

Linear

One sample: $y = w_1x_1 + w_2x_2 + w_3x_3 + b \quad y \rightarrow \mathcal{L}$

$$y = \begin{bmatrix} x_1 & x_2 & x_3 \end{bmatrix} \begin{bmatrix} w_1 \\ w_2 \\ w_3 \end{bmatrix} + [b]$$

$$\textcircled{1} \quad \frac{d\mathcal{L}}{dw} = \frac{dy}{dy} \cdot \frac{dy}{dw} = \begin{bmatrix} \frac{dy}{dw_1} \\ \frac{dy}{dw_2} \\ \frac{dy}{dw_3} \end{bmatrix} \frac{d\mathcal{L}}{dy}$$

$$= \begin{bmatrix} x_1 \\ x_2 \\ x_3 \end{bmatrix} \frac{d\mathcal{L}}{dy} = X^T \frac{d\mathcal{L}}{dy}$$

↙
a constant
1x1 matrix
side don't matt

$$\textcircled{2} \quad \frac{d\mathcal{L}}{db} = \frac{dy}{dy} \cdot \frac{dy}{db} = \frac{d\mathcal{L}}{dy}$$

$$\textcircled{3} \quad \frac{d\mathcal{L}}{dx} = \frac{d\mathcal{L}}{dy} \cdot \frac{dy}{dx} = \frac{d\mathcal{L}}{dy} [w_1 \ w_2 \ w_3] = \frac{d\mathcal{L}}{dy} W^T$$

Batch Dimension

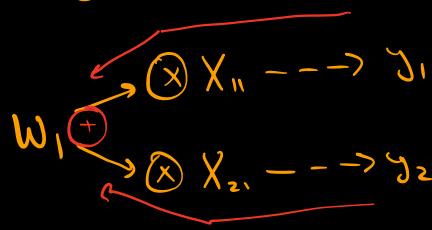
$$\begin{bmatrix} y_1 \\ y_2 \end{bmatrix} = \begin{bmatrix} x_{11} & x_{12} & x_{13} \\ x_{21} & x_{22} & x_{23} \end{bmatrix} \begin{bmatrix} w_1 \\ w_2 \\ w_3 \end{bmatrix} + [b]$$

↙ broadcasting

$$\frac{d\mathcal{L}}{dw} = \begin{bmatrix} \frac{d\mathcal{L}}{dw_1} \\ \frac{d\mathcal{L}}{dw_2} \\ \frac{d\mathcal{L}}{dw_3} \end{bmatrix}$$

[2x1] [1x1]

① $\frac{d\mathcal{L}}{dw_1} = \frac{d\mathcal{L}}{dy} \cdot \frac{dy}{dw_1}$



$$= \sum_{i=1}^N \frac{d\mathcal{L}}{dy_i} \cdot \frac{dy_i}{dw_1} = \frac{d\mathcal{L}}{dy_1} \cdot X_{11} + \frac{d\mathcal{L}}{dy_2} \cdot X_{21}$$

$$\frac{dy}{dw} = \begin{bmatrix} \frac{dL}{dy_1} \cdot X_{11} + \frac{dL}{dy_2} X_{21} \\ \frac{dL}{dy_1} \cdot X_{12} + \frac{dL}{dy_2} X_{22} \\ \frac{dL}{dy_1} \cdot X_{13} + \frac{dL}{dy_2} X_{23} \end{bmatrix}$$

↓ sum over batches contained in the matrix!

$$= \begin{bmatrix} X_{11} & X_{21} \\ X_{12} & X_{22} \\ X_{13} & X_{23} \end{bmatrix} \begin{bmatrix} dL/dy_1 \\ dL/dy_2 \end{bmatrix} = \boxed{\overbrace{X^T \frac{dL}{dy}}^{\sim}}$$

② $\frac{dL}{db}$

b \oplus $\begin{array}{c} \oplus w_1 X_{11} + w_2 X_{12} + w_3 X_{13} \longrightarrow y_1 \\ \oplus w_1 X_{21} + w_2 X_{22} + w_3 X_{23} \longrightarrow y_3 \end{array}$

but $\frac{dy}{db} = \frac{dL}{db} \cdot \frac{dy}{dL} = \frac{dL}{db} \cdot 2 \Rightarrow \boxed{\sum_{i=1}^N \frac{dL}{dy_i}}$

$$\textcircled{3} \quad \frac{dL}{dx} = \frac{dL}{dy} W^T \quad \leftarrow \begin{array}{l} \text{same as before, all the} \\ \text{batch grad contributions} \\ \text{contained inside sum.} \end{array}$$

Softmax (Slow)

More complex as it's a vector function:

$$\vec{o} = [o_1 \ o_2 \ o_3 \ \dots \ o_n]$$

$$\text{Softmax}(\vec{o}) = [s_1 \ s_2 \ s_3 \ \dots \ s_n]$$

Where $s_i = \frac{e^{o_i}}{\sum_{k=1}^n e^{o_k}}$ interdependency b/w some i and j so two cases.

if $i=j$

$$\begin{aligned} \frac{ds_i}{do_j} &= \frac{\left[\sum_{k=1}^n e^{o_k} \right] e^{o_i} - e^{o_i} e^{o_j}}{\left(\sum_{k=1}^n e^{o_k} \right)^2} \\ &= \frac{e^{o_i} \left[\sum_{k=1}^n e^{o_k} - e^{o_j} \right]}{\left(\sum_{k=1}^n e^{o_k} \right)^2} \end{aligned}$$

$$= \frac{e^{O_i}}{\sum_{k=1}^n e^{O_k}} \cdot \frac{\sum_{k=1}^n e^{O_k} - e^{O_j}}{\sum_{k=1}^n e^{O_k}}$$

$\sim s_i (1 - s_j)$

if $i \neq j$

$$\frac{dS_i}{dO_j} = \frac{\left[\sum_{k=1}^n e^{O_k} \right] O - e^{O_i} e^{O_j}}{\left[\sum_{k=1}^n e^{O_k} \right]^2}$$

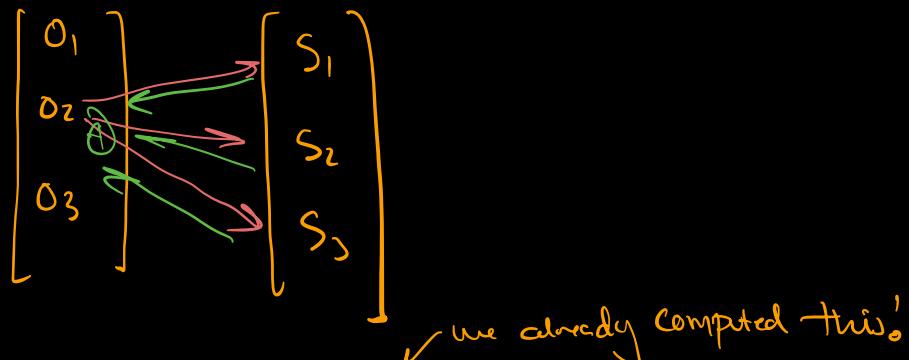
$$= - \frac{e^{O_i}}{\sum_{k=1}^n e^{O_k}} \frac{e^{O_j}}{\sum_{k=1}^n e^{O_k}} = - s_i s_j$$

$$J = \begin{bmatrix} s_1(1-s_1) & -s_1 s_2 & -s_1 s_3 \\ -s_2 s_1 & s_2(1-s_2) & -s_2 s_3 \\ -s_3 s_1 & -s_3 s_2 & s_3(1-s_3) \end{bmatrix}$$

↑

Numerical Stability	
$\frac{e^{x-c}}{\sum e^{x-c}} = \frac{e^x/e^c}{(\sum e^x)/e^c}$	Per sample jacobian! if N is large this is a massive N^2 matrix!
Can we do better?	

What we want: $\frac{dL}{d\theta_j}$ and θ_j has grad contributions from all s_j



$$\frac{dL}{d\theta_j} = \sum_i \frac{dL}{ds_i} \cdot \underbrace{\frac{ds_i}{d\theta_j}}_{\text{we already computed this!}}$$

$$= \sum_i \frac{dL}{ds_i} \left[s_i [\delta_{ij} - s_j] \right]$$

split sum
into 2
cases

Where $\delta_{ij} = 1$ if $i=j$ else 0.

$$= \frac{dL}{ds_j} [s_j(1-s_j)] + \sum_{\substack{i \\ i \neq j}} \frac{dL}{ds_i} (-s_i s_j)$$

$$= \frac{dL}{ds_j} s_j - \frac{dL}{ds_j} s_j s_j - \sum_{\substack{i \\ i \neq j}} \frac{dL}{ds_i} (-s_i s_j)$$

$$= \frac{dL}{ds_j} s_j - \sum_i \frac{dL}{ds} (s_i s_j) \quad \begin{matrix} \leftarrow & \text{sum does} \\ & \text{depend on } s_j \end{matrix}$$

$$= S_j \left[\frac{dL}{ds_j} - \sum_i \underbrace{\frac{dL}{ds_i}}_{\text{dot prod.}} s_i \right]$$



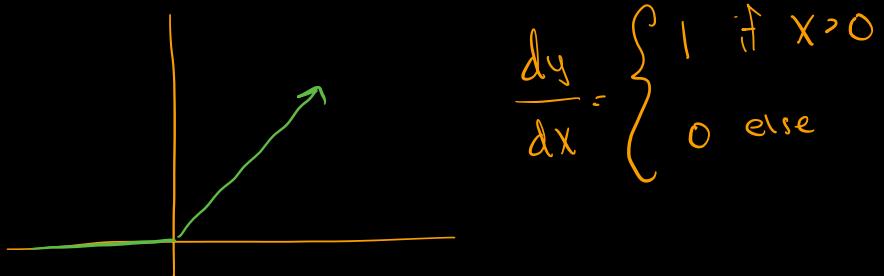
Vector notation

$$S \cdot \left[\frac{dL}{ds} - \left[\frac{dL}{ds} \bullet S \right] \right]$$



Relu

$$\begin{bmatrix} -2 \\ 3 \\ 8 \end{bmatrix} \rightarrow \text{Relu}(x) \rightarrow \begin{bmatrix} 0 \\ 3 \\ 8 \end{bmatrix}$$



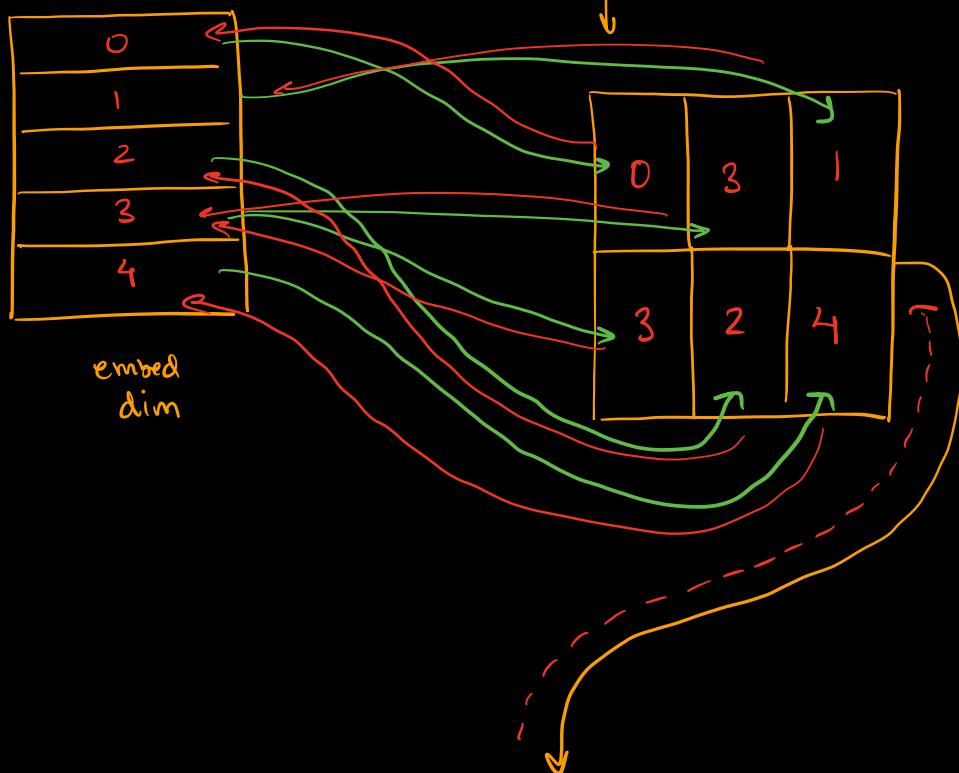
Embeddings

3

$$\begin{bmatrix} 0 & 3 & 1 \\ 3 & 2 & 4 \end{bmatrix}$$

Words 5

embed dim



Transformer

↑ ↓
L

Layer Norm

Data: $B \times S \times E$

normalize along last dimension E .

$$M_{b,s} = \frac{1}{E} \sum_{e=1}^E X_{b,s,e} \quad \leftarrow \text{w/ } X \text{ interdependency}$$

$$\text{Var}_{b,s} = \frac{1}{E} \sum_{e=1}^E (X_{b,s,e} - M_{b,s})^2 \quad \leftarrow$$

$$\hat{X}_{b,s,e} = \frac{X_{b,s,e} - M_{b,s}}{\sqrt{\text{Var}_{b,s} + \epsilon}}$$

$$y_{b,s,e} = \gamma_e \cdot \hat{X}_{b,s,e} + \beta_e$$

↙ learnable scale param ↙ learnable shift param .

given upstream grad $\frac{\partial L}{\partial y_{b,s,e}}$ we need

$$\frac{\partial L}{\partial \gamma_e}, \frac{\partial L}{\partial \beta_e}, \frac{\partial L}{\partial \hat{X}_{b,s,e}}$$

Remember, our B, S indexes are independent,
 \therefore the grad are accumulated. For simplicity
I will ignore them.

$$\bar{M} = \frac{1}{E} \sum_{k=1}^E X_k$$

$$U \stackrel{\Delta}{=} X_k - \bar{M} \leftarrow \text{mean centered.}$$

$$\sigma^2 = \frac{1}{E} \sum_{k=1}^E U_k^2$$

$$S = \sqrt{\sigma^2 + \epsilon} \quad \begin{matrix} \downarrow \\ \text{has } X \text{ in it} \end{matrix}$$

$$\hat{X}_j = \frac{U_j}{S} = \frac{X_j - \bar{M}}{S} \quad \begin{matrix} \downarrow \\ \text{has } X \text{ in it} \end{matrix}$$

$$y_j = \gamma_j \hat{X}_j + \beta_j$$

$$\nabla_j = \frac{\partial L}{\partial y_j} \leftarrow \text{upstream gradient.}$$

$$\frac{\partial L}{\partial x_i} = \sum_{j=1}^E \frac{\partial L}{\partial y_j} \cdot \frac{\partial y_j}{\partial x_i} \quad \begin{matrix} \leftarrow & \text{sum up all} \\ & \text{those contributions.} \end{matrix}$$

$$= \sum_{j=1}^E \nabla_j \underbrace{\frac{d\gamma_j}{d\hat{x}_j}}_{\delta_j} \cdot \frac{d\hat{x}_j}{dx_i}$$

$$= \sum_{j=1}^E \nabla_j \delta_j \underbrace{\frac{d\hat{x}_j}{dx_i}}_{\text{solve this:}}$$

$$\mathcal{M} = \frac{1}{E} \sum_{k=1}^E x_k \rightarrow \frac{d\mathcal{M}}{dx_i} = \frac{1}{E}$$

$$u_j = x_j - \mu \rightarrow \frac{du_j}{dx_i} = \delta_{ij} - \frac{1}{E}$$

1 if $i=j$ else 0.

$$\sigma^2 = \frac{1}{E} \sum_{k=1}^E u_k^2$$

$$\frac{d\sigma^2}{dx_i} = \frac{1}{E} \sum_{k=1}^E 2u_k \cdot \frac{du_k}{dx_i}$$

$$= \frac{2}{E} \sum_{k=1}^E u_k \left\{ \delta_{ik} - \frac{1}{E} \right\}$$

$$= \frac{2}{E} \left[\sum_{k=1}^E u_k \delta_{ik} - \underbrace{\sum_{k=1}^E u_k \frac{1}{E}}_{0 \text{ if } i \neq k} \right]$$

$$= \frac{2}{E} \left[u_i - \underbrace{\frac{1}{E} \sum_{k=1}^E u_k}_{\text{mean}} \right]$$

$$\frac{1}{E} \sum_{k=1}^E u_k = \frac{1}{E} \sum_{k=1}^E x_k - \bar{M} = 0$$

Mean centered data.
So mean of mean
centered is 0!

$$\therefore \boxed{\frac{\partial \sigma^2}{\partial x_i} = \frac{2}{E} u_i}$$

$$S = \sqrt{\sigma^2 + \epsilon} = (\sigma^2 + \epsilon)^{1/2} \quad \text{constant}.$$

$$\frac{dS}{dx_i} = \frac{1}{2} (\sigma^2 + \epsilon)^{-\frac{1}{2}} \cdot \frac{d(\sigma^2 + \epsilon)}{dx_i}$$

$$= \frac{1}{2\sqrt{\sigma^2 + \epsilon}} \sum_{E} u_i$$

$$= \frac{1}{S} \sum_{E} u_i$$

$$\therefore \left\{ \begin{array}{l} \frac{dS}{dx_i} = \frac{u_i}{SE} \\ \end{array} \right.$$

$$\hat{x}_j = \frac{u_j}{S} = u_j \left(\frac{1}{S} \right) \quad \text{Product rule.}$$

$$\frac{d\hat{x}_j}{dx_i} = u_j \underbrace{\frac{d(1/S)}{dx_i}}_1 + \underbrace{\frac{du_j}{dx_i}}_2 \frac{1}{S}$$

$$\textcircled{1} \quad u_j \frac{d\left(\frac{1}{s}\right)}{dx_i} = u_j \frac{d}{dx_i} s^{-1} = -s^{-2} \frac{ds}{dx_i} u_j$$

$$= -\frac{1}{s^2} u_j \frac{u_i}{sE} = -\frac{u_i u_j}{s^3 E}$$

$$\textcircled{2} \quad \frac{du_i}{dx_i} \frac{1}{s} = \left(\delta_{ij} - \frac{1}{E} \right) \left(\frac{1}{s} \right)$$

$$\therefore \frac{d\hat{x}_i}{dx_i} = -\frac{u_i}{s^3 E} + \left(\delta_{ij} - \frac{1}{E} \right) \left(\frac{1}{s} \right)$$

$$= \frac{1}{s} \left[\delta_{ij} - \frac{1}{E} - \frac{u_i u_j}{s^2 E} \right] \hat{x}_i \frac{u}{s}$$

$$= \frac{1}{s} \left[\delta_{ij} - \frac{1}{E} - \underbrace{\frac{u_i u_j}{s^2} \frac{1}{E}}_{\hat{x}_i \hat{x}_j} \right]$$

$$= \frac{1}{s} \left[\delta_{ij} - \frac{1}{E} - \frac{\hat{x}_i \hat{x}_j}{E} \right]$$

Now remember that: we just solved this.

$$\frac{\partial L}{\partial x_i} = \sum_{j=1}^E \nabla_j \delta_j \frac{\partial \hat{x}_j}{\partial x_i}$$

$$= \sum_{j=1}^E \nabla_j \delta_j - \frac{1}{S} \left[\delta_{ij} - \frac{1}{E} - \frac{\hat{x}_i \hat{x}_j}{E} \right]$$

$$= \frac{1}{S} \left[\sum_{j=1}^E \nabla_j \delta_j \delta_{ij} \stackrel{j}{\downarrow} 1 \text{ if } i=j \right]$$

$$- \frac{1}{E} \sum_{j=1}^E \nabla_j \delta_j$$

$$- \frac{1}{E} \sum_{j=1}^E \nabla_j \delta_j \hat{x}_i \hat{x}_j \right]$$

$$= \frac{1}{S} \left[\nabla_i \delta_i - \frac{1}{E} \sum_{j=1}^E \nabla_j \delta_j - \frac{\hat{x}_i}{E} \sum_{j=1}^E \nabla_j \delta_j \hat{x}_j \right]$$

Now remember

$$\nabla = \frac{dL}{dy} \quad \text{and} \quad \frac{d\mathcal{Y}}{d\hat{x}} = \mathcal{D}$$

so for simplicity, our upstream grad to \hat{x}

(not just y) is:

$$\frac{\partial L}{\partial y} \cdot \frac{\partial y}{\partial \hat{x}} = \nabla \mathcal{D} \triangleq \tilde{\nabla} \quad \begin{matrix} \leftarrow \text{let just give it} \\ \text{a new symbol} \end{matrix}$$

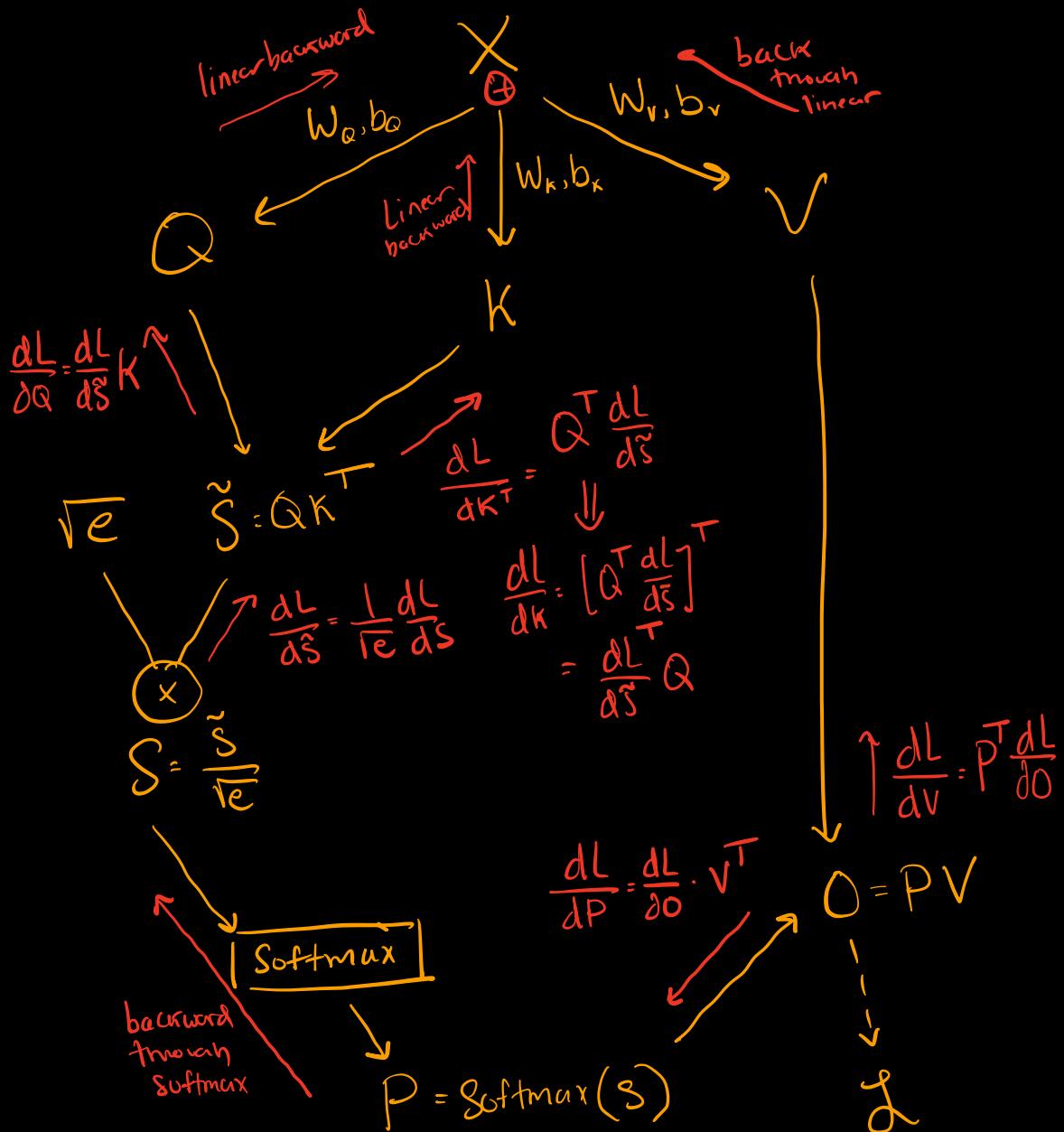
$$= \frac{1}{S} \left[\tilde{\nabla} - \frac{1}{E} \sum_{j=1}^E \tilde{\nabla}_j - \frac{\hat{x}_i}{E} \underbrace{\sum_{j=1}^E \tilde{\nabla}_j \hat{x}_j}_{\text{dot product}} \right]$$

$$= \boxed{\frac{1}{\sqrt{\sigma^2 + \epsilon}} \left[\tilde{\nabla} - \text{mean}(\tilde{\nabla}) - \hat{x} \cdot \text{mean}(\tilde{\nabla} \cdot \hat{x}) \right]}$$

$$\frac{\partial L}{\partial \beta_j} = \sum_{y,s} \frac{\partial L}{\partial y} \cdot \frac{\partial y}{\partial \beta_j} = \sum_{y,s} \nabla$$

$$\frac{\partial L}{\partial \delta_j} = \sum_{y,s} \frac{\partial L}{\partial y} \cdot \frac{\partial y}{\partial \delta_j} = \sum_{y,s} \nabla \hat{x}_j$$

Multihead Attention



Cross Entropy Loss (For OHE Labels)

$$J(P, y) = - \sum_{i=1}^n y_i \log(P_i)$$

$\underbrace{\qquad\qquad\qquad}_{P_i \text{ is a probability (softmaxed)}}$

y is one hot:

$$\begin{bmatrix} 0 & 0 & 1 & 0 & 0 \end{bmatrix}$$

$$J(\cdot, y) = - \left[0 \log(P_1) + 0 \log(P_2) + 1 \log(P_3) + 0 \log(P_4) + 0 \log(P_5) \right]$$

$$= -\log(P_3)$$

$$\therefore \frac{dJ}{dP_i} = -\frac{y_i}{P_i} = \begin{cases} 0 & \text{if } i \text{ is not correct} \\ -\frac{1}{P_i} & \text{if } i \text{ is correct} \end{cases}$$

logits \rightarrow Softmax \rightarrow loss



Problem. Previous derivation expects already softmaxed outputs. This means in the backward pass we have to backprop through softmax.

Can we avoid that?

$$CE = - \sum_j y_j \log(P_j)$$

$$\text{Where } P_j = \frac{e^{x_j}}{\sum_i e^{x_i}}$$

but y is OHE so

$$CE = -\log(P_{\text{label}})$$

$$= -\log \left[\frac{e^{x_{\text{label}}}}{\sum_i e^{x_i}} \right]$$

$$= - \left[\log(e^{x_{\text{label}}}) - \log \sum_i e^{x_i} \right]$$

$$= \log \sum_i e^{x_i} - \log(e^{x_{\text{true}}})$$

$$= \log \sum_i e^{x_i} - x_{\text{true}}$$

$$\frac{\partial L}{\partial x_j} \left[\log \sum_i e^{x_i} - x_{\text{label}} \right]$$

$$= \frac{1}{\sum_i e^{x_i}} \frac{\partial}{\partial x_j} \sum_i e^{x_i} - \begin{cases} 1 & \text{if } j = \text{label} \\ 0 & \text{else} \end{cases}$$

$$= \frac{e^{x_j}}{\sum_i e^{x_i}} - \begin{cases} 1 & \text{if } j = \text{label} \\ 0 & \text{else} \end{cases}$$

\sim

Our grad uses our probs!

$$= p_j - \delta_j \quad \text{where} \quad \delta_j = \begin{cases} 1 & \text{if } j = \text{label} \\ 0 & \text{else.} \end{cases}$$

logits \rightarrow loss