Example:

$$y_1 = f(x) = x W_1, \quad x \in R^{1 \times D}, \quad W_1 \in R^{D \times F}$$

$$y_2 = g(x) = f(x) W_2, \quad f(x) \in R^{1 \times F}, \quad W_2 \in R^{F \times G}$$

$$y_2 \in R^{1 \times G}$$

$$L = loss = \sum_i y_{2i} \in R$$

How to compute $\frac{\partial L}{\partial W_2} \in R^{F \times G}$, if we already have $\frac{\partial L}{\partial y_2}$

$$\frac{\partial L}{\partial W_2} = f(x)^T \frac{\partial L}{\partial y_2}$$

$$R^{F \times 1} \qquad \frac{\partial L}{\partial y_2} \in R^{1 \times G} \quad [0 \; 0 \; \overset{k-term}{\underset{\downarrow}{0}} \; 0 \; 0] \qquad \text{Sum over rows} \; \downarrow$$

Why? Because:

chain rule: $L(y_1, y_2, y_3, \ldots)$

$$y_{2k} = \sum_r f(x)_r W_{2rk}$$

$$\frac{\partial L}{\partial W_{2ij}} = \sum_k \frac{\partial L}{\partial y_{2k}} \cdot \boxed{\frac{\partial y_{2k}}{\partial W_{2ij}}} \implies \frac{\partial y_{2k}}{\partial W_{2ij}} = \begin{cases} f(x)_i & j=k \\ \\ 0 & j \neq k \end{cases}$$

$$= \sum_k \frac{\partial L}{\partial y_{2k}} \cdot (j=k) \cdot f(x)_i$$

$$\rightarrow \text{only 1 when } j=k \qquad \text{Zero because } W_{2ij} \text{ not in this column}$$

$$= \frac{\partial L}{\partial y_{2j}} f(x)_i$$

Putting together :

$$\frac{\partial L}{\partial W_{2ij}} = \frac{\partial L}{\partial y_{2j}} f(x)_i$$

to matrix

$$\Rightarrow \frac{\partial L}{\partial W_2} = f(X)^T \boxed{\frac{\partial L}{\partial y_2}} \quad \text{from upstream}$$

$$\underbrace{\begin{bmatrix} | \\ | \end{bmatrix}}_{\in R^{F \times 1}} \underbrace{[\quad\quad\quad]}_{\in R^{1 \times G}}$$

---

Extension : What if $f(X) \in R^{B \times F}$ is a batch?
(each row an example)

$$\frac{\partial L}{\partial W_2} = f(\underbrace{X_i}_{\text{row } i})^T \underbrace{\frac{\partial L}{\partial Y_{2i}}}_{\text{row } i} \quad \text{for example } i$$

We need to take the mean of $\frac{\partial L}{\partial i}$ gradient

which is:

$$\frac{\partial L}{\partial W_2} = \frac{1}{B} \sum_i \frac{\partial L}{\partial X_i} = \frac{1}{B} \sum f(X_i)^T \frac{\partial L}{\partial Y_{2i}}$$

$$= \frac{1}{B} f(X)^T \frac{\partial L}{\partial Y_2} \in R^{B \times G}$$

$$\underbrace{\begin{bmatrix} | & | & | & | \end{bmatrix}}_{\text{batch}} \underset{\text{batch}}{\begin{bmatrix} \equiv \end{bmatrix}}$$

To clean up the notation a bit:

$$y = xW \quad, \quad x \in R^{1 \times F} \quad, \quad W \in R^{F \times D}$$
$$, \quad y \in R^{1 \times D}$$

$$L(y) = \text{----.} \quad \text{Scalar}$$
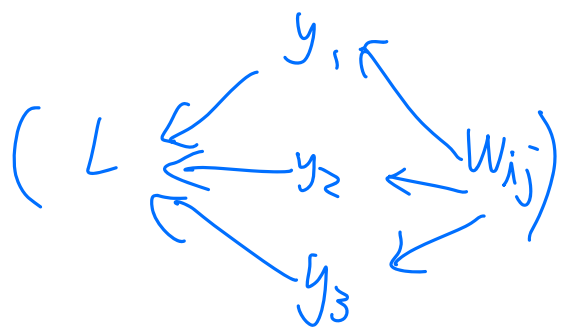
Assume: $\dfrac{\partial L}{\partial y}$ is computed already

$$\frac{\partial L}{\partial W} = \begin{bmatrix} \text{----} \\ \text{---} \frac{\partial L}{\partial w_{ij}} \text{...} \\ \text{----} \end{bmatrix}$$

by chain rule: $L$ is affected by $y_1, y_2 y_3, \text{----} y_D$

$$\frac{\partial L}{\partial w_{ij}} = \sum_{k} \frac{\partial L}{\partial y_k} \cdot \frac{\partial y_k}{\partial w_{ij}}$$

$$\left( L \Leftarrow \begin{array}{c} y_1 \\ y_2 \\ y_3 \end{array} \Leftarrow W_{ij} \right)$$

$$= \sum_{k} \frac{\partial L}{\partial y_k} \cdot \underbrace{(k=j)} \cdot x_r \qquad y_k = \sum_{r} x_r W_{rk}$$

$$= \frac{\partial L}{\partial y_k} \cdot x_r \qquad \Rightarrow \frac{\partial y_k}{\partial w_{ij}} = \begin{cases} x_r & k=j \\ 0 & k \neq j \end{cases}$$

Putting it to vector form

$$\frac{\partial L}{\partial W} = X^T \frac{\partial L}{\partial y}$$

---

If we introduce batch dimension:

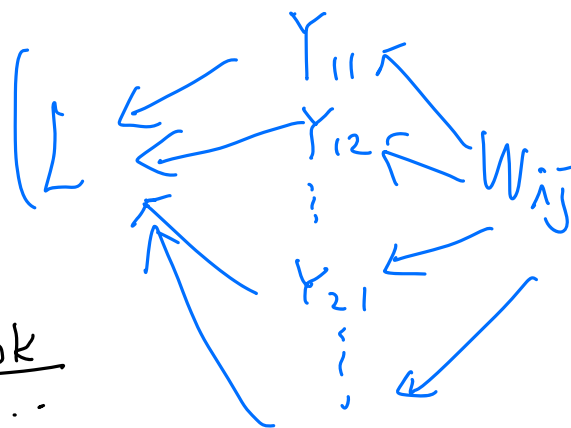$$Y = X W, \quad X \in R^{B \times F}, \quad W \in R^{F \times D},$$
$$L(Y) = scalar \qquad Y \in R^{B \times D}$$

Assume $\frac{\partial L}{\partial Y} \in R^{B \times D}$ is already computed

Compute:

$$\frac{\partial L}{\partial W} = \left[ \begin{array}{c} \cdots \cdots \frac{\partial L}{\partial W_{ij}} \\ \end{array} \right] \quad \text{known}$$

$$\frac{\partial L}{\partial W_{ij}} \overset{\text{chain rule}}{=} \sum_{b=1}^{B} \sum_{k=1}^{D} \boxed{\frac{\partial L}{\partial Y_{bk}}} \frac{\partial Y_{bk}}{\partial W_{ij}}$$

First compute $\frac{\partial Y_{bk}}{\partial W_{ij}}$, to compute this, expand $Y_{bk}$

$$Y_{bk} = X_{b,:} \text{ @ } W_{:,k}$$

$$= \sum_{r=1}^{F} X_{br} W_{rk}$$

$$\frac{\partial Y_{bk}}{\partial W_{ij}} = \frac{\partial}{\partial W_{ij}} \sum_{r=1}^{F} X_{br} W_{rk}$$

$$= \begin{cases} X_{bi} \frac{\partial W_{ij}}{\partial W_{ij}} & k = j \\ 0 & k \neq j \end{cases}$$

Plug this back:

$$\frac{\partial L}{\partial W_{ij}} = \sum_{b=1}^{B} \sum_{k=1}^{D} \frac{\partial L}{\partial Y_{bk}} \cdot (k=j) \cdot X_{bi}$$

when $k \neq j$, terms are zero

$$= \sum_{b=1}^{B} \frac{\partial L}{\partial Y_{bj}} \cdot X_{bi}$$

right    left dim

$B$ is the contracting dim

Putting this to vector: $\dfrac{\partial L}{\partial W} = X^T \dfrac{\partial L}{\partial Y}$

$$X^T \in R^{F \times B}, \frac{\partial L}{\partial Y} \in R^{B \times D}$$