

Binary classification

- * There will be output 0 or 1
- * But it doesn't mean it always gives you True or not.

* It will be used to differentiate between pictures, words, to translate means wider applications other than regression which include numbers

$$(n, y) \quad n \in \mathbb{R}^{nx}, y \in \{0, 1\}$$

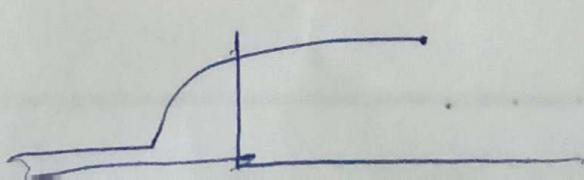
$$X = \left[x^{(1)}, x, \dots, x^n \right]$$

- * Logistic Regression (Supervised) (used to predict cat or not like task)

Given x, y

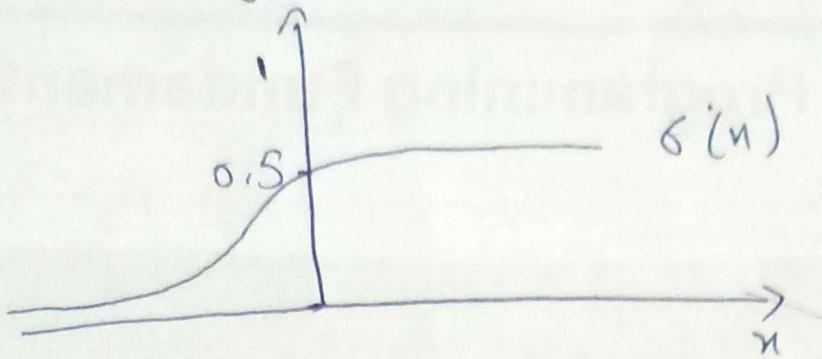
Given x , want $\hat{y} = P(y=1 | x)$

Parameter $w \in \mathbb{R}^n, b \in \text{Real num}$
output $\hat{y} = g(w^T x + b)$



sigmoid

thus will form yes or no



σ = sigmoid

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

* If z large $\sigma(z) \approx \frac{1}{1} = 1$

* If z is small $= \frac{1}{1+e^{-z}} \approx \frac{1}{1+\text{Bignum}} \approx 0$

$$\hat{y} = \sigma(\theta^T x)$$

$$\theta = \begin{bmatrix} \theta_0 \\ \theta_1 \\ \theta_2 \\ \vdots \\ \theta_{Ax} \end{bmatrix} \quad b \leftarrow *$$

Just to understand

Logistic Regression Cost Function

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

loss error $L(\hat{y}, y) = \frac{1}{2} (\hat{y} - y)^2$

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

Condition

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

If $y=0$: $L(\hat{y}, y) = -\log(1-\hat{y}) \leftarrow$ want $\log(1-\hat{y})$ large

cost function $= \frac{1}{m} \sum_{i=1}^m L(\hat{y}^{(i)}, y^{(i)}) = \left[-\frac{1}{m} \sum_{i=1}^m \left[y^{(i)} \log \hat{y}^{(i)} + (1-y^{(i)}) \log(1-\hat{y}^{(i)}) \right] \right]$

$$\cancel{\frac{1}{m}} = -\frac{1}{m} \sum_{i=1}^m [y^{(i)} \log \hat{y}^{(i)} + (1-y^{(i)}) \log(1-\hat{y}^{(i)})]$$

Gradient Descent

Gradient Descent

* Cost and loss function difference

The loss function computes the error for single training example,

* Cost function

The loss function computes the average of loss function of the entire training set

I mean, Cost (Loss)

Gradient descent

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

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

Gradient Descent

Gradient Descent

* Cost and loss function difference

The loss function computes the error for
single training example,

* Cost function

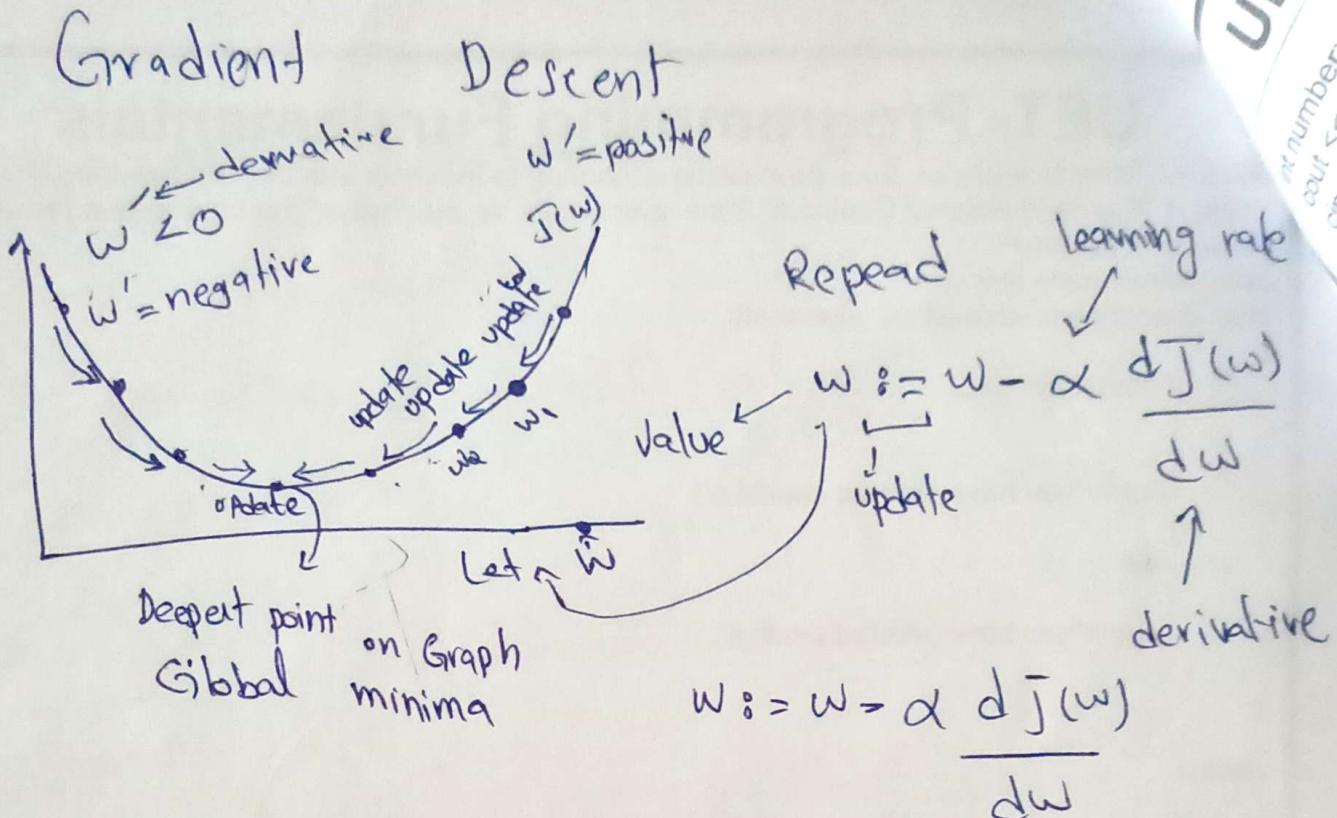
The loss function computes the ~~error~~
average of loss function of the entire
training set

I mean, cost (loss)

Gradient descent

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

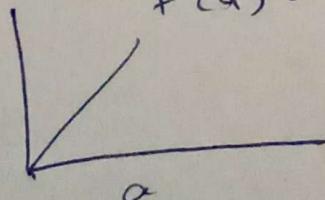
$$\text{cost function } J(w, b) = \frac{1}{m} \sum_{i=1}^m L(g(x^{(i)}), y^{(i)}) = \frac{1}{m} \sum y^{(i)} \log \frac{\hat{y}^{(i)}}{1 - \hat{y}^{(i)}}$$



⇒ Hence, we get the accurate value by reaching global minima.

$$b := b - \alpha \frac{dJ(w, b)}{db} \quad \xrightarrow{\text{partial derivative}}$$

learning derivative



$$f(a) = 3a$$

$$\text{if } a = 2 \Rightarrow f(a) = 6$$

\Rightarrow Slope = $\frac{\text{height}}{\text{width}}$ of triangle

\Rightarrow derivatives are small values
0.0000000
0.001 it should
be tiny

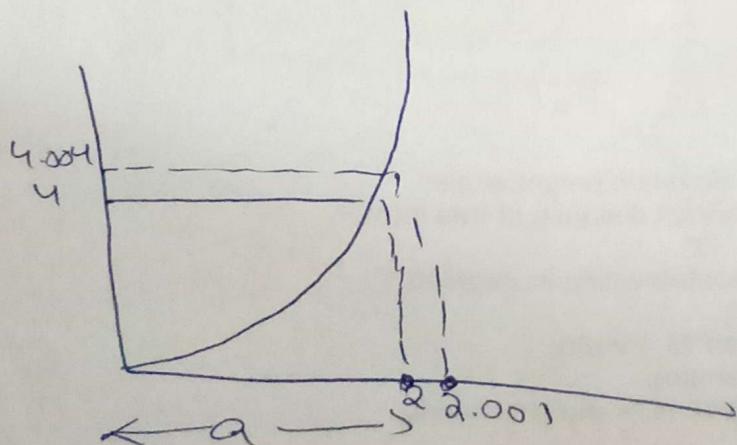
Example : 2

$$f(a) = a^2$$

$$a = 2$$

$$a = 2.001$$

$$\begin{aligned} f(a) &= 4 \\ f(2.001) &= 4.004 \end{aligned}$$



$$\frac{d}{da} a^2 \Rightarrow 2a$$

$$a = 5 \quad f(a) = a^2 \Rightarrow 25$$

$$a = 5.001 \Rightarrow f(a) \Rightarrow 25.010$$

we got the 25.010 value in $\frac{d}{da} a^2$
derivative we got the 25.010 value in $\frac{d}{da} a^2$

$$f(a) = a^3$$

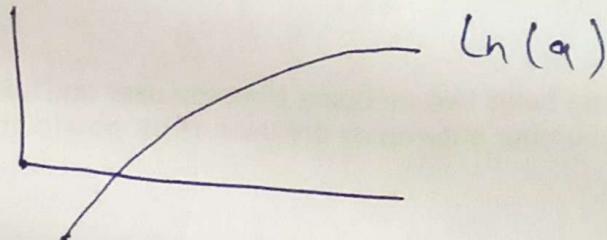
(let $a = 2$)
 $a = 2.0001$

$$\frac{d}{da} f(a) = 3a^2$$

~~\approx~~

$$f(a) = 2.0018$$

$$f(a) = \log_e(a)$$
$$= \ln(a)$$



$$a = 2$$

$$a = 2.001$$

$$\frac{df(a)}{da} = \frac{d}{da} \log a$$

$$df = \frac{1}{a}$$

$$f(a)' = \frac{1}{2} = 0.5$$

$$f(a) = 2.0005$$

Q.

Computation Graph

$$J(a, b, c) = 3(a + bc) \quad J = 33$$

$U = bc = 6, a=5, b=3, c=2$

$V = a+U = 11$

$J = 3V = 33$

$$\frac{dJ}{da} = ? \quad \because J = 3V$$

$\checkmark \rightarrow V = 11.001$

$\Rightarrow J = 33.003$

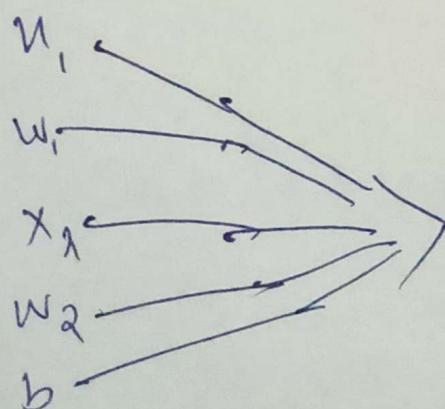
$J' = 3$

Example $a = 5.001$

$\frac{dJ}{da} = J'$

$$\text{derivative} = a + a'$$

Use case



$$Z = w_1x_1 + w_2x_2 + b \Rightarrow q = g(Z) \rightarrow \text{loss}$$

$$dq = \frac{\partial L(q, y)}{\partial q}$$

$$= -\frac{y}{q} + \frac{1-y}{1-q}$$

$$dz = \frac{df}{dz} = \frac{dL(\alpha, y)}{dz}$$

$$= \alpha \hat{y}$$

$$\stackrel{(a-y)}{=} \frac{df}{da} \cdot \frac{da}{dz} \xrightarrow{\quad} a(1-a)$$

$$= -\frac{\hat{y}}{a} + \underbrace{\frac{1-\hat{y}}{1-a}}$$

$$J\omega = (y^{(i)}, y^{(0)})$$

$$J(w, b) \leftarrow \frac{1}{m} \sum_{i=1}^m l(a^{(i)}, y^{(i)})$$

$$a^{(i)} = \hat{y}^{(i)} = \sigma(z^{(i)}) = \sigma(w^T x^{(i)} + b)$$

$$\frac{\partial}{\partial w_1} J(w, b) = \frac{1}{m} \sum_{i=1}^m \frac{\partial}{\partial w_1} l(a^{(i)}, y^{(i)})$$

}

$$dw_1^{(i)} = (a^{(i)}, y^{(i)})$$

$$J=0, dw_1=0, dw_a=0, db=0$$

For $i=1$ to m

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

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

$$J_t = -[y^{(i)} \log a^{(i)} + (1-y^{(i)}) \log (1-a^{(i)})]$$

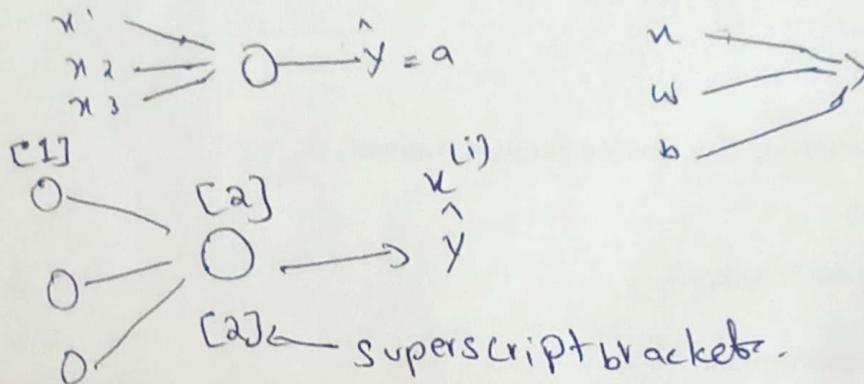
$$d\hat{z}^{(i)} = \alpha^{(i)} - y^{(i)}$$

$$dw_{1,t} = n_1^{(i)} d\hat{z}^{(i)}$$

$$dw_{2,t} = n_2^{(i)} d\hat{z}^{(i)}$$

$$db_t = d\hat{z}^{(i)}$$

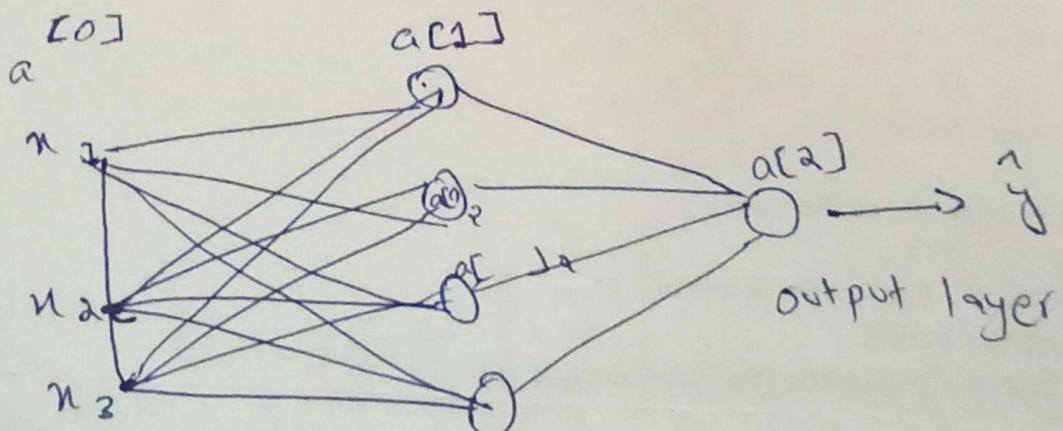
Shallow + Neural network



$$w^T n + b \rightarrow a = \sigma(z) \rightarrow L(a, y)$$

$$z^{[0]} = w^{[0]} n + b^{[0]} \rightarrow a^{[1]} = \sigma(z^{[1]}) \rightarrow \sum^{[2]} = w^{[2]} a^{[1]} + b^{[2]}$$

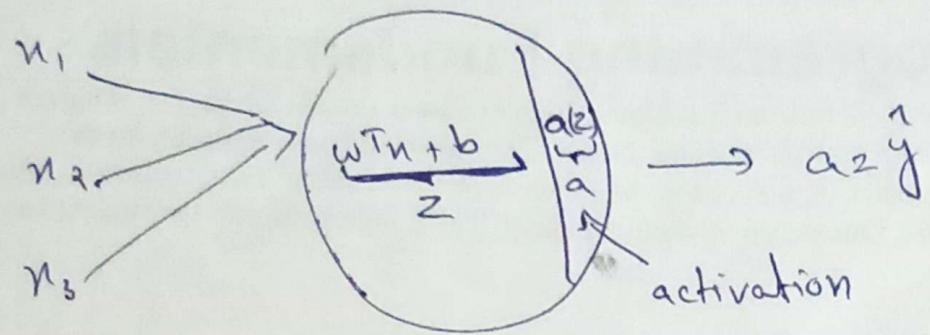
Neural Network representation



Hidden layer

$$a^{[1]} = \begin{bmatrix} a_1^{[1]} \\ a_2^{[1]} \\ \vdots \\ a_y^{[1]} \end{bmatrix}$$

Representative of Neural network



$$z^{[1]} = w_1^{[1]T} \times b_1^{[1]}$$

$$a_1^{[1]} = g(z_i^{[1]})$$

u_1, u_2, u_3

$$z^{[2]} = w_2^{[1]T} u + b_2^{[1]}$$

$$a_2^{[2]} = g(z_2^{[2]})$$

u_1, u_2, u_3

transpose

$$w^T = \begin{bmatrix} w_1^{[1]T} & w_2^{[1]T} & \dots \\ w_a^{[1]} & w_a^{[2]} & \dots \\ w_3^{[1]T} & w_3^{[2]T} & \dots \\ w_m^{[1]T} + b & w_m^{[2]T} + b & \dots \end{bmatrix} = \begin{bmatrix} z_1^{[1]} \\ z_2^{[1]} \\ z_3^{[1]} \\ \vdots \end{bmatrix}$$

$$\alpha = \begin{bmatrix} a_1^{[1]} \\ a_2^{[1]} \\ \vdots \\ a_n^{[1]} \end{bmatrix} = g(z)^{[1]}$$

Vectoring across multiple example.

layer working on example
 $\underline{a}^{[2]}, \underline{(i)}$

layer $\rightarrow [a]$ working on example
 $\underline{a}^{[1]}, \underline{(i)}$

for $l = 1$ to m_1

$$z^{[1](i)} = w^{[1]}_n z^{(i)} + b^{[1]}$$

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

$$z^{[2](i)} = w^{[2]}_n a^{[1](i)} + b$$

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

activation

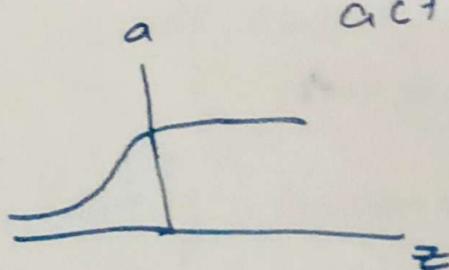
$$z^{[2](i)} = w^{[2]}_n a^{[1](i)} + b^{[2]}$$

for activation

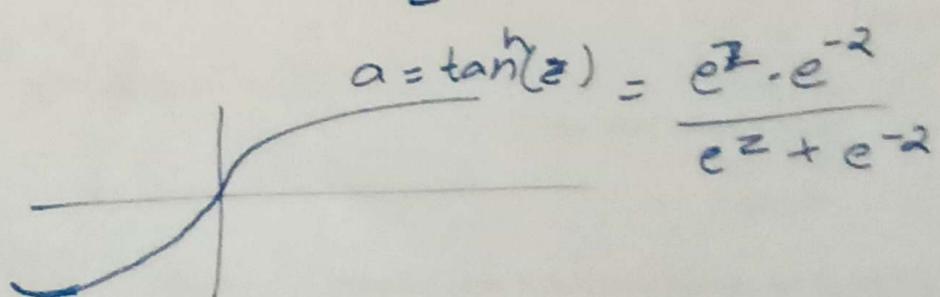
$$A^{[1]} = g(z^{[1]})$$

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

One hidden activation Layer. Choice of activation function:-



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

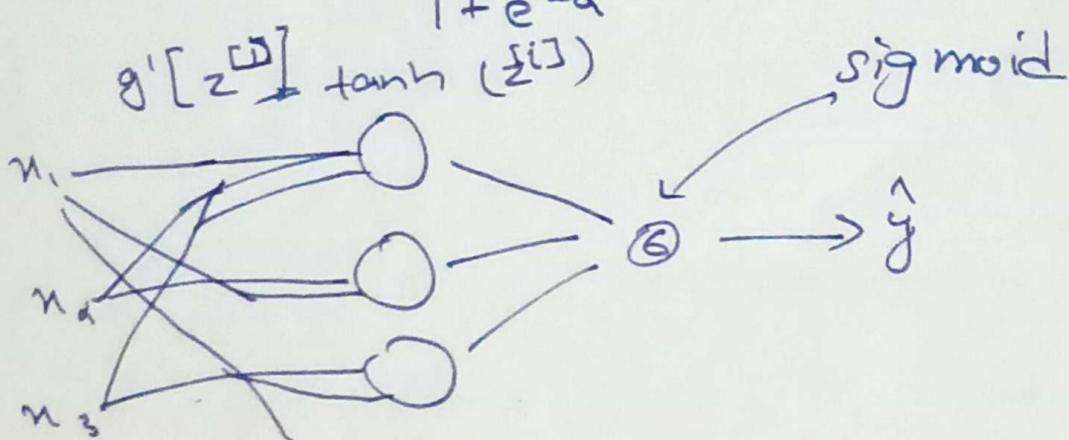


$$a = \tanh(z) = \frac{e^z - e^{-z}}{e^z + e^{-z}}$$

$g(z^{(i)}) = \tanh(z^{(i)})$ is superior

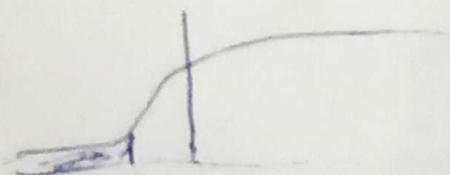
than $a = 1$

$$g'(z^{(i)}) = \tanh(z^{(i)})$$



Limitation of σ function and tanh func

\Rightarrow When value is very Large or very small
it cause gradient very large or small



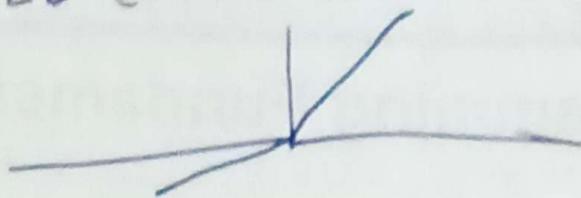
Solution:- Rectify Linear Unit (ReLU)

$$z^+ \Rightarrow z = 1$$

$$\therefore a = m^+(0, z)$$

$$z^- \Rightarrow z = 0$$

Leaky RLV (another formula)



$$\max(0.01z, 0)$$

- ② Why we need activation function in neural network

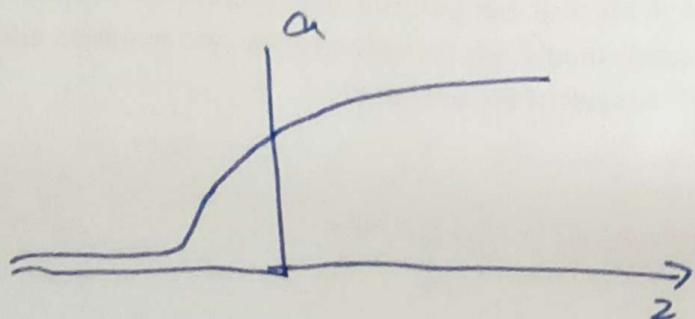
$$a = \sigma(z)$$

$$\Rightarrow a^{[2]} = w^{[2]}(w^{[1]}\mathbf{n} + b^{[1]}) + b^{[2]}$$

$$\Rightarrow (w^2 w^1) \mathbf{n} + [w^2 b^1 + b^{[2]}]$$
$$= w^1 \mathbf{n} + b$$

$$= w^1 \mathbf{n} + b$$

Derivative of activation function



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

$$\frac{d}{dt} g(z) = \text{slope of } g(z) \text{ at } z$$

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

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

$$= g(z)(1 - g(z)) \Rightarrow z = 10, g(z) \approx 1$$

$$\frac{d}{dz}g(z) \approx 1(1-1) = 0$$

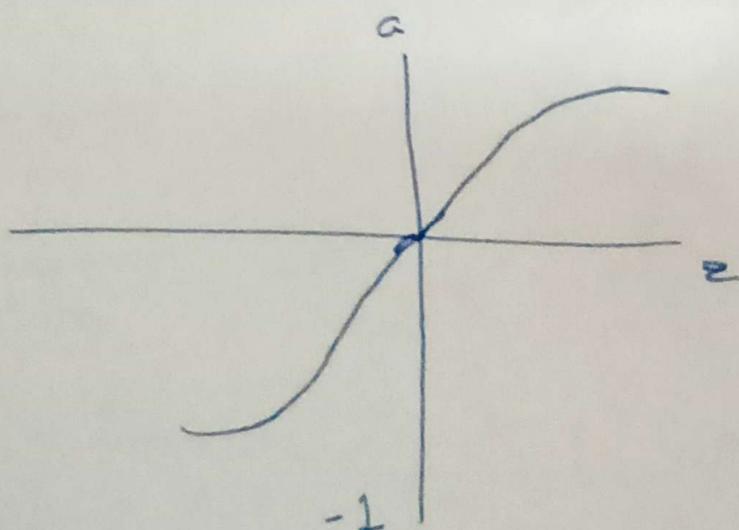
$$z = -10, g(z) \approx 0$$

$$\frac{d}{dz}g(z) = 0(1-0) = 0$$

$$z=0, g(z) = \frac{1}{2}$$

$$\frac{d}{dz}g(z) = \frac{1}{2}(1-\frac{1}{2}) = \frac{1}{4}$$

derivative of $\tan(h)$ activation function:-



$$g'(z) = \frac{d}{dz}g(z) = \text{slope of } g(z) \text{ at } z$$

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

ReLU and leaky ReLU

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

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

ReLU

Leaky ReLU

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

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

Topic:- Gradient descent for neural network

Parameters : $w^{[1]}, b^{[1]}, w^{[2]}, b^{[2]}$

Cost function = $\frac{1}{m} \sum \ell(\hat{y}, y)$ (Loss function)
 $\uparrow a^{[2]}$

Repeat $\left\{ (\hat{r}^{(i)}, i: \dots, n)$

$$\frac{\partial \hat{r}}{\partial w} = w^{[i]} - \alpha dw^{[i]}$$

$$\frac{\partial \hat{r}}{\partial b} = b^{[i]} - \alpha db^{[i]}$$

formula for propagation in Python

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

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

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

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

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

$$d z^{[1]} = w^{[2]T} d z^{[2]} x^T$$

Backpropagation (Topic) (goes backward)
to find w, b from $L(a, y)$

$$w \xrightarrow{\longrightarrow} z = w^T x + b \rightarrow a = \sigma(z) \rightarrow L(a, y)$$

lose

chain rule

$$\# \frac{d L(a, y)}{da} = -\frac{y}{a} - \frac{(1-y)}{1-a}$$

$$dz = da \cdot g'(z)$$

$$= -\frac{y}{a} + \frac{1-y}{1-a}$$

$$\frac{d L}{dz} = \frac{d h}{da} \cdot \frac{d a}{d z}$$

$$"dz" = "da" \hookrightarrow d$$

Q# Back propagation

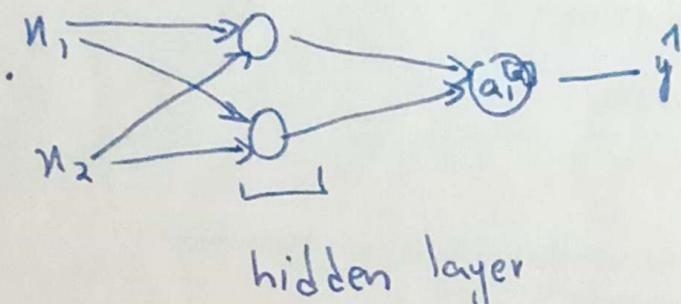
$$L(a, y)$$

$$dz^{[2]} = a^{[2]} - y$$

$$dw = dz^{[2]} a^{[1]T}$$

$$db^{[2]} = dz^{[2]}$$

One hidden layer



$$z = w \cdot x + b$$



where w is randomly initialized

$$w = np.random.rand(n_outputs, n_inputs) \times 0.01$$

Deep neural network

* more than 3 hidden layers

Shallow neural layer

* 2 hidden layer

Top) Forward propagation

$$\text{first layer } z^{[1]} = w^{[1]}x + b^{[1]}$$

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

$$\text{2nd layer.} = z^{[2]} = w^{[2]} \cdot h^{[1]} + b^{[2]}$$

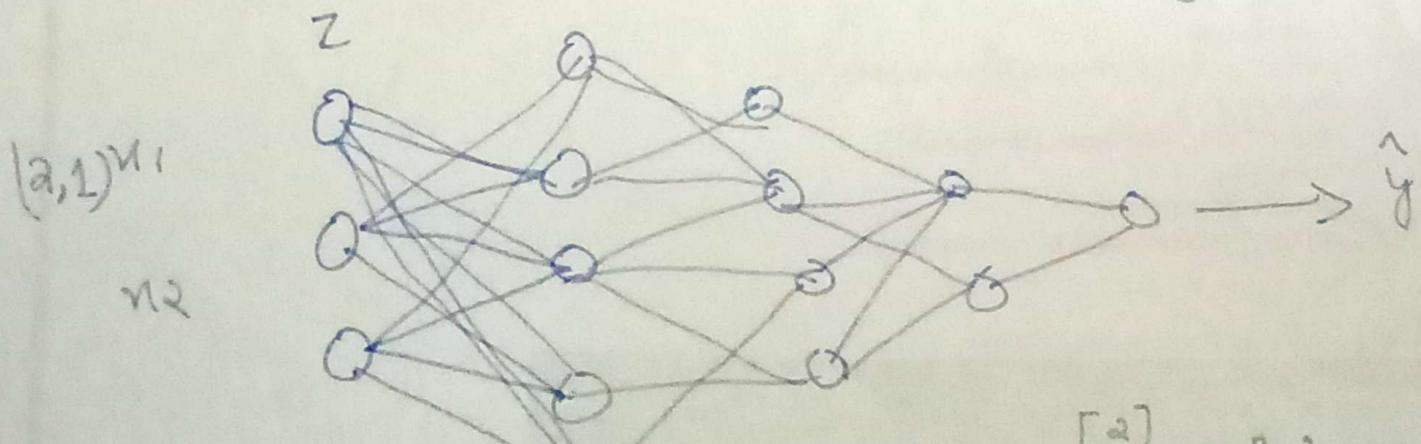
$$z^{[4]} = w^{[4]} x^{[3]} + b^{[4]}$$

$$a^{\left[\cdot \right]} \rightarrow S\left[z^{\left[\cdot \right]} \right]$$

$$z^l = \omega^{l_n} \bar{t}^{-1} + b^{[l]}$$

$$a^{[\ell]} = g^{[\ell]}(z^{[\ell]})$$

Top) Dimension Getting or dimension right.



$$z^{[2]} = w^{[1]} x + b^{[2]}$$

$$(3, 1) = \begin{pmatrix} 3, 2 \\ 2, 1 \end{pmatrix} \begin{pmatrix} n, 2 \end{pmatrix}$$

$$Z^{[2]} = W^{[2]} \overset{a^{[1]}}{\underset{[5,3]}{\uparrow}} + b^{[3]}$$

$$w^{[l]} = [n^{\text{in}}, n^{[l+1]}]$$

$$b^{[l]} = [n^{[l]}, 1]$$

Deep Neural Networks Representation

Forward propagation:-

Input $a^{[L-1]}$

output $a^{[l]}$, cache ($z^{[l]}$)

use to store ~~result~~ of intermediate computation
one layer in back propagation

Back propagation

Input $da^{[l]}$

output $da^{[l-1]}, dw^{[l]}, db^{[l]}$

$$\rightarrow dz^{[l]} = da^{[l]} * g^{[l]}(z)$$

$$\rightarrow dw^{[l]} = dz^{[l]} \cdot a^{[l-1]T}$$

$$\rightarrow db^{[l]} = dz^{[l]}$$

$$\rightarrow da^{[l-1]} = w^{[l]T} \cdot dz^{[l]}$$

$$\rightarrow dz^{[l]} = w^{[l+1]T} \cdot dz^{[l+1]} * g^{[l]}(z^{[l]})$$

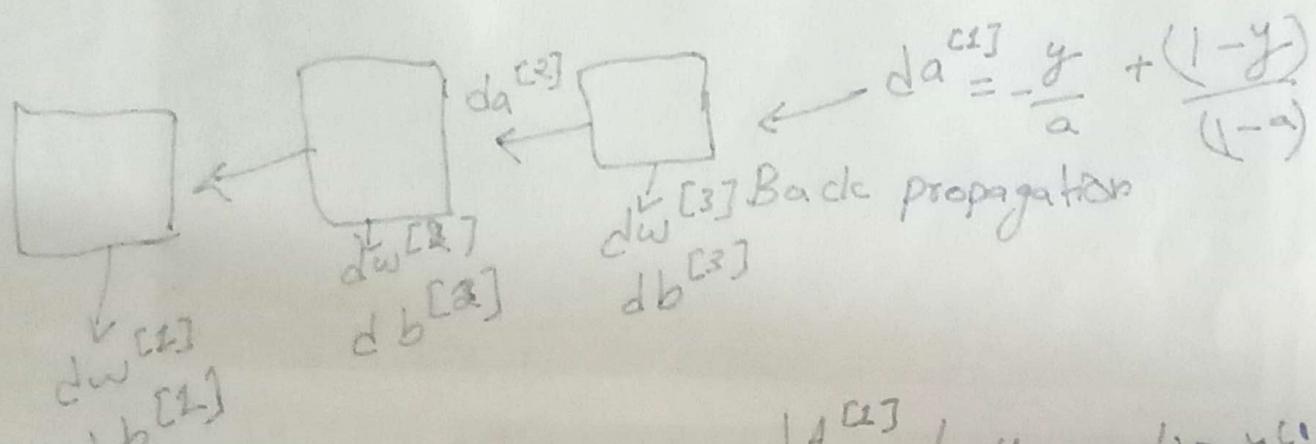
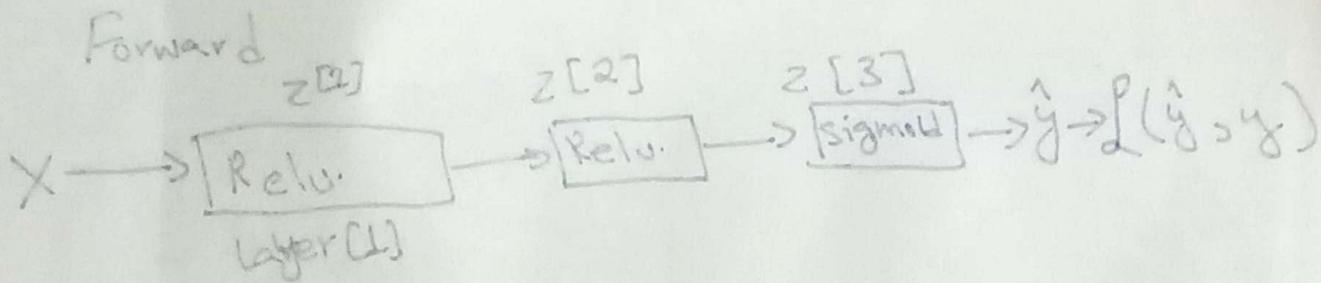
$$dz^{[l]} = dA^{[l-1]} \star g^{[l]}(z^{[l]})$$

$$dw^{[l]} = \frac{1}{m} dz^{[l]} \cdot A^{[l-1]T}$$

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

$$dA^{[l+1]} = w^{[l]T} \cdot dz^{[l]}$$

Summary:



$$dA^{[l]} = \left(\frac{-y}{a^{[l]}} + \frac{(1-y^{[l]})}{1-a^{[l]}} \right)$$

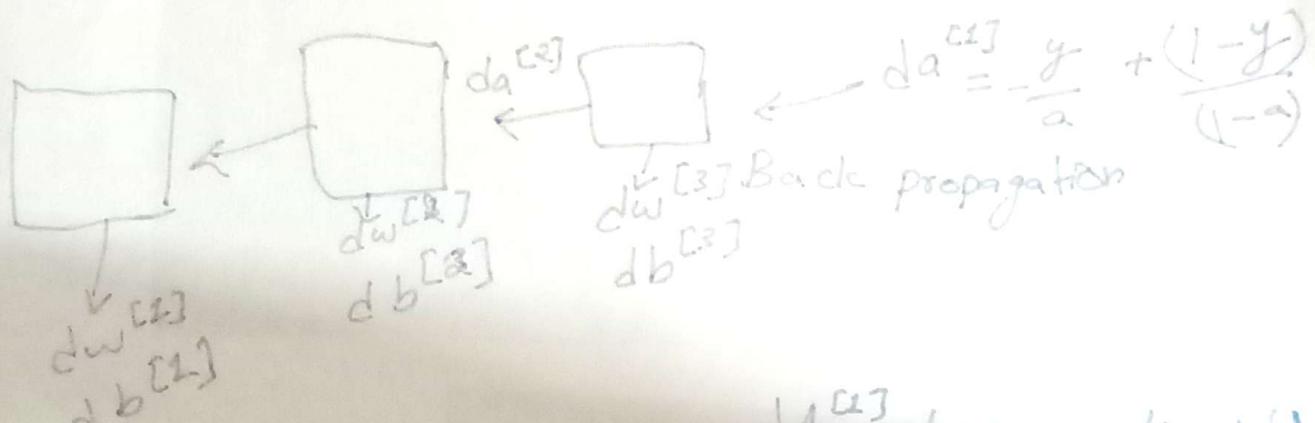
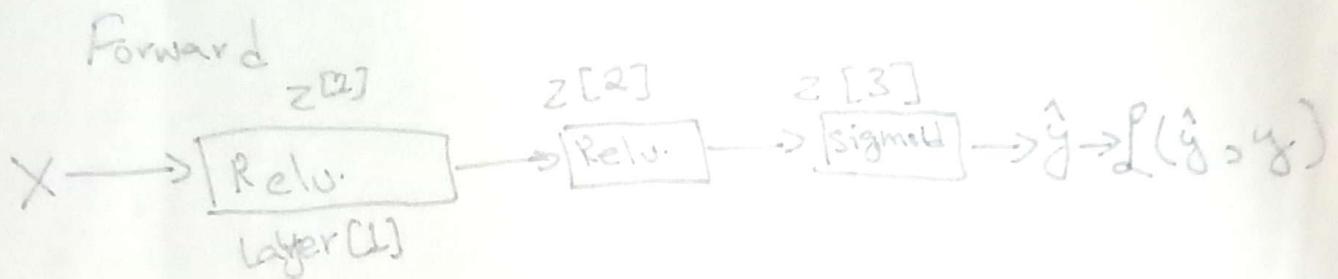
$$dz^{[l]} = dA^{[l-1]} \star g^{[l-1]}(z^{[l]})$$

$$dw^{[l]} = \frac{1}{m} dz^{[l]} \cdot A^{[l-1]T}$$

$$db^{[l]} = \frac{1}{m} \text{np.sum}(dz^{[l]}, \text{axis}=1, \text{keepdims}=1)$$

$$dA^{[l+1]} = w^{[l]T} \cdot dz^{[l]}$$

Summary:



$$da^{[l]} = -\frac{y}{a} + \frac{(1-y)}{(1-a)}$$

$$dA^{[l]} = \left(\frac{-y}{a^{[l]}} + \frac{(1-y)}{1-a^{[l]}} \right)$$