

The Expectation-Maximization (EM) Algorithm

Instructor: Sham Kakade

1 The EM algorithm: the general case.

The EM algorithm [1] is a general procedure to estimate the parameters in a model with latent (unobserved) factors. *EM improves the log likelihood function at every step and will converge.* However, it may not converge to the global optima. Think of it as a more general (and probabilistic) adaptation of the K -means algorithm.

Let x_1, \dots, x_N be all our data. Let h_1, \dots, h_N be the unobserved (“hidden”) variables (we do not observe these). We suppose we have a model which specifies the distribution $\Pr(x, h)$, and we assume that each (x_i, h_i) is sampled independently (and we do not observe the h_i ’s).

The maximum likelihood estimation problem is:

$$\arg \max_{\theta} \sum_{n=1}^N \log \Pr(x_n | \theta).$$

Note that

$$\Pr(x_n | \theta) = \sum_h \Pr(x_n, h_n = h | \theta)$$

where the sum is over all the unobserved variables. Using this, we have that the optimization problem is:

$$\arg \max_{\theta} \sum_{n=1}^N \log \sum_h \Pr(x_n, h_n = h | \theta).$$

which shows how the log likelihood function depends on the parameters θ .

Initialization. Initialize the parameters to some θ . Then alternate between the E-step and the M-step below.

The E step. Compute the *posterior* distribution of h_n given the x_n ’s, for our given parameter. For every n and every possible value of h , set:

$$z_n(h) = \Pr(h_n = h | x_n, \theta)$$

The M step Set the new parameters as the solution of the following optimization problem:

$$\theta \leftarrow \arg \min_{\theta'} \sum_n \sum_h z_n(h) \log \Pr(x_n, h_n = h | \theta')$$

Remark: the key is that in many natural models, this M-step is very easy to solve for in closed form, similar to the case of the mixture of Gaussians problem.

1.1 Convergence

For a general class of latent variable models — models which have unobserved random variables — we can say EM only increases the value of the objective function until a local minima (or a saddle point) is reached. To be precise, we have that:

Lemma 1.1. *Let θ be the parameter at the start of an iteration and let θ' be parameter at the end of the iteration. We have that:*

$$\log \Pr(x_1, \dots, x_N | \theta') \geq \log \Pr(x_1, \dots, x_N | \theta)$$

1.2 (local) Convergence

- If the algorithm has not converged, then, after every M step, the negative log likelihood function decreases in value.
- The algorithm will converge in the limit (to some point, under mild assumptions). Unfortunately, this point may *not* be the global minima. This is related to the that the log likelihood objective function (for these latent variable models) is typically not convex.

2 Another example: the problem of document clustering/topic modeling

Suppose we have N documents x_1, \dots, x_n . Each document is is of length T , and we only keep track of the word count in each document. Let us say $\text{Count}^{(n)}(w)$ is the number of times word w appeared in the n -th document.

We are interested in a “soft” grouping of the documents along with estimating a model for document generation. Let us start with a simple model.

3 A generative model for documents

For a moment, put aside the document clustering problem. Let us instead posit a (probabilistic) procedure which underlies how our documents were generated.

3.1 “Bag of words” model: a (single) topic model

Random variables: a “hidden” (or *latent* topic) $i \in \{1 \dots k\}$ and T -word outcomes w_1, w_2, \dots, w_T which take on some discrete values (these T outcomes constitute a document).

Parameters: the *mixing weights* $\pi_i = \Pr(\text{topic} = i)$, the *topics* $b_{wi} = \Pr(\text{word} = w | \text{topic} = i)$

The generative model for a T -word document, where every document is only about one topic, is specified as follows:

1. sample a topic i , which has probability π_i
2. generate T words w_1, w_2, \dots, w_T , independently. in particular, we choose word w_t as the t -th word with probability $b_{w_t i}$.

Note this generative model ignores the word order, so it is not a particularly faithful generative model.

The conditional independencies implied by the generative procedure imply that we can write the *joint* probability of the outcome topic i occurring with a document containing the words w_1, w_2, \dots, w_T as:

$$\begin{aligned} \Pr(\text{topic} = i \text{ and } w_1, w_2, \dots, w_T) &= \Pr(\text{topic} = i) \Pr(w_1, w_2, \dots, w_T | \text{topic} = i) \\ &= \Pr(\text{topic} = i) \Pr(w_1 | \text{topic} = i) \Pr(w_2 | \text{topic} = i) \Pr(w_T | \text{topic} = i) \\ &= \pi_i b_{w_1 i} b_{w_2 i} \dots b_{w_T i} \end{aligned}$$

where the second to last step follows due to the fact that the words are generated independently given the topic i .

Inference

Suppose we were given a document with w_1, w_2, \dots, w_T . One *inference* question would be: what is the probability the underlying topic is i ? By Bayes rule, we have:

$$\begin{aligned} \Pr(\text{topic} = i | w_1, w_2, \dots, w_T) &= \frac{1}{\Pr(w_1, w_2, \dots, w_T)} \Pr(\text{topic} = i \text{ and } w_1, w_2, \dots, w_T) \\ &= \frac{1}{Z} \pi_i b_{w_1 i} b_{w_2 i} \dots b_{w_T i} \end{aligned}$$

where Z is a number chosen so that the probabilities sum to 1. Critically, note that Z is not a function of i .

3.2 Back to our topic modeling/document clustering problem: Maximum Likelihood estimation

Given the N documents, we could estimate the parameters as follows:

$$\hat{b}, \hat{\pi} = \arg \max_{b, \pi} \log \Pr(x_1, \dots, x_n | b, \pi)$$

How can we do this efficiently?

4 Another EM example: the topic modeling case

The EM algorithm is an *alternating minimization* algorithm. We start at some initialization and then alternate between the E and M steps as follows:

Initialization. Initialize with some \hat{b} and $\hat{\pi}$ (which is not symmetric so that the topic vector are all not identical).

The E step. Estimate the *posterior* probabilities, i.e. the soft assignments, of each document:

$$\widehat{Pr}(\text{topic } i | x_n) = \frac{1}{Z} \hat{\pi}_i \hat{b}_{w_1 i} \hat{b}_{w_2 i} \dots \hat{b}_{w_T i}$$

The M step. Note that $\text{Count}^{(n)}(w)/T$ is the empirical frequency of word w in the n -th document.

Given the posterior probabilities (which we can view as “soft” assignments), we go back and re-estimate the topic probabilities and the mixing weights as follows

$$\hat{b}_{wi} = \frac{\sum_{n=1}^N \widehat{Pr}(\text{topic } i|x_n) \text{Count}^{(n)}(w)/T}{\sum_{n=1}^N \widehat{Pr}(\text{topic } i|x_n)}$$

and

$$\hat{\pi}_i = \frac{1}{N} \sum_{n=1}^N \widehat{Pr}(\text{topic } i|x_n)$$

Now got back to the E -step.

References

- [1] A. P. Dempster, N. M. Laird, and D. B. Rubin. Maximum likelihood from incomplete data via the em algorithm. *JOURNAL OF THE ROYAL STATISTICAL SOCIETY, SERIES B*, 39(1):1–38, 1977.