

Title: EM algorithm and GMM model

URL: https://en.wikipedia.org/wiki/EM_algorithm_and_GMM_model

PageID: 64563432

Categories: Category:Machine learning, Category:Regression models

Source: Wikipedia (CC BY-SA 4.0). Content may require attribution.

In statistics, EM (expectation maximization) algorithm handles latent variables, while GMM is the Gaussian mixture model.

Background

In the picture below, are shown the red blood cell hemoglobin concentration and the red blood cell volume data of two groups of people, the Anemia group and the Control Group (i.e. the group of people without Anemia). As expected, people with Anemia have lower red blood cell volume and lower red blood cell hemoglobin concentration than those without Anemia.

x is a random vector such as $x := (\text{red blood cell volume}, \text{red blood cell hemoglobin concentration})$, and from medical studies [citation needed] it is known that x are normally distributed in each group, i.e. $x \sim N(\mu, \Sigma)$.

z is denoted as the group where x belongs, with $z_i = 0$ when x_i belongs to Anemia Group and $z_i = 1$ when x_i belongs to Control Group. Also $z \sim \text{Categorical}(k, \phi)$ where $k = 2$, $\phi_j \geq 0$, and $\sum_{j=1}^k \phi_j = 1$. See Categorical distribution.

The following procedure can be used to estimate ϕ, μ, Σ .

A maximum likelihood estimation can be applied:

As the z_i for each x_i are known, the log likelihood function can be simplified as below:

Now the likelihood function can be maximized by making partial derivative over μ, Σ, ϕ , obtaining:

If z_i is known, the estimation of the parameters results to be quite simple with maximum likelihood estimation. But if z_i is unknown it is much more complicated.

Being z a latent variable (i.e. not observed), with unlabeled scenario, the Expectation Maximization Algorithm is needed to estimate z as well as other parameters. Generally, this problem is set as a GMM since the data in each group is normally distributed. [circular reference]

In machine learning, the latent variable z is considered as a latent pattern lying under the data, which the observer is not able to see very directly. x_i is the known data, while ϕ, μ, Σ are the parameter of the model. With the EM algorithm, some underlying pattern z in the data x_i can be found, along with the estimation of the parameters. The wide application of this circumstance in machine learning is what makes EM algorithm so important.

EM algorithm in GMM

The EM algorithm consists of two steps: the E-step and the M-step. Firstly, the model parameters and the $z^{(i)}$ can be randomly initialized. In the E-step, the algorithm tries to

guess the value of $z(i)$ based on the parameters, while in the M-step, the algorithm updates the value of the model parameters based on the guess of $z(i)$ of the E-step. These two steps are repeated until convergence is reached.

The algorithm in GMM is:

Repeat until convergence:

With Bayes Rule, the following result is obtained by the E-step:

$$p(z(i) = j | x(i); \phi, \mu, \Sigma) = \frac{p(x(i) | z(i) = j; \mu, \Sigma) p(z(i) = j; \phi)}{\sum_{l=1}^K p(x(i) | z(i) = l; \mu, \Sigma) p(z(i) = l; \phi)}$$

According to GMM setting, these following formulas are obtained:

$$p(x(i) | z(i) = j; \mu, \Sigma) = \frac{1}{(2\pi)^{n/2} |\Sigma_j|^{1/2}} \exp \left(-\frac{1}{2} (x(i) - \mu_j)^T \Sigma_j^{-1} (x(i) - \mu_j) \right)$$

$$p(z(i) = j; \phi) = \phi_j$$

In this way, a switch between the E-step and the M-step is possible, according to the randomly initialized parameters.

References