

# The derivative of LLK score function with the Multinomial Logistic function

reda.dehak

October 2022

## 1 The Softmax Function:

$$\text{SoftMax}_{\Theta}(\mathbf{x}, k) = \frac{e^{\theta_{.k}^T \mathbf{x}}}{\sum_{c=1}^K e^{\theta_{.c}^T \mathbf{x}}} \quad (1)$$

where  $\theta_{.k}$  is the column number  $k$  of the matrix  $\Theta$ ,  $K$  is the number of classes

$$\frac{\partial \text{SoftMax}_{\Theta}(\mathbf{x}, k)}{\partial \theta_{ij}} = \begin{cases} \frac{x_i e^{\theta_{.k}^T \mathbf{x}} \sum_{c=1}^K e^{\theta_{.c}^T \mathbf{x}} - x_i (e^{\theta_{.k}^T \mathbf{x}})^2}{\left( \sum_{c=1}^K e^{\theta_{.c}^T \mathbf{x}} \right)^2} & \text{if } j = k \\ -\frac{x_i e^{\theta_{.k}^T \mathbf{x}} e^{\theta_{.j}^T \mathbf{x}}}{\left( \sum_{c=1}^K e^{\theta_{.c}^T \mathbf{x}} \right)^2} & \text{Otherwise} \end{cases} \quad (2)$$

$$= \begin{cases} x_i \left( \frac{e^{\theta_{.k}^T \mathbf{x}}}{\sum_{c=1}^K e^{\theta_{.c}^T \mathbf{x}}} - \left( \frac{e^{\theta_{.k}^T \mathbf{x}}}{\sum_{c=1}^K e^{\theta_{.c}^T \mathbf{x}}} \right)^2 \right) & \text{if } j = k \\ -x_i \frac{e^{\theta_{.k}^T \mathbf{x}}}{\sum_{c=1}^K e^{\theta_{.c}^T \mathbf{x}}} \frac{e^{\theta_{.j}^T \mathbf{x}}}{\sum_{c=1}^K e^{\theta_{.c}^T \mathbf{x}}} & \text{Otherwise} \end{cases} \quad (3)$$

$$= \begin{cases} x_i \left( \text{SoftMax}_{\Theta}(\mathbf{x}, k) - \text{SoftMax}_{\Theta}(\mathbf{x}, k)^2 \right) & \text{if } j = k \\ -x_i \text{SoftMax}_{\Theta}(\mathbf{x}, j) \text{SoftMax}_{\Theta}(\mathbf{x}, k) & \text{Otherwise} \end{cases} \quad (4)$$

$$= \begin{cases} x_i \text{SoftMax}_{\Theta}(\mathbf{x}, k) (1 - \text{SoftMax}_{\Theta}(\mathbf{x}, k)) & \text{if } j = k \\ -x_i \text{SoftMax}_{\Theta}(\mathbf{x}, j) \text{SoftMax}_{\Theta}(\mathbf{x}, k) & \text{Otherwise} \end{cases} \quad (5)$$

$$= \begin{cases} x_i \text{SoftMax}_{\Theta}(\mathbf{x}, k) (1 - \text{SoftMax}_{\Theta}(\mathbf{x}, k)) & \text{if } j = k \\ x_i \text{SoftMax}_{\Theta}(\mathbf{x}, j) (0 - \text{SoftMax}_{\Theta}(\mathbf{x}, k)) & \text{Otherwise} \end{cases} \quad (6)$$

$$\frac{\partial \text{SoftMax}_{\Theta}(\mathbf{x}, k)}{\partial \theta_{ij}} = x_i \text{SoftMax}_{\Theta}(\mathbf{x}, j) \left( \boxed{\begin{cases} 1 & \text{if } j = k \\ 0 & \text{Otherwise} \end{cases}} - \text{SoftMax}_{\Theta}(\mathbf{x}, k) \right) \quad (7)$$

$$= x_i \text{SoftMax}_\Theta(\mathbf{x}, j) (\mathbb{1}_{(j=k)} - \text{SoftMax}_\Theta(\mathbf{x}, k)) \quad (8)$$

So we can generalise to vectors

$$\frac{\partial \text{SoftMax}_\Theta(\mathbf{x}, k)}{\partial \theta_{.j}} = \mathbf{x} \text{SoftMax}_\Theta(\mathbf{x}, j) (\mathbb{1}_{(j=k)} - \text{SoftMax}_\Theta(\mathbf{x}, k)) \quad (9)$$

## 2 NLLK Cost Function

$$\text{NLLK} = -\frac{1}{N} \sum_{n=1}^N \sum_{k=1}^K \mathbb{1}_{(y_n=k)} \log(\text{SoftMax}_\Theta(\mathbf{x}_n, k)) \quad (10)$$

$$\frac{\partial \text{NLLK}}{\partial \theta_{.j}} = -\frac{1}{N} \sum_{n=1}^N \sum_{k=1}^K \mathbb{1}_{(y_n=k)} \frac{\frac{\partial \text{SoftMax}_\Theta(\mathbf{x}_n, k)}{\partial \theta_{.j}}}{\text{SoftMax}_\Theta(\mathbf{x}_n, k)} \quad (11)$$

$$= -\frac{1}{N} \sum_{n=1}^N \sum_{k=1}^K \mathbb{1}_{(y_n=k)} \frac{\mathbf{x}_n \text{SoftMax}_\Theta(\mathbf{x}_n, j) (\mathbb{1}_{(j=k)} - \text{SoftMax}_\Theta(\mathbf{x}_n, k))}{\text{SoftMax}_\Theta(\mathbf{x}_n, k)} \quad (12)$$

$$= -\frac{1}{N} \sum_{n=1}^N \mathbf{x}_n \text{SoftMax}_\Theta(\mathbf{x}_n, j) \left( \sum_{k=1}^K \frac{\mathbb{1}_{(y_n=k)} \mathbb{1}_{(j=k)}}{\text{SoftMax}_\Theta(\mathbf{x}_n, k)} - \sum_{k=1}^K \frac{\mathbb{1}_{(y_n=k)} \cancel{\text{SoftMax}_\Theta(\mathbf{x}_n, k)}}{\cancel{\text{SoftMax}_\Theta(\mathbf{x}_n, k)}} \right)$$

$$= -\frac{1}{N} \sum_{n=1}^N \mathbf{x}_n \text{SoftMax}_\Theta(\mathbf{x}_n, j) \left( \frac{\mathbb{1}_{(y_n=j)}}{\text{SoftMax}_\Theta(\mathbf{x}_n, j)} - \sum_{k=1}^K \mathbb{1}_{(y_n=k)} \right) \quad (13)$$

$$= -\frac{1}{N} \sum_{n=1}^N \mathbf{x}_n \cancel{\text{SoftMax}_\Theta(\mathbf{x}_n, j)} \left( \frac{\mathbb{1}_{(y_n=j)}}{\cancel{\text{SoftMax}_\Theta(\mathbf{x}_n, j)}} - 1 \right) \quad (14)$$

$$= -\frac{1}{N} \sum_{n=1}^N \mathbf{x}_n (\mathbb{1}_{(y_n=j)} - \text{SoftMax}_\Theta(\mathbf{x}_n, j)) \quad (15)$$

$$\frac{\partial \text{NLLK}}{\partial \theta_{.j}} = -\frac{1}{N} \sum_{n=1}^N \mathbf{x}_n e_{nj} \quad (16)$$

where  $e_{nj}$  represents the output error for sample  $n$  and class  $j$

$$e_{nj} = \mathbb{1}_{(y_n=j)} - \text{SoftMax}_\Theta(\mathbf{x}_n, j) \quad (17)$$