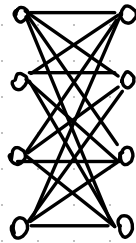


GRADIENT OF SOFTMAX + CROSS-ENTROPY

Let us consider the last layer of classification and the cost.



N size
of dataset



K K

logits outputs

$\rightarrow L = \text{cross-entropy}$

- $x \in \mathbb{R}^{N \times K}$ inputs
- $\tilde{y} \in \mathbb{R}^{N \times K}$ true targets (one-hot encoded)

$$L = -\frac{1}{N} \sum_{i=0}^{N-1} \sum_{k=0}^{K-1} \tilde{y}_{ik} \log \sigma_k(x_i)$$

Recall that

$$\sigma_k(x_{i.}) = \frac{e^{x_{ik}}}{\sum_{h=0}^{k-1} e^{x_{ih}}}$$

We want to compute

$$\frac{\partial L}{\partial x} \in \mathbb{R}^{K \times N}$$

$$\frac{\partial}{\partial x_{jl}} \left(\sigma_k(x_{i.}) \right) = \text{First of all :}$$

0 if $j \neq i$

Next, we distinguish two cases:

Case 1: $k=l$

$$\frac{\partial}{\partial x_{ik}} \left(\frac{e^{x_{ik}}}{\sum_{h=0}^{k-1} e^{x_{ih}}} \right) = \frac{e^{x_{ik}}}{\sum_{h=0}^{k-1} e^{x_{ih}}} - \frac{(e^{x_{ik}})^2}{\left(\sum_{h=0}^{k-1} e^{x_{ih}} \right)^2} =$$

$$= \sigma_k(x_{i.}) - \sigma_k(x_{i.})^2 =$$

$$= \sigma_k(x_{i.}) (1 - \sigma_k(x_{i.}))$$

Case 2: $k \neq l$

$$\begin{aligned}\frac{\partial}{\partial x_{il}} \left(\frac{e^{x_{ik}}}{\sum_{h=0}^{K-1} e^{x_{ih}}} \right) &= - \frac{e^{x_{ik}}}{\left(\sum_{h=0}^{K-1} e^{x_{ih}} \right)^2} e^{x_{il}} = \\ &= - \sigma_k(x_{i \cdot}) \sigma_l(x_{i \cdot})\end{aligned}$$

In a compact way:

$$\frac{\partial}{\partial x_{jl}} (\sigma_k(x_{i \cdot})) = \text{Id}_{ij} \sigma_k(x_{i \cdot}) (\text{Id}_{kl} - \sigma_l(x_{i \cdot}))$$

We are ready to compute the differential of the loss:

$$\begin{aligned}\frac{\partial L}{\partial x_{jl}} &= \frac{\partial}{\partial x_{jl}} \left(-\frac{1}{N} \sum_{i=0}^{N-1} \sum_{k=0}^{K-1} \tilde{y}_{ik} \log \sigma_k(x_{i \cdot}) \right) = \\ &= -\frac{1}{N} \sum_{i=0}^{N-1} \sum_{k=0}^{K-1} \tilde{y}_{ik} \frac{1}{\sigma_k(x_{i \cdot})} \text{Id}_{ij} \sigma_k(x_{i \cdot}) (\text{Id}_{kl} - \sigma_l(x_{i \cdot})) \\ &= -\frac{1}{N} \sum_{k=0}^{K-1} \tilde{y}_{jk} (\text{Id}_{kl} - \sigma_l(x_{j \cdot})) = \\ &= -\frac{1}{N} \tilde{y}_{jl} + \frac{1}{N} \sum_{k=0}^{K-1} \tilde{y}_{jk} \sigma_l(x_{j \cdot}) =\end{aligned}$$

$$= -\frac{1}{N} \tilde{y}_{j,l} + \frac{1}{N} \sigma_l(x_{j,:})$$

$$= \frac{1}{N} (\sigma_l(x_{j,:}) - \tilde{y}_{j,l})$$

Setting $q \in \mathbb{R}^{N \times K}$, $q_{j,l} = \sigma_l(x_{j,:})$,

$$\frac{\partial L}{\partial x} = \frac{1}{N} (q - \tilde{y})^T \in \mathbb{R}^{K \times N}.$$