

Generalized Expectation Maximization¹

1 Maximum Likelihood Estimation (MLE)

First let's recall the definition of the maximum-likelihood estimation problem. We have a density function $f(x | \Theta)$, that is governed by the set of parameters Θ (e.g., f might be a Gaussian and Θ could be the means (vector) and covariance (matrix)). We also have a data set of size N , supposedly drawn from this distribution with density f , i.e., $\mathcal{X} = \{x_1, \dots, x_N\}$. That is, we assume that these data vectors are independent and identically distributed (IID) with distribution f . Therefore, the resulting joint density for the samples is:

$$p(\mathbf{X} | \Theta) = \prod_{i=1}^N p(x_i | \Theta) = L(\Theta | \mathbf{X})$$

$L(\Theta | \mathbf{X})$ is called the [likelihood function](#) of the parameters (Θ) given the data, or just the likelihood. The likelihood is thought of as a function of the parameters (Θ) where the data \mathcal{X} is fixed (observed). In the [maximum likelihood problem](#), our goal is to find a parameter vector Θ that maximizes $L(\Theta | \mathbf{X})$. In other words, we look for Θ^* , where

$$\Theta^* = \underset{\Theta}{\text{ArgMax}} L(\Theta | \mathbf{X})$$

Oftentimes we choose to maximize $\text{Log}(L(\Theta | \mathbf{X}))$ instead because it is analytically easier or computationally appealing.

Depending on the form of $f(x | \Theta)$ this problem can be easy or hard. For example, if $f(x | \Theta)$ is simply a single Gaussian distribution where the parameter vector $\Theta = (\mu, \sigma^2)$, then we can solve the maximum likelihood problem of determining estimates (MLE) of μ & σ^2 by setting the partial derivatives of $\text{Log}(L(\Theta | \mathbf{X}))$ to zero (in fact, this is how the familiar formulas for the population mean and variance are obtained). For many problems, however, it is not possible to find such analytical expressions, and we must resort to more elaborate techniques (e.g., EM technique).

¹ This set of class notes is based on previous work by Bilmer, Bishop, Dempster, Ghahramani, Jordan, Rabiner, Redner, Wu and Xu see references at the end of the document.

Example 1: Suppose $\{X_1, \dots, X_n\}$ IID $N(\mu, \sigma^2)$, where μ is unknown. We estimate μ by:
 $MLE(\mu) = \hat{\mu} = \underset{\mu}{\text{ArgMax}} L(\mu | \{X_1, \dots, X_n\})$.

$$MLE(\mu) = \text{Log} \left(\prod_{i=1}^n \frac{e^{-\frac{(x_i - \mu)^2}{2\sigma^2}}}{\sqrt{2\pi\sigma^2}} \right) = L(\mu)$$

$$0 = L'(\hat{\mu}) = \frac{1}{(2\pi\sigma^2)^{n/2}} \left(e^{-\sum_{i=1}^n \frac{(x_i - \hat{\mu})^2}{2\sigma^2}} \right) \frac{\sum_{i=1}^n 2(x_i - \hat{\mu})}{2\sigma^2}$$

$$\Leftrightarrow 0 = 2\sum_{i=1}^n (x_i - \hat{\mu}) \Leftrightarrow \hat{\mu} = \frac{\sum_{i=1}^n x_i}{n}.$$

Similarly can show that : $MLE(\sigma) = \hat{\sigma} = \frac{\sum_{i=1}^n (x_i - \mu)^2}{n - 1}$.

Example 2: Suppose $\{X_1, \dots, X_n\}$ IID Poisson(λ), where (the population mean, and standard deviation) λ is unknown. Estimate λ by:

$$MLE(\lambda) = \hat{\lambda} = \underset{\lambda}{\text{ArgMax}} L(\lambda | \{X_1, \dots, X_n\})$$

$$MLE(\lambda) = \text{Log} \left(\prod_{i=1}^n \frac{e^{-\lambda} \lambda^{x_i}}{(x_i)!} \right) = L(\lambda)$$

$$0 = L'(\hat{\lambda}) = \frac{\partial}{\partial \lambda} \text{Log} \left(\frac{e^{-n\lambda} \lambda^{\sum_{i=1}^n x_i}}{\prod_{i=1}^n (x_i)!} \right) =$$

$$= \frac{\partial}{\partial \lambda} \left(-n\lambda + \text{Log}(\lambda) \sum_{i=1}^n x_i \right) = -n + \frac{1}{\lambda} \sum_{i=1}^n x_i \Leftrightarrow \hat{\lambda} = \frac{\sum_{i=1}^n x_i}{n}.$$

2 General Expectation Maximization (GEM) Algorithms

The EM algorithm is one such technique, which allows estimating parameter vectors in the cases when such analytic solutions to the likelihood minimization are difficult or impossible. The EM algorithm [see references at the end] is a general method of finding the maximum-likelihood

estimates of the parameters of an underlying distribution from a given data set when the data is incomplete or has missing values. There are two main applications of the EM algorithm.

Application 1: The first application of EM algorithm occurs when the data indeed has missing values, due to problems with or limitations of the observation process.

Application 2: The second occurs when optimizing the likelihood function is analytically intractable, however, we still need to assume the likelihood function can be simplified by assuming the existence of and values for additional but *missing* (or *hidden*) parameters. This second application is more common in the computational pattern recognition community.

As before, we assume that data \mathbf{X} is observed and is generated by some distribution. We call \mathbf{X} the *incomplete data*. We assume that a complete data set exists $\mathbf{Z}=(\mathbf{X};\mathbf{Y})$ and also assume (or specify) a joint density function:

$$p(\mathbf{Z}|\Theta) = p(\mathbf{X},\mathbf{Y}|\Theta) = p(\mathbf{Y}|\mathbf{X},\Theta) p(\mathbf{X}|\Theta)$$

Recall that for each probability measure, P , $P(A,B)=P(A|B)P(B)$ and hence, in terms of the conditional probability, $P(A,B|C) = P(A|B,C)P(B|C)$. This joint density often times arises from the marginal density function $p(\mathbf{X}|\Theta)$ and the assumption of hidden variables and parameter value guesses (e.g., mixture-densities, this example is coming up; and Baum-Welch algorithm). In other cases (e.g., missing data values in samples of a distribution), we must assume a joint relationship between the missing and observed values.

With this new density function, we can define a new likelihood function,

$$L(\Theta|\mathbf{Z}) = L(\Theta|\mathbf{X},\mathbf{Y}) = p(\mathbf{X},\mathbf{Y}|\Theta),$$

called the [complete-data likelihood](#). Note that this function is in fact a random variable since the missing information \mathbf{Y} is unknown, i.e., random, and presumably governed by an underlying distribution. That is, we can think of $L(\Theta|\mathbf{X},\mathbf{Y}) = h_{\mathbf{X},\Theta}(\mathbf{Y})$ for some function $h_{\mathbf{X},\Theta}(\mathbf{Y})$, where \mathbf{X} and Θ are constant and \mathbf{Y} is a random variable. The original likelihood $L(\Theta|\mathbf{X})$ is referred to as the [incomplete-data likelihood function](#).

The Expectation Maximization algorithm then proceeds in two steps – expectation followed by its maximization.

Step 1 of EM (Expectation): The EM algorithm needs to first find the expected value of the complete-data log-likelihood $\text{Log}(p(\mathbf{X},\mathbf{Y}|\Theta))$ with respect to the unknown data \mathbf{Y} given the observed data \mathbf{X} and the current parameter estimates Θ^i , the *E-Step*. Below we define the [expectation of \$\Theta^{\(i-1\)}\$](#) , the second argument represents the parameters that we use to evaluate the expectation. The first argument Θ simply indicates the parameter (vector) that ultimately will be optimized in an attempt to maximize the likelihood:

$$Q(\Theta, \Theta^{(i-1)}) = E[\text{Log } p(\mathbf{X},\mathbf{Y}|\Theta) | \mathbf{X}, \Theta^{(i-1)}]. \quad (1)$$

Where $\Theta^{(i-1)}$ are the current parameters estimates that we used to evaluate the expectation and Θ are the new parameters that we optimize to increase (maximize) Q . Note that the expression above is a *conditional expectation* w.r.t. $(\mathbf{X} \ \& \ \Theta^{(i-1)})$, i.e.,

$$E[h(Y) | X=x] := E[h(Y) | X = x] := \int_{\mathcal{Y}} h(y) \times f_{Y|X}(y | x) dy.$$

The key thing to understand is that \mathbf{X} and $\Theta^{(i-1)}$ are *given* constants, Θ is a random variable that we wish to adjust/estimate, and \mathbf{Y} is a random variable governed by the distribution $f(\mathbf{Y} | \mathbf{X}, \Theta^{(i-1)})$. The right side of equation (1) can therefore be expressed as:

$$E[\text{Log } \rho(\mathbf{X}, \mathbf{Y} | \Theta) | \mathbf{X}, \Theta^{(i-1)}] = \int_{\mathcal{Y} \in \mathcal{Y}} \text{Log } \rho(\mathbf{X}, y | \Theta) f(y | \mathbf{X}, \Theta^{(i-1)}) dy \quad (2)$$

The $f(\mathbf{Y} | \mathbf{X}, \Theta^{(i-1)})$ is the **marginal** distribution of the unobserved data (\mathbf{Y}) and is dependent on: the observed data \mathbf{X} , on the current parameters $\Theta^{(i-1)}$, and \mathcal{Y} is the space of values \mathbf{y} can take on. In the best-case situations, this marginal distribution is a simple analytical expression of the assumed parameters $\Theta^{(i-1)}$, and perhaps the observed data (\mathbf{X}). In the worst-case scenario, this density might be very hard to obtain. In fact, sometimes the actually used density is:

$$f(\mathbf{y}, \mathbf{X} | \Theta^{(i-1)}) = f(\mathbf{y} | \mathbf{X}, \Theta^{(i-1)}) f(\mathbf{X} | \Theta^{(i-1)}),$$

but this doesn't effect subsequent steps since the extra factor, $f(\mathbf{X} | \Theta^{(i-1)})$ is not dependent on Θ .

As an analogy, suppose we have a function $h(\cdot; \cdot)$ of two variables. Consider $h(\theta; \mathbf{Y})$ where θ is a constant and \mathbf{Y} is a random variable governed by some distribution $f_{\mathbf{Y}}(y)$. Then

$$q(\theta) = E_{\mathbf{Y}} [h(\theta; \mathbf{Y})] = \int_{\mathcal{Y}} h(\theta; \mathbf{Y}) f_{\mathbf{Y}}(y) dy$$

is now a deterministic function that could be maximized if desired, w.r.t. θ . The evaluation of this expectation is called the **E-step** of the algorithm. Notice the meaning of the two arguments in the function $Q(\Theta, \Theta^{(i-1)})$. The first argument Θ corresponds to the parameters that ultimately will be optimized in an attempt to maximize the likelihood. The second argument, $\Theta^{(i-1)}$, corresponds to the parameters that we use to evaluate the expectation at each iteration ($\underline{\mathbf{E}} \rightarrow \underline{\mathbf{M}} \rightarrow \underline{\mathbf{E}} \rightarrow \underline{\mathbf{M}} \rightarrow \dots$).

Step 2 of EM (Maximization): The second step (the **M-step**) of the EM algorithm is to maximize the expectation we computed in the first step. That is, we iteratively compute:

$$\Theta^{(i)} = \underset{\Theta}{\text{ArgMax}} \left(Q \left(\Theta, \Theta^{(i-1)} \right) \right)$$

That is we *maximize the expectation of the log-likelihood function*. These two steps are repeated as necessary. Each iteration is guaranteed to increase the log-likelihood and the **algorithm is guaranteed to converge** to a local maximum of the likelihood function. There are many theoretical and empirical rate-of-convergence papers (see references below).

A modified form of the M-step is to, instead of maximizing the rather difficult function $Q(\Theta, \Theta^{(i-1)})$, we **find some $\Theta^{(i)}$** such that $Q(\Theta^{(i)}, \Theta^{(i-1)}) > Q(\Theta, \Theta^{(i-1)})$. This form of the algorithm is called Generalized EM (GEM) and is **also guaranteed to converge**. This description of the GEM does not yield a direct computer implementation scheme (it's not constructive (the coding algorithm is not explicit). This is the way, however, that the algorithm is presented in its most general form. The details of the steps required to compute the given quantities are very dependent on the particular application, so they are not discussed when the algorithm is presented in this abstract form, but rather detailed for each specific application.

3 EM Application: Finding Maximum Likelihood Mixture Densities Parameters via EM

The mixture-density parameter estimation problem is probably one of the most widely used applications of the EM algorithm in the computational pattern recognition community. In this case, we assume the following (mixture-density) model:

$$p(x | \Theta) = \sum_{i=1}^M a_i p_i(x | \Theta) \quad M = \# \text{Mixture Distributions} \quad (3)$$

where the parameter vector is $\Theta = (\alpha_1, \dots, \alpha_M; \theta_1, \dots, \theta_M)$, where the mixture-model weights satisfy $\sum_{i=1}^M \alpha_i = 1$ and each p_i is a density function parameterized, in general, by its own parameter vector θ_i . In other words, we assume we have M component densities mixed together with M mixing coefficients α_i .

The incomplete-data log-likelihood expression for this density from the data $\mathbf{X} = \{x_1, \dots, x_N\}$, $N = \# \text{ Observations}$, is given by:

$$\text{Log}(L(\Theta | \mathbf{X})) = \text{Log}\left[\prod_{i=1}^N p(x_i | \Theta)\right] = \sum_{i=1}^N \text{Log}\left[\sum_{j=1}^M a_j p_j(x | \theta_j)\right]$$

which is difficult to optimize because it contains the logarithm function of the sum [if the sum and the log were interchanged then optimizing the outside logarithm would have been equivalent to optimizing its argument – the sum – as the log function has always a positive derivative over its domain $(0; \infty)$]. If we consider X as incomplete, however, and assume the existence of unobserved data items $Y = \{y_i\}_{i=1}^N$, whose values inform us which component density “generated” each data item, the likelihood expression is significantly simplified. That is, if we assume that $y_i \in \{1, 2, \dots, M\}$ for each $1 \leq i \leq N$, and $y_i = k$ if the i^{th} sample, x_i , was generated by the k^{th} mixture component p_k . If we know the values of \mathbf{Y} , the likelihood becomes:

$$\text{Log}(L(\Theta | \mathbf{X}, \mathbf{Y})) = \text{Log}[P(\mathbf{X}, \mathbf{Y} | \Theta)] = \sum_{i=1}^N \text{Log}[P(x_i | y_i)P(y_i)] = \sum_{i=1}^N \text{Log}\left[\alpha_{y_i} p_{y_i}(x_i | \theta_{y_i})\right]$$

which, given a particular form of the component densities, can be optimized using a variety of techniques. The problem, of course, is that we do not know the values of \mathbf{Y} . If we assume \mathbf{Y} is a random vector, however, we can proceed. We first must derive an expression for the distribution of the unobserved (missing) data, \mathbf{Y} . Let's first guess at parameters for the mixture density, i.e., we guess that $\Theta^g = (\alpha_1^g, \dots, \alpha_M^g; \theta_1^g, \dots, \theta_M^g)$, are the appropriate parameters for the likelihood $L(\Theta^g | \mathbf{X}, \mathbf{Y})$. Given Θ^g , we can compute $p_j(x_i | \theta_j^g)$ for each i and j . In addition, the mixing parameters, α_i , can be thought of as *prior* probabilities of each mixture component, that is $\alpha_i = p(\text{component } i)$. Therefore, using Bayes's rule – $P(Y_i | X_i) = P(X_i | Y_i)P(Y_i) / P(X_i)$, we can compute:

$$p(y_i | x_i, \Theta^g) = \frac{\alpha_{y_i}^g p_{y_i} \left(\begin{matrix} x_i | \theta_{y_i}^g \\ y_i \end{matrix} \right)}{p(x_i | \Theta^g)} = \frac{\alpha_{y_i}^g p_{y_i} \left(\begin{matrix} x_i | \theta_{y_i}^g \\ y_i \end{matrix} \right)}{\sum_{k=1}^M \alpha_k^g p_k(x_i | \Theta_k^g)}$$

and therefore, $p(y | X, \Theta^g) = \prod_{i=1}^N p(y_i | x_i, \Theta^g)$, where $y = (y_1, \dots, y_N)$ is an instance of the unobserved data independently drawn. When we now look at equation (2), we see that in this case we have obtained the desired marginal density (of \mathbf{Y}), $f(y | \mathbf{X}, \Theta^{(g)})$, by assuming the existence of the hidden variables and making a guess at the initial parameters of their distribution. In this case, equation (1) takes the specific form:

$$\begin{aligned} Q(\Theta, \Theta^{(g)}) &= \sum_{y \in \mathbf{Y}} \log(L(\Theta | X, y)) p(y | X, \Theta^g) \\ &= \sum_{y \in \mathbf{Y}} \sum_{i=1}^N \text{Log}(\alpha_{y_i} p_{y_i}(x_i | \theta_{y_i})) \prod_{j=1}^N p(y_j | x_j, \Theta^g) \\ &= \sum_{y_1=1}^M \sum_{y_2=1}^M \dots \sum_{y_N=1}^M \sum_{i=1}^N \text{Log}(\alpha_{y_i} p_{y_i}(x_i | \theta_{y_i})) \prod_{j=1}^N p(y_j | x_j, \Theta^g) \\ &= \sum_{y_1=1}^M \sum_{y_2=1}^M \dots \sum_{y_N=1}^M \sum_{i=1}^N \sum_{l=1}^M \delta_{l, y_i} \text{Log}(\alpha_l p_l(x_i | \theta_l)) \prod_{j=1}^N p(y_j | x_j, \Theta^g) \\ &= \sum_{l=1}^M \sum_{i=1}^N \text{Log}(\alpha_l p_l(x_i | \theta_l)) \sum_{y_1=1}^M \sum_{y_2=1}^M \dots \sum_{y_N=1}^M \delta_{l, y_i} \prod_{j=1}^N p(y_j | x_j, \Theta^g) \quad (4) \end{aligned}$$

where $\delta_{l, y_i} = \begin{cases} 1, & \text{if } l = y_i \\ 0, & \text{otherwise} \end{cases}$. In this form, $Q(\Theta, \Theta^{(g)})$ appears computationally challenging,

however, it can be simplified, since for $l \in \{1, 2, \dots, M\}$,

$$\begin{aligned} \sum_{y_1=1}^M \sum_{y_2=1}^M \dots \sum_{y_N=1}^M \delta_{l, y_i} \prod_{j=1}^N p(y_j | x_j, \Theta^g) &= \\ &= \left(\sum_{y_1=1}^M \dots \sum_{y_{i-1}=1}^M \sum_{y_{i+1}=1}^M \dots \sum_{y_N=1}^M \prod_{j=1, j \neq i}^N p(y_j | x_j, \Theta^g) \right) p(l | x_i, \Theta^g) \\ &= \prod_{j=1, j \neq i}^N \left[\sum_{y_j=1}^M p(y_j | x_j, \Theta^g) \right] p(l | x_i, \Theta^g) = p(l | x_i, \Theta^g) \quad (5) \end{aligned}$$

This is because $\sum_{i=1}^M p(i | x_j, \Theta^g) = 1$. Using equation (5), we can write equation (4) as:

$$\begin{aligned}
 Q(\Theta, \Theta^{(g)}) &= \sum_{l=1}^M \sum_{i=1}^N \text{Log}(a_l p_l(x_i | \theta_l)) p(l | x_i, \Theta^g) \\
 &= \sum_{l=1}^M \sum_{i=1}^N \text{Log}(a_l) \text{Log}(p_l(x_i | \Theta^g)) + \sum_{l=1}^M \sum_{i=1}^N \text{Log}(p_l(x_i | \theta_l)) p(l | x_i, \Theta^g) \quad (6) \\
 &= \Phi(\{\alpha_l\}) + \Psi(\{\theta_l\}).
 \end{aligned}$$

Now, to maximize this expression, we can maximize independently Φ and Ψ , the terms containing α_l and θ_l , since they are not related. To find the expression for α_l , which maximizes $Q(\Theta, \Theta^{(g)})$, we introduce the Lagrange multiplier λ with the constraint that $g = \sum_l \alpha_l - 1 = 0$, and solve the following equation:

$$\frac{\partial}{\partial \alpha_l} \left[\sum_{l=1}^M \sum_{i=1}^N \text{Log}(\alpha_l) \text{Log}(p_l(x_i | \Theta^g)) + \lambda \left(\sum_{l=1}^M \alpha_l - 1 \right) \right] = 0$$

This yields: $\sum_{i=1}^N \frac{1}{\alpha_l} p_l(x_i | \Theta^g) + \lambda = 0$

Summing over l we get:

$$\sum_{l=1}^M \left[\sum_{i=1}^N \frac{1}{\alpha_l} p_l(x_i | \Theta^g) + \lambda \right] = \sum_{i=1}^N \sum_{l=1}^M \frac{1}{\alpha_l} p_l(x_i | \Theta^g) + M \lambda = 0$$

Therefore, $\lambda = -N$, which yields that $\alpha_l = \frac{1}{N} \sum_{i=1}^N p_l(x_i | \Theta^g)$. This is how we determine the mixture parameters $\{\alpha_l\}$, most of the time.

Now let us try to estimate the second part of the parameter vector $\Theta = (\alpha_1, \dots, \alpha_M; \theta_1, \dots, \theta_M)$, i.e., the distribution specific parameters $(\theta_1, \dots, \theta_M)$. Clearly, these need to be estimated in a case by case manner, as different distributions have different number and type of parameters. We consider again a couple of cases that illustrate the basic strategy for estimating $(\theta_1, \dots, \theta_M)$ using EM approaches. For some distributions, it will be possible to get analytic expressions for θ_l directly, as functions of all other variables.

Example 1: Suppose that the mixture model in equation (3) involves Poisson(λ_l) distributions, $1 \leq l \leq M$. Then the

$$\begin{aligned}
 Q(\Theta, \Theta^{(g)}) &= \sum_{l=1}^M \sum_{i=1}^N \text{Log}(p_l(x_i | \mu_l, \Sigma_l)) p(l | x_i, \Theta^g) \\
 &= \sum_{l=1}^M \sum_{i=1}^N \left[-\lambda_l + x_i \text{Log}(\lambda_l) - \text{Log}((x_i)!) \right] p(l | x_i, \Theta^g)
 \end{aligned}$$

Taking the derivatives w.r.t. λ_l and setting these equal to zero yields,

$$0 = \sum_{i=1}^N \left[-1 + \frac{x_i}{\lambda_l} \right] p(l | x_i, \Theta^g) \Rightarrow -N + \frac{1}{\lambda_l} \sum_{i=1}^N x_i p(l | x_i, \Theta^g) = 0$$

$$\lambda_l = \frac{\sum_{i=1}^N x_i p(l | x_i, \Theta^g)}{\sum_{i=1}^N p(l | x_i, \Theta^g)}.$$

Therefore, we have explicit expressions for iterative calculation of the estimates of the mixture parameters, $\Theta=(\alpha_1, \dots, \alpha_M)$, and the *Poisson* distribution parameters, $(\theta_1, \dots, \theta_M)=(\lambda_1, \dots, \lambda_M)$.

Example 2: If we assume d -dimensional Gaussian component distributions with a mean vector μ and covariance matrix Σ , i.e., $\theta = (\mu, \Sigma)$ then the probability density function is

$$p_l(x | \theta) = p_l(x | \mu_l, \Sigma_l) = \frac{1}{(2\pi)^{d/2} |\Sigma_l|^{1/2}} \text{Exp} \left(-\frac{1}{2} (x - \mu_l)^T \Sigma_l^{-1} (x - \mu_l) \right) \quad (7)$$

Then we may derive the update equations [equations (1), (4), (6)] for this specific distribution, we need to recall some results from matrix algebra. The **trace** of a square matrix $tr(A)$ is equal to the sum of A 's diagonal elements. In 1D the trace of a scalar equals that scalar. Also, $tr(A + B) = tr(A) + tr(B)$, and $tr(AB) = tr(BA)$, which implies that if $B = \sum_i x_i x_i^T \Rightarrow \sum_i x_i^T A x_i = tr(AB)$. Also note that $|A|$ indicates the determinant of the matrix A , and $|A^{-1}| = 1/|A|$. Differentiation of a function of a matrix $f(A)$ is accomplished by differentiating with respect to elements of that matrix. Therefore, we define $df(A)/dA$ to be the matrix with (i,j) -th entry equal to $[df(A)/da_{i,j}]$, where $A=(a_{i,j})$. This definition also applies taking derivatives with respect to a vector. First, $d(x^T A x)/dx = (A + A^T)x$. Second, it can be shown that when A is a symmetric matrix:

$$\frac{\partial |A|}{\partial a_{i,j}} = \begin{cases} A_{i,j}, & \text{if } i = j \\ 2A_{i,j}, & \text{if } i \neq j \end{cases}$$

where $A_{i,j}$ is the (i,j) -the cofactor of A . Given the above, we see that:

$$\frac{\partial \text{Log } |A|}{\partial A} = \left\{ \begin{array}{l} A_{i,j} / |A|, \text{ if } i = j \\ 2A_{i,j} / |A|, \text{ if } i \neq j \end{array} \right\} = 2A^{-1} - \text{diag}(A^{-1})$$

by the definition of the inverse of a matrix. Finally, it can be shown that $dtr(AB)/dA = B + B^T - \text{diag}(B)$.

Now, for the d -dimensional Gaussian distribution example, if we take a log of equation (7), ignoring any constant terms (since they disappear after taking derivatives), and substituting into the right side of equation (6), we get:

$$\begin{aligned} Q(\Theta, \Theta^{(g)}) &= \sum_{l=1}^M \sum_{i=1}^N \text{Log}(p_l(x_i | \mu_l, \Sigma_l)) p(l | x_i, \Theta^{(g)}) \\ &= \sum_{l=1}^M \sum_{i=1}^N \left(-\frac{1}{2} \text{Log}(|\Sigma_l|) - \frac{1}{2} (x - \mu_l)^T \Sigma_l^{-1} (x - \mu_l) \right) p(l | x_i, \Theta^{(g)}) \end{aligned} \quad (8)$$

Taking the derivative of equation (8) with respect to μ_l and setting it equal to zero, we get:

$$\sum_{i=1}^N \left(\Sigma_l^{-1} (x_i - \mu_l) p(l | x_i, \Theta^{(g)}) \right) = 0$$

which solving for μ_l yields:

$$\mu_l = \frac{\sum_{i=1}^N x_i p(l | x_i, \Theta^{(g)})}{\sum_{i=1}^N p(l | x_i, \Theta^{(g)})}$$

To find Σ_l , note that we can write equation (8) as:

$$\begin{aligned} \sum_{l=1}^M \left[\frac{1}{2} \text{Log} |\Sigma_l^{-1}| \sum_{i=1}^N p(l | x_i, \Theta^{(g)}) - \frac{1}{2} \sum_{i=1}^N p(l | x_i, \Theta^{(g)}) \text{tr} \left[\Sigma_l^{-1} (x - \mu_l) (x - \mu_l)^T \right] \right] &= \\ = \sum_{l=1}^M \left[\frac{1}{2} \text{Log} |\Sigma_l^{-1}| \sum_{i=1}^N p(l | x_i, \Theta^{(g)}) - \frac{1}{2} \sum_{i=1}^N p(l | x_i, \Theta^{(g)}) \text{tr} \left[\Sigma_l^{-1} N_{l,i} \right] \right] \end{aligned} \quad (9)$$

where $N_{l,i} = (x - \mu_l)(x - \mu_l)^T$. In equation (9) we now take the derivative with respect to the matrix Σ_l^{-1} , and we obtain:

$$\begin{aligned} \frac{1}{2} \sum_{i=1}^N p(l | x_i, \Theta^{(g)}) \left(2\Sigma_l - \text{diag}(\Sigma_l) \right) - \frac{1}{2} \sum_{i=1}^N p(l | x_i, \Theta^{(g)}) \left(2N_{l,i} - \text{diag}(N_{l,i}) \right) &= \\ = \frac{1}{2} \sum_{i=1}^N p(l | x_i, \Theta^{(g)}) \left(2M_{l,i} - \text{diag}(M_{l,i}) \right) = 2S - \text{diag}(S), \end{aligned} \quad (10)$$

where $N_{l,i} = (x - \mu_l)(x - \mu_l)^T$, $M_{l,i} = \Sigma_l - N_{l,i}$ and $S = \frac{1}{2} \sum_{i=1}^N p(l | x_i, \Theta^{(g)}) M_{l,i}$. To find the extreme values (maxima) of the , equation (9) we set the derivative to zero [equation (10)], i.e., $2S - \text{diag}(S) = 0$. This implies that $S=0 \rightarrow \sum_{i=1}^N p(l | x_i, \Theta^{(g)}) (\Sigma_l - N_{l,i}) = 0$. So, we obtain an exact expression (variance-covariance matrix estimate, $1 \leq l \leq M$) for Σ_l .

$$\Sigma_l = \frac{\sum_{i=1}^N p(l | x_i, \Theta^g) N_{l,i}}{\sum_{i=1}^N p(l | x_i, \Theta^g)} = \frac{\sum_{i=1}^N \left[p(l | x_i, \Theta^g) \left((x_i - \mu_l)(x_i - \mu_l)^T \right) \right]}{\sum_{i=1}^N p(l | x_i, \Theta^g)}$$

Summarizing, the estimates of the new parameters in terms of the old parameters (guessed parameters super-indexed by g , $\Theta^g = (\alpha_1^g, \dots, \alpha_M^g; \theta_1^g, \dots, \theta_M^g)$) are as follows:

$$\begin{aligned} \alpha_l^{new} &= \frac{1}{N} \sum_{i=1}^N p(l | x_i, \Theta^g) \\ \mu_l^{new} &= \frac{\sum_{i=1}^N x_i p(l | x_i, \Theta^g)}{\sum_{i=1}^N p(l | x_i, \Theta^g)} \end{aligned} \quad (11)$$

$$\Sigma_l^{new} = \frac{\sum_{i=1}^N \left[p(l | x_i, \Theta^g) \left((x_i - \mu_l^{new})(x_i - \mu_l^{new})^T \right) \right]}{\sum_{i=1}^N p(l | x_i, \Theta^g)}$$

Note that the above equations (11) perform both the expectation step and the maximization step simultaneously. The algorithm proceeds by using the newly derived parameters as the guess for the next iteration.

See Online [SOCR](http://socr.stat.ucla.edu/Applets.dir/MixtureEM.html) Demo at:
<http://socr.stat.ucla.edu/Applets.dir/MixtureEM.html>

References:

A.P.Dempster, N.M. Laird, and D.B. Rubin. Maximum-likelihood from incomplete data via the EM algorithm. *J. Royal Statist. Soc. Ser. B.*, 39, 1977.

C. Bishop. *Neural Networks for Pattern Recognition*. Clarendon Press, Oxford, 1995.

J.A. Bilmes. A Gentle Tutorial of the EM Algorithm and its Application to Parameter Estimation for Gaussian Mixture and Hidden Markov Models. Department of Electrical Engineering and Computer Science, U.C. Berkeley, Berkeley, CA 94704, *TR-97-021*, April 1998.

Z. Ghahramani and M. Jordan. Learning from incomplete data. Technical Report AI Lab Memo No. 1509, CBCL Paper No. 108, MIT AI Lab, August 1995.

M. Jordan and R. Jacobs. Hierarchical mixtures of experts and the EM algorithm. *Neural Computation*, 6:181–214, 1994.

M. Jordon and L. Xu. Convergence results for the EM approach to mixtures of experts architectures. *Neural Networks*, 8:1409–1431, 1996.

L. Rabiner and B.-H. Juang. *Fundamentals of Speech Recognition*. Prentice Hall Signal Processing Series, 1993.

R. Redner and H. Walker. Mixture densities, maximum likelihood and the EM algorithm. *SIAM Review*, 26(2), 1984.

C.F.J. Wu. On the convergence properties of the EM algorithm. *The Annals of Statistics*, 11(1):95–103, 1983.

L. Xu and M.I. Jordan. On convergence properties of the EM algorithm for Gaussian mixtures. *Neural Computation*, 8:129–151, 1996.