

MIXTURE MODELS & EM ALGORITHM

General set-up:

Description of distribution:

- There are K component distributions specified by densities

$$f_j(x; \theta_j) \quad j=1, 2, \dots, K$$

- A "selecting" random variable Z takes values in $\{1, 2, \dots, K\}$, with probability distribution π_j , so

$$\mathbb{P}(Z=j) = \pi_j, \quad j=1, 2, \dots, K$$

Sometimes we will represent Z by a K -dimensional vector $z = (z_1, \dots, z_K)$ with

$$z_j = \begin{cases} 1 & \text{if } j \text{ is selected} \\ 0 & \text{otherwise} \end{cases}$$

So, e.g., if 3 is selected,
 $z = (0, 0, 1, 0, \dots, 0)$

- Given that $Z=j$, the random variable X is sampled from $f_j(\cdot; \theta_j)$ above

Thus the overall density of X is,
by conditioning on Z , given by

$$f_X(x; \theta_1, \dots, \theta_K, \pi) = \sum_{j=1}^K f_j(x; \theta_j) \mathbb{I}(Z=j)$$

$$f_x(x; \theta_1, \dots, \theta_j, \pi) = \sum_{j=1}^k f_j(x; \theta_j) \pi_j$$

The pair (x, z) has joint density

$$f_{x,z}(x, z; \theta_1, \dots, \theta_j, \pi) = \prod_{j=1}^k [\pi_j f_j(x; \theta_j)]^{z_j}$$

because if $z = (0, 0, \dots, 1, \dots, 0)$ in its
 $\overset{\text{1}}{d_0} \overset{\text{1}}{d_0}$ component

vector representation, then

$$\begin{aligned} f_{x,z}(x, (0, \dots, 1, \dots, 0); \theta_1, \dots, \theta_j, \pi) \\ = \prod_j \pi_j f_j(x; \theta_j) \end{aligned}$$

since $z_j = 1$ only for $j = j_0$

Now, suppose $(x_1, z_1), (x_2, z_2), \dots, (x_n, z_n)$
 are sampled from above distribution.

NOTE Now z_i in its vector representation is
 $z_i = (z_{i,1}, \dots, z_{i,k})$ with
 $z_{i,j} = \begin{cases} 1 & \text{if } x_i \text{ is sampled from } f_j(\cdot; \theta_j) \\ 0 & \text{otherwise} \end{cases}$

Note: The z_i 's are HIDDEN, or LATENT
 VARIABLES

We have already seen an example where each $f_j(x; \theta_j)$ is $N(\mu_j, \Sigma_j)$ a Normal (or GAUSSIAN) density

The goal is to estimate $\theta_1, \dots, \theta_k$ and π using only x_1, \dots, x_n , the observed data. For each x_i , we do not know which distribution f_j generated x_i . (The Z_i variable would tell us this information.)

The E-M algorithm provides method for maximizing the likelihood function for the joint observations (x_1, \dots, x_n)

$$\begin{aligned}
 L(\theta_1, \dots, \theta_k, \pi; x_1, \dots, x_n) &= f(x_1, \dots, x_n; \theta_1, \dots, \theta_k, \pi) \\
 &= \prod_{i=1}^n f_x(x_i; \theta_1, \dots, \theta_k, \pi) \quad \text{by INDEPENDENCE} \\
 \text{Here } f_x(x_i; \theta_1, \dots, \theta_k, \pi) &= \sum_{j=1}^k \pi_j f_j(x_i; \theta_j)
 \end{aligned}$$

as above

The E-M algorithm starts by writing down the complete data log-likelihood, that is,

$$\log f_{x,z}(x_1, \dots, x_n, z_1, \dots, z_n; \theta_1, \dots, \theta_K, \pi)$$

$$f_{x,z}(x_1, \dots, x_n, z_1, \dots, z_n; \theta)$$

$$\begin{matrix} \uparrow & \nwarrow \\ (x_1, \dots, x_n) & (z_1, \dots, z_n) \end{matrix} \quad \text{write } \theta = (\theta_1, \dots, \theta_K, \pi)$$

$$= \prod_{i=1}^n f_{x,z}(x_i, z_i; \theta)$$

$$= \prod_{i=1}^n \prod_{j=1}^K [f_j(x_i; \theta_j) \pi_j]^{z_{ij}}$$

since each pair (x_i, z_i) has density

$$\prod_{j=1}^K [f_j(x_i; \theta_j) \pi_j]^{z_{ij}}$$

Taking logarithms shows

$$\log f_{x,z}(x_1, \dots, x_n, z_1, \dots, z_n; \theta)$$

$$= \sum_{i=1}^n \sum_{j=1}^K z_{ij} \log \pi_j + \sum_{i=1}^n \sum_{j=1}^K z_{ij} \log f_j(x_i; \theta_j)$$

In the E-step of the algorithm, the "Q-function" is calculated. This is defined as

$$Q(\theta, \theta_0) = \mathbb{E}[\log f_{x,z}(x_1, \dots, x_n, z_1, \dots, z_n; \theta)]$$

where the expectation \mathbb{E} is calculated using

The distribution $P_{Z|X}(z_1, \dots, z_n | x_1, \dots, x_n; \theta_0)$

That is, the x_1, \dots, x_n are held constant, and the function $\log f_{x,z}(x_1, \dots, x_n, z_1, \dots, z_n; \theta)$ is averaged with respect to the probability distribution over (z_1, \dots, z_n) given by

$P_{Z|X}(z_1, \dots, z_n | x_1, \dots, x_n; \theta_0)$

So using our formula for $\log f_{x,z}(x_i, z_i; \theta)$ above

$$x = (x_1, \dots, x_n) \quad z = (z_1, \dots, z_n)$$

$$\begin{aligned} Q(\theta, \theta_0) &= \sum_{i=1}^n \sum_{j=1}^k \mathbb{E}[z_{i,j} | x_1, \dots, x_n; \theta_0] \log \pi_j \\ &\quad + \sum_{i=1}^n \sum_{j=1}^k \mathbb{E}[z_{i,j} | x_1, \dots, x_n; \theta_0] \log f_j(z_i; \theta_j) \end{aligned}$$

So all we need is

$$\begin{aligned} \mathbb{E}[z_{i,j} | x_1, \dots, x_n; \theta_0] &= \mathbb{E}[z_{i,j} | x_i; \theta_0] \\ &\quad \uparrow z_i \text{ only depends on } x_i \\ &= P(z_{i,j}=1 | x_i; \theta_0) \\ &\stackrel{\text{DEF}}{=} \gamma^{(i)}(j) \end{aligned}$$

This is the conditional prob., given observation x_i , that this observation came from distribution j

Bayes Theorem allows calculation

$$\gamma^{(i)}(j) = \frac{\pi_j f_j(x_i; \theta_j)}{\sum_{l=1}^K \pi_l f_l(x_i; \theta_l)}$$

Note that $\sum_{j=1}^K \gamma^{(i)}(j) = 1$

The distribution $(\gamma^{(1)}(1), \dots, \gamma^{(1)}(K))$ gives
The conditional prob.'s that x_i came from
each of the K distributions $f_j(\cdot; \theta_j)$

To conclude,

$$Q(\theta, \theta_0) = \sum_{j=1}^K \sum_{i=1}^n \gamma^{(i)}(j) \log \pi_j + \sum_{j=1}^K \sum_{i=1}^n \gamma^{(i)}(j) \log f_j(x_i; \theta_j)$$

Let $N_j = \sum_{i=1}^n \gamma^{(i)}(j) = \mathbb{E}[\#x_i's \text{ from } j^{\text{th}} \text{ component}]$

$$Q(\theta, \theta_0) = \underbrace{\sum_{j=1}^K N_j \log \pi_j}_{T_1} + \underbrace{\sum_{j=1}^K \sum_{i=1}^n \gamma^{(i)}(j) \log f_j(x_i; \theta_j)}_{T_2}$$

We are done with the E-step; this
amounts to calculation of γ 's as above

In the M-step, θ_0 is held fixed
and $Q(\theta, \theta^0)$ is maximized over θ

The maximization of $T_1 = \sum_{j=1}^k N_j \log \pi_j$

is always the same, it does not
depend on the form of f_j 's

The solution is always

$$\pi_j = \frac{N_j}{n}$$

The maximization of T_2 depends on form
of f_j 's

Let us try when f_j is poisson(λ_j)

$$f_j(x) = e^{-\lambda_j} \lambda_j^x / x! \quad x=0, 1, \dots$$

$$\log f_j(x) = -\lambda_j + x_i \log \lambda_j - \log x_i!$$

$$\text{Thus } T_2 = \sum_{j=1}^k \sum_{i=1}^n \gamma^{(i)}(j) [-\lambda_j + \pi_i \log \lambda_j - \log x_i!]$$

we can maximize each inner sum independently

$$\begin{aligned} & \frac{\partial}{\partial \lambda_j} \sum_{i=1}^n \gamma^{(i)}(j) [-\lambda_j + \pi_i \log \lambda_j - \log x_i!] \\ &= \sum_{i=1}^n \gamma^{(i)}(j) [-1 + \pi_i / \lambda_j] \end{aligned}$$

$$\begin{aligned}
 &= -\sum_i \gamma^{(t)}(j) + \frac{1}{N_j} \sum_{i=1}^n \gamma^{(t)}(j) x_i \\
 &= -N_j + \frac{1}{N_j} \sum_{i=1}^n \gamma^{(t)}(j) x_i
 \end{aligned}$$

Solving for $\hat{\gamma}_j$ after setting to 0 yields

$$\hat{\gamma}_j = \frac{1}{N_j} \sum_{i=1}^n \gamma^{(t)}(j) x_i$$

Thus given the posterior prob. $\gamma^{(t)}(j)$ that x_i came from component j , the estimate of $\hat{\gamma}_j$ weights the data point x_i by $\gamma^{(t)}(j)$, and the weighted average is the estimate of the mean $\hat{\gamma}_j$ of the j^{th} component

So for Poisson, the "M-step" amounts to writing down

$$\hat{\gamma}_j = \frac{1}{N_j} \sum_{i=1}^n \gamma^{(t)}(j) x_i$$

The EM-algorithm then iterates; once $\hat{\gamma}_j$'s are calculated as above, then and the π_j 's

The γ 's are recomputed

This converges to a local max of the likelihood function.

Above can be applied where each f_j is a BERNoulli vector

e.g. digit recognition

First we describe the "single-variable" dist'n.
let $X = (X_1, \dots, X_D)$ be a collection
of independent BITS, where

$$P(X_l = 1) = \mu_l$$

$$P(X_l = 0) = 1 - \mu_l$$

thus $P(X_1 = x_1, \dots, X_D = x_D) = \prod_{l=1}^D \mu_l^{x_l} (1 - \mu_l)^{1-x_l}$

$$\text{For } x_l = 0 \text{ or } 1 \quad l=1, \dots, D$$

This distribution has parameters (μ_1, \dots, μ_D)

Now suppose we have K such distributions

The j^{th} distribution has parameters

$$(\mu_{j,1}, \dots, \mu_{j,D})$$

Now we mix these with mixing probabilities

$$\pi_1, \dots, \pi_j$$

So the set-up is as above, but now

$$f_j(x_i; \mu_j = (\mu_{j,1}, \dots, \mu_{j,D})) \\ (x_{i,1}, \dots, x_{i,D}) = \prod_{l=1}^D \mu_{j,l}^{x_{i,l}} (1 - \mu_{j,l})^{1-x_{i,l}}$$

i = "which data point"

j = "which component distribution"

l = "which bit out of D bits"

T_1 term in $Q(\theta, \theta_0)$ does not change

T_2 term is

$$\sum_{j=1}^k \sum_{i=1}^n \gamma^{(i)}(j) \left[\sum_{l=1}^D x_{i,l} \log \mu_{j,l} + (1 - x_{i,l}) \log (1 - \mu_{j,l}) \right]$$

Fixing j , we can differentiate with respect to
 μ_{j,l_0} on inside sum

$$\frac{\partial}{\partial \mu_{j,l_0}} = \sum_{i=1}^n \gamma^{(i)}(j) \left[x_{i,l_0} \frac{1}{\mu_{j,l_0}} - (1 - x_{i,l_0}) \frac{1}{1 - \mu_{j,l_0}} \right]$$

Setting to 0 and solving (multiply both sides by
 $\mu_{j,l_0} (1 - \mu_{j,l_0})$)

$$0 = \sum_{i=1}^n \gamma^{(i)}(j) x_{i,l_0} (1 - \mu_{j,l_0}) - \mu_{j,l_0} \left(N_j - \underbrace{\sum_{i=1}^n \gamma^{(i)}(j) x_{i,l_0}}_{S_{j,l_0}(j)} \right)$$

$$0 = S_{j,l_0}(j) - \mu_{j,l_0} S_{j,l_0}(j) - \mu_{j,l_0} N_j + \mu_{j,l_0} \overline{(S_{j,l_0}(j))}$$

$$\mu_{j,l_0} = \frac{S_{j,l_0}(j)}{N_j}$$

In vector notation

$$\mu_j = \frac{s(j)}{N_j} \leftarrow s(j) = (s_1(j), \dots, s_D(j))$$

\uparrow

$$(\mu_{j,1}, \dots, \mu_{j,D})$$

Thus to apply E-M in this case, during M-step, set

$$\mu_j = \frac{s(j)}{N_j}, \quad s(j) = \sum_{l=1}^n \gamma^{(t)}(j) \chi_{i,l}$$

HIGH-LEVEL SUMMARY

For mixtures, the E-step is always the same: compute γ 's

The M-step depends on form of f_j 's