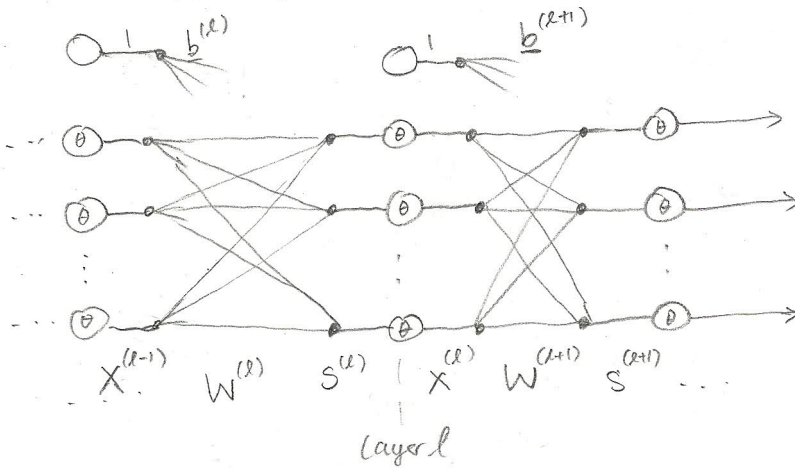


VECTORIZING THE NEURAL NETWORK

Model of L-layer NN



$$\text{Let } X^{(l)} = \begin{bmatrix} X_1^{(l)T} \\ \vdots \\ X_N^{(l)T} \end{bmatrix} \quad Y = \begin{bmatrix} y_1^T \\ \vdots \\ y_N^T \end{bmatrix}$$

$$\underline{b} = [b_1, \dots, b_k]$$

one-hot encoded

X_{ij} corresponds to the j^{th} feature of the i^{th} datapoint

$\dot{\underline{X}}$ indicates broadcasting.

FORWARD PROPAGATION VECTORIZATION

$$S_{nj}^{(l)} = \sum_i X_{ni}^{(l-1)} W_{ij}^{(l)} + b_j^{(l)}$$

$$X_{nj}^{(l)} = \theta(S_{nj}^{(l)})$$

$$\boxed{S^{(l)} = X^{(l-1)} W^{(l)} + \dot{\underline{b}}^{(l)}}$$

$$\boxed{X^{(l)} = \theta(S^{(l)})}$$

activation function of the l^{th} layer. (ie. ReLU or Softmax)

GRADIENT VECTORIZATION

Let $\delta_{nj}^{(l)} = \frac{\partial E_{in}}{\partial S_{nj}^{(l)}}$ the sensitivity of the i^{th} datapoint at the node j layer l .

$$\frac{\partial E_{in}}{\partial W_{ij}^{(l)}} = \sum_n \frac{\partial E_{in}}{\partial S_{nj}^{(l)}} \cdot \frac{\partial S_{nj}^{(l)}}{\partial W_{ij}^{(l)}}$$

(This is b/c $E_{in} = \frac{1}{N} \sum_n E_n = f(S_1^{(l)}, S_2^{(l)}, \dots, S_N^{(l)})$ which means we need to apply a branching chain rule)

$$= \sum_n \delta_{nj}^{(l)} \cdot X_{ni}^{(l-1)} = \sum_n (X_{in}^{(l-1)})^T \cdot \delta_{nj}^{(l)} \Rightarrow$$

$$\boxed{\frac{\partial E_{in}}{\partial W^{(l)}} = [X^{(l-1)}]^T \delta^{(l)}}$$

$$\frac{\partial E_{in}}{\partial b_j^{(l)}} = \sum_n \frac{\partial E_{in}}{\partial S_{nj}^{(l)}} \cdot \frac{\partial S_{nj}^{(l)}}{\partial b_j^{(l)}}$$

$$= \sum_n \delta_{nj}^{(l)} \cdot 1$$

\Rightarrow

$$\boxed{\frac{\partial E_{in}}{\partial \underline{b}^{(l)}} = \underline{1}^T \delta^{(l)}}$$

BACKPROPAGATION VECTORIZATION

$$\begin{aligned}\delta^{(L)} &= \frac{\partial E_{in}}{\partial S_{nj}^{(L)}} = \frac{\partial E_{in}}{\partial X_{nj}^{(L)}} \cdot \frac{\partial X_{nj}^{(L)}}{\partial S_{nj}^{(L)}} \\ &= \left(\sum_k \frac{\partial E_{in}}{\partial S_{nk}^{(L+1)}} \cdot \frac{\partial S_{nk}^{(L+1)}}{\partial X_{nj}^{(L)}} \right) \cdot \theta'(S_{nj}^{(L)}) \\ &= \left(\sum_k \delta_{nk}^{(L+1)} w_{jk}^{(L+1)} \right) \cdot \theta'(S_{nj}^{(L)}) = \left(\sum_k \delta_{nk}^{(L+1)} (w_{kj}^{(L+1)})^T \right) \cdot \theta'(S_{nj}^{(L)})\end{aligned}$$

$$\Rightarrow \boxed{\delta^{(L)} = \left(\delta^{(L+1)} [W^{(L+1)}]^T \right) \otimes \theta'(S^{(L)})}$$

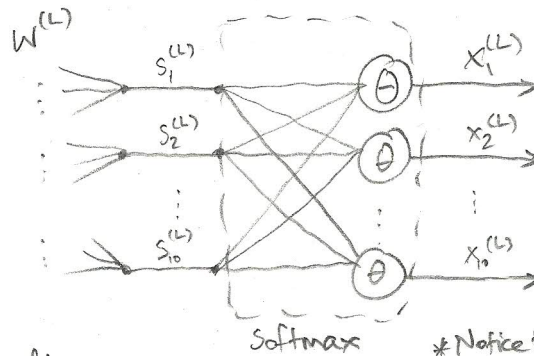
in the hidden layer for Relu.
 $\theta'(S^{(L)}) = \text{sign}(S^{(L)})$

SEEDING the SENSITIVITY

What is $\delta^{(L)}$? $L = \text{output layer}$.

$$\delta_{nj}^{(L)} = \frac{\partial E_{in}}{\partial S_{nj}^{(L)}} = \sum_k \frac{\partial E_{in}}{\partial X_{nk}^{(L)}} \cdot \frac{\partial X_{nk}^{(L)}}{\partial S_{nj}^{(L)}}$$

(Note: unlike for hidden layer activation, the softmax depends on all $S_{nj}^{(L)}$ $\rightarrow \therefore$ we need branching chain rule)



* Notice that the softmax activation depends on all $S_j^{(L)}$

$$E_{in} = -\frac{1}{N} \sum_{n=1}^N \sum_{k=1}^K y_{nk} \ln(X_{nk}^{(L)}) \quad \text{for part 1.1.5}$$

$$\frac{\partial E_{in}}{\partial X_{nk}^{(L)}} = -\frac{1}{N} \frac{y_{nk}}{X_{nk}^{(L)}} \Rightarrow \frac{\partial E_{in}}{\partial X^{(L)}} = -\frac{1}{N} \frac{Y}{X^{(L)}}$$

$$\frac{\partial X_{nk}^{(L)}}{\partial S_{nj}^{(L)}} = \frac{\partial}{\partial S_{nj}^{(L)}} \left(\sigma(S_{nk}^{(L)}) \right) = X_{nk}^{(L)} (\delta_{kj} - X_{nj}^{(L)})$$

\uparrow dirac delta.

$$\delta_{nj}^{(L)} = \sum_k \frac{\partial E_{in}}{\partial X_{nk}^{(L)}} \cdot \frac{\partial X_{nk}^{(L)}}{\partial S_{nj}^{(L)}}$$

$$= -\frac{1}{N} \sum_k \frac{y_{nk}}{X_{nk}^{(L)}} \cdot X_{nk}^{(L)} (\delta_{kj} - X_{nj}^{(L)})$$

$$= -\frac{1}{N} \sum_k (y_{nk} \delta_{kj} - y_{nk} X_{nj}^{(L)})$$

$$\Rightarrow = -\frac{1}{N} \left(\sum_k y_{nk} \delta_{kj} - X_{nj}^{(L)} \sum_k y_{nk} \right)$$

$$\Rightarrow \delta^{(L)} = -\frac{1}{N} (Y \cdot I - X^{(L)})$$

is always = 1 b/c
one-hot encoded.

$$\boxed{\delta^{(L)} = \frac{1}{N} (X^{(L)} - Y)}$$

For the specific 2-Layer Neural Network we are asked to implement...

$$\left\{ \begin{array}{l} \frac{\partial L}{\partial w_0} = \frac{\partial E_{in}}{\partial w^{(2)}} = [X^{(1)}]^T \delta^{(2)} \\ \frac{\partial L}{\partial b_0} = \frac{\partial E_{in}}{\partial b^{(2)}} = \underline{1}^T \delta^{(2)} \end{array} \right.$$

$$\left\{ \begin{array}{l} \frac{\partial L}{\partial w_n} = \frac{\partial E_{in}}{\partial w^{(1)}} = [X^{(0)}]^T \delta^{(1)} \\ \frac{\partial L}{\partial b_n} = \frac{\partial E_{in}}{\partial b^{(1)}} = \underline{1}^T \delta^{(1)} \end{array} \right.$$

$$\left\{ \begin{array}{l} \delta^{(2)} = \frac{1}{N} (X^{(2)} - Y) \\ \delta^{(1)} = (\delta^{(2)} [W^{(2)}]^T) \otimes \text{Sign}(s^{(1)}) \end{array} \right.$$