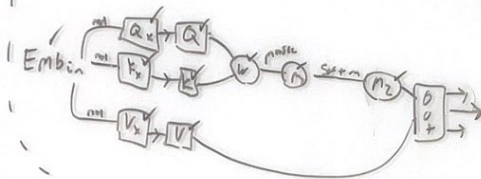


# Deriving Self-Attention Head ( $Q_x, K_x, V_x, Emb_{in}$ )

$$\begin{aligned} \frac{\partial L}{\partial V} &= \left( M^T \frac{\partial L}{\partial out} \right) \cdot \text{Soft}(A) \\ \frac{\partial L}{\partial M_2} &= \frac{\partial L}{\partial out} V^T \\ \frac{\partial L}{\partial M_1} &= \frac{row_x}{\text{Soft} + \text{max}_i \sum \text{Soft} + \text{max}_i} \uparrow \text{col}_x \\ \frac{\partial L}{\partial w} &\rightarrow \text{Self-Attention} \text{ learn } w \end{aligned}$$

## 1 Head Self-Attention



$$\begin{aligned} \frac{\partial L}{\partial K} &= \frac{\partial L}{\partial w} \cdot \frac{\partial w}{\partial Q} = d_{out}^T \frac{\partial Q}{\partial K} \\ \frac{\partial L}{\partial Q} &= \frac{\partial L}{\partial w} \cdot \frac{\partial w}{\partial K} = d_{out}^T \frac{\partial Q}{\partial Q} \end{aligned} \quad Q K^T = w$$

$$\frac{\partial L}{\partial Q_x} = \frac{\partial L}{\partial Q} \cdot \frac{\partial Q}{\partial Q_x} \rightarrow \text{Same for } K_x, V_x$$

$$Q = Emb_{in} \cdot Q_x$$

For this,  $EL_{I,J}$  affects  $j$ th col out by dot + w/  $Emb_{in}$  +  $Q_x$

$$\frac{\partial L}{\partial Q_x} = Emb_{in}^T \frac{\partial L}{\partial Q}$$

$$\frac{\partial L}{\partial K_x} = Emb_{in}^T \frac{\partial L}{\partial K} \rightarrow \text{Affects all Outches, so sum in = 0}$$

$$\frac{\partial L}{\partial V_x} = Emb_{in}^T \frac{\partial L}{\partial V}$$

Embedding affects entire  $i$ th row of output by ext to each element  $j$ th row of matrix. (L-shaped for all  $j$ , so summed)

$$\sum_x \frac{\partial L}{\partial out_{I,x}} \cdot \text{matrix} \uparrow Q_x, K_x, V_x = \frac{\partial L}{\partial V} \cdot V_x^T + \frac{\partial L}{\partial K} \cdot K_x^T + \frac{\partial L}{\partial Q} \cdot Q_x^T$$

$$\begin{bmatrix} Q_1 \\ Q_2 \\ \vdots \end{bmatrix} \cdot \begin{bmatrix} K_1^T \\ K_2^T \\ \vdots \end{bmatrix} = \begin{bmatrix} w_1 \\ w_2 \\ \vdots \end{bmatrix}$$

(B)

$EL_{I,J}$  of  $K$  is  $(EL_{J,I})$  affects  $j$ th col of out by  $I$ th col of  $Q$

$$\frac{\partial L}{\partial EL_{I,J}} = \sum_x Q_{xJ} \cdot \frac{\partial L}{\partial out_{Ix}}$$

dot +  $j$ th col of  $Q$  w/  $i$ th col  $\frac{\partial L}{\partial out}$

$$d_{out}^T \cdot Q$$

Each  $EL_{I,J}$  of  $Q$  affects  $i$ th row of  $\frac{\partial L}{\partial out}$  by  $j$ th col

$$\sum_x K_{Jx}^T \cdot \frac{\partial L}{\partial out_{Ix}}$$

$$d_{out} \cdot w(K^T)^T \cdot K$$