

Student Name: Sakshi  
Roll Number: 180653  
Date: December 19, 2020

---

We've to find the eigenvectors of  $\mathbf{S}$ :

$$\mathbf{S} = \frac{1}{N} \mathbf{X}^T \mathbf{X}$$

And we're given the eigenvectors of  $\mathbf{R}$ :

$$\mathbf{R} = \frac{1}{N} \mathbf{X} \mathbf{X}^T$$

Let's look at the singular value decomposition of  $\frac{1}{\sqrt{N}} \mathbf{X}$ :

$$\frac{1}{\sqrt{N}} \mathbf{X} = \mathbf{U} \mathbf{\Sigma} \mathbf{V}^T = \sum_{n=1}^N \sigma_n \mathbf{u}_n \mathbf{v}_n^T$$

Here,  $\mathbf{U}$  &  $\mathbf{V}$  are  $N \times N$  and  $D \times D$  orthogonal matrices respectively.  $\mathbf{\Sigma}$  is  $N \times D$  diagonal matrix of singular values. Thus,  $\mathbf{S}$  and  $\mathbf{R}$  can be written as:

$$\mathbf{S} = \frac{1}{N} \mathbf{X}^T \mathbf{X} = \mathbf{V} \mathbf{\Sigma}^T \mathbf{\Sigma} \mathbf{V}^T = \mathbf{V} \mathbf{\Sigma}_S \mathbf{V}^T$$

$$\mathbf{R} = \frac{1}{N} \mathbf{X} \mathbf{X}^T = \mathbf{U} \mathbf{\Sigma} \mathbf{\Sigma}^T \mathbf{U}^T = \mathbf{U} \mathbf{\Sigma}_R \mathbf{U}^T$$

It can be seen that  $\mathbf{\Sigma}_R$  and  $\mathbf{\Sigma}_S$  have equal non-zero entries. Thus, every column of  $\mathbf{V}$  is an eigenvector of  $\mathbf{S}$  and every column of  $\mathbf{U}$  is an eigenvector of  $\mathbf{R}$ . We have:

$$\begin{aligned} (\mathbf{v}'_i)^T &= \frac{1}{\sqrt{N}} \mathbf{u}_i^T \mathbf{X} \\ &= \sum_{n=1}^N \sigma_n (\mathbf{u}_i^T \mathbf{u}_n) \mathbf{v}_n^T \\ &= \sigma_i \mathbf{v}_i^T \end{aligned}$$

$$\mathbf{v}_i = \frac{\mathbf{v}'_i}{\|\mathbf{v}'_i\|}$$

Thus, we've have the eigenvector of  $\mathbf{S}$  using that of  $\mathbf{R}$ . It takes  $\mathcal{O}(D^2)$  complexity in power method. Here it takes  $\mathcal{O}(ND)$  and is better as  $N < D$ .

Student Name: Sakshi  
 Roll Number: 180653  
 Date: December 19, 2020

**Estimation Maximization for Poisson Mixture Model:**

1. Initialize:  $\Theta = \{\pi_l, \lambda_l\}_{l=1}^L$  as  $\Theta^{(0)}$  with  $t = 1$
2. Next is the Estimation step in which we will compute the conditional posterior  $\mathbf{p}(\mathbf{Z}|\mathbf{K}, \Theta^{(t-1)})$   
 As the observations are independent and identically distributed :

$$\begin{aligned} \mathbf{p}(\mathbf{z}_n = l | \mathbf{k}_n, \Theta^{(t-1)}) &\propto \mathbf{p}(\mathbf{z}_n = l, \Theta^{(t-1)}) \mathbf{p}(\mathbf{k}_n | \mathbf{z}_n = l, \Theta^{(t-1)}) \\ \implies \mathbf{p}(\mathbf{z}_n = l | \mathbf{k}_n, \Theta^{(t-1)}) &= \pi_l^{(t-1)} \prod_{m=1}^M \text{Poisson}(\mathbf{k}_{n,m} | \lambda_l) \\ &= \mathbf{p}(\mathbf{z}_n = l | \mathbf{k}_n, \Theta^{(t-1)}) = \pi_l^{(t-1)} \prod_{m=1}^M \frac{(\lambda_l^{(t-1)})^{(\mathbf{k}_{n,m})} e^{-(\lambda_l^{(t-1)})}}{\mathbf{k}_{n,m}} \end{aligned}$$

Hence,  $\gamma_{nl}^{(t)} = \frac{\pi_l^{(t-1)} \prod_{m=1}^M \text{Poisson}(\mathbf{k}_{n,m} | \lambda_l)}{\sum_{i=1}^L \pi_i^{(t-1)} \prod_{m=1}^M \text{Poisson}(\mathbf{k}_{n,m} | \lambda_i)}$

3. Estimation is followed by Maximization step which maximizes the expected complete data log likelihood. Update equation is as follows:

$$\begin{aligned} \Theta^{(t)} &= \underset{\Theta}{\operatorname{argmax}} \mathbf{E}_{\mathbf{p}(\mathbf{Z}^{(t-1)} | \mathbf{K}, \Theta^{(t-1)})} \left[ \log \mathbf{p}(\mathbf{Z}^{(t-1)}, \mathbf{K} | \Theta) \right] \\ &= \underset{\Theta}{\operatorname{argmax}} \sum_{n=1}^N \mathbf{E}_{\mathbf{p}(\mathbf{Z}^{(t-1)} | \mathbf{K}, \Theta^{(t-1)})} \left[ \log \mathbf{p}(\mathbf{z}_n^{(t-1)}, \mathbf{k}_n | \Theta) \right] \\ &= \underset{\Theta}{\operatorname{argmax}} \sum_{n=1}^N \mathbf{E}_{\mathbf{p}(\mathbf{Z}^{(t-1)} | \mathbf{K}, \Theta^{(t-1)})} \left[ \sum_{l=1}^L \log \mathbf{p}(\mathbf{z}_n = l | \Theta)^{\mathbf{z}_{nl}^{t-1}} + \log \mathbf{p}(\mathbf{k}_n | \mathbf{z}_n = l, \Theta)^{\mathbf{z}_{nl}^{(t-1)}} \right] \\ &= \underset{\Theta}{\operatorname{argmax}} \sum_{n=1}^N \sum_{l=1}^L \mathbf{E} \left[ \mathbf{z}_{nl}^{t-1} \right] \left[ \log \mathbf{p}(\mathbf{z}_n = l | \Theta) + \log \mathbf{p}(\mathbf{k}_n | \mathbf{z}_n = l, \Theta) \right] \\ &= \underset{\Theta}{\operatorname{argmax}} \sum_{n=1}^N \sum_{l=1}^L \gamma_{nl}^{(t)} \left[ \log \pi_l + \sum_{m=1}^M \log \mathbf{p}(\mathbf{k}_{n,m} | \mathbf{z}_n = l, \Theta) \right] \\ &= \underset{\Theta}{\operatorname{argmax}} \sum_{n=1}^N \sum_{l=1}^L \gamma_{nl}^{(t)} \left[ \log \pi_l + \left( \sum_{m=1}^M \mathbf{k}_{n,m} \right) \log \lambda_l - M \lambda_l \right] \end{aligned}$$

Optimization constraints:  $\sum_{l=1}^L \pi_l = 1$  and  $\lambda_l > 0, l = 1, 2, \dots, L$ . Using Lagrangian operator, we get the required updates:  $\pi_l^{(t)} = \frac{N_l}{N}$  &  $\lambda_l^{(t)} = \frac{1}{MN_l} \sum_{n=1}^N \sum_{m=1}^M \gamma_{nl}^{(t)} k_{n,m}$

4. Iterate :  $t=t+1$  until convergence

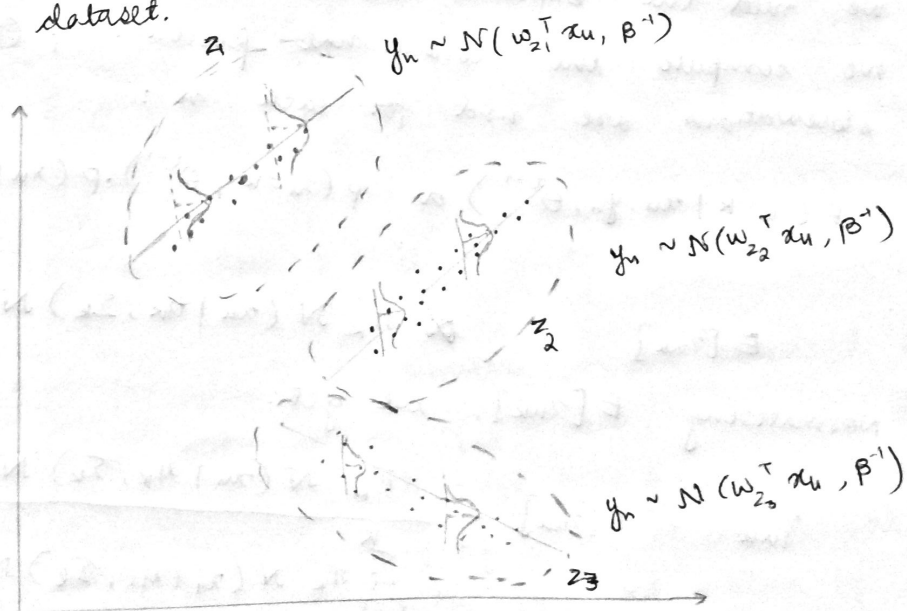
### Question 3

#### Part 1:

- The latent variable model for regression is doing linear regression and clustering simultaneously. We are learning  $K$ -linear regressions for each cluster. The inputs  $x_n$  are assumed to belong to one of the  $K$ -clusters. Each cluster  $z_n$  is modelled as a Gaussian distribution  $N(\mu_{z_n}, \Sigma_{z_n})$ . The output  $y_n$  is modelled as a cluster dependent linear transformation of input added with Gaussian noise.

$$y_n = w_{z_n}^T x_n + \epsilon; \quad \epsilon \sim N(0, \beta^{-1})$$

Thus, this model can be used for non-linearly distributed dataset. In the standard linear regression problem, a single linear transformation is used for the whole dataset. Thus we learn  $\mu$  &  $\Sigma$  for the entire dataset. Here, we first find the clusters and then learn  $\mu_{z_k}$  and  $\Sigma_{z_k}$  for each cluster. Standard linear regression works well only for linearly separable distributed dataset.



• CLL:

$$\theta = \{\mu_k, \Sigma_k, w_k\}_{k=1}^K$$

$$CLL(\theta) = \log \left[ \prod_{n=1}^N p(x_n, y_n, z_n | \theta) \right]$$

$$= \log \prod_{n=1}^N \prod_{k=1}^K \left\{ p(x_n | z_n=k, \theta) p(y_n | z_n=k, x_n, \theta) p(z_n=k) \right\}^{z_{nk}}$$

$$= \sum_{n=1}^N \sum_{k=1}^K z_{nk} \left\{ \log p(x_n | z_n=k, \theta) + \log p(y_n | z_n=k, x_n, \theta) + \log p(z_n=k | \theta) \right\}$$

\*  $\log p(x_n | z_n=k, \theta) : \mathcal{N}(x_n | \mu_k, \Sigma_k)$

\*  $\log p(y_n | z_n=k, \theta) : \mathcal{N}(y_n | w_k^T x_n, \beta^{-1})$

\*  $\log p(z_n=k) : \pi_k$

$$CLL(\theta) \propto \sum_{n=1}^N \sum_{k=1}^K z_{nk} \left\{ -\frac{1}{2} \log |\Sigma_k| - \frac{1}{2} (x_n - \mu_k)^T \Sigma_k^{-1} (x_n - \mu_k) - \frac{\beta}{2} (y_n - w_k^T x_n)^2 + \log \pi_k \right\}$$

• EM for Variable Regression

• E-step:

We need the expected CLL.

We compute the conditional posterior  $p(z | x, y, \theta^{(t-1)})$ . Since observations are iid for each  $n$ :

$$p(z_n=k | x_n, y_n, \theta^{(t-1)}) \propto p(z_n=k | \theta^{(t-1)}) \cdot p(x_n | z_n=k, \theta^{(t-1)}) \cdot p(y_n | x_n, z_n=k, \theta^{(t-1)})$$

$$\therefore E[z_{nk}] \propto \pi_k \mathcal{N}(x_n | \mu_k, \Sigma_k) \mathcal{N}(y_n | w_k^T x_n, \beta^{-1})$$

Normalizing  $E[z_{nk}]$ , we get:

$$Y_{nk} = E[z_{nk}] = \frac{\pi_k \mathcal{N}(x_n | \mu_k, \Sigma_k) \mathcal{N}(y_n | w_k^T x_n, \beta^{-1})}{\sum_{l=1}^K \pi_l \mathcal{N}(x_n | \mu_l, \Sigma_l) \mathcal{N}(y_n | w_l^T x_n, \beta^{-1})}$$

M-step: Updating  $\theta^{t-1}$  by maximizing the Expected LL.

$$\theta^t = \arg \max_{\theta} E_{P(z^{t-1} | x, y, \theta^{t-1})} [\log P(z^{t-1}, x, y | \theta)]$$

As we have already seen, this comes out to be:

$$\theta^t = \arg \max_{\theta} \sum_{n=1}^N \sum_{k=1}^K \gamma_{nk}^{(t-1)} \left[ \log \pi_k + \log N(x_n | \mu_k, \Sigma_k) - \frac{1}{2} (x_n - \mu_k^T x_n)^2 - \frac{1}{2} \log |\Sigma_k| \right]$$

This is to be done with respect to constraints:

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

As seen in class, we can use Lagrangian method to get:

$$\pi_k^{(t)} = \frac{1}{N} \sum_{n=1}^N \gamma_{nk}^{(t)} = \frac{N_k}{N} \quad N_k = \sum_{n=1}^N \gamma_{nk}^{(t)}$$

$$\Sigma_k^{(t)} = \frac{1}{N_k} \sum_{n=1}^N \gamma_{nk}^{(t)} (x_n - \mu_k^t) (x_n - \mu_k^t)^T$$

$$\mu_k^t = \frac{1}{N_k} \sum_{n=1}^N \gamma_{nk}^{(t)} x_n \quad ; \quad \hat{w}_k = \left( \sum_n \gamma_{nk} x_n x_n^T \right) \left( \sum_n \gamma_{nk} x_n \right)^{-1}$$

Thus, the overall algorithm is:

1. Initialize  $\theta = \theta^0$  ;  $t = 1$
2. Perform the E-step
3. Perform the M-step
4.  $t = t + 1$  ; Go to step 2 if not converged.

### Update Equation:

Our update equation for  $w_k$  is similar to the standard linear regression one only with extra  $\gamma_{nk}$ . Here, only those points are considered for  $k^{th}$  cluster which have  $\gamma_{nk} = 1$ .  $\gamma_{nk} = 0$  to  $\gamma_{nk} = 1$  determines the contribution of a point in a cluster. Thus, linear regression happens with expected value of points determining how much a point contributes in a cluster.



• ALT-OPT Algorithm:

1. Initialize  $\Theta = \{\mu_k, \Sigma_k, w_k\}_{k=1}^K$ ; Set  $t=1$ ;
2. Here we'll take point estimate of  $z_n$  instead of  $E[z_n]$ .  
Let's compute the most probable value of  $z_n$  as:

$$\hat{z}_n = \underset{k}{\operatorname{argmax}} \pi_k N(x_n | \mu_k^{(t-1)}, \Sigma_k^{(t-1)}) N(y_n | (w_k^{(t-1)})^T x_n, \beta^{-1})$$

Given all  $\pi_k = \frac{1}{K}$ , we can neglect it.

3. We can solve the MLE problem for  $\Theta$  using  $z_n$ .  
Here, we don't need expected value as we have the actual  $z_n$ s. So we can just replace  $\gamma_{nk}$ s by  $z_{nk}$ s, got in the ALT-OPT method.

Thus, we get:

$$\mu_k^{(t)} = \frac{1}{N} \sum_{n=1}^N z_{nk}$$

$$N_k = \sum_{n=1}^N z_{nk}$$

$$\Sigma_k^{(t)} = \frac{1}{N_k} \sum_{n=1}^N z_{nk} (x_n - \hat{\mu}_k) (x_n - \hat{\mu}_k)^T$$

$$w_k^{(t)} = \left( \sum_{n: z_{nk}=1} x_n x_n^T \right) \left( \sum_{n: z_{nk}=1} y_n x_n \right)$$

4.  $t=t+1$ ; Go to step 2 if not converged.