

EM algorithm

A training set $D = \{x^{(1)}, x^{(2)}, \dots, x^{(m)}\}$ m i.i.d. samples.

For each $x^{(i)}$, $i \in [m]$, there is a latent variable $y^{(i)}$ associated with it.

The log-likelihood of D is given by

$$\ell(\theta) = \sum_{i=1}^m \log p(x^{(i)}; \theta)$$

$$= \sum_{i=1}^m \log \sum_y p(x^{(i)}, y; \theta)$$

Explicitly find $\theta^* = \max_{\theta} \ell(\theta)$ may be hard. Instead in EM, we optimize over a lower bound of $\ell(\theta)$.

Iterate until convergence

(E-step) For each $i \in [m]$, set

$$p(y|i) = p(y^{(i)} = y | x^{(i)}; \theta)$$

(M-step) Set

$$\theta := \arg \max_{\theta} \sum_i \sum_y p(y|i) \log \frac{p(x^{(i)}, y; \theta)}{p(y|i)}$$

Let $\hat{\ell}(\theta) = \sum_i \sum_y Q(y|i) \log \frac{p(x^{(i)}, y; \theta)}{Q(y|i)}$. We'll show

where $Q(\cdot|i)$ is some dist over $y^{(i)}$

① $\ell(\theta) \geq \hat{\ell}(\theta)$, i.e. $\hat{\ell}(\theta)$ is a lower bound on $\ell(\theta)$ and $\ell(\theta) = \hat{\ell}(\theta)$ if $Q(y|i) = p(y|i)$.

② Let $\theta^{(t)}$ be parameters obtained after the t -th iteration of the EM algorithm. Then $\ell(\theta^{(t+1)}) \geq \ell(\theta^{(t)}) \geq \ell(\theta^{(t-1)}) \dots \geq \ell(\theta^{(0)})$. That is, the log-likelihood $\ell(\theta)$ is monotonically increasing as EM iterates.

We have

$$\ell(\theta) = \sum_{i=1}^m \log \sum_y Q(y|i) \frac{p(x^{(i)}, y; \theta)}{Q(y|i)}$$

$$\geq \sum_{i=1}^m \sum_y Q(y|i) \log \frac{p(x^{(i)}, y; \theta)}{Q(y|i)}$$

$$= \hat{\ell}(\theta)$$

which gives ①. The equality holds if

$$\frac{p(x^{(i)}, y; \theta)}{Q(y|i)} = \text{const}$$

$$\Rightarrow Q(y|i) \propto p(x^{(i)}, y; \theta)$$

$$\Rightarrow Q(y|i) = \frac{p(x^{(i)}, y; \theta)}{\sum_y p(x^{(i)}, y; \theta)} = \frac{p(x^{(i)}, y; \theta)}{p(x^{(i)}; \theta)} = p(y|i).$$

Hence $\ell(\theta) = \hat{\ell}(\theta)$ if $Q(y|i) = p(y|i)$

we used the fact
 $\log E[X] \geq E[\log X]$
 Check Jensen's inequality

$E[f(X)] \geq f(E[X])$ for convex

Here $-\log x$ is convex, then

$$E[-\log X] \geq -\log(E[X])$$

$$\Rightarrow E[\log X] \leq \log(E[X])$$

Now we show ②. In this case $q(y|i) = p(y^{(i)} = y | x^{(i)}; \theta^{(t)})$.

$$\begin{aligned} \ell(\theta^{(t+1)}) &\geq \sum_i \sum_y p(y|i) \log \frac{p(x^{(i)}, y; \theta^{(t+1)})}{p(y|i)} \quad (\text{by ①}) \\ &\geq \sum_i \sum_y p(y|i) \log \frac{p(x^{(i)}, y; \theta^{(t)})}{p(y|i)} \quad (\text{by the maximization in M step}) \\ &= \ell(\theta^{(t)}) \quad (\text{by } \ell(\theta^{(t)}) = \hat{\ell}(\theta^{(t)}) \text{ if } q(y|i) = p(y|i)) \\ &\Rightarrow \ell(\theta^{(t+1)}) \geq \ell(\theta^{(t)}). \end{aligned}$$

Hence $\ell(\theta^{(t)})$ is monotonically increasing and bounded by 0 from above, $\ell(\theta^{(t)})$ is guaranteed to converge. However, convergence to the global max of $\ell(\theta)$ is not guaranteed.

In the application, run EM several times using randomly initialized μ 's and then choose the best one.

Example: problem 2 and 3 of part 2 in project 3 (show Lagrange multipliers)

Per all

$$Pr(\mathcal{X}_{\text{obs}}, \mathcal{Z} | \pi, \alpha) = \prod_{i=1}^n Pr(\mathcal{Z}^{(i)} | \pi) \prod_{d=1}^D Pr(\mathcal{X}_d^{(i)} | \alpha_d, \mathcal{Z}^{(i)}) \quad [[\alpha_d^{(i)} \text{ is not missing}]]$$

E step

$$p(\mathcal{Z}^{(i)} | x^{(i)}, \pi, \alpha) = \frac{\pi_{\mathcal{Z}^{(i)}} \prod_{d=1}^D p(\alpha_d, \mathcal{Z}^{(i)}, x_d^{(i)})}{\sum_{j=1}^K \pi_j \prod_{d=1}^D p(\alpha_d, j, x_d^{(i)})} \quad [[x_d^{(i)} \text{ is not missing}]]$$

M step:

$$J(\alpha, \pi) = \sum_k p(\mathcal{Z} = k | \mathcal{X}_{\text{obs}}) \log(\mathcal{X}_{\text{obs}} | \alpha, \pi)$$

$$= \sum_{i=1}^n \sum_{k=1}^K p(\mathcal{Z} = k | x^{(i)}) \log p(\mathcal{Z} = k | \pi) +$$

$$\sum_{i=1}^n \sum_{k=1}^K \sum_{d=1}^D p(\mathcal{Z} = k | x^{(i)}) \log p(x_d^{(i)} | \mathcal{Z} = k, \alpha_d, \pi) \quad [[x_d^{(i)} \text{ is not missing}]]$$

$$= \sum_{i=1}^n \sum_{k=1}^K p(\mathcal{Z} = k | x^{(i)}) \log \pi_k + \sum_{i=1}^n \sum_{k=1}^K \sum_{d=1}^D p(\mathcal{Z} = k | x^{(i)}) \log \alpha_{d,k} x_d^{(i)} \quad [[x_d^{(i)} \text{ is observed}]]$$

in M step, use gradient

$$\max_{\alpha, \pi} J(\alpha, \pi)$$

$$\alpha, \pi$$

$$\text{s.t. } \sum_{k=1}^K \pi_k = 1 \quad \dots \textcircled{1}$$

$$\sum_{\ell=1}^L \alpha_{d,k,\ell} = 1, \forall d \in [D], k \in [K] \quad \dots \textcircled{2}$$

$$\Leftrightarrow \max_{\alpha, \pi, \beta, \eta_{d,k}} J(\alpha, \pi) + \beta \left(\sum_{k=1}^K \pi_k - 1 \right) + \sum_{d=1}^D \sum_{k=1}^K \eta_{d,k} \left(\sum_{\ell=1}^L \alpha_{d,k,\ell} - 1 \right)$$

$$F(\alpha, \pi, \beta, \eta_{d,k}; \forall d \in [D], k \in [K])$$

necessary condition.

$$\frac{\partial F}{\partial \alpha_{d,k,\ell}} = 0 \quad \forall d \in [D], k \in [K], \ell \in [L] \quad \textcircled{3}$$

$$\frac{\partial F}{\partial \pi_k} = 0 \quad \forall k \in [K] \quad \textcircled{4}$$

$$\frac{\partial F}{\partial \beta} = 0 \Rightarrow \textcircled{1}$$

$$\frac{\partial F}{\partial \eta_{d,k}} = 0 \quad \forall d \in [D], k \in [K] \Rightarrow \textcircled{2}$$

From $\textcircled{4}$ and $\textcircled{1}$, we have

$$\left. \begin{aligned} \sum_{i=1}^n \frac{P(z^{(i)}=k | x^{(i)})}{\pi_k} + \beta &= 0 \quad \forall k \in [K] \\ \sum_{k=1}^K \pi_k &= 1 \end{aligned} \right\} \Rightarrow \pi_k = \frac{1}{n} \sum_{i=1}^n P(z^{(i)}=k | x^{(i)})$$

From $\textcircled{3}$ and $\textcircled{2}$, we have

$$\sum_{i=1}^n \frac{P(z^{(i)}=k | x^{(i)})}{\alpha_{d,k,\ell}} [[x_d^{(i)} = \ell]] + \eta_{d,k} = 0 \quad \text{for } \forall \ell \in [L]$$

$$\sum_{\ell=1}^L \alpha_{d,k,\ell} = 1$$

$$\Rightarrow \alpha_{d,k,\ell} = \frac{\sum_{i=1}^n P(z^{(i)}=k | x^{(i)}) [[x_d^{(i)} = \ell]]}{\sum_{i=1}^n P(z^{(i)}=k | x^{(i)}) [[x_d^{(i)} \text{ is not missing}]]}$$