

Градиент функции потерь (softmax loss) :

$$p_c = \frac{\exp(s_c)}{\sum_{k=1}^C \exp(s_k)}$$

$$L = - \sum_{c=1}^C [y_c \log(p_c)] = - \sum_{c=1}^C \left[y_c \log \left(\frac{\exp(s_c)}{\sum_{k=1}^C \exp(s_k)} \right) \right]$$

$$\left(\frac{\partial p_c}{\partial s_j} \right)_{c \neq j} = \frac{\partial}{\partial s_j} \left[\frac{\exp(s_c)}{\sum_{k=1}^C \exp(s_k)} \right] = - \frac{\exp(s_c) \exp(s_j)}{\left(\sum_{k=1}^C \exp(s_k) \right)^2} = -p_c p_j$$

$$\frac{\partial p_j}{\partial s_j} = \frac{\partial}{\partial s_j} \left[\frac{\exp(s_j)}{\sum_{k=1}^C \exp(s_k)} - \frac{(\exp(s_j))^2}{\left(\sum_{k=1}^C \exp(s_k) \right)^2} \right] = p_j - (p_j)^2 = p_j (1 - p_j)$$

Градиент сумматора второго слоя (в разрезе одного наблюдения) :

$$\frac{\partial L}{\partial s_{2j}} = - \frac{\partial}{\partial s_{2j}} \left(\sum_{c=1}^C [y_c \log(p_c)] \right) = - \sum_{c=1}^C \left[y_c \frac{1}{p_c} \frac{\partial p_c}{\partial s_{2j}} \right] =$$

$$= -y_j \frac{1}{p_j} p_j (1 - p_j) - \sum_{c \neq j} \left[y_c \frac{1}{p_c} (-p_c p_j) \right] =$$

$$= -y_j + y_j p_j + \sum_{c \neq j} [y_c p_j] = -y_j + \sum_{c=1}^C [y_c p_j] =$$

$$= -y_j + p_j \sum_{c=1}^C y_c = p_j - y_j$$

Градиент порога первого слоя :

$$\frac{\partial L}{\partial \mathbf{o1}} = \frac{\partial L}{\partial \mathbf{s2}} \times \frac{\partial \mathbf{s2}}{\partial \mathbf{o1}} = \frac{\partial L}{\partial \mathbf{s2}} \times \mathbf{w2}^T = (\mathbf{p_{ic}} - \mathbf{y_{ic}}) \times \mathbf{w2}^T$$

Градиент по весам второго слоя :

$$\frac{\partial L}{\partial \mathbf{w2}} = \frac{\partial L}{\partial \mathbf{s2}} \times \frac{\partial \mathbf{s2}}{\partial \mathbf{w2}} = \frac{\partial \mathbf{s2}}{\partial \mathbf{w2}} \times \frac{\partial L}{\partial \mathbf{s2}} = \mathbf{o1}^T \times (\mathbf{p_{ic}} - \mathbf{y_{ic}})$$

Градиент по входным данным :

$$\begin{aligned} \frac{\partial L}{\partial \mathbf{x}} &= \left(\frac{\partial L}{\partial \mathbf{o1}} \cdot \frac{\partial \mathbf{o1}}{\partial \mathbf{s1}} \right) \times \frac{\partial \mathbf{o1}}{\partial \mathbf{x}} = \left(\frac{\partial L}{\partial \mathbf{o1}} \cdot \frac{\partial \mathbf{o1}}{\partial \mathbf{s1}} \right) \times \left(\frac{\partial \mathbf{s1}}{\partial \mathbf{x}} \right)^T = \left(\frac{\partial L}{\partial \mathbf{o1}} \cdot \frac{\partial \mathbf{o1}}{\partial \mathbf{s1}} \right) \times \left(\frac{\mathbf{w1}}{[D \times H]} \right)^T = \\ &= \left(\frac{\partial L}{\partial \mathbf{o1}} \cdot [\mathbf{s1} > \mathbf{0}] \right) \times \left(\frac{\mathbf{w1}}{[D \times H]} \right)^T = \left(\left((\mathbf{p_{ic}} - \mathbf{y_{ic}}) \times \left(\frac{\mathbf{w2}^T}{[C \times H]} \right) \right) \cdot [\mathbf{s1} > \mathbf{0}] \right) \times \left(\frac{\mathbf{w1}}{[D \times H]} \right)^T \end{aligned}$$

Градиент по весам первого слоя :

$$\begin{aligned}
 \frac{\partial L}{\frac{\partial \mathbf{w}_1}{[D \times H]}} &= \left(\frac{\partial L}{\frac{\partial \mathbf{o}_1}{[N \times H]} \cdot \frac{\partial \mathbf{s}_1}{[N \times H]}} \right) \times \frac{\partial \mathbf{o}_1}{\frac{\partial \mathbf{w}_1}{[D \times N]}} = \\
 &= \frac{\partial \mathbf{o}_1}{\frac{\partial \mathbf{w}_1}{[D \times N]}} \times \left(\frac{\partial L}{\frac{\partial \mathbf{o}_1}{[N \times H]} \cdot \frac{\partial \mathbf{s}_1}{[N \times H]}} \right) = \frac{\partial \mathbf{o}_1}{\frac{\partial \mathbf{w}_1}{[D \times N]}} \times \left(\frac{\partial L}{\frac{\partial \mathbf{o}_1}{[N \times H]} \cdot \frac{\partial \mathbf{s}_1}{[N \times H]}} \right) = \\
 &= \left(\frac{\mathbf{x}}{\frac{[N \times D]}{[D \times N]}} \right)^T \times \left(\frac{\partial L}{\frac{\partial \mathbf{o}_1}{[N \times H]} \cdot \frac{\partial \mathbf{s}_1}{[N \times H]}} \right) = \left(\frac{\mathbf{x}}{\frac{[N \times D]}{[D \times N]}} \right)^T \times \left(\left(\mathbf{p}_{ic} - \mathbf{y}_{ic} \right) \times \left(\mathbf{w}_2^T \right) \cdot \left[\mathbf{s}_1 > 0 \right] \right)
 \end{aligned}$$