

First calculate the gradient:

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

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

$$\frac{\partial}{\partial \beta_k} \log p_k(x_i) = \frac{\partial}{\partial \beta_k} \left[ x_i \beta_k - \log \left( \sum_{l=0}^{K-1} e^{x_i \beta_l} \right) \right]$$

$$= \left[ x_i - \frac{1 \cdot [e^{x_i \beta_k}] x_i}{\sum_{l=0}^{K-1} e^{x_i \beta_l}} \right]$$

$$= [x_i - p_k(x_i) x_i] = [1 - p_k(x_i)] x_i \quad \text{--- (2)}$$

Substitute result of (2) in (1)

$$\nabla_k f(\beta_k) = \left[ -\sum_{i=1}^n 1(y_i=k) [1 - p_k(x_i)] x_i + \lambda \beta_k \right] \quad \text{--- (3)}$$

$$= \left[ \sum_{i=1}^n \left\{ -1(y_i=k) + 1(y_i=k) p_k(x_i) \right\} x_i + \lambda \beta_k \right]$$

$$\nabla_k f(\beta_k) = \left[ X^T \{ p_k - 1(y=k) \} + \lambda \beta_k \right] \quad \text{--- (4)}$$

Next calculate the Hessian

from ③

$$\nabla_k^2 f(\beta_k) = \left[ -\sum_{i=1}^n 1(y_i=k) \left[ -\frac{\partial}{\partial \beta_k} p_k(x_i) \right] x_i + \lambda \right] - \textcircled{5}$$

$$\frac{\partial}{\partial \beta_k} p_k(x_i) = \frac{\left( \sum_{l=0}^{K-1} e^{x_i \beta_l} \right) (e^{x_i \beta_k}) x_i - (e^{x_i \beta_k}) (e^{x_i \beta_k}) x_i}{\left( \sum_{l=0}^{K-1} e^{x_i \beta_l} \right)^2}$$

$$= x_i \left[ \frac{e^{x_i \beta_k}}{\sum_{l=0}^{K-1} e^{x_i \beta_l}} - \frac{(e^{x_i \beta_k})^2}{\left( \sum_{l=0}^{K-1} e^{x_i \beta_l} \right)^2} \right]$$

$$= x_i p_k(x_i) [1 - p_k(x_i)] = x_i w_{kii} - \textcircled{6}$$

Substitute the result from ⑥ in ⑤

$$\nabla_k^2 f(\beta_k) = \left[ \sum_{i=1}^n 1(y_i=k) [x_i w_{kii}] x_i + \lambda \right]$$

$$\nabla_k^2 f(\beta_k) = X^T W_k X + \lambda I - \textcircled{7}$$

Damped Newton's update :

$$\beta_k^{(t+1)} = \beta_k^{(t)} - \eta [\nabla_k^2 f(\beta_k)]^{-1} [\nabla_k f(\beta_k)]$$

substituting results from ④ and ⑦

$$\beta_k^{(t+1)} = \beta_k^{(t)} - \eta [X^T W_k X + \lambda I]^{-1} [X^T \{ p_k - 1(y=k) \} + \lambda \beta_k^{(t)}]$$

for  $k = 0, 1, 2, \dots, K-1$