

$$P(y_i = k | x_i) = p_k(x_i; \beta), \quad \sum_{k=0}^{K-1} p_k(x_i; \beta) = 1$$

$$p_k(x_i; \beta) = \frac{e^{x_i^T \beta_k}}{\sum_{l=0}^{K-1} e^{x_i^T \beta_l}}$$

$$f(\beta) = \left[-\sum_{i=1}^n \left\{ \sum_{k=0}^{K-1} 1(y_i = k) \log p_k(x_i; \beta) \right\} + \frac{n}{2} \sum_{k=0}^{K-1} \sum_{j=1}^p \beta_{kj}^2 \right]$$

Damped Newton's update with learning rate $\eta > 0$ has the form;

$$\beta_k^{(t+1)} = \beta_k^{(t)} - \eta \left[\nabla^2 f(\beta_k^{(t)}) \right]^{-1} \nabla f(\beta_k^{(t)}) \quad (*)$$

We want to find $\nabla f(\beta_k)$ and $\nabla^2 f(\beta_k)$.

- Compute the gradient.

$$\begin{aligned} \frac{\partial f(\beta)}{\partial \beta_k} &= -\sum_{i=1}^n 1(y_i = k) \frac{\partial [\log p_k(x_i; \beta)]}{\partial \beta_k} + \frac{n}{2} \times 2 \beta_k \\ &= -\sum_{i=1}^n 1(y_i = k) \frac{\partial [\log p_k(x_i; \beta)]}{\partial \beta_k} + n \beta_k \quad \text{--- (1)} \end{aligned}$$

Consider $\frac{\partial [\log p_k(x_i; \beta)]}{\partial \beta_k}$

$$\log p_k(x_i; \beta) = \log \frac{e^{x_i^T \beta_k}}{\sum_{l=0}^{K-1} e^{x_i^T \beta_l}} = x_i^T \beta_k - \log \sum_{l=0}^{K-1} e^{x_i^T \beta_l}$$

$$\begin{aligned} \frac{\partial}{\partial \beta_k} (x_i^T \beta_k - \log \sum_{l=0}^{K-1} e^{x_i^T \beta_l}) &= x_i - \frac{1}{\sum_{l=0}^{K-1} e^{x_i^T \beta_l}} \cdot \frac{\partial}{\partial \beta_k} \left(\sum_{l=0}^{K-1} e^{x_i^T \beta_l} \right) \quad \text{--- (2)} \end{aligned}$$

Consider $\frac{\partial}{\partial \beta_k} \sum_{l=0}^{K-1} e^{x_i^T \beta_l}$

If $l = k \Rightarrow \frac{\partial}{\partial \beta_k} \sum_{l=0}^{K-1} e^{x_i^T \beta_l} = e^{x_i^T \beta_k} \cdot x_i$

If $l \neq k \Rightarrow \frac{\partial}{\partial \beta_k} \sum_{l=0}^{K-1} e^{x_i^T \beta_l} = 0$

\therefore From (2);

$$\frac{\partial}{\partial \beta_k} \log p_k(x_i; \beta) = x_i - \frac{e^{x_i^T \beta_k}}{\sum_{l=0}^{K-1} e^{x_i^T \beta_l}} \cdot x_i$$

$$= x_i \left[1 - p_k(x_i; \beta) \right]$$

\therefore From (1);

$$\begin{aligned} \frac{\partial f(\beta)}{\partial \beta_k} &= - \sum_{i=1}^n 1(y_i = k) \cdot x_i (1 - p_k(x_i; \beta)) + n \beta_k \\ &= - \sum_{i=1}^n x_i (1(y_i = k) (1 - p_k(x_i; \beta)) + n \beta_k \\ &= - x^T (1(Y=k) - p_k) + n \beta_k \\ &= x^T (p_k - 1(Y=k)) + n \beta_k \quad \text{--- (3)} \end{aligned}$$

$$\frac{\partial^2 f(\beta)}{\partial \beta_k^2} = \frac{\partial}{\partial \beta_k} \left[\sum_{i=1}^n x_i (p_k(x_i; \beta) - 1(y_i = k)) \right] + n \quad \text{--- (4)}$$

Consider $\frac{\partial}{\partial \beta_k} \left[\sum_{i=1}^n x_i p_k(x_i; \beta) \right]$

$$= \sum_{i=1}^n x_i \cdot \frac{\partial}{\partial \beta_k} [p_k(x_i; \beta)] \quad \text{--- (5)}$$

⇒ Consider $\frac{\partial}{\partial \beta_k} p_k(x_i; \beta)$

$$= \frac{\partial}{\partial \beta_k} \frac{e^{x_i^T \beta_k}}{\sum_{l=0}^{K-1} e^{x_i^T \beta_l}}$$

$$\begin{aligned} \frac{\partial}{\partial \beta_k} \frac{e^{x_i^T \beta_k}}{\sum_{l=0}^{K-1} e^{x_i^T \beta_l}} &= \frac{\sum_{l=0}^{K-1} e^{x_i^T \beta_l} \frac{\partial}{\partial \beta_k} e^{x_i^T \beta_k} - e^{x_i^T \beta_k} \frac{\partial}{\partial \beta_k} \sum_{l=0}^{K-1} e^{x_i^T \beta_l}}{\left(\sum_{l=0}^{K-1} e^{x_i^T \beta_l} \right)^2} \\ &= \frac{e^{x_i^T \beta_k}}{\sum_{l=0}^{K-1} e^{x_i^T \beta_l}} x_i - \frac{e^{x_i^T \beta_k}}{\left(\sum_{l=0}^{K-1} e^{x_i^T \beta_l} \right)^2} \frac{\partial}{\partial \beta_k} \sum_{l=0}^{K-1} e^{x_i^T \beta_l} \end{aligned}$$

If $l=k \Rightarrow \frac{e^{x_i^T \beta_k}}{\sum_{l=0}^{K-1} e^{x_i^T \beta_l}} x_i - \frac{e^{x_i^T \beta_k}}{\left(\sum_{l=0}^{K-1} e^{x_i^T \beta_l} \right)^2} \cdot e^{x_i^T \beta_k} \cdot x_i$

$$= p_k(x_i; \beta) \cdot x_i - [p_k(x_i; \beta)]^2 x_i$$

$$= p_k(x_i; \beta) x_i (1 - p_k(x_i; \beta))$$

If $l \neq k \Rightarrow \frac{e^{x_i^T \beta_k}}{\sum_{l=0}^{K-1} e^{x_i^T \beta_l}} x_i - \frac{e^{x_i^T \beta_k}}{\sum_{l=0}^{K-1} e^{x_i^T \beta_l}} \cdot \sum_{l=0}^{K-1} \frac{e^{x_i^T \beta_l}}{\sum_{l=0}^{K-1} e^{x_i^T \beta_l}} x_i$

$$= p_k(x_i; \beta) x_i - p_k(x_i; \beta) \cdot \sum_{l=0}^{K-1} p_l(x_i; \beta) \cdot x_i$$

$$= p_k(x_i; \beta) x_i - p_k(x_i; \beta) \cdot x_i \quad ; \quad \sum_{l=0}^{K-1} p_l(x_i; \beta) = 1$$

$$= 0$$

∴ From (5)

$$\frac{\partial}{\partial \beta_k} \left[\sum_{i=1}^n x_i p_k(x_i; \beta) \right] = \sum_{i=1}^n x_i \cdot p_k(x_i; \beta) (1 - p_k(x_i; \beta)) x_i$$

$$(W_k)_{ii} = p_k(x_i; \beta) (1 - p_k(x_i; \beta))$$

From (4);

$$\nabla^2 f = \sum_{i=1}^n x_i p_k(x_i; \beta) (1 - p_k(x_i; \beta)) x_i + n$$

$$= X^T \left(\sum_{i=1}^n W_{kii} x_i x_i^T \right) + n I$$

$$= X^T W_k X + n I \quad \text{--- (5)}$$

From (*);

$$\beta_k^{(t+1)} = \beta_k^{(t)} - \eta (X^T W_k X + n I)^{-1} [X^T \{p_{k-1}(Y=k)\} + n \beta_k^{(t)}]$$

; $k = 0, \dots, K-1$