Solution 1

Suppose $\vec{v}$ is an eigen vector of matrix $S' = \frac{1}{N} XX^T$. By definition,     $S'v = \lambda v$ , where $\lambda$ is corresponding eigen value

$\therefore$    $\frac{1}{N}(XX^T)v = \lambda v$

Premultiplying by $X^T$ both sides, we get

$$\frac{1}{N} X^T X X^T v = \lambda(X^T v).$$

substituting $X^T v = u$    $[X^T \to D \times N, v \in R^N \therefore u \in R^D]$.

$$\frac{1}{N}(X^T X) u = \lambda(u)$$

From above expression, it is clear that

$u = \underline{X^T v}$  is eigen vector for $\frac{1}{N}(X^T X)$ ie S.

In normal way, time complexity to compute eigen vector of

S   is $O(KD^2)$ $\longrightarrow T_1$

For this way, time complexity will be

$$\underbrace{O(KN^2)}_{\substack{eigen\ vectors \\ for\ S'}} + \underbrace{O(KND)}_{\substack{Matrix \\ Multiplication}} \longrightarrow T_2$$

since, we are given $D > N$ in the question, thus $T_1 > T_2$.

i.e we can say that computing eigen vector of S through S' is more efficient than computing them directly.

Solution 2.

Given: Poisson distribution, N webservers monitored for M minutes.

$k_{n,m} \rightarrow$ no. of hits to $n^{th}$ webserver in minute $m$.

As shown in the hint, the complete data likelihood is as follows.

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

$$1[z_n = l] = \begin{cases} 1 & \text{if } z_n = l \\ 0 & \text{otherwise} \end{cases}, \quad \text{and also } p(z_n = l) = \pi_l.$$

Taking log both sides.

$$CLL = \cancel{1[\text{const}]} \sum_{n=1}^{N} \sum_{l=1}^{L} 1[z_n = l] \left[ \log(\pi_l) + \log\left( \prod_{m=1}^{M} \text{Poisson}(\cdots) \right) \right].$$

$$= \sum_{n=1}^{N} \sum_{l=1}^{L} z_{nl} \left[ \log(\pi_l) + \log\left( \prod_{m=1}^{M} \frac{1}{e^{\lambda_l}} \frac{\lambda_l^{k_{n,m}}}{(k_{m,m})!} \right) \right]$$

$$= \sum_{n=1}^{N} \sum_{l=1}^{L} z_{nl} \left[ \log(\pi_l) + \sum_{m=1}^{M} \left( k_{n,m} \log \lambda_l - \lambda_l - \log(k_{n,m}!) \right) \right]$$

Thus; the complete data log likelihood is given as.

$$\sum_{n=1}^{N} \sum_{l=1}^{L} z_{nl} \left[ \log(\pi_l) + \sum_{m=1}^{M} \left( k_{n,m} \log \lambda_l - \lambda_l - \log(k_{n,m}!) \right) \right]$$

$z_{nl} = 1$ and all other components of one hot vector $z_n = 0$.

Estimating $z_{n\ell}$. $\rightarrow$ E-step

$$E[z_{n\ell}] = 1 \times p(z_{n\ell}=1 \mid k_n, \theta) + 0 \times \cancel{p(z_{n\ell}=0 \mid k_n)}$$

$$= p(z_{n\ell}=1 \mid k_n)$$

$$\propto p(z_{n\ell}=1) \, p(k_n \mid z_{n\ell}=1, \theta)$$

$$\propto \cancel{\prod_k \mathcal{N}(\cdots)} \cdot \pi_\ell \left[ \prod_{m=1}^{M} \text{Poisson}(k_{nm} \mid \lambda_\ell) \right]$$

$$\propto \pi_\ell \left( \frac{e^{-\lambda_\ell M} \, \lambda^{\sum_{m=1}^{\ } k_{nm}}}{\prod_{m=1}^{M} (k_{nm}!)} \right)$$

c.) M-step.

Taking $E[z_{n\ell}] = \gamma_{n\ell}$.

Expected complete Data log-likelihood is.

$$EL = \sum_{n=1}^{N} \sum_{\ell=1}^{L} \gamma_{n\ell} \left[ \log(\pi_\ell) + \sum_{m=1}^{M} \left( \log(\lambda_\ell) \cdot k_{n,m} - \lambda_\ell - \log(k_{n,m}!) \right) \right]$$

Taking derivative w.r.t $\lambda_\ell$., ( $\gamma_{ni} = 0 \ \forall i \neq \ell$)

$$\frac{\partial EL}{\partial \lambda_\ell} = \sum_{n=1}^{N} \gamma_{n\ell} \left[ \cancel{\frac{\partial}{\partial \lambda_\ell} \log(\pi_\ell)} + \frac{\partial}{\partial \lambda_\ell} \left( \sum_{m=1}^{M} (k_{n,m} \log \lambda_\ell - \lambda_\ell - \log(s)) \right) \right]$$

$$= \sum_{n=1}^{N} \gamma_{n\ell} \left[ \sum_{m=1}^{M} \left( \frac{k_{n,m}}{\lambda_\ell} - 1 \right) \right]$$

Now, equating it with 0., we get

$$\sum_{n=1}^{N} \gamma_{n\ell} \left[ \sum_{m=1}^{M} \left( \frac{k_{n,m}}{\lambda_\ell} - 1 \right) \right] = 0$$

$$\Rightarrow \sum_{n=1}^{N} \left( \gamma_{n\ell} \sum_{m=1}^{M} \frac{k_{n,m}}{\lambda_\ell} - \gamma_{n\ell} M \right) = 0$$

$$\Rightarrow \sum_{n=1}^{N} \gamma_{n\ell} \sum_{m=1}^{M} \frac{k_{n,m}}{\lambda_\ell} = \sum_{n=1}^{N} \gamma_{n\ell} M$$

$\rightarrow$ where $\gamma_{n\ell} = E[z_{n\ell}]$

$$\Rightarrow \boxed{\lambda_\ell = \frac{\sum_{n=1}^{N} \gamma_{n\ell} \sum_{m=1}^{M} k_{n,m}}{\sum_{n=1}^{N} \gamma_{n\ell} M}} \quad \checkmark$$

## Estimating $\pi_\ell$

$E[CLL]$ ↙ $E.L = \sum_{n=1}^{N} \sum_{\ell=1}^{L} \gamma_{n\ell} \left[ \log \pi_\ell + \sum_{m=1}^{M} (\text{some terms}) \right]$

We need to max. above eq$^n$, but there is constraint $\sum_{\ell=1}^{L} \pi_\ell = 1$, so, we will use Langrangian method,

$$\ell = \sum_{n=1}^{N} \sum_{\ell=1}^{L} \gamma_{n\ell} \log (\pi_\ell) + \alpha * \left( 1 - \sum_{\ell=1}^{L} \pi_\ell \right). \quad \begin{bmatrix} \text{ignored} \\ \text{ind. terms} \end{bmatrix}$$

$\dfrac{\partial \ell}{\partial \pi_\ell} = \sum_{n=1}^{N} \dfrac{\gamma_{n\ell}}{\pi_\ell} - \alpha$, equating to zero gives us.

$$\sum_{n=1}^{N} \dfrac{\gamma_{n\ell}}{\pi_\ell} = \alpha \quad \Rightarrow \quad \pi_\ell = \dfrac{\sum_{n=1}^{N} \gamma_{n\ell}}{\alpha}, \quad \text{but } \alpha = ?$$

$\dfrac{\partial \ell}{\partial \alpha} = 0 \Rightarrow \sum_{\ell=1}^{L} \pi_\ell = 1 \quad \Rightarrow \quad \dfrac{\sum_{\ell=1}^{L} \sum_{n=1}^{N} \gamma_{n\ell}}{\alpha} = 1$

Hence $\alpha = \sum_{n=1}^{N} \sum_{\ell=1}^{L} \gamma_{n\ell}$, therefore

$$\boxed{\pi_\ell = \dfrac{\sum_{n=1}^{N} \gamma_{n\ell}}{\sum_{n=1}^{N} \sum_{\ell=1}^{L} \gamma_{n\ell}}} \quad ✓ \quad \underline{\text{Final Ans}}$$

Now since $\sum_{n=1}^{N} \sum_{\ell=1}^{L} \gamma_{n\ell} = N$. $\quad \therefore \quad \boxed{\pi_\ell = \dfrac{\sum_{n=1}^{N} \gamma_{n\ell}}{N}}$
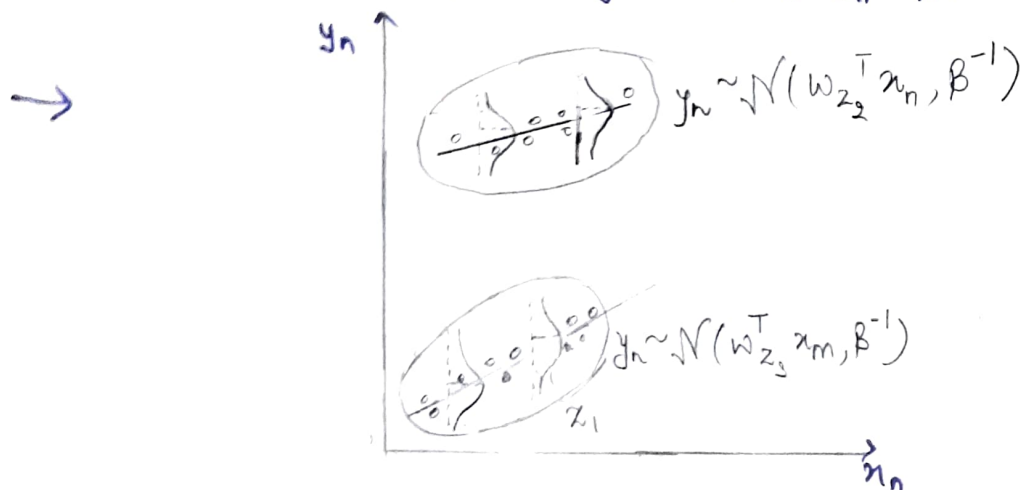
# Solution 3.  Part 1

Generative story as given in question.

(i) $z_n \sim$ multinomial $(\pi_1, \pi_2, \cdots \pi_k)$

(ii) Generate inputs $x_n \sim N(\mu_{z_n}, \Sigma_{z_n})$

(iii) Outputs $y_n \sim N(w_{z_n}^T x_n, \beta^{-1})$.



$y_n \sim N(w_{z_2}^T x_n, \beta^{-1})$

$y_n \sim N(w_{z_3}^T x_m, \beta^{-1})$

Latent variable model

The model will learn a combination of $k$-linear regressions, as depicted by the above graph. In resp. cluster, the input-output relationship is linear. However, in standard linear regression, it will only learn a single linear boundary. Also, the LVM model will separate the outliers in clustering, therefore also help in reducing their effect on predictions of model.

Now A.T.Q. $p(z_n|\theta) = $ multinoulli $(\pi_1, \cdots \pi_k)$

$$p(x_n|z_n, \theta) = N(\mu_{z_n}, \Sigma_{z_n})$$

$$p(y_n|x_n, z_n, \theta) = N(w_{z_n}^T x_n, \beta^{-1}).$$

## 3.1.2.

### Deriving EM algorithm.

~~Q~~ $= \prod$  $\quad CLL = \sum\limits_{n=1}^{N} \log p(x_n, y_n, z_n | \theta)$

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

$= \sum\limits_{n=1}^{N} \log p(y_n | x_n, z_n, \theta) + \log p(x_n | z_n, \theta) + \log (z_n | \theta).$

$= \sum\limits_{n=1}^{N} \sum\limits_{k=1}^{K} ~~z_{nk} [\log N(\mu_k, \Sigma_k) ] + \log \pi_k + \log N(y_n | w~~$

$= \sum\limits_{n=1}^{N} \sum\limits_{k=1}^{K} z_{nk} \left[ \log N(w_k^T x_n, \beta^{-1}) + \log N(\mu_k, \Sigma_k) + \log (\pi_k) \right]$

### E-step.

$p(z | x, Y, \theta) = \prod\limits_{n=1}^{N} p(z_n | x_n, y_n, \theta)$

$p(z_n | x_n, y_n, \theta) = \dfrac{p(y_n | x_n, z_n, \theta) \, p(x_n | z_n, \theta) \, p(z_n | \theta)}{p(x_n, y_n | \theta).} \rightarrow$ Independent of $z$.

$p(z_n = k | x_n, y_n, \theta) \propto N(w_k^T x_n, \beta^{-1}) \, N(\mu_k, \Sigma_k) \, \pi_k.$

$p(z_n = k | x_n, y_n, \theta) = \dfrac{N(w_k^T x_n, \beta^{-1}) \, N(\mu_k, \Sigma_k) \, \pi_k}{\sum\limits_{l=1}^{K} N(w_l^T x_n, \beta^{-1}) \, N(\mu_l, \Sigma_l) \, \pi_l.}$

$\gamma_{nk} = E[z_{nk}] = 1 \times p(z_n = k) \not= = ~~\#~~ \cdot$ same as above ("").

## M-step.

$$E[CLL] = \sum_{n=1}^{N} \sum_{k=1}^{K} E[z_{nk}] \left( \log \pi_k + \log \mathcal{N}(\mu_k, z_k) + \log \mathcal{N}(w_k^T x_n, B^{-1}) \right)$$

For maximization, we will differentiate this w.r.t to $\mu_k, \Sigma_k$ and $\pi_k$ and $w_k$.

This will be same as done in class. for GMM.

$$\hat{\pi}_k = \frac{N_k}{N}$$

$$\hat{\mu}_k = \frac{\sum_{n=1}^{N} \gamma_{nk} x_n}{N_k}$$

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

$$\hat{w}_k = \left( \sum_n \gamma_{nk} x_n x_n^T \right)^{-1} \left( \sum_n \gamma_{nk} x_n y_n \right)$$

→ shown on next page.

here $\gamma_{nk} = E[z_{nk}]$, calculated in E step

$N_k = \sum_{n=1}^{N} \gamma_{nk}$ = effective no. of pts in cluster $k$

## Overall EM algorithm

1. Initialize $\theta = \theta^0$, set $t = 1$

2. <u>E-step</u> :- Compute $E[z_{nk}]$

3. <u>M-step</u> :- Maximize $E[CLL]$ and update parameters as described above.

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

Calculation of $\hat{w_k}$,

$$\frac{\partial}{\partial w_k}\left(E[CLL]\right) = 0$$

$$E[CLL] = \sum_{n=1}^{N}\sum_{k=1}^{K} E[z_{nk}]\left(\underbrace{\log \pi_k + \log \mathcal{N}(\mu_k, \Sigma_k)}_{\text{independent of } w_k} + \log \mathcal{N}(w_k^T x_n, \beta^{-1})\right)$$

→ ③

ignoring constants in expression of ③

$$= \sum_{n=1}^{N}\sum_{k=1}^{K} \gamma_{nk}\left[\text{indep. terms} + \frac{-\beta}{2}\left(y_n - w_k^T x_n\right)^2\right]$$

$$\frac{\partial E[CLL]}{\partial w_k} = \cdot \sum_{n=1}^{N} \beta \gamma_{nk} \left(y_n - w_k^T x_n\right) x_n$$

equating to zero, we get

$$\sum_{n=1}^{N} \cdot \gamma_{nk}\left(y_n - \hat{w_k}^T x_n\right) x_n = 0$$

$$\sum_{n=1}^{N} \gamma_{nk} x_n y_n = \sum_{n=1}^{N} \gamma_{nk} x_n x_n^T \hat{w_k}$$

↓ doesn't depend on $n$.

$$\therefore \hat{w_k} = \left(\sum_{n=1}^{N} \gamma_{nk} x_n x_n^T\right)^{-1}\left(\sum_{n=1}^{N} \gamma_{nk} x_n y_n\right)$$

# Intution of update eqⁿ of $\hat{w}_k$

in normal regression $w = \left[\sum\limits_{n=1}^{N} (x_n x_n^T)\right]^{-1} \left(\sum y_n x_n\right)$.

Our update eqⁿ is similar, rather it is specific to every cluster, ie $w_k$.. Note our eqⁿ only considers points belonging to that cluster. This property is governed by $\gamma_{nk}$ in the expression

## ALT-OPT Algorithm.

Instead of $E[z_{nk}]$, we find $\hat{z}_n$

$$\hat{z}_n = \max_{k \in [1,K]} \pi_k \, \mathcal{N}(w_k^T x_n, \beta^{-1}) \, \mathcal{N}(\mu_k, \Sigma_k)$$

since $\pi_k = 1/k$

$$\hat{z}_n = \max_{k \in [1,k]} \mathcal{N}(w_k^T x_n, \beta^{-1}) \, \mathcal{N}(\mu_k, \Sigma_k)$$

Simply replace $\gamma_{nk}$ with $\hat{z}_{nk}$ in all the update parameters.

$$\hat{\mu_k} = \frac{1}{N_k} \sum_{n=1}^{N} \hat{z}_{nk} \, x_n$$

$$\hat{\Sigma_k} = \frac{1}{N_k} \sum_{n=1}^{N} \hat{z}_{nk} \, (x_n - \hat{\mu_k})(x_n - \hat{\mu_k})^T$$

$$\hat{w_k} = \left(\sum_{n=1}^{N} \hat{z}_n \, x_n x_n^T\right)^{-1} \left(\sum_{n=1}^{N} \hat{z}_n \, y_n x_n\right).$$

Note here $N_k = \sum_{n=1}^{N} \hat{z}_{nk}$.

## Overall algorithm ALT-OPT

1. Intialize $\theta = \{\overline{\phantom{xx}}, \mu_k, \Sigma_k, w_k\}$ as $\theta_0$, set $t = 1$

2. For each $n$, compute $\hat{z}_n$

$$\hat{z}_n = \underset{k \in [1,K]}{\text{argmax}} \left[ N_k(x_n \mid \mu_k, \Sigma_k) \, N(y_n \mid w_k^T x_n, \beta^{-1}) \right]$$

3. Solve MLE problem using updates given in last page. $\forall k$.

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

Part 2

$x_n \to$ given (not modeled)

$$\pi_k(x_n) \to \frac{\exp(\eta_k^T x_n)}{\sum_{l=1}^{K} e^{\eta_l^T x_n}}$$

$$CLL = \sum_{n=1}^{N} \log p(y_n, z_n | x_n, \theta)$$

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

$$= \sum_{n=1}^{N} \sum_{k=1}^{K} z_{nk} \left[ \log \mathcal{N}(y_n | w_k^T x_n, \beta^{-1}) + \log \pi_k(x_n) \right]$$

$$EL = E[CLL] = \sum_{n=1}^{N} \sum_{k=1}^{K} \gamma_{nk} \left[ -\frac{\beta}{2}(y_n - w_k^T x_n)^2 + \log \pi_k(x_n) \right]$$

E-step

Calculation of $\gamma_{nk} \to$ ?

$$p(z_n = k | x_n, y_n, \theta) = \frac{p(z_n = k | \theta, x_n)\, p(y_n | z_n = k, \theta, x_n)}{\sum_{l=1}^{K} p(z_n = l | \theta, x_n)\, p(y_n | z_n = l)}$$

$$\alpha \, \underset{\alpha\, e^{\eta_k^T x_n}}{\pi_k(x_n)} \cdot \mathcal{N}(y_n | w_k^T x_n, \beta^{-1})$$

$$p(z_n = k | x_n, y_n, \theta) = \frac{\mathcal{N}(y_n | w_k^T x_n, \beta^{-1})\, e^{\eta_k^T x_n}}{\sum_{l=1}^{K} \mathcal{N}(y_n | w_l^T x_n, \beta^{-1})\, e^{\eta_l^T x_n}}$$

$$\gamma_{nk} = E[z_n = k] = 1 \times p(z_n = k | x_n, y_n, \theta) + 0.$$

$$= p(z_n = k | x_n, y_n, \theta)$$

update for $W_k \longrightarrow$ remains same. as part 1.

$$\hat{W_k} = \left( \sum_{n=1}^{N} \gamma_{nk} \, x_n \, x_n^T \right)^{-1} \left( \sum_{n=1}^{N} \gamma_{nk} \, y_n \, x_n \right) \quad \forall \, k \in [1, K]$$

Now, since $\eta_k$ is also a parameter in $E[CLL]$

$$\frac{\partial E[CLL]}{\partial \eta_k} = 0.$$

$$\Longrightarrow \frac{\partial}{\partial \eta_k} \sum_{n=1}^{N} \sum_{k=1}^{K} \gamma_{nk} \, \log \pi_k (x_n) = 0$$

$$\Longrightarrow \frac{\partial}{\partial \eta_k} \sum_{n=1}^{N} \sum_{k=1}^{K} \gamma_{nk} \left( \eta_k^T x_n - \log \sum_{l=1}^{K} e^{\eta_l^T x_n} \right) = 0$$

$$\hookrightarrow \frac{\partial}{\partial \eta_k} \sum_{n=1}^{N} \cancel{\sum_k} \; \gamma_{nk} \left( \eta_k^T x_n - \underbrace{\log \sum_{l=1}^{K} e^{\eta_l^T x_n}}_{\text{log of sum}} \right) = 0$$

No, we won't be able to find its closed form solution because of the summation term inside the log. As discussed in class, we will need gradient based optimisation to get the point of estimate.