

The EM Algorithm
for Maximum Likelihood
Estimation of Gaussian
Mixture Models.

Gaussian Mixture Models

LDA and QDA assume a Gaussian model for $X|Y=0$ and $X|Y=1$, estimate the model parameters via MLE, and plug in to the Bayes classifier.

Advantages

- Simple
- Mahalanobis distance interpretation

Disadvantages

- Data is often non Gaussian.

A Gaussian mixture model (GMM) extends the basic Gaussian model and affords a more flexible representation:

$$g(x|\theta) = \sum_{k=1}^K w_k \phi(x|\mu_k, \Sigma_k),$$

where $\theta = (w_1, \dots, w_K, \mu_1, \Sigma_1, \dots, \mu_K, \Sigma_K)$,

$\sum w_k = 1$, and ϕ denotes the Gaussian pdf.

A classifier is obtained by estimating a GMM for each class, $\hat{\theta}_0$ and $\hat{\theta}_1$, and plugging in to the Bayes classifier:

$$\frac{g(x | \hat{\theta}_1)}{g(x | \hat{\theta}_0)} \underset{0}{\overset{1}{>}} \frac{\hat{\pi}_0}{\hat{\pi}_1}$$

In this lecture we will focus on an algorithm, called the EM algorithm, for computing the MLE of a GMM.

GMMs: A generative view

Suppose

$$X \sim g(x | \theta) = \sum_{k=1}^K w_k \phi(x | \mu_k, \Sigma_k)$$

and θ is given.

How could you generate a realization of X using basic random number generators (e.g., uniform and Gaussian)?

Let $S \in \{1, 2, \dots, K\}$ be a discrete RV such that

$$\Pr \{S = k\} = w_k.$$

Generate X as follows:

1. Generate a realization s of S .
2. Generate $X \sim \phi(x | \mu_s, \Sigma_s)$.

Then the density of X generated in this way is

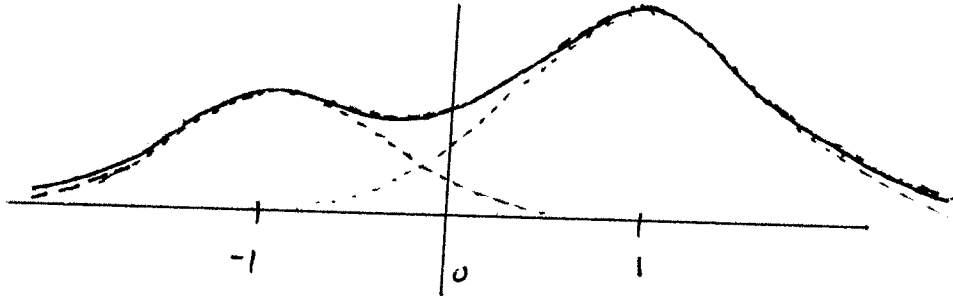
$$\begin{aligned} q(x) &= \sum_{k=1}^K q(x | S=k) \cdot \Pr \{S=k\} \\ &= \sum_{k=1}^K w_k \phi(x | \mu_k, \Sigma_k) \end{aligned}$$

as desired.

The variable S is called a (hidden) state variable. We will imagine that every realization from a GMM is associated with a specific realization of a state variable.

Example

$$X \sim \frac{1}{3} \phi(x|-1,1) + \frac{2}{3} \phi(x|1,1)$$



↑ simulate this $\frac{1}{3}$ of the time
↑ simulate this $\frac{2}{3}$ of the time

Maximum Likelihood Estimation

Consider the likelihood function

$$l(\theta | \underline{x}) = \prod_{i=1}^n g(x_i | \theta)$$

based on a random sample $\underline{x} = (x_1, \dots, x_n)$.

Recall each $x_i \in \mathbb{R}^d$ is a vector.

For the GMM we have

$$l(\theta | \underline{x}) = \prod_{i=1}^n \left(\sum_{k=1}^K w_k \phi(x_i | \mu_k, \Sigma_k) \right)$$

Unfortunately, solving

$$\frac{\partial \ell(\theta | \underline{x})}{\partial \theta} = 0 \quad \text{or} \quad \frac{\partial \log \ell(\theta | \underline{x})}{\partial \theta} = 0$$

is not straightforward.

Complete data

Let $\underline{z} = (s_1, \dots, s_n)$ be the state variables corresponding to $\underline{x} = (x_1, \dots, x_n)$.

We call \underline{z} the unobserved, or latent, or hidden data.

We call $\underline{z} = (\underline{x}, \underline{s})$ the complete data.

Why? Because with knowledge of \underline{z} , the MLE is easy to compute.

MLE of θ given complete data

Introduce the notation

$$I_k = \{i: s_i = k\}$$

$$n_k = |I_k|.$$

Exercise 1

Find the likelihood of the complete data

$$\bar{l}(\theta | \underline{z}) =$$

Solution

$$\bar{l}(\theta | \underline{z}) = \prod_{i=1}^n g(z_i | \theta)$$

$$= \prod_{i=1}^n g(x_i, s_i | \theta)$$

$$= \prod_{i=1}^n g(x_i | \theta, S_i = s_i) \cdot \Pr(S_i = s_i)$$

$$= \prod_{i=1}^n w_{s_i} \cdot \prod_{i=1}^n \phi(x_i | \mu_{s_i}, \Sigma_{s_i})$$

$$= \prod_{k=1}^K w_k^{n_k} \cdot \prod_{k=1}^K \prod_{i \in I_k} \phi(x_i | \mu_k, \Sigma_k).$$

⇒ The likelihood factors into independent terms

⇒ Can maximize w.r.t. (μ_k, Σ_k)

independently for each k .

⇒

$$\hat{\mu}_k = \frac{1}{n_k} \sum_{i \in I_k} x_i$$

$$\hat{\Sigma}_k = \frac{1}{n_k} \sum_{i \in I_k} (x_i - \hat{\mu}_k)(x_i - \hat{\mu}_k)^T$$

What about $\hat{w}_1, \dots, \hat{w}_k$? Because of the constraint $\sum w_k = 1$ we cannot optimize each term separately. However, this ^{is the} MLE for a multinomial distribution:

$$\hat{w}_k = \frac{n_k}{n}$$

Use Lagrange multipliers

In summary, the complete data makes the problem easy. Unfortunately, we don't observe \underline{z} .

The EM Algorithm

Introduce the indicator variable

$$\Delta_{i,k} = \begin{cases} 1 & \text{if } s_i = k \\ 0 & \text{if } s_i \neq k \end{cases}$$

Observe that the complete log-likelihood may be written

$$\begin{aligned} \log \bar{\ell}(\theta | \underline{z}) &= \sum_{i=1}^n \log \left(\sum_{k=1}^K \Delta_{i,k} w_k \phi(x_i | \mu_k, \Sigma_k) \right) \\ &= \sum_{i=1}^n \sum_{k=1}^K \Delta_{i,k} \left[\log w_k + \log \phi(x_i | \mu_k, \Sigma_k) \right] \end{aligned}$$

EM Algorithm

Initialize $\theta^{(0)}$

Repeat

E-Step : Compute

$$Q(\theta, \theta^{(j)}) = E \left[\log \bar{\ell}(\theta | \underline{z}) \mid \theta^{(j)}, \underline{x} \right]$$

M-Step :

$$\theta^{(j+1)} = \arg \max_{\theta} Q(\theta, \theta^{(j)})$$

Until termination criterion satisfied.

E-Step for GMM

$$Q(\theta, \theta^{(j)}) = \sum_{i=1}^n \sum_{k=1}^K \gamma_{i,k}(\theta^{(j)}) \cdot \left[\log w_k - \frac{1}{2} \log |\Sigma| - \frac{d}{2} \log 2\pi - \frac{1}{2} (x_i - \mu_k)^T \Sigma_k^{-1} (x_i - \mu_k) \right]$$

where

$$\gamma_{i,k}(\theta^{(j)}) = E[\Delta_{i,k} \mid x_i, \theta^{(j)}].$$

Exercise Express $\gamma_{i,k}(\theta^{(j)})$ in terms of $x_i, \theta^{(j)}$.

Solution

$$\begin{aligned}\gamma_{i,k}(\theta^{(j)}) &= E[\Delta_{i,k} \mid x_i, \theta^{(j)}] \\ &= \Pr\{\Delta_{i,k} = 1 \mid x_i, \theta^{(j)}\}\end{aligned}$$

$$= \Pr\{s_i = k \mid x_i, \theta^{(j)}\}$$

Bayes rule →

$$= \frac{\Pr\{s_i = k \mid \theta^{(j)}\} \cdot g(x_i \mid s_i = k, \theta^{(j)})}{g(x_i \mid \theta^{(j)})}$$

$$= \frac{w_k^{(j)} \cdot \phi(x_i \mid \mu_k^{(j)}, \Sigma_k^{(j)})}{\sum_{l=1}^K w_l^{(j)} \phi(x_i \mid \mu_l^{(j)}, \Sigma_l^{(j)})}$$

M-Step for GMM

$$\begin{aligned}Q(\theta, \theta^{(j)}) &= \sum_{i=1}^n \sum_{k=1}^K \gamma_{i,k}^{(j)} \left[\log w_k - \frac{d}{2} \log 2\pi \right. \\ &\quad \left. - \frac{1}{2} \log |\Sigma| - \frac{1}{2} (x_i - \mu_k)^T \Sigma_k^{-1} (x_i - \mu_k) \right]\end{aligned}$$

Maximizing w.r.t. θ yields ...

$$\mu_k^{(j+1)} = \frac{\sum_{i=1}^n \gamma_{i,k}^{(j)} \cdot x_i}{\sum_{i=1}^n \gamma_{i,k}^{(j)}}, \quad k = 1, \dots, K$$

$$\sum_k^{(j+1)} = \frac{\sum_{i=1}^n \gamma_{i,k}^{(j)} \cdot (x_i - \mu_k^{(j+1)}) \cdot (x_i - \mu_k^{(j+1)})^T}{\sum_{i=1}^n \gamma_{i,k}^{(j)}}, \quad k = 1, \dots, K$$

$$W_k^{(j+1)} = \frac{1}{n} \sum_{i=1}^n \gamma_{i,k}^{(j)}, \quad k = 1, \dots, K.$$

The number $\gamma_{i,k}$ is called the responsibility of component k for observation i .

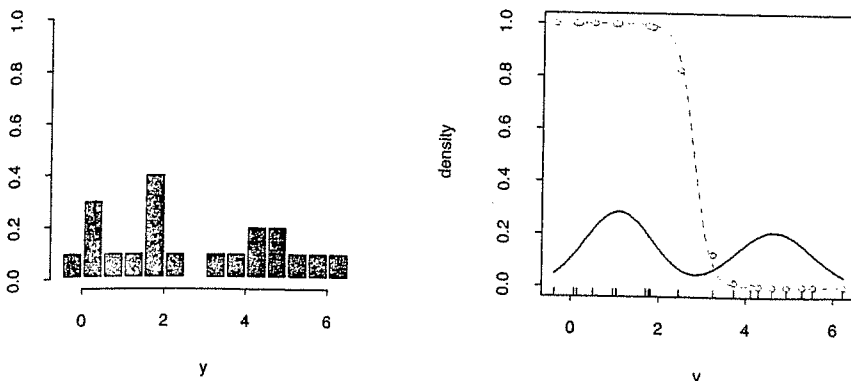


FIGURE 8.5. Mixture example. Left panel: histogram of data. Right panel: maximum likelihood fit of Gaussian densities (solid red) and responsibility (dotted green) of the left component density for observation y , as a function of y .

TABLE 8.1. 20 fictitious data points used in the two-component mixture example in Figure 8.5.

-0.39	0.12	0.94	1.67	1.76	2.44	3.72	4.28	4.92	5.53
0.06	0.48	1.01	1.68	1.80	3.25	4.12	4.60	5.28	6.22

← Hastie, Tibshirani & Friedman, The Elements of Statistical Learning

Initialization

In general, the expected complete-data log-likelihood has several local maxima. Therefore initialization of the EM algorithm is critical.

In fact, a global maximum is obtained by putting $\mu_i = x_i$ for some i , $\Sigma_i = 0$, $w_i = 1$, which is not a useful solution. Hence, we are actually seeking a local maxima.

A good initialization for the GMM is

$$\mu_k^{(0)} = \text{random } x_i \text{ (distinct)}$$

$$\Sigma_k^{(0)} = \text{sample covariance}$$

$$w_k^{(0)} = \frac{1}{K}$$

In practice, it may be beneficial to initialize the algorithm and run it several times, and select the final estimate with largest expected complete-data log-likelihood.

Termination

The EM algorithm can be terminated when

$$|Q(\theta^{(j+1)}, \theta^{(j)}) - Q(\theta^{(j)}, \theta^{(j-1)})| \leq \epsilon$$

for some user defined tolerance ϵ , or when the estimated parameters do not change by more than a predetermined amount.

Question | Why is the EM algorithm couched in terms of the log likelihood as opposed to just the likelihood?

The EM Algorithm in General

The EM algorithm is a general class of algorithms for maximum likelihood estimation.

It is useful whenever MLE given the observed data \underline{x} is intractable, but there exists some unobserved data \underline{z} that makes the MLE tractable. Applications are numerous.

The EM algorithm also applies to maximum a posteriori (MAP) estimation.

Convergence Properties

Theorem: For each $j = 1, 2, \dots$

$$l(\theta^{(j)} | \underline{x}) \geq l(\theta^{(j-1)} | \underline{x})$$

where $l(\theta | \underline{x})$ is the likelihood of θ given the observed data \underline{x} .