (12)$$

because $\ln(\cdot)$ is a monotonically increasing function.

Then, parameter estimation becomes an optimization problem. We will use the notation $L(\theta)$ to denote the log-likelihood, that is,

$$L(\theta) = \ln p(\mathcal{X}|\theta). \quad (13)$$

Recent developments in optimization tells that we can generally consider a minimization problem as “easy” if it is *convex*, but non-convex problems are usually difficult to solve. Equivalently, a *concave* maximization problem is generally considered easy, while non-concave maximization is usually difficult, because the negative of a convex function is a concave one, and vice versa.

Unfortunately, the log-likelihood is non-concave in most cases. Take the Gaussian mixture model as an example, the likelihood $p(\mathcal{X}|\boldsymbol{\theta})$ is

$$p(\mathcal{X}|\boldsymbol{\theta}) = \prod_{j=1}^M \left(\sum_{i=1}^N \frac{\alpha_i}{(2\pi)^{d/2} |\Sigma_i|^{1/2}} \exp \left(-\frac{1}{2} (\mathbf{x}_j - \boldsymbol{\mu}_i)^T \Sigma_i^{-1} (\mathbf{x}_j - \boldsymbol{\mu}_i) \right) \right). \quad (14)$$

The log-likelihood has the following form:

$$\sum_{j=1}^M \ln \left(\sum_{i=1}^N \frac{\alpha_i}{(2\pi)^{d/2} |\Sigma_i|^{1/2}} \exp \left(-\frac{1}{2} (\mathbf{x}_j - \boldsymbol{\mu}_i)^T \Sigma_i^{-1} (\mathbf{x}_j - \boldsymbol{\mu}_i) \right) \right), \quad (15)$$

This equation is non-concave with respect to the joint optimization variables $\{\alpha_i, \boldsymbol{\mu}_i, \Sigma_i\}_{i=1}^N$. In other words, this is a difficult maximization problem.

We have two sets of random variables X and Z . The log-likelihood in Equation 15 is called the *incomplete data* log-likelihood because Z is not in that equation.

4.2 (Possibly) Concave complete data log-likelihood

The complete data log-likelihood is

$$\ln p(\mathcal{X}, \mathcal{Z}|\boldsymbol{\theta}). \quad (16)$$

Let us use GMM as an example once more. In GMM, the \mathbf{z}_j vectors (which form \mathcal{Z}) is an N -dimensional vector with $N-1$ 0's and only one dimension with value 1. Hence, the complete data likelihood is

$$p(\mathcal{X}, \mathcal{Z}|\boldsymbol{\theta}) = \prod_{j=1}^M \prod_{i=1}^N \left[\frac{\alpha_i}{(2\pi)^{d/2} |\Sigma_i|^{1/2}} \exp \left(-\frac{1}{2} (\mathbf{x}_j - \boldsymbol{\mu}_i)^T \Sigma_i^{-1} (\mathbf{x}_j - \boldsymbol{\mu}_i) \right) \right]^{z_{ij}}. \quad (17)$$

This equation can be explained using the two-step sampling process. Let us assume \mathbf{x}_j is generated by the i' -th Gaussian component. Then, if $i \neq i'$, we know that $z_{ij} \neq 1$, otherwise $z_{ij} = z_{i'j} = 1$. In other words, the term inside $[\cdot]$ will equal 1 for $N-1$ times when $z_{ij} = 0$, and the remaining one entry will be evaluated to $\alpha_{i'} N(\mathbf{x}; \boldsymbol{\mu}_{i'}, \Sigma_{i'})$, which exactly matches the 2-step sampling procedure.⁵

Then, the complete data log-likelihood is

$$\sum_{j=1}^M \sum_{i=1}^N z_{ij} \left(\frac{1}{2} (\ln |\Sigma_i^{-1}| - (\mathbf{x}_j - \boldsymbol{\mu}_i)^T \Sigma_i^{-1} (\mathbf{x}_j - \boldsymbol{\mu}_i)) + \ln \alpha_i \right) + \text{const}. \quad (18)$$

Let us consider the scenario when the hidden variable z_{ij} is known, but $\alpha_i, \boldsymbol{\mu}_i$ and Σ_i are unknown. Here we suppose Σ_i is invertible for $1 \leq i \leq N$. Instead

⁵The first step has probability $\alpha_{i'}$, and the second step has density $N(\mathbf{x}; \boldsymbol{\mu}_{i'}, \Sigma_{i'})$. These two steps are independent of each other, hence the product rule applies.

of considering parameters $(\boldsymbol{\mu}_i, \Sigma_i)$, we consider $(\boldsymbol{\mu}_i, \Sigma_i^{-1})$.⁶ It is well known that the log-determinant function $\ln |\cdot|$ is concave. It is also easy to prove that the quadratic term $(\mathbf{z}_j - \boldsymbol{\mu}_i)^T \Sigma_i^{-1} (\mathbf{z}_j - \boldsymbol{\mu}_i)$ is jointly convex with respect to variables $(\boldsymbol{\mu}_i, \Sigma_i^{-1})$, which directly implies that its negative is concave.⁷ Hence, this sub-problem can be efficiently solved.

From this optimization perspective, we can understand the EM algorithm from a different point of view. Although the original maximum likelihood parameter estimation problem is difficult to solve (jointly non-concave), the EM algorithm can usually (but not always) make concave subproblems, hence becoming efficiently solvable.

4.3 The general EM derivation

Now we talk about EM in the general sense. We have observable variables X and samples \mathcal{X} . We also have hidden variables Z and unobservable samples \mathcal{Z} . The overall system parameters are denoted by $\boldsymbol{\theta}$.

The parameter learning problem tries to find optimal parameters $\hat{\boldsymbol{\theta}}$ by maximizing the incomplete data log-likelihood

$$\hat{\boldsymbol{\theta}} = \arg \max_{\boldsymbol{\theta}} \ln p(\mathcal{X}|\boldsymbol{\theta}). \quad (19)$$

We assume Z is discrete, and hence

$$p(\mathcal{X}|\boldsymbol{\theta}) = \sum_{\mathcal{Z}} p(\mathcal{X}, \mathcal{Z}|\boldsymbol{\theta}). \quad (20)$$

However, this assumption is mainly for notational simplicity. If Z is continuous, we can replace the summation with an integral.

Although we have mentioned previously that we can use the posterior of Z , i.e., $p(\mathcal{Z}|\mathcal{X}, \boldsymbol{\theta})$ as our guess, it is also interesting to observe what will happen to the complete data likelihood if we use a *arbitrary* distribution for Z (and hence understand why the posterior is special and why we should use it).

Let q be any valid probability distribution for Z , we can measure how different is q to the posterior using the classic Kullback-Leibler (KL) divergence measure, as

$$\text{KL}(q||p) = - \sum_{\mathcal{Z}} q(\mathcal{Z}) \ln \left(\frac{p(\mathcal{Z}|\mathcal{X}, \boldsymbol{\theta})}{q(\mathcal{Z})} \right). \quad (21)$$

The probability theory tells us that

$$p(\mathcal{X}|\boldsymbol{\theta}) = \frac{p(\mathcal{X}, \mathcal{Z}|\boldsymbol{\theta})}{p(\mathcal{Z}|\mathcal{X}, \boldsymbol{\theta})} \quad (22)$$

$$= \frac{p(\mathcal{X}, \mathcal{Z}|\boldsymbol{\theta})}{q(\mathcal{Z})} \frac{q(\mathcal{Z})}{p(\mathcal{Z}|\mathcal{X}, \boldsymbol{\theta})}. \quad (23)$$

⁶It is more natural to understand this choice as using the canonical parameterization of a normal distribution. Please refer to my note on properties of normal distributions.

⁷For knowledge about convexity, please refer to the book *Convex optimization* by Stephen Boyd and Lieven Vandenberghe, Cambridge University Press. The PDF version of this book is available at <http://stanford.edu/~boyd/cvxbook/>.

Hence,

$$\ln p(\mathcal{X}|\boldsymbol{\theta}) = \left(\sum_{\mathcal{Z}} q(\mathcal{Z}) \right) \ln p(\mathcal{X}|\boldsymbol{\theta}) \quad (24)$$

$$= \sum_{\mathcal{Z}} q(\mathcal{Z}) \ln p(\mathcal{X}|\boldsymbol{\theta}) \quad (25)$$

$$= \sum_{\mathcal{Z}} q(\mathcal{Z}) \ln \left(\frac{p(\mathcal{X}, \mathcal{Z}|\boldsymbol{\theta})}{q(\mathcal{Z})} \frac{q(\mathcal{Z})}{p(\mathcal{Z}|\mathcal{X}, \boldsymbol{\theta})} \right) \quad (26)$$

$$= \sum_{\mathcal{Z}} \left(q(\mathcal{Z}) \ln \frac{p(\mathcal{X}, \mathcal{Z}|\boldsymbol{\theta})}{q(\mathcal{Z})} - q(\mathcal{Z}) \ln \frac{p(\mathcal{Z}|\mathcal{X}, \boldsymbol{\theta})}{q(\mathcal{Z})} \right) \quad (27)$$

$$= \sum_{\mathcal{Z}} q(\mathcal{Z}) \ln \frac{p(\mathcal{X}, \mathcal{Z}|\boldsymbol{\theta})}{q(\mathcal{Z})} + \text{KL}(q\|p) \quad (28)$$

$$= \mathcal{L}(q, \boldsymbol{\theta}) + \text{KL}(q\|p). \quad (29)$$

We have decomposed the incomplete data log-likelihood into two terms. The first term is $\mathcal{L}(q, \boldsymbol{\theta})$, defined as

$$\mathcal{L}(q, \boldsymbol{\theta}) = \sum_{\mathcal{Z}} q(\mathcal{Z}) \ln \frac{p(\mathcal{X}, \mathcal{Z}|\boldsymbol{\theta})}{q(\mathcal{Z})}. \quad (30)$$

The second term is a KL-divergence between q and the posterior

$$\text{KL}(q\|p) = - \sum_{\mathcal{Z}} q(\mathcal{Z}) \ln \left(\frac{p(\mathcal{Z}|\mathcal{X}, \boldsymbol{\theta})}{q(\mathcal{Z})} \right), \quad (31)$$

which was copied from Equation 21.

There are some nice properties of the KL-divergence. For example,

$$D(q\|p) \geq 0 \quad (32)$$

always holds, and the quality sign is true if and only if $q = p$.⁸ One direct consequence of this property is that

$$\mathcal{L}(q, \boldsymbol{\theta}) \leq \ln p(\mathcal{X}|\boldsymbol{\theta}) \quad (33)$$

always holds, and

$$\mathcal{L}(q, \boldsymbol{\theta}) = \ln p(\mathcal{X}|\boldsymbol{\theta}) \text{ if and only if } q(\mathcal{Z}) = p(\mathcal{Z}|\mathcal{X}, \boldsymbol{\theta}). \quad (34)$$

In other words, we have found a lower bound of $\ln p(\mathcal{X}|\boldsymbol{\theta})$. Hence, in order to maximize $\ln p(\mathcal{X}|\boldsymbol{\theta})$, we can perform two steps.

⁸For more properties of the KL-divergence, please refer to the book *Elements of information theory* by Thomas M. Cover and Joy A. Thomas, John Wiley & Sons, Inc.

- The first step is to make the lower bound $\mathcal{L}(q, \boldsymbol{\theta})$ equal to $\ln p(\mathcal{X}|\boldsymbol{\theta})$. As aforementioned, we know the equality hold if and only if $\hat{q}(\mathcal{Z}) = p(\mathcal{Z}|\mathcal{X}, \boldsymbol{\theta})$. Now we have

$$\ln p(\mathcal{X}|\boldsymbol{\theta}) = \mathcal{L}(\hat{q}, \boldsymbol{\theta}), \quad (35)$$

and \mathcal{L} only depends on $\boldsymbol{\theta}$ now. This is the Expectation step (E-step) in the EM algorithm.

- In the second step, we can maximize $\mathcal{L}(\hat{q}, \boldsymbol{\theta})$ with respect to $\boldsymbol{\theta}$. Since $\ln p(\mathcal{X}|\boldsymbol{\theta}) = \mathcal{L}(\hat{q}, \boldsymbol{\theta})$, an increase of $\mathcal{L}(\hat{q}, \boldsymbol{\theta})$ also means an increase of the log-likelihood $\ln p(\mathcal{X}|\boldsymbol{\theta})$. And, because we are maximizing $\mathcal{L}(\hat{q}, \boldsymbol{\theta})$ in this step, *the log-likelihood will always increase* if we are not already at a local minimum of the log-likelihood. This is the Maximization step (M-step) in the EM algorithm.

4.4 The E- & M-steps

In the E-step, we already know that we should set

$$\hat{q}(\mathcal{Z}) = p(\mathcal{Z}|\mathcal{X}, \boldsymbol{\theta}), \quad (36)$$

which is straightforward (at least in its mathematical form).

Then, how shall we maximize $\mathcal{L}(\hat{q}, \boldsymbol{\theta})$? We can substitute \hat{q} into the definition of \mathcal{L} . We will find the optimal $\boldsymbol{\theta}$ that maximizes \mathcal{L} after plugging in \hat{q} . However, note that \hat{q} involves $\boldsymbol{\theta}$ too. Hence, we need some more notations.

Suppose we are in the t -th iteration. In the E-step, \hat{q} is computed using the current parameter, as

$$\hat{q}(\mathcal{Z}) = p(\mathcal{Z}|\mathcal{X}, \boldsymbol{\theta}^{(t)}). \quad (37)$$

Then, \mathcal{L} becomes

$$\mathcal{L}(\hat{q}, \boldsymbol{\theta}) = \sum_{\mathcal{Z}} \hat{q}(\mathcal{Z}) \ln \frac{p(\mathcal{X}, \mathcal{Z}|\boldsymbol{\theta})}{\hat{q}(\mathcal{Z})} \quad (38)$$

$$= \sum_{\mathcal{Z}} \hat{q}(\mathcal{Z}) \ln p(\mathcal{X}, \mathcal{Z}|\boldsymbol{\theta}) - \hat{q}(\mathcal{Z}) \ln \hat{q}(\mathcal{Z}) \quad (39)$$

$$= \sum_{\mathcal{Z}} p(\mathcal{Z}|\mathcal{X}, \boldsymbol{\theta}^{(t)}) \ln p(\mathcal{X}, \mathcal{Z}|\boldsymbol{\theta}) + \text{const}, \quad (40)$$

in which $\text{const} = -\hat{q}(\mathcal{Z}) \ln \hat{q}(\mathcal{Z})$ does not involve the variable $\boldsymbol{\theta}$, hence can be ignored.

The term remaining is in fact an expectation, which we denote as $\mathcal{Q}(\boldsymbol{\theta}, \boldsymbol{\theta}^{(t)})$,

$$\mathcal{Q}(\boldsymbol{\theta}, \boldsymbol{\theta}^{(t)}) = \sum_{\mathcal{Z}} p(\mathcal{Z}|\mathcal{X}, \boldsymbol{\theta}^{(t)}) \ln p(\mathcal{X}, \mathcal{Z}|\boldsymbol{\theta}) \quad (41)$$

$$= \mathbb{E}_{\mathcal{Z}|\mathcal{X}, \boldsymbol{\theta}^{(t)}} [\ln p(\mathcal{X}, \mathcal{Z}|\boldsymbol{\theta})]. \quad (42)$$

That is, in the E-step, we compute the posterior of Z . In the M-step, we compute the expectation of the complete data log-likelihood $\ln p(\mathcal{X}, \mathcal{Z}|\boldsymbol{\theta})$

with respect to the posterior distribution $p(\mathcal{Z}|\mathcal{X}, \boldsymbol{\theta}^{(t)})$, and we maximize the expectation to get a better parameter estimate:

$$\boldsymbol{\theta}^{(t+1)} = \arg \max_{\boldsymbol{\theta}} \mathcal{Q}(\boldsymbol{\theta}, \boldsymbol{\theta}^{(t)}) = \arg \max_{\boldsymbol{\theta}} \mathbb{E}_{\mathcal{Z}|\mathcal{X}, \boldsymbol{\theta}^{(t)}} [\ln p(\mathcal{X}, \mathcal{Z}|\boldsymbol{\theta})] . \quad (43)$$

Thus, three computations are involved in EM: 1) posterior, 2) expectation, 3) maximization. We treat 1) as the E-step, and 2)+3) as the M-step. Some researchers prefer to treat 1)+2) as the E-step, and 3) as the M-step. However, no matter how the computations are attributed to different steps, the EM algorithm does not change.

4.5 The EM algorithm

Now we are ready to write down the EM algorithm.

Algorithm 1 The Expectation-Maximization Algorithm

- 1: $t \leftarrow 0$
- 2: Initialize the parameters to $\boldsymbol{\theta}^{(0)}$
- 3: The **E**(expectation)-step: Find $p(\mathcal{Z}|\mathcal{X}, \boldsymbol{\theta}^{(t)})$
- 4: The **M**(aximization)-step.1: Find the expectation

$$\mathcal{Q}(\boldsymbol{\theta}, \boldsymbol{\theta}^{(t)}) = \mathbb{E}_{\mathcal{Z}|\mathcal{X}, \boldsymbol{\theta}^{(t)}} [\ln p(\mathcal{X}, \mathcal{Z}|\boldsymbol{\theta})] \quad (44)$$

- 5: The **M**(aximization)-step.2: Find a new parameter estimate

$$\boldsymbol{\theta}^{(t+1)} = \arg \max_{\boldsymbol{\theta}} \mathcal{Q}(\boldsymbol{\theta}, \boldsymbol{\theta}^{(t)}) \quad (45)$$

- 6: $t \leftarrow t + 1$
 - 7: If the log-likelihood has not converged, go to the E-step again (Line 3)
-

4.6 Will EM converge?

The analysis of EM's convergence property is a complex topic. However, it is easy to show that the EM algorithm will help achieve higher likelihood and converge to a local minimum.

Let us consider two time steps $t - 1$ and t . From Equation 35, we get that:

$$\mathcal{L}(\hat{q}^{(t)}, \boldsymbol{\theta}^{(t)}) = \ln p(\mathcal{X}|\boldsymbol{\theta}^{(t)}), \quad (46)$$

$$\mathcal{L}(\hat{q}^{(t-1)}, \boldsymbol{\theta}^{(t-1)}) = \ln p(\mathcal{X}|\boldsymbol{\theta}^{(t-1)}). \quad (47)$$

Note that we have added the time index to the superscript of \hat{q} to emphasize that \hat{q} also changes among iterations.

Now because at the $(t - 1)$ -th iteration

$$\boldsymbol{\theta}^{(t)} = \arg \max_{\boldsymbol{\theta}} \mathcal{L}(\hat{q}^{(t-1)}, \boldsymbol{\theta}), \quad (48)$$

we have

$$\mathcal{L}(\hat{q}^{(t-1)}, \boldsymbol{\theta}^{(t)}) \geq \mathcal{L}(\hat{q}^{(t-1)}, \boldsymbol{\theta}^{(t-1)}). \quad (49)$$

Similarly, at the t -th iteration, based on Equation 33 and Equation 35, we have

$$\mathcal{L}(\hat{q}^{(t-1)}, \boldsymbol{\theta}^{(t)}) \leq \ln p(\mathcal{X}|\boldsymbol{\theta}^{(t)}) = \mathcal{L}(\hat{q}^{(t)}, \boldsymbol{\theta}^{(t)}). \quad (50)$$

Putting these equations together, we get

$$\ln p(\mathcal{X}|\boldsymbol{\theta}^{(t)}) = \mathcal{L}(\hat{q}^{(t)}, \boldsymbol{\theta}^{(t)}) \quad [\text{Use (46)}] \quad (51)$$

$$\geq \mathcal{L}(\hat{q}^{(t-1)}, \boldsymbol{\theta}^{(t)}) \quad [\text{Use (50)}] \quad (52)$$

$$\geq \mathcal{L}(\hat{q}^{(t-1)}, \boldsymbol{\theta}^{(t-1)}) \quad [\text{Use (49)}] \quad (53)$$

$$= \ln p(\mathcal{X}|\boldsymbol{\theta}^{(t-1)}). \quad [\text{Use (47)}] \quad (54)$$

Hence, EM will converge to a local minimum of the likelihood. However, the analysis of its convergence rate is very complex and beyond the scope of this introductory note.

5 EM for GMM

Now we can apply the EM algorithm to GMM.

The first thing is to compute the posterior. Using the Bayes' theorem, we have

$$p(z_{ij}|\mathbf{x}_j, \boldsymbol{\theta}^{(t)}) = \frac{p(\mathbf{x}_j, z_{ij}|\boldsymbol{\theta}^{(t)})}{p(\mathbf{x}_j|\boldsymbol{\theta}^{(t)})}, \quad (55)$$

in which z_{ij} can be 0 or 1, and $z_{ij} = 1$ is true if and only if \mathbf{x}_j is generated by the i -th Gaussian component.

Next, we will compute the \mathcal{Q} function, which is the expectation of the complete data log-likelihood $\ln p(\mathcal{X}, \mathcal{Z}|\boldsymbol{\theta})$ with respect to the posterior distribution we just found. The GMM complete data log-likelihood was already computed in Equation 18. For easier reference, we copy this equation here:

$$\sum_{j=1}^M \sum_{i=1}^N z_{ij} \left(\frac{1}{2} (\ln |\Sigma_i^{-1}| - (\mathbf{x}_j - \boldsymbol{\mu}_i)^T \Sigma_i^{-1} (\mathbf{x}_j - \boldsymbol{\mu}_i)) + \ln \alpha_i \right) + \text{const}, \quad (56)$$

The expectation of Equation 56 with respect to Z is

$$\sum_{j=1}^M \sum_{i=1}^N \gamma_{ij} \left(\frac{1}{2} (\ln |\Sigma_i^{-1}| - (\mathbf{x}_j - \boldsymbol{\mu}_i)^T \Sigma_i^{-1} (\mathbf{x}_j - \boldsymbol{\mu}_i)) + \ln \alpha_i \right), \quad (57)$$

where the constant term is ignored and γ_{ij} is the expectation of $z_{ij}|\mathbf{x}_j, \boldsymbol{\theta}^{(t)}$. In other words, we need to compute the expectation of the conditional distribution defined by Equation 55.

In Equation 55, the denominator does not depend on Z , and $p(\mathbf{x}_j|\boldsymbol{\theta}^{(t)})$ equals $\sum_{i=1}^N \alpha_i^{(t)} N(\mathbf{x}_j; \boldsymbol{\mu}_i^{(t)}, \Sigma_i^{(t)})$. For the numerator, we can directly compute its expectation, as

$$\mathbb{E} \left[p(\mathbf{x}_j, z_{ij}|\boldsymbol{\theta}^{(t)}) \right] = \mathbb{E} \left[p(z_{ij}|\boldsymbol{\theta}^{(t)}) p(\mathbf{x}_j|z_{ij}, \boldsymbol{\theta}^{(t)}) \right]. \quad (58)$$

Note that when $z_{ij} = 0$, we always have $p(\mathbf{x}_j|z_{ij}, \boldsymbol{\theta}^{(t)}) = 0$. Thus,

$$\mathbb{E} \left[p(z_{ij}|\boldsymbol{\theta}^{(t)}) p(\mathbf{x}_j|z_{ij}, \boldsymbol{\theta}^{(t)}) \right] = \Pr(z_{ij} = 1) p(\mathbf{x}_j|\boldsymbol{\mu}_i^{(t)}, \Sigma_i^{(t)}) \quad (59)$$

$$= \alpha_i^{(t)} N(\mathbf{x}_j; \boldsymbol{\mu}_i^{(t)}, \Sigma_i^{(t)}). \quad (60)$$

Hence, we have

$$\gamma_{ij} = \mathbb{E} \left[z_{ij}|\mathbf{x}_j, \boldsymbol{\theta}^{(t)} \right] \propto \alpha_i^{(t)} N(\mathbf{x}_j; \boldsymbol{\mu}_i^{(t)}, \Sigma_i^{(t)}), \quad (61)$$

or,

$$\gamma_{ij} = \mathbb{E} \left[z_{ij}|\mathbf{x}_j, \boldsymbol{\theta}^{(t)} \right] = \frac{\alpha_i^{(t)} N(\mathbf{x}_j; \boldsymbol{\mu}_i^{(t)}, \Sigma_i^{(t)})}{\sum_{k=1}^N \alpha_k^{(t)} N(\mathbf{x}_j; \boldsymbol{\mu}_k^{(t)}, \Sigma_k^{(t)})} \quad (62)$$

for $1 \leq i \leq N, 1 \leq j \leq M$.

After γ_{ij} is computed, Equation 57 is completely specified. We start the optimization from α_i . Because there is a constraint that $\sum_{i=1}^N \alpha_i = 1$, we use the Lagrange multiplier method, remove irrelevant terms, and get

$$\sum_{j=1}^M \sum_{i=1}^N \gamma_{ij} \ln \alpha_i + \lambda \left(\sum_{i=1}^N \alpha_i - 1 \right). \quad (63)$$

Setting the derivative to 0 gives us that for any $1 \leq i \leq N$,

$$\frac{\sum_{j=1}^M \gamma_{ij}}{\alpha_i} + \lambda = 0 \quad (64)$$

or, $\alpha_i = -\frac{\sum_{j=1}^M \gamma_{ij}}{\lambda}$. Because $\sum_{i=1}^N \alpha_i = 1$, we know that $\lambda = -\sum_{j=1}^M \sum_{i=1}^N \gamma_{ij}$.

Hence, $\alpha_i = \frac{\sum_{j=1}^M \gamma_{ij}}{\sum_{j=1}^M \sum_{i=1}^N \gamma_{ij}}$.

For notational simplicity, we define

$$m_i = \sum_{j=1}^M \gamma_{ij}. \quad (65)$$

From the definition of γ_{ij} , it is easy to prove that

$$\sum_{i=1}^N m_i = \sum_{i=1}^N \sum_{j=1}^M \gamma_{ij} = \sum_{j=1}^M \left(\sum_{i=1}^N \gamma_{ij} \right) = \sum_{j=1}^M 1 = M. \quad (66)$$

Then, we get the updating rule for α_i :

$$\alpha_i^{(t+1)} = \frac{m_i}{M}. \quad (67)$$

Furthermore, using similar steps in deriving the single Gaussian equations,⁹ it is easy to show that for any $1 \leq i \leq N$,

$$\boldsymbol{\mu}_i^{(t+1)} = \frac{\sum_{j=1}^M \gamma_{ij} \mathbf{x}_j}{m_i}, \quad (68)$$

$$\Sigma_i^{(t+1)} = \frac{\sum_{j=1}^M \gamma_{ij} \left(\mathbf{x}_j - \boldsymbol{\mu}_i^{(t+1)} \right) \left(\mathbf{x}_j - \boldsymbol{\mu}_i^{(t+1)} \right)^T}{m_i}. \quad (69)$$

Putting these results together, we have the complete set of updating rules for GMM. If at iteration t , the parameter are estimated as $\alpha_i^{(t)}$, $\boldsymbol{\mu}_i^{(t)}$, and $\Sigma_i^{(t)}$ for $1 \leq i \leq N$, the EM algorithm updates these parameters as (for $1 \leq i \leq N$, $1 \leq j \leq M$)

$$\gamma_{ij} = \frac{\alpha_i N(\mathbf{x}_j; \boldsymbol{\mu}_i^{(t)}, \Sigma_i^{(t)})}{\sum_{k=1}^N \alpha_k N(\mathbf{x}_j; \boldsymbol{\mu}_k^{(t)}, \Sigma_k^{(t)})}, \quad (70)$$

$$m_i = \sum_{j=1}^M \gamma_{ij}, \quad (71)$$

$$\boldsymbol{\mu}_i^{(t+1)} = \frac{\sum_{j=1}^M \gamma_{ij} \mathbf{x}_j}{m_i}, \quad (72)$$

$$\Sigma_i^{(t+1)} = \frac{\sum_{j=1}^M \gamma_{ij} \left(\mathbf{x}_j - \boldsymbol{\mu}_i^{(t+1)} \right) \left(\mathbf{x}_j - \boldsymbol{\mu}_i^{(t+1)} \right)^T}{m_i}. \quad (73)$$

Exercises

1. Derive the updating equations for Gaussian Mixture Models by yourself. You should not refer to Section 5 during your derivation. If you have just finished reading Section 5, wait for at least 2 to 3 hours before working on this problem.
2. In this problem, we will use the Expectation-Maximization method to learn parameters in a hidden Markov model (HMM). As will be shown in this problem, the Baum-Welch algorithm is indeed performing EM updates. To work out the solution for this problem, you will also need knowledge and facts learned in the HMM and information theory notes.
We will use the notations in the HMM note. For your convenience, the notations are repeated as follows.

⁹Please refer to my note on properties of normal distributions.

- There are N discrete states, denoted by symbols S_1, S_2, \dots, S_N .
- There are M output discrete symbols, denoted by V_1, V_2, \dots, V_M .
- Assuming one sequence with T time steps, whose hidden state is Q_t and whose observed output is O_t at time t ($1 \leq t \leq T$). We use q_t and o_t to denote the indexes for state and output symbols at time t , respectively, i.e., $Q_t = S_{q_t}$ and $O_t = V_{o_t}$.
- The notation $1 : t$ denotes all the ordered time steps between 1 and t . For example, $o_{1:T}$ is the sequence of all observed output symbols.
- An HMM has parameters $\lambda = (\boldsymbol{\pi}, A, B)$, where $\boldsymbol{\pi} \in \mathbb{R}^N$ specifies the initial state distribution, $A \in \mathbb{R}^{N \times N}$ is the state transition matrix, and $B \in \mathbb{R}^{N \times M}$ is the observation probability matrix. Note that $A_{ij} = \Pr(Q_t = S_j | Q_{t-1} = S_i)$ and $b_j(k) = \Pr(O_t = V_k | Q_t = S_j)$ are elements of A and B , respectively.
- In this problem, we use a variable r to denote the index of EM iterations. Hence, $\lambda^{(1)}$ are the initial parameters.
- Various probabilities have been defined in the HMM note, denoted by $\alpha_t(i)$, $\beta_t(i)$, $\gamma_t(i)$, $\delta_t(i)$ and $\xi_t(i, j)$. In this problem, we assume that at the r -th iterations, $\lambda^{(r)}$ are known and these probabilities are computed using $\lambda^{(r)}$.

The purpose of this problem is to use the EM algorithm to find $\lambda^{(r+1)}$ using a training sequence $o_{1:T}$ and $\lambda^{(r)}$, by treating Q and O as the hidden and observed random variables, respectively.

(a) Suppose the hidden variables can be observed as $S_{q_1}, S_{q_2}, \dots, S_{q_T}$. Show that the complete data log-likelihood is

$$\ln \pi_{q_1} + \sum_{t=1}^{T-1} \ln A_{q_t q_{t+1}} + \sum_{t=1}^T \ln b_{q_t}(o_t). \quad (74)$$

(b) The expectation of Equation 74 with respect to the hidden variables Q_t (conditioned on $o_{1:T}$ and $\lambda^{(r)}$) forms an auxiliary function $\mathcal{Q}(\lambda, \lambda^{(r)})$ (the E-step). Show that the expectation of the first term in Equation 74 equals $\sum_{i=1}^N \gamma_1(i) \ln \pi_i$, i.e.,

$$\mathbb{E}_{Q_{1:T}}[\ln \pi_{Q_1}] = \sum_{i=1}^N \gamma_1(i) \ln \pi_i. \quad (75)$$

(c) Because the parameter $\boldsymbol{\pi}$ only hinges on Equation 75, the update rule for $\boldsymbol{\pi}$ can be found by maximizing this equation. Prove that we should set $\pi_i^{(r+1)} = \gamma_1(i)$ in the M-step. Note that $\gamma_1(i)$ is computed using $\lambda^{(r)}$

as parameter values. (Hint: The right hand side of Equation 75 is related to the cross entropy.)

(d) The second part of the E-step calculates the expectation of the middle term in Equation 74. Show that

$$\mathbb{E}_{Q_{1:T}} \left[\sum_{t=1}^{T-1} \ln A_{q_t q_{t+1}} \right] = \sum_{i=1}^N \sum_{j=1}^N \left(\sum_{t=1}^{T-1} \xi_t(i, j) \right) \ln A_{ij}. \quad (76)$$

(e) For the M-step relevant to A , prove that we should set

$$A_{ij}^{(r+1)} = \frac{\sum_{t=1}^{T-1} \xi_t(i, j)}{\sum_{t=1}^{T-1} \gamma_t(i)}. \quad (77)$$

(f) The final part of the E-step calculates the expectation of the last term in Equation 74. Show that

$$\mathbb{E}_{Q_{1:T}} \left[\sum_{t=1}^T \ln b_{q_t}(o_t) \right] = \sum_{j=1}^M \sum_{k=1}^M \sum_{t=1}^T \llbracket o_t = k \rrbracket \gamma_t(j). \quad (78)$$

(g) For the M-step relevant to B , prove that we should set

$$b_j^{(r+1)}(k) = \frac{\sum_{t=1}^T \llbracket o_t = k \rrbracket \gamma_t(j)}{\sum_{t=1}^T \gamma_t(j)}, \quad (79)$$

in which $\llbracket \cdot \rrbracket$ is the indicator function.

(h) Are these results obtained using EM the same as those in the Baum-Welch?