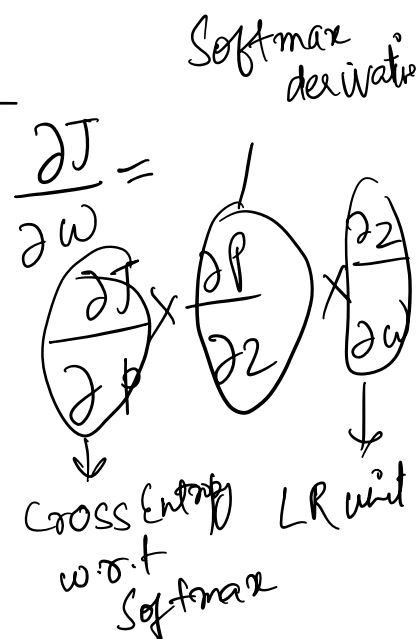


Backpropagation

① Softmax derivative

$$\text{softmax}_i = \frac{e^{x_i}}{\sum_j e^{x_j}}$$



What would be derivative of softmax w.r.t inputs

Jacobian →

$$\nabla_{\text{Softmax}} = \begin{bmatrix} \frac{\partial S_1}{\partial x_1} & \frac{\partial S_1}{\partial x_2} & \dots & \frac{\partial S_1}{\partial x_n} \\ \frac{\partial S_2}{\partial x_1} & \frac{\partial S_2}{\partial x_2} & \dots & \frac{\partial S_2}{\partial x_n} \\ \vdots & \vdots & \ddots & \vdots \\ \frac{\partial S_n}{\partial x_1} & \dots & \dots & \frac{\partial S_n}{\partial x_n} \end{bmatrix}$$

Lets look at s_i

$$s_i = \frac{e^{x_i}}{\sum_j e^{x_j}}$$

Lets take \log both sides

$$\log s_i = \log \left(\frac{e^{x_i}}{\sum e^{x_j}} \right)$$

$$= \log(e^{x_i}) - \log(\sum e^{x_j})$$

$$\log s_i = x_i - \log(\sum e^{x_j})$$

Differentiate w.r.t any x_k

$$\frac{\partial \log s_i}{\partial x_k} = \frac{\partial x_i}{\partial x_k} - \frac{\partial \log (\sum e^{x_i})}{\partial x_k}$$



$$\frac{\partial x_i}{\partial x_k} = \begin{cases} 1 & x_i = x_k \\ 0 & \text{otherwise} \end{cases} = \mathbb{1}(x_i = x_k)$$

$$= \delta_{ik} \text{ where } \mathbb{1}(x_i = x_k) \text{ (Indicator fun}^n)$$

2nd term

$$\neq \frac{\partial \log (\sum e^{x_i})}{\partial x_k} \left[\begin{array}{l} \frac{\partial \log g(x)}{\partial x} \\ = \frac{1}{g(x)} \frac{\partial g(x)}{\partial x} \end{array} \right]$$

$$\Rightarrow \frac{1}{\sum e^{x_i}} \left(\frac{\partial \sum e^{x_i}}{\partial x_k} \right)$$

$$= \frac{1}{\sum e^{x_i}} \sum \frac{\partial e^{x_i}}{\partial x_k}$$

$$\delta_{jk} = \begin{cases} 1 & x_i = x_k \\ 0 & \text{otherwise} \end{cases}$$

$$\Rightarrow \frac{1}{\sum e^{x_i}} \sum_{i=1}^C e^{x_i} \delta_{jk}$$

$$\Rightarrow \frac{e^{x_k}}{\sum e^{x_i}} = S_k$$

So final equation becomes

$$\frac{\partial \log S_i}{\partial x_k} = \delta_{ik} - S_k$$

$\delta_{ik} = \text{Indicator function}$
 $\delta = \begin{cases} 1 & i=k \\ 0 & i \neq k \end{cases}$

$$\frac{1}{s_i} \frac{\partial s_i}{\partial x_k} = \delta_{ik} - s_k$$

$$\frac{\partial s_i}{\partial x_k} = s_i (\delta_{ik} - s_k)$$

So our final Jacobian becomes

$$J_{\text{softmax}} = \begin{pmatrix} s_1(1-s_1) & -s_1s_2 & \dots & -s_1s_n \\ -s_2s_1 & s_2(1-s_2) & \dots & -s_2s_n \\ \vdots & \vdots & \ddots & \vdots \\ -s_ns_1 & \dots & \dots & s_n(1-s_n) \end{pmatrix}$$

② Differentiating Cross-Entropy loss

$$L(s, y) = - \sum_{i=1}^C y_i \log s_i$$

$$\frac{\partial L(s, y)}{\partial x_k} = - \sum_{i=1}^C y_i \frac{\partial \log s_i}{\partial x_k}$$

$$\Rightarrow - \sum_{i=1}^C \left(\frac{y_i}{s_i} \frac{\partial s_i}{\partial x_k} \right)$$

$$= - \sum_{i=1}^C \left(\frac{y_i}{s_i} s_i (\delta_{ik} - s_k) \right)$$

$$= - \sum_{i=1}^C \left(y_i (\delta_{ik} - s_k) \right)$$

$$= - \sum_{i=1}^C (y_i \delta_{ik} - y_i S_k)$$

\Rightarrow

When $i=k$, only then
first term will become y_k

$$\frac{\partial L(S, y)}{\partial x_k} = -y_k + \sum_{i=1}^C y_i S_k$$

$$\Rightarrow \sum_{i=1}^C y_i = 1 \quad (\text{Since it's one-hot})$$

Hence

$$\frac{\partial L(S, y)}{\partial x_k} \Rightarrow -y_k + S_k$$

0.1.2

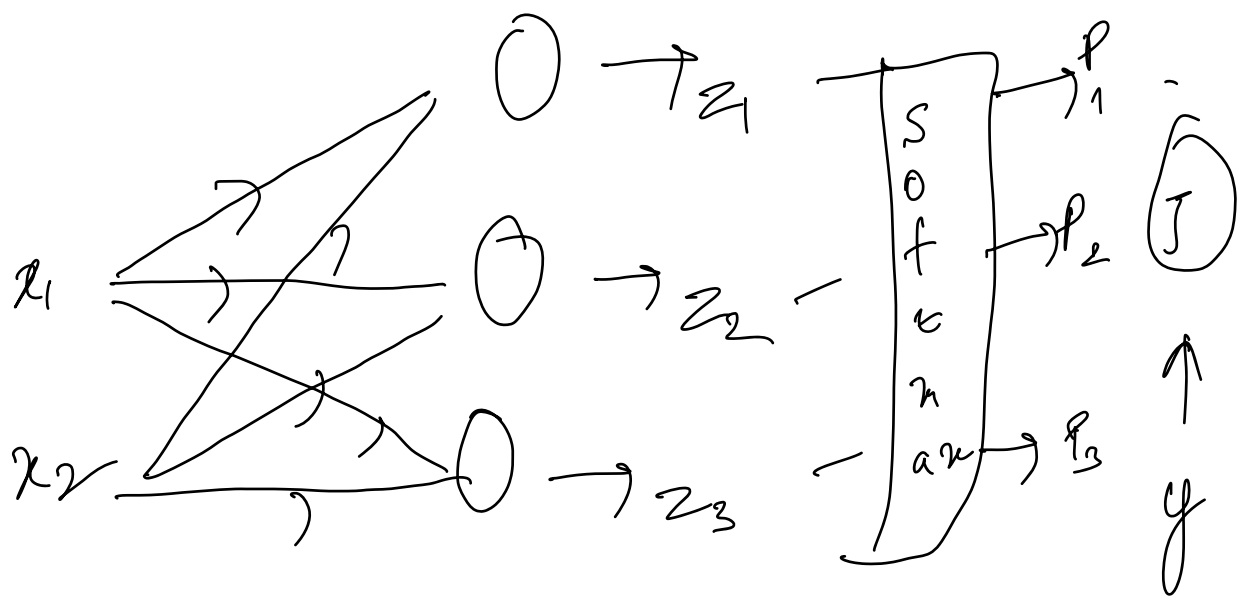
$$\frac{\partial L(s, y)}{\partial x_k} = s_k - y_k$$

In vectorized form this becomes

$$[s_1 - y_1, s_2 - y_2, \dots]$$

$$\frac{\partial L(x, y)}{\partial x} = s - y$$

→ So now in our original chain equation



$$\frac{\partial J}{\partial w} = \frac{\partial J}{\partial p} \frac{\partial p}{\partial z} \left(\frac{\partial z}{\partial w} \right) ?$$

$$\frac{\partial z}{\partial w} = x$$