

# 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:

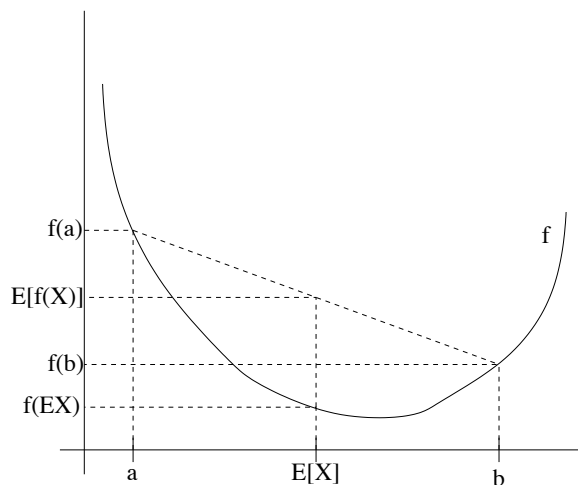
**Theorem.** Let  $f$  be a convex function, and let  $X$  be a random variable. Then:

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

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  $E[f(X)] \geq f(EX)$ .

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 ( $E[f(X)] \leq f(EX)$ , 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

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

marginalizing over joint distribution

But, explicitly finding the maximum likelihood estimates of the parameters  $\theta$  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  $\ell$  (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:<sup>1</sup>

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

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

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

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

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

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  $\theta$  of the parameters, it seems

---

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

natural to try to **make the lower-bound tight at that value of  $\theta$** . I.e., we'll make the inequality above hold with equality at our particular value of  $\theta$ . (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  $\theta$ , 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  $\theta$ . **posterior distribution -> tight bound**

Now, for this choice of the  $Q_i$ 's, Equation (3) gives a lower-bound on the loglikelihood  $\ell$  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  $\theta$ '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

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

**choose theta that maximizes the lower bound**

}

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  $\theta$ , 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  $\theta$ .

### 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  $\phi$ ,  $\mu$  and  $\Sigma$  in a mixture of Gaussians. For the sake of brevity, we carry out the derivations for the M-step updates only for  $\phi$  and  $\mu_j$ , and leave the updates for  $\Sigma_j$  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  $\phi, \mu, \Sigma$ , the quantity

$$\begin{aligned} & \sum_{i=1}^m \sum_{z^{(i)}} Q_i(z^{(i)}) \log \frac{p(x^{(i)}, z^{(i)}; \phi, \mu, \Sigma)}{Q_i(z^{(i)})} \\ &= \sum_{i=1}^m \sum_{j=1}^k Q_i(z^{(i)} = j) \log \frac{p(x^{(i)} | z^{(i)} = j; \mu, \Sigma) p(z^{(i)} = j; \phi)}{Q_i(z^{(i)} = j)} \\ &= \sum_{i=1}^m \sum_{j=1}^k w_j^{(i)} \log \frac{\frac{1}{(2\pi)^{n/2} |\Sigma_j|^{1/2}} \exp\left(-\frac{1}{2}(x^{(i)} - \mu_j)^T \Sigma_j^{-1} (x^{(i)} - \mu_j)\right) \cdot \phi_j}{w_j^{(i)}} \end{aligned}$$

Let's maximize this with respect to  $\mu_l$ . If we take the derivative with respect to  $\mu_l$ , we find

$$\begin{aligned}
\nabla_{\mu_l} \sum_{i=1}^m \sum_{j=1}^k w_j^{(i)} \log \frac{\frac{1}{(2\pi)^{n/2} |\Sigma_j|^{1/2}} \exp \left( -\frac{1}{2} (x^{(i)} - \mu_j)^T \Sigma_j^{-1} (x^{(i)} - \mu_j) \right) \cdot \phi_j}{w_j^{(i)}} \\
= -\nabla_{\mu_l} \sum_{i=1}^m \sum_{j=1}^k w_j^{(i)} \frac{1}{2} (x^{(i)} - \mu_j)^T \Sigma_j^{-1} (x^{(i)} - \mu_j) \\
= \frac{1}{2} \sum_{i=1}^m w_l^{(i)} \nabla_{\mu_l} 2\mu_l^T \Sigma_l^{-1} x^{(i)} - \mu_l^T \Sigma_l^{-1} \mu_l \\
= \sum_{i=1}^m w_l^{(i)} (\Sigma_l^{-1} x^{(i)} - \Sigma_l^{-1} \mu_l)
\end{aligned}$$

Setting this to zero and solving for  $\mu_l$  therefore yields the update rule

$$\mu_l := \frac{\sum_{i=1}^m w_l^{(i)} x^{(i)}}{\sum_{i=1}^m w_l^{(i)}},$$

which was what we had in the previous set of notes.

Let's do one more example, and derive the M-step update for the parameters  $\phi_j$ . Grouping together only the terms that depend on  $\phi_j$ , we find that we need to maximize

$$\sum_{i=1}^m \sum_{j=1}^k w_j^{(i)} \log \phi_j.$$

However, there is an additional constraint that the  $\phi_j$ 's sum to 1, since they represent the probabilities  $\phi_j = p(z^{(i)} = j; \phi)$ . To deal with the constraint that  $\sum_{j=1}^k \phi_j = 1$ , we construct the Lagrangian

$$\mathcal{L}(\phi) = \sum_{i=1}^m \sum_{j=1}^k w_j^{(i)} \log \phi_j + \beta \left( \sum_{j=1}^k \phi_j - 1 \right),$$

where  $\beta$  is the Lagrange multiplier.<sup>2</sup> Taking derivatives, we find

$$\frac{\partial}{\partial \phi_j} \mathcal{L}(\phi) = \sum_{i=1}^m \frac{w_j^{(i)}}{\phi_j} + 1$$

---

<sup>2</sup>We don't need to worry about the constraint that  $\phi_j \geq 0$ , because as we'll shortly see, the solution we'll find from this derivation will automatically satisfy that anyway.

Setting this to zero and solving, we get

$$\phi_j = \frac{\sum_{i=1}^m w_j^{(i)}}{-\beta}$$

I.e.,  $\phi_j \propto \sum_{i=1}^m w_j^{(i)}$ . Using the constraint that  $\sum_j \phi_j = 1$ , we easily find that  $-\beta = \sum_{i=1}^m \sum_{j=1}^k w_j^{(i)} = \sum_{i=1}^m 1 = m$ . (This used the fact that  $w_j^{(i)} = Q_i(z^{(i)} = j)$ , and since probabilities sum to 1,  $\sum_j w_j^{(i)} = 1$ .) We therefore have our M-step updates for the parameters  $\phi_j$ :

$$\phi_j := \frac{1}{m} \sum_{i=1}^m w_j^{(i)}.$$

The derivation for the M-step updates to  $\Sigma_j$  are also entirely straightforward.