

Expectation Maximization for mixture models

Seong-Wook Joo

1. Motivation

Suppose measurement data $\mathbf{X} = \{\mathbf{x}_1, \dots, \mathbf{x}_N\}$ is i.i.d. from a distribution parameterized by Θ , i.e., $p(\mathbf{x}|\Theta)$. Then Θ can be estimated by Maximum-Likelihood:

$$\Theta^* = \operatorname{argmax}_{\Theta} [p(\mathbf{X}|\Theta) = \prod_{i=1:N} p(\mathbf{x}_i|\Theta)]$$

But for many problems, it is not possible to find an analytical expression (e.g., mixture model ?)

2. Abstract form of EM

EM can be applied to mainly two types of *missing data problems*: (1) some terms in a data vector are missing (2) (the analytically intractable) likelihood function can be simplified by assuming some missing variables (e.g., classification labels). We concentrate on the latter.

A basic intuition is this: Guess an initial Θ . (E) Using current Θ , obtain an expectation of the complete data likelihood. (M) Find (and update) Θ^+ that maximizes the expectation. Here Θ from the (E) step is a constant. (Fix parameter and get the missing data (and expectation wrt missing data) – fix missing data and find the best parameter - ...). Note that “missing data” means the “distribution” of the missing data determined (parameterized) by Θ .

Let \mathbf{X} be the observed data, \mathbf{Y} be the missing data, $\mathbf{Z} = (\mathbf{X}, \mathbf{Y})$ be the complete data. The complete-data likelihood function is defined as

$$L(\Theta|\mathbf{Z}) = L(\Theta|\mathbf{X}, \mathbf{Y}) = p(\mathbf{X}, \mathbf{Y}|\Theta)$$

Initialization: Guess a value of Θ .

E-Step:

Find the expected value of the complete-data likelihood with respect to missing \mathbf{Y} , given observed \mathbf{X} and current parameter Θ^{i-1} .

$$Q(\Theta, \Theta^{i-1}) = (E[L(\Theta|\mathbf{X}, \mathbf{Y}) | \Theta^{i-1}]) = E [p(\mathbf{X}, \mathbf{Y}|\Theta) | \mathbf{X}, \Theta^{i-1}] = \int_{\mathbf{y}} p(\mathbf{X}, \mathbf{y}|\Theta) f(\mathbf{y} | \mathbf{X}, \Theta^{i-1}) d\mathbf{y}$$

Note that $f(\mathbf{y} | \mathbf{X}, \Theta^{i-1})$ is the distribution of the missing data and is dependent on the observed \mathbf{X} and current parameter Θ^{i-1} .

M-Step:

Maximize the expectation from the E-Step wrt Θ . i.e.,

$$\Theta = \operatorname{argmax}_{\Theta} Q(\Theta, \Theta^{i-1})$$

Repeat the two steps until convergence ($Q(\cdot)$ changes little). The algorithm is guaranteed to converge to a *local* maximum.

3. EM for mixture model – image segmentation

Assume each pixel of an image is associated with a d -dimensional feature vector \mathbf{y} (e.g., the d features include color, texture, etc.) We model our image as a mixture of g image segments each of which has its density and weight. That is, in terms of how to produce the image (“generative model”), each pixel vector in the image is obtained by selecting the l th segment with probability (mixing weight) π_l and then drawing a sample from the l th density:

$$p(\mathbf{x}) = \sum \pi_l p(\mathbf{y}|\theta_l)$$

We would like to determine (a) the parameters of each density (b) the mixing weights (c) the segment label (the missing data). Assuming each component density $p(\mathbf{y}|\theta_l)$ is a d -dimensional Gaussian with parameter $\theta_l = (\mu_l, \Sigma_l)$, we encapsulate all the parameters into a vector $\Theta = (\pi_1, \dots, \pi_g, \theta_1, \dots, \theta_g)$. The log-likelihood for the entire image of N pixels is (iid pixels assumed)

$$\sum_{j=1:N} \log(\sum_{l=1:g} \pi_l p(\mathbf{y}_j|\theta_l))$$

(As can be seen, the \sum inside the log makes it hard to optimize this analytically.) If we knew the missing data, it would be easy to estimate the parameters (via ML) and the mixing weights. If we knew the parameters, we could determine the missing data from the densities. The difficulty is that we know neither the missing data nor the parameters.

To make things easy, let's make the missing variables linear wrt the complete data-log likelihood by the following. Let \mathbf{x} be a vector of g components. If the label of the corresponding pixel is l , then the l th component is 1 and the rest are zeros. Then the complete data-log likelihood is

$$\sum_{j=1:N} (\sum_{l=1:g} \mathbf{x}_{lj} \log(\pi_l p(\mathbf{y}_j|\theta_l)))$$

where \mathbf{x}_{lj} is the l th component of the j th vector. Note the inner \sum is simply selecting one of g segments' log-likelihood. (Since the likelihood function is linear in \mathbf{x}_{lj} we need only take the expectation of \mathbf{x}_{lj} instead of the whole function.)

E-Step:

$$I_{lj} \equiv E(\mathbf{x}_{lj}) = P(j\text{th pixel is from } l\text{th component}) = \pi_l p(\mathbf{y}_j|\theta_l^{(s)}) / \sum_{k=1:g} \pi_k p(\mathbf{y}_j|\theta_k^{(s)})$$

M-Step:

$$\begin{aligned} \pi_l^{(s+I)} &= 1/N \sum_{j=1:N} I_{lj} \\ \mu_l^{(s+I)} &= \sum_{j=1:N} \mathbf{x}_j I_{lj} / (\sum_{j=1:N} I_{lj}) \\ \Sigma_l^{(s+I)} &= \sum_{j=1:N} \{(\mathbf{x}_j - \mu_l^{(s)})(\mathbf{x}_j - \mu_l^{(s)})^T\} I_{lj} / (\sum_{j=1:N} I_{lj}) \end{aligned}$$

Note $\theta_l^{(s+I)} = (\mu_l^{(s+I)}, \Sigma_l^{(s+I)})$ is the weighted mean and covariance.

References

- [1] Bilmes, J., “A gentle tutorial on the EM algorithm and its application to parameter estimation for Gaussian mixture and hidden Markov models” Technical Report ICSI-TR-97-021, University of Berkeley (1998)
- [2] Forsyth and Ponce, “Computer Vision: A Modern Approach” Prentice Hall, 2003, Chapter 16.