

EM algorithm

Saki Kuzushima

June 27, 2019

1 Intro

This memo sketches the proof of the monotonicity of the EM algorithm.

2 Set up

We assume a simple mixture model. Let ϕ, γ be the parameters and γ is the latent parameter. Let y be data. We want to maximize the marginal probability of ϕ given y

$$\max p(\phi|y) \tag{1}$$

This often corresponds to the MAP (maximum a posteriori) estimate in the Bayesian framework. However, the functional form of (1) is often intractable. EM is useful when (1) is intractable but the joint distribution of ϕ and γ , $p(\phi, \gamma|y)$ is tractable. One of the typical cases is mixture models.

3 EM algorithm

(1) Start with a rough estimate $\phi = \phi^0$

(2) While the increase in the maximand $p(\phi|y)$ becomes less than the arbitrary threshold, Do;

(2-1): E step. Obtain the Q function. i.e. the expectation of the complete log likelihood.

$$Q = \mathbb{E}_{\gamma}[\log p(\phi, \gamma|y)] = \int \log p(\phi, \gamma|y) p(\gamma|\phi^t, y) d\gamma \tag{2}$$

(2-2): M step. Maximize the Q function with respect to ϕ .

$$\phi^{t+1} = \operatorname{argmax}_{\phi} \mathbb{E}_{\gamma}[\log p(\phi, \gamma|y)] = \operatorname{argmax}_{\phi} \int \log p(\phi, \gamma|y) p(\gamma|\phi^t, y) d\gamma \tag{3}$$

4 Monotonicity proof

We want to show

$$p(\phi^{t+1}|y) > p(\phi^t|y) \quad (4)$$

where t is the current iteration and $t + 1$ is the next iteration. First, we can write (1)

$$p(\phi|y) = \frac{p(\phi, \gamma|y)}{p(\gamma|\phi, y)} \quad (5)$$

$$\iff \log p(\phi|y) = \log p(\phi, \gamma|y) - \log p(\gamma|\phi, y) \quad (6)$$

Taking the expectation with respect to γ in the current iteration, $p(\gamma|\phi^t, y)$,

$$\log p(\phi|y) = \mathbb{E}_\gamma[\log p(\phi, \gamma|y)] - \mathbb{E}_\gamma[\log p(\gamma|\phi, y)] \quad (7)$$

Showing (4) is equivalent to showing

$$\log p(\phi^{t+1}|y) > \log p(\phi^t|y) \quad (8)$$

Using (7),

$$\log p(\phi^{t+1}|y) - \log p(\phi^t|y) \quad (9)$$

$$= \mathbb{E}_\gamma[\log p(\phi^{t+1}, \gamma|y)] - \mathbb{E}_\gamma[\log p(\phi^t, \gamma|y)] \quad (10)$$

$$- \{ \mathbb{E}_\gamma[\log p(\gamma|\phi^{t+1}, y)] - \mathbb{E}_\gamma[\log p(\gamma|\phi^t, y)] \} \quad (11)$$

Observe that the last line is a KL divergence.

$$- \{ \mathbb{E}_\gamma[\log p(\gamma|\phi^{t+1}, y)] - \mathbb{E}_\gamma[\log p(\gamma|\phi^t, y)] \} \quad (12)$$

$$= \mathbb{E}_\gamma[\log p(\gamma|\phi^t, y)] - \mathbb{E}_\gamma[\log p(\gamma|\phi^{t+1}, y)] \quad (13)$$

$$= \int \log p(\gamma|\phi^t, y) p(\gamma|\phi^t, y) d\gamma - \int \log p(\gamma|\phi^{t+1}, y) p(\gamma|\phi^t, y) d\gamma \quad (14)$$

$$= \int \log \frac{p(\gamma|\phi^t, y)}{p(\gamma|\phi^{t+1}, y)} p(\gamma|\phi^t, y) d\gamma \quad (15)$$

$$= KL(p(\gamma|\phi^t, y) \| p(\gamma|\phi^{t+1}, y)) \quad (16)$$

$$\geq KL(p(\gamma|\phi^t, y) \| p(\gamma|\phi^t, y)) = 0 \quad (17)$$

Therefore,

$$\log p(\phi^{t+1}|y) > \log p(\phi^t|y) \quad (18)$$

$$\iff \mathbb{E}_\gamma[\log p(\phi^{t+1}, \gamma|y)] > \mathbb{E}_\gamma[\log p(\phi^t, \gamma|y)] \quad (19)$$

As long as we set ϕ^{t+1} so that it satisfies (19), the EM algorithm is guaranteed to increase the maximand (1). Usually, we choose $\phi^{t+1} = \operatorname{argmax}_\gamma \mathbb{E}_\gamma[\log p(\phi^t, \gamma|y)]$.

5 Refernece

Gelman, A., Stern, H.S., Carlin, J.B., Dunson, D.B., Vehtari, A. and Rubin, D.B., 2013. Bayesian Data Analysis. Chapman and Hall/CRC., Chapter 13 p.320-321.