

Jupyter/Python Notebooks

→ Shift + Enter to Execute program

- Restart kernel if get any errors
- submit assignment

Explanation of Logistic Regression cost Function

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

Interpret

$$\begin{cases} \hat{y} = P(y=1/x) \\ \text{if } y=1 : P(y/x) = \hat{y} \\ y=0 : P(y/x) = 1-\hat{y} \end{cases}$$

$$P(y/x) = \hat{y}^y (1-\hat{y})^{(1-y)}$$

$$\log P(Y/x) = \log y \log \hat{y} + (1-y) \log (1-\hat{y})$$

$$= -L(\hat{y}, y)$$

cost on m examples

$$\log P(\text{labels in training set}) = \sum_{i=1}^m P(y^{(i)}/x^{(i)})$$

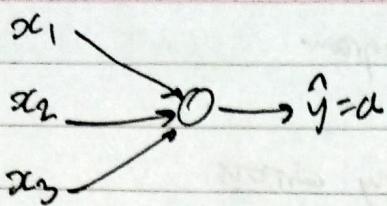
$$\log P(\dots) = \underbrace{\sum_{i=1}^m \log P(y^{(i)}/x^{(i)})}_{-L(\hat{y}^{(i)}, y^{(i)})}$$

maximum likelihood estimation

$$\text{cost } J(w, b) = -\sum_{i=1}^m L(\hat{y}^{(i)}, y^{(i)})$$

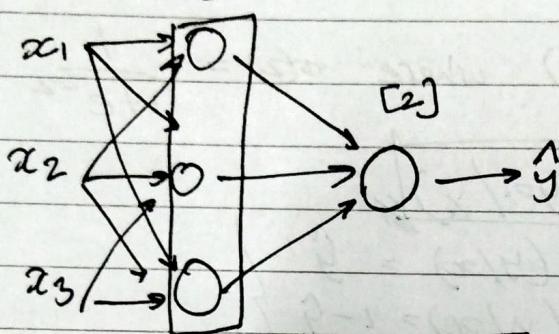
~~cost~~

$$\rightarrow \left[\begin{array}{c|c} \vdots & \vdots \\ \vdots & \vdots \\ \vdots & \vdots \end{array} \right] \quad \left[\begin{array}{c} \hat{y}^{(1)} \\ \vdots \\ \hat{y}^{(m)} \end{array} \right]$$

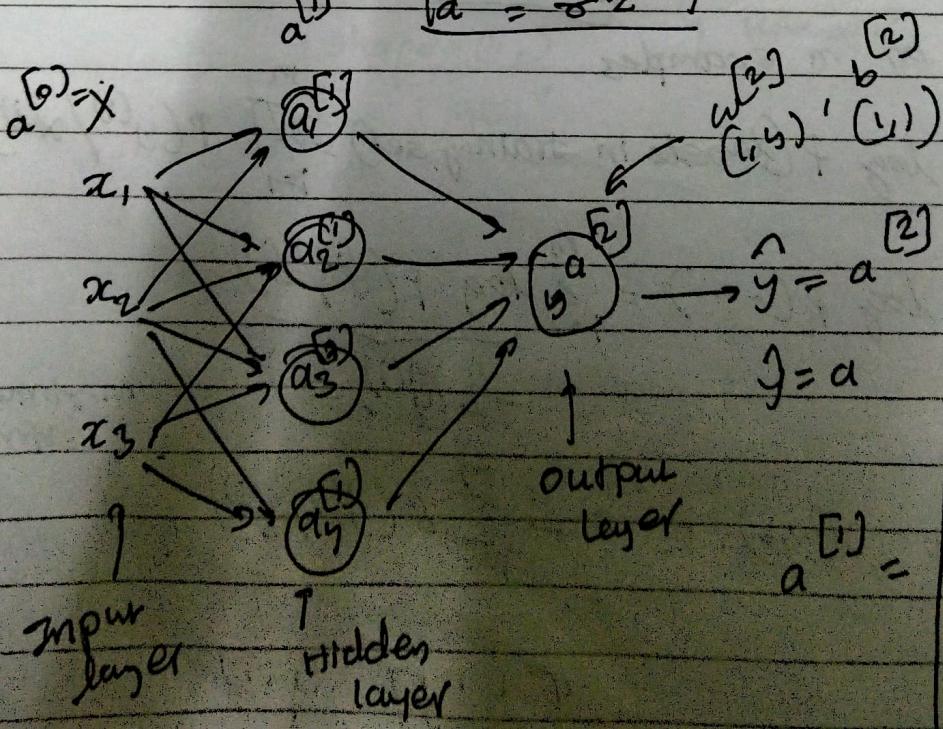


$$x \rightarrow [z = w^T x + b] \rightarrow [a = \sigma(z)] \rightarrow (a, y)$$

[1]



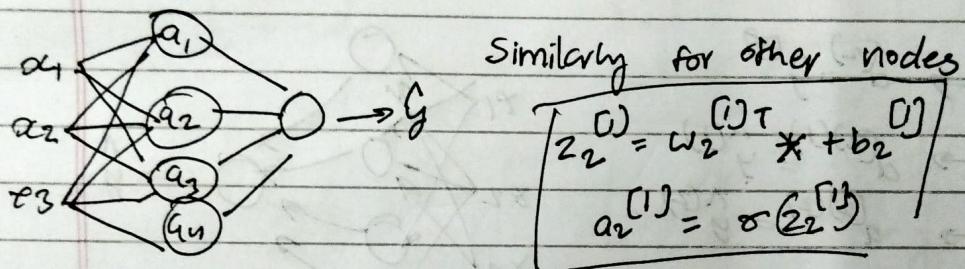
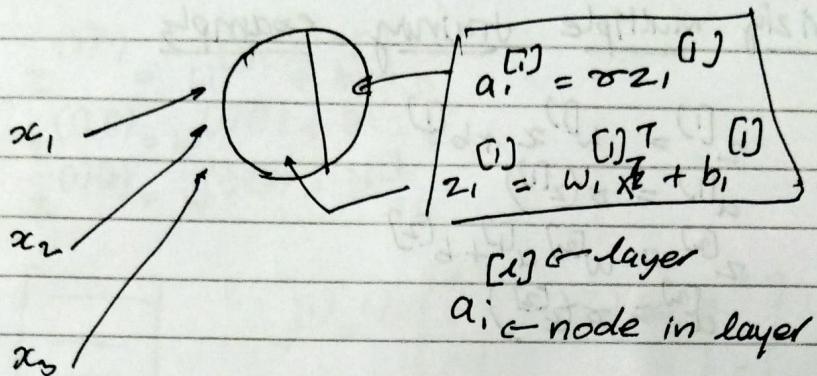
$$\begin{aligned} w^{(1)} &\rightarrow z^{(1)} = w^{(1)} x + b^{(1)} \rightarrow a^{(1)} = \sigma(z^{(1)}) \\ w^{(2)} &\rightarrow z^{(2)} = w^{(2)} a^{(1)} + b^{(2)} \end{aligned}$$



Output layer

$$a^{(1)} = \begin{bmatrix} a_1^{(1)} \\ a_2^{(1)} \\ a_3^{(1)} \\ a_4^{(1)} \end{bmatrix}$$

computing a Neural Network's output



$$\begin{bmatrix} z_1^{[1]} \\ z_2^{[1]} \\ z_3^{[1]} \\ z_4^{[1]} \end{bmatrix} = \begin{bmatrix} (w_1^{[1]})^T \\ w_2^{[1] T} \\ w_3^{[1] T} \\ w_4^{[1] T} \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \\ x_3 \end{bmatrix} + \begin{bmatrix} b_1^{[1]} \\ b_2^{[1]} \\ b_3^{[1]} \\ b_4^{[1]} \end{bmatrix} = \begin{bmatrix} w_1^{[1] T} x + b_1^{[1]} \\ w_2^{[1] T} x + b_2^{[1]} \\ w_3^{[1] T} x + b_3^{[1]} \\ w_4^{[1] T} x + b_4^{[1]} \end{bmatrix}$$

$n \times 3$

$w^{[1]}$

$b^{[1]} (a_i^{[1]})$

$$a^{[1]} = \begin{bmatrix} a_1^{[1]} \\ \vdots \\ a_n^{[1]} \end{bmatrix} = \sigma(z^{[1]})$$

$$\rightarrow z^{[2]} = w^{[2] T} x + b^{[2]}$$

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

$$\rightarrow z^{[2]} = w^{[2] T} a^{[1]} + b^{[2]}$$

$$\rightarrow a^{[2]} = \underbrace{\sigma(z^{[2]})}_{(1,1)} \quad (1,1)$$

$$\rightarrow a^{[2]} = \sigma(z^{[2]}) \quad (1,1)$$

Vectorizing multiple training examples

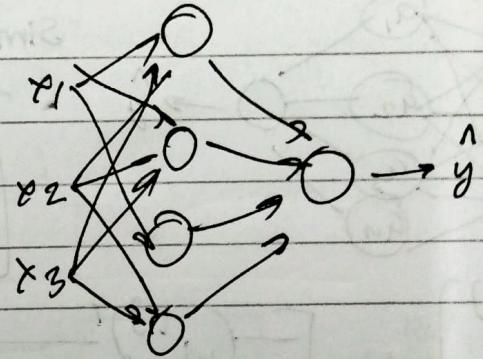
$$z^{(1)} = w^{(1)} x + b^{(1)}$$

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

$$z^{(2)} = w^{(2)} a^{(1)} + b^{(2)}$$

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

$$\begin{aligned} x &\rightarrow a^{(2)} = \hat{y} \\ x^{(1)} &\rightarrow a^{(2)(1)} = \hat{y}^{(1)} \\ x^{(2)} &\rightarrow a^{(2)(2)} = \hat{y}^{(2)} \\ \vdots & \\ x^{(m)} &\rightarrow a^{(2)(m)} = \hat{y}^{(m)} \end{aligned}$$



$a^{(2)(i)}$ example i
layer 2

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

$$\text{for } i=1 \text{ to } m: z^{(1)(i)} = w^{(1)} x^{(i)} + b^{(1)} \rightarrow z^{(1)} = w^{(1)} x + b^{(1)}$$

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

$$z^{(2)(i)} = w^{(2)} a^{(1)(i)} + b^{(2)} \rightarrow z^{(2)} = w^{(2)} A^{(1)} + b^{(2)} \rightarrow A^{(2)} = \sigma(z^{(2)})$$

$$X = \left[\begin{array}{c|c|c|c} x^{(0)} & x^{(1)} & \cdots & x^{(m)} \\ \hline x^{(0)} & x^{(1)} & \cdots & x^{(m)} \end{array} \right]^{(n_x, m)}$$

$$z^{(0)} = \left[\begin{array}{c|c|c|c} z^{(0)(1)} & z^{(0)(2)} & \cdots & z^{(0)(m)} \\ \hline z^{(0)(1)} & z^{(0)(2)} & \cdots & z^{(0)(m)} \end{array} \right]$$

training examples

$$A^{(0)} = \left[\begin{array}{c|c|c|c} a^{(0)(1)} & a^{(0)(2)} & \cdots & a^{(0)(m)} \\ \hline a^{(0)(1)} & a^{(0)(2)} & \cdots & a^{(0)(m)} \end{array} \right]$$

hidden units
nodes in
neural network

Justification for vectorized implementation

$$z^{(1)(1)} = w^{(1)(1)} + b^{(1)}_0$$

$$z^{(1)(2)} = w^{(1)(2)} + b^{(1)}_1$$

$$z^{(1)(3)} = w^{(1)(3)} + b^{(1)}_2$$

$$w^{(1)} = \begin{bmatrix} \dots \\ \dots \\ \dots \end{bmatrix} \quad w^{(1)}x^{(1)} = \begin{bmatrix} \vdots \\ \vdots \\ \vdots \end{bmatrix} \quad w^{(1)}x^{(2)} = \begin{bmatrix} \vdots \\ \vdots \\ \vdots \end{bmatrix} \quad w^{(1)}x^{(3)} = \begin{bmatrix} \vdots \\ \vdots \\ \vdots \end{bmatrix}$$

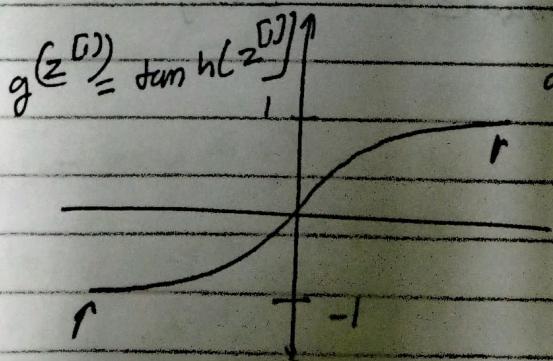
$$w^{(1)} \begin{bmatrix} x_1 & x_2 & x_3 \\ 1 & 1 & 1 \end{bmatrix} = \begin{bmatrix} \vdots & \vdots & \vdots \\ \vdots & \vdots & \vdots \\ \vdots & \vdots & \vdots \end{bmatrix} = \begin{bmatrix} z^{(1)(1)} \\ z^{(1)(2)} \\ z^{(1)(3)} \end{bmatrix} = \begin{bmatrix} w^{(1)(1)}x^{(1)} + b^{(1)}_0 \\ w^{(1)(2)}x^{(2)} + b^{(1)}_1 \\ w^{(1)(3)}x^{(3)} + b^{(1)}_2 \end{bmatrix}$$

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

ACTIVATION FUNCTIONS

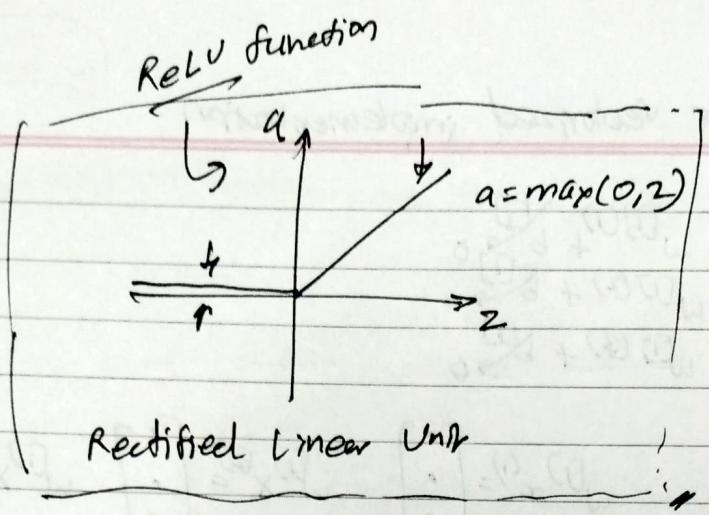
$$\sigma \rightarrow \text{Activation function} \quad \sigma = \frac{1}{1+e^{-z}}$$

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

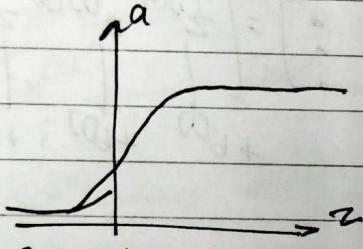


$$a = \tanh(z) = \frac{e^z - e^{-z}}{e^z + e^{-z}}$$

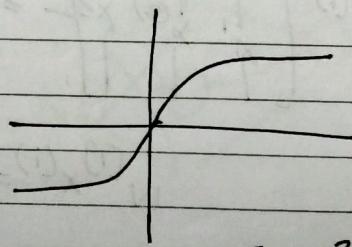
$y \in [0, 1] \leftarrow \text{sigmoid activation function}$



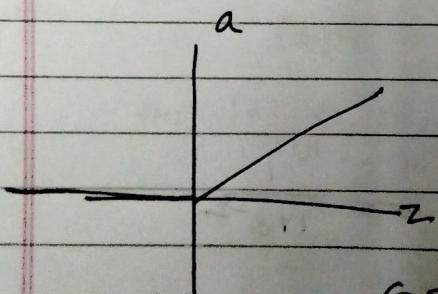
Pros & cons of activation functions



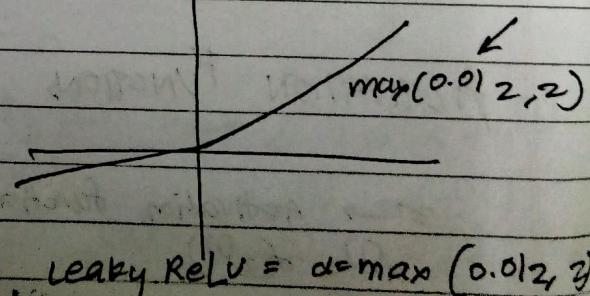
Sigmoid - $a = \frac{1}{1+e^{-z}}$



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



ReLU $a = \max(0, z)$



Leaky ReLU = $a = \max(0.01z, z)$

Why do you need non-linear activation function?

$$\alpha^{(l)} = g^{(l)}(z^{(l)})$$

$g^{(l)}(z) = z \rightarrow$ linear activation function

hidden layer doesn't at all come into picture
as it just outputs a linear like

Hidden units should not use linear activation function
layers

linear activation function are ok for last layer

Slope & functions

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

$$\frac{d}{dz} g(z) = g(z)(1-g(z))$$

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

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

ReLU

$$g(z) = \max(0, z)$$

$$g'(z) = \begin{cases} 0 & \text{if } z < 0 \\ 1 & \text{if } z \geq 0 \end{cases}$$

leaky ReLU

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

$$g'(z) = \begin{cases} 0.01 & \text{if } z < 0 \\ 1 & \text{if } z \geq 0 \end{cases}$$

Date _____
Page _____

~~classmate~~ drawing examples

$z \& A$ vectors have
matrices have shapes $\begin{pmatrix} y, m \end{pmatrix}$

THIS is no. of neurons in that layer
itself

Gradient Descent for neural networks

Parameters $\begin{pmatrix} w^{[1]} \\ n^{[0]} \end{pmatrix}, \begin{pmatrix} b^{[1]} \\ n^{[0]} \end{pmatrix}, \begin{pmatrix} w^{[2]} \\ b^{[2]} \end{pmatrix}, \begin{pmatrix} n^{[1]}, n^{[0]} \\ n^{[2]}, n^{[1]} \end{pmatrix}, \underline{\begin{pmatrix} n^{[2]}, 1 \end{pmatrix}}$

$$n_p = n^{[0]}, n^{[1]}, n^{[2]} = 1$$

cost function

$$J(w^{[1]}, b^{[1]}, w^{[2]}, b^{[2]})$$

$$= \frac{1}{m} \sum_{i=1}^m L(\hat{y}_i, y_i)$$

gradient descent:

Repeat {

compute predictions $(\hat{y}^{(i)}, i=1 \dots m)$

$$\frac{d w^{[1]}}{d w^{[1]}} = \frac{d J}{d w^{[1]}}, \frac{d b^{[1]}}{d b^{[1]}} = \frac{d J}{d b^{[1]}}, \dots$$

$$w^{[1]} = w^{[1]} - \alpha d w^{[1]}$$

$$b^{[1]} = b^{[1]} - \alpha d b^{[1]}$$

$$w^{[2]} = w^{[2]} -$$

Forward propagation

$$z^{[1]} = w^{[1]} x + b^{[1]}$$

$$A^{[1]} = g^{[1]}(z^{[1]})$$

$$z^{[2]} = w^{[2]} x + b^{[2]}$$

$$A^{[2]} = g^{[2]}(z^{[2]})$$

Back propagation

$$d z^{[2]} = A^{[2]} - y$$

$$d w^{[2]} = \frac{1}{m} d z^{[2]} A^{[1]T}$$

$$d b^{[2]} = \frac{1}{m} \text{np.sum}(d z^{[2]}, \text{axis}=1, \text{keepdims}=1)$$

$\text{axis}=1$ horizontally

keepdims → to save from horizontal funny outputs →

$$d\hat{z}^{[1]} = w^{[2]} + d\hat{z}^{[2]} \quad \# \quad g^{[1]}'(\hat{z}^{[0]}) \quad (n^{[1]}, m)$$

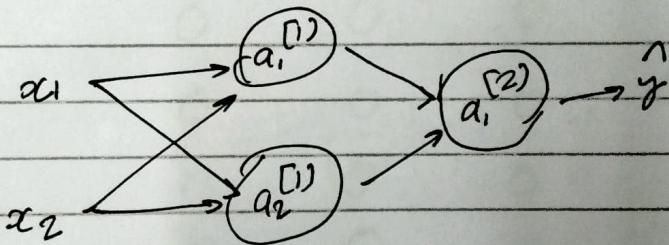
$\underbrace{(n^{[1]}, m)}_{\text{elementwise product}}$

$$dw^{[1]} = \frac{1}{m} d\hat{z}^{[0]} x^T$$

$$db^{[1]} = \frac{1}{m} \text{np.sum}(d\hat{z}^{[1]}, \text{axis}=1, \text{keepdims=True})$$

$(n^{[1]}, 1) \quad \text{not } (n^{[1]}, 0)$
so keepdims used

RANDOM INITIALIZATION



If you initialise weights to zero.

$$w^{[1]} = \begin{bmatrix} 0 & 0 \\ 0 & 0 \end{bmatrix} \quad b^{[1]} = \begin{bmatrix} 0 \\ 0 \end{bmatrix}$$

$$a_1^{[1]} = a_2^{[1]} \quad d\hat{z}_1^{[1]} = d\hat{z}_2^{[1]}$$

$$w^{[2]} = \begin{bmatrix} 0 & 0 \end{bmatrix}$$

$$\therefore dw = \begin{bmatrix} u & v \\ u & v \end{bmatrix}$$

$$w^{[1]} = w^{[1]} - \alpha dw$$

$$w^{[1]} = \begin{bmatrix} - & - \\ - & - \end{bmatrix}$$

$\therefore w^{(1)} = \text{np. random. randn}((2, 2)) \neq 0.0$

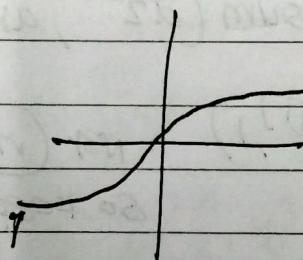
$b^{(1)} = \text{np. zeros}((2, 1))$

why?

$w^{(2)} = \text{np. } - - -$

$b^{(2)} = 0$

So we want small
random values not
huge as they will affect



↓ As you will end up at
the large values of z
saturating & slowing
learning algorithm