

- 1) The probability of being correct for example t is:

$$P(y = i | x_t) = \text{softmax}(wx_t + b)_i = \frac{e^{w_i x_t + b_i}}{\sum_{j=1}^k e^{w_j x_t + b_j}}$$

- 2) The Loss function is:

$$L(w; b) = \arg \min \sum_t -\log\left(\frac{e^{w_i x_t + b_i}}{\sum_{j=1}^k e^{w_j x_t + b_j}}\right)$$

- 3) The update rule will be according to the derivative of the loss function w.r.t- w , and w.r.t- b . For SGD we will calculate the derivative for a single example t :

$$\begin{aligned} \frac{\partial L}{\partial w} &= \frac{\partial \left(-\log \left(\frac{e^{w_i x_t + b_i}}{\sum_{j=1}^k e^{w_j x_t + b_j}} \right) \right)}{\partial w} = \frac{\partial (-\log(e^{w_i x_t + b_i}) + \log(\sum_{j=1}^k e^{w_j x_t + b_j}))}{\partial w} \\ &= \frac{\partial (-w_i x_t - b_i + \log(\sum_{j=1}^k e^{w_j x_t + b_j}))}{\partial w} = \frac{\partial (\log(\sum_{j=1}^k e^{w_j x_t + b_j}) - w_i x_t - b_i)}{\partial w} \end{aligned}$$

When we derive w.r.t the vector of the correct class i :

$$\begin{aligned} \frac{\partial (\log(\sum_{j=1}^k e^{w_j x_t + b_j}) - w_i x_t - b_i)}{\partial w_i} &= \frac{1}{\sum_{j=1}^k e^{w_j x_t + b_j}} e^{w_i x_t + b_i} * x_t - x_t \\ &= \frac{e^{w_i x_t + b_i}}{\sum_{j=1}^k e^{w_j x_t + b_j}} * x_t - x_t = x_t * \text{softmax}(wx_t + b)_i - x_t \\ &= x_t (\text{softmax}(wx_t + b)_i - 1) \end{aligned}$$

When we derive all the other w_a vectors, when $a \neq i$, we get:

$$\begin{aligned} \frac{\partial (\log(\sum_{j=1}^k e^{w_j x_t + b_j}) - w_i x_t - b_i)}{\partial w_a} &= \frac{1}{\sum_{j=1}^k e^{w_j x_t + b_j}} e^{w_a x_t + b_a} * x_t = \frac{e^{w_a x_t + b_a}}{\sum_{j=1}^k e^{w_j x_t + b_j}} x_t \\ &= x_t (\text{softmax}(wx_t + b)_a) \end{aligned}$$

Now let's derive the loss w.r.t b : When we are in the b of the correct class, b_i :

$$\begin{aligned} \frac{\partial (\log(\sum_{j=1}^k e^{w_j x_t + b_j}) - w_i x_t - b_i)}{\partial b_i} &= \frac{1}{\sum_{j=1}^k e^{w_j x_t + b_j}} e^{w_i x_t + b_i} - 1 \\ &= \text{softmax}(wx_t + b)_i - 1 \end{aligned}$$

And w.r.t all others b_a :

$$\begin{aligned} \frac{\partial (\log(\sum_{j=1}^k e^{w_j x_t + b_j}) - w_i x_t - b_i)}{\partial b_a} &= \frac{1}{\sum_{j=1}^k e^{w_j x_t + b_j}} e^{w_a x_t + b_a} \\ &= \text{softmax}(wx_t + b)_a \end{aligned}$$

