

Let  $X$  be training data matrix with additional column of 1s on the first column,  $\beta$  be the coefficient matrix.  $\vec{x}_i = (1, x_{i1}, x_{i2}, \dots, x_{im})$  is a datum on row  $i$  in  $X$ .  $\vec{\beta}_k = (\beta_{0k}, \beta_{1k}, \dots, \beta_{mk})$  is a vector in  $\beta$  of class  $k$  in which there are  $K$  classes in response  $y$ . Let  $\hat{y}_{ik}$  be the prediction (or probability) of  $\vec{x}_i$  on class  $k$ , then  $\hat{y}_k = (\hat{y}_{1k}, \hat{y}_{2k}, \dots, \hat{y}_{nk})^T$ . Since

$$\begin{aligned} \hat{y}_{ik} &= P[Y = k | X = \vec{x}_i] = \frac{e^{\vec{\beta}_k \cdot \vec{x}_i}}{\sum_{l=1}^K e^{\vec{\beta}_l \cdot \vec{x}_i}} \\ &\rightarrow \ln(\hat{y}_{ik}) = \vec{\beta}_k \cdot \vec{x}_i - \ln\left(\sum_{l=1}^K e^{\vec{\beta}_l \cdot \vec{x}_i}\right) \\ &\rightarrow \begin{cases} \nabla_{\vec{\beta}_k} \ln(\hat{y}_{ik}) = \vec{x}_i - \frac{\nabla_{\vec{\beta}_k} [\sum_{l=1}^K e^{\vec{\beta}_l \cdot \vec{x}_i}]}{\sum_{l=1}^K e^{\vec{\beta}_l \cdot \vec{x}_i}} = \vec{x}_i - \frac{e^{\vec{\beta}_k \cdot \vec{x}_i}}{\sum_{l=1}^K e^{\vec{\beta}_l \cdot \vec{x}_i}} \vec{x}_i = \vec{x}_i - \hat{y}_{ik} \vec{x}_i = (1 - \hat{y}_{ik}) \vec{x}_i \\ \nabla_{\vec{\beta}_j} \ln(\hat{y}_{ik}) = -\frac{e^{\vec{\beta}_j \cdot \vec{x}_i}}{\sum_{l=1}^K e^{\vec{\beta}_l \cdot \vec{x}_i}} \vec{x}_i = -\hat{y}_{ij} \vec{x}_i \end{cases} \end{aligned}$$

Now we use likelihood function  $l$  to build loss function and try to maximize it to find matrix  $\beta$  in which

$$\begin{aligned} l &= \prod_{i_1:Y=1} \hat{y}_{i_1 1} \prod_{i_2:Y=2} \hat{y}_{i_2 2} \dots \prod_{i_K:Y=K} \hat{y}_{i_K K} \\ &\rightarrow \ln(l) = \sum_{i_1} \ln(\hat{y}_{i_1 1}) + \sum_{i_2} \ln(\hat{y}_{i_2 2}) + \dots + \sum_{i_K} \ln(\hat{y}_{i_K K}) \end{aligned}$$

To maximize  $l$  is to maximize  $\ln(l)$ , then the goal is to find  $\beta$  for which

$$\nabla_{\vec{\beta}_k} \ln(l) = 0 \quad \text{for } k = 1 \dots K$$

Let  $I_k$  be a  $n \times 1$  array where

$$I_k = \begin{cases} 1 & Y=k \\ 0 & \text{otherwise} \end{cases}$$

Then

$$\begin{aligned} \nabla_{\vec{\beta}_k} \ln(l) &= -\sum_{i_1} \hat{y}_{i_1 k} \vec{x}_{i_1} - \sum_{i_2} \hat{y}_{i_2 k} \vec{x}_{i_2} - \dots + \sum_{i_k} (1 - \hat{y}_{i_k k}) \vec{x}_{i_k} - \dots - \sum_{i_K} \hat{y}_{i_K k} \vec{x}_{i_K} \\ &= \sum_{i_k} \vec{x}_{i_k} - \sum_{i_1} \hat{y}_{i_1 k} \vec{x}_{i_1} - \sum_{i_2} \hat{y}_{i_2 k} \vec{x}_{i_2} - \dots - \sum_{i_K} \hat{y}_{i_K k} \vec{x}_{i_K} \\ &= \sum_{i_k} \vec{x}_{i_k} - \sum_i^n \hat{y}_{ik} \vec{x}_i \\ &= \sum_i^n (I_k - \hat{y}_{ik}) \vec{x}_i \\ &= X^T \cdot (I_k - \hat{y}_k) \end{aligned}$$

So

$$\left\{ \begin{array}{l} \nabla_{\vec{\beta}_1} \ln(l) = \sum_i^n (I_1 - \hat{y}_{i1}) \vec{x}_i = X^T \cdot (I_1 - \hat{y}_1) \\ \nabla_{\vec{\beta}_2} \ln(l) = \sum_i^n (I_2 - \hat{y}_{i2}) \vec{x}_i = X^T \cdot (I_2 - \hat{y}_2) \\ \vdots \\ \nabla_{\vec{\beta}_K} \ln(l) = \sum_i^n (I_K - \hat{y}_{iK}) \vec{x}_i = X^T \cdot (I_K - \hat{y}_K) \end{array} \right.$$

$$\rightarrow \nabla_{\vec{\beta}} \ln(l) = X^T \cdot (I - \hat{y})$$

Now we know which direction to renew matrix  $\beta$  via  $\nabla_{\vec{\beta}} \ln(l)$ , and with the help of learning rate  $\alpha$ , we then have

$$\beta_{new} = \beta_{old} + \alpha X^T \cdot (I - \hat{y})$$

Below are matrices definition used in this document

$$X = \left[ \begin{array}{c} \left[ \begin{array}{c} 1 \\ 1 \\ \vdots \\ 1 \end{array} \right] [x_1] [x_2] \dots [x_m] \end{array} \right]$$

$$I = \left[ [I_1] [I_2] \dots [I_K] \right]$$

$$\beta = \left[ \begin{array}{c} \left[ \begin{array}{c} \beta_{01} \\ \beta_{11} \\ \vdots \\ \beta_{m1} \end{array} \right] [\beta_2] \dots [\beta_K] \end{array} \right]$$

$$\hat{y} = \left[ [\hat{y}_1] [\hat{y}_2] \dots [\hat{y}_K] \right]$$