

We derive the Expectize-Maximize algorithm (EM) for Gaussian Mixture modeling (GMM). In general, any random space can be modeled by a mixture (weighed linear combination) of normal densities. Then, such an arbitrary density can be represented as:

$$f(\phi; \mathbf{y}) = \sum_{k=1}^K \lambda_k f_k(\theta_k; \mathbf{y}) \quad (1)$$

where

- K is the total number of densities
- \mathbf{y} is the data vector
- $\phi = [\lambda_1 \dots \lambda_K, \theta_1 \dots \theta_K]$ and $\theta_k = [\mu_k, \Sigma_k]$ represent the mean vector and covariance matrix of the k th density
- $f_k \sim \mathcal{N}(\mu_k, \Sigma_k)$
- $\{\lambda_k\}$ represent the mixing weights with constraint $\sum_{k=1}^K \lambda_k = 1$

From now onwards, we drop the bold face notation for vectors and matrices for convenience. To derive the EM algorithm, we need to identify the missing information. Often identifying this missing information can lead to different results affecting convergence. In the mixture modeling problem, if know the label of the data, i.e., if we know the density index that could have possibly produced the observation, we could easily estimate the mixture parameters by simply collecting all the data corresponding to that particular density. In other words, if we associate each observation to one of the K densities, estimating each density parameters is trivial. Hence, we choose the label as the missing information.

Once we identify the missing information, we need to define the complete data likelihood and marginal unobserved data density conditional on the observed data. A very good explanation of the EM algorithm is given in [1, pp.66-87]

Let $u(n) \in \{1 \dots K\}$ be the missing label corresponding to the n th observation $y(n)$. The marginal density (in fact pmf) of this unobserved data is then given as:

$$p(u(n) = k) = \lambda_k \quad (2)$$

The density of the observed data given the unobserved data is:

$$p(y(n)|u(n)) = f_{u(n)}(y(n); \theta_{u(n)}) \quad (3)$$

Then we obtain the complete data density as:

$$p(y(n), u(n); \phi) = f_{u(n)}(y(n); \theta_{u(n)}) \lambda_{u(n)} \quad (4)$$

The complete data density of N independent observations is then:

$$p(y, u; \phi) = \prod_{n=1}^N p(y(n), u(n); \phi) \quad (5)$$

Marginal unobserved data density given the observed data is:

$$p(u|y; \phi) = p(y, u; \phi) / p(y; \phi) \quad (6)$$

where $p(y; \phi)$ is the incomplete or observed data density (likelihood function of the observed data).

The EM algorithm proceeds as follows:

- E-Step: Given an initial estimate of the parameter $\hat{\phi}$, we calculate the expectation of the complete log-likelihood. The expectation is taken over the marginal density of the unobserved data given the observed data. We use $E_{u|y}$ to emphasize the distribution over which the expectation is taken.
- M-step: We maximize this function w.r.t ϕ to obtain an improved (hopefully) estimate of ϕ

E-Step

Since we have independent observations

$$E_{u|y} [\ln p(u, y; \phi)] = \sum_{n=1}^N E_{u(n)|y(n)} [\ln p(u(n), y(n); \theta_{u(n)})] \quad (7)$$

Upon expanding the expectation operator and letting $u(n) = k$,

$$E_{k|y(n)} [\ln p(k, y(n); \theta_k)] = \sum_{k=1}^K [\ln p(k, y(n); \theta_k)] \frac{p(y(n), k; \hat{\theta}_k)}{p(y(n); \hat{\phi})} \quad (8)$$

The dependence of $u(n)$ on k becomes obvious if we recollect that $u(n) \in \{1 \dots K\}$. Also, the denominator is only a function of the data and hence can be ignored while performing the M-step. Hence, we restate the above equation as:

$$E_{k|y(n)} [\ln p(k, y(n); \theta_k)] \propto \sum_{k=1}^K [\ln p(k, y(n); \theta_k)] p(y(n), k; \hat{\theta}_k) \quad (9)$$

Then,

$$E_{u|y} [\ln p(u, y; \phi)] \propto \sum_{n=1}^N \sum_{k=1}^K \hat{\lambda}_k f_k(y(n); \hat{\theta}_k) \left[\ln \lambda_k - \frac{1}{2} \exp \left[(y(n) - \mu_k)^T \Sigma_k^{-1} (y(n) - \mu_k) \right] - \ln |\Sigma_k| \right] \quad (10)$$

The above function has to be maximized subject to the constraint $\sum_{k=1}^K \lambda_k = 1$ and define it as $Q(\phi; \hat{\phi})$.

M-Step

$$\hat{\phi}_{new} = \arg \max_{\phi} Q(\phi; \hat{\phi}) \quad (11)$$

Let φ be the Lagrange multiplier. The then function to be maximized is:

$$U(\phi; \hat{\phi}) = Q(\phi; \hat{\phi}) - \varphi \left(\sum_{k=1}^K \lambda_k - 1 \right) \quad (12)$$

To obtain $\hat{\lambda}_{k,new}$, we set the derivative of $U(\phi; \hat{\phi})$ w.r.t λ_k to zero and solve for the roots and similarly for μ_k and Σ_k . After some straightforward algebraic manipulations, we obtain the expressions for the estimates as [2, Chapter 5]:

$$\hat{\lambda}_{k,new} = \frac{1}{N} \sum_{n=1}^N \frac{\hat{\lambda}_k f_k(y(n); \hat{\theta}_k)}{f(y(n); \hat{\phi})} \quad (13)$$

$$\hat{\mu}_{k,new} = \frac{1}{N \hat{\lambda}_{k,new}} \sum_{n=1}^N \frac{\hat{\lambda}_k f_k(y(n); \hat{\theta}_k)}{f(y(n); \hat{\phi})} y(n) \quad (14)$$

$$\hat{\Sigma}_{k,new} = \frac{1}{N \hat{\lambda}_{k,new}} \sum_{n=1}^N \frac{\hat{\lambda}_k f_k(y(n); \hat{\theta}_k)}{f(y(n); \hat{\phi})} [(y(n) - \hat{\mu}_k)^T (y(n) - \hat{\mu}_k)] \quad (15)$$

References

- [1] Aleksander Dogandžić, “Detection and estimation theory,” Spring, 2005, class notes #3, ECpE Dept., Iowa State University, url: home.eng.iastate.edu/~ald/EE527.html.
- [2] Soma Sekhar Dhavala, *Time-frequency representations: Analysis, Synthesis and Implementation*, M. S, Indian Institute of Technology, Madras, March 2000, url: www.asankhya.org/me/msthesisitmt.html.