

1. Linear param. $W = \begin{bmatrix} & \# \text{ in} \\ \# \text{ out} & \end{bmatrix}$ $b = \begin{bmatrix} 1 \\ \# \text{ out} \end{bmatrix}$

input, $\begin{bmatrix} & \text{batch-size} \\ \# \text{ in} & \end{bmatrix}$

$$\text{output} = \begin{bmatrix} W \\ \text{batch-size} \end{bmatrix} \cdot \text{input} + \begin{bmatrix} b, b, b \dots b \end{bmatrix}$$

Compute Gradient

① w.r.t "w"

$$f(w) = W \cdot \overset{\text{input}}{x} + b \Rightarrow \text{grad. } w = \frac{dv_output \cdot \text{input}^T}{\text{batch-size}}$$

$$\therefore \frac{\partial}{\partial w} f(w) = x^T \quad \text{recall } \frac{\partial}{\partial w} aw = a, \quad \frac{\partial}{\partial w} wx = \frac{\partial}{\partial w} (x^T w^T)^T$$

② w.r.t "b"

$$f(b) = w \cdot x + b$$

$$\Rightarrow \frac{\partial f}{\partial b} = dv_output \cdot \begin{bmatrix} 1 \\ \# \text{ batch-size} \end{bmatrix}$$

$$\frac{\partial}{\partial b} f(b) = 1$$

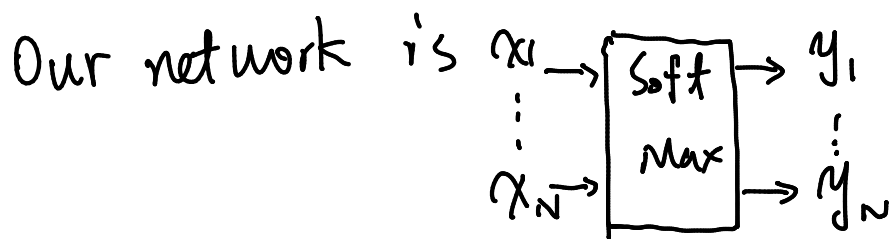
$$\Leftrightarrow \text{sum}(dv_output, 2)$$

$$\text{③ } dv_input = \left(\frac{\partial f}{\partial x} \right)^T dv_output = w^T \cdot dv_output$$

2. Softmax

$$y_i = \frac{e^{x_i}}{\sum_{j=1}^N e^{x_j}}$$

Because Softmax is the last layer, so, it only has $\frac{\partial f(x)}{\partial x_i}$, or say the original dv-output.



We know $\frac{\partial L}{\partial y_i}$, we want $\frac{\partial L}{\partial x_i}$

According to chain rule.

$$\frac{\partial L}{\partial x_i} = \sum_{k=1}^N \frac{\partial L}{\partial y_k} \cdot \frac{\partial y_k}{\partial x_i}$$

for $\frac{\partial}{\partial x_i} y_k$ let $\Gamma = \sum_{j=1}^N e^{x_j}$

* if $i \neq k$

$$\frac{\partial}{\partial x_i} y_k = \frac{\partial}{\partial x_i} \cdot \frac{e^{x_k}}{\sum_{j \neq i} e^{x_j} + e^{x_i}} = - \frac{e^{x_i} \cdot e^{x_k}}{\left(\sum_{j \neq i} e^{x_j} + e^{x_i}\right)^2} = - \frac{e^{x_i + x_k}}{\Gamma^2}$$

* if $i = k$

$$\frac{\partial}{\partial x_i} y_i = \frac{\partial}{\partial x_i} \frac{e^{x_i}}{\sum_{j \neq i} e^{x_j} + e^{x_i}} = \frac{e^{x_i} \Gamma - e^{x_i} \cdot e^{x_i}}{\Gamma^2} = \frac{e^{x_i}}{\Gamma} - \frac{e^{x_i + x_k}}{\Gamma^2}$$

$$\begin{aligned}
\therefore \frac{\partial L}{\partial x_i} &= \sum_{k=1}^N \frac{\partial L}{\partial y_k} \frac{\partial y_k}{\partial x_i} \\
&= \sum_{k \neq i} \frac{\partial L}{\partial y_k} \cdot \left(-\frac{e^{x_i + x_k}}{P^2} \right) + \frac{\partial L}{\partial y_i} \cdot \left(\frac{e^{x_i}}{P} - \frac{e^{x_i + x_i}}{P^2} \right) \\
&= \sum_{k=1}^N \frac{\partial L}{\partial y_k} \left(-\frac{e^{x_i}}{P} \cdot \frac{e^{x_k}}{P} \right) + \frac{\partial L}{\partial y_i} \frac{e^{x_i}}{P}
\end{aligned}$$

3. Cross entropy

$$f(\underline{x}) = - \frac{\sum_{t \in T} \log x_t}{N}$$

T is the set of target label,
or ground truth
 N is batch size.

Also, we don't have $\frac{\partial}{\partial w} f(\underline{x})$

$$\frac{\partial}{\partial x_t} f(\underline{x}) = \frac{\partial}{\partial x_t} - \frac{\sum_{t \in T} \log x_t}{N} = - \frac{1}{x_t \cdot N}$$

4. ReLU

$$\frac{\partial}{\partial x} L = \frac{\partial L}{\partial y} \cdot \mathbb{1}(x \geq 0)$$