

Motivation of Expectation Maximization algorithm

[Ask Question](#)


In the EM algorithm approach we use Jensen's inequality to arrive at

19



$$\log p(x|\theta) \geq \int \log p(z, x|\theta) p(z|x, \theta^{(k)}) dz - \int \log p(z|x, \theta) p(z|x, \theta^{(k)}) dz$$

$$\log p(x|\theta) \geq \int \log p(z, x|\theta) p(z|x, \theta^{(k)}) dz - \int \log p(z|x, \theta) p(z|x, \theta^{(k)}) dz$$



15

and define $\theta^{(k+1)}$ $\theta^{(k+1)}$ by

$$\theta^{(k+1)} = \arg \max_{\theta} \int \log p(z, x|\theta) p(z|x, \theta^{(k)}) dz$$

$$\theta^{(k+1)} = \arg \max_{\theta} \int \log p(z, x|\theta) p(z|x, \theta^{(k)}) dz$$

Everything I read EM just plops it down but I've always felt uneasy by not having an explanation of why the EM algorithm arises naturally. I understand that loglog likelihood is typically dealt with to deal with addition instead of multiplication but the appearance of loglog in the definition of $\theta^{(k+1)}$ $\theta^{(k+1)}$ feels unmotivated to me. Why should one consider loglog and not other monotonic functions? For various reasons I suspect that the "meaning" or "motivation" behind expectation maximization has some kind of explanation in terms of information theory and sufficient statistics. If there were such an explanation that would be much more satisfying than just an abstract algorithm.

[mixture](#)
[expectation-maximization](#)
[share](#)
[cite](#)
[improve this question](#)

edited Jul 20 '13 at 20:02


[jpmuc](#)

9,531 19 46

asked Jul 18 '13 at 11:05


[user782220](#)

412 4 12

migrated from [math.stackexchange.com](#) Jul 20 '13 at 15:49

This question came from our site for people studying math at any level and professionals in related fields.

- 3 [What is the expectation maximization algorithm?](#), *Nature Biotechnology* **26**:897–899 (2008) has a nice picture that illustrates how the algorithm works. – [chl](#) ♦ Jul 20 '13 at 20:58

@chl: I have seen that article. The point I'm asking is that notice that nowhere does it explain why a non-log approach can't work – [user782220](#) Jul 21 '13 at 1:11

[add a comment](#)

5 Answers

[active](#)
[oldest](#)
[votes](#)



The EM algorithm has different interpretations and can arise in different forms in different applications.

10



It all starts with the likelihood function $p(x|\theta)$ or equivalently, the log-likelihood function $\log p(x|\theta)$ we would like to maximize. (We generally use logarithm as it simplifies the calculation: It is strictly monotone, concave, and $\log(ab) = \log a + \log b$.) In an ideal world, the value of p depends only on the **model parameter** θ , so we can search through the space of θ and find one that maximizes p .

However, in many interesting real-world applications things are more complicated, because not all the variables are observed. Yes, we might directly observe x , but some other variables z are unobserved. Because of the **missing variables** z , we are in a kind of chicken-and-eggs situation: Without z we cannot estimate the parameter θ and without θ we cannot infer what the value of z may be.

It is where the EM algorithm comes into play. We start with an initial guess of the model parameters θ and derive the expected values of the missing variables z (i.e., the E step). When we have the values of z , we can maximize the likelihood w.r.t. the parameters θ (i.e., the M step, corresponding to the $\arg \max$ equation in the problem statement). With this θ we can derive the new expected values of z (another E step), so on and so forth. In another word, in each step we assume one of the both, z and θ , is known. We repeat this iterative process until the likelihood cannot be increased anymore.

This is the EM algorithm in a nutshell. It is well-known that the likelihood will never decrease during this iterative EM process. But keep in mind that EM algorithm doesn't guarantee global optimum. That is, it might end up with a local optimum of the likelihood function.

The appearance of \log in the equation of $\theta^{(k+1)}$ is inevitable, because here the function you would like to maximize is written as a log-likelihood.

share cite improve this answer edited Jan 28 '14 at 12:13

answered Jul 20 '13 at 20:17



Weiwei

570 2 10

I don't see how this answers the question. – broncoAbiertoApr 14 '17 at 18:21

add a comment



Likelihood vs. log-likelihood

9

As has already been said, the \log is introduced in maximum likelihood simply because it is generally easier

to optimize sums than products. The reason we don't consider other monotonic functions is that the [logarithm is the unique function](#) with the property of turning products into sums.

Another way to motivate the logarithm is the following: Instead of maximizing the probability of the data under our model, we could equivalently try to minimize the [Kullback-Leibler divergence](#) between the data distribution, $p_{\text{data}}(x)$, and the model distribution, $p(x | \theta)$.

$$D_{\text{KL}}[p_{\text{data}}(x) \parallel p(x | \theta)] = \int p_{\text{data}}(x) \log \frac{p_{\text{data}}(x)}{p(x | \theta)} dx = \text{const} - \int p_{\text{data}}(x) \log p(x | \theta) dx.$$

$$\text{DKL}[p_{\text{data}}(x) \parallel p(x|\theta)] = \int p_{\text{data}}(x) \log p_{\text{data}}(x) p(x|\theta) dx = \text{const} - \int p_{\text{data}}(x) \log p(x|\theta) dx.$$

The first term on the right-hand side is constant in the parameters. If we have N samples from the data distribution (our data points), [we can approximate the second term](#) with the average log-likelihood of the data,

$$\int p_{\text{data}}(x) \log p(x | \theta) dx \approx \frac{1}{N} \sum_n \log p(x_n | \theta).$$

$$\int p_{\text{data}}(x) \log p(x|\theta) dx \approx \frac{1}{N} \sum_n \log p(x_n|\theta).$$

An alternative view of EM

I am not sure this is going to be the kind of explanation you are looking for, but I found the following view of expectation maximization much more enlightening than its motivation via Jensen's inequality (you can find a detailed description in [Neal & Hinton \(1998\)](#) or in Chris Bishop's PRML book, Chapter 9.3).

It is not difficult to show that

$$\log p(x | \theta) = \int q(z | x) \log \frac{p(x, z | \theta)}{q(z | x)} dz + D_{\text{KL}}[q(z | x) \parallel p(z | x, \theta)]$$

$$\log p(x|\theta) = \int q(z|x) \log p(x,z|\theta) q(z|x) dz + \text{DKL}[q(z|x) \parallel p(z|x,\theta)]$$

for any $q(z | x)$. If we call the first term on the right-hand side $F(q, \theta)$, this implies that

$$F(q, \theta) = \int q(z | x) \log \frac{p(x, z | \theta)}{q(z | x)} dz = \log p(x | \theta) - D_{\text{KL}}[q(z | x) \parallel p(z | x, \theta)].$$

$$F(q,\theta) = \int q(z|x) \log p(x,z|\theta) q(z|x) dz = \log p(x|\theta) - \text{DKL}[q(z|x) \parallel p(z|x,\theta)].$$

Because the [KL divergence is always positive](#), $F(q, \theta)$ is a lower bound on the log-likelihood for every fixed q . Now, EM can be viewed as alternately maximizing F with respect to q and θ . In particular, by setting $q(z | x) = p(z | x, \theta)$ in the E-step, we minimize the KL divergence on the right-hand side and thus maximize F .

share cite improve this answer answered Jul 20 '13 at 22:56



Lucas

4,155 16 30

Thanks for the post! Though [the given document](#) doesn't say

logarithm is the unique function turning products into sums. It says logarithm is the only function that fulfills all three listed properties *at the same time*. – Weiwei Jul 22 '13 at 3:47

@Weiwei: Right, but the first condition mainly requires that the function is invertible. Of course, $f(x) = 0$ also implies $f(x + y) = f(x)f(y)$, but this is an uninteresting case. The third condition asks that the derivative at 1 is 1, which is only true for the logarithm to base e . Drop this constraint and you get logarithms to different bases, but still logarithms. – Lucas Jul 22 '13 at 7:49

[add a comment](#)

4

The paper that I found clarifying with respect to expectation-maximization is [Bayesian K-Means as a "Maximization-Expectation" Algorithm \(pdf\)](#) by Welling and Kurihara.

Suppose we have a probabilistic model $p(x, z, \theta)$ with x observations, z hidden random variables, and a total of θ parameters. We are given a dataset D and are forced (by higher powers) to establish $p(z, \theta | D)$.

1. Gibbs sampling

We can approximate $p(z, \theta | D)$ by sampling. Gibbs sampling gives $p(z, \theta | D)$ by alternating:

$$\begin{aligned}\theta &\sim p(\theta | z, D) \\ z &\sim p(z | \theta, D) \\ \theta &\sim p(\theta | z, D) \quad z \sim p(z | \theta, D)\end{aligned}$$

2. Variational Bayes

Instead, we can try to establish a distribution $q(\theta)$ and $q(z)$ and minimize the difference with the distribution we are after $p(\theta, z | D)$. The difference between distributions has a convenient fancy name, the KL-divergence. To minimize $KL[q(\theta)q(z) || p(\theta, z | D)]$ we update:

$$\begin{aligned}q(\theta) &\propto \exp(E_{q(z)}[\log p(\theta, z, D)]) \\ q(z) &\propto \exp(E_{q(\theta)}[\log p(\theta, z, D)]) \\ q(\theta) &\propto \exp(E[\log p(\theta, z, D)]_{q(z)}) \\ q(z) &\propto \exp(E[\log p(\theta, z, D)]_{q(\theta)})\end{aligned}$$

3. Expectation-Maximization

To come up with full-fledged probability distributions for both z and θ might be considered extreme. Why don't we instead consider a point estimate for one of these and keep the other nice and nuanced. In EM the parameter θ is established as the one being unworthy of a full distribution, and set to its MAP (Maximum A Posteriori) value, θ^* .

$$\begin{aligned}\theta^* &= \arg\max_{\theta} E_{q(z)}[\log p(\theta, z, D)] \\ q(z) &= p(z | \theta^*, D)\end{aligned}$$

$$\theta^* = \operatorname{argmax}_{\theta} E[\log p(\theta, z, D)] q(z) q(z) = p(z|\theta^*, D)$$

Here $\theta^* \in \operatorname{argmax}_{\theta}$ $\theta^* \in \operatorname{argmax}$ would actually be a better notation: the argmax operator can return multiple values. But let's not nitpick. Compared to variational Bayes you see that correcting for the loglog by $\exp \exp$ doesn't change the result, so that is not necessary anymore.

4. Maximization-Expectation

There is no reason to treat z as a spoiled child. We can just as well use point estimates z^* for our hidden variables and give the parameters θ the luxury of a full distribution.

$$z^* = \operatorname{argmax}_z E[\log p(\theta, z, D)]_{q(\theta)}$$

$$q(\theta) = p(\theta|z^*, D)$$

$$z^* = \operatorname{argmax}_z E[\log p(\theta, z, D)] q(\theta) q(\theta) = p(\theta|z^*, D)$$

If our hidden variables z are indicator variables, we suddenly have a computationally cheap method to perform inference on the number of clusters. This is in other words: model selection (or automatic relevance detection or imagine another fancy name).

5. Iterated conditional modes

Of course, the poster child of approximate inference is to use point estimates for both the parameters θ as well as the observations z .

$$\theta^* = \operatorname{argmax}_{\theta} p(\theta, z^*, D)$$

$$z^* = \operatorname{argmax}_z p(\theta^*, z, D)$$

$$\theta^* = \operatorname{argmax}_{\theta} p(\theta, z^*, D) \quad z^* = \operatorname{argmax}_z p(\theta^*, z, D)$$

To see how Maximization-Expectation plays out I highly recommend the article. In my opinion, the strength of this article is however not the application to a k -means alternative, but this lucid and concise exposition of approximation.

share cite improve this answer answered Dec 23 '14 at 19:06



Anne van Rossum

308 3 13

(+1) this is a beautiful summary of all methods. – kedarps Apr 20 '17 at 17:54

add a comment



4



There is a useful optimisation technique underlying the EM algorithm. However, it's usually expressed in the language of probability theory so it's hard to see that at the core is a method that has nothing to do with probability and expectation.

Consider the problem of maximising

$$g(x) = \sum_i \exp(f_i(x))$$

$$g(x) = \sum_i \exp(f_i(x))$$

(or equivalently $\log g(x)$ with respect to x). If you write down an expression for $g'(x)$ and set it equal to zero you will often end up with a transcendental equation to solve. These can be nasty.

Now suppose that the f_i play well together in the sense that linear combinations of them give you something easy to optimise. For example, if all of the $f_i(x)$ are quadratic in x then a linear combination of the $f_i(x)$ will also be quadratic, and hence easy to optimise.

Given this supposition, it'd be cool if, in order to optimise $\log g(x) = \log \sum_i \exp(f_i(x))$ we could somehow shuffle the \log past the \sum so it could meet the \exp s and eliminate them. Then the f_i could play together. But we can't do that.

Let's do the next best thing. We'll make another function h that is similar to g . And we'll make it out of linear combinations of the f_i .

Let's say x_0 is a guess for an optimal value. We'd like to improve this. Let's find another function h that matches g and its derivative at x_0 , i.e.

$g(x_0) = h(x_0)$ and $g'(x_0) = h'(x_0)$. If you plot a graph of h in a small neighbourhood of x_0 it's going to look similar to g .

You can show that

$$g'(x) = \sum_i f_i'(x) \exp(f_i(x)).$$

$$g'(x) = \sum_i f_i'(x) \exp(f_i(x)).$$

We want something that matches this at x_0 . There's a natural choice:

$$h(x) = \text{constant} + \sum_i f_i(x) \exp(f_i(x_0)).$$

$$h(x) = \text{constant} + \sum_i f_i(x) \exp(f_i(x_0)).$$

You can see they match at $x = x_0$. We get

$$h'(x) = \sum_i f_i'(x) \exp(f_i(x_0)).$$

$$h'(x) = \sum_i f_i'(x) \exp(f_i(x_0)).$$

As x_0 is a constant we have a simple linear combination of the f_i whose derivative matches g . We just have to choose the constant in h to make $g(x_0) = h(x_0)$.

So starting with x_0 , we form $h(x)$ and optimise that. Because it's similar to $g(x)$ in the neighbourhood of x_0 we hope the optimum of h is similar to the

optimum of g . Once you have a new estimate, construct the next h and repeat.

I hope this has motivated the choice of h . This is exactly the procedure that takes place in EM.

But there's one more important point. Using Jensen's inequality you can show that $h(x) \leq g(x)$ $h(x) \leq g(x)$. This means that when you optimise $h(x)$ you always get an x that makes g bigger compared to $g(x_0)$. So even though h was motivated by its *local* similarity to g , it's safe to *globally* maximise h at each iteration. The hope I mentioned above isn't required.

This also gives a clue to when to use EM: when linear combinations of the arguments to the \exp function are easier to optimise. For example when they're quadratic - as happens when working with mixtures of Gaussians. This is particularly relevant to statistics where many of the standard distributions are from [exponential families](#).

share cite improve this answer answered Oct 20 '16 at 1:51



Dan Piponi
142 6

[add a comment](#)

3

As you said, I will not go into technical details. There are quite a few very nice tutorials. One of my favourites are Andrew Ng's [lecture notes](#). Take a look also at the references [here](#).

1. EM is naturally motivated in mixture models and models with hidden factors in general. Take for example the case of Gaussian mixture models (GMM). Here we model the density of the observations as a weighted sum of K Gaussians:

$$p(x) = \sum_{i=1}^K \pi_i N(x|\mu_i, \Sigma_i)$$

$$p(x) = \sum_{i=1}^K \pi_i \text{tr} N(x|\mu_i, \Sigma_i)$$

where π_i is the probability that the sample x was caused/generated by the i th component, μ_i is the mean of the distribution, and Σ_i is the covariance matrix. The way to understand this expression is the following: each data sample has been generated/caused by one component, but we do not know which one. The approach is then to express the uncertainty in terms of probability (π_i represents the chances that the i th component can account for that sample), and take the weighted sum. As a concrete example, imagine you want to cluster text documents. The idea is to assume that each document belongs to a topic (science, sports,...) which you do not know beforehand!. The possible topics are hidden variables. Then you are given a bunch of documents, and by counting n -grams or whatever features you extract, you want to then find those clusters and see to which cluster each document

belongs to. EM is a procedure which attacks this problem step-wise: the expectation step attempts to improve the assignments of the samples it has achieved so far. The maximization step you improve the parameters of the mixture, in other words, the form of the clusters.

2. The point is not using monotonic functions but convex functions. And the reason is the Jensen's inequality which ensures that the estimates of the EM algorithm will improve at every step.
-