

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 strictly positive semi-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:

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

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

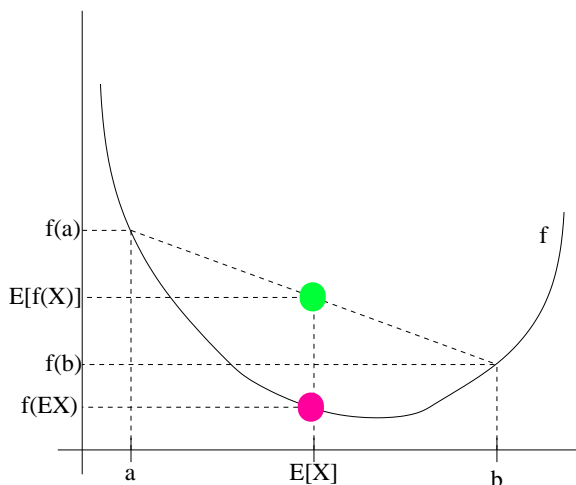
Recall our convention of occasionally dropping the parentheses when writing expectations, so in the theorem above, $f(EX) = f(E[X])$.

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

这里补充一点儿，凸函数并不要求 f 一定可导，但是如果存在二阶导数，那么二阶导数一定大于等于0

凸函数
严格凸函数

如果是严格凸函数，想要等号成立，即 $E[f(X)] = f(E[X])$ ，当且仅当 $p(X=E[X]) = 1$ 时成立。



f 是一个凸函数

X 是随机变量，取到 a, b 的概率都是 0.5 ，所以 $E[X] = (a + b) / 2$

$f(E[X])$ 对应图中的红点

$E[f(X)] = (f(a) + f(b)) / 2$ ，对应图中的绿点

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

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

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

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

2 The EM algorithm

Suppose we have an estimation problem in which we have a training set $\{x^{(1)}, \dots, x^{(m)}\}$ consisting of m independent examples. We wish to fit the parameters of a model $p(x, z)$ to the data, where the likelihood is given by

$$\ell(\theta) = \sum_{i=1}^m \log p(x; \theta)$$

极大对数似然，拟合参数 θ 使画红线等式最大

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

我们考虑 x 和 z (样本特征和标签) 的联合概率分布，引入 z

$$= p(x; \theta)$$

但是这里对参数 θ 的估计是很困难的，因为 z_i 是潜在随机变量（“潜在”也就是说我们并不知道），如果我们知道了 z_i ，那么对 θ 的估计就容易多了

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

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

EM算法采用的策略

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

Q_i 是 Z_i 的一个概率分布，圈起来的部分有点复杂，但是整体上也可以看的很简单（就把这一部分看成一个 z_i 的函数），那么绿色的圈整体上是对于这个函数求期望 $E[X]$ ，接下来就是 $\log(E[X])$ ，根据Jensen不等式，有以下不等式成立： $E[\log(X)] \geq \log(E[X])$ ，注意不等号的方向，不理解可以想一下 \log 函数的大致形状

$$\begin{aligned} \sum_i \log p(x^{(i)}; \theta) &= \sum_i \log \sum_{z^{(i)}} p(x^{(i)}, z^{(i)}; \theta) & (1) \\ &= \sum_i \log \sum_{z^{(i)}} Q_i(z^{(i)}) \frac{p(x^{(i)}, z^{(i)}; \theta)}{Q_i(z^{(i)})} & (2) \\ &\geq \sum_i \sum_{z^{(i)}} Q_i(z^{(i)}) \log \frac{p(x^{(i)}, z^{(i)}; \theta)}{Q_i(z^{(i)})} & (3) \end{aligned}$$

到这里其实我们已经构建了最大似然函数的一个下界了
并且在这个下届中，将对数函数中的连加符号移到了函数外面，使得问题可解

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

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

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 θ of the parameters, it seems

怎么去选择 Q_i 的分布？

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

如果是严格凸函数，想要等号成立，即 $E[f(X)] = f(E[X])$ ，当且仅当 $p(X=E[X]) = 1$ 时成立。
 我们就是要通过设置 Q_i 来让这个等号成立
 为什么要让等号成立，这是个很直观的问题，在后面有相关的证明。
 总的来说就是保证当 $\text{lowbound}(\theta + 1) > \text{lowbound}(\theta)$ 的时候
 有 $l(\theta+1) > l(\theta)$

4

natural to try to make the lower-bound tight at that value of θ . I.e., we'll make the inequality above hold with equality at our particular value of θ . (We'll see later how this enables us to prove that $\ell(\theta)$ increases monotonically with successive iterations of EM.)

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

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

$E[\log(X)] = \log(E[X])$ ，让这个等式成立其实很简单，让 X 只能取一个点（ X 只有一个取值）即可~，也就是这里的 c

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

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

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

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

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

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

初始化 θ

Repeat until convergence {

(E-step) For each i , set

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

类比下高斯混合分布，初始化一些参数，对于当前参数和 x ，求使得上面推出的等号成立的 Q_i （整体上是 Q 的分布）

(M-step) Set

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

关于这个M-step，为了方便自己理解，做个笔记，圈红圈的地方是计算一个联合概率分布，但是 $z^{(i)}$ 我们不知道啊，联合概率分布可以通过贝叶斯公式求解，通常我们都是假设 $z^{(i)}$ 服从什么分布（比如服从多项式分布）， $p(x_i, z_i) = p(x_i | z_i) * p(z_i)$ ， $p(x_i | z_i)$ 是可以求的，就是说 z_i 为多少时， x_i 是分布是知道的（比如， z_i 已知时，设 x_i 服从某些参数下的高斯分布，那就能求出 z_i 为各个值时的 $p(x_i | z_i)$ 了，只不过求出来是带参数的），所以也就大致得到我想要的结论了，首先是能求，只不过求出来带参数，那就可以对这些参数求导，让导为0，找到使蓝圈最大的

}

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

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

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

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

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

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

This first inequality comes from the fact that

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

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

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

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

证明
EM算法
能收敛

这其实也是为什么要让
 $\ell(\theta)$ 和lowbound(θ)相等的原因

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

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

如果定义损失函数为 $J(Q, \theta)$ ，那么 EM 算法可以看成是这个损失函数的坐标上升法

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

3 Mixture of Gaussians revisited MoG的EM推导

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

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

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

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

Next, in the M-step, we need to maximize, with respect to our parameters ϕ, μ, Σ , the quantity

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

!

对公式1的结果求偏导让偏导为0获得极大值

Lets maximize this with respect to μ_l . If we take the derivative with respect to μ_l , we find

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

Setting this to zero and solving for μ_l therefore yields the update rule

$$\mu_l := \frac{\sum_{i=1}^m w_l^{(i)} x^{(i)}}{\sum_{i=1}^m w_l^{(i)}}, \quad \text{u的更新规则}$$

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

Lets do one more example, and derive the M-step update for the parameters ϕ_j . Grouping together only the terms that depend on ϕ_j , we find that we need to maximize

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

约束条件

However, there is an additional constraint that the ϕ_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

$$\mathcal{L}(\phi) = \sum_{i=1}^m \sum_{j=1}^k w_j^{(i)} \log \phi_j + \beta \left(\sum_{j=1}^k \phi_j - 1 \right), \quad \text{构造拉格朗日函数}$$

where β is the Lagrange multiplier.² Taking derivatives, we find

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

²We 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 ϕ_j :

$$\phi_j := \frac{1}{m} \sum_{i=1}^m w_j^{(i)}. \quad \text{\color{violet} \(\phi\) 的更新规则}$$

The derivation for the M-step updates to Σ_j are also entirely straightforward.