

Consider $L(Y) \in \mathbb{R}$, for $Y \in \mathbb{R}^{D \times N}$

Suppose you have $\nabla_Y L \in \mathbb{R}^{D \times N}$.

Let $Y = H \cdot W$ for $H \in \mathbb{R}^{D \times M}$ & $W \in \mathbb{R}^{M \times N}$

Show that

$$\nabla_H L = \nabla_Y L \cdot W^T \quad (1)$$

$D \times M \quad D \times N \quad N \times M$

and

$$\nabla_W L = H^T \cdot \nabla_Y L \quad (2)$$

$M \times N \quad M \times D \quad D \times N$

Pf (of (1))

Consider $\frac{\partial L}{\partial H_{dm}}$

$$= \sum_{n=1}^N \frac{\partial L}{\partial Y_{dn}} \frac{\partial Y_{dn}}{\partial H_{dm}}$$

$$= \left[\frac{\partial L}{\partial Y_{d1}} \dots \frac{\partial L}{\partial Y_{dN}} \right] \begin{bmatrix} \frac{\partial Y_{d1}}{\partial H_{dm}} \\ \vdots \\ \frac{\partial Y_{dN}}{\partial H_{dm}} \end{bmatrix}$$

$$\therefore \frac{\partial L}{\partial H_{dm}} = \nabla_Y L \cdot \begin{bmatrix} W_{m1} \\ \vdots \\ W_{mN} \end{bmatrix}$$

This is where the W^T comes from?

this is the transpose of the m^{th} row of W

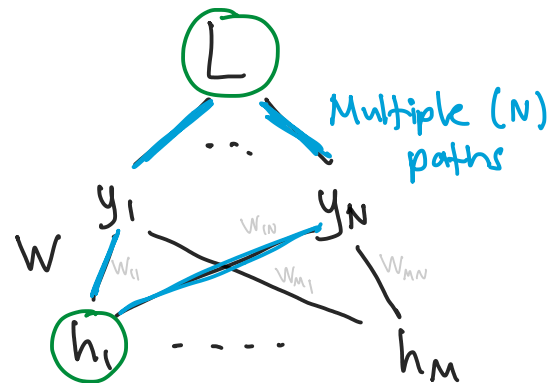
Put all of these, $m=1, \dots, M$, in one row vector...

$$\therefore \left[\frac{\partial L}{\partial H_{d1}} \dots \frac{\partial L}{\partial H_{dM}} \right] = \nabla_Y L \cdot \left[W_{1\cdot} \dots W_{M\cdot} \right] = \nabla_Y L \cdot W^T$$

$m=1 \quad m=M$

M^{th} col. of W^T .

Stacking for $d=1, \dots, D$



Thinking through just 1 node of 1 sample.

$$y_n = \sum_{m=1}^M h_m W_{mn}$$

$$\therefore \frac{\partial y_n}{\partial h_m} = W_{mn}$$

$$\begin{aligned}
 d=1 \rightarrow & \begin{bmatrix} \frac{\partial L}{\partial H_{11}} & \dots & \frac{\partial L}{\partial H_{1M}} \\ \vdots & & \vdots \\ \frac{\partial L}{\partial H_{d1}} & \dots & \frac{\partial L}{\partial H_{dM}} \end{bmatrix} \rightarrow \begin{bmatrix} \nabla_{Y_1} L \\ \vdots \\ \nabla_{Y_d} L \end{bmatrix} \cdot W^T \\
 d=D \rightarrow & \underbrace{\begin{bmatrix} \frac{\partial L}{\partial H_{11}} & \dots & \frac{\partial L}{\partial H_{1M}} \\ \vdots & & \vdots \\ \frac{\partial L}{\partial H_{d1}} & \dots & \frac{\partial L}{\partial H_{dM}} \end{bmatrix}}_{\nabla_H L} = \underbrace{\begin{bmatrix} \nabla_{Y_1} L \\ \vdots \\ \nabla_{Y_d} L \end{bmatrix}}_{\nabla_Y L} \cdot W^T
 \end{aligned}$$

This is how the D samples factor into the formula?

$\therefore \nabla_H L = \nabla_Y L \cdot W^T$

PS (of ②)

Need to show that $\nabla_w L = H^T \cdot \nabla_Y L$

Consider $\frac{\partial L}{\partial w_{ij}}$. w_{ij} is the weight from H_{di} to Y_{dj}

$$Y_{dj} = \sum_{i=1}^M H_{di} w_{ij} \quad \therefore \frac{\partial Y_{dj}}{\partial w_{ij}} = H_{di} \quad \forall j=1, \dots, N$$

$$\text{So, } \frac{\partial L}{\partial w_{ij}} = \sum_{d=1}^D \frac{\partial L}{\partial Y_{dj}} \frac{\partial Y_{dj}}{\partial w_{ij}}$$

Consider a term for a single sample (d):

$$\begin{aligned}
 \begin{bmatrix} \frac{\partial L}{\partial w_{11}} & \dots & \frac{\partial L}{\partial w_{1N}} \\ \vdots & & \vdots \\ \frac{\partial L}{\partial w_{d1}} & \dots & \frac{\partial L}{\partial w_{dN}} \end{bmatrix} &= \begin{bmatrix} \frac{\partial L}{\partial Y_{d1}} \frac{\partial Y_{d1}}{\partial w_{11}} & \dots & \frac{\partial L}{\partial Y_{d1}} \frac{\partial Y_{d1}}{\partial w_{1N}} \\ \vdots & & \vdots \\ \frac{\partial L}{\partial Y_{dj}} \frac{\partial Y_{dj}}{\partial w_{d1}} & \dots & \frac{\partial L}{\partial Y_{dj}} \frac{\partial Y_{dj}}{\partial w_{dN}} \end{bmatrix} = \begin{bmatrix} \frac{\partial L}{\partial Y_{d1}} H_{d1} & \dots & \frac{\partial L}{\partial Y_{d1}} H_{dN} \\ \vdots & & \vdots \\ \frac{\partial L}{\partial Y_{dj}} H_{dj1} & \dots & \frac{\partial L}{\partial Y_{dj}} H_{djN} \end{bmatrix} \\
 &= \begin{bmatrix} H_{d1} \\ H_{d2} \\ \vdots \\ H_{dM} \end{bmatrix} \begin{bmatrix} \frac{\partial L}{\partial Y_{d1}} & \frac{\partial L}{\partial Y_{d2}} & \dots & \frac{\partial L}{\partial Y_{dN}} \end{bmatrix} = h^T \cdot \nabla_Y L \\
 &= (H_d)^T \cdot \nabla_{Y_d} L
 \end{aligned}$$

It's ok if they just took this from the lectures.

They could point out this outer-product for a single sample?
 (h ...)

Adding over d ...

$$\nabla_W L = \sum_{d=1}^D (H_{d.})^T \cdot \nabla_{Y_d.} L$$

(rank-1)

Then combine all the rank-1 products together?

$$= \begin{bmatrix} (H_{1.})^T & \dots & (H_{D.})^T \end{bmatrix} \begin{bmatrix} \nabla_{Y_1.} L \\ \vdots \\ \nabla_{Y_D.} L \end{bmatrix}$$

$$\therefore \nabla_W L = H^T \cdot \nabla_Y L$$



Alternatively, consider the function $\bar{L}(Y^T) = L(Y)$. It's just a version of L that operates on the transpose of Y . Everything is transposed for \bar{L} .

Thus, $\nabla_{Y^T} \bar{L} = (\nabla_Y L)^T$, and $\nabla_{W^T} \bar{L} = (\nabla_W L)^T$.

$$\text{And } \nabla_{W^T} \bar{L}(Y^T) = \nabla_{W^T} \bar{L}(W^T H^T)$$

$$\text{Let } A = W^T, B = H^T.$$

According to the proof of (1), we know

$$\begin{aligned} \nabla_A \bar{L}(AB) &= \nabla_{Y^T} \bar{L} \cdot B^T \\ &= (\nabla_Y L)^T \cdot H \end{aligned}$$

$$\therefore \nabla_{W^T} \bar{L} = (\nabla_Y L)^T \cdot H$$

$$\text{But } \nabla_W L = (\nabla_{W^T} \bar{L})$$

$$= \left[(\nabla_Y L)^T H \right]^T = H^T \cdot \nabla_Y L \text{ as required. } \blacksquare$$

A note about proof design.

The solution must prove (1) and (2).

One of them must be proven outright, but the proof of the second one can use the first (already proven) result.