

# EXPECTATION MAXIMIZATION: A GENTLE INTRODUCTION

MORITZ BLUME

## 1. INTRODUCTION

This tutorial was basically written for students/researchers who want to get into first touch with the EXPECTATION MAXIMIZATION (EM) Algorithm. The main motivation for writing this tutorial was the fact that I did not find any text that fitted my needs. I started with the great book “Artificial Intelligence: A modern Approach” of Russel and Norvig [6], which provides lots of intuition, but I was disappointed that the general explication of the EM algorithm did not directly fit the shown examples. There is a very good tutorial written by Jeff Bilmes [1]. However, I felt that some parts of the mathematical derivation could be shortened by another choice of the hidden variable (explained later on). Finally I had a look at the book “The EM Algorithm and Extensions” by McLachlan and Krishnan [4]. I can definitely recommend this book to anybody involved with the EM algorithm - though, for a beginner, some steps could be more elaborate.

So, I ended up by writing my own tutorial. Any feedback is greatly appreciated (error corrections, spelling/grammar mistakes, etc.).

## 2. SHORT REVIEW: MAXIMUM LIKELIHOOD

The goal of Maximum Likelihood (ML) estimation is to find parameters that maximize the probability of having received certain measurements of a random variable distributed by some probability density function (p.d.f.). For example, we look at a random variable  $Y$  and a measurement vector  $\mathbf{y} = (y_1, \dots, y_N)^T$ . The probability of receiving some measurement  $y_i$  is given by the p.d.f.

$$(1) \quad p(y_i|\Theta) ,$$

where the p.d.f. is governed by the parameter  $\Theta$ . The probability of having received the whole series of measurements is then

$$(2) \quad p(\mathbf{y}|\Theta) = \prod_{i=1}^N p(y_i|\Theta) ,$$

if the measurements are independent. The likelihood function is defined as a function of  $\Theta$ :

$$(3) \quad L(\Theta) = p(\mathbf{y}|\Theta) .$$

The ML estimate of  $\Theta$  is found by maximizing  $L$ . Often, it is easier to maximize the log-likelihood

$$(4) \quad \log L(\boldsymbol{\Theta}) = \log p(\mathbf{y}|\boldsymbol{\Theta})$$

$$(5) \quad = \log \prod_{i=1}^N p(y_i|\boldsymbol{\Theta})$$

$$(6) \quad = \sum_{i=1}^N \log p(y_i|\boldsymbol{\Theta}) \quad .$$

Since the logarithm is a strictly increasing function, the maximum of  $L$  and  $\log(L)$  is the same.

It is important to note that we do not include any *a priori* knowledge of the parameter by calculating the ML estimate. Instead, we assume that each choice of the parameter vector is equally likely and therefore that the p.d.f. of the parameter vector is a uniform distribution. If such a prior p.d.f. for the parameter vector is available then methods of Bayesian parameter estimation should be preferred. In this tutorial we will only look at ML estimates.

In some cases a closed form can be derived by just setting the derivative with respect to  $\boldsymbol{\Theta}$  to zero. However, things are not always as easy, as the following example will show.

**Example: mixture distributions.** We assume that a data point  $y_j$  is produced as follows: first, one of  $I$  distributions/components which produces the measurement is chosen. Then, according to the parameters of this distribution, the actual measurement is produced. There are  $I$  components. The probability density function (p.d.f.) of the mixture distribution is

$$(7) \quad p(y_j|\boldsymbol{\Theta}) = \sum_{i=1}^I \alpha_i p(y_j|C=i, \boldsymbol{\beta}_i) \quad ,$$

where  $\boldsymbol{\Theta}$  is composed as  $\boldsymbol{\Theta} = (\alpha_1^\top, \dots, \alpha_I^\top, \beta_1^\top, \dots, \beta_M^\top)^\top$ . The parameters  $\alpha_i$  define the weight of component  $i$  and must sum up to 1, and the parameter vectors  $\boldsymbol{\beta}_i$  are associated to the p.d.f. of component  $i$ .

Based on some measurements  $\mathbf{y} = (y_1, \dots, y_J)^\top$  of the underlying random variable  $Y$ , the goal of ML (Maximum Likelihood) estimation is to find the parameters  $\boldsymbol{\Theta}$  that maximize the probability of having received these measurements. This involves maximization of the log-likelihood for  $\boldsymbol{\Theta}$ :

$$(8) \quad \log L(\boldsymbol{\Theta}) = \log p(\mathbf{y}|\boldsymbol{\Theta})$$

$$(9) \quad = \log \{p(y_1|\boldsymbol{\Theta}) \cdot \dots \cdot p(y_J|\boldsymbol{\Theta})\}$$

$$(10) \quad = \sum_{j=1}^J \log p(y_j|\boldsymbol{\Theta})$$

$$(11) \quad = \sum_{j=1}^J \log \left\{ \sum_{i=1}^I \alpha_i p(y_j|C=i, \boldsymbol{\beta}_i) \right\} \quad .$$

Since the logarithm is outside the sum, optimization of this term is not an easy task. There does not exist a closed form for the optimal value of  $\boldsymbol{\Theta}$ .

## 3. THE EM ALGORITHM

The goal of the EM algorithm is to facilitate maximum likelihood parameter estimation by introducing so-called *hidden* random variables which are not observed and therefore define the unobserved data. Let us denote this hidden random variable by some  $Z$ , and the unobserved data by a vector  $\mathbf{z}$ . Note that it does not play a role if this unobserved data is technically unobservable or just not observed due to other reasons. Often, it is just an artificial construct in order to enable an EM scheme.

We define the complete data  $\mathbf{x}$  to consist of the incomplete observed data  $\mathbf{y}$  and the unobserved data  $\mathbf{z}$ :

$$(12) \quad \mathbf{x} = (\mathbf{y}^\top, \mathbf{z}^\top)^\top .$$

Then, instead of looking at the log-likelihood for the *observed data*, we consider the log-likelihood function of the *complete data*:

$$(13) \quad \log L_C(\boldsymbol{\Theta}) = \log p(\mathbf{x}|\boldsymbol{\Theta}) .$$

Remember that in the ML case we just had a likelihood function and wanted to find the parameter that maximizes the probability of having obtained the observed data. The problem in the present case of the complete data likelihood function is that we did not observe the complete data ( $\mathbf{z}$  is missing!).

It may be surprising, but this problem of dealing with unobserved data in the end facilitates calculation of the ML parameter estimate.

The EM algorithm faces this problem by first finding an estimate for the likelihood function, and then maximizing the whole term. In order to find this estimate, it is important to note that each entry of  $\mathbf{z}$  is a realization of a hidden random variable. However, since these realizations do not exist in reality, we have to consider each entry of  $\mathbf{z}$  as a random variable itself. So, the whole likelihood function is nothing else than a function of a random variable and therefore by itself is a random variable (a quantity which depends on a random variable is a random variable). Its expectation given the observed data  $\mathbf{y}$  is:

$$(14) \quad \mathbb{E} \left[ \log L_C(\boldsymbol{\Theta}) | \mathbf{y}, \boldsymbol{\Theta}^{(i)} \right] .$$

The EM algorithm maximizes this expectation:

$$(15) \quad \boldsymbol{\Theta}^{(i+1)} = \underset{\boldsymbol{\Theta}}{\operatorname{argmax}} \mathbb{E} \left[ \log L_C(\boldsymbol{\Theta}) | \mathbf{y}, \boldsymbol{\Theta}^{(i)} \right] .$$

**Continuation of the example.** In the case of our mixture distributions, a good choice for  $\mathbf{z}$  is a matrix whose entry  $z_{ij}$  is equal to one if and only if component  $i$  produced measurement  $j$  - otherwise it is zero. Note that each column contains exactly one entry equal to one.

By this choice, we introduce additional information into the process:  $\mathbf{z}$  tells us the probability density function that underlies a certain measurement. Nobody knows if we can really measure this hidden variable, and it is not important. We just assume that we have it, do some calculus, and see what happens...

Accordingly, the log-likelihood function of the complete data is:

$$(16) \quad \log L_C(\Theta) = \log p(\mathbf{x}|\Theta)$$

$$(17) \quad = \log \prod_{j=1}^J \sum_{i=1}^I z_{ij} \alpha_i p(y_j | C = i, \beta_i)$$

$$(18) \quad = \sum_{j=1}^J \log \sum_{i=1}^I z_{ij} \alpha_i p(y_j | C = i, \beta_i) .$$

This can now be rewritten as

$$(19) \quad \log L_C(\Theta) = \sum_{j=1}^J \sum_{i=1}^I z_{ij} \log \alpha_i p(y_j | C = i, \beta_i) ,$$

since  $z_{ij}$  is zero for all but one term in the inner sum!

The choice of this hidden variable might look like magic. In fact this is the most difficult part for defining an EM scheme. The intuition behind the choice of  $z_{ij}$  might be that if the information about the assigned components were available, then it would be easily possible to estimate the parameters - one just has to do a simple ML estimation for each component.

One also should note that even if he decides to incorporate the right information (here: the assignment of measurements to components), the simplicity of the further steps strongly depends on how this is modeled. For example, instead of indicator variables  $z_{ij}$  one could define random variables  $z_j$  that represent the value of the respective component for measurement  $j$ . However, this heavily complicates further calculations (as can be seen in [1]).

Compared to the ML approach which involves just maximizing a log likelihood, the EM algorithm makes one step in between: the calculation of the expectation. This is denoted as the expectation step. Please note that in order to calculate the expectation, an initial parameter estimate  $\Theta^{(i)}$  is needed. But still, the whole term will depend on a parameter  $\Theta$ . In the maximization step, a new update  $\Theta^{(i+1)}$  for  $\Theta$  is found in such manner, that it maximizes the whole term. This whole procedure is repeated until some stopping criteria, e. g. until the difference of change between the parameter updates becomes very small.

**EM for mixture distributions.** Now, back to our example. The parameter update is found by solving

$$(20) \quad \Theta^{(i+1)} = \operatorname{argmax}_{\Theta} \mathbb{E} \left[ \log L_C(\Theta) | \mathbf{y}, \Theta^{(i)} \right]$$

$$(21) \quad = \operatorname{argmax}_{\Theta} \mathbb{E} \left[ \sum_{j=1}^J \sum_{i=1}^I z_{ij} \log \{ \alpha_i p(y_j) | C = i, \beta_i \} \middle| \mathbf{y}, \Theta^{(i)} \right] .$$

**3.1. Expectation Step.** The expectation step consists of calculating the expected value of the complete data likelihood function. Due to linearity of expectation we have

$$(22) \quad \boldsymbol{\Theta}^{(i+1)} = \underset{\boldsymbol{\Theta}}{\operatorname{argmax}} \sum_{j=1}^J \sum_{i=1}^I \mathbb{E} \left[ z_{ij} | \mathbf{y}, \boldsymbol{\Theta}^{(i)} \right] \log \alpha_i p(y_j | C = i, \beta_i) .$$

In order to complete this step, we have to see how  $\mathbb{E} [z_{ij} | \mathbf{y}, \boldsymbol{\Theta}^{(i)}]$  looks like:

$$(23) \quad \mathbb{E} \left[ z_{ij} | \mathbf{y}, \boldsymbol{\Theta}^{(i)} \right] = 0 \cdot p(z_{ij} = 0 | \boldsymbol{\Theta}^{(i)}) + 1 \cdot p(z_{ij} = 1 | \boldsymbol{\Theta}^{(i)}) = p(z_{ij} = 1 | \mathbf{y}, \boldsymbol{\Theta}^{(i)}) .$$

This is the probability that class  $i$  produced measurement  $j$ . In other words: The probability that class  $i$  was active while measurement  $j$  was produced:

$$(24) \quad p(z_{ij} = 1) = p(C = i | x_j, \boldsymbol{\Theta}^{(i)}) .$$

Application of Bayes theorem leads to

$$(25) \quad p(C = i | x_j, \boldsymbol{\Theta}^{(i)}) = \frac{p(x_j | C = i, \boldsymbol{\Theta}^{(i)}) p(C = i | \boldsymbol{\Theta}^{(i)})}{p(x_j | \boldsymbol{\Theta}^{(i)})} .$$

All the terms on the right side can be easily calculated.  $p(x_j | C = i, \boldsymbol{\Theta}^{(i)})$  just comes from the  $i$ -th components density function.  $p(C = i | \boldsymbol{\Theta}^{(i)})$  is the probability of choosing the  $i$ th component, which is nothing else than the parameter  $\alpha_i$  (which are already given by  $\boldsymbol{\Theta}^{(i)}$ ). Finally,  $p(x_j | \boldsymbol{\Theta}^{(i)})$  is the whole distributions probability density function, which, based on  $\boldsymbol{\Theta}^{(i)}$ , is easy to calculate.

**3.2. Maximization Step.** Remember that our parameter  $\boldsymbol{\Theta}$  is composed of the  $\alpha_k$  and  $\beta_k$ . We will perform the optimization by setting the respective partial derivatives to zero.

**3.2.1. Optimization with respect to  $\alpha_k$ .** Optimization with respect to  $\alpha_k$  in fact is a constrained optimization, since we suppose that the  $\alpha_i$  sum up to one. So, the subproblem is: maximize (22) with respect to some  $\alpha_k$  subject to  $\sum_{i=1}^I \alpha_i = 1$ . We do this by using the method of Lagrange multipliers (a very good tutorial explaining constrained optimization is [3]):

$$(26) \quad \frac{\partial}{\partial \alpha_k} \sum_{j=1}^J \sum_{i=1}^I \mathbb{E} \left[ z_{ij} | \mathbf{y}, \boldsymbol{\Theta}^{(i)} \right] \log \{ \alpha_i p(y_j | C = i, \beta_i) \} + \lambda \left( \sum_{i=1}^I \alpha_i - 1 \right) = 0$$

$$(27) \quad \Leftrightarrow \sum_{j=1}^J \mathbb{E} \left[ z_{kj} | \mathbf{y}, \boldsymbol{\Theta}^{(i)} \right] + \alpha_k \lambda = 0$$

$\lambda$  is found by getting rid of  $\alpha_k$ , that is, by inserting the constraint. Therefore, we sum the whole equation up over  $i$  running from 1 to  $I$ :

$$(28) \quad \sum_{k=1}^M \sum_{j=1}^J \mathbb{E} [z_{ij} | \mathbf{y}, \boldsymbol{\Theta}^{(i)}] + \lambda \sum_{k=1}^I \alpha_k = 0$$

$$(29) \quad \Leftrightarrow \lambda = - \sum_{k=1}^I \sum_{j=1}^J \mathbb{E} [z_{ij} | \mathbf{y}, \boldsymbol{\Theta}^{(i)}]$$

$$(30) \quad \Leftrightarrow \lambda = -J \ .$$

Now, (27) looks like

$$(31) \quad \alpha_k = \frac{1}{J} \sum_{j=1}^J \mathbb{E} [z_{ij} | \mathbf{y}, \boldsymbol{\Theta}^{(i)}] \ ,$$

and  $\alpha_k$  is our new estimate. As can be seen, it depends on the expectation of  $z_{ij}$ , which was already found in the expectation step!

**3.2.2. Maximization with respect to  $\beta_k$ .** Maximization with respect to  $\beta_k$  depends on the components p.d.f. In general, it is not possible to derive a closed form. For some, a closed form is sought by setting the first derivative to zero.

In this tutorial, we will take a closer look at Gaussian mixture distributions. In this case,  $\beta_k$  consists of the  $k$ -th mean value and the  $k$ -th variance:

$$(32) \quad \boldsymbol{\beta}_k = (\mu_k, \sigma_k)^\top$$

Their maxima are found by setting the derivative with respect to  $\mu_k$  and  $\sigma_k$  equal to zero:

$$(33) \quad \frac{\partial}{\partial \beta_k} \sum_{j=1}^J \sum_{i=1}^I \mathbb{E} [z_{ij} | \mathbf{y}, \boldsymbol{\Theta}^{(i)}] \log \{ \alpha_i p(y_j | C = i, \beta_i) \} \ .$$

$$(34) \quad = \frac{\partial}{\partial \beta_k} \sum_{j=1}^J \sum_{i=1}^I \mathbb{E} [z_{ij} | \mathbf{y}, \boldsymbol{\Theta}^{(i)}] \log \alpha_i + \mathbb{E} [z_{ij} | \mathbf{y}, \boldsymbol{\Theta}^{(i)}] \log p(y_j | C = i, \beta_i) \ .$$

$$(35) \quad = \sum_{j=1}^J \mathbb{E} [z_{kj} | \mathbf{y}, \boldsymbol{\Theta}^{(i)}] \frac{\partial}{\partial \beta_k} \log p(y_j | C = k, \beta_k) \ .$$

By inserting the definition of a Gaussian distribution we arrive at

$$(36) \quad \sum_{j=1}^J \mathbb{E} [z_{kj} | \mathbf{y}, \boldsymbol{\Theta}^{(i)}] \cdot \frac{\partial}{\partial \beta_k} \log \left[ \frac{1}{\sigma_k \sqrt{2\pi}} \exp \left( -\frac{1}{2} \frac{(y_j - \mu_k)^2}{\sigma_k^2} \right) \right]$$

$$(37) \quad = \sum_{j=1}^J \mathbb{E} [z_{kj} | \mathbf{y}, \boldsymbol{\Theta}^{(i)}] \cdot \frac{\partial}{\partial \beta_k} \left[ -\log(\sigma_k) - \log \sqrt{2\pi} - \frac{1}{2} \frac{(y_j - \mu_k)^2}{\sigma_k^2} \right]$$

**Maximization with respect to  $\mu_k$ .** For any maxima, the following condition has to be fulfilled:

$$(38) \quad \sum_{j=1}^J \mathbb{E} \left[ z_{kj} | \mathbf{y}, \boldsymbol{\Theta}^{(i)} \right] \cdot \frac{\partial}{\partial \mu_k} \left[ -\log(\sigma_k) - \log \sqrt{2\pi} - \frac{1}{2} \frac{(y_j - \mu_k)^2}{\sigma_k^2} \right] = 0 \quad .$$

The derivation for  $\mu_k$  is straightforward:

$$(39) \quad \sum_{j=1}^J \mathbb{E} \left[ z_{kj} | \mathbf{y}, \boldsymbol{\Theta}^{(i)} \right] \cdot \frac{\sigma_k^2 (\mu_k - y_j)}{\sigma_k^4} = 0$$

$$(40) \quad \Leftrightarrow \sum_{j=1}^J \mathbb{E} \left[ z_{kj} | \mathbf{y}, \boldsymbol{\Theta}^{(i)} \right] \cdot (\mu_k - y_j) = 0$$

$$(41) \quad \Leftrightarrow \mu_k = \frac{\sum_{j=1}^J \mathbb{E} \left[ z_{kj} | \mathbf{y}, \boldsymbol{\Theta}^{(i)} \right] y_j}{\sum_{j=1}^J \mathbb{E} \left[ z_{kj} | \mathbf{y}, \boldsymbol{\Theta}^{(i)} \right]}$$

**Maximization with respect to  $\sigma_k$ .** For any maxima, the following condition has to be fulfilled:

$$(42) \quad \sum_{j=1}^J \mathbb{E} \left[ z_{kj} | \mathbf{y}, \boldsymbol{\Theta}^{(i)} \right] \cdot \frac{\partial}{\partial \sigma_k} \left[ -\log(\sigma_k) - \log \sqrt{2\pi} - \frac{1}{2} \frac{(y_j - \mu_k)^2}{\sigma_k^2} \right] = 0 \quad .$$

Again, the derivation is straight forward:

$$(43) \quad \sum_{j=1}^J \mathbb{E} \left[ z_{kj} | \mathbf{y}, \boldsymbol{\Theta}^{(i)} \right] \cdot \left[ -\frac{1}{\sigma_k} - \frac{-(y_j - \mu_k)^2}{\sigma_k^3} \right] = 0$$

$$(44) \quad \Leftrightarrow \sum_{j=1}^J \mathbb{E} \left[ z_{kj} | \mathbf{y}, \boldsymbol{\Theta}^{(i)} \right] \cdot \left[ -\sigma_k^2 + \frac{1}{2} (y_j - \mu_k)^2 \right] = 0$$

$$(45) \quad \Leftrightarrow \sigma_k^2 = \frac{1}{2} \frac{\sum_{j=1}^J \mathbb{E} \left[ z_{kj} | \mathbf{y}, \boldsymbol{\Theta}^{(i)} \right] (y_j - \mu_k)^2}{\sum_{j=1}^J \mathbb{E} \left[ z_{kj} | \mathbf{y}, \boldsymbol{\Theta}^{(i)} \right]}$$

#### 4. CONCLUSION

As can be seen, the EM Algorithm depends very much on the application. In fact, it has been criticized by some reviewers of the original paper by Dempster [2] that the term ‘‘Algorithm’’ is not adequate for this iteration scheme, since its formulation is too general. The most difficult task in the identification of the EM Algorithm for a new application is to find the hidden variable.

The reader is encouraged to look at different applications of the EM Algorithm in order to develop a better sense for where and how it can be applied.

#### REFERENCES

1. J. Bilmes, *A gentle tutorial on the em algorithm and its application to parameter estimation for gaussian mixture and hidden markov models*, 1998.

2. A. P. Dempster, N. M. Laird, and D. B. Rubin, *Maximum likelihood from incomplete data via the em algorithm*, Journal of the Royal Statistical Society. Series B (Methodological) **39** (1977), no. 1, 1–38.
3. Dan Klein, *Lagrange multipliers without permanent scarring*, University of California at Berkeley, Computer Science Division.
4. Geoffrey J. McLachlan and Thriyambakam Krishnan, *The EM algorithm and extensions*, Wiley, 1996.
5. Thomas P. Minka, *Old and new matrix algebra useful for statistics*, 2000.
6. Stuart J. Russell and Peter Norvig, *Artificial intelligence: A modern approach (2nd edition)*, Prentice Hall, 2002.

MORITZ BLUME, TECHNISCHE UNIVERSITÄT MÜNCHEN, INSTITUT FÜR INFORMATIK / I16, BOLTZMANNSTRASSE 3, 85748 GARCHING B. MÜNCHEN, GERMANY

*E-mail address:* `blume@cs.tum.edu`

*URL:* <http://www.navab.cs.tum.edu/Main/MoritzBlume>