

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

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

$$-\sum_{i=1}^n \sum_{k=0}^{K-1} \mathbb{1}(y_i = k) \log p_k(x_i; \beta)$$

## Gradient

$$-\sum_{i=1}^n \mathbb{1}(y_i = k) \log p_k(x_i; \beta)$$

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

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

$$\frac{\partial}{\partial \beta_k} (x_i^\top \beta_k) = x_i$$

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

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

So:

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

$$\frac{\partial}{\partial \beta_{ik}} \log \left( \sum_{l=0}^{k-1} e^{\mathbf{x}_i^\top \boldsymbol{\beta}_l} \right) = \frac{e^{\mathbf{x}_i^\top \boldsymbol{\beta}_k}}{\sum_{l=0}^{k-1} e^{\mathbf{x}_i^\top \boldsymbol{\beta}_l}} \mathbf{x}_i = p_k(\mathbf{x}_i; \boldsymbol{\beta}) \cdot \mathbf{x}_i$$

$$\frac{\partial}{\partial \beta_{ik}} (-\mathbb{1}(y_i=k) \log p_k(\mathbf{x}_i; \boldsymbol{\beta})) = -\mathbb{1}(y_i=k) \mathbf{x}_i + p_k(\mathbf{x}_i; \boldsymbol{\beta}) \mathbf{x}_i$$

Sum over all  $i$ :

$$\nabla_{\beta_k} L(\boldsymbol{\beta}) = - \sum_{i=1}^n (\mathbb{1}(y_i=k) - p_k(\mathbf{x}_i; \boldsymbol{\beta})) \mathbf{x}_i$$

### Regularization term

$$\frac{1}{2} \sum_{k=0}^{K-1} \sum_{j=1}^p \beta_{kj}^2$$

$$\frac{\partial}{\partial \beta_k} \left( \frac{1}{2} \sum_{k=0}^{K-1} \sum_{j=1}^p \beta_{kj}^2 \right) = \frac{1}{2} \cdot 2 \sum_{j=1}^p \beta_{kj} = \lambda \sum_{j=1}^p \beta_{kj}$$

$$= \lambda \beta_k$$

Total Gradienten:

$$\nabla_{\beta_k} L(\boldsymbol{\beta}) = - \sum_{i=1}^n \underbrace{(\mathbb{1}(y_i=k) - p_k(\mathbf{x}_i; \boldsymbol{\beta}))}_{\text{red bracket}} \mathbf{x}_i$$

$$\mathbb{1}(y=k) = \begin{bmatrix} \mathbb{1}(y_1=k) \\ \vdots \\ \mathbb{1}(y_n=k) \end{bmatrix} \quad p_k = \begin{bmatrix} p_k(\mathbf{x}_1; \boldsymbol{\beta}) \\ \vdots \\ p_k(\mathbf{x}_n; \boldsymbol{\beta}) \end{bmatrix}$$

$$\nabla_{\beta_k} L(\beta) = \sum_{i=1}^n x_i (p_{ik}(x_i; \beta) - 1(y_i=k))$$

$$\nabla_{\beta_k} L(\beta) = X^T (P_k - 1(y=k))$$

Add regularization term:

$$\nabla_{\beta_k} f(\beta) = X^T (P_k - 1(y=k)) + \lambda \beta_k$$

## Hessian Calculation

$$\nabla_{\beta_k} L(\beta) = - \sum_{i=1}^n (1(y_i=k) - p_{ik}(x_i; \beta)) x_i$$

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

Apply quotient rule:

$$\frac{d}{dx} \left( \frac{f(x)}{g(x)} \right) = \frac{g(x)f'(x) - f(x)g'(x)}{g(x)^2}$$

$$f(x) = e^{x_i^T \beta_k} \quad g(x) = \sum_{l=0}^{k-1} e^{x_i^T \beta_l}$$

$$f'(x) = e^{x_i^T \beta_k} x_i \quad g'(x) = e^{x_i^T \beta_k} x_i$$

$$\frac{\partial}{\partial \beta_k} p_{ik}(x_i; \beta) = \frac{\left( \sum_{l=0}^{k-1} e^{x_i^T \beta_l} \cdot e^{x_i^T \beta_k} x_i \right) - e^{x_i^T \beta_k} \cdot e^{x_i^T \beta_k} x_i}{\left( \sum_{l=0}^{k-1} e^{x_i^T \beta_l} \right)^2}$$

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

$$\frac{\partial}{\partial \beta_k} p_k(x_i; \beta) = p_k(x_i; \beta) (1 - p_k(x_i; \beta)) x_i$$

Sub into simplified product rule expression:

$$\frac{\partial}{\partial \beta_k} (p_k(x_i; \beta) x_i) = (p_k(x_i; \beta) (1 - p_k(x_i; \beta)) x_i) x_i$$

$$\frac{\partial}{\partial \beta_k} (p_k(x_i; \beta) x_i) = p_k(x_i; \beta) (1 - p_k(x_i; \beta)) x_i x_i^\top$$

Hessian:  $\sum_{i=1}^n p_k(x_i; \beta) (1 - p_k(x_i; \beta)) x_i x_i^\top$

$x_i \in \mathbb{R}^p$  is a  
feature vector

$X \in \mathbb{R}^{n \times p}$  is the matrix of input features where each row corresponds to feature vector  $x_i^\top$

$$X = \begin{bmatrix} x_1^\top \\ x_2^\top \\ \vdots \\ x_n^\top \end{bmatrix}$$

Weights  $p_k(x_i; \beta) (1 - p_k(x_i; \beta))$  apply to each outer product individually. Use diagonal matrix  $W_k \in \mathbb{R}^{n \times n}$  so:

$$W_K = \begin{bmatrix} p_K(x_1; \beta) (1 - p_K(x_1; \beta)) & 0 & \dots & \dots & 0 \\ 0 & \ddots & \ddots & \ddots & \vdots \\ \vdots & \ddots & \ddots & \ddots & 0 \\ 0 & \dots & \dots & \dots & 0 & p_K(x_n; \beta) (1 - p_K(x_n; \beta)) \end{bmatrix}$$

## Hessian Regularization Term

$$\frac{\partial}{\partial \beta_k} \left( \frac{1}{2} \sum_{j=1}^p \beta_{kj}^2 \right) = \lambda \beta_k \Rightarrow \frac{\partial}{\partial \beta_k} \lambda \beta_k = \lambda I$$

Thus Hessian matrix:

$$H_{\beta_k} = X^T W_K X + \lambda I$$

## Newton's Update Formula

$$\beta_k^{(t+1)} = \beta_k^{(t)} - \eta H_{\beta_k}^{-1} \nabla_{\beta_k} f(\beta)$$

$$\beta_k^{(t+1)} = \beta_k^{(t)} - \eta (X^T W_K X + \lambda I)^{-1} (X^T (p_k - \mathbb{1}(y_i=k)) + \lambda \beta_k^{(t)})$$

Case 1 :  $t=K$  (when  $y_i = K$ )

For  $t=K$ , the indicator function  $\mathbb{1}(y_i=k)=1$ :

$$\nabla_{\beta_k} f(\beta) = X^T (p_k - 1) + \lambda \beta_k$$

Apply Newton's Update:

$$\beta_k^{(t+1)} = \beta_k^{(t)} - \eta \underbrace{(X^T W_K X + \lambda I)^{-1}}_{\text{Hessian, capturing the loss function}} \underbrace{(X^T (p_k - 1) + \lambda \beta_k^{(t)})}_{\text{gradient}}$$

Update accounts for both the predicted probabilities & the true labels

$$\nabla_{\beta_k} f(\beta) = X^T (P_k - 1) + \lambda \beta_k$$

Case 2:  $t \neq k$  (when  $y_i \neq k$ )

For  $t \neq k$ , indicator function  $\mathbb{1}(y_i = t) = \emptyset$

The gradient for this case simplifies to

$$\nabla_{\beta_k} f(\beta) = X^T P_k + \lambda \beta_k$$

Newton's update:

$$\beta_k^{(t+1)} = \beta_k^{(t)} - \eta (X^T W_k X + \lambda I)^{-1} (X^T P_k + \lambda \beta_k^{(t)})$$

Update is simpler involving only the predicted probabilities as

$$\nabla_{\beta_k} f(\beta) = X^T P_k + \lambda \beta_k$$

