

Derivation of EM algorithm

• Notation

* Data: $X = \{x_1, \dots, x_n\}$, $x_i \in \mathbb{R}^d \quad \forall 1 \leq i \leq n$

$$x_i = \{x_i^1, \dots, x_i^d\} \quad x_i^j \in \mathbb{R} \quad \forall 1 \leq j \leq d$$

* $\Pi = \{\pi_1, \dots, \pi_K\}$ $\pi_k \in \mathbb{R}$, $0 \leq \pi_k \leq 1 \quad \forall 1 \leq k \leq K$

such that,

$$\sum_{k=1}^K \pi_k = 1$$

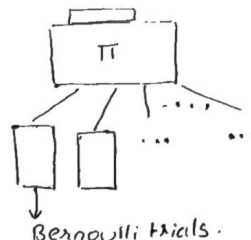


Fig. Schematic of the model

* $\Theta = \begin{bmatrix} p_1^1 & \dots & p_1^d \\ \vdots & & \vdots \\ p_K^1 & \dots & p_K^d \end{bmatrix} \quad \Theta \in \mathbb{R}^{K \times d}$
 $0 \leq p_k^j \leq 1 \quad \forall 1 \leq j \leq d, 1 \leq k \leq K$

Assume $\{x_1, \dots, x_n\}$ drawn from i.i.d Bernoulli trials (each x_i^j is a Bernoulli trial)

• Likelihood

$$L(\Theta; X) = \prod_{i=1}^n P(x_i; \Theta) \quad (\text{By i.i.d})$$

$$= \prod_{i=1}^n \left[\sum_{k=1}^K \pi_k P(x_i; \Theta_k) \right] \quad (\text{By Lm of total probability and } \Theta_k \text{ is } k^{\text{th}} \text{ row of } \Theta)$$

$$= \prod_{i=1}^n \left[\sum_{k=1}^K \pi_k \left(\prod_{j=1}^d (p_k^j)^{x_i^j} (1-p_k^j)^{1-x_i^j} \right) \right]$$

[Note: $P(x_i; \Theta_k) = \prod_{j=1}^d (p_k^j)^{x_i^j} (1-p_k^j)^{1-x_i^j}$ [Bernoulli trials]]

* Introduce λ_{ik}^j $1 \leq i \leq n, 1 \leq k \leq K$

$$\log L(\Theta; X) = \sum_{i=1}^n \log \left[\sum_{k=1}^K \lambda_{ik}^j \left(\frac{\pi_k \prod_{j=1}^d (p_k^j)^{x_i^j} (1-p_k^j)^{1-x_i^j}}{\lambda_{ik}^j} \right) \right]$$

* log is concave, applying Jensen's inequality.

$$\begin{aligned} \text{mod-log-} L(\Theta; X) &= \sum_{i=1}^n \sum_{k=1}^K \lambda_{ik}^j \log \left(\frac{\pi_k P(x_i; \Theta_k)}{\lambda_{ik}^j} \right) \\ &= \sum_{i=1}^n \sum_{k=1}^K \lambda_{ik}^j \left[\log \left(\frac{\pi_k}{\lambda_{ik}^j} \right) + \sum_{j=1}^d (x_i^j \log p_k^j + (1-x_i^j) \log (1-p_k^j)) \right] \end{aligned}$$

Here

$$\log L(\Theta) \geq \text{mod-log-} L(\Theta, \lambda) \quad (\text{since log is concave})$$

Note: constant on λ is $\sum_{k=1}^K \lambda_{ik}^j = 1 \quad \forall i \quad 0 \leq \lambda_{ik}^j \leq 1 \quad \forall k$

• Obtaining Maximum Likelihood estimators (θ, π)

$$\frac{\partial}{\partial p_k^j} \text{mod} \cdot \log L(\theta) = \sum_{i=1}^n \lambda_k^i \left(\frac{x_i^j}{p_k^j} - \frac{(1-x_i^j)}{1-p_k^j} \right) = 0$$

$$\therefore \sum_{i=1}^n (\lambda_k^i (x_i^j - x_i^j p_k^j - p_k^j + p_k^j x_i^j)) = 0$$

$$\Rightarrow \boxed{p_k^j = \frac{\sum_{i=1}^n \lambda_k^i x_i^j}{\sum_{i=1}^n \lambda_k^i}} \quad \forall 1 \leq j \leq d, 1 \leq k \leq K \quad \text{--- (1)}$$

Now,

$$\begin{aligned} \max_{\pi_1, \dots, \pi_K} \text{mod} \log L(\theta) &= \min_{\pi \in \mathbb{R}^K} - \sum_{i=1}^n \sum_{k=1}^K \lambda_k^i \left(\log\left(\frac{\pi_k}{\lambda_k^i}\right) + \log P(x_i; \theta_k) \right) \\ \text{s.t. } \sum_{k=1}^K \pi_k &= 1 \end{aligned}$$

Reparametrize $\pi_k = e^{\gamma_k}, \sum e^{\gamma_k} = 1$

$$\nabla(-\text{mod} \log L(\theta))_k = -\mathcal{L} e^{\gamma_k} \quad (\mathcal{L} \rightarrow \text{Lagrange multiplier})$$

$$\therefore \sum_{i=1}^n \lambda_k^i = \mathcal{L} e^{\gamma_k} \Rightarrow \pi_k = \frac{\sum_{i=1}^n \lambda_k^i}{\mathcal{L}}$$

$$\text{now } \sum_{k=1}^K \pi_k = 1 = \frac{\sum_{i=1}^n \sum_{k=1}^K (\lambda_k^i)}{\mathcal{L}} = \frac{n}{\mathcal{L}} \Rightarrow n = \mathcal{L}$$

$$\therefore \boxed{\pi_k = \frac{\sum_{i=1}^n \lambda_k^i}{n}} \quad \forall 1 \leq k \leq K \quad \text{--- (2)}$$

Now Maximize $\text{mod} \log L(\theta)$ over λ with (θ, π) as constants.

$$\begin{aligned} \max_{\lambda_1^i, \dots, \lambda_K^i} \sum_{k=1}^K \lambda_k^i \left(\log\left(\frac{\pi_k}{\lambda_k^i}\right) + \log(P(x_i; \theta_k)) \right) &= \min_{\lambda_1^i, \dots, \lambda_K^i} - \sum_{k=1}^K \lambda_k^i \left(\log\left(\frac{\pi_k}{\lambda_k^i}\right) + \log(P(x_i; \theta_k)) \right) \\ \text{s.t. } \sum_{i=1}^n \lambda_k^i &= 1 \end{aligned}$$

Reparametrize $\lambda_k^i = e^{\gamma_k}$

$$\nabla \left(- \sum_{k=1}^K \lambda_k^i \left(\log\left(\frac{\pi_k}{\lambda_k^i}\right) + \log(P(x_i; \theta_k)) \right) \right)_k = -\mathcal{L} e^{\gamma_k} \quad (\mathcal{L} \text{ is lagrange multiplier})$$

$$\therefore e^{\gamma_k} (\log \pi_k + \log(P(x_i; \theta_k)) - \gamma_k - 1) = \mathcal{L} e^{\gamma_k}$$

$$\Rightarrow \gamma_k = \log \pi_k P(x_i; \theta_k) - \mathcal{L} - 1 = \log \lambda_k^i$$

$$\therefore (e^{-\mathcal{L}-1}) \pi_k P(x_i; \theta_k) = \lambda_k^i \quad (\text{take exp on both sides})$$

$$\text{Now } \sum_{k=1}^K \lambda_k^i = 1 \therefore e^{-\mathcal{L}-1} = \frac{1}{\sum_{k=1}^K \pi_k P(x_i; \theta_k)} \Rightarrow \boxed{\lambda_k^i = \frac{\pi_k P(x_i; \theta_k)}{\sum_{k=1}^K \pi_k P(x_i; \theta_k)}} \quad \text{--- (3)}$$