

What is the expectation maximization algorithm?

Chuong B Do & Serafim Batzoglou

The expectation maximization algorithm arises in many computational biology applications that involve probabilistic models. What is it good for, and how does it work?

Probabilistic models, such as hidden Markov models or Bayesian networks, are commonly used to model biological data. Much of their popularity can be attributed to the existence of efficient and robust procedures for learning parameters from observations. Often, however, the only data available for training a probabilistic model are incomplete. Missing values can occur, for example, in medical diagnosis, where patient histories generally include results from a limited battery of tests. Alternatively, in gene expression clustering, incomplete data arise from the intentional omission of gene-to-cluster assignments in the probabilistic model. The expectation maximization algorithm enables parameter estimation in probabilistic models with incomplete data.

A coin-flipping experiment

As an example, consider a simple coin-flipping experiment in which we are given a pair of coins A and B of unknown biases, θ_A and θ_B , respectively (that is, on any given flip, coin A will land on heads with probability θ_A and tails with probability $1-\theta_A$ and similarly for coin B). Our goal is to estimate $\theta = (\theta_A, \theta_B)$ by repeating the following procedure five times: randomly choose one of the two coins (with equal probability), and perform ten independent coin tosses with the selected coin. Thus, the entire procedure involves a total of 50 coin tosses (Fig. 1a).

During our experiment, suppose that we keep track of two vectors $x = (x_1, x_2, \dots, x_5)$ and

$z = (z_1, z_2, \dots, z_5)$, where $x_i \in \{0, 1, \dots, 10\}$ is the number of heads observed during the i th set of tosses, and $z_i \in \{A, B\}$ is the identity of the coin used during the i th set of tosses. Parameter estimation in this setting is known as the complete data case in that the values of all relevant random variables in our model (that is, the result of each coin flip and the type of coin used for each flip) are known.

Here, a simple way to estimate θ_A and θ_B is to return the observed proportions of heads for each coin:

$$\hat{\theta}_A = \frac{\text{\# of heads using coin A}}{\text{total \# of flips using coin A}} \quad (1)$$

and

$$\hat{\theta}_B = \frac{\text{\# of heads using coin B}}{\text{total \# of flips using coin B}}$$

This intuitive guess is, in fact, known in the statistical literature as maximum likelihood estimation (roughly speaking, the maximum likelihood method assesses the quality of a statistical model based on the probability it assigns to the observed data). If $\log P(x, z; \theta)$ is the logarithm of the joint probability (or log-likelihood) of obtaining any particular vector of observed head counts x and coin types z , then the formulas in (1) solve for the parameters $\hat{\theta} = (\hat{\theta}_A, \hat{\theta}_B)$ that maximize $\log P(x, z; \theta)$.

Now consider a more challenging variant of the parameter estimation problem in which we are given the recorded head counts x but not the identities z of the coins used for each set of tosses. We refer to z as hidden variables or latent factors. Parameter estimation in this new setting is known as the incomplete data case. This time, computing proportions of heads for each coin is no longer possible, because we

don't know the coin used for each set of tosses. However, if we had some way of completing the data (in our case, guessing correctly which coin was used in each of the five sets), then we could reduce parameter estimation for this problem with incomplete data to maximum likelihood estimation with complete data.

One iterative scheme for obtaining completions could work as follows: starting from some initial parameters, $\hat{\theta}^{(i)} = (\hat{\theta}_A^{(i)}, \hat{\theta}_B^{(i)})$, determine for each of the five sets whether coin A or coin B was more likely to have generated the observed flips (using the current parameter estimates). Then, assume these completions (that is, guessed coin assignments) to be correct, and apply the regular maximum likelihood estimation procedure to get $\hat{\theta}^{(i+1)}$. Finally, repeat these two steps until convergence. As the estimated model improves, so too will the quality of the resulting completions.

The expectation maximization algorithm is a refinement on this basic idea. Rather than picking the single most likely completion of the missing coin assignments on each iteration, the expectation maximization algorithm computes probabilities for each possible completion of the missing data, using the current parameters $\hat{\theta}^{(i)}$. These probabilities are used to create a weighted training set consisting of all possible completions of the data. Finally, a modified version of maximum likelihood estimation that deals with weighted training examples provides new parameter estimates, $\hat{\theta}^{(i+1)}$. By using weighted training examples rather than choosing the single best completion, the expectation maximization algorithm accounts for the confidence of the model in each completion of the data (Fig. 1b).

In summary, the expectation maximization algorithm alternates between the steps

Chuong B. Do and Serafim Batzoglou are in the Computer Science Department, Stanford University, 318 Campus Drive, Stanford, California 94305-5428, USA.
e-mail: chuong@cs.stanford.edu

of guessing a probability distribution over completions of missing data given the current model (known as the E-step) and then re-estimating the model parameters using these completions (known as the M-step). The name 'E-step' comes from the fact that one does not usually need to form the probability distribution over completions explicitly, but rather need only compute 'expected' sufficient statistics over these completions. Similarly, the name 'M-step' comes from the fact that model reestimation can be thought of as 'maximization' of the expected log-likelihood of the data.

Introduced as early as 1950 by Ceppellini *et al.*¹ in the context of gene frequency estimation, the expectation maximization algorithm

was analyzed more generally by Hartley² and by Baum *et al.*³ in the context of hidden Markov models, where it is commonly known as the Baum-Welch algorithm. The standard reference on the expectation maximization algorithm and its convergence is Dempster *et al.*⁴.

Mathematical foundations

How does the expectation maximization algorithm work? More importantly, why is it even necessary?

The expectation maximization algorithm is a natural generalization of maximum likelihood estimation to the incomplete data case. In particular, expectation maximization attempts to find the parameters $\hat{\theta}$ that maximize the

log probability $\log P(x; \theta)$ of the observed data. Generally speaking, the optimization problem addressed by the expectation maximization algorithm is more difficult than the optimization used in maximum likelihood estimation. In the complete data case, the objective function $\log P(x; \theta)$ has a single global optimum, which can often be found in closed form (e.g., equation 1). In contrast, in the incomplete data case the function $\log P(x; \theta)$ has multiple local maxima and no closed form solution.

To deal with this, the expectation maximization algorithm reduces the difficult task of optimizing $\log P(x; \theta)$ into a sequence of simpler optimization subproblems, whose objective functions have unique global maxima that can often be computed in closed form. These subproblems are chosen in a way that guarantees their corresponding solutions $\hat{\theta}^{(1)}, \hat{\theta}^{(2)}, \dots$ and will converge to a local optimum of $\log P(x; \theta)$.

More specifically, the expectation maximization algorithm alternates between two phases. During the E-step, expectation maximization chooses a function g_t that lower bounds $\log P(x; \theta)$ everywhere, and for which $g_t(\hat{\theta}^{(t)}) = \log P(x; \hat{\theta}^{(t)})$. During the M-step, the expectation maximization algorithm moves to a new parameter set $\hat{\theta}^{(t+1)}$ that maximizes g_t . As the value of the lower-bound g_t matches the objective function at $\hat{\theta}^{(t)}$, it follows that $\log P(x; \hat{\theta}^{(t)}) = g_t(\hat{\theta}^{(t)}) \leq g_t(\hat{\theta}^{(t+1)}) = \log P(x; \hat{\theta}^{(t+1)})$ — so the objective function monotonically increases during each iteration of expectation maximization! A graphical illustration of this argument is provided in **Supplementary Figure 1** online, and a concise mathematical derivation of the expectation maximization algorithm is given in **Supplementary Note 1** online.

As with most optimization methods for nonconcave functions, the expectation maximization algorithm comes with guarantees only of convergence to a local maximum of the objective function (except in degenerate cases). Running the procedure using multiple initial starting parameters is often helpful; similarly, initializing parameters in a way that breaks symmetry in models is also important. With this limited set of tricks, the expectation maximization algorithm provides a simple and robust tool for parameter estimation in models with incomplete data. In theory, other numerical optimization techniques, such as gradient descent or Newton-Raphson, could be used instead of expectation maximization; in practice, however, expectation maximization has the advantage of being simple, robust and easy to implement.

Applications

Many probabilistic models in computational biology include latent variables. In some

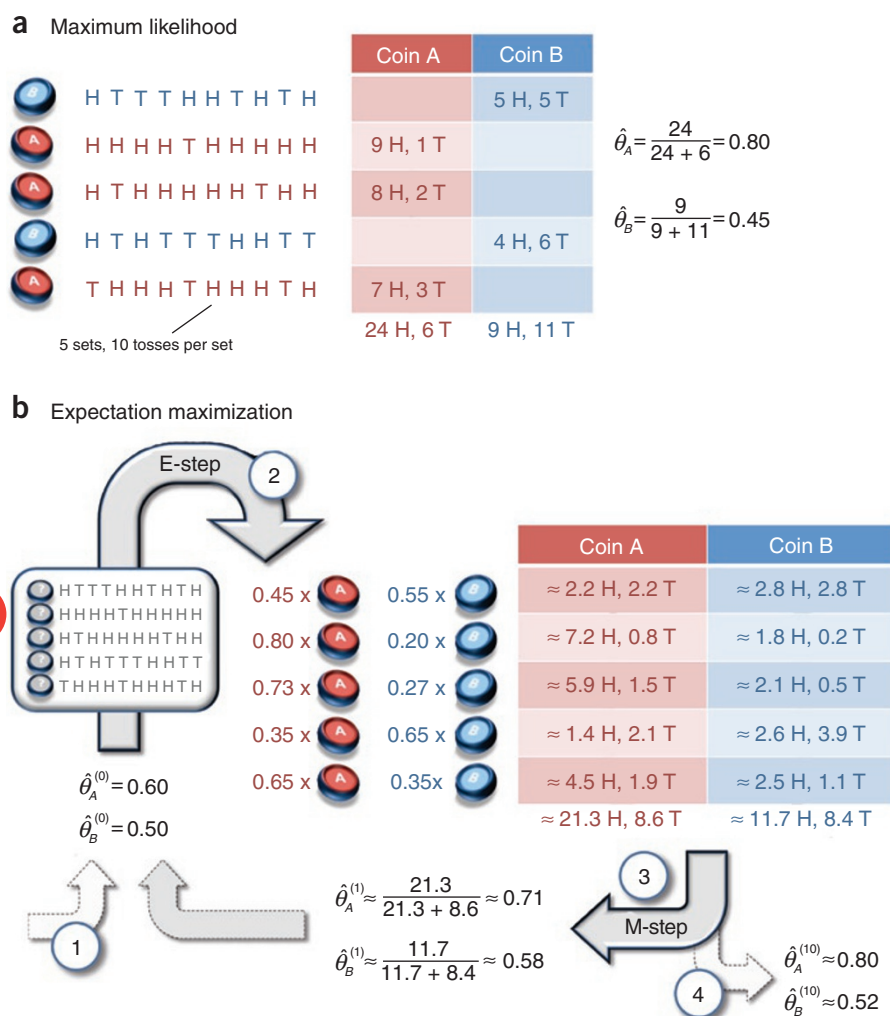


Figure 1 Parameter estimation for complete and incomplete data. **(a)** Maximum likelihood estimation. For each set of ten tosses, the maximum likelihood procedure accumulates the counts of heads and tails for coins A and B separately. These counts are then used to estimate the coin biases. **(b)** Expectation maximization. 1. EM starts with an initial guess of the parameters. 2. In the E-step, a probability distribution over possible completions is computed using the current parameters. The counts shown in the table are the expected numbers of heads and tails according to this distribution. 3. In the M-step, new parameters are determined using the current completions. 4. After several repetitions of the E-step and M-step, the algorithm converges.

cases, these latent variables are present due to missing or corrupted data; in most applications of expectation maximization to computational biology, however, the latent factors are intentionally included, and parameter learning itself provides a mechanism for knowledge discovery.

For instance, in gene expression clustering⁵, we are given microarray gene expression measurements for thousands of genes under varying conditions, and our goal is to group the observed expression vectors into distinct clusters of related genes. One approach is to model the vector of expression measurements for each gene as being sampled from a multivariate Gaussian distribution (a generalization of a standard Gaussian distribution to multiple correlated variables) associated with that gene's cluster. In this case, the observed data x correspond to microarray measurements, the unobserved latent factors z are the assignments of genes to clusters, and the parameters θ include the means and covariance matrices of the multivariate Gaussian distributions representing the expression patterns for each cluster. For parameter learning, the expectation maximization algorithm alternates between computing probabilities for assignments of each gene to each cluster (E-step) and updating the cluster means and covariances based on the set of genes predominantly belonging to that cluster (M-step). This can be thought of as a 'soft' variant of the popular k-means clustering algorithm, in which one alternates between 'hard' assignments of genes to clusters and reestimation of cluster means based on their assigned genes.

In motif finding⁶, we are given a set of unaligned DNA sequences and asked to identify a pattern of length W that is present (though possibly with minor variations) in every sequence from the set. To apply the expectation maximization algorithm, we model the instance of the motif in each sequence as having each letter sampled independently from a position-specific distribution over letters,

and the remaining letters in each sequence as coming from some fixed background distribution. The observed data x consist of the letters of sequences, the unobserved latent factors z include the starting position of the motif in each sequence and the parameters θ describe the position-specific letter frequencies for the motif. Here, the expectation maximization algorithm involves computing the probability distribution over motif start positions for each sequence (E-step) and updating the motif letter frequencies based on the expected letter counts for each position in the motif (M-step).

In the haplotype inference problem⁷, we are given the unphased genotypes of individuals from some population, where each unphased genotype consists of unordered pairs of single-nucleotide polymorphisms (SNPs) taken from homologous chromosomes of the individual. Contiguous blocks of SNPs inherited from a single chromosome are called haplotypes. Assuming for simplicity that each individual's genotype is a combination of two haplotypes (one maternal and one paternal), the goal of haplotype inference is to determine a small set of haplotypes that best explain all of the unphased genotypes observed in the population. Here, the observed data x are the known unphased genotypes for each individual, the unobserved latent factors z are putative assignments of unphased genotypes to pairs of haplotypes and the parameters θ describe the frequencies of each haplotype in the population. The expectation maximization algorithm alternates between using the current haplotype frequencies to estimate probability distributions over phasing assignments for each unphased genotype (E-step) and using the expected phasing assignments to refine estimates of haplotype frequencies (M-step).

Other problems in which the expectation maximization algorithm plays a prominent role include learning profiles of protein domains⁸ and RNA families⁹, discovery of

transcriptional modules¹⁰, tests of linkage disequilibrium¹¹, protein identification¹² and medical imaging¹³.

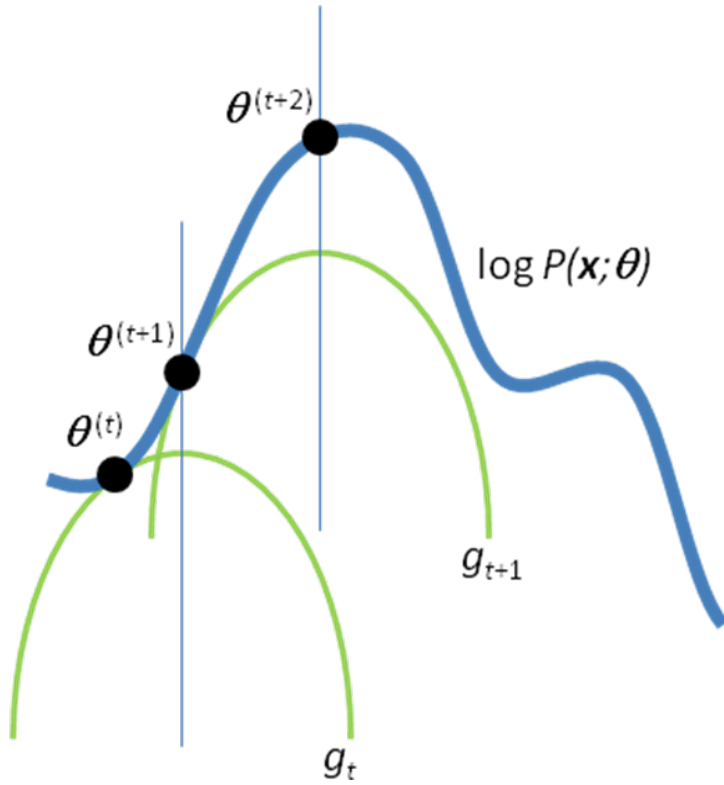
In each case, expectation maximization provides a simple, easy-to-implement and efficient tool for learning parameters of a model; once these parameters are known, we can use probabilistic inference to ask interesting queries about the model. For example, what cluster does a particular gene most likely belong to? What is the most likely starting location of a motif in a particular sequence? What are the most likely haplotype blocks making up the genotype of a specific individual? By providing a straightforward mechanism for parameter learning in all of these models, expectation maximization provides a mechanism for building and training rich probabilistic models for biological applications.

Note: Supplementary information is available on the Nature Biotechnology website.

ACKNOWLEDGMENTS

C.B.D. is supported in part by an National Science Foundation (NSF) Graduate Fellowship. S.B. wishes to acknowledge support by the NSF CAREER Award. We thank four anonymous referees for helpful suggestions.

1. Ceppellini, R., Siniscalco, M. & Smith, C.A. *Ann. Hum. Genet.* **20**, 97–115 (1955).
2. Hartley, H. *Biometrics* **14**, 174–194 (1958).
3. Baum, L.E., Petrie, T., Soules, G. & Weiss, N. *Ann. Math. Stat.* **41**, 164–171 (1970).
4. Dempster, A.P., Laird, N.M. & Rubin, D.B. *J. R. Stat. Soc. Ser. B* **39**, 1–38 (1977).
5. D'haeseleer, P. *Nat. Biotechnol.* **23**, 1499–1501 (2005).
6. Lawrence, C.E. & Reilly, A.A. *Proteins* **7**, 41–51 (1990).
7. Excoffier, L. & Slatkin, M. *Mol. Biol. Evol.* **12**, 921–927 (1995).
8. Krogh, A., Brown, M., Mian, I.S., Sjölander, K. & Haussler, D. *J. Mol. Biol.* **235**, 1501–1543 (1994).
9. Eddy, S.R. & Durbin, R. *Nucleic Acids Res.* **22**, 2079–2088 (1994).
10. Segal, E., Yelensky, R. & Koller, D. *Bioinformatics* **19**, i273–i282 (2003).
11. Slatkin, M. & Excoffier, L. *Heredity* **76**, 377–383 (1996).
12. Nesvizhskii, A.I., Keller, A., Kolker, E. & Aebersold, R. *Anal. Chem.* **75**, 4646–4658 (2003).
13. De Pierro, A.R. *IEEE Trans. Med. Imaging* **14**, 132–137 (1995).



Supplementary Figure 1 Convergence of the EM algorithm. Starting from initial parameters $\theta^{(t)}$, the E-step of the EM algorithm constructs a function g_t that lower-bounds the objective function $\log P(x; \theta)$. In the M-step, $\theta^{(t+1)}$ is computed as the maximum of g_t . In the next E-step, a new lower-bound g_{t+1} is constructed; maximization of g_{t+1} in the next M-step gives $\theta^{(t+2)}$, etc.

Supplementary Note 1

Here, we provide a short derivation of the EM algorithm based on the idea of bound maximization. A more general (though somewhat subtle) argument, which leads to a number of variants of the EM algorithm, was described in:

Neal RM and Hinton GE. 1998. A view of the EM algorithm that justifies incremental, sparse, and other variants. In MI Jordan, ed., *Learning in Graphical Models*, Kluwer Academic Publishers, 355-358.

Mathematically, the EM algorithm derives from the fact that for any probability distribution $Q(\mathbf{z})$,

$$\log \left(\sum_{\mathbf{z}} P(\mathbf{x}, \mathbf{z}; \boldsymbol{\theta}) \right) = \log \left(\sum_{\mathbf{z}} Q(\mathbf{z}) \cdot \frac{P(\mathbf{x}, \mathbf{z}; \boldsymbol{\theta})}{Q(\mathbf{z})} \right) \geq \sum_{\mathbf{z}} Q(\mathbf{z}) \log \left(\frac{P(\mathbf{x}, \mathbf{z}; \boldsymbol{\theta})}{Q(\mathbf{z})} \right), \quad (3)$$

where the inequality is tight (i.e. holds with equality) whenever $Q(\mathbf{z}) = P(\mathbf{z}|\mathbf{x}; \boldsymbol{\theta})$. The derivation of (3) relies on a tool from mathematical analysis known as **Jensen's inequality** for concave functions.¹

Now, consider the update rule, $\hat{\boldsymbol{\theta}}^{(t+1)} = \arg \max_{\boldsymbol{\theta}} g_t(\boldsymbol{\theta})$, where

$$g_t(\boldsymbol{\theta}) = \sum_{\mathbf{z}} P(\mathbf{z}|\mathbf{x}; \hat{\boldsymbol{\theta}}^{(t)}) \log \left(\frac{P(\mathbf{x}, \mathbf{z}; \boldsymbol{\theta})}{P(\mathbf{z}|\mathbf{x}; \hat{\boldsymbol{\theta}}^{(t)})} \right). \quad (4)$$

Applying the tightness conditions of equation (3), $g_t(\hat{\boldsymbol{\theta}}^{(t)}) = \log P(\mathbf{x}; \hat{\boldsymbol{\theta}}^{(t)})$. However, observe that $g_t(\hat{\boldsymbol{\theta}}^{(t)}) \leq g_t(\hat{\boldsymbol{\theta}}^{(t+1)})$ by definition of the update rule. Furthermore, $g_t(\hat{\boldsymbol{\theta}}^{(t+1)}) \leq \log P(\mathbf{x}; \hat{\boldsymbol{\theta}}^{(t+1)})$ since (3) guarantees that $g_t(\boldsymbol{\theta})$ is a lower bound on $\log P(\mathbf{x}; \boldsymbol{\theta})$ for any parameter vector $\boldsymbol{\theta}$. Therefore, the update rule results in monotonic improvement of the maximum likelihood objective for incomplete data.

To see the connection between (4) and the description of the EM algorithm given in the text, consider the following equivalent update rule:

$$\hat{\boldsymbol{\theta}}^{(t+1)} = \arg \max_{\boldsymbol{\theta}} \sum_{\mathbf{z}} P(\mathbf{z}|\mathbf{x}; \hat{\boldsymbol{\theta}}^{(t)}) \log P(\mathbf{x}, \mathbf{z}; \boldsymbol{\theta}). \quad (5)$$

Here, equivalence follows from the fact that the objective function of (5) differs from $g_t(\boldsymbol{\theta})$ by a constant offset which does not depend on $\boldsymbol{\theta}$. In this final form, we see that the EM update rule effectively maximizes the log-likelihood of a dataset expanded to contain all possible completions \mathbf{z}

¹ In brief, Jensen's inequality states that a concave function (e.g., $\log(\cdot)$) of the expectation of a random variable (e.g., \mathbf{y}) is never less than expected value of that concave function applied to the random variable (i.e., $\log(E[\mathbf{y}]) \geq E[\log(\mathbf{y})]$); furthermore, equality holds if and only if the random variable is constant with probability 1 (i.e., $\mathbf{y} = E[\mathbf{y}]$). Equation

(3) follows from Jensen's inequality, using the random variable $y = \frac{P(\mathbf{x}, \mathbf{z}; \boldsymbol{\theta})}{Q(\mathbf{z})}$.

of the unobserved variables, where each completion is weighted by the posterior probability, $P(\mathbf{z}|\mathbf{x}; \hat{\boldsymbol{\theta}}^{(t)})$.

CS229 Lecture notes

Andrew Ng

Part IX

The EM algorithm

In the previous set of notes, we talked about the EM algorithm as applied to fitting a mixture of Gaussians. In this set of notes, we give a broader view of the EM algorithm, and show how it can be applied to a large family of estimation problems with latent variables. We begin our discussion with a very useful result called **Jensen's inequality**

1 Jensen's inequality

Let f be a function whose domain is the set of real numbers. Recall that f is a convex function if $f''(x) \geq 0$ (for all $x \in \mathbb{R}$). In the case of f taking vector-valued inputs, this is generalized to the condition that its hessian H is positive semi-definite ($H \geq 0$). If $f''(x) > 0$ for all x , then we say f is **strictly** convex (in the vector-valued case, the corresponding statement is that H must be positive definite, written $H > 0$). Jensen's inequality can then be stated as follows:

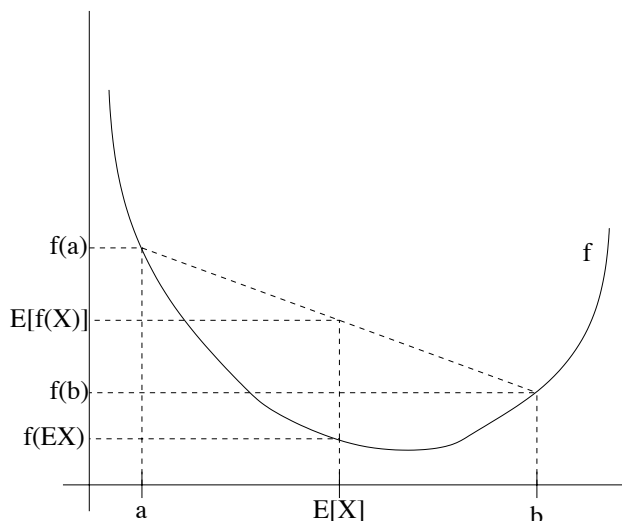
Theorem. Let f be a convex function, and let X be a random variable. Then:

$$\mathbb{E}[f(X)] \geq f(\mathbb{E}X).$$

Moreover, if f is strictly convex, then $\mathbb{E}[f(X)] = f(\mathbb{E}X)$ holds true if and only if $X = \mathbb{E}[X]$ with probability 1 (i.e., if X is a constant).

Recall our convention of occasionally dropping the parentheses when writing expectations, so in the theorem above, $f(\mathbb{E}X) = f(\mathbb{E}[X])$.

For an interpretation of the theorem, consider the figure below.



Here, f is a convex function shown by the solid line. Also, X is a random variable that has a 0.5 chance of taking the value a , and a 0.5 chance of taking the value b (indicated on the x -axis). Thus, the expected value of X is given by the midpoint between a and b .

We also see the values $f(a)$, $f(b)$ and $f(E[X])$ indicated on the y -axis. Moreover, the value $E[f(X)]$ is now the midpoint on the y -axis between $f(a)$ and $f(b)$. From our example, we see that because f is convex, it must be the case that $E[f(X)] \geq f(E[X])$.

Incidentally, quite a lot of people have trouble remembering which way the inequality goes, and remembering a picture like this is a good way to quickly figure out the answer.

Remark. Recall that f is [strictly] concave if and only if $-f$ is [strictly] convex (i.e., $f''(x) \leq 0$ or $H \leq 0$). Jensen's inequality also holds for concave functions f , but with the direction of all the inequalities reversed ($E[f(X)] \leq f(E[X])$, etc.).

2 The EM algorithm

Suppose we have an estimation problem in which we have a training set $\{x^{(1)}, \dots, x^{(m)}\}$ consisting of m independent examples. We wish to fit the

parameters of a model $p(x, z)$ to the data, where the likelihood is given by

$$\begin{aligned}\ell(\theta) &= \sum_{i=1}^m \log p(x; \theta) \\ &= \sum_{i=1}^m \log \sum_z p(x, z; \theta).\end{aligned}$$

But, explicitly finding the maximum likelihood estimates of the parameters θ may be hard. Here, the $z^{(i)}$'s are the latent random variables; and it is often the case that if the $z^{(i)}$'s were observed, then maximum likelihood estimation would be easy.

In such a setting, the EM algorithm gives an efficient method for maximum likelihood estimation. Maximizing $\ell(\theta)$ explicitly might be difficult, and our strategy will be to instead repeatedly construct a lower-bound on ℓ (E-step), and then optimize that lower-bound (M-step).

For each i , let Q_i be some distribution over the z 's ($\sum_z Q_i(z) = 1$, $Q_i(z) \geq 0$). Consider the following:¹

$$\sum_i \log p(x^{(i)}; \theta) = \sum_i \log \sum_{z^{(i)}} p(x^{(i)}, z^{(i)}; \theta) \quad (1)$$

$$= \sum_i \log \sum_{z^{(i)}} Q_i(z^{(i)}) \frac{p(x^{(i)}, z^{(i)}; \theta)}{Q_i(z^{(i)})} \quad (2)$$

$$\geq \sum_i \sum_{z^{(i)}} Q_i(z^{(i)}) \log \frac{p(x^{(i)}, z^{(i)}; \theta)}{Q_i(z^{(i)})} \quad (3)$$

The last step of this derivation used Jensen's inequality. Specifically, $f(x) = \log x$ is a concave function, since $f''(x) = -1/x^2 < 0$ over its domain $x \in \mathbb{R}^+$. Also, the term

$$\sum_{z^{(i)}} Q_i(z^{(i)}) \left[\frac{p(x^{(i)}, z^{(i)}; \theta)}{Q_i(z^{(i)})} \right]$$

in the summation is just an expectation of the quantity $[p(x^{(i)}, z^{(i)}; \theta)/Q_i(z^{(i)})]$ with respect to $z^{(i)}$ drawn according to the distribution given by Q_i . By Jensen's inequality, we have

$$f \left(\mathbb{E}_{z^{(i)} \sim Q_i} \left[\frac{p(x^{(i)}, z^{(i)}; \theta)}{Q_i(z^{(i)})} \right] \right) \geq \mathbb{E}_{z^{(i)} \sim Q_i} \left[f \left(\frac{p(x^{(i)}, z^{(i)}; \theta)}{Q_i(z^{(i)})} \right) \right],$$

¹If z were continuous, then Q_i would be a density, and the summations over z in our discussion are replaced with integrals over z .

where the “ $z^{(i)} \sim Q_i$ ” subscripts above indicate that the expectations are with respect to $z^{(i)}$ drawn from Q_i . This allowed us to go from Equation (2) to Equation (3).

Now, for *any* set of distributions Q_i , the formula (3) gives a lower-bound on $\ell(\theta)$. There’re many possible choices for the Q_i ’s. Which should we choose? Well, if we have some current guess θ of the parameters, it seems natural to try to make the lower-bound tight at that value of θ . I.e., we’ll make the inequality above hold with equality at our particular value of θ . (We’ll see later how this enables us to prove that $\ell(\theta)$ increases monotonically with successive iterations of EM.)

To make the bound tight for a particular value of θ , we need for the step involving Jensen’s inequality in our derivation above to hold with equality. For this to be true, we know it is sufficient that the expectation be taken over a “constant”-valued random variable. I.e., we require that

$$\frac{p(x^{(i)}, z^{(i)}; \theta)}{Q_i(z^{(i)})} = c$$

for some constant c that does not depend on $z^{(i)}$. This is easily accomplished by choosing

$$Q_i(z^{(i)}) \propto p(x^{(i)}, z^{(i)}; \theta).$$

Actually, since we know $\sum_z Q_i(z^{(i)}) = 1$ (because it is a distribution), this further tells us that

$$\begin{aligned} Q_i(z^{(i)}) &= \frac{p(x^{(i)}, z^{(i)}; \theta)}{\sum_z p(x^{(i)}, z; \theta)} \\ &= \frac{p(x^{(i)}, z^{(i)}; \theta)}{p(x^{(i)}; \theta)} \\ &= p(z^{(i)} | x^{(i)}; \theta) \end{aligned}$$

Thus, we simply set the Q_i ’s to be the posterior distribution of the $z^{(i)}$ ’s given $x^{(i)}$ and the setting of the parameters θ .

Now, for this choice of the Q_i ’s, Equation (3) gives a lower-bound on the loglikelihood ℓ that we’re trying to maximize. This is the E-step. In the M-step of the algorithm, we then maximize our formula in Equation (3) with respect to the parameters to obtain a new setting of the θ ’s. Repeatedly carrying out these two steps gives us the EM algorithm, which is as follows:

Repeat until convergence {

(E-step) For each i , set

$$Q_i(z^{(i)}) := p(z^{(i)}|x^{(i)}; \theta).$$

(M-step) Set

$$\theta := \arg \max_{\theta} \sum_i \sum_{z^{(i)}} Q_i(z^{(i)}) \log \frac{p(x^{(i)}, z^{(i)}; \theta)}{Q_i(z^{(i)})}.$$

}

How do we know if this algorithm will converge? Well, suppose $\theta^{(t)}$ and $\theta^{(t+1)}$ are the parameters from two successive iterations of EM. We will now prove that $\ell(\theta^{(t)}) \leq \ell(\theta^{(t+1)})$, which shows EM always monotonically improves the log-likelihood. The key to showing this result lies in our choice of the Q_i 's. Specifically, on the iteration of EM in which the parameters had started out as $\theta^{(t)}$, we would have chosen $Q_i^{(t)}(z^{(i)}) := p(z^{(i)}|x^{(i)}; \theta^{(t)})$. We saw earlier that this choice ensures that Jensen's inequality, as applied to get Equation (3), holds with equality, and hence

$$\ell(\theta^{(t)}) = \sum_i \sum_{z^{(i)}} Q_i^{(t)}(z^{(i)}) \log \frac{p(x^{(i)}, z^{(i)}; \theta^{(t)})}{Q_i^{(t)}(z^{(i)})}.$$

The parameters $\theta^{(t+1)}$ are then obtained by maximizing the right hand side of the equation above. Thus,

$$\ell(\theta^{(t+1)}) \geq \sum_i \sum_{z^{(i)}} Q_i^{(t)}(z^{(i)}) \log \frac{p(x^{(i)}, z^{(i)}; \theta^{(t+1)})}{Q_i^{(t)}(z^{(i)})} \quad (4)$$

$$\geq \sum_i \sum_{z^{(i)}} Q_i^{(t)}(z^{(i)}) \log \frac{p(x^{(i)}, z^{(i)}; \theta^{(t)})}{Q_i^{(t)}(z^{(i)})} \quad (5)$$

$$= \ell(\theta^{(t)}) \quad (6)$$

This first inequality comes from the fact that

$$\ell(\theta) \geq \sum_i \sum_{z^{(i)}} Q_i(z^{(i)}) \log \frac{p(x^{(i)}, z^{(i)}; \theta)}{Q_i(z^{(i)})}$$

holds for any values of Q_i and θ , and in particular holds for $Q_i = Q_i^{(t)}$, $\theta = \theta^{(t+1)}$. To get Equation (5), we used the fact that $\theta^{(t+1)}$ is chosen explicitly to be

$$\arg \max_{\theta} \sum_i \sum_{z^{(i)}} Q_i(z^{(i)}) \log \frac{p(x^{(i)}, z^{(i)}; \theta)}{Q_i(z^{(i)})},$$

and thus this formula evaluated at $\theta^{(t+1)}$ must be equal to or larger than the same formula evaluated at $\theta^{(t)}$. Finally, the step used to get (6) was shown earlier, and follows from $Q_i^{(t)}$ having been chosen to make Jensen's inequality hold with equality at $\theta^{(t)}$.

Hence, EM causes the likelihood to converge monotonically. In our description of the EM algorithm, we said we'd run it until convergence. Given the result that we just showed, one reasonable convergence test would be to check if the increase in $\ell(\theta)$ between successive iterations is smaller than some tolerance parameter, and to declare convergence if EM is improving $\ell(\theta)$ too slowly.

Remark. If we define

$$J(Q, \theta) = \sum_i \sum_{z^{(i)}} Q_i(z^{(i)}) \log \frac{p(x^{(i)}, z^{(i)}; \theta)}{Q_i(z^{(i)})},$$

then we know $\ell(\theta) \geq J(Q, \theta)$ from our previous derivation. The EM can also be viewed a coordinate ascent on J , in which the E-step maximizes it with respect to Q (check this yourself), and the M-step maximizes it with respect to θ .

3 Mixture of Gaussians revisited

Armed with our general definition of the EM algorithm, let's go back to our old example of fitting the parameters ϕ , μ and Σ in a mixture of Gaussians. For the sake of brevity, we carry out the derivations for the M-step updates only for ϕ and μ_j , and leave the updates for Σ_j as an exercise for the reader.

The E-step is easy. Following our algorithm derivation above, we simply calculate

$$w_j^{(i)} = Q_i(z^{(i)} = j) = P(z^{(i)} = j | x^{(i)}; \phi, \mu, \Sigma).$$

Here, " $Q_i(z^{(i)} = j)$ " denotes the probability of $z^{(i)}$ taking the value j under the distribution Q_i .

Next, in the M-step, we need to maximize, with respect to our parameters