# TASK 1:

**Gradient:** firstly rewrite $f(\beta) = \left[ -\sum_{i=1}^{n}\sum_{k=0}^{K-1} \mathbb{1}_{\{y_i=k\}} \log p(\vec{x}_i; \vec{\beta}_k) + \frac{\lambda}{2}\sum_{k=0}^{K-1} \|\vec{\beta}_k\|_2^2 \right]$

$$\nabla_{\vec{\beta}_k} f(\beta) = -\sum_{i=1}^{n}\sum_{j=0}^{K-1} \mathbb{1}_{\{y_i=j\}} \nabla_{\vec{\beta}_k}^{\log} p(\vec{x}_i; \vec{\beta}_j) + \frac{\lambda}{2}\cdot 2\vec{\beta}_k$$

Since $\nabla_{\vec{\beta}_k}^{\log} p(\vec{x}_i; \vec{\beta}_k) \overset{\text{chain}}{\underset{\text{Rule}}{=}} \nabla_{\vec{\beta}_k}(\vec{x}_i^T \vec{\beta}_k) \cdot \frac{\partial \log P}{\partial \vec{x}_i^T \vec{\beta}_k} \overset{C-R}{=} \nabla_{\vec{\beta}_k}(\vec{x}_i^T \vec{\beta}_k) \frac{\partial \exp(\vec{x}_i^T \vec{\beta}_k)}{\partial \vec{x}_i^T \vec{\beta}_k} \cdot \frac{\partial P}{\partial e^{\vec{x}_i^T \vec{\beta}_k}} \cdot \frac{\partial \log P}{\partial P}$

$$= \frac{\vec{x}_i \cdot \exp(\vec{x}_i^T \vec{\beta}_k)}{p(\vec{x}_i; \vec{\beta}_k)} \frac{\sum_{\ell \neq k} e^{\vec{x}_i^T \vec{\beta}_\ell}}{(\sum_0^{K-1} e^{\vec{x}_i^T \vec{\beta}_\ell})^2} = \frac{\vec{x}_i \cdot p(\vec{x}_i; \vec{\beta}_k)\left[1 - p(\vec{x}_i; \vec{\beta}_k)\right]}{p(\vec{x}_i; \vec{\beta}_k)}$$

when $j \neq k$ $\qquad = \vec{x}_i \cdot \left[1 - p(\vec{x}_i; \vec{\beta}_k)\right]$

Similarly $\nabla_{\vec{\beta}_k} \log p(\vec{x}_j; \vec{\beta}_j) \overset{C-R}{=} \vec{x}_i \exp(\vec{x}_i^T \vec{\beta}_k) \cdot \frac{1}{P}\left(-\frac{\exp(\vec{x}_i^T \vec{\beta}_j)}{(\sum_{\ell=0}^{K-1} e^{\vec{x}_i^T \vec{\beta}_\ell})^2}\right) = -\frac{\vec{x}_i \, e^{\vec{x}_i^T \vec{\beta}_k}}{\sum_{\ell=0}^{K-1} e^{\vec{x}_i^T \vec{\beta}_\ell}}$

Thus, $\sum_{j=0}^{K-1} \mathbb{1}_{\{y_i=j\}} \nabla_{\vec{\beta}_k} \log p(\vec{x}_i; \vec{\beta}_j) = \vec{x}_i \mathbb{1}_{\{y_i=k\}} - \vec{x}_i \sum_j \mathbb{1}_{\cdots} \cdot p \qquad = -\vec{x}_i \cdot p(\vec{x}_i; \vec{\beta}_k)$

$$= \vec{x}_i \left(\mathbb{1}_{\{y_i=k\}} - p(\vec{x}_i; \vec{\beta}_k)\right)$$

$\therefore \nabla_{\vec{\beta}_k} f(\beta) = -\sum_i \vec{x}_i \left(\mathbb{1}_{(y_i=k)} - p(\vec{x}_i; \vec{\beta}_k)\right) + \lambda\vec{\beta}_k = X^T\left(P_k(\vec{x}_i) - \mathbb{1}_{\{\vec{y}=k\}}\right) + \lambda\vec{\beta}_k$

**Hessian:** With the results of $\nabla_{\vec{\beta}_k} f(\beta)$, we can directly derive $H_{\vec{\beta}_k} f(\beta) = \frac{\partial}{\partial \vec{\beta}_k^T} X^T\big($

$P_k(\vec{x}_i) - \mathbb{1}_{\{\vec{y}=k\}}\big) + \lambda I$, with the former $= X^T W_k X$, since for the $i^{th}$ row,

$\frac{\partial}{\partial \vec{\beta}_k^T} P(\vec{x}_i; \vec{\beta}_k) = \vec{x}_i^T \cdot P(\vec{x}_i; \vec{\beta}_k)\left[1 - P(\vec{x}_i; \vec{\beta}_k)\right]$. (Collect to get $W_k X$)

Therefore, $-\eta \underset{\substack{\wedge \\ \text{damped} \\ \text{rate}}}{\left[\lambda I + X^T W_k X\right]^{-1}} \left[X^T\{P_k(\vec{x}_i) - \mathbb{1}_{\{\vec{y}=k\}} + \lambda\vec{\beta}_k^{(t)}\right]$ is a damped Newton

step.