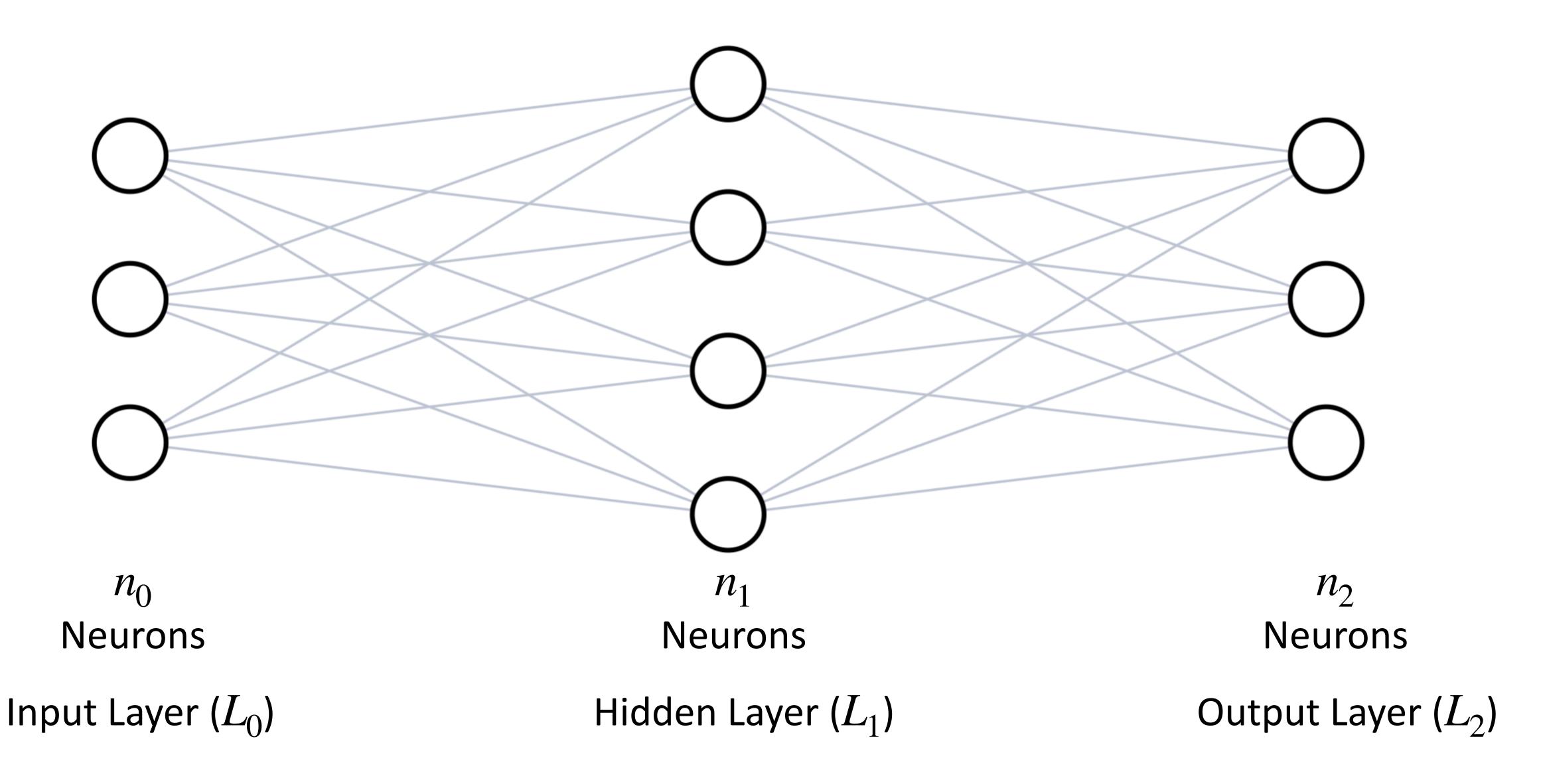
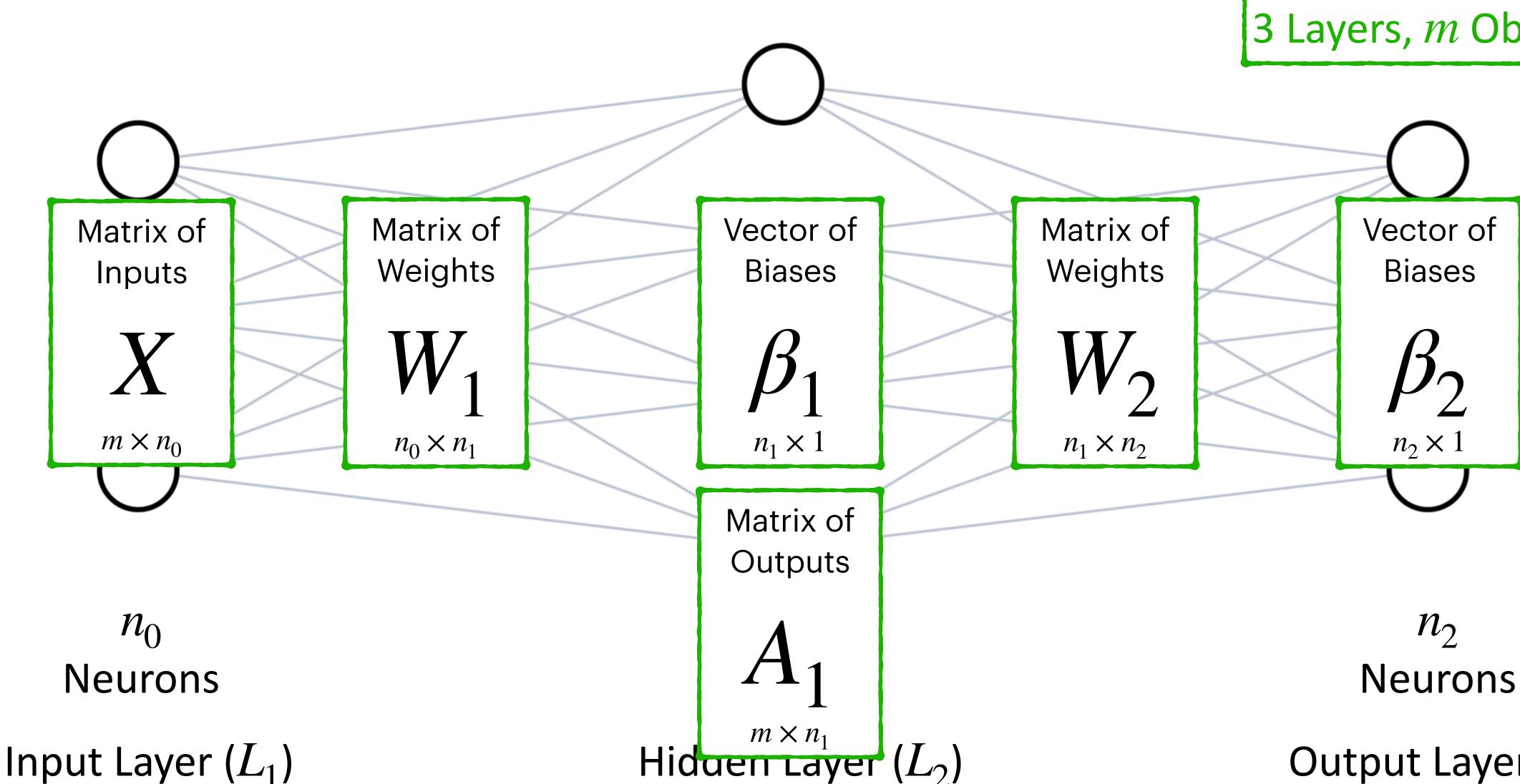
## **Back Propagation Equations for Binary Classification**

Rahul Singh rsingh@arrsingh.com



3 Layers, m Observations

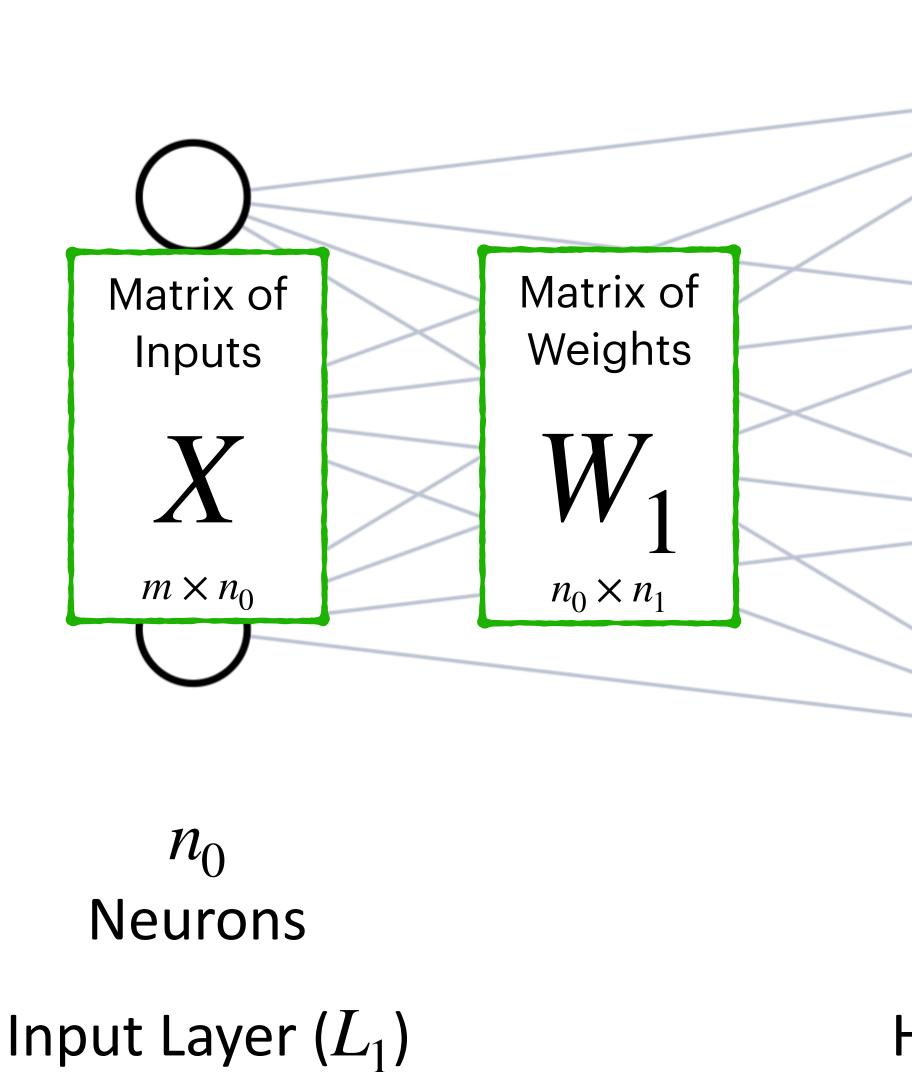


Matrix of Outputs

 $m \times n_2$ 

Output Layer  $(L_3)$ 

3 Layers, m Observations



Vector of Biases

 $\beta_{1}$   $n_1 \times 1$ 

Matrix of Weights

 $\frac{W_2}{n_1 \times n_2}$ 

Vector of Biases

 $\beta_2$   $n_2 \times 1$ 

Matrix of Outputs



 $m \times n_2$ 

Matrix of Outputs

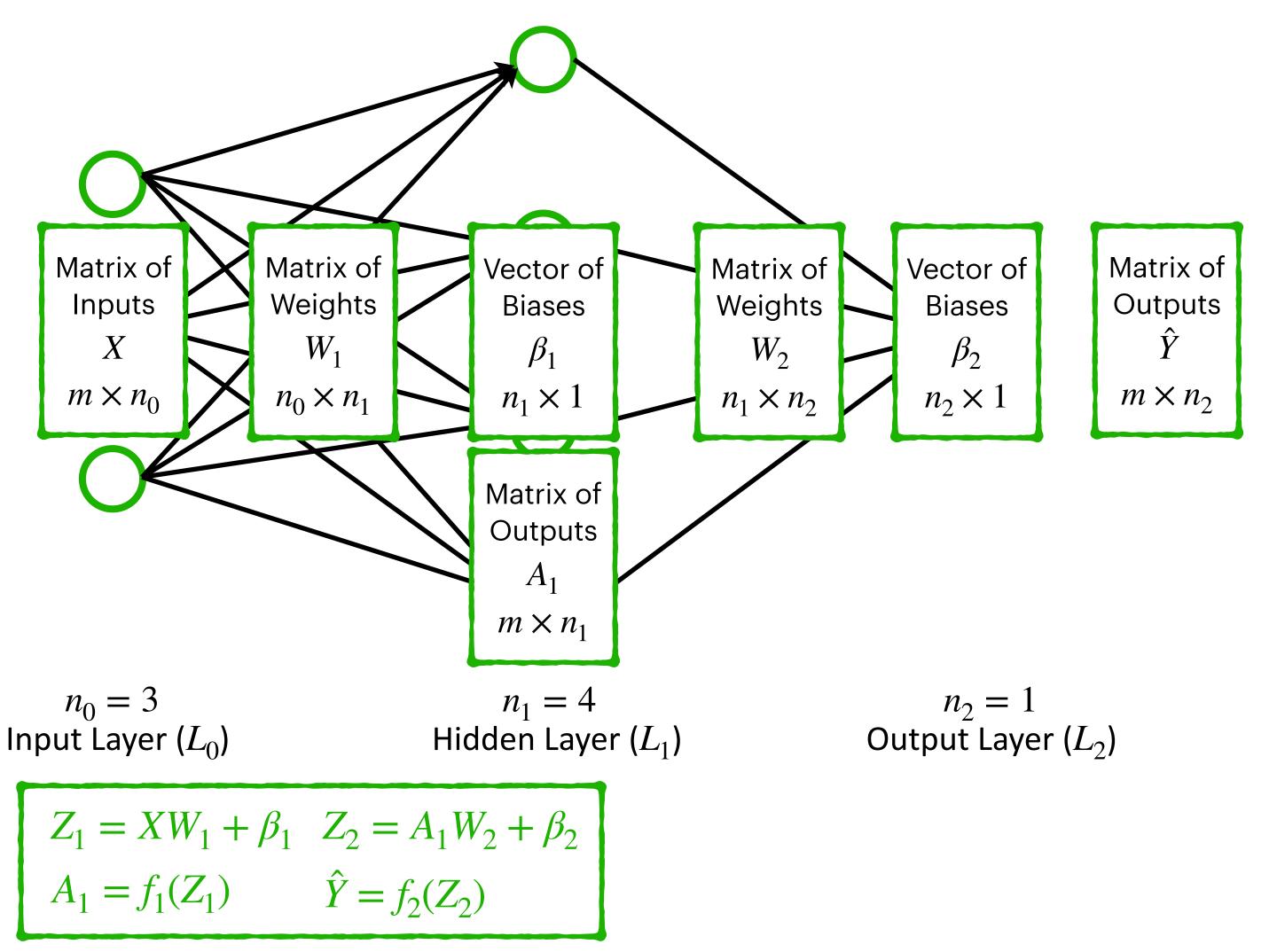
 $A_1$ 

liduen Laye

Matrix Equations to compute the Predicted values  $\hat{Y}$  given inputs X, Weights  $W_1$  and  $W_2$  and Biases  $\beta_1$  and  $\beta_2$ 

$$Z_1 = XW_1 + \beta_1$$
  $Z_2 = A_1W_2 + \beta_2$   
 $A_1 = f(Z_1)$   $\hat{Y} = f_2(Z_2)$ 

 $A_1=f(Z_1) \qquad \hat{Y}=f_2(Z_2)$   $f_1(g) \ {\rm and} \ f_2(g) \ {\rm are \ the \ activation \ functions \ in \ Layers} \ L_1 \ {\rm and} \ L_2$ 



### **Neural Networks**

Cost function is the **Binary Cross Entropy**:

$$-\frac{1}{n}\sum_{i=1}^{n} [y \log_{e} \hat{y} + (1 - y) \log_{e} (1 - \hat{y})]$$

Loss function:

$$L = -[y \log_e \hat{y} + (1 - y) \log_e (1 - \hat{y})]$$

$$\hat{y} = f_2(x) = \frac{1}{1 + e^{-x}}$$

## **Neural Networks**

Cost function is the **Binary Cross Entropy**:

$$-\frac{1}{n} \sum_{i=1}^{n} \left[ y \log_e \hat{y} + (1 - y) \log_e (1 - \hat{y}) \right]$$

Loss function:

## Problem Statement: Calculate the four partial derivatives:

$$\frac{\partial}{\partial \beta_2} cost$$
,  $\frac{\partial}{\partial W_2} cost$ ,  $\frac{\partial}{\partial \beta_1} cost$ ,  $\frac{\partial}{\partial W_1} cost$ 

 $n_0 = 3$  Input Layer ( $L_0$ )

$$n_1 = 4$$
 Hidden Layer ( $L_1$ )

Outputs

 $m \times n_1$ 

$$n_2 = 1$$
  
Output Layer ( $L_2$ )

$$Z_1 = XW_1 + \beta_1$$
  $Z_2 = A_1W_2 + \beta_2$   
 $A_1 = f_1(Z_1)$   $\hat{Y} = f_2(Z_2)$ 

$$\frac{\partial}{\partial W_2} L = \frac{\partial}{\partial \hat{y}} L \frac{\partial}{\partial z_2} \hat{y} \frac{\partial}{\partial W_2} z_2$$
 Chain Rule:  
  $L$  depends on  $\hat{y}, z_2$  in that order

 $\frac{d}{dx}log_e x = \frac{1}{x}$ 

 $+y\hat{y}-y\hat{y}$  cancels out

First let's calculate 
$$\frac{\partial}{\partial \hat{y}} L$$

$$\Rightarrow \frac{\partial}{\partial \hat{\mathbf{y}}} L = \frac{\partial}{\partial \hat{\mathbf{y}}} [-y \log_e \hat{\mathbf{y}} - (1 - y) \log_e (1 - \hat{\mathbf{y}})]$$

$$\Rightarrow \frac{\partial}{\partial \hat{\mathbf{y}}} L = -\frac{\mathbf{y}}{\hat{\mathbf{y}}} - \frac{(1-\mathbf{y})}{(1-\hat{\mathbf{y}})} (-1)$$

$$\Rightarrow \frac{\partial}{\partial \hat{y}} L = \frac{-y(1-\hat{y}) + \hat{y}(1-y)}{\hat{y}(1-\hat{y})}$$

$$\Rightarrow \frac{\partial}{\partial \hat{\mathbf{y}}} L = \frac{-\mathbf{y} + \mathbf{y}\hat{\mathbf{y}} + \hat{\mathbf{y}} - \mathbf{y}\hat{\mathbf{y}}}{\hat{\mathbf{y}}(1 - \hat{\mathbf{y}})}$$

$$\Rightarrow \frac{\partial}{\partial \hat{\mathbf{y}}} L = \frac{-\mathbf{y} + \hat{\mathbf{y}}}{\hat{\mathbf{y}}(1 - \hat{\mathbf{y}})}$$

$$\Rightarrow \frac{\partial}{\partial \hat{y}} L = \frac{\hat{y} - y}{\hat{y}(1 - \hat{y})}$$

## **Neural Networks**

Cost function is the **Binary Cross Entropy**:

$$-\frac{1}{n} \sum_{i=1}^{n} \left[ y_i \log_e \hat{y}_i + (1 - y_i) \log_e (1 - \hat{y}_i) \right]$$

Loss function:

$$L = -[y \log_e \hat{y} + (1 - y) \log_e (1 - \hat{y})]$$

$$\hat{y} = f_2(x) = \frac{1}{1 + e^{-x}}$$

Next lets calculate 
$$\frac{\partial}{\partial z_2} \hat{y}$$

$$\Rightarrow \frac{\partial}{\partial z_2} \hat{y} = \frac{\partial}{\partial z_2} \left( \frac{1}{1 + e^{-z_2}} \right)$$

$$\Rightarrow \frac{\partial}{\partial z_2} \hat{y} = \frac{\partial}{\partial z_2} (1 + e^{-z_2})^{-1}$$

$$\Rightarrow \frac{\partial}{\partial z_2} \hat{y} = (-1)(1 + e^{-z_2})^{-2} \frac{\partial}{\partial z_2} (1 + e^{-z_2})$$

$$\Rightarrow \frac{\partial}{\partial z_2} \hat{y} = (-1)(1 + e^{-z_2})^{-2}(-1)e^{-z_2}$$

$$\Rightarrow \frac{\partial}{\partial z_2} \hat{y} = \frac{e^{-z_2}}{(1 + e^{-z_2})^2}$$

$$\Rightarrow \frac{\partial}{\partial z_2} \hat{y} = \frac{\hat{y}^2 (1 - \hat{y})}{\hat{y}}$$

$$\Rightarrow \frac{\partial}{\partial z_2} \hat{y} = \hat{y}(1 - \hat{y})$$

#### Power Rule:

$$\frac{d}{dx}x^n = n \cdot x^{(n-1)}$$

$$\frac{d}{dx}e^{-x} = -e^{-x}$$

$$\hat{y} = \frac{1}{1 + e^{-z^2}}$$

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

$$\hat{y} = \frac{1}{1 + e^{-z^2}}$$

$$\Rightarrow e^{-z_2} = \frac{(1 - \hat{y})}{\hat{y}}$$

### **Neural Networks**

Cost function is the **Binary Cross Entropy**:

$$-\frac{1}{n}\sum_{i=1}^{n} [y_i \log_e \hat{y}_i + (1 - y_i) \log_e (1 - \hat{y}_i)]$$

Loss function:

$$L = -[y \log_e \hat{y} + (1 - y) \log_e (1 - \hat{y})]$$

$$\hat{y} = f_2(x) = \frac{1}{1 + e^{-x}}$$

Next lets calculate 
$$\frac{\partial}{\partial W_2} z_2$$

$$\Rightarrow \frac{\partial}{\partial W_2} z_2 = \frac{\partial}{\partial W_2} (a_1 W_2 + \beta_2)$$

$$\Rightarrow \frac{\partial}{\partial W_2} z_2 = a_1$$

Putting all three terms together:  $\frac{\partial}{\partial W_2} L = \frac{\partial}{\partial \hat{y}} L \frac{\partial}{\partial z_2} \hat{y} \frac{\partial}{\partial W_2} z_2$ 

$$\Rightarrow \frac{\partial}{\partial W_2} L = \frac{\hat{y} - y}{\hat{y}(1 - \hat{y})} \hat{y}(1 - \hat{y}) a_1$$

$$\Rightarrow \frac{\partial}{\partial W_2} L = (\hat{y} - y) a_1$$

Derivative of the cost over m observations:

$$\Rightarrow \frac{\partial}{\partial W_2} cost = \frac{1}{m} A_1^T (\hat{Y} - Y)$$

### **Neural Networks**

Cost function is the **Binary Cross Entropy**:

$$-\frac{1}{n} \sum_{i=1}^{n} \left[ y_i \log_e \hat{y}_i + (1 - y_i) \log_e (1 - \hat{y}_i) \right]$$

Loss function:

$$L = -[y \log_e \hat{y} + (1 - y) \log_e (1 - \hat{y})]$$

Sigmoid Activation function

$$\hat{y} = f_2(x) = \frac{1}{1 + e^{-x}}$$

 $A_1$  is an  $m \times n_1$  matrix  $\hat{Y}$  is a  $m \times n_2$  matrix Y is a  $m \times n_2$  matrix

$$\Rightarrow A_1^T$$
 is a  $n_1 \times m$  matrix

$$\Rightarrow \frac{1}{m} A_1^T (\hat{Y} - Y) \text{ is a } n_1 \times n_2 \text{ matrix}$$

$$\frac{\partial}{\partial \beta_2} L = \frac{\partial}{\partial \hat{y}} L \frac{\partial}{\partial z_2} \hat{y} \frac{\partial}{\partial \beta_2} z_2$$

Chain Rule:

L depends on  $\hat{y}, z_2$  in that order

$$\frac{\partial}{\partial \hat{y}} L = \frac{\hat{y} - y}{\hat{y}(1 - \hat{y})} \qquad \frac{\partial}{\partial z_2} \hat{y} = \hat{y}(1 - \hat{y})$$

Let's calculate  $\frac{\partial}{\partial \beta_2} z_2$ 

$$\Rightarrow \frac{\partial}{\partial \beta_2} z_2 = \frac{\partial}{\partial \beta_2} (a_1 W_2 + \beta_2) = 1$$

Putting all three terms together:  $\frac{\partial}{\partial \beta_2} L = \frac{\partial}{\partial \hat{v}} L \frac{\partial}{\partial z_2} \hat{y} \frac{\partial}{\partial \beta_2} z_2$ 

$$\Rightarrow \frac{\partial}{\partial \beta_2} L = \frac{\hat{y} - y}{\hat{y}(1 - \hat{y})} \hat{y}(1 - \hat{y})$$

$$\Rightarrow \frac{\partial}{\partial \beta_2} L = (\hat{y} - y)$$

Derivative of the cost over m observations:

$$\Rightarrow \frac{\partial}{\partial \beta_2} cost = \frac{1}{m} \sum_{i=1}^{m} (\hat{Y} - Y)$$

## **Neural Networks**

Cost function is the **Binary Cross Entropy**:

$$-\frac{1}{n}\sum_{i=1}^{n} \left[ y_i \log_e \hat{y}_i + (1 - y_i) \log_e (1 - \hat{y}_i) \right]$$

Loss function:

$$L = -[y \log_e \hat{y} + (1 - y) \log_e (1 - \hat{y})]$$

Sigmoid Activation function

$$\hat{y} = f_2(x) = \frac{1}{1 + e^{-x}}$$

Y is a  $m \times n_2$  matrix Y is a  $m \times n_2$  matrix

$$\Rightarrow \frac{1}{m} \sum_{i=1}^{m} (\hat{Y} - Y) \text{ is a } 1 \times n_2 \text{ vector}$$

Summation from 1..m collapses the *m* dimension

$$\frac{\partial}{\partial W_1} L = \frac{\partial}{\partial \hat{y}} L \frac{\partial}{\partial z_2} \hat{y} \frac{\partial}{\partial a_1} z_2 \frac{\partial}{\partial z_1} a_1 \frac{\partial}{\partial W_1} z_1$$
Chain Rule:
$$L \text{ depends on } \hat{y}, z_2, a_1, z_1 \text{ in that order}$$

$$\frac{\partial}{\partial \hat{y}} L = \frac{\hat{y} - y}{\hat{y}(1 - \hat{y})} \qquad \frac{\partial}{\partial z_2} \hat{y} = \hat{y}(1 - \hat{y})$$

First let's calculate 
$$\frac{\partial}{\partial a_1} z_2$$

$$\Rightarrow \frac{\partial}{\partial a_1} z_2 = \frac{\partial}{\partial a_1} (a_1 W_2 + \beta_2) = W_2$$

Next let's calculate  $\frac{\partial}{\partial z_1}a_1$ 

$$\Rightarrow \frac{\partial}{\partial z_1} a_1 = \frac{\partial}{\partial z_1} f_1(z_1)$$

Next let's calculate  $\frac{\partial}{\partial W_1} z_1$ 

$$\Rightarrow \frac{\partial}{\partial W_1} z_1 = \frac{\partial}{\partial W_1} (xW_1 + \beta_1) = x$$

## **Neural Networks**

Cost function is the **Binary Cross Entropy**:

$$-\frac{1}{n}\sum_{i=1}^{n} [y_i \log_e \hat{y}_i + (1 - y_i) \log_e (1 - \hat{y}_i)]$$

Loss function:

$$L = -[y \log_e \hat{y} + (1 - y) \log_e (1 - \hat{y})]$$

$$\hat{y} = f_2(x) = \frac{1}{1 + e^{-x}}$$

$$\frac{\partial}{\partial W_1} L = \frac{\partial}{\partial \hat{y}} L \frac{\partial}{\partial z_2} \hat{y} \frac{\partial}{\partial a_1} z_2 \frac{\partial}{\partial z_1} a_1 \frac{\partial}{\partial W_1} z_1$$

Putting all five terms together:

$$\Rightarrow \frac{\partial}{\partial W_1} L = \frac{\partial}{\partial \hat{y}} L \frac{\partial}{\partial z_2} \hat{y} \frac{\partial}{\partial a_1} z_2 \frac{\partial}{\partial z_1} a_1 \frac{\partial}{\partial W_1} z_1$$

$$\Rightarrow \frac{\partial}{\partial W_1} L = \frac{\partial}{\partial z_2} L \frac{\partial}{\partial a_1} z_2 \frac{\partial}{\partial z_1} a_1 \frac{\partial}{\partial W_1} z_1$$

$$\Rightarrow \frac{\partial}{\partial W_1} L = \left[ \frac{\partial}{\partial z_1} L \right] x \longleftarrow$$

Derivative of the cost over m observations:

$$\Rightarrow \frac{\partial}{\partial W_1} cost = \frac{1}{m} X^T \left[ \frac{\partial}{\partial Z_1} cost \right]$$

$$X^{T} \text{ is a } n_{0} \times m \text{ matrix}$$

$$\Rightarrow \frac{1}{m} X^{T} \left[ \frac{\partial}{\partial Z_{1}} cost \right] \text{ is a } n_{0} \times n_{1} \text{ matrix}$$

Chain Rule:

L depends on  $\hat{y}, z_2, a_1, z_1$  in that order

$$\frac{\partial}{\partial \hat{y}} L \frac{\partial}{\partial z_2} \hat{y} = \frac{\partial}{\partial z_2} L$$

$$\Rightarrow \frac{\partial}{\partial z_1} L = \frac{\partial}{\partial \hat{y}} L \frac{\partial}{\partial z_2} \hat{y} \frac{\partial}{\partial a_1} z_2 \frac{\partial}{\partial z_1} a_1$$

$$= \frac{\partial}{\partial z_2} L \frac{\partial}{\partial a_1} z_2 \frac{\partial}{\partial z_1} a_1$$

$$= \left( \left[ \frac{\partial}{\partial z_2} L \right] W_2^T \right) \odot \left[ \frac{\partial}{\partial z_1} f_1(z_1) \right]$$

$$\frac{\partial}{\partial Z_1} cost = \left( \left[ \frac{\partial}{\partial Z_2} cost \right] W_2^T \right) \odot \left[ \frac{\partial}{\partial Z_1} f_1(Z_1) \right]$$

$$\Rightarrow \frac{\partial}{\partial Z_1} cost = \left( (\hat{Y} - Y) W_2^T \right) \odot \left[ \frac{\partial}{\partial Z_1} f_1(Z_1) \right]$$

 $\frac{\partial}{\partial Z_1} cost$  is an  $m \times n_1$  matrix

#### **Neural Networks**

Cost function is the **Binary Cross Entropy**:

$$-\frac{1}{n}\sum_{i=1}^{n} \left[ y_i \log_e \hat{y}_i + (1 - y_i) \log_e (1 - \hat{y}_i) \right]$$

Loss function:

$$L = -[y \log_e \hat{y} + (1 - y) \log_e (1 - \hat{y})]$$

Sigmoid Activation function

$$\hat{y} = f_2(x) = \frac{1}{1 + e^{-x}}$$

$$\Rightarrow \frac{\partial}{\partial z_2} L = \frac{\partial}{\partial \hat{y}} L \frac{\partial}{\partial z_2} \hat{y}$$

$$= \frac{\hat{y} - y}{\hat{y}(1 - \hat{y})} \hat{y}(1 - \hat{y})$$

$$= \hat{y} - y$$

Vectorized over *m* observations

$$\Rightarrow \frac{\partial}{\partial Z_2} cost = \hat{Y} - Y$$

 $\hat{Y} - Y$  is an  $m \times n_2$  matrix

$$\frac{\partial}{\partial \beta_1} L = \frac{\partial}{\partial \hat{y}} L \frac{\partial}{\partial z_2} \hat{y} \frac{\partial}{\partial a_1} z_2 \frac{\partial}{\partial z_1} a_1 \frac{\partial}{\partial \beta_1} z_1$$
Chain Rule:
$$L \text{ depends on } \hat{y}, z_2, a_1, z_1 \text{ in that order}$$

Chain Rule:

#### Putting all five terms together:

$$\Rightarrow \frac{\partial}{\partial \beta_1} L = \frac{\partial}{\partial \hat{y}} L \frac{\partial}{\partial z_2} \hat{y} \frac{\partial}{\partial a_1} z_2 \frac{\partial}{\partial z_1} a_1 \frac{\partial}{\partial \beta_1} z_1$$

$$\Rightarrow \frac{\partial}{\partial \beta_1} L = \frac{\partial}{\partial z_2} L \frac{\partial}{\partial a_1} z_2 \frac{\partial}{\partial z_1} a_1 \frac{\partial}{\partial \beta_1} z_1$$

$$\Rightarrow \frac{\partial}{\partial \beta_1} L = \left[ \frac{\partial}{\partial z_1} L \right] \frac{\partial}{\partial \beta_1} z_1$$

$$\Rightarrow \frac{\partial}{\partial \beta_1} L = \left[ \frac{\partial}{\partial z_1} L \right] \frac{\partial}{\partial \beta_1} (xW_1 + \beta_1)$$

$$\Rightarrow \frac{\partial}{\partial \beta_1} L = \frac{\partial}{\partial z_1} L$$

$$\frac{\partial}{\partial \hat{y}} L \frac{\partial}{\partial z_2} \hat{y} = \frac{\partial}{\partial z_2} L$$

$$\Rightarrow \frac{\partial}{\partial z_1} L = \frac{\partial}{\partial \hat{y}} L \frac{\partial}{\partial z_2} \hat{y} \frac{\partial}{\partial a_1} z_2 \frac{\partial}{\partial z_1} a_1$$

$$= \frac{\partial}{\partial z_2} L \frac{\partial}{\partial a_1} z_2 \frac{\partial}{\partial z_1} a_1$$

$$= \left( \left[ \frac{\partial}{\partial z_2} L \right] W_2^T \right) \odot \left[ \frac{\partial}{\partial z_1} f_1(z_1) \right]$$

#### Derivative of the cost over m observations:

$$\Rightarrow \frac{\partial}{\partial \beta_1} cost = \frac{1}{m} \sum_{1=1}^{m} \frac{\partial}{\partial z_1} L$$

#### **Neural Networks**

Cost function is the **Binary Cross Entropy**:

$$-\frac{1}{n}\sum_{i=1}^{n} \left[ y_i \log_e \hat{y}_i + (1 - y_i) \log_e (1 - \hat{y}_i) \right]$$

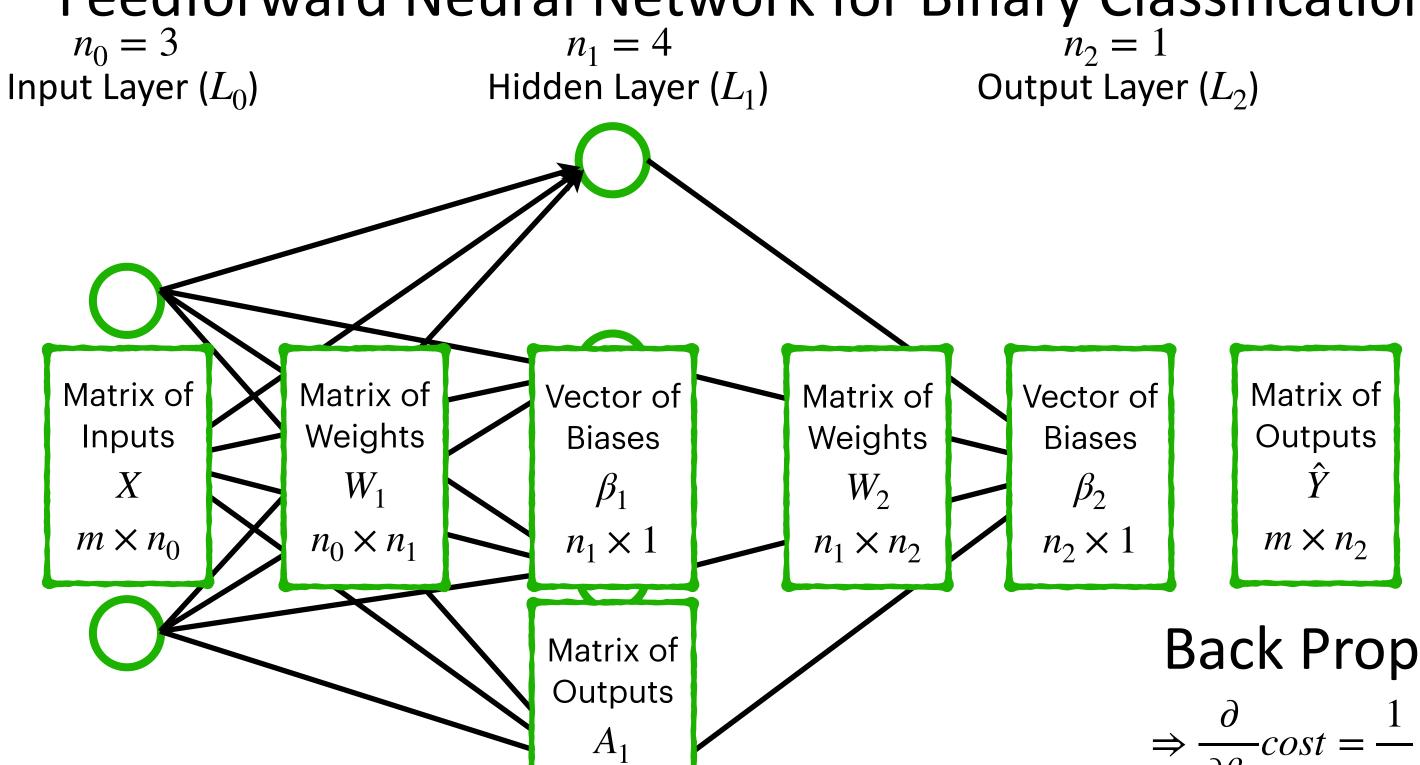
Loss function:

$$L = -[y \log_e \hat{y} + (1 - y) \log_e (1 - \hat{y})]$$

$$\hat{y} = f_2(x) = \frac{1}{1 + e^{-x}}$$

$$\frac{\partial}{\partial z_1} L \text{ is a } n_1 \times 1 \text{ vector}$$

$$\frac{1}{m} \sum_{1=1}^{m} \frac{\partial}{\partial z_1} L \text{ is an } n_1 \times 1 \text{ vector}$$



## **Neural Networks**

Cost function is the **Binary Cross Entropy**:

$$-\frac{1}{n}\sum_{i=1}^{n} [y \log_{e} \hat{y} + (1-y) \log_{e} (1-\hat{y})]$$

Loss function:

$$L = -[y \log_e \hat{y} + (1 - y) \log_e (1 - \hat{y})]$$

Sigmoid Activation function

$$\hat{y} = f_2(x) = \frac{1}{1 + e^{-x}}$$

**Back Propagation (Summary):** 

$$\Rightarrow \frac{\partial}{\partial \beta_2} cost = \frac{1}{m} \sum_{i=1}^{m} (\hat{Y} - Y)$$

$$\Rightarrow \frac{\partial}{\partial W_2} cost = \frac{1}{m} A_1^T (\hat{Y} - Y)$$

$$\Rightarrow \frac{\partial}{\partial \beta_1} cost = \frac{1}{m} \sum_{i=1}^{m} \frac{\partial}{\partial z_i} L$$

$$\Rightarrow \frac{\partial}{\partial W_1} cost = \frac{1}{m} X^T \left[ \frac{\partial}{\partial Z_1} cost \right]$$

$$\Rightarrow \frac{\partial}{\partial \beta_{2}} cost = \frac{1}{m} \sum_{i=1}^{m} (\hat{Y} - Y)$$

$$\Rightarrow \frac{\partial}{\partial \beta_{2}} cost = \frac{1}{m} \sum_{i=1}^{m} (\hat{Y} - Y)$$

$$\Rightarrow \frac{\partial}{\partial W_{2}} cost = \frac{1}{m} A_{1}^{T} (\hat{Y} - Y)$$

$$\Rightarrow \frac{\partial}{\partial \beta_{1}} cost = \frac{1}{m} \sum_{i=1}^{m} \frac{\partial}{\partial z_{1}} L$$

$$\Rightarrow \frac{\partial}{\partial W_{2}} cost = \frac{1}{m} \sum_{i=1}^{m} \frac{\partial}{\partial z_{1}} L$$

$$\Rightarrow \frac{\partial}{\partial W_{2}} cost = \frac{1}{m} \sum_{i=1}^{m} \frac{\partial}{\partial z_{1}} L$$

$$\Rightarrow \frac{\partial}{\partial W_{2}} cost = \frac{1}{m} \sum_{i=1}^{m} \frac{\partial}{\partial z_{1}} L$$

$$\Rightarrow \frac{\partial}{\partial W_{2}} cost = \frac{1}{m} \sum_{i=1}^{m} \frac{\partial}{\partial z_{1}} L$$

$$\Rightarrow \frac{\partial}{\partial W_{2}} cost = \frac{1}{m} \sum_{i=1}^{m} \frac{\partial}{\partial z_{1}} L$$

$$\Rightarrow \frac{\partial}{\partial W_{2}} cost = \frac{1}{m} \sum_{i=1}^{m} \frac{\partial}{\partial z_{1}} L$$

$$\Rightarrow \frac{\partial}{\partial W_{2}} cost = \frac{1}{m} \sum_{i=1}^{m} \frac{\partial}{\partial z_{1}} cost$$

$$\Rightarrow \frac{\partial}{\partial W_{2}} cost = \frac{1}{m} \sum_{i=1}^{m} \frac{\partial}{\partial z_{1}} cost$$

$$\Rightarrow \frac{\partial}{\partial W_{2}} cost = \frac{1}{m} \sum_{i=1}^{m} \frac{\partial}{\partial z_{1}} cost$$

$$\Rightarrow \frac{\partial}{\partial W_{2}} cost = \frac{1}{m} \sum_{i=1}^{m} \frac{\partial}{\partial z_{1}} cost$$

$$\Rightarrow \frac{\partial}{\partial W_{2}} cost = \frac{1}{m} \sum_{i=1}^{m} \frac{\partial}{\partial z_{1}} cost$$

$$\Rightarrow \frac{\partial}{\partial W_{2}} cost = \frac{1}{m} \sum_{i=1}^{m} \frac{\partial}{\partial z_{1}} cost$$

$$\Rightarrow \frac{\partial}{\partial W_{2}} cost = \frac{1}{m} \sum_{i=1}^{m} \frac{\partial}{\partial z_{1}} cost$$

$$\Rightarrow \frac{\partial}{\partial W_{2}} cost = \frac{1}{m} \sum_{i=1}^{m} \frac{\partial}{\partial z_{1}} cost$$

$$\Rightarrow \frac{\partial}{\partial W_{2}} cost = \frac{1}{m} \sum_{i=1}^{m} \frac{\partial}{\partial z_{1}} cost$$

$$\Rightarrow \frac{\partial}{\partial W_{2}} cost = \frac{1}{m} \sum_{i=1}^{m} \frac{\partial}{\partial z_{1}} cost$$

$$\Rightarrow \frac{\partial}{\partial W_{2}} cost = \frac{1}{m} \sum_{i=1}^{m} \frac{\partial}{\partial z_{1}} cost$$

$$\Rightarrow \frac{\partial}{\partial W_{2}} cost = \frac{1}{m} \sum_{i=1}^{m} \frac{\partial}{\partial z_{1}} cost$$

 $m \times n_1$ 

#### Back Propagation (Summary for a 3 layer network):

$$\Rightarrow \frac{\partial}{\partial Z_2} cost = dZ_2 = \hat{Y} - Y$$

$$\Rightarrow \frac{\partial}{\partial \beta_2} cost = dB_2 = \frac{1}{m} \sum_{i=1}^{m} dZ_2$$

$$\Rightarrow \frac{\partial}{\partial W_2} cost = dW_2 = \frac{1}{m} A_1^T dZ_2$$

$$\beta_{2} = \beta_{2} - learning\_rate \times \frac{\partial}{\partial \beta_{2}} cost$$

$$W_{2} = W_{2} - learning\_rate \times \frac{\partial}{\partial W_{2}} cost$$

$$D_{3} = \beta_{1} - learning\_rate \times \frac{\partial}{\partial \beta_{1}} cost$$

$$W_{4} = W_{1} - learning\_rate \times \frac{\partial}{\partial \beta_{1}} cost$$

$$W_{5} = f_{2}(x) = \frac{1}{1 + e^{-x}}$$

$$\hat{y} = f_{2}(x) = \frac{1}{1 + e^{-x}}$$

## **Neural Networks**

Cost function is the **Binary Cross Entropy**:

$$-\frac{1}{n}\sum_{i=1}^{n} [y \log_{e} \hat{y} + (1 - y) \log_{e} (1 - \hat{y})]$$

Loss function:

$$L = -[y \log_e \hat{y} + (1 - y) \log_e (1 - \hat{y})]$$

Sigmoid Activation function

$$\hat{y} = f_2(x) = \frac{1}{1 + e^{-x}}$$

$$\Rightarrow \frac{\partial}{\partial \beta_1} cost = dB_1 = \frac{1}{m} \sum_{1=1}^m dZ_1$$

 $\Rightarrow \frac{\partial}{\partial W_1} cost = dW_1 = \frac{1}{m} X^T dZ_1$ 

 $\Rightarrow \frac{\partial}{\partial Z_1} cost = dZ_1 = \left( dZ_2 W_2^T \right) \odot \left| \frac{\partial}{\partial Z_1} f_1(Z_1) \right|$ 

#### Back Propagation (Summary for a 4 layer network):

$$\Rightarrow \frac{\partial}{\partial Z_{3}} cost = dZ_{3} = \hat{Y} - Y$$

$$\Rightarrow \frac{\partial}{\partial \beta_{3}} cost = dB_{3} = \frac{1}{m} \sum_{i=1}^{m} dZ_{3}$$

$$\Rightarrow \frac{\partial}{\partial W_{3}} cost = dW_{3} = \frac{1}{m} A_{2}^{T} dZ_{3}$$

$$\Rightarrow \frac{\partial}{\partial Z_{2}} cost = dZ_{2} = (dZ_{3} W_{3}^{T}) \odot \left[ \frac{\partial}{\partial Z_{2}} f_{2}(Z_{2}) \right]$$

$$\Rightarrow \frac{\partial}{\partial W_{2}} cost = dW_{2} = \frac{1}{m} A_{1}^{T} dZ_{2}$$

$$\Rightarrow \frac{\partial}{\partial \beta_{2}} cost = dB_{2} = \frac{1}{m} \sum_{i=1}^{m} dZ_{2}$$

$$\Rightarrow \frac{\partial}{\partial Z_{1}} cost = dZ_{1} = (dZ_{2} W_{2}^{T}) \odot \left[ \frac{\partial}{\partial Z_{1}} f_{1}(Z_{1}) \right]$$

$$\Rightarrow \frac{\partial}{\partial W_{1}} cost = dW_{1} = \frac{1}{m} X^{T} dZ_{1}$$

$$\Rightarrow \frac{\partial}{\partial \beta_{1}} cost = dB_{1} = \frac{1}{m} \sum_{i=1}^{m} dZ_{1}$$

$$\beta_{3} = \beta_{3} - learning\_rate \times \frac{\partial}{\partial \beta_{3}} cost$$

$$W_{3} = W_{3} - learning\_rate \times \frac{\partial}{\partial W_{3}} cost$$

$$\beta_{2} = \beta_{2} - learning\_rate \times \frac{\partial}{\partial \beta_{2}} cost$$

$$W_{2} = W_{2} - learning\_rate \times \frac{\partial}{\partial W_{2}} cost$$

$$\beta_{1} = \beta_{1} - learning\_rate \times \frac{\partial}{\partial \beta_{1}} cost$$

$$W_{1} = W_{1} - learning\_rate \times \frac{\partial}{\partial W_{1}} cost$$

### **Neural Networks**

Cost function is the **Binary Cross Entropy**:

$$-\frac{1}{n}\sum_{i=1}^{n} [y \log_{e} \hat{y} + (1 - y) \log_{e} (1 - \hat{y})]$$

Loss function:

$$L = -[y \log_e \hat{y} + (1 - y) \log_e (1 - \hat{y})]$$

$$\hat{y} = f_2(x) = \frac{1}{1 + e^{-x}}$$

#### Related Tutorials & Textbooks

#### **Neural Networks**

An introduction to Neural Networks starting from a foundation of linear regression, logistic classification and multi class classification models along with the matrix representation of a neural network generalized to I layers with n neurons

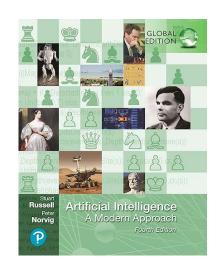
#### Forward and Back Propagation in Neural Networks

A deep dive into how Neural Networks are trained using Gradient Descent. Output predictions, are compared to observations to calculate loss and Backward propagation then computes gradients by working backward through the network

#### **Gradient Descent for Multiple Regression**

Gradient Descent algorithm for multiple regression and how it can be used to optimize k + 1 parameters for a Linear model in multiple dimensions.

#### **Recommended Textbooks**



<u>Artificial Intelligence: A Modern Approach</u>

by Peter Norvig, Stuart Russell

For a complete list of tutorials see:

https://arrsingh.com/ai-tutorials