

Week 2

- Binary Classification (output values b/w 0 and 1)
  - Logistic Regression
- Notations used

$(x, y) \rightarrow$  single training example where  
 $x \in \mathbb{R}^{n_x}$ ,  $y \in \{0, 1\}$

$M =$  training examples :  $\{(x^{(1)}, y^{(1)}), (x^{(2)}, y^{(2)}) \dots (x^m, y^m)\}$

$M_{\text{test}} =$  test examples

$$X = \begin{bmatrix} | & | & | & | \\ x^{(1)} & x^{(2)} & \dots & x^{(m)} \\ | & | & \dots & | \end{bmatrix} \quad \begin{matrix} 1 \\ n_x \\ m \end{matrix} \quad (\text{this structure makes the implementation easier})$$

$X \in \mathbb{R}^{n_x \times m} \rightarrow$  dimensional we matrix of  $n_x \times m$   
 $X.\text{shape} \rightarrow (n_x, m)$   
 gives output

$$Y = [y^{(1)}, y^{(2)}, y^{(3)}, \dots, y^{(m)}]$$

$Y \in \mathbb{R}^{1 \times m}; Y.\text{shape} \rightarrow (1, m)$

→ Logistic Regression

- learning algorithm that you use when the output labels  $Y$  in a supervised learning problem are all either zero or one
- Input given  $x$ , we want  $y = P(y=1/x)$

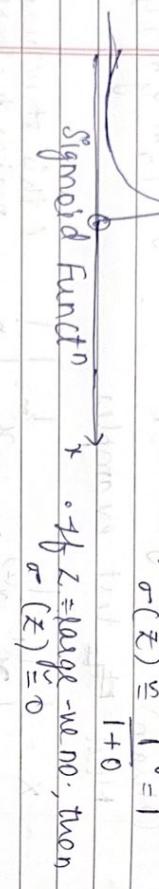
where  $0 \leq y \leq 1$   
 $b \Rightarrow$  real number

- Parameters:  $w \in \mathbb{R}^n$ ,  $b \in \mathbb{R}$
- Output =  $\hat{y} = \sigma(w^T x + b)$

$$\sigma(z) =$$

$$\frac{1}{1+e^{-z}}$$

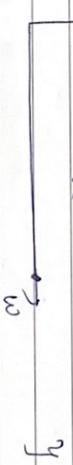
→ Gradient Descent



• If  $z$  is large, then

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

$y$



→ Cost function

$$\hat{y}^{(i)} = \sigma(w^T x^{(i)} + b), \text{ where } \sigma(z^{(i)}) = \frac{1}{1+e^{-z}}$$

Given  $(x^{(1)}, y^{(1)}), \dots, (x^{(m)}, y^{(m)})$ , want  $\hat{y}^{(1)} \approx y^{(1)}$

- Loss (error) function:  $J(\hat{y}, y) = \frac{1}{2} (\hat{y} - y)^2$   
 (optimization problem  $\Rightarrow$  you're with multiple local optimum), where

$$J(w, b) = w - \alpha \frac{dJ(w, b)}{dw}$$

$$b := b - \alpha \frac{dJ(w, b)}{db}$$

→ Logistic Regression Gradient descent

$$z = w^T x + b$$

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

$$J(w, b) = -(y \log(\hat{y}) + (1-y) \log(1-\hat{y}))$$

- $J(\hat{y}, y) = - (y \log \hat{y} + (1-y) \log (1-\hat{y}))$   
 (optimization problem  $\Rightarrow$  you're with multiple local optimum), where

If  $y=1 \Rightarrow J(\hat{y}, y) = -\log \hat{y} \leftarrow$  want  $\hat{y}$  large,

If  $y=0 \Rightarrow J(\hat{y}, y) = -\log(1-\hat{y}) \leftarrow$  want  $1-\hat{y}$  large, want  $\hat{y}$  to be small.

$$\begin{aligned} x_1 &\rightarrow [x_1 = w_1 x_1 + w_2 x_2 + b] \rightarrow \hat{y} = a = \sigma(z) \\ w_1 &\nearrow \\ b &\nearrow \end{aligned}$$

classmate  
Date \_\_\_\_\_  
Page \_\_\_\_\_

classmate  
Date \_\_\_\_\_  
Page \_\_\_\_\_

$$\frac{dL(a_i|y)}{da} = -\frac{y}{a} + \frac{1-y}{1-a}$$

Let  $\hat{J}=0$ ,  $dw_1=0$ ,  $dw_2=0$ ;  $db=0$   $\forall i=1 \dots m$   
For  $i=1 \dots m$

$$\frac{dz}{da} = \frac{dk}{dz} = \frac{dL(a_i|y)}{dz}$$

$$= a-y$$

$$= \frac{dL}{da} \cdot \frac{da}{dz}$$

$$= \sigma\left(\frac{-y}{a} + \frac{1-y}{1-a}\right) \left(a(1-a)\right)$$

$$dZ^{(i)} = a^{(i)} - y^{(i)}$$

$$dw_1 + = \alpha_1^{(i)} dz^{(i)}$$

$$dw_2 + = \alpha_2^{(i)} dz^{(i)}$$

$$db + = dz^{(i)}$$

$$J = -\left[\sum_{i=1}^m \log(a^{(i)}) + (1-y^{(i)}) \log(1-a^{(i)})\right]$$

$$dw_2 = X_2 \circ dz$$

$$db = dz$$

$$w_1 := w_1 - \alpha dw_1$$

$$w_2 := w_2 - \alpha dw_2$$

$$b := b - \alpha db$$

$$J / s = m$$

$$dw_1 / = m \Rightarrow dw_2 / = m \cdot db / = m$$

$$w_1 := w_1 - \alpha dw_1$$

$$w_2 := w_2 - \alpha dw_2$$

$$b := b - \alpha db$$

→ Gradient descent on m examples

→ Vectorization  
art of eliminating for loops in the code

Vectorized implementation is faster, computing  $w^T z$  in Logistic regression

$$z = np.dot(w, x)$$

import numpy → create array A → create array B →  
tic time → c = np.dot(A, B) → toc time →  
print vectorized version

eq.

$$v_i = \sum_j A_{ij} v_j$$

 $v = np.zeros((n, 1))$ 

for i

for j

 $v[j] += A[i][j] * v[j]$ 

 Non vectorised  
 Version 1

$$z = w^\top X + b$$

 $v = np.dot(A, v) \Rightarrow \text{vectorized}$ 

$$\text{eq: } v = \begin{bmatrix} v_1 \\ \vdots \\ v_n \end{bmatrix} \rightarrow v = \begin{bmatrix} e^{v_1} \\ \vdots \\ e^{v_n} \end{bmatrix}$$

 $\text{import numpy as np}$   
 $y = np.exp(v)$ 

→ Vectorizing logistic regression's gradient computation

$$\frac{d z^{(1)}}{d x^{(2)}} = a^{(1)} - y^{(1)}$$

$$\frac{d z^{(2)}}{d x^{(3)}} = a^{(2)} - y^{(2)}$$

$$\frac{d z^{(m)}}{d x^{(m)}} = a^{(m)} - y^{(m)}$$

$$\frac{d L}{d x^{(1)}} = \frac{d z^{(1)}}{d x^{(2)}} \dots \frac{d z^{(m)}}{d x^{(m)}}$$

$$A = \begin{bmatrix} a^{(1)} & \dots & a^{(m)} \end{bmatrix}$$

$$y = \begin{bmatrix} y^{(1)} & \dots & y^{(m)} \end{bmatrix}$$

$$dL = A - y = \begin{bmatrix} a^{(1)} - y^{(1)} & a^{(2)} - y^{(2)} & \dots & a^{(m)} - y^{(m)} \end{bmatrix}$$

$$z^{(0)} = w^\top x^{(0)} + b \quad z^{(1)} = w^\top x^{(1)} + b \quad z^{(2)} = w^\top x^{(2)} + b \quad z^{(3)} = w^\top x^{(3)} + b$$

$$a^{(1)} = \sigma(z^{(0)}) \quad a^{(2)} = \sigma(z^{(1)}) \quad a^{(3)} = \sigma(z^{(2)})$$

$$X = \begin{bmatrix} 1 & x^{(1)} & x^{(2)} & \dots & x^{(m)} \end{bmatrix}$$

$$Z = [Z^{(0)} Z^{(1)} Z^{(2)} \dots Z^{(m)}] = w^\top X + [b \ b \ \dots \ b]$$

$$= [w^\top x^{(0)} + b \ w^\top x^{(1)} + b \ \dots \ w^\top x^{(m)} + b]$$

$$db = \frac{1}{m} \sum_{i=1}^m d z^{(1)}$$

$$= \frac{1}{m} np.sum(d z^{(1)})$$

$$dw = \frac{1}{m} X dz^T + (X, b) \text{ for } q_1 - \epsilon$$

$$= \frac{1}{m} \left[ x_1^T, x_2^T, \dots, x_m^T \right] \begin{bmatrix} dz^{(1)} \\ dz^{(2)} \\ \vdots \\ dz^{(m)} \end{bmatrix}$$

$$= \frac{1}{m} [x^{(1)} dz^{(1)} + x^{(2)} dz^{(2)} + x^{(3)} dz^{(3)} + \dots + x^{(m)} dz^{(m)}]$$

\* Efficient code

$$z = w^T X + b$$

$$= np.\text{dot}(w.T, X) + b$$

$$A = \sigma(z)$$

$$dz = A - Y$$

$$dw = \frac{1}{m} X dz^T$$

$$db = \frac{1}{m} np.\text{sum}(dz)$$

$$w := w - \alpha dw$$

$$b := b - \alpha db$$

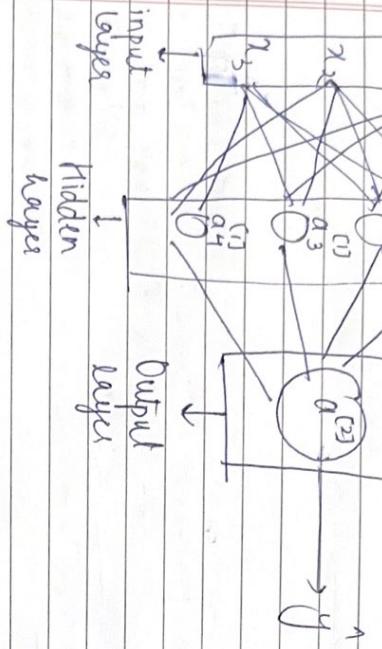
$$0 = db$$

WEEK 3:

→ Neural Network Representation

$$a^{[0]} = x$$

2 layer NN



→ Computing a Neural Network's Output

- for  $i = 1$  to  $m$ :

$$z^{(0)(i)} = w^{(0)}x^{(i)} + b^{(0)}$$

$$a^{(0)(i)} = \sigma(z^{(0)(i)})$$

$$z^{(1)(i)} = w^{(1)}a^{(0)(i)} + b^{(1)}$$

$$a^{(1)(i)} = \sigma(z^{(1)(i)})$$

$$a^{(0)} = \sigma(z^{(0)})$$

$$a^{(1)} = \sigma(z^{(1)})$$

$$a^{(2)} = \sigma(z^{(2)})$$

$$a^{(3)} = \sigma(z^{(3)})$$

$$a^{(4)} = \sigma(z^{(4)})$$

$$a^{(5)} = \sigma(z^{(5)})$$

similarly for mode 2

$$z^{(0)} = w^{(0)}x + b^{(0)}$$

$$a^{(0)} = \sigma(z^{(0)})$$

$$z^{(1)} = w^{(1)}a^{(0)} + b^{(1)}$$

$$a^{(1)} = \sigma(z^{(1)})$$

→ Vectorizing Across Multiple Examples

$$\begin{aligned} X^{(0)} &\longrightarrow a^{(2)} = \hat{y}^{(1)} \\ X^{(1)} &\longrightarrow a^{(2)(1)} = \hat{y}^{(1)} \\ &\dots \\ X^{(m)} &\longrightarrow a^{(2)(m)} = \hat{y}^{(m)} \end{aligned}$$

example  $i$ :  $\hat{y}^{(i)} \rightarrow$  layer 2

$$X = \begin{bmatrix} x^{(1)} & x^{(2)} & \dots & x^{(m)} \end{bmatrix}$$

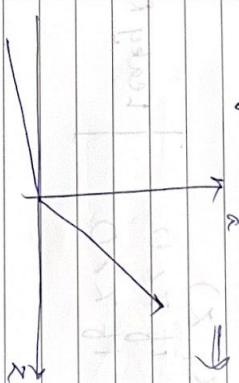
$$\begin{aligned} z^{(0)} &= w^{(0)}X + b^{(0)} \\ a^{(0)} &= \sigma(z^{(0)}) \\ z^{(1)} &= w^{(1)}a^{(0)} + b^{(1)} \\ a^{(1)} &= \sigma(z^{(1)}) \end{aligned}$$

Sigmoid range:  $0 \leq y \leq 1$   
 Tanh( $z$ ) range:  $-1 \leq y \leq 1$   
 Default values are for Rectified Linear Unit (ReLU)

$$A^{(l)} = \begin{bmatrix} 1 & a^{(l)(1)} & a^{(l)(2)} & \dots & a^{(l)(m)} \end{bmatrix} \quad \text{Hidden units}$$

training examples  $\rightarrow$

ReLU



vector of slope of  
 $z$  doesn't become 0,  
 it takes a slight slope

→ Activation function can be different for different layers

Sigmoid activatn functn is the most common, but other choices can work better

- Tanh( $z$ ) activatn function is better than sigmoid, as it ranges b/w -1 and 1 and has zero mean

→ Need of Non-linear Activation Function  
 Neural networks need a non-linear activation function to compute interesting functions

$$a = \tanh(z) = \frac{e^z - e^{-z}}{e^z + e^{-z}}$$

derivatives of Activation functions

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

$$g'(z) = g(z)(1-g(z))$$

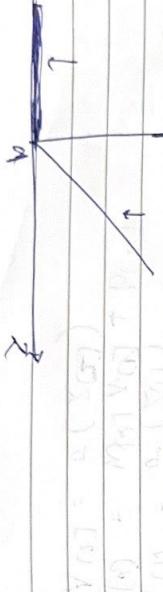
ReLU function

Tanh functn

$$g(z) = \tanh(z) = a$$

$$g'(z) = 1 - (\tanh(z))^2$$

$$g'(z) = 1 - a^2$$



- ReLU and leaky ReLU using sigmoid

$$f(x) = \max(0, x)$$

if  $x < 0$

ReLU

$$1 \quad \text{if } x \geq 0$$

sigmoid

$$g(x) = \max(0.01x, x)$$

$$f'(x) = \begin{cases} 0.01 & \text{if } x < 0 \\ 1 & \text{if } x > 0 \end{cases}$$

Leaky ReLU



## WEEK 4

→ Deep L.-layer Neural Network

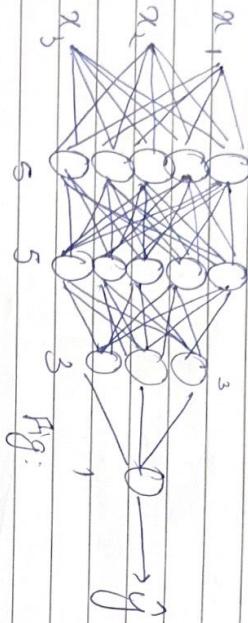


Fig:

$$L = 4 \quad (\# \text{ layers})$$

$$n^{(x)} = \# \text{ units in layer } x, n^{(1)} = 5, n^{(2)} = 3, n^{(3)} = n^{(4)} = 1$$

$$a^{(c_j)} = \text{activation}_j \text{ in layer } c$$

$$a^{(c_j)} = g^{(c_j)}(z^{(c_j)})$$

$$w^{(c_j)} = \text{weights for } z^{(c_j)}$$

→ Forward propagation in a Deep Network

for above figure

$$x: \mathbb{R}^5, z^{(0)} = W^{(0)}x + b^{(0)}$$

$$a^{(0)} = g^{(0)}(z^{(0)})$$

$$z^{(1)} = W^{(1)}a^{(0)} + b^{(1)}$$

$$a^{(1)} = g^{(1)}(z^{(1)})$$

$$\vdots$$

$$z^{(4)} = W^{(4)}a^{(3)} + b^{(4)}$$

$$a^{(4)} = g^{(4)}(z^{(4)})$$

## Vectorized Version

$$z^{(c_j)} = w^{(c_j)}a^{(c_{j-1})} + b^{(c_j)}$$

$$a^{(c_j)} = g^{(c_j)}(z^{(c_j)})$$

$$y = g(z^{(c_4)}) = A^{(c_4)}$$

→ Backpropagation for layer  $c$

Input:  $dz^{(c+1)}$

Output:  $da^{(c-1)}, dw^{(c)}, db^{(c)}$

$$dz^{(c)} = da^{(c)} * g^{(c)}'(z^{(c)})$$

$$dw^{(c)} = dz^{(c)} * A^{(c-1)}$$

$$db^{(c)} = dz^{(c)}$$

$$da^{(c-1)} = w^{(c)}.dz^{(c)}$$

$$dz^{(c)} = w^{(c+1)}dz^{(c+1)} + g^{(c)}(z^{(c)})$$

$$dA^{(c)} = dA^{(c)} * g^{(c)}'(z^{(c)})$$

$$dw^{(c)} = \frac{1}{m} dA^{(c)} * A^{(c-1)T}$$

$$db^{(c)} = \frac{1}{m} \text{np.sum}(dA^{(c)}, axis=1, keepdims=True)$$