

Neural Network

Inputs : (x, y) ; $x \in \mathbb{R}^d$, $y \in \mathbb{R}^m$ where ' d ' is the number of features and ' m ' is the number of classes in the final output.

$$x = [x_1, x_2, \dots, x_m - x_n] \quad y = [y_1, y_2, \dots, y_n]$$

There are N number of examples [tuple (x, y)] but for simplicity we will just focus on one example i.e $(x, y) \leftarrow$ one data point.

Notation

For inputs $x_i^l \leftarrow$ layer
 $x_i^l \leftarrow$ it value in vector

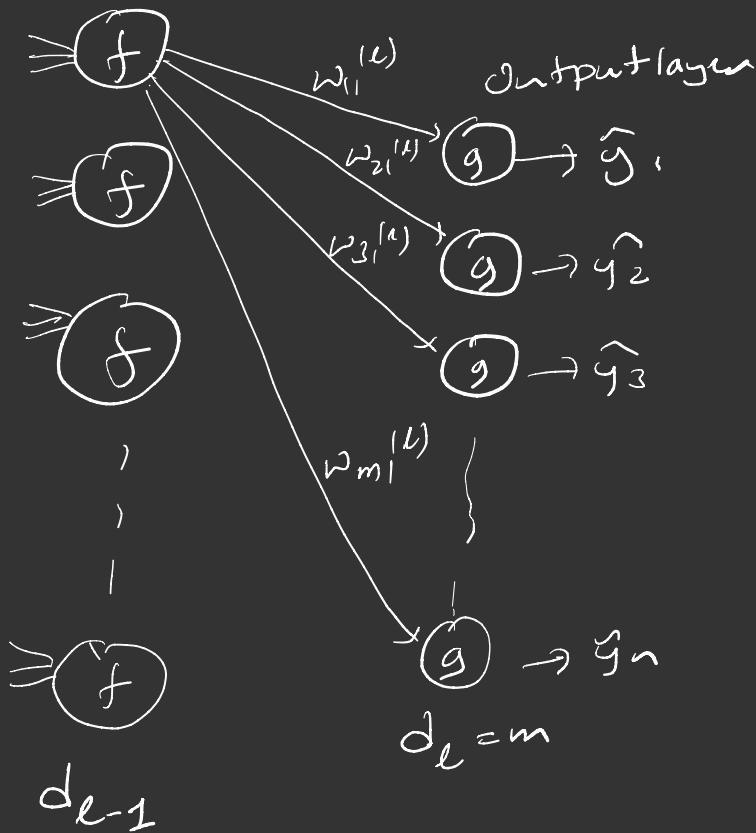
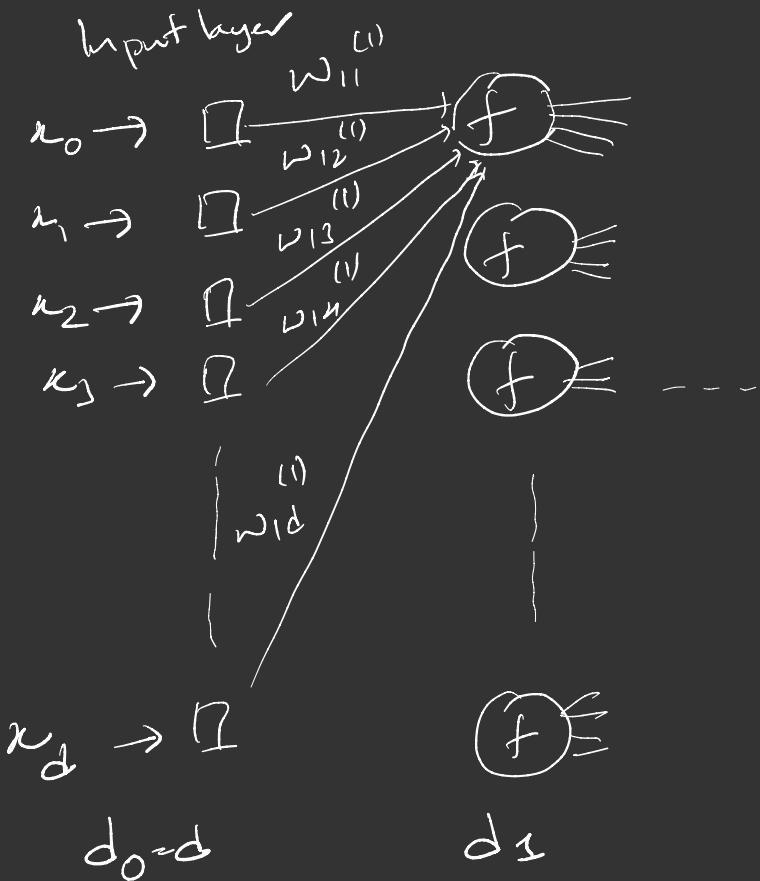
$$1 \leq l \leq N \quad 0 \leq i \leq d_l$$

for weights

$w_{ji}^{(l)} \leftarrow$ layer
 $i \rightarrow$ from i^{th} neuron of $(l-1)^{\text{th}}$ layer
 $j \rightarrow$ going to j^{th} neuron of l^{th} layer

(x_i, y)

Hidden layer



$d_l = \text{number of neurons in the layer } l'$

* More Notation

For targets : y_i \leftarrow itth value of vector y . y is the actual output

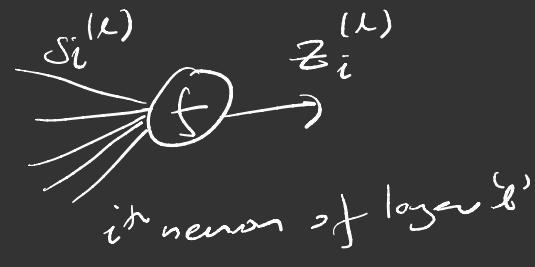
for output : \hat{y}_i \leftarrow itth value of vector \hat{y} . \hat{y} is the predicted

output of example 'n'

for hidden neurons

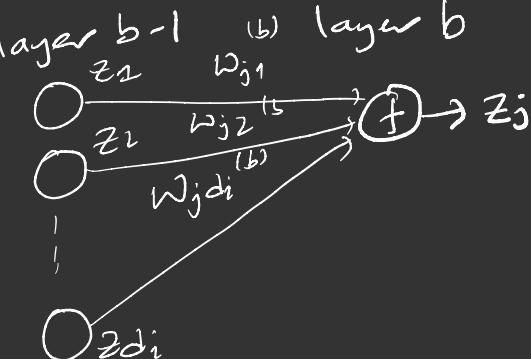
Weighted sum : $s_i^{(l)}$ \leftarrow layer
 $i \leftarrow$ itth neurons.

Output of a neuron : $z_i^{(l)}$ \leftarrow layer
 $i \leftarrow$ itth neurons.



Then, $z_j^{(b)} = f(s_j^{(b)})$ where f is an activation function.

And, $s_j^{(b)} = \sum_{i=0}^{d_{b-1}} w_{ji} z_i^{(b-1)}$



Then, in vector notation.

$$W_{i \in i^{\text{th}} \text{ layer}} = \left[\begin{array}{c} \boxed{\quad} \\ \vdots \\ \boxed{\quad} \end{array} \right] \left. \begin{array}{l} \text{di rows} \\ \text{di} \times \text{di-1} \end{array} \right\} \text{di columns} \left. \begin{array}{l} \text{or neurons} \\ \text{jth neuron} \end{array} \right\} \begin{array}{l} \text{each row has weights} \\ \text{of jth neuron in layer} \\ i \\ \text{layer-i} \end{array}$$

Similarly,

$$X_{i \in i^{\text{th}} \text{ layer}} = \left[\begin{array}{c} \boxed{\quad} \\ \vdots \\ \boxed{\quad} \end{array} \right] \left. \begin{array}{l} \text{layer i} \\ O \rightarrow z_1 \\ O \rightarrow z_L \\ \vdots \\ O \rightarrow z_i \end{array} \right\}$$

N number of examples in N rows.

Each row contains di values as the output at ith layer.

$$O \rightarrow z_{di}$$

$$\text{Then, } X_i = f(X_{i-1} W_i^T)$$

\downarrow \downarrow
 $N \times d_{i-1}$ $d_{i-1} \times d_i$
 $\underbrace{\hspace{10em}}$
 $N \times d_i$

* Functions

- i) Logistic: $f(s_i) = \frac{1}{1+e^{-s_i}}$ $f'(s_i) = f(s_i)[1-f(s_i)]$
- ii) Softmax: normalizing function: $R^N \rightarrow R^N$

$$g(s_i) = \frac{e^{s_i}}{\sum_{i=1}^m e^{s_i}} = \hat{y}_i$$

$$\begin{bmatrix} \frac{\partial \hat{y}_1}{\partial s_1} & \frac{\partial \hat{y}_1}{\partial s_2} & \dots & \frac{\partial \hat{y}_1}{\partial s_n} \\ \vdots & \searrow & & \vdots \\ \frac{\partial \hat{y}_n}{\partial s_1} & \dots & \dots & \frac{\partial \hat{y}_n}{\partial s_n} \end{bmatrix}_{ij}$$

Now, finding the derivative of softmax.

$$\begin{aligned}\frac{\partial g(s_i)}{\partial s_j} &= \frac{\sum_{i=1}^m e^{s_i} \times e^{s_i} \{i=j\} - e^{s_i} \times e^{s_j}}{\left(\sum_{i=1}^m e^{s_i}\right)^2} \\ &= \frac{\sum \times e^{s_i} \{i=j\} - e^{s_i} \times e^{s_j}}{\sum^2} \\ &= \frac{e^{s_i} (\sum \{i=j\} - e^{s_j})}{\sum \times \sum} \\ &= \frac{e^{s_i}}{\sum} \times \frac{(\sum \{i=j\} - e^{s_j})}{\sum} \\ &= g(s_i) \times [1 \{i=j\} - g(s_j)]\end{aligned}$$

So, when $i=j$, $g'(s_i) = g(s_i) [1 - g(s_i)]$

when $i \neq j$, $g'(s_j) = g(s_i) [-g(s_j)]$

$$= -g(s_i)g(s_j)$$

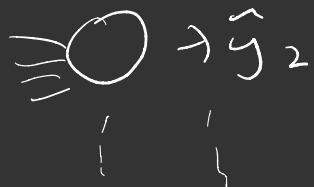
* Forward Propagation (Inference).

layer-l



$$\hat{y}_k = g(z_k^{(l)})$$

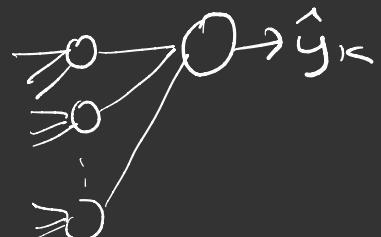
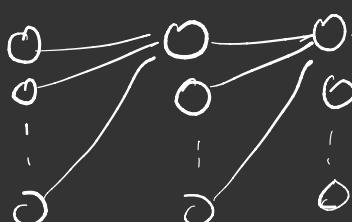
$$= g\left(\sum_{j=0}^{d_{l-1}} w_{kj}^{(l)} z_j^{l-1}\right)$$



$$= g\left[\sum_{j=0}^{d_{l-1}} w_{kj}^{(l)} f\left(\dots f\left(\sum_{p=0}^d w_{qp}^{(l)} x_p\right)\right)\right]$$



m' output neurons



* Cost function: cross entropy

$$C(\omega) = -\sum_{n=1}^N \sum_{k=1}^m y_k \log(\hat{y}_k) \Leftarrow \text{for all examples.}$$

+ for a single example. (x, y)

$$C_i(\omega) = -\sum_{k=1}^m y_k \log(\hat{y}_k) \quad 1 \leq k \leq m$$

* Now, computing derivatives using back propagation.

$$\frac{\partial C(\omega)}{\partial w_{cj}^{(l)}} = \frac{\partial C}{\partial s_c^{(l)}} \times \frac{\partial s_c^{(l)}}{\partial w_{cj}^{(l)}} \quad 1 \leq c \leq m$$

Here, $\frac{\partial s_c^{(l)}}{\partial w_{cj}^{(l)}} = z_j^{(l-1)}$ $\left[\dots, s_c^{(l)} = \sum_{j=0}^{d-1} w_{cj}^{(l)} z_j^{(l-1)} \right]$

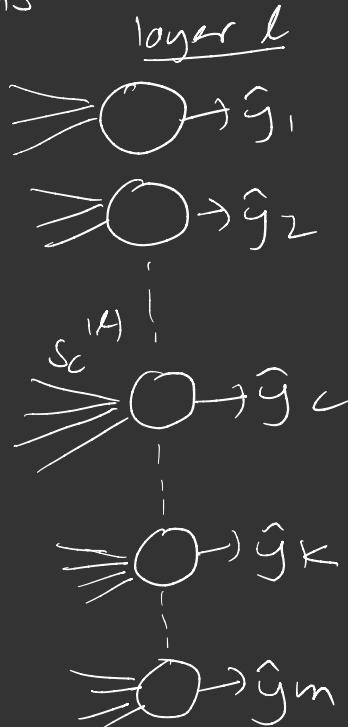
Then, computing $\frac{\partial \mathcal{L}}{\partial s_c^{(l)}}$ w.r.t weighted sum.

Fixing 'c' between $1 \leq c \leq m$ and 'k' changes between $1 \leq k \leq n$
 They both represent the output layer neurons

$$\frac{\partial C}{\partial s_c^{(l)}} = \frac{\partial}{\partial s_c^{(l)}} \left[- \sum_{k=1}^m y_k \log(\hat{y}^k) \right]$$

$$= - \sum_{k=1}^m y_k \frac{\partial}{\partial s_c^{(l)}} \log(\hat{y}^k)$$

$$= - \sum_{k=1}^m y_k \frac{\partial}{\partial \hat{y}^k} (\log \hat{y}^k) \times \frac{\partial \hat{y}^k}{\partial s_c^{(l)}}$$



$$= - \sum_{K=1}^m y_K \times \frac{1}{\hat{y}_K} \times \hat{y}_L (1 \{ K=c \} - \hat{y}_c)$$

$\therefore \hat{y}_K = g(s_c^{(w)})$ and

$$\frac{\partial g(s_i)}{\partial s_j} = g'(s_i) (1 \{ i=j \} - g(s_j))$$

$$= - \sum_{K=1}^m y_K (1 \{ K=c \} - \hat{y}_c)$$

$$= \left[- \sum_{K \neq c}^m y_K (-\hat{y}_c) \right] - y_c (1 - \hat{y}_c)$$

$$= \sum_{K \neq c}^m y_K \hat{y}_c - y_c + y_c \hat{y}_c$$

$$= \sum_{K=1}^m y_K \hat{y}_c - y_c$$

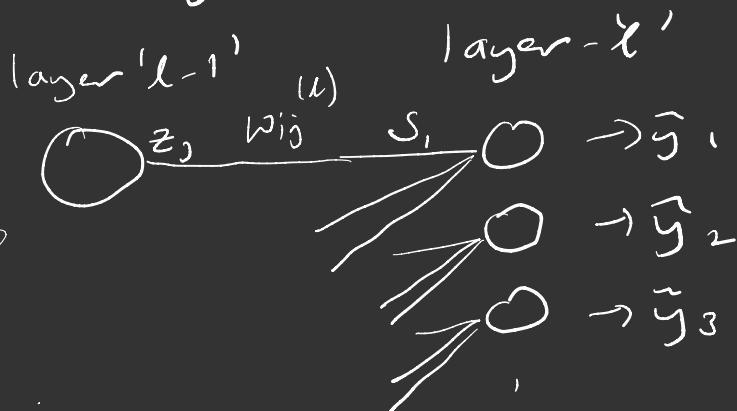
$$= \hat{y}_c \sum_{k=1}^m y_k - y_c$$

$$= \hat{y}_c - y_c \quad [\sum_{k=1}^m y_k = 1; \text{sum of probabilities is 1}]$$

$$\therefore \frac{\partial L}{\partial w_{ij}^{(l)}} = (\hat{y}_c - y_c) \times z_j^{(l-1)}$$

\Downarrow
weights of last layer 'l'

between layers 'l' and 'l-1'



$$\text{So, } w_l = w_l - \frac{\alpha}{|B|} (\hat{y} - y)^T X_{l-1}$$

$$\downarrow \qquad \downarrow \qquad \downarrow$$

$$m \times d_{l-1} \qquad m \times B \qquad B \times d_{l-1}$$

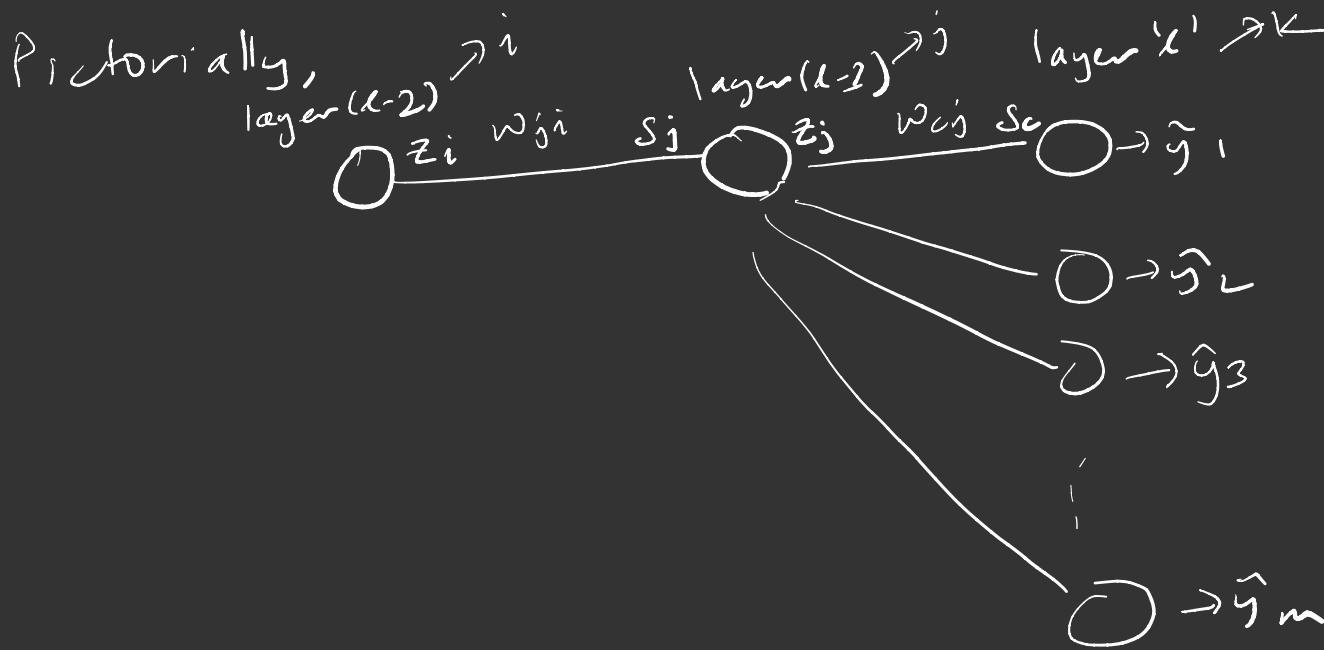
$$\Rightarrow \vec{D} \rightarrow \hat{y}_m$$

* Now, computing gradients $\frac{\partial C}{\partial w_{ji}^{(l-1)}}$ for the second last layer weights between layer ' $l-2$ ' and ' $l-1$ '.

$$\begin{aligned}\frac{\partial C}{\partial w_{ji}^{(l-1)}} &= \sum_{c=1}^m \frac{\partial C}{\partial s_c^{(l)}} \times \frac{\partial s_c^{(l)}}{\partial z_j^{(l-1)}} \times \frac{\partial z_j^{(l-1)}}{\partial s_j^{(l-1)}} \times \frac{\partial s_j^{(l-1)}}{\partial w_{ji}^{(l-2)}} \\ &= \sum_{c=1}^m (\hat{y}_c - y_c) w_{cj}^{(l)} \times [z_j^{(l-1)} \times (1 - z_j^{(l-1)})] \times z_i^{(l-2)}\end{aligned}$$

We used $\frac{\partial C}{\partial s_c^{(l)}}$ from the previous computation and then

$$\begin{aligned}\frac{\partial z_j^{(l-1)}}{\partial s_j^{(l-1)}} &= z_j^{(l-1)} (1 - z_j^{(l-1)}) \text{ is the derivation of sigmoid activation function. } z_j^{(l-1)} = f(s_j^{(l-1)}) \text{ and } s_j^{(l-1)} = \sum_{i=0}^n w_{ji} z_i^{(l-2)}.\end{aligned}$$



So, in general, if we define the error in neuron ' i '

in the layer ' l ' as $\delta_i^{(l)} = \frac{\partial C}{\partial s_i^{(l)}}$, then

$$\delta_i^{(l)} = \frac{\partial C}{\partial z_i^{(l)}} \times \frac{\partial z_i^{(l)}}{\partial s_i^{(l)}} = \frac{\partial C}{\partial z_i^{(l)}} \times f'(s_i^{(l)})$$

which can be written as.

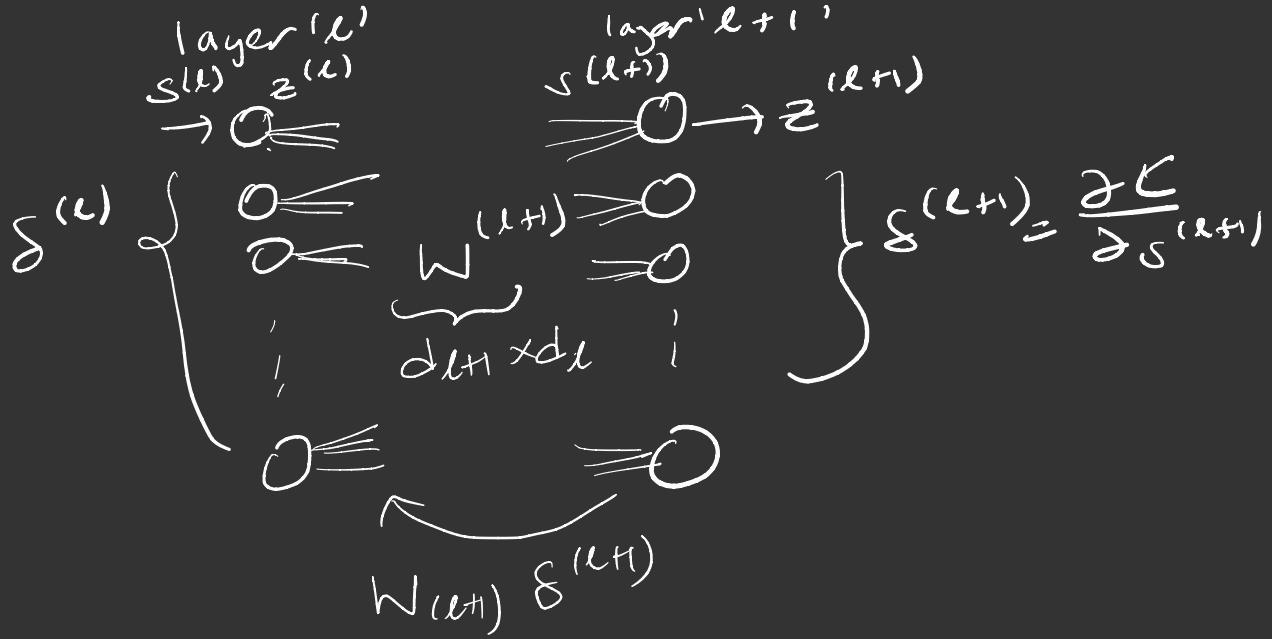
$$\delta^l = \nabla_z C \odot f'(s^{(l)}) \quad \text{where } \odot \text{ is the element wise product.}$$

Now, the error δ^l in terms of error in the next layer

δ^{l+1} is given by

$$\delta^l = ((w^{l+1})^T \delta^{l+1}) \odot f'(s^{(l)})$$

The way to understand this is suppose we know the error in the $(l+1)$ th layer i.e δ^{l+1} . Then applying $(w^{l+1})^T$ is intuitively moving back the error through the network, giving some measure of error δ^l at the output of layer l .



Then, applying $\odot f'(s^{(l)})$ moves the error backward through the activation function in layer l , giving us the error δ^l in the weighted input ($s^{(l)}$) to layer l

$$\delta^{(l)} = (W^{(l+1)} \delta^{(l+1)}) \oplus f'(s^{(l)})$$

$$\text{Now, } \frac{\partial C}{\partial w_{ji}^{(l)}} = \frac{\partial C}{\partial s_j^{(l)}} \times \frac{\partial s_j^{(l)}}{\partial w_{ji}^{(l)}}$$

$$= \delta_j^l \times z_i^{(l-1)}$$

So, generally,

i) Input : (x_i, y_i) for $1 \leq i \leq N$ N examples.

ii) Feed forward: for each example, each layer $l = 1, 2, 3, \dots$

compute $s^l = w^l z^{l-1}$ and $z^{l-1} = f(s^{l-1})$

iii) Output error: $\delta^l = \nabla_z C \odot f'(s^l)$

iv) Back propagate error: For each layer $l = 1, 2, 3, \dots$

$$\delta^{(l)} = ((w^{l+1}) \quad \delta^{l+1}) \odot f'(s^l)$$

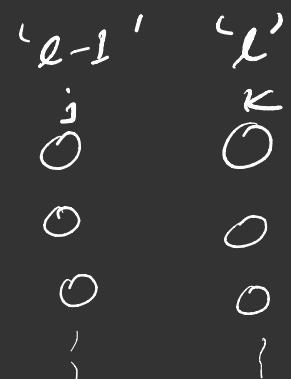
v) Gradient of loss function :

$$\frac{\partial C}{\partial w_{ji}^{(l)}} = z_i^{(l-1)} \times \delta_j^{(l)}$$

For our case, logistic function in the hidden layer and softmax on the output layer.

i) For last layer ' l '

$$\frac{\partial C}{\partial w_{kj}^{(l)}} = \delta_k^{(l)} \times z_j^{(l-1)}$$



$$= \frac{\partial C}{\partial s_k^{(l)}} \times z_j^{(l-1)}$$

$$= (\hat{y}_k - y_k) \times z_j^{(l-1)}$$



i) for second last layer ($\ell-1$)

$$\frac{\partial C}{\partial w_{ji}^{(\ell-1)}} = \delta_j^{(\ell-1)} \times z_i^{(\ell-2)}$$

$$= [w_{kj}^{(u)} \delta_k^{(\ell)} \odot f'(s_j^{(\ell-1)})] z_i^{(\ell-2)}$$

$$\begin{array}{ccccccccc} i & j & k & & & & & & \\ \ell-2 & \ell-1 & \ell & & & & & & \\ 0 & 0 & 0 & & & & & & \\ 0 & 0 & 0 & & & & & & \\ 0 & 0 & 0 & & & & & & \\ 0 & 0 & 0 & & & & & & \\ | & | & | & & & & & & \\ 0 & 0 & 0 & & & & & & \end{array} = \sum_{k=1}^m (\hat{y}_k - y_k) \underbrace{w_{kj}^{(u)} \times [z_j^{(\ell-1)} (1 - z_j^{(\ell-1)})]}_{\delta_k^{(u)}} \times z_i^{(\ell-2)}$$

confusion - for the third last layer and onward.

$$\frac{\partial L}{\partial w_{ih}^{(l-2)}} = \delta_i^{(l-2)} \times z_n^{(l-3)}$$

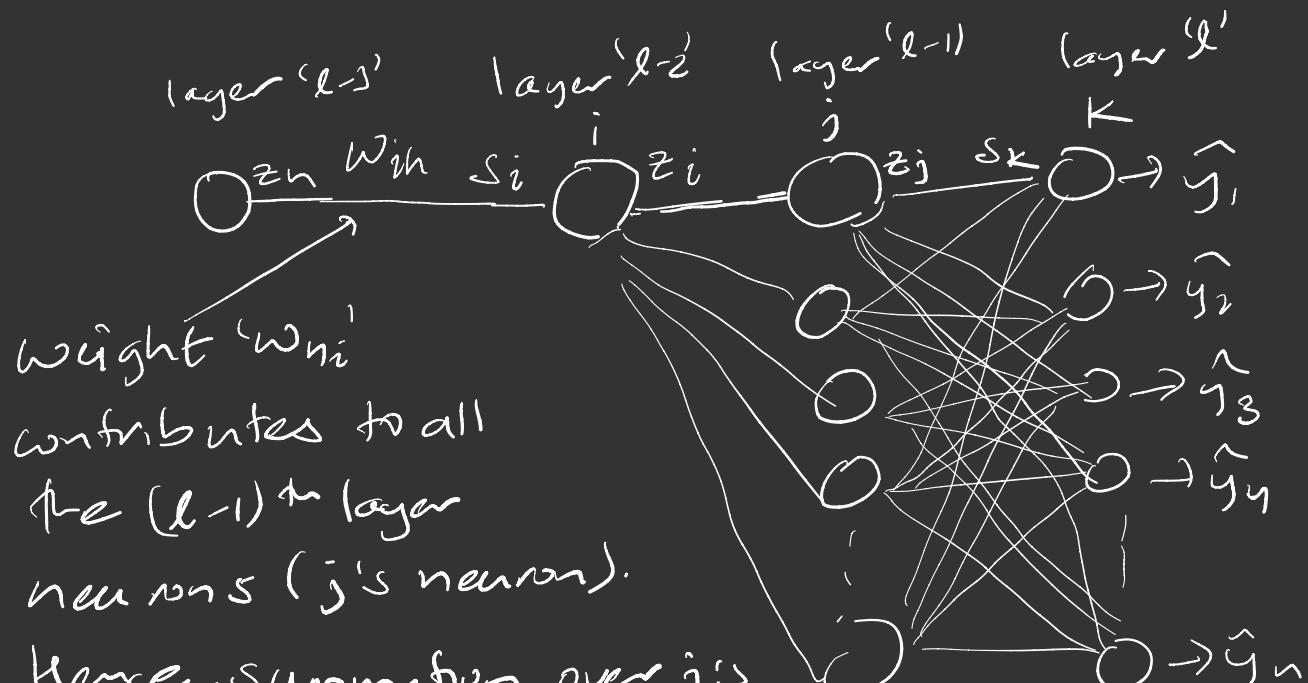
$$= [w_{ji}^{(l-1)} \delta_j^{(l-1)} \odot f'(s_i^{(l-2)})] \times z_n^{(l-3)}$$

$$\begin{matrix} h & i & j & k \\ l-3 & l-2 & l-1 & l \\ 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 \\ \vdots & \vdots & \vdots & \vdots \\ 0 & 0 & 0 & 0 \end{matrix} = \delta_j^{(l-1)} \times w_{ji}^{(l-1)} \times f'(s_i^{(l-2)}) \times z_n^{(l-3)}$$

$$\begin{aligned} &= \delta_j^{(l-1)} \times w_{ji}^{(l-1)} \times [z_i^{(l-2)} (1 - z_i^{(l-2)})] \times z_n^{(l-3)} \\ &\quad \xrightarrow{\sum} \sum_{k=1}^m (\hat{y}_k - y_k) w_{kj}^{(l)} \times z_j^{(l-1)} [1 - z_j^{(l-1)}] \times w_{ji}^{(l-1)} \\ &\quad \times [z_i^{(l-2)} (1 - z_i^{(l-2)})] \times z_n^{(l-3)} \end{aligned}$$

Should have summation over j or not? $\sum_{j=0}^{d_j}$

We should be having because without it what would
the 'j' terms in the derivative equation mean.



Hence, summation over j 's makes sense. But need to confirm this.