

We wish to show:

$$\underbrace{\beta_k^{(t+1)} = \beta_k^{(t)} - \eta \left(\underbrace{x^T W_k x}_{P \times P} + \lambda I \right)^{-1} \left[\underbrace{x^T \{P_k - \mathbb{1}_{\{y_i=k\}}\} \cdot \beta_k^t}_{P \times 1} \right]}_{P \times 1}, \quad k=0, \dots, K-1$$

is given by

$$\beta^{t+1} = \beta^t - \eta [\nabla f(\beta^t)]^{-1} \nabla f(\beta^t)$$

$$\text{When } f(\beta) = \left[-\sum_{i=1}^n \sum_{k=0}^{K-1} \mathbb{1}_{\{y_i=k\}} \log(p_k(x_i; \beta)) + \frac{\lambda}{2} \sum_{k=0}^{K-1} \sum_{j=1}^p \beta_{kj}^2 \right], \text{ where } p_k(x_i; \beta) = \frac{\exp(x_i^T \beta_k)}{\sum_{e=0}^{K-1} \exp(x_i^T \beta_e)}$$

And W_k is a diagonal matrix with elements $(W_k)_{ii} = p_k(x_i; \beta)(1 - p_k(x_i; \beta))$ and P_k is an \mathbb{R}^p vector with $(P_k)_i = p_k(x_i; \beta)$

$$\log\left(\frac{\exp(x_i^T \beta_k)}{\sum_{e=0}^{K-1} \exp(x_i^T \beta_e)}\right) = (\log(\exp(x_i^T \beta_k)) - \log(\sum_{e=0}^{K-1} \exp(x_i^T \beta_e))) = x_i^T \beta_k - \log(\sum_{e=0}^{K-1} \exp(x_i^T \beta_e))$$

So

$$\begin{aligned} f(\beta) &= \left[-\sum_{i=1}^n \sum_{k=0}^{K-1} \mathbb{1}_{\{y_i=k\}} (x_i^T \beta_k - \log(\sum_{e=0}^{K-1} \exp(x_i^T \beta_e))) + \frac{\lambda}{2} \sum_{k=0}^{K-1} \sum_{j=1}^p \beta_{kj}^2 \right] \\ &= \left[-\sum_{i=1}^n \sum_{k=0}^{K-1} \mathbb{1}_{\{y_i=k\}} (x_i^T \beta_k) + \sum_{i=1}^n \sum_{k=0}^{K-1} \mathbb{1}_{\{y_i=k\}} \log(\sum_{e=0}^{K-1} \exp(x_i^T \beta_e)) + \frac{\lambda}{2} \sum_{k=0}^{K-1} \sum_{j=1}^p \beta_{kj}^2 \right] \\ &= \left[-\sum_{i=1}^n \sum_{k=0}^{K-1} \mathbb{1}_{\{y_i=k\}} (x_i^T \beta_k) + \sum_{i=1}^n \log(\sum_{e=0}^{K-1} \exp(x_i^T \beta_e)) + \frac{\lambda}{2} \sum_{k=0}^{K-1} \sum_{j=1}^p \beta_{kj}^2 \right] \end{aligned}$$

Since only one term of $\sum_{k=0}^{K-1} \mathbb{1}_{\{y_i=k\}} \log(\sum_{e=0}^{K-1} \exp(x_i^T \beta_e))$ will be non-zero.

I will need the following calculations to take the gradient of $f(\beta)$ above

$$(a) \nabla_{\beta_j} (\mathbb{1}_{\{y_i=k\}} (x_i^T \beta_k)) = \begin{cases} x_i & \text{if } j = k \\ 0 & \text{if } j \neq k \end{cases} = \prod_{j=0}^{K-1} \mathbb{1}_{\{j=k\}} \cdot x_i$$

This holds since if $j=k$, $\nabla_{\beta_j} (\mathbb{1}_{\{y_i=k\}} (x_i^T \beta_k)) = \nabla_{\beta_k} (x_i^T \beta_k) = x_i$, and $\nabla_{\beta_j} (\mathbb{1}_{\{y_i=k\}} (x_i^T \beta_k)) = \nabla_{\beta_j} (0) = 0$ if $j \neq k$

Moreover

$$\nabla_{\beta_j} \log(\sum_{e=0}^{K-1} \exp(x_i^T \beta_e)) = (\sum_{e=0}^{K-1} \exp(x_i^T \beta_e))^{-1} \nabla_{\beta_j} (\sum_{e=0}^{K-1} \exp(x_i^T \beta_e))$$

Since for $f: \mathbb{R} \rightarrow \mathbb{R}$ and $g: \mathbb{R}^n \rightarrow \mathbb{R}$, $\nabla_t f(g(t)) = f'(g(t)) \nabla_t g(t)$ (*)

However $\nabla_{\beta_j} (\sum_{e=0}^{K-1} \exp(x_i^T \beta_e)) = \sum_{e=0}^{K-1} \nabla_{\beta_j} \exp(x_i^T \beta_e)$ where $\nabla_{\beta_j} \exp(x_i^T \beta_e) = \exp(x_i^T \beta_e) \nabla_{\beta_j} (x_i^T \beta_e)$ by applying (*) again

$$\text{Hence } \nabla_{\beta_j} \exp(x_i^T \beta_e) = \begin{cases} x_i \exp(x_i^T \beta_e) & \text{if } j = e \\ 0 & \text{if } j \neq e \end{cases} = \prod_{j=0}^{K-1} \mathbb{1}_{\{j=e\}} \cdot x_i \exp(x_i^T \beta_e)$$

So

$$(b) \nabla_{\beta_j} \log(\sum_{e=0}^{K-1} \exp(x_i^T \beta_e)) = (\sum_{e=0}^{K-1} \exp(x_i^T \beta_e))^{-1} \sum_{e=0}^{K-1} \prod_{j=0}^{K-1} \mathbb{1}_{\{j=e\}} x_i \exp(x_i^T \beta_e) = \frac{\exp(x_i^T \beta_e)}{(\sum_{e=0}^{K-1} \exp(x_i^T \beta_e))} \cdot x_i = x_i \rho_j(x_i, \beta)$$

Finally

$$\nabla_{\beta_j} \left(\frac{\lambda}{2} \sum_{k=0}^{K-1} \sum_{i=1}^n \beta_{k,i}^2 \right) = \nabla_{\beta_j} \left(\frac{\lambda}{2} \sum_{k=0}^{K-1} \|\beta_k\|^2 \right) = \frac{\lambda}{2} \sum_{k=0}^{K-1} \nabla_{\beta_j} \|\beta_k\|^2$$

However $\nabla_{\beta_j} \|\beta_k\|^2 = \begin{cases} 2\beta_j & \text{if } j=k \\ 0 & \text{if } j \neq k \end{cases}$

So

$$(c) \quad \nabla_{\beta_j} \left(\frac{\lambda}{2} \sum_{k=0}^{K-1} \sum_{i=1}^n \beta_{k,i}^2 \right) = \frac{\lambda}{2} \sum_{k=0}^{K-1} \mathbb{1}_{\{j=k\}} (2\beta_j) = \frac{\lambda}{2} \cdot 2\beta_j = \lambda\beta_j$$

Using these results to calculate $\nabla_{\beta_j} f(\beta)$ for some fixed $j=0, \dots, K-1$
we have (a), (b) and (c) giving us

$$\begin{aligned} \nabla_{\beta_j} (f(\beta)) &= \nabla_{\beta_j} \left[-\frac{\lambda}{2} \sum_{k=0}^{K-1} \mathbb{1}_{\{j=k\}} (\beta_j \beta_k) + \frac{\lambda}{2} \log \left(\sum_{k=0}^{K-1} \exp(\beta_k^T \beta_k) \right) + \frac{\lambda}{2} \sum_{k=0}^{K-1} \sum_{i=1}^n \beta_{k,i}^2 \right] \\ &= -\sum_{i=1}^n \sum_{k=0}^{K-1} \nabla_{\beta_j} \left(\mathbb{1}_{\{j=k\}} (\beta_j \beta_k) \right) + \sum_{i=1}^n \nabla_{\beta_j} \left[\log \left(\sum_{k=0}^{K-1} \exp(\beta_k^T \beta_k) \right) \right] + \nabla_{\beta_j} \left[\sum_{k=0}^{K-1} \sum_{i=1}^n \beta_{k,i}^2 \right] \\ &= -\sum_{i=1}^n \sum_{k=0}^{K-1} x_i \mathbb{1}_{\{j=k\}} + \sum_{i=1}^n x_i p_j(x_i, \beta) + \lambda\beta_j \\ &= -\sum_{i=1}^n x_i + \sum_{i=1}^n x_i p_j(x_i, \beta) + \lambda\beta_j \\ &= \sum_{i=1}^n x_i (p_j(x_i, \beta) - 1) + \lambda\beta_j \\ &= X^T (P_j - \mathbb{1}_{\{j=j\}}) + \lambda\beta_j, \quad \text{for every } j=0, \dots, K-1 \end{aligned}$$

Now to calculate the Hessian of f denoted by $\nabla_{\beta_j}^2 f(\beta)$ we observe
that

$$\nabla_{\beta_j}^2 f(\beta) = \frac{\partial^2}{\partial \beta_j \partial \beta_k} f(\beta) = \nabla_{\beta_k} [\nabla_{\beta_j} f(\beta)]$$

where we've shown $\nabla_{\beta_j} f(\beta) = \sum_{i=1}^n x_i (p_j(x_i, \beta) - 1) + \lambda\beta_j$

It remains to calculate $\nabla_{\beta_k} \left[\sum_{i=1}^n x_i (p_k(x_i, \beta) - 1) + \lambda\beta_k \right]$

I'll need the following calculations

$$\nabla_{\beta_j} [x_i (p_k(x_i, \beta) - 1)] = \nabla_{\beta_j} [x_i p_k(x_i, \beta) - x_i] = x_i \nabla_{\beta_j} p_k(x_i, \beta)$$

$$\text{However } \nabla_{\beta_j} p_k(x_i, \beta) = \nabla_{\beta_j} \left[\frac{\exp(x_i^T \beta_k)}{\sum_{k=0}^{K-1} \exp(x_i^T \beta_k)} \right] = \frac{\left(\sum_{k=0}^{K-1} \exp(x_i^T \beta_k) \right) \nabla_{\beta_j} [\exp(x_i^T \beta_k)] - \exp(x_i^T \beta_k) \nabla_{\beta_j} \left[\sum_{k=0}^{K-1} \exp(x_i^T \beta_k) \right]}{\left(\sum_{k=0}^{K-1} \exp(x_i^T \beta_k) \right)^2}$$

We've shown $\nabla_{\beta_j} [\exp(x_i^T \beta_k)] = \mathbb{1}_{\{j=k\}} \cdot x_i \exp(x_i^T \beta_j)$

and we've shown $\nabla_{\beta_j} \left[\sum_{k=0}^{K-1} \exp(x_i^T \beta_k) \right] = \sum_{k=0}^{K-1} \mathbb{1}_{\{j=k\}} \cdot x_i \exp(x_i^T \beta_j) = x_i \exp(x_i^T \beta_j)$

$$\text{So } \nabla_{\beta_k} P_K(x_i, \beta) = \frac{\left(\sum_{e=0}^{k-1} \exp(x_i^T \beta_e) \right) \left(\prod_{\{j=k\}} x_i^T \exp(x_i^T \beta_j) \right) - \exp(x_i^T \beta_k) x_i^T \exp(x_i^T \beta_k)}{\left(\sum_{e=0}^{k-1} \exp(x_i^T \beta_e) \right)^2}$$

$$= \frac{\prod_{\{j=k\}} x_i^T \exp(x_i^T \beta_j)}{\sum_{e=0}^{k-1} \exp(x_i^T \beta_e)} - \frac{\exp(x_i^T \beta_k)}{\left(\sum_{e=0}^{k-1} \exp(x_i^T \beta_e) \right)} \cdot \frac{x_i^T \exp(x_i^T \beta_k)}{\left(\sum_{e=0}^{k-1} \exp(x_i^T \beta_e) \right)}$$

$$= \prod_{\{j=k\}} x_i^T P(x_i, \beta) - P_K(x_i, \beta) (x_i^T P(x_i, \beta))$$

We consider only the case when $j=k$ since we're calculating $\nabla_{\beta_k} [\nabla_{\beta_k} f(\beta)]$

So we have $\nabla_{\beta_k} [\nabla_{\beta_k} P_K(x_i, \beta)] = x_i^T P_K(x_i, \beta) - x_i^T P_K^2(x_i, \beta) = x_i^T P_K(x_i, \beta)(1 - P_K(x_i, \beta))$

$$\begin{aligned} \text{So } \nabla_{\beta_k} [\nabla f(\beta)] &= \sum_{i=1}^n \nabla_{\beta_k} [x_i^T (P_K(x_i, \beta) - 1)] + \nabla_{\beta_k} [\lambda \beta_k] \\ &= \sum_{i=1}^n x_i^T (P_K(x_i, \beta))(1 - P_K(x_i, \beta)) x_i + \lambda \\ &= X^T W_K X + \lambda I \end{aligned}$$

Where $(W_K)_{ii} = P_K(x_i, \beta)(1 - P_K(x_i, \beta))$ is a diagonal matrix

Hence by the damped Newton's Method:

$$\begin{aligned} \beta_k^{(t+1)} &= \beta_k^{(t)} - \eta [\nabla^2 f(\beta)]^{-1} [\nabla f(\beta)] \\ &= \beta_k^{(t)} - \eta [X^T W_K X + \lambda I]^{-1} (X^T (P_K - \prod_{\{y=k\}}) + \lambda \beta_k) \end{aligned}$$