

*Student Name:* Deeksha Arora  
*Roll Number:* 20111017  
*Date:* December 19, 2020

---

Let  $\mathbf{v} \in \mathbb{R}^N$  is the eigen vector of the covariance matrix  $\frac{1}{N}\mathbf{X}\mathbf{X}^T$ , where  $\mathbf{X}$  is a  $N \times D$  matrix ( $D > N$ ).

We know that,  $\mathbf{A}\mathbf{v} = \lambda\mathbf{v}$ , where  $\mathbf{v}$  is the eigen vector of matrix  $\mathbf{A}$ . Therefore,

$$\left(\frac{1}{N}\mathbf{X}\mathbf{X}^T\right)\mathbf{v} = \lambda\mathbf{v} \quad (1)$$

Multiplying by  $\mathbf{X}^T$  we get:

$$\mathbf{X}^T \left(\frac{1}{N}\mathbf{X}\mathbf{X}^T\right)\mathbf{v} = \lambda\mathbf{X}^T\mathbf{v} \quad (2)$$

$$\left(\frac{1}{N}\mathbf{X}^T\mathbf{X}\right)(\mathbf{X}^T\mathbf{v}) = \lambda(\mathbf{X}^T\mathbf{v}) \quad (3)$$

Let  $\mathbf{u} = \mathbf{X}^T\mathbf{v}$ . Therefore,

$$\left(\frac{1}{N}\mathbf{X}^T\mathbf{X}\right)\mathbf{u} = \lambda\mathbf{u} \quad (4)$$

So,  $\mathbf{u} = \mathbf{X}^T\mathbf{v}$  is an eigen vector of the matrix  $\frac{1}{N}\mathbf{X}^T\mathbf{X}$ . So for every eigen vector  $\mathbf{v}$  of  $\frac{1}{N}\mathbf{X}\mathbf{X}^T$  there exists an eigen vector  $\mathbf{u} = \mathbf{X}^T\mathbf{v}$  for the matrix  $\frac{1}{N}\mathbf{X}^T\mathbf{X}$ .

Therefore, instead of doing eigen decomposition of  $\frac{1}{N}\mathbf{X}^T\mathbf{X}$ ,  $D \times D$  matrix, we can find eigen vectors by doing eigen decomposition of  $\frac{1}{N}\mathbf{X}\mathbf{X}^T$  which is a  $N \times N$  matrix. This way of obtaining eigen vectors is advantageous because  $D > N$  and hence it is computationally more efficient to find eigen vectors of  $N \times N$  matrix instead of  $D \times D$  matrix.

Student Name: Deeksha Arora  
 Roll Number: 20111017  
 Date: December 19, 2020

Name - Deeksha Arora  
 Roll No - 20111017

Q2. Poisson distribution has form  $p(k|\lambda) = \frac{e^{-\lambda} \lambda^k}{k!}$

Therefore,  $Poisson(k_{n,m}|\lambda_e) = \frac{e^{-\lambda_e} \lambda_e^{k_{n,m}}}{(k_{n,m})!} \quad \text{--- (1)}$

$p(z_n=l)$  tells the probability that web server  $n$  belongs to cluster  $l$ . Therefore,

$$p(z_n=l) = \pi_e \quad \text{--- (2)}$$

Complete data log likelihood is given by:

$$p(\vec{k}, \vec{z} | \vec{\lambda}, \vec{\pi}) = \prod_{n=1}^N \prod_{l=1}^L \left[ p(z_n=l) \prod_{m=1}^M \text{Poisson}(k_{n,m}|\lambda_e) \right]^{1[z_n=l]}$$

Substituting values from eq (1) & (2)

$$p(\vec{k}, \vec{z} | \vec{\lambda}, \vec{\pi}) = \prod_{n=1}^N \prod_{l=1}^L \left[ \pi_e \prod_{m=1}^M \frac{e^{-\lambda_e} \lambda_e^{k_{n,m}}}{(k_{n,m})!} \right]^{1[z_n=l]}$$

Part 1: Complete Data Log Likelihood for the model is:

$$\log p(\vec{k}, \vec{z} | \vec{\lambda}, \vec{\pi}) = \sum_{n=1}^N \sum_{l=1}^L z_{n,l} \left[ \log \pi_e + \sum_{m=1}^M \{-\lambda_e + k_{n,m} \log \lambda_e - \log(k_{n,m}!) \} \right]$$

$z_{n,l} = 1$ , if web server  $n$  belongs to cluster  $l$ , else 0.

Name- Deeksha Aurora  
Roll No - 20111017

classmate

Date \_\_\_\_\_  
Page \_\_\_\_\_

### Part 2: Estimation Step

$$E[z_{ne}] = Y_{ne} = p(z_n = e | \vec{R}_n, \vec{\pi}, \vec{\lambda})$$

$$\propto p(z_n = e | \vec{\pi}) p(\vec{R}_n | \vec{\lambda})$$

$$\propto \pi_e \prod_{m=1}^M \frac{e^{-\lambda_e} \lambda_e^{k_{n,m}}}{(k_{n,m})!}$$

Part 3: Maximization Step : update  $\theta$  by maximizing the expected complete data log-likelihood

$$\hat{\theta} = \arg \max_{\theta} E_{p(z| \vec{R}_n, \vec{\pi}, \vec{\lambda})} \log (p(\vec{R}, \vec{z} | \vec{x}, \vec{\pi}))$$

$$= \arg \max_{\theta} \sum_{n=1}^N E_{p(z_n | \vec{R}_n, \vec{\pi}, \vec{\lambda})} \log (p(\vec{R}_n, \vec{z}_n | \vec{x}, \vec{\pi}))$$

(4)

Differentiating the above equation w.r.t.  
 ~~$\pi_e$~~  and equating it to 0 gives:

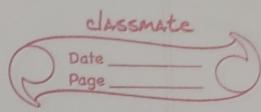
$$\pi_e = \frac{\sum_{n=1}^N E[z_{ne}]}{N}$$

$$\text{Let } Y_{ne} = E[z_{ne}] \quad \text{and } N_e = \sum_{n=1}^N Y_{ne}$$

$$\therefore \boxed{\pi_e = \frac{N_e}{N}}$$

Therefore, the MLE estimate for the mixing

Name - Deepika Aurora  
Roll No - 20111017



proportions are given by the ratio of sample sizes from each cluster

⇒ Differentiating eq ④ w.r.t.  $\lambda_e$  and equating to 0 gives :-

$$\frac{\partial \log p(\vec{X}, \vec{Z} | \vec{\lambda}, \vec{\pi})}{\partial \lambda_e} = \sum_{n=1}^N \left[ E[z_{ne}] \sum_{m=1}^M (-1) + \sum_{m=1}^M \frac{k_{n,m}}{\lambda_e} \right]$$

$$\therefore \lambda_e = \frac{\sum_{n=1}^N \sum_{m=1}^M k_{n,m}}{M \sum_{n=1}^N E[z_{ne}]}$$

$$\text{Since } Y_{ne} = E[z_{ne}] \text{ and } N_e = \sum_{n=1}^N Y_{ne}$$

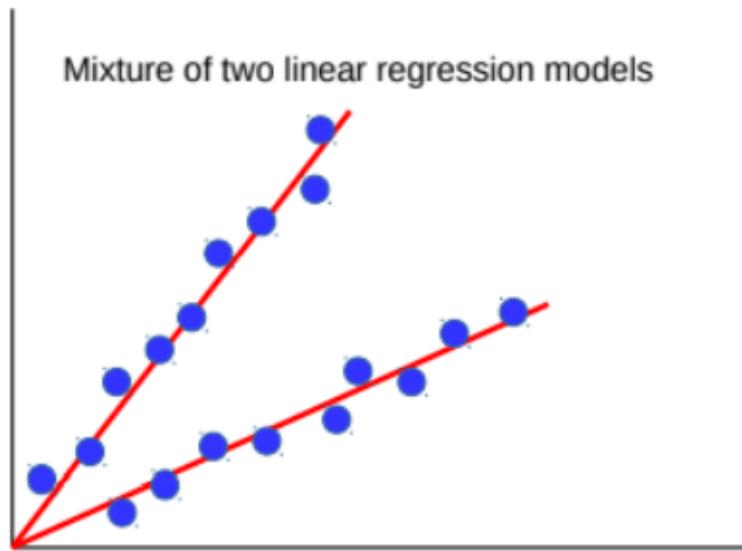
$$\therefore \lambda_e = \frac{\sum_{n=1}^N \sum_{m=1}^M k_{n,m}}{M \cdot N_e}$$

$$\boxed{\lambda_e = \frac{\sum_{n=1}^N k_n}{M \cdot N_e}}, \text{ where } k_n = \{k_{n1}, \dots, k_{nM}\}$$

Therefore,  $\lambda_e$  for cluster 'l' is the sample average for cluster l.

This model will learn more than one weight vector to model the given data. In many cases a single linear curve may not fit the data efficiently and we need more than one linear curve to model the data. In such cases this model will group the data into different clusters and will learn a different linear model for data present within each cluster. At test time, first the cluster of the test input is identified and then prediction is made using the linear model of that class. This method also reduces the impact of outliers on the linear curve because outliers may get separate out due to clustering. Hence this model is expected to work better than a standard probabilistic model.

The following image shows a possible scenario where latent variable model for regression is useful:



#### EM Algorithm:

Given :  $z_n \sim \text{multinoulli}(\pi_1, \dots, \pi_K)$ ,  $\mathbf{x}_n \sim \mathcal{N}(\mu_{z_n}, \Sigma_{z_n})$  and  $y_n \sim \mathcal{N}(\mathbf{w}_{z_n}^\top \mathbf{x}_n, \beta^{-1})$ .  
Therefore for the latent variable model:

**Prior Probability,**  $p(z_n = k) = \pi_k$

**Likelihood,**  $p(y_n | z_n, \theta) = \mathcal{N}(\mathbf{w}_{z_n}^\top \mathbf{x}_n, \beta^{-1})$  and  $p(\mathbf{x}_n | z_n, \theta) = \mathcal{N}(\mu_{z_n}, \Sigma_{z_n})$   
Here, we have two likelihood terms unlike GMM.

**Conditional Posterior Probability,**

$$p(z_n = k | \mathbf{x}_n, y_n, \theta) = p(z_n = k) p(\mathbf{x}_n | z_n = k) p(y_n | z_n = k, \mathbf{x}_n) \quad (5)$$

$$= \pi_k \mathcal{N}(\mathbf{x}_n | \boldsymbol{\mu}_k, \Sigma_k) \mathcal{N}(y_n | w_n^\top \mathbf{x}_n, \beta^{-1}) \quad (6)$$

**Computing  $\mathbb{E}[z_{nk}]$  :**

$$\mathbb{E}[z_{nk}] \propto \pi_k p(\mathbf{x}_n | z_n = k) p(y_n | z_n = k, \mathbf{x}_n) \quad (7)$$

$$\propto \pi_k \mathcal{N}(\mathbf{x}_n | \boldsymbol{\mu}_k, \Sigma_k) \mathcal{N}(y_n | w_n^\top \mathbf{x}_n, \beta^{-1}) \quad (8)$$

$$= \frac{\pi_k \mathcal{N}(\boldsymbol{\mu}_k, \Sigma_k) \mathcal{N}(\mathbf{w}_k^\top \mathbf{x}_n, \beta^{-1})}{\sum_{l=1}^K \pi_l \mathcal{N}(\boldsymbol{\mu}_k, \Sigma_k) \mathcal{N}(\mathbf{w}_l^\top \mathbf{x}_n, \beta^{-1})} \quad (9)$$

Therefore,

$$\mathbb{E}[z_{nk}] = \frac{\pi_k \exp\left(\frac{-1}{2} (\mathbf{x}_n - \boldsymbol{\mu}_k)^\top \Sigma_k (\mathbf{x}_n - \boldsymbol{\mu}_k)\right) \exp\left(\frac{-1}{2} (y_n - \mathbf{w}_k^\top \mathbf{x}_n)^\top \beta (y_n - \mathbf{w}_k^\top \mathbf{x}_n)\right)}{\sum_{l=1}^K \pi_l \exp\left(\frac{-1}{2} (\mathbf{x}_n - \boldsymbol{\mu}_k)^\top \Sigma_l (\mathbf{x}_n - \boldsymbol{\mu}_l)\right) \exp\left(\frac{-1}{2} (y_n - \mathbf{w}_l^\top \mathbf{x}_n)^\top \beta (y_n - \mathbf{w}_l^\top \mathbf{x}_n)\right)} \quad (10)$$

**Complete data log likelihood (CLL) :**

$$\log p(X, Y, Z | \Theta) = \log \prod_{n=1}^N p(\mathbf{x}_n, y_n, z_n | \theta) \quad (11)$$

$$= \log \prod_{n=1}^N \prod_{k=1}^K [p(\mathbf{x}_n | z_n = k) p(y_n | z_n = k, \mathbf{x}_n) p(z_n = k)]^{z_{nk}} \quad (12)$$

$$= \sum_{n=1}^N \sum_{k=1}^K z_{nk} \left[ -\frac{1}{2} \log |\Sigma_k| - \frac{1}{2} (\mathbf{x}_n - \boldsymbol{\mu}_k)^\top \Sigma_k^{-1} (\mathbf{x}_n - \boldsymbol{\mu}_k) - \frac{\beta}{2} (y_n - \mathbf{w}_k^\top \mathbf{x}_n)^2 + \log \pi_k \right] \quad (13)$$

Therefore, the expression for CLL (after ignoring constants) is given by equation 12.

### The EM Algorithm

1. Initialize  $\Theta$  as  $\Theta^{(0)}$ , set  $t=1$ .

2. **The Expectation Step :** Compute the expectation of each  $z_n$  using equation 8.

$$\begin{aligned} \mathbb{E}[z_{nk}^{(t)}] &= \frac{\pi_k^{(t-1)} \mathcal{N}(\boldsymbol{\mu}_k^{(t-1)}, \Sigma_k^{(t-1)}) \mathcal{N}(\mathbf{x}_n^\top \mathbf{w}_k^{(t-1)}, \beta^{-1})}{\sum_{l=1}^K \pi_l^{(t-1)} \mathcal{N}(\boldsymbol{\mu}_k^{(t-1)}, \Sigma_k^{(t-1)}) \mathcal{N}(\mathbf{x}_n^\top \mathbf{w}_l^{(t-1)}, \beta^{-1})} \\ &= \frac{\pi_k^{(t-1)} \exp\left(\frac{-1}{2} (\mathbf{x}_n - \boldsymbol{\mu}_k^{(t-1)})^\top \Sigma_k^{(t-1)} (\mathbf{x}_n - \boldsymbol{\mu}_k^{(t-1)})\right) \exp\left(\frac{-1}{2} (y_n - \mathbf{w}_k^{(t-1)T} \mathbf{x}_n)^\top \beta (y_n - \mathbf{w}_k^{(t-1)T} \mathbf{x}_n)\right)}{\sum_{l=1}^K \pi_l^{(t-1)} \exp\left(\frac{-1}{2} (\mathbf{x}_n - \boldsymbol{\mu}_k^{(t-1)})^\top \Sigma_l^{(t-1)} (\mathbf{x}_n - \boldsymbol{\mu}_l^{(t-1)})\right) \exp\left(\frac{-1}{2} (y_n - \mathbf{w}_l^{(t-1)T} \mathbf{x}_n)^\top \beta (y_n - \mathbf{w}_l^{(t-1)T} \mathbf{x}_n)\right)} \end{aligned}$$

3. **The Maximization Step :** Update  $\Theta$  by maximizing the Complete data log likelihood (CLL)

$$\hat{\Theta} = \operatorname{argmax}_{\Theta} \mathbb{E}_{p(\mathbf{Z}|\mathbf{X}, \mathbf{Y}, \hat{\Theta})} [\log p(\mathbf{X}, \mathbf{Y}, \mathbf{Z} | \Theta)] \quad (14)$$

$$= \operatorname{argmax}_{\Theta} \sum_{n=1}^N \mathbb{E}_{p(z_n | \mathbf{x}_n, y_n, \hat{\Theta})} [\log p(\mathbf{x}_n, y_n, z_n | \Theta)] \quad (15)$$

Differentiating the above equation gives the following MLE estimate of  $\Theta$ :

Suppose  $\gamma_{nk} = \mathbb{E}[z_{nk}]$ , and  $N_k = \sum_{n=1}^N \gamma_{nk}$

$$\pi_k = N_k/N \quad (16)$$

$$\boldsymbol{\mu}_k = \frac{1}{N_k} \sum_{n=1}^N \gamma_{nk} \mathbf{x}_n \quad (17)$$

$$\boldsymbol{\Sigma}_k = \frac{1}{N_K} \sum_{n=1}^N \gamma_{nk} (\mathbf{x}_n - \boldsymbol{\mu}_k) (\mathbf{x}_n - \boldsymbol{\mu}_k)^T \quad (18)$$

$$\hat{\mathbf{w}}_k = \left( \sum_{n=1}^N \gamma_{nk} \mathbf{x}_n \mathbf{x}_n^\top \right)^{-1} \left( \sum_{n=1}^N \gamma_{nk} y_n \mathbf{x}_n \right) \quad (19)$$

4. Set  $t=t+1$  and go to step 2 if not yet converged.

#### Intuitive sense of update equation of each weight vector:

In the update expression of  $\mathbf{w}_k$ , each training input is multiplied by  $\gamma_{nk}$  which denotes the expectation of training input  $n$  belonging to cluster  $k$ . Therefore, the contribution of each input in the updation of weight vector of class  $k$  depends on the probability of input  $n$  belonging to class  $k$  and logically this makes sense because the linear model of a class will depend more on the values of points belonging to the class.

#### The ALT-OPT Algorithm :

In Alt-Opt algorithm, we assume a "hard" (most probable) guess of  $z_n$ .

1. Initialize  $\Theta$  as  $\hat{\Theta}$ .
2. For  $n = 1, \dots, N$  compute the most probable value of  $z_n$ . Since  $\pi_k = \frac{1}{K}$  for all  $k \in \{1, \dots, K\}$ , therefore, we don't need to compute  $\pi_k$ . Thus, the update equation for  $z_n$  can be written as :

$$\hat{z}_n = \operatorname{argmin}_{k=1, \dots, K} [-\log p(\mathbf{x}_n | z_n = k) - \log p(y_n | z_n = k, \mathbf{x}_n)] \quad (20)$$

$$= \operatorname{argmin}_{k=1, \dots, K} \left[ \frac{1}{2} \log |\boldsymbol{\Sigma}| + \frac{1}{2} (\mathbf{x}_n - \boldsymbol{\mu}_k)^T \boldsymbol{\Sigma}_k^{-1} (\mathbf{x}_n - \boldsymbol{\mu}_k) + \frac{1}{2} (y_n - \mathbf{w}_k^T \mathbf{x}_n)^2 \right] \quad (21)$$

The optimal values of  $z_n$  can be used to create one-hot vector  $\gamma_n = [\gamma_{n1}, \dots, \gamma_{nk}]$

3. Re-estimate  $\Theta$  using MLE, given  $\mathbf{Z} = \{\hat{z}_1, \dots, \hat{z}_n\}$ ,

$$\hat{\Theta} = \operatorname{argmax}_{\Theta} \log p(\mathbf{X}, \mathbf{Y}, \hat{\mathbf{Z}} | \Theta) \quad (22)$$

$$= \operatorname{argmax}_{\Theta} \sum_{n=1}^N \sum_{k=1}^K \hat{z}_{nk} [\log p(\mathbf{x}_n | z_n = k) + \log p(y_n | z_n = k, \mathbf{x}_n)] \quad (23)$$

The updates will be same as EM except that here  $\gamma_n$  is a one-hot vector and only the points with  $\gamma_{nk} = 1$  will contribute to the updates. Therefore,

$$\boldsymbol{\mu}_k = \frac{1}{N_k} \sum_{n=1}^N \gamma_{nk} \mathbf{x}_n \quad (24)$$

$$\boldsymbol{\Sigma}_k = \frac{1}{N_K} \sum_{n=1}^N \gamma_{nk} (\mathbf{x}_n - \boldsymbol{\mu}_k) (\mathbf{x}_n - \boldsymbol{\mu}_k)^T \quad (25)$$

$$\hat{\mathbf{w}}_k = \left( \sum_{n=1}^N \gamma_{nk} \mathbf{x}_n \mathbf{x}_n^\top \right)^{-1} \left( \sum_{n=1}^N \gamma_{nk} y_n \mathbf{x}_n \right) \quad (26)$$

4. Go to step 2 if not yet converged