

$$\text{out} = Xw + b \Rightarrow \text{out}_{ij} = X_i w_j + b_j$$

$$\frac{\partial \text{out}_{ij}}{\partial X_i} = w_j$$

$$\frac{\partial \text{out}_{ij}}{\partial w_j} = X_i$$

$$\frac{\partial \text{out}_{ij}}{\partial b_j} = 1$$

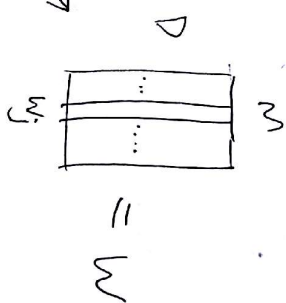
$$\frac{\partial L}{\partial \text{out}} \quad \text{given}$$

$$\frac{\partial L}{\partial X} = \frac{\partial L}{\partial \text{out}} \cdot \frac{\partial \text{out}}{\partial X}$$

We know  $\text{out}_i$  is only computed from  $X_i$ ,  
not other  $X$ 's.

For one data point  $i$

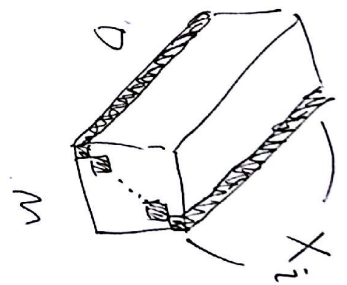
$$\frac{\partial L}{\partial X_i} = \frac{\partial L}{\partial \text{out}_i} \cdot \frac{\partial \text{out}_i}{\partial X_i}$$



To vectorize for all data points,

$$\frac{\partial L}{\partial X} = \frac{\partial L}{\partial \text{out}} \cdot \frac{\partial \text{out}}{\partial X}$$

$$\begin{aligned} \frac{\partial L}{\partial W} &= \frac{\partial L}{\partial out} \cdot \frac{\partial out}{\partial W} \\ \text{For one point } i, \frac{\partial L}{\partial W} &= \frac{\partial L}{\partial out_i} \cdot \frac{\partial out_i}{\partial W} \\ &= \frac{\partial L}{\partial out_i} \cdot \frac{\partial out_i}{\partial W} \cdot \frac{\partial out_i}{\partial W} \end{aligned}$$



$$= \frac{\partial L}{\partial out_i} \times \begin{bmatrix} M \end{bmatrix} \times \begin{bmatrix} D \end{bmatrix}$$

For all data points,

$$\begin{aligned} \frac{\partial L}{\partial W} &= \frac{\partial L}{\partial out} \times \begin{bmatrix} N \times D \times M \end{bmatrix} \\ &\quad \times \begin{bmatrix} M \end{bmatrix} \times \begin{bmatrix} D \end{bmatrix} \end{aligned}$$

take average

$$\frac{\partial L}{\partial W} = \frac{\partial L}{\partial out} \cdot \frac{\partial out}{\partial W}$$

$$\frac{\partial \mathcal{L}}{\partial \mathbf{b}} = \frac{\partial \mathcal{L}}{\partial \mathbf{out}} \frac{\partial \mathbf{out}}{\partial \mathbf{b}}$$

$\boxed{M}$ 
 $\boxed{N \times M}$ 
 $\boxed{(N \times M) \times M}$

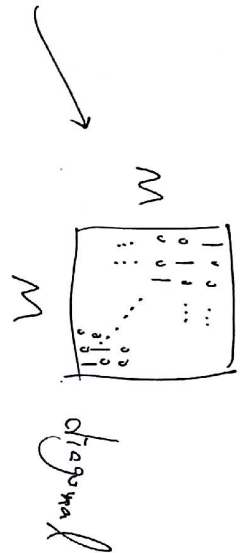
For one data point  $i$ ,

$$\frac{\partial \mathcal{L}}{\partial \mathbf{b}} = \frac{\partial \mathcal{L}}{\partial \mathbf{out}_i} \frac{\partial \mathbf{out}_i}{\partial \mathbf{b}}$$

$\boxed{M}$ 
 $\boxed{M}$ 
 $\boxed{M \times M}$

$$= \frac{\partial \mathcal{L}}{\partial \mathbf{out}_i}$$

$\boxed{M}$



To vectorize for all data points,

$$\frac{\partial \mathcal{L}}{\partial \mathbf{b}} = \frac{\partial \mathcal{L}}{\partial \mathbf{out}}$$

$N \times M$ 
 $N \times M$