

EXPECTATION MAXIMIZATION: FROM THE HORSE'S MOUTH

MAKSIM LEVENTAL

Expectation¹ maximization (EM) is a way to iteratively approximate the maximum likelihood estimators (MLEs) for a parametric family when solving the MLE equations analytically is intractable. Recall that finding the MLEs for a parametric family

$$\mathbf{Y} \sim f_{\mathbf{Y}}(\cdot; \boldsymbol{\theta}) \quad \boldsymbol{\theta} \in \Theta$$

is tantamount to maximizing the likelihood function

$$L(\boldsymbol{\theta}; \mathbf{y}) = f_{\mathbf{Y}}(\mathbf{y}; \boldsymbol{\theta})$$

as a function of $\boldsymbol{\theta}$. The maximum likelihood estimators are estimators of/for $\boldsymbol{\theta}$, usually denoted $\hat{\boldsymbol{\theta}}$. The actual maximization is effected by finding critical points of L with respect to $\boldsymbol{\theta}$ and testing concavity, i.e. solving

$$\nabla L(\boldsymbol{\theta}; \mathbf{y}) = 0$$

and checking the negative definiteness of Hessian of L . In general $\partial_{\theta_i} L(\boldsymbol{\theta}; \mathbf{y})$ might be highly non-linear in θ_i and hence finding each might be very difficult. What to do?

1. EXPECTATION MAXIMIZATION BETA

Suppose² there exists another random, unobserved, vector \mathbf{X} such that were \mathbf{X} in fact observed, maximizing the joint likelihood for \mathbf{X} and \mathbf{Y}

$$L(\boldsymbol{\theta}; \mathbf{x}, \mathbf{y}) = f_{\mathbf{X}, \mathbf{Y}}(\mathbf{x}, \mathbf{y}; \boldsymbol{\theta})$$

is easier than the original problem. \mathbf{X} is called a latent or hidden variable. Now a priori there's absolutely no reason to assume that having more information in hand and complicating the model should make things easier. On the contrary, for example, a good physical model makes simplifying assumptions and thereby becomes tractable. But indeed for some very useful models such as Hidden Markov Models and Gaussian Mixture Models this ansatz does simplify computing the MLEs.

All is going swimmingly except \mathbf{X} is unobserved - the only observed data are \mathbf{Y} . What to do? Estimate \mathbf{X} of course. How? Using the best estimator $E[\mathbf{X}|\mathbf{y}; \boldsymbol{\theta}]$ of \mathbf{X} , based on observed \mathbf{y} , that minimizes the risk for the quadratic loss function³, i.e. minimizes the mean square error. Note two things. Firstly, it's important that you can actually compute *this* in closed form, otherwise you've traded one intractable problem for another. Secondly, since $E[\mathbf{X}|\mathbf{y}; \boldsymbol{\theta}]$ is implicitly a function of $\boldsymbol{\theta}$, which is unknown to begin with⁴ you need some initial guess for it too, otherwise you can't

¹These notes are basically a transcription of notes (taken dutifully by Chris Gianelli) from lectures given by Dr. Kshitij Khare at UF in Spring 15. He's not a horse - it's just an expression.

²Suppose!

³ $l(\mathbf{x}, \hat{\mathbf{x}}) = (\mathbf{x} - \hat{\mathbf{x}})^2$

⁴Don't forget the point of all this is actually to estimate $\boldsymbol{\theta}$.

compute the expectation. Hence the expectation computed is actually $E[\mathbf{X}|\mathbf{y}; \boldsymbol{\theta}^{(r)}]$ where $\boldsymbol{\theta}^{(r)}$ is the current guess for $\boldsymbol{\theta}$. Then once you have this estimate for \mathbf{X} just maximize

$$L(\boldsymbol{\theta}; E[\mathbf{X}|\mathbf{y}; \boldsymbol{\theta}^{(r)}], \mathbf{y}) = f_{\mathbf{X}, \mathbf{Y}}(E[\mathbf{X}|\mathbf{y}; \boldsymbol{\theta}^{(r)}], \mathbf{y}; \boldsymbol{\theta})$$

The procedure alternates between estimating \mathbf{x} using $E[\mathbf{X}|\mathbf{y}; \boldsymbol{\theta}^{(r)}]$ and maximizing $L(\boldsymbol{\theta}; E[\mathbf{X}|\mathbf{y}; \boldsymbol{\theta}^{(r)}], \mathbf{y})$. Compute expectation, then perform maximization, compute expectation, then perform maximization, compute expectation,... hence **Expectation Maximization** algorithm. Just to be clear

Definition 1. Expectation Algorithm Beta. Given some observed data \mathbf{y} , in order to perform the intractable maximization of $L(\boldsymbol{\theta}; \mathbf{y})$, posit the existence of some latent \mathbf{x} such that $L(\boldsymbol{\theta}; \mathbf{x}, \mathbf{y})$ is easier. Set $\boldsymbol{\theta}^{(0)}$ to be some initial guess for $\boldsymbol{\theta}$ then

- (1) E-step: Set $\mathbf{x}^{(r)} = E[\mathbf{X}|\mathbf{y}; \boldsymbol{\theta}^{(r)}]$
- (2) M-step: Set $\boldsymbol{\theta}^{(r+1)} = \arg \max_{\boldsymbol{\theta}} L(\boldsymbol{\theta}; \mathbf{x}^{(r)}, \mathbf{y})$
- (3) Go to step 1 unless $|\boldsymbol{\theta}^{(r+1)} - \boldsymbol{\theta}^{(r)}| < \varepsilon$ for some ε of your choosing.

The equations in step 1 and 2 are called update equations because specify how to update the current estimate for $\mathbf{x}^{(i)}$ and $\boldsymbol{\theta}^{(i)}$.

Example 2. Here's an example: let Y_i be iid such that

$$f_{Y_i}(y_i; \theta) = \theta g_1(y_i) + (1 - \theta) g_0(y_i)$$

where $\theta \in (0, 1)$ and g_0 and g_1 are Gaussians with known means (μ_0, μ_1) and variances (σ_0^2, σ_1^2) . We want the MLE for θ . Note that

$$L(\theta; \mathbf{y}) = \prod_{i=1}^n (\theta g_1(y_i) + (1 - \theta) g_0(y_i))$$

Quite messy⁵. What to do? EM to the rescue! Suppose the Y_i are drawn by a process where a θ -biased coin is flipped and either g_0 or g_1 generates the y_i depending on whether the coin lands heads up or down. The latent variable here then is which Gaussian generated y_i . Hence let X_i be Bernoulli random variables where

$$\begin{aligned} P(X_i = x_i) &= \begin{cases} \theta & \text{if } x_i = 1 \\ 1 - \theta & \text{if } x_i = 0 \end{cases} \\ &= \theta^{x_i} (1 - \theta)^{1-x_i} \end{aligned}$$

and $f_{Y_i|X_i}(y_i|x_i = 1; \theta) = g_1(y_i)$ and $f_{Y_i|X_i}(y_i|x_i = 0; \theta) = g_0(y_i)$. Then

$$f_{\mathbf{X}, \mathbf{Y}}(x_i, y_i; \theta) = f_{Y_i|X_i}(y_i|x_i; \theta) P(X_i = x_i)$$

⁵Even if you think you're clever and try to maximize log-likelihood you're still going to have a rough go at it.

and so

$$\begin{aligned}
 L(\theta; \mathbf{x}, \mathbf{y}) &= \prod_{i=1}^n f_{Y_i|X_i}(y_i|x_i; \theta) P(X_i = x_i) \\
 &= \prod_{i=1}^n (g_1(y_i))^{x_i} (g_0(y_i))^{1-x_i} \theta^{x_i} (1-\theta)^{1-x_i} \\
 &= \theta^{\sum_i x_i} (1-\theta)^{n-\sum_i x_i} \prod_{i=1}^n (g_1(y_i))^{x_i} (g_0(y_i))^{1-x_i}
 \end{aligned}$$

Note the trick in writing $f_{Y_i|X_i}(y_i|x_i; \theta) = (g_1(y_i))^{x_i} (g_0(y_i))^{1-x_i}$ - it comes up a lot for a class of models called mixture models. For the E-step, to compute $E[\mathbf{X}|\mathbf{y}; \theta']$ we use i.i.d-ness and compute $E[X_i|y_i; \theta']$ instead, which will hold for each i . To compute the conditional expectation $E[X_i|y_i; \theta]$ we need the conditional distribution $f_{X_i|Y_i}(x_i|y_i)$. By Bayes' Theorem

$$\begin{aligned}
 f_{X_i|Y_i}(x_i|y_i; \theta) &= \frac{f_{X_i, Y_i}(x_i, y_i; \theta)}{f_{Y_i}(y_i; \theta)} \\
 &= \frac{f_{X_i, Y_i}(x_i, y_i; \theta)}{\sum_{x_i} f_{X_i, Y_i}(x_i, y_i; \theta)} \\
 &= \frac{(g_1(y_i))^{x_i} (g_0(y_i))^{1-x_i} \theta^{x_i} (1-\theta)^{1-x_i}}{\sum_{x_i} (g_1(y_i))^{x_i} (g_0(y_i))^{1-x_i} \theta^{x_i} (1-\theta)^{1-x_i}} \\
 &= \frac{(g_1(y_i))^{x_i} (g_0(y_i))^{1-x_i} \theta^{x_i} (1-\theta)^{1-x_i}}{g_1(y_i) \theta + g_0(y_i) (1-\theta)} \\
 &= \left(\frac{g_1(y_i) \theta}{g_1(y_i) \theta + g_0(y_i) (1-\theta)} \right)^{x_i} \left(\frac{g_0(y_i) (1-\theta)}{g_1(y_i) \theta + g_0(y_i) (1-\theta)} \right)^{1-x_i}
 \end{aligned}$$

So $X_i|Y_i$ is still Bernoulli just renormalized. Hence

$$E[X_i|y_i; \theta] = \frac{g_1(y_i) \theta}{g_1(y_i) \theta + g_0(y_i) (1-\theta)}$$

For the M-step, since

$$L(\theta; \mathbf{x}, \mathbf{y}) = \left[\theta^{\sum_i x_i} (1-\theta)^{n-\sum_i x_i} \right] \left[\prod_{i=1}^n (g_1(y_i))^{x_i} (g_0(y_i))^{1-x_i} \right]$$

and the second term is independent of θ we can just maximize the first term. But this is just the the joint distribution for n i.i.d Bernoulli random variables and the MLE $\hat{\theta}$ is

$$\hat{\theta} = \frac{\sum_i x_i}{n}$$

Therefore, finally, the update equations are

$$\begin{aligned}
 x_i^{(r)} &= \frac{g_1(y_i) \theta^{(r)}}{g_1(y_i) \theta^{(r)} + g_0(y_i) (1-\theta^{(r)})} \\
 \theta^{(r+1)} &= \frac{1}{n} \sum_{i=1}^n x_i^{(r)}
 \end{aligned}$$

■

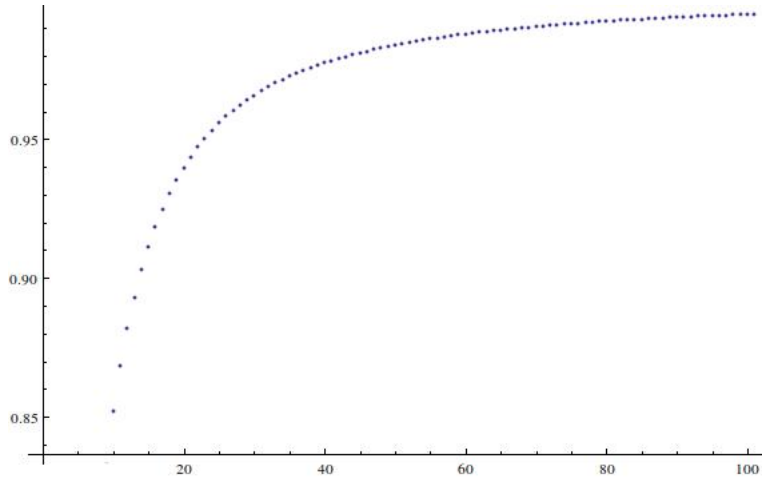


FIGURE 1.1. EM Beta

So why did I dub it EM **Beta** instead of just EM? Because this isn't the standard EM algorithm. But why alter this algorithm at all? What's wrong with it as is? Well there are no convergence guarantees. Indeed the iterates in Example 1 don't converge to the right answer: figure ?? shows the first 100 iterates for $g_1 \sim n(1, 2)$, $g_0 \sim n(3, 4)$, $\theta = 2/3$. So they converge but to $1 \neq 2/3$. Why did I present this algorithm first? Standard EM is slightly more complicated and much less intuitive but legend has it that it was in fact conceived in this way first and then manipulated to get convergence.

2. EXPECTATION MAXIMIZATION FOR REAL

Recall that the whole point of this procedure is to actually maximize the likelihood for \mathbf{Y} . This is equivalent to maximizing the log-likelihood for \mathbf{Y}

$$l(\theta; \mathbf{y}) = \log L(\theta; \mathbf{y}) = \log f_{\mathbf{Y}}(\mathbf{y}; \theta)$$

The ansatz here is the same: suppose there exists another random, unobserved, vector \mathbf{X} such that were \mathbf{X} in fact observed, maximizing the joint log-likelihood for \mathbf{X} and \mathbf{Y}

$$l(\theta; \mathbf{x}, \mathbf{y}) = \log L(\theta; \mathbf{x}, \mathbf{y}) = \log f_{\mathbf{X}, \mathbf{Y}}(\mathbf{x}, \mathbf{y}; \theta)$$

is easier than the original problem. The ease of maximizing $l(\theta; \mathbf{x}, \mathbf{y})$ over $l(\theta; \mathbf{y})$ isn't immediately apparent but first using Bayes' theorem we see that

$$(2.1) \quad l(\theta; \mathbf{y}) = \log f_{\mathbf{Y}}(\mathbf{y}; \theta)$$

$$(2.2) \quad = \log \left(\frac{f_{\mathbf{X}, \mathbf{Y}}(\mathbf{x}, \mathbf{y}; \theta)}{f_{\mathbf{X}|\mathbf{Y}}(\mathbf{x}|\mathbf{y}; \theta)} \right)$$

$$(2.3) \quad = \log(f_{\mathbf{X}, \mathbf{Y}}(\mathbf{x}, \mathbf{y}; \theta)) - \log(f_{\mathbf{X}|\mathbf{Y}}(\mathbf{x}|\mathbf{y}; \theta))$$

If we can easily maximize the first term in ??, with respect to θ and the second term in ?? doesn't spoil things somehow, then we'll indeed maximize $l(\theta; \mathbf{y})$. This seems like a workflow different from that of the beta algorithm⁶ but the analogy follows. Since \mathbf{X} is unobserved we face the same

⁶Here we're explicitly maximizing a function of the likelihood for \mathbf{Y} by maximizing a lowerbound.

difficulty as in the beta algorithm: the terms in ?? can't be computed and must be estimated. We use the same estimator as before $E[\cdot | \mathbf{y}; \boldsymbol{\theta}^{(r)}]$

$$(2.4) \quad l(\boldsymbol{\theta}; \mathbf{y}) = E \left[\log(f_{\mathbf{X}, \mathbf{Y}}(\mathbf{X}, \mathbf{y}; \boldsymbol{\theta})) | \mathbf{y}; \boldsymbol{\theta}^{(r)} \right] - E \left[\log(f_{\mathbf{X}|\mathbf{Y}}(\mathbf{X} | \mathbf{y}; \boldsymbol{\theta})) | \mathbf{y}; \boldsymbol{\theta}^{(r)} \right]$$

Note some subtle things:

- $\mathbf{x} \rightarrow \mathbf{X}$ because each of the terms on the right-hand side of ?? are estimates of the random variable \mathbf{X} as a function of observed data \mathbf{y} .
- The expectations are again computed with respect to the conditional distribution of \mathbf{X} given \mathbf{Y} for fixed (iterate) values of $\boldsymbol{\theta}^{(r)}$, i.e. $f_{\mathbf{X}|\mathbf{Y}}(\mathbf{x} | \mathbf{y}; \boldsymbol{\theta}^{(r)})$.
- The estimators, for which the expectations are computed, are functions of $\boldsymbol{\theta}$.
- Equality is maintained because since $l(\boldsymbol{\theta}; \mathbf{y})$ is independent of \mathbf{X}

$$E \left[l(\boldsymbol{\theta}; \mathbf{y}) | \mathbf{y}; \boldsymbol{\theta}^{(r)} \right] = l(\boldsymbol{\theta}; \mathbf{y})$$

The algorithm then is

Definition 3. Expectation Algorithm. Given some observed data \mathbf{y} , in order to perform the intractable maximization of $L(\boldsymbol{\theta}; \mathbf{y})$, posit the existence of some latent \mathbf{X} such that $L(\boldsymbol{\theta}; \mathbf{x}, \mathbf{y})$ is easier. Set $\boldsymbol{\theta}^{(0)}$ to be some initial guess for $\boldsymbol{\theta}$ then

- (1) E-step: Compute $E \left[\log(f_{\mathbf{X}, \mathbf{Y}}(\mathbf{X}, \mathbf{y}; \boldsymbol{\theta})) | \mathbf{y}; \boldsymbol{\theta}^{(r)} \right]$
- (2) M-step: Set $\boldsymbol{\theta}^{(r+1)} = \arg \max_{\boldsymbol{\theta}} E \left[\log(f_{\mathbf{X}, \mathbf{Y}}(\mathbf{X}, \mathbf{y}; \boldsymbol{\theta})) | \mathbf{y}; \boldsymbol{\theta}^{(r)} \right]$
- (3) Go to step 1 unless $|\boldsymbol{\theta}^{(r+1)} - \boldsymbol{\theta}^{(r)}| < \varepsilon$ for some ε of your choosing.

The only thing to remains is to prove that maximizing $E \left[\log(f_{\mathbf{X}, \mathbf{Y}}(\mathbf{X}, \mathbf{y}; \boldsymbol{\theta})) | \mathbf{y}; \boldsymbol{\theta}^{(r)} \right]$ is sufficient to maximize $l(\boldsymbol{\theta}; \mathbf{y})$. I won't but I'll prove a thing on the way to that result, namely that $l(\boldsymbol{\theta}^{(r)}; \mathbf{y}) \geq l(\boldsymbol{\theta}^{(r-1)}; \mathbf{y})$ and since $l(\boldsymbol{\theta}; \mathbf{y})$ is bounded above (it's a density) the sequence of $l(\boldsymbol{\theta}^{(r)}; \mathbf{y})$ must converge. The last hurdle would be proving that convergence of $l(\boldsymbol{\theta}^{(r)}; \mathbf{y})$ implies the convergence of the iterates $\boldsymbol{\theta}^{(r)}$ themselves[?].

Theorem 4. (Monotonic EM Sequence) The sequence $\{\boldsymbol{\theta}^{(r)}\}$ satisfies $l(\boldsymbol{\theta}^{(r+1)}; \mathbf{y}) \geq l(\boldsymbol{\theta}^{(r)}; \mathbf{y})$.

Proof. Start with

$$l(\boldsymbol{\theta}; \mathbf{y}) = \log(f_{\mathbf{X}, \mathbf{Y}}(\mathbf{x}, \mathbf{y}; \boldsymbol{\theta})) - \log(f_{\mathbf{X}|\mathbf{Y}}(\mathbf{x} | \mathbf{y}; \boldsymbol{\theta}))$$

Take conditional expectation with respect to $\mathbf{X} | \mathbf{Y}; \boldsymbol{\theta}^{(r)}$ of both sides

$$\begin{aligned} \int l(\boldsymbol{\theta}; \mathbf{y}) f_{\mathbf{X}|\mathbf{Y}}(\mathbf{x} | \mathbf{y}; \boldsymbol{\theta}^{(r)}) d\mathbf{x} &= \int \log(f_{\mathbf{X}, \mathbf{Y}}(\mathbf{x}, \mathbf{y}; \boldsymbol{\theta})) f_{\mathbf{X}|\mathbf{Y}}(\mathbf{x} | \mathbf{y}; \boldsymbol{\theta}^{(r)}) d\mathbf{x} \\ &\quad - \int \log(f_{\mathbf{X}|\mathbf{Y}}(\mathbf{x} | \mathbf{y}; \boldsymbol{\theta})) f_{\mathbf{X}|\mathbf{Y}}(\mathbf{x} | \mathbf{y}; \boldsymbol{\theta}^{(r)}) d\mathbf{x} \end{aligned}$$

Since $l(\boldsymbol{\theta}; \mathbf{y})$ is independent of \mathbf{x}

$$\begin{aligned} \int l(\boldsymbol{\theta}; \mathbf{y}) f_{\mathbf{X}|\mathbf{Y}}(\mathbf{x} | \mathbf{y}; \boldsymbol{\theta}^{(r)}) d\mathbf{x} &= l(\boldsymbol{\theta}; \mathbf{y}) \int f_{\mathbf{X}|\mathbf{Y}}(\mathbf{x} | \mathbf{y}; \boldsymbol{\theta}^{(r)}) d\mathbf{x} \\ &= l(\boldsymbol{\theta}; \mathbf{y}) \end{aligned}$$

and so

$$\begin{aligned}
 l(\boldsymbol{\theta}; \mathbf{y}) &= \int \log(f_{\mathbf{X}, \mathbf{Y}}(\mathbf{x}, \mathbf{y}; \boldsymbol{\theta})) f_{\mathbf{X}|\mathbf{Y}}(\mathbf{x}|\mathbf{y}; \boldsymbol{\theta}^{(r)}) d\mathbf{x} \\
 &\quad - \int \log(f_{\mathbf{X}|\mathbf{Y}}(\mathbf{x}|\mathbf{y}; \boldsymbol{\theta})) f_{\mathbf{X}|\mathbf{Y}}(\mathbf{x}|\mathbf{y}; \boldsymbol{\theta}^{(r)}) d\mathbf{x} \\
 &= E \left[\log(f_{\mathbf{X}, \mathbf{Y}}(\mathbf{X}, \mathbf{y}; \boldsymbol{\theta})) | \mathbf{y}; \boldsymbol{\theta}^{(r)} \right] - E \left[\log(f_{\mathbf{X}|\mathbf{Y}}(\mathbf{X}|\mathbf{y}; \boldsymbol{\theta})) | \mathbf{y}; \boldsymbol{\theta}^{(r)} \right]
 \end{aligned}$$

Let

$$Q(\boldsymbol{\theta}, \boldsymbol{\theta}^{(r)}) = E \left[\log(f_{\mathbf{X}, \mathbf{Y}}(\mathbf{X}, \mathbf{y}; \boldsymbol{\theta})) | \mathbf{y}; \boldsymbol{\theta}^{(r)} \right]$$

and $K(\boldsymbol{\theta}, \boldsymbol{\theta}^{(r)}) = E \left[\log(f_{\mathbf{X}|\mathbf{Y}}(\mathbf{X}|\mathbf{y}; \boldsymbol{\theta})) | \mathbf{y}; \boldsymbol{\theta}^{(r)} \right]$. Then

$$l(\boldsymbol{\theta}; \mathbf{y}) = Q(\boldsymbol{\theta}, \boldsymbol{\theta}^{(r)}) - K(\boldsymbol{\theta}, \boldsymbol{\theta}^{(r)})$$

If we can show that

$$\begin{aligned}
 Q(\boldsymbol{\theta}^{(r)}, \boldsymbol{\theta}^{(r)}) &\leq Q(\boldsymbol{\theta}^{(r+1)}, \boldsymbol{\theta}^{(r)}) \\
 &\quad \text{and} \\
 K(\boldsymbol{\theta}^{(r)}, \boldsymbol{\theta}^{(r)}) &\geq K(\boldsymbol{\theta}^{(r+1)}, \boldsymbol{\theta}^{(r)})
 \end{aligned}$$

then it will follow that $l(\boldsymbol{\theta}^{(r+1)}; \mathbf{y}) \geq l(\boldsymbol{\theta}^{(r)}; \mathbf{y})$. Well by definition of

$$\boldsymbol{\theta}^{(r+1)} = \arg \max_{\boldsymbol{\theta}} E \left[\log(f_{\mathbf{X}, \mathbf{Y}}(\mathbf{X}, \mathbf{y}; \boldsymbol{\theta})) | \mathbf{y}; \boldsymbol{\theta}^{(r)} \right] = \arg \max_{\boldsymbol{\theta}} Q(\boldsymbol{\theta}, \boldsymbol{\theta}^{(r)})$$

so

$$Q(\boldsymbol{\theta}, \boldsymbol{\theta}^{(r)}) \leq Q(\boldsymbol{\theta}, \arg \max_{\boldsymbol{\theta}} Q(\boldsymbol{\theta}, \boldsymbol{\theta}^{(r)})) = Q(\boldsymbol{\theta}, \boldsymbol{\theta}^{(r+1)})$$

To see that $K(\boldsymbol{\theta}^{(r)}, \boldsymbol{\theta}^{(r)}) \geq K(\boldsymbol{\theta}^{(r+1)}, \boldsymbol{\theta}^{(r)})$ inspect

$$\int \log \left(\frac{f_{\mathbf{X}|\mathbf{Y}}(\mathbf{x}|\mathbf{y}; \boldsymbol{\theta}^{(r+1)})}{f_{\mathbf{X}|\mathbf{Y}}(\mathbf{x}|\mathbf{y}; \boldsymbol{\theta}^{(r)})} \right) f_{\mathbf{X}|\mathbf{Y}}(\mathbf{x}|\mathbf{y}; \boldsymbol{\theta}^{(r)}) d\mathbf{x}$$

On the one hand

$$\begin{aligned}
 \int \log \left(\frac{f_{\mathbf{X}|\mathbf{Y}}(\mathbf{x}|\mathbf{y}; \boldsymbol{\theta}^{(r+1)})}{f_{\mathbf{X}|\mathbf{Y}}(\mathbf{x}|\mathbf{y}; \boldsymbol{\theta}^{(r)})} \right) f_{\mathbf{X}|\mathbf{Y}}(\mathbf{x}|\mathbf{y}; \boldsymbol{\theta}^{(r)}) d\mathbf{x} &= \int \log(f_{\mathbf{X}|\mathbf{Y}}(\mathbf{x}|\mathbf{y}; \boldsymbol{\theta}^{(r+1)})) f_{\mathbf{X}|\mathbf{Y}}(\mathbf{x}|\mathbf{y}; \boldsymbol{\theta}^{(r)}) d\mathbf{x} \\
 &\quad - \int \log(f_{\mathbf{X}|\mathbf{Y}}(\mathbf{x}|\mathbf{y}; \boldsymbol{\theta}^{(r)})) f_{\mathbf{X}|\mathbf{Y}}(\mathbf{x}|\mathbf{y}; \boldsymbol{\theta}^{(r)}) d\mathbf{x} \\
 &= K(\boldsymbol{\theta}^{(r+1)}, \boldsymbol{\theta}^{(r)}) - K(\boldsymbol{\theta}^{(r)}, \boldsymbol{\theta}^{(r)})
 \end{aligned}$$

on the other hand, by Jensen's⁷ inequality⁸

$$\begin{aligned} \int \log \left(\frac{f_{\mathbf{X}|\mathbf{Y}}(\mathbf{x}|\mathbf{y}; \boldsymbol{\theta}^{(r+1)})}{f_{\mathbf{X}|\mathbf{Y}}(\mathbf{x}|\mathbf{y}; \boldsymbol{\theta}^{(r)})} \right) f_{\mathbf{X}|\mathbf{Y}}(\mathbf{x}|\mathbf{y}; \boldsymbol{\theta}^{(r)}) d\mathbf{x} &\leq \log \left[\int \left(\frac{f_{\mathbf{X}|\mathbf{Y}}(\mathbf{x}|\mathbf{y}; \boldsymbol{\theta}^{(r+1)})}{f_{\mathbf{X}|\mathbf{Y}}(\mathbf{x}|\mathbf{y}; \boldsymbol{\theta}^{(r)})} \right) f_{\mathbf{X}|\mathbf{Y}}(\mathbf{x}|\mathbf{y}; \boldsymbol{\theta}^{(r)}) d\mathbf{x} \right] \\ &= \log \left[\int f_{\mathbf{X}|\mathbf{Y}}(\mathbf{x}|\mathbf{y}; \boldsymbol{\theta}^{(r+1)}) d\mathbf{x} \right] \\ &= 0 \end{aligned}$$

Hence

$$K(\boldsymbol{\theta}^{(r+1)}, \boldsymbol{\theta}^{(r)}) - K(\boldsymbol{\theta}^{(r)}, \boldsymbol{\theta}^{(r)}) \leq 0$$

which completes the proof. \square

3. APPLICATIONS

3.1. Gaussian Mixture Models. EM works really well for mixture models, e.g. \mathbf{Y}_i is distributed iid such that

$$f_{\mathbf{Y}_1}(\mathbf{y}_1) = \sum_{j=1}^m \alpha_j g_j(\mathbf{y}_1; \boldsymbol{\theta}_j)$$

where g_j are (in general multivariate) densities with parameter vectors $\boldsymbol{\theta}_j$ (in general, distinct for distinct j) and $\sum_j \alpha_j = 1$. If $g_j \sim \mathcal{N}(\mu_j, \sigma_j^2)$ then the mixture model is called a Gaussian mixture model (GMM). Naively if you wanted to find the MLEs you would maximize

$$L((\mathbf{y}_1, \dots, \mathbf{y}_n); \boldsymbol{\theta}) = \prod_{i=1}^n f_{\mathbf{Y}_i}(\mathbf{y}_i; \boldsymbol{\theta}) = \prod_{i=1}^n \sum_{j=1}^m \alpha_j g_j(\mathbf{y}_i; \boldsymbol{\theta}_j)$$

Even taking the log won't help you because of the interior sum. EM to the rescue! Let $\mathbf{X} = (X_1, \dots, X_n)$ be the mixture components, i.e. each X_i is a categorical random variable that indicates which component of the mixture density produced the correspondent \mathbf{Y}_i

$$P(X_i = j) = \begin{cases} \alpha_0 & \text{if } x_i = 0 \\ \alpha_1 & \text{if } x_i = 1 \\ \vdots & \\ \alpha_m & \text{if } x_i = m \end{cases}$$

Hence

$$P(X_i = j) = a_j$$

And the $\mathbf{Y}_i|X_i$ takes the form

$$f_{\mathbf{Y}_i|X_i}(\mathbf{y}_i|x_i = j; \boldsymbol{\theta}) = g_j(\mathbf{y}_i; \boldsymbol{\theta}_j)$$

⁷Pronounce Yen-sen you uncultured boor!

⁸For convex $g(X)$ it's the case that $g(E(X)) \leq Eg(X)$ and for concave $g(X)$ the inequality is reversed.

The x_i in $g_{x_i}(\mathbf{y}_i; \boldsymbol{\theta}_{x_i})$ “picks” which mixture component the \mathbf{y}_i is generated by. Then the log-likelihood becomes

$$\begin{aligned}
 l((x_1, \dots, x_n), (\mathbf{y}_1, \dots, \mathbf{y}_n); \boldsymbol{\theta}) &= \log \left[\prod_{i=1}^n f_{X_i, \mathbf{Y}_i}(x_i, \mathbf{y}_i; \boldsymbol{\theta}) \right] \\
 &= \log \left[\prod_{i=1}^n [f_{\mathbf{Y}_i | X_i}(\mathbf{y}_i | x_i; \boldsymbol{\theta}) P(X_i = j)] \right] \\
 &= \log \left[\prod_{i=1}^n [g_j(\mathbf{y}_i; \boldsymbol{\theta}_j) \alpha_j] \right] \\
 &= \sum_{i=1}^n \log (g_j(\mathbf{y}_i; \boldsymbol{\theta}_j) \alpha_j) \\
 &= \sum_{i=1}^n [\log \alpha_j + \log (g_j(\mathbf{y}_i; \boldsymbol{\theta}_j))]
 \end{aligned}$$

There are a lot of indices and subscripts and superscripts to keep track of: $\boldsymbol{\theta}$ is all of the parameters of all of the mixture components, \mathbf{y}_i are observed samples (n of them), the X_i are the unobserved data (the mixture components), the $\boldsymbol{\theta}_j$ are the parameters of the j th mixture density, and the α_j are the mixture coefficients (i.e. in what proportion the j th density contributes). Much as we did for Example 1 we need to compute $f_{X_i | \mathbf{Y}_i}(x_i | \mathbf{y}_i; \boldsymbol{\theta}^{(r)}) = P(X_i = j | \mathbf{y}_i; \boldsymbol{\theta}^{(r)})$ in order to compute the conditional expectations

$$\begin{aligned}
 P(X_i = j | \mathbf{y}_i; \boldsymbol{\theta}^{(r)}) &= \frac{f_{X_i, \mathbf{Y}_i}(j, \mathbf{y}_i; \boldsymbol{\theta}^{(r)})}{f_{\mathbf{Y}_i}(\mathbf{y}_i; \boldsymbol{\theta})} \\
 &= \frac{g_j(\mathbf{y}_i; \boldsymbol{\theta}_j^{(r)}) \alpha_j^{(r)}}{\sum_{k=1}^m (g_k(\mathbf{y}_i; \boldsymbol{\theta}_k^{(r)}) \alpha_k^{(r)})}
 \end{aligned}$$

So again $X_i | \mathbf{Y}_i$ is categorical random variable. By i.i.d

$$P((X_1 = j_1, \dots, X_n = j_n) | (\mathbf{y}_1, \dots, \mathbf{y}_n); \boldsymbol{\theta}^{(r)}) = \prod_{i=1}^n \frac{g_{j_i}(\mathbf{y}_i; \boldsymbol{\theta}_{j_i}^{(r)}) \alpha_{j_i}^{(r)}}{\sum_{k=1}^m (g_k(\mathbf{y}_i; \boldsymbol{\theta}_k^{(r)}) \alpha_k^{(r)})}$$

Now we just need to take the expectation of the log-likelihood against this conditional density

$$\begin{aligned}
 E \left[l((X_1, \dots, X_n), (\mathbf{y}_1, \dots, \mathbf{y}_n); \boldsymbol{\theta}) | (\mathbf{y}_1, \dots, \mathbf{y}_n); \boldsymbol{\theta}^{(r)} \right] &= \\
 \sum_{j_1=1}^m \sum_{j_2=1}^m \cdots \sum_{j_n=1}^m \left(\left[\sum_{k=1}^n \{ \log \alpha_{j_k} + \log [g_{j_k}(\mathbf{y}_k; \boldsymbol{\theta}_{j_k})] \} \right] \prod_{i=1}^n \frac{g_{j_i}(\mathbf{y}_i; \boldsymbol{\theta}_{j_i}^{(r)}) \alpha_{j_i}^{(r)}}{\sum_{k=1}^m (g_k(\mathbf{y}_i; \boldsymbol{\theta}_k^{(r)}) \alpha_k^{(r)})} \right)
 \end{aligned}$$

Pretty ugly right? Suffice it to say you're not actually taking this expectation. So let

$$\gamma_{ij}^{(r)} := P(X_i = j | \mathbf{y}_i; \boldsymbol{\theta}^{(r)}) = \frac{g_j(\mathbf{y}_i; \boldsymbol{\theta}_j^{(r)}) \alpha_j^{(r)}}{\sum_{k=1}^m (g_k(\mathbf{y}_i; \boldsymbol{\theta}_k^{(r)}) \alpha_k^{(r)})}$$

and a helpful lemma:

Lemma 5. *For i.i.d incomplete samples Y_i , with completions X_i*

$$E \left[l((X_1, \dots, X_n), (\mathbf{y}_1, \dots, \mathbf{y}_n); \boldsymbol{\theta}) \mid (\mathbf{y}_1, \dots, \mathbf{y}_n); \boldsymbol{\theta}^{(r)} \right] = \sum_{i=1}^n E \left[l(X_i | \mathbf{y}_i; \boldsymbol{\theta}) \mid \mathbf{y}_i; \boldsymbol{\theta}^{(r)} \right]$$

where \mathbf{y}_i is the i th sample.

Proof. By i.i.d

$$\begin{aligned} E \left[l((X_1, \dots, X_n), (\mathbf{y}_1, \dots, \mathbf{y}_n); \boldsymbol{\theta}) \mid (\mathbf{y}_1, \dots, \mathbf{y}_n); \boldsymbol{\theta}^{(r)} \right] &= \\ E \left[\log \left[\prod_{i=1}^n f(X_i | (\mathbf{y}_1, \dots, \mathbf{y}_n); \boldsymbol{\theta}) \right] \mid (\mathbf{y}_1, \dots, \mathbf{y}_n); \boldsymbol{\theta}^{(r)} \right] &= \\ E \left[\sum_{i=1}^n \log [f(X_i | (\mathbf{y}_1, \dots, \mathbf{y}_n); \boldsymbol{\theta})] \mid (\mathbf{y}_1, \dots, \mathbf{y}_n); \boldsymbol{\theta}^{(r)} \right] &= (\text{ since } X_i \perp X_j \text{ for } i \neq j) \\ \sum_{i=1}^n E \left[\log [f(X_i | (\mathbf{y}_1, \dots, \mathbf{y}_n); \boldsymbol{\theta})] \mid (\mathbf{y}_1, \dots, \mathbf{y}_n); \boldsymbol{\theta}^{(r)} \right] &= (\text{ since } P(X_i = x_i | (\mathbf{Y}_1, \dots, \mathbf{Y}_n)) = P(X_i = x_i | \mathbf{Y}_i)) \\ \sum_{i=1}^n E \left[\log [f(X_i | \mathbf{y}_i; \boldsymbol{\theta})] \mid \mathbf{y}_i; \boldsymbol{\theta}^{(r)} \right] &= \\ \sum_{i=1}^n E \left[l(X_i | \mathbf{y}_i; \boldsymbol{\theta}) \mid \mathbf{y}_i; \boldsymbol{\theta}^{(r)} \right] & \end{aligned}$$

□

Therefore we only need to compute $E \left[l(X_i | \mathbf{y}_i; \boldsymbol{\theta}) \mid \mathbf{y}_i; \boldsymbol{\theta}^{(r)} \right]$ to perform the E-step. Let $j_i \rightarrow j$ for convenience and (μ_j, Σ_j) be the parameters for the j th Gaussian. Actually let's take the 1 dimensional Gaussian case⁹ so $(\mu_j, \Sigma_j) = (\mu_j, \sigma_j^2)$ and $\mathbf{y}_i \rightarrow y_i$.

$$\begin{aligned} E \left[l(X_i | y_i; \boldsymbol{\theta}) \mid y_i; \boldsymbol{\theta}^{(r)} \right] &= \sum_{j=1}^m \left([\log \alpha_j + \log (g_j(y_i; \boldsymbol{\theta}_j))] \gamma_{ij}^{(r)} \right) \\ &= \sum_{j=1}^m \left(\left[\log \alpha_j + \left(-\frac{1}{2} (\log(\sigma_j^2) + \log(2\pi)) - \frac{1}{2\sigma_j^2} (y_i - \mu_j)^2 \right) \right] \gamma_{ij}^{(r)} \right) \end{aligned}$$

Then

$$\begin{aligned} E \left[l((X_1, \dots, X_n) \mid (y_1, \dots, y_n); \boldsymbol{\theta}) \mid (y_1, \dots, y_n); \boldsymbol{\theta}^{(r)} \right] &= \\ \sum_{i=1}^n \sum_{j=1}^m \left([\log \alpha_j + \log (g_j(y_i; \boldsymbol{\theta}_j))] \gamma_{ij}^{(r)} \right) &= \\ \sum_{i=1}^n \sum_{j=1}^m \left(\left[\log \alpha_j + \left(-\frac{1}{2} (\log(\sigma_j^2) + \log(2\pi)) - \frac{1}{2\sigma_j^2} (y_i - \mu_j)^2 \right) \right] \gamma_{ij}^{(r)} \right) \end{aligned}$$

⁹[?] is a good place for the general case.

Define

$$n_j^{(r)} = \sum_{i=1}^n \gamma_{ij}^{(r)} = \sum_{i=1}^n P(X_i = j | \mathbf{y}_i; \boldsymbol{\theta}^{(r)})$$

This is something like the portion of the samples that were generated by the j th component of the density¹⁰. The expectation becomes

$$\begin{aligned} E \left[l((X_1, \dots, X_n) | (y_1, \dots, y_n); \boldsymbol{\theta}) | (y_1, \dots, y_n); \boldsymbol{\theta}^{(r)} \right] &= \\ \sum_{i=1}^n \sum_{j=1}^m \left(\left[\log \alpha_j + \left(-\frac{1}{2} (\log(\sigma_j^2) + \log(2\pi)) - \frac{1}{2\sigma_j^2} (y_i - \mu_j)^2 \right) \right] \gamma_{ij}^{(r)} \right) &= \\ \sum_{i=1}^n \sum_{j=1}^m \log \alpha_j \gamma_{ij}^{(r)} - \sum_{i=1}^n \sum_{j=1}^m \frac{1}{2} (\log(\sigma_j^2) + \log(2\pi)) \gamma_{ij}^{(r)} - \sum_{i=1}^n \sum_{j=1}^m \frac{1}{2\sigma_j^2} (y_i - \mu_j)^2 \gamma_{ij}^{(r)} &= \\ \sum_{j=1}^m n_j^{(r)} \log \alpha_j - \frac{1}{2} \sum_{j=1}^m n_j^{(r)} (\log(\sigma_j^2) + \log(2\pi)) - \frac{1}{2} \sum_{i=1}^n \sum_{j=1}^m \frac{(y_i - \mu_j)^2}{\sigma_j^2} \gamma_{ij}^{(r)} \end{aligned}$$

This is the expression that needs to be maximized with respect to $\alpha_j, \mu_j, \sigma_j^2$ ¹¹. The fully specified maximization problem is

$$\begin{aligned} \max_{\alpha_j, \mu_j, \sigma_j^2} E \left[l((X_1, \dots, X_n) | (y_1, \dots, y_n); \boldsymbol{\theta}) | (y_1, \dots, y_n); \boldsymbol{\theta}^{(r)} \right] \\ \iff \\ \max_{\alpha_j, \mu_j, \sigma_j^2} \sum_{j=1}^m n_j^{(r)} \log \alpha_j - \frac{1}{2} \sum_{j=1}^m n_j^{(r)} (\log(\sigma_j^2) + \log(2\pi)) - \frac{1}{2} \sum_{i=1}^n \sum_{j=1}^m \frac{(y_i - \mu_j)^2}{\sigma_j^2} \gamma_{ij}^{(r)} \\ \text{subject to } \sum_{j=1}^m \alpha_j = 1, \alpha_j \geq 0 \end{aligned}$$

This looks tough but because it's a bunch of uncoupled (somewhat) terms summed together so it's easier than the original log-likelihood

$$\log \left[\sum_{j=1}^m \alpha_j g_j(\mathbf{y}_i; \boldsymbol{\theta}_j) \right]$$

The local maxima for $\alpha_j, \mu_j, \sigma_j^2$ become the new iterates, ie. $\alpha_j^{(r+1)}, \mu_j^{(r+1)}, (\sigma_j^2)^{(r+1)}$. To perform the constrained maximization (only the α_j) we use Lagrange multipliers; form the Lagrangian

$$\mathcal{L}(\boldsymbol{\alpha}, \lambda) = \sum_{j=1}^m n_j^{(r)} \log \alpha_j + \lambda \left(\sum_{j=1}^m \alpha_j - 1 \right)$$

¹⁰Why? Because $\sum_{j=1}^m n_j = \sum_{j=1}^m \sum_{i=1}^n \gamma_{ij} = \sum_{i=1}^n 1 = n$.

¹¹Pay attention to the difference between the arguments that come from the $\log \alpha_j, \mu_j, \sigma_j^2$ and the current estimates of the parameters, those with iterate superscripts (r) that come from conditional density, with respect to which the expectation was taken.

Then computing derivatives and setting to zero to find α_j

$$\frac{\partial \mathcal{L}}{\partial \alpha_j} = \frac{n_j^{(r)}}{\alpha_j} + \lambda$$

Therefore the critical¹² α_j satisfy

$$\alpha_j = -\frac{n_j^{(r)}}{\lambda}$$

To eliminate λ use the equality constraint, i.e. $\sum_{j=1}^m \alpha_j = 1$ implies $\lambda = -\sum_{j=1}^m n_j^{(r)}$ and hence

$$\alpha_j^{(r+1)} = \frac{n_j^{(r)}}{\sum_{j=1}^m n_j^{(r)}} = \frac{n_j^{(r)}}{n}$$

So $\alpha_j^{(r+1)}$ is just the current best estimate for how many of the samples were generated by the j th component of the mixture density¹³. To compute the updates $\mu_j^{(r+1)}$ we just need to do the standard maximization of $E \left[l((X_1, \dots, X_n) | (y_1, \dots, y_n); \boldsymbol{\theta}) | (y_1, \dots, y_n); \boldsymbol{\theta}^{(r)} \right]$ since the μ_j are unconstrained. Therefore

$$\frac{\partial}{\partial \mu_j} E \left[l((X_1, \dots, X_n) | (y_1, \dots, y_n); \boldsymbol{\theta}) | (y_1, \dots, y_n); \boldsymbol{\theta}^{(r)} \right] = \sum_{i=1}^n \frac{(y_i - \mu_j)}{\sigma_j^2} \gamma_{ij}^{(r)} = 0$$

Note the sum of j disappears because only the j th term is non-zero. Then simplifying further

$$\begin{aligned} \sum_{i=1}^n \frac{(y_i - \mu_j)}{\sigma_j^2} \gamma_{ij}^{(r)} &= \frac{1}{\sigma_j^2} \left(\sum_{i=1}^n \gamma_{ij}^{(r)} y_i - \mu_j \sum_{i=1}^n \gamma_{ij}^{(r)} \right) \\ &= \frac{1}{\sigma_j^2} \left(\sum_{i=1}^n \gamma_{ij}^{(r)} y_i - \mu_j n_j^{(r)} \right) = 0 \\ \Rightarrow \\ \mu_j^{(r+1)} &= \frac{1}{n_j^{(r)}} \sum_{i=1}^n \gamma_{ij}^{(r)} y_i \end{aligned}$$

Finally doing the same kind of thing for $(\sigma_j^2)^{(r+1)}$

$$\frac{\partial}{\partial \sigma_j^2} E \left[l((X_1, \dots, X_n) | (y_1, \dots, y_n); \boldsymbol{\theta}) | (y_1, \dots, y_n); \boldsymbol{\theta}^{(r)} \right] = -\frac{1}{2} \frac{n_j^{(r)}}{\sigma_j^2} + \frac{1}{2} \sum_{i=1}^n \frac{(y_i - \mu_j^{(r+1)})^2}{(\sigma_j^2)^2} \gamma_{ij}^{(r)} = 0$$

where μ_j now gets a superscript because it's already been updated. Hence

$$(\sigma_j^2)^{(r+1)} = \frac{1}{n_j^{(r+1)}} \sum_{i=1}^n (y_i - \mu_j^{(r+1)})^2 \gamma_{ij}^{(r)}$$

¹²You need to do the convex analysis (second-derivative test) to determine whether α_j are maxima but they are and I'm not going to.

¹³Pretty much the epitome of the MLE for a parameter.

In summary the update equations for a univariate GMM

$$\begin{aligned}\alpha_j^{(r+1)} &= \frac{n_j^{(r)}}{n} \\ \mu_j^{(r+1)} &= \frac{1}{n_j^{(r)}} \sum_{i=1}^n \gamma_{ij}^{(r)} y_i \\ (\sigma_j^2)^{(r+1)} &= \frac{1}{n_j^{(r+1)}} \sum_{i=1}^n \left(y_i - \mu_j^{(r+1)} \right)^2 \gamma_{ij}^{(r)}\end{aligned}$$

where

$$\begin{aligned}\gamma_{ij}^{(r)} &= \frac{g_j(\mathbf{y}_i; \boldsymbol{\theta}_j^{(r)}) \alpha_j^{(r)}}{\sum_{k=1}^m \left(g_k(\mathbf{y}_i; \boldsymbol{\theta}_k^{(r)}) \alpha_k^{(r)} \right)} \\ n_j^{(r)} &= \sum_i \gamma_{ij}^{(r)}\end{aligned}$$

For multivariate GMM the only difference is $\sigma_j^2 \rightarrow \Sigma_j$ and so you need to use matrix derivatives. Consult [?] for the full derivation

$$\Sigma_j^{(r+1)} = \frac{1}{n_j^{(r+1)}} \sum_{i=1}^n \left(\mathbf{y}_i - \boldsymbol{\mu}_j^{(r+1)} \right) \left(\mathbf{y}_i - \boldsymbol{\mu}_j^{(r+1)} \right)^T \gamma_{ij}^{(r)}$$