

$$J = \frac{-1}{N} \sum_{i=1}^N \sum_{k=1}^K y_{ik} \log \left( \frac{\exp f_k}{\sum_{c=1}^K \exp f_c} \right) + \lambda \sum_{j=1}^d w_k^2$$

Rearranging the terms,

$$J = \frac{-1}{N} \sum_{i=1}^N \sum_{k=1}^K y_{ik} \left( \log \exp f_k - \log \sum_{c=1}^K \exp f_c \right) + \lambda \sum_{j=1}^d w_k^2$$

Since  $\log e^a = a$

$$J = \frac{-1}{N} \sum_{i=1}^N \sum_{k=1}^K y_{ik} f_k + \frac{1}{N} \sum_{i=1}^N \sum_{k=1}^K \log \sum_{c=1}^K e^{f_c} + \lambda \sum_{j=1}^d w_k^2$$

For gradient we need to differentiate the function  $J$  wrt to  $w_k$

$$\frac{dJ}{dw_k} = \frac{d}{dw_k} \left[ \frac{-1}{N} \sum_{i=1}^N \sum_{k=1}^K y_{ik} f_k \right] + \frac{d}{dw_k} \left[ \frac{1}{N} \sum_{i=1}^N \sum_{k=1}^K \log \sum_{c=1}^K e^{f_c} \right]$$

+  $2w_k$

$$= \frac{-1}{N} \sum_{i=1}^N y_{ik} f'_k + \frac{1}{N} \sum_{i=1}^N \frac{e^{f_k}}{\sum_{c=1}^K e^{f_c}} \cdot f'_k + 2w_k$$

Taking common terms out.

$$\frac{1}{N} \sum_{i=1}^N \left( \frac{e^{f_k}}{\sum_{c=1}^K e^{f_c}} - y_{ik} \right) f'(k) + 2w_k$$

According to the definition of Softmax

$$\frac{e^{f_k}}{\sum_{c=1}^K e^{f_c}} = P(C_i | x)$$

The above equation becomes

$$\frac{1}{N} \sum_{i=1}^N \left[ P(C_i | x) - y_{ik} \right] f'(k) + 2w_k$$