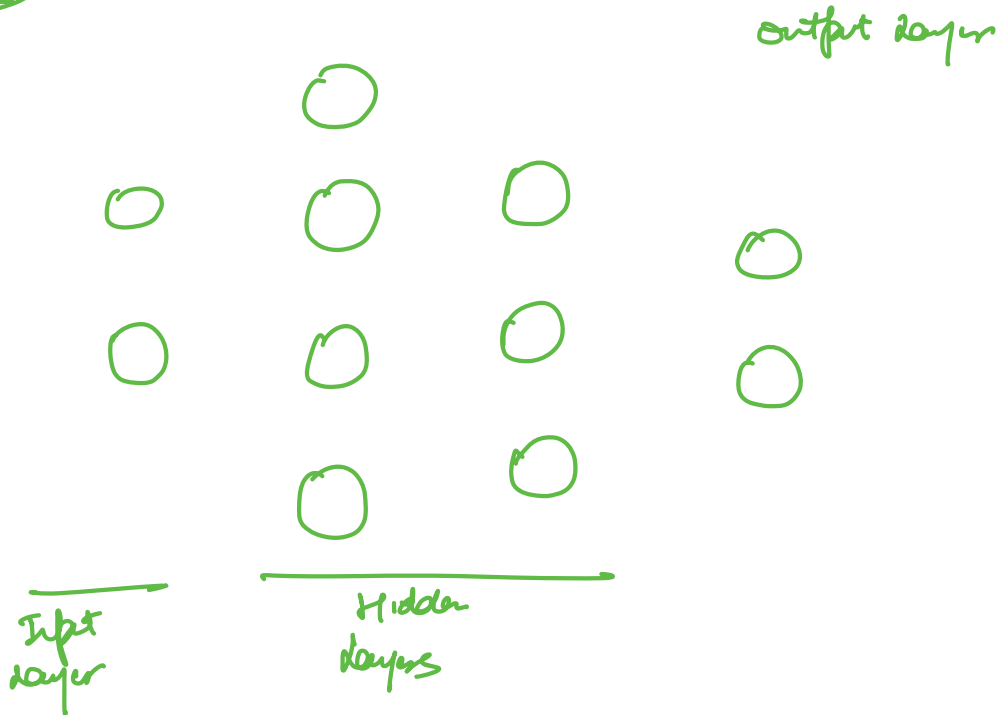
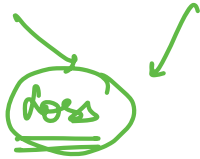


weights bias



2

a) Output (L)

$$\frac{\partial L}{\partial w_{i,j}^L}$$

↓

∂L wrt w in last layer

$$\frac{\partial L}{\partial b^L}$$

↓

∂L wrt bias in last layer

b) Hidden (L)

$$\frac{\partial L}{\partial w_{i,j}^L}$$

↓

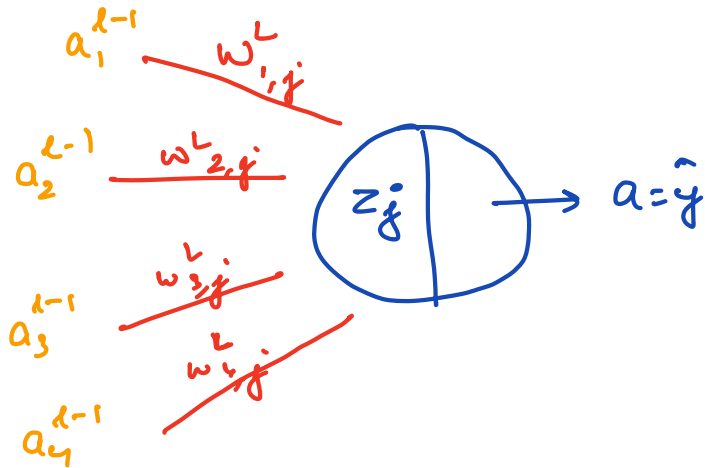
∂L wrt w in hidden layers

$$\frac{\partial L}{\partial b^L}$$

↓

∂L wrt bias in hidden layers

Output layer



$$z_j^L = \sum_i w_{i,j}^L a_i^{L-1} + b_j^L$$

i iterates
over all neurons
in l-1 layer

$$\frac{\partial L}{\partial \omega}$$

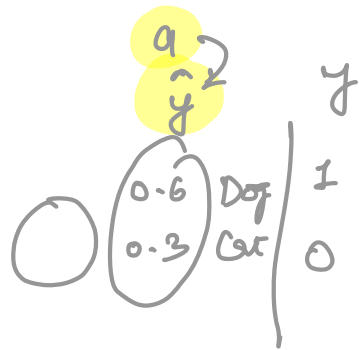
$$w_{i,j} \rightarrow z_j \rightarrow a_j \rightarrow L(a, y)$$

3

$$\frac{\partial L}{\partial w_{i,j'}} = \underbrace{\frac{\partial L}{\partial a_j}}_{(1)} \cdot \underbrace{\frac{\partial a_j}{\partial z_j}}_{(2)} \underbrace{\frac{\partial z_j}{\partial w_{i,j'}}}_{(3)}$$

$$\downarrow$$

$$\delta_j^L = \frac{\partial L}{\partial z_j}$$



$$\textcircled{1} \quad \frac{\partial L}{\partial a_j}$$

$$L = \frac{1}{2} \sum_i (y_i^L - a_i^L)^2$$

<u>L</u>	a^L	y^L
○	0.3	0
○	0.5	1

$$\frac{\partial L}{\partial a_j} = \frac{1}{2} \cdot 2 \cdot (y_j - a_j) (-1)$$

$$0 \rightarrow 0.2 \mid 0$$

$$\frac{\partial L}{\partial a_j} = -(y_j - a_j)$$

2 $\frac{\partial a_j}{\partial z_j}$

$$a_j = \sigma(z_j)$$

$$\frac{\partial a_j}{\partial z_j} = \sigma'(z_j)$$

$$\frac{\partial a_j}{\partial z_j} = (1 - \sigma(z_j)) \sigma(z_j)$$

(Sigmoid function)

$$\sigma(z) = \frac{1}{1 + e^{-z}}$$

$$\frac{\partial \sigma(z)}{\partial z} = \frac{\partial (1 + e^{-z})^{-1}}{\partial z}$$

$$= \frac{-1 (e^{-z}) (-1)}{(1 + e^{-z})^2}$$

$$= \frac{e^{-z}}{(1 + e^{-z})^2}$$

$$= \left(\frac{1}{1 + e^{-z}} \right) \left(1 - \frac{1}{1 + e^{-z}} \right)$$

$$= \sigma(z) (1 - \sigma(z))$$

3 $\frac{\partial z_j^l}{\partial w_{i,j}}$

$$z_j^l = \sum_i w_{i,j} a_i^{l-1} + b_j^l$$

$$\frac{\partial z_j^l}{\partial w_{i,j}} = a_i^{l-1}$$

$$\frac{\partial (w_{1,j} a_1^{l-1} + w_{2,j} (a_2^{l-1}) + \dots)}{\partial w_{2,j}}$$

Combine

$$\frac{\partial L}{\partial \omega_{i,j}} = \delta_j^L \cdot \frac{\partial z_j}{\partial \omega_{i,j}}$$

$$\frac{\partial L}{\partial \omega_{i,j}} = \delta_j^L \cdot a_i^{L-1}$$

$$\delta_j^L = -(y_j - a_j) \sigma'(z_j)$$

Bias

$$\frac{\partial L}{\partial b_j} = \underbrace{\frac{\partial L}{\partial a_j} \cdot \frac{\partial a_j}{\partial z_j}}_{\delta_j^L} \cdot \frac{\partial z_j}{\partial b_j}$$

$$z_j^L = \sum_i \omega_{i,j} a_i^{L-1} + b_j^L$$

$$\frac{\partial z_j}{\partial b_j} = 1$$

$$\frac{\partial L}{\partial b_j} = \delta_j^L \cdot 1$$

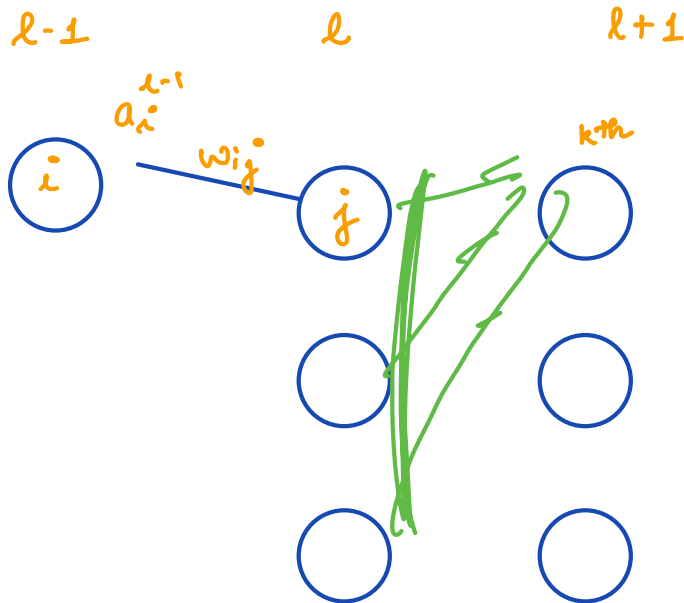
for output layer

$$\frac{\partial L}{\partial \omega_{i,j}} = \delta_j^L \cdot a_i^{L-1}$$

$$\frac{\partial L}{\partial b_j} = \delta_j^L$$

$$\delta_j^L = (a_j - y_j) \sigma'(z_j)$$

Hidden Layer



$$\frac{\partial L}{\partial w_{ij}^l} = \frac{\partial L}{\partial z_k^{l+1}} \cdot \underbrace{\frac{\partial z_k^{l+1}}{\partial a_j^l}}_{(1)} \cdot \underbrace{\frac{\partial a_j^l}{\partial z_j^l}}_{(2)} \cdot \underbrace{\frac{\partial z_j^l}{\partial w_{ij}^l}}_{(3)}$$

$$(1) \quad \frac{\partial z_k^{l+1}}{\partial a_j^l}$$

$$z_k^{l+1} = \sum_j w_{jk} a_j^l + b^{l+1}$$

$$\frac{\partial z_k^{l+1}}{\partial a_j^l} = w_{jk}$$

$$(2) \quad \frac{\partial a_j^l}{\partial z_j^l}$$

$$a_j^l = \sigma(z_j^l)$$

$$\frac{\partial a_j}{\partial z_j} = \sigma'(z_j)$$

$$(3) \quad \frac{\partial z_j}{\partial w_{ij}}$$

$$z_j = w_{ij} a_i^{l-1}$$

$$\frac{\partial z_j}{\partial w_{ij}} = a_i^{l-1}$$

$$\frac{\partial L}{\partial w_{ij}^l} = \sum_k \frac{\partial L}{\partial z_k^{l+1}} \cdot \frac{\partial z_k^{l+1}}{\partial a_j} \cdot \frac{\partial a_j}{\partial z_j} \cdot \frac{\partial z_j}{\partial w_{ij}}$$

$$= \boxed{\sum_k \delta_k^{l+1} \cdot w_{jk} \cdot \sigma'(z_j^l)} \cdot a_i^{l-1}$$

$$\downarrow$$

$$\delta_j^l$$

$$\frac{\partial L}{\partial w_{ij}^l} = \delta_j^l \cdot a_i^{l-1}$$

Bias Update Rule

$$\frac{\partial L}{\partial b} = \boxed{\sum_k \frac{\partial L}{\partial z_k^{l+1}} \cdot \frac{\partial z_k^{l+1}}{\partial a_j} \cdot \frac{\partial a_j}{\partial z_j}} \cdot \frac{\partial z_j}{\partial b}$$

$$\downarrow$$

$$\delta_j^l$$

$$\downarrow$$

$$1$$

$$\frac{\partial L}{\partial b} = \delta_j^l \cdot 1$$

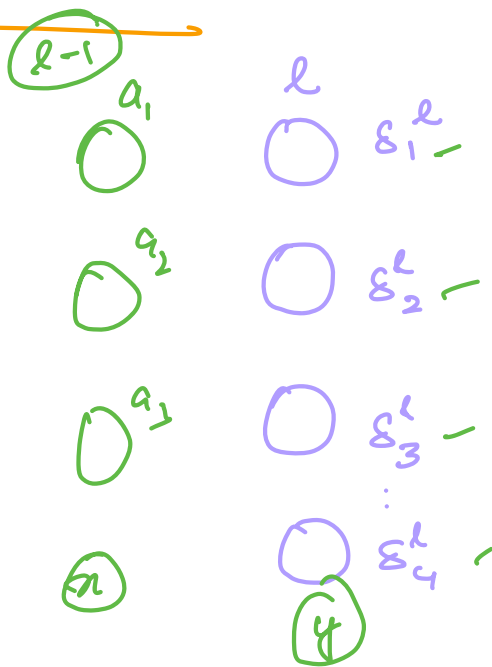
Final Result:

$$\frac{\partial L}{\partial \omega_{ij}^l} = \delta_j^l \cdot a_i^{l-1}$$

$$\frac{\partial L}{\partial b} = \delta_j^l \cdot 1$$

$$\delta_j^l = \sum_k \frac{\partial L}{\partial z_k^{l+1}} \cdot \frac{\partial z_k^{l+1}}{\partial a_j} \cdot \frac{\partial a_j}{\partial z_j}$$

Matrix Representation:



$$\frac{\partial L}{\partial \omega_{ij}^l} \rightarrow \frac{\partial L}{\partial \omega^l}$$

$$\frac{\partial L}{\partial \omega^l} = \underbrace{a_i^{l-1}}_{(n \times 1)} \cdot \underbrace{(\delta^l)^T}_{(1 \times n)}$$

$n \times n$

$$\begin{bmatrix} a_1^{L-1} \\ a_2^{L-1} \\ \vdots \\ a_n^{L-1} \end{bmatrix} \cdot [\delta_1^L \quad \delta_2^L \quad \delta_3^L \quad \delta_4^L \quad \dots \quad \delta_y^L]$$

$$\frac{\partial L}{\partial w_{ij}} \leftarrow \begin{bmatrix} a_1^{L-1} \delta_1^L & a_1^{L-1} \delta_2^L & a_1^{L-1} \delta_3^L & \dots \\ a_2^{L-1} \delta_1^L & a_2^{L-1} \delta_2^L & \dots & \dots \end{bmatrix}$$

$$\frac{\partial L}{\partial w^L} = a^{L-1} \cdot (\delta^L)^T$$

For bias: $\frac{\partial L}{\partial b_j^L} = \delta_j^L$

$$\frac{\partial L}{\partial b^L} = \delta^L$$

→ ○

→ ○

→ ○

$$b = b - \eta \frac{\partial L}{\partial b}$$

$$\begin{bmatrix} b_1^L \\ b_2^L \\ b_3^L \end{bmatrix}$$

$$\frac{\partial L}{\partial b} = \begin{bmatrix} \delta_1^L \\ \delta_2^L \\ \delta_3^L \end{bmatrix}$$

Loss fnⁿ: Square Error

Binary Classification:

output layer:

Dog/Cat

(z)

$$z \rightarrow a = \sigma(z) = \hat{y}$$

$< 0.5 \rightarrow \text{Dog}$

$> 0.5 \rightarrow \text{Cat}$

k classes $k > 2$

x neurms in output layer.

Binary Classifcats:

↳ Binary Cross Entropy:

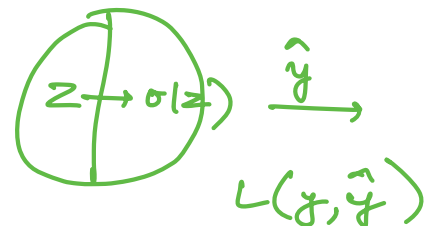
$$= - \sum_{i=1}^m (y_i \log \hat{y}_i + (1-y_i) \log (1-\hat{y}_i))$$

y_i = True Label

\hat{y}_i = Predicted
Label

$$\delta_L = \frac{\partial L}{\partial z^2} = \frac{\partial L}{\partial a} \cdot \frac{\partial a}{\partial z}$$

\downarrow
 \hat{y}



$$\delta_L = \left(\frac{-y_i}{\hat{y}_i} + \frac{1-y_i}{1-\hat{y}_i} \right)$$

?

$$a = \sigma(z)$$

$$\begin{aligned}\frac{\partial a}{\partial z} &= \sigma(z) (1 - \sigma(z)) \\ &= a (1 - a) \\ &= \hat{y} (1 - \hat{y})\end{aligned}$$

$$\delta_L = \left(\frac{-y_i}{\hat{y}_i} + \frac{1-y_i}{1-\hat{y}_i} \right) (\hat{y}_i (1-\hat{y}_i))$$

$$= \left(\frac{-\cancel{y_i} + \cancel{y_i} \hat{y}_i + \hat{y}_i - y_i \hat{y}_i}{\cancel{\hat{y}_i} (1-\cancel{\hat{y}_i})} \right) (\cancel{\hat{y}_i} (1-\cancel{\hat{y}_i}))$$

$$\delta_L = (\hat{y}_i - y_i)$$

↓
 δ_L in cross
 Entropy loss

Multiple Examples

C : no. of classes in o/p layer

$$a_c = \begin{bmatrix} a_1 \\ a_2 \\ \vdots \\ a_c \end{bmatrix}$$

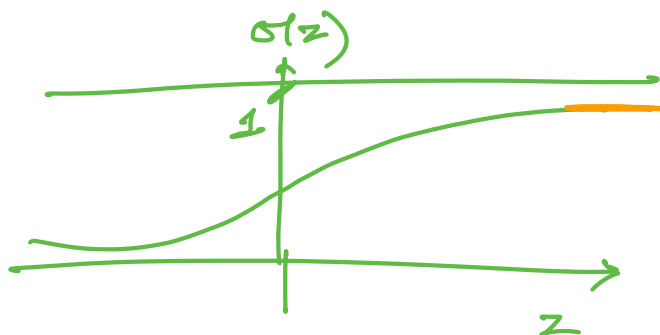


$$\begin{bmatrix} \text{---} a^{(1)} \text{---} \\ \text{---} a^{(2)} \text{---} \\ \text{---} a^{(3)} \text{---} \\ \vdots \\ \text{---} a^{(m)} \text{---} \end{bmatrix}$$

$$\begin{bmatrix} 1^{\text{st}} \\ 2^{\text{nd}} \\ \vdots \\ m^{\text{th}} \end{bmatrix} \begin{matrix} 0.1 & 0.8 & 0.2 \end{matrix}$$

$$m \times c$$

Sigmoid



$$\sigma'(z) = \sigma(z)(1 - \sigma(z))$$

Z is large

$$\sigma'(z) = \sigma(z) \underbrace{(1 - \sigma(z))}_{\substack{1 \\ 0}}$$

2 is small

$$\sigma'(z) = \underbrace{\sigma(z)}_p (1 - \sigma(z))$$

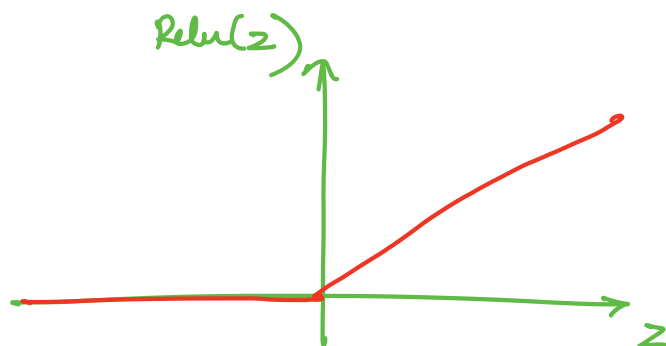
20

$w = w - \eta \left(\frac{\partial L}{\partial w} \right) \rightarrow 0$

gradient

no change in costs.

Relu



$$\text{Relu}(z) = \begin{cases} z & z \geq 0 \\ \underline{\underline{0}} & \underline{\underline{z < 0}} \end{cases}$$

① Overfit
Dropout }