

# Derivation

(1)

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

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

We will calculate the gradient and Hessian of this  $f(\beta)$ .

## Gradient

$$\frac{\lambda}{2} \cdot 2\beta_k$$

We will take derivative with respect to  $\beta_k$  ( $\nabla_{\beta_k}$ ):

$$\nabla_{\beta_k} f(\beta) = - \sum_{i=1}^n \left\{ \sum_{t=0}^{K-1} \mathbb{1}(y_i=t) \cdot \nabla_{\beta_k} \log p_t(x_i; \beta) \right\} + \lambda \cdot \beta_k$$

$$\text{and for } \nabla_{\beta_k} \log p_t(x_i; \beta) = \nabla_{\beta_k} \log \left( \frac{e^{x_i^T \beta_k}}{\sum_{l=0}^{K-1} e^{x_i^T \beta_l}} \right) =$$

If  $t=k$ :

Derivative of numerator of  $p_k; \beta$  wrt  $\beta_k$ :  $e^{x_i^T \beta_k} \cdot x_i$

Derivative of denominator of  $p_k; \beta$  wrt  $\beta_k$  ( $k=l$ ):  $e^{x_i^T \beta_k} \cdot x_i$

$$\text{By Quotient rule: } \nabla_{\beta_k} p_k(x_i; \beta) = \frac{\left( \sum_{l=0}^{K-1} e^{x_i^T \beta_l} \right) (e^{x_i^T \beta_k} \cdot x_i) - (e^{x_i^T \beta_k} \cdot e^{x_i^T \beta_k} \cdot x_i)}{\left( \sum_{l=0}^{K-1} e^{x_i^T \beta_l} \right)^2}$$

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

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

$$\nabla_{\beta_k} p_k(x_i; \beta) = p_k(x_i; \beta) \cdot x_i \cdot (1 - p_k(x_i; \beta))$$

②

By Quotient rule:  $\nabla_{\beta_k} p_k(x_i; \beta) = \frac{\left( \sum_{l=0}^{K-1} e^{x_i^T \beta_l} \right) \cdot 0 - e^{x_i^T \beta_k} \cdot e^{x_i^T \beta_k} \cdot x_i}{\left( \sum_{l=0}^{K-1} e^{x_i^T \beta_l} \right)^2}$

$$\nabla_{\beta_k} p_k(x_i; \beta) = -p_t(x_i; \beta) \cdot p_k(x_i; \beta) \cdot x_i$$

$$\nabla_{\beta_K} f(\beta) = - \sum_{i=1}^n \left\{ \sum_{t=0}^{K-1} \mathbb{1}(y_i=t) \cdot \nabla_{\beta_K} \log p_t(x_i; \beta) \right\} + \lambda \cdot \beta_K$$

$$= - \sum_{i=1}^n \left\{ \sum_{t=0}^{K-1} \mathbb{1}_{(y_i=t)} \cdot \frac{\nabla_{\beta_k} p_t(x_i; \beta)}{p_t(x_i; \beta)} \right\} + \lambda \beta_k$$

$$= - \sum_{i=1}^n \left\{ \sum_{t=0}^{K-1} \mathbb{1}_{(y_i=t)} \left[ \frac{P_K(X_i; \beta) \cdot X_i \cdot (1 - P_K(X_i; \beta))}{P_K(X_i; \beta)} + \frac{-X_i \cdot P_t(X_i; \beta) \cdot P_K(X_i; \beta)}{P_t(X_i; \beta)} \right] \right\}$$

$$= - \sum_{i=1}^n \left\{ \mathbb{1}_{(y_i=k)} x_i (1 - p_k(x_i; \beta)) - \sum_{t \neq k}^{K-1} \mathbb{1}_{(y_i=t)} x_i p_t(x_i; \beta) \right\} + \lambda \cdot \beta_k$$

$$= - \sum_{i=1}^n \left\{ \mathbb{1}_{(y_i=k)} x_i - \left( \mathbb{1}_{(y_i=k)} \frac{p_k(x_i; \beta)}{\sum_{t \neq k} p_t(x_i; \beta)} x_i + \sum_{t \neq k} \mathbb{1}_{(y_i=t)} x_i \frac{p_t(x_i; \beta)}{\sum_{t \neq k} p_t(x_i; \beta)} \right) \right\} + \lambda \cdot \beta_k$$

$$= - \sum_{i=1}^n \left\{ \mathbb{1}_{(y_i=k)} x_i - \left( \left[ p_k(x_i; \beta) \cdot x_i \right] \left( \mathbb{1}_{(y_i=k)} + \sum_{t \neq k}^{K-1} \mathbb{1}_{(y_i=t)} \right) \right) \right\} + \lambda \cdot \beta_k$$

$$= - \sum_{i=1}^n \left\{ \mathbb{1}_{(y_i=k)} \cdot x_i - p_k(x_i; \beta) \cdot x_i \right\} + \frac{1}{\lambda \cdot \beta_k}$$

$$= - \sum_{i=1}^n \{ x_i ( \mathbb{1}_{(y_i=k)} - p_k(x_i; \beta) ) \} + \lambda \cdot \beta_k.$$

$$= -X^T(\mathbb{1}(Y=k) - p_k) + \lambda \beta_k = X^T(p_k - \mathbb{1}(Y=k)) + \lambda \beta_k$$

$\mathbb{1}(y=k)$  is 1 if  $y=k$  (class), 0 if  $y \neq k$  (class),  $P_k$  is vector for each  $p_k(x_i; \beta)$ ,  $X$  is design matrix,  $\lambda \beta_k$  is regularization

$X$  is design matrix,  $\lambda \beta_K$  is regularization term

# Hessian

(3)

To derive the Hessian we need to differentiate the gradient  $\nabla_{\beta_k} f(\beta)$

$$\nabla_{\beta_k} f(\beta) = - \sum_{i=1}^n \left\{ x_i \left( \mathbb{1}_{(y_i=k)} - p_k(x_i; \beta) \right) \right\} + \lambda \beta_k$$

Since  $\mathbb{1}_{(y_i=k)}$  does not depend on  $\beta_k$ , derivative  $\rightarrow 0$ .

So we focus on  $x_i \cdot p_k(x_i; \beta)$  for the two cases when  $t=k$  (same class) and  $t \neq k$  (different class).

$t=k$  case:

$$\begin{aligned} \frac{d^2 p_k(x_i; \beta)}{d\beta_k^2} &= \frac{d}{d\beta_k} \nabla_{\beta_k} f(\beta) = x_i \cdot p_k(x_i; \beta) (1 - p_k(x_i; \beta)) x_i' \\ &\quad + p_k(x_i; \beta) \cdot 0 + \lambda \\ &= p_k(x_i; \beta) (1 - p_k(x_i; \beta)) \cdot x_i \cdot x_i' + \lambda \end{aligned}$$

$t \neq k$  case:

$$\begin{aligned} \frac{d^2 p_k(x_i; \beta)}{d\beta_k^2} &= \frac{d}{d\beta_k} \nabla_{\beta_k} f(\beta) = x_i \cdot -p_t(x_i; \beta) \cdot p_k(x_i; \beta) \cdot x_i' + \\ &\quad p_k(x_i; \beta) \cdot 0 \\ &= -p_t(x_i; \beta) \cdot p_k(x_i; \beta) x_i \cdot x_i' \end{aligned}$$

Now summing over all samples:

$$\text{Hessian}(\beta) = \begin{cases} \sum_{i=1}^n \underbrace{p_k(x_i; \beta) \cdot (1 - p_k(x_i; \beta))}_{W_{kk}} \cdot x_i x_i' + \lambda I, & \text{if } t=k \text{ (diagonal entries of Hessian matrix)} \\ \sum_{i=1}^n -p_t(x_i; \beta) \cdot p_k(x_i; \beta) x_i x_i', & \text{if } t \neq k \text{ (off-diagonal entries of Hessian matrix)} \end{cases}$$

So for diagonal entries ( $t=k$ ) we can write  $\text{Hessian}(\beta)$  as

$$H_{kk} = X^T W_k X + \lambda I, \quad W_{kii} = p_k(x_i; \beta) (1 - p_k(x_i; \beta)) \text{ as given}$$

Now let us plug in Gradient and Hessian into <sup>Damped</sup> Newton's update.



## Damped Newton's Update

(4)

$$\beta_k^{(t+1)} = \beta_k^{(t)} - \eta (H_{kk})^{-1} [\nabla_{\beta_k} f(\beta)], \quad k = 0, \dots, K-1.$$

Substituting in:

$$\beta_k^{(t+1)} = \beta_k^{(t)} - \eta (X^T W_k X + \lambda I)^{-1} [X^T (P_k - \mathbf{1}(Y=k)) + \lambda \beta_k^{(t)}].$$

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

---

Final Result