# WEEK 3

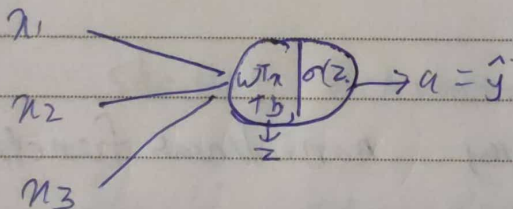$$X^{[i]} \rightarrow \text{For } X \text{ in } i^{th} \text{ layer.} \qquad [X]^{(i)} \rightarrow X \text{ in } i^{th} \text{ training example}$$

Neural network is basically multiple Logistic layers one after other.

# # Representation of NN

— repetitive, refer ML notes

# # Computing a neural networks output.



every neuron performs 2 computations $\rightarrow z, a$.

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

# # Vectorizing across multiple training examples.

//careful

usually for loop on for $i=1$ to $i=m$ : ↓
m training data

$$z^{[1](i)} = w^{[1]} x^{(i)} + b^{[1]}.$$
$$a^{[1](i)} = \sigma(z^{[1](i)})$$
$$z^{[2](i)} = w^{[2]} a^{[1](i)} + b^{[2]}$$
$$a^{[2](i)} = \sigma(z^{[2](i)})$$

horizontally $\rightarrow$ different training examples

column wise $\rightarrow$ vertically $\rightarrow$ one neuron/feature of one training example $\underset{\phantom{a}}{\underline{\underline{\phantom{aaa}}}}$

we combine all training example.

$$X = \begin{bmatrix} \vdots & \vdots & & \vdots \\ x^1 & x^2 & \cdots & x^m \\ \vdots & \vdots & & \vdots \end{bmatrix} \longrightarrow Z \text{ becomes } \begin{bmatrix} \vdots & \vdots & & \vdots \\ z^1 & z^2 & \cdots & z^m \\ \vdots & \vdots & & \vdots \end{bmatrix}$$

n×m Matrix

W remains same across all training examples

a becomes

$$a = \begin{bmatrix} \vdots & \vdots & & \vdots \\ a^1 & a^2 & \cdots & a^m \\ \vdots & \vdots & & \vdots \end{bmatrix}$$

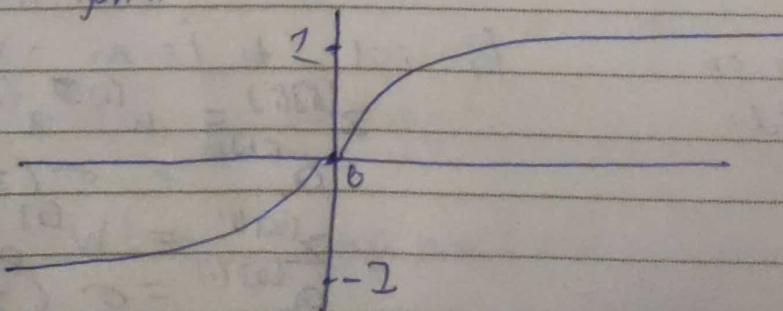// same code, just remove for loop !!

# Activation functions
↳ our choice.
↳ can be any non-linear function.

eg:- Sigmoid :- $\sigma(z) = \dfrac{1}{1 + e^{-z}}$

b/w 0 and 1
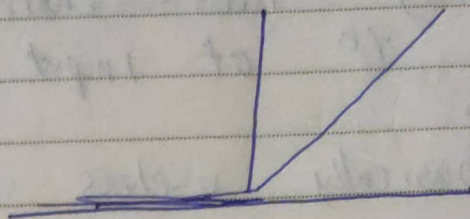↓
useful in binary classification



tanh



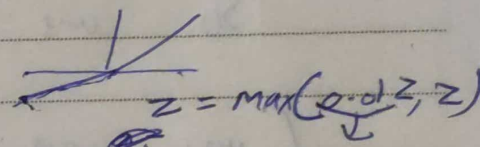$$\tanh(z) = \dfrac{e^z - e^{-z}}{e^z + e^{-z}}$$

Relu→ReLU $z = max(0, z)$



tanh is superior

- use sigmoid only if you're binary classificat<sup>n</sup>
- use ReLU mostly
- ReLU is becoming very common
- disadvantage of ReLU - derivative is 0 for -ve no
- ↓ leaky ReLU: cons it ⟶ $z = max(0.01z, z)$

can change this constant

Why should Activation f<sup>n</sup> be Non-Linear ?

Why should we even use Activation f<sup>n</sup> ?

lets not use $\xrightarrow[to]{equivalent}$ using a linear activation fn. (with slope 1)

Our algo becomes:

Given $x$:
$$z^{[1]} = W^{[1]} x + b^{[1]}$$
$$a^{[1]} = z^{[1]}$$
$$z^{[2]} = W^{[2]} a^{[1]} + b^{[2]}$$
$$a^{[2]} = z^{[2]}$$

then.
$$a^{[2]} = W^{[2]} (W^{[1]} x + b^{[1]}) + b^{[2]}$$
$$a^{[2]} \to \text{our output is a linear f}^n \text{ of } x !!$$

Most problems can't be solved by just a
linear fn

So basically our $\overbrace{\text{entire fn}}^{\text{entire mode}}$ is just
a linear fn of input
↓
its basically useless
If you put multiple layers but they're
all linear, the final output will just
be still a linear fn of input

So all the middle layers were useless

So we need to induce non-linearity

[ you can use Linear fn for Regression
problems btw ] → only in final output layer.


Derivatives of Activation function

$$\frac{d\sigma(z)}{dz} = \sigma(z)[1-\sigma(z)]$$


for large $(z) \longrightarrow \sigma(z) \approx 1$
$$\frac{d\sigma(z)}{dz} \approx 1(1-1) \approx 0$$


for small $(z)$  ⊕ $\sigma(z) \approx 0$
$$\frac{d\sigma(z)}{dz} \approx 0(1-0) \approx 0$$


for $z=0 \rightarrow \sigma(z) = \frac{1}{2}$
$$\frac{d}{dz}\sigma(z) = \frac{1}{4}$$

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

for $z = 10$      $\tanh(z) \approx 1$
                derivative $\approx 0$

for $z = -10$      $\tanh(z) \approx -10$
                derivative $\approx 0$

for $z = 0$      $\tanh(z) \approx 0$
                derivative $\approx 1$

$$\frac{d}{dz} \text{ReLU}(z) = \begin{cases} 0 & \text{if } z < 0 \\ 1 & \text{if } z > 0 \\ \text{undefined} & \text{if } z = 0 \end{cases}$$

$$\frac{d}{dz} Q(\text{Leaky ReLU}(z)) = \begin{cases} 0.01 & z < 0 \\ 1 & z > 0 \\ \text{undefine} & z \leq 0 \end{cases}$$

[undefined for $z = 0$ doesn't matter b/z
derivative almost never becomes exactly zero)

Gradient Descent for Neural Networks

example taken is of 1 layer
so basically same as logistic regression

# Forward + backward propogation for layes l

Input: $a^{[l-1]}$
Output: $a^{[l]}$, cach $z^{[l]}$

forward prop
$$z = w \cdot a^T + b$$
$$a = sigmoid(z)$$

using vectoriz$^n$
$$Z = WA + b$$
$$A^{[l]} = g^{[l]}(Z^{[l]})$$

Backward part

Input: $da^{[l]}$

Output: $da^{[l-1]}, dW^{[l]}, db^{[l]}$

$$dz^{[l]} = da^{[l]} \times g'^{[l]}(z^{[l]})$$
$$dW^{[l]} = dz^{[l]} \cdot a$$
$$db^{[l]} = dz^{[l]}$$
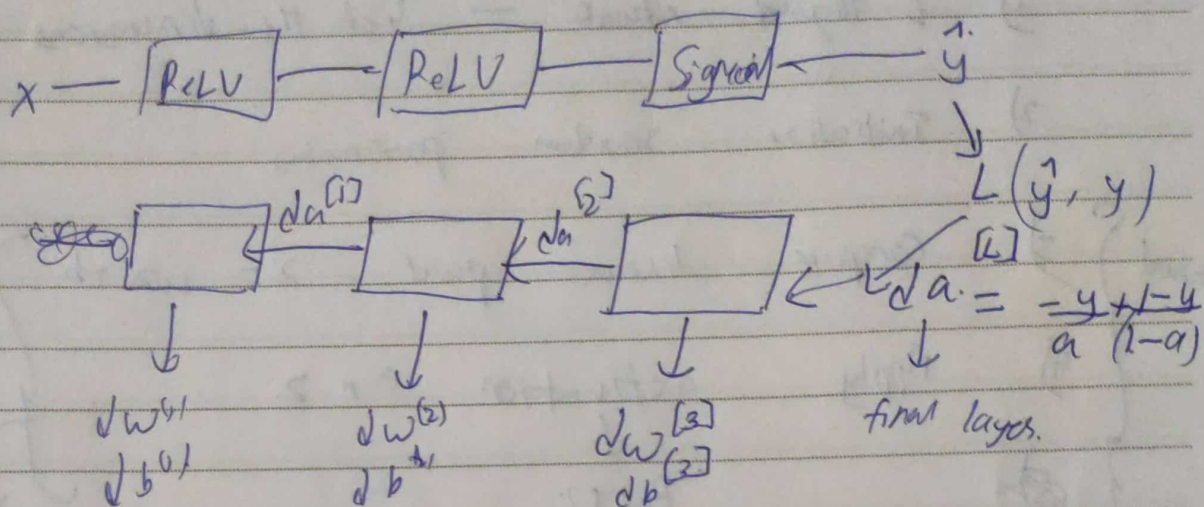$$da^{[l-1]} = W^{[l]T} \cdot dz^{[l]}$$

basically
d(param)
$= g'(z) d(param\ output)$

comes from partial derivatives

$$dz^{[l]} = W^{[l+1]T} * g'^{[l]}(z^{[l]})$$

$$dz^{[l]} = d$$

# Same for vectorize v.

Summarising :-

$$x \longrightarrow \boxed{ReLU} \longrightarrow \boxed{ReLU} \longrightarrow \boxed{Sigmoid} \longrightarrow \hat{y}$$

$$\downarrow$$
$$L(\hat{y}, y)$$

$$\boxed{\phantom{xx}} \xleftarrow{da^{[1]}} \boxed{\phantom{xx}} \xleftarrow{da^{[2]}} \boxed{\phantom{xx}} \xleftarrow{} \sqrt{da^{[L]}} \quad da = \frac{-y + 1 - y}{a(1-a)}$$

$$\downarrow \qquad\qquad \downarrow \qquad\qquad \downarrow \qquad\qquad \downarrow$$

$$\begin{array}{cccc} dw^{[1]} & dw^{[2]} & dw^{[3]} & \text{final layer.} \\ db^{[1]} & db^{[2]} & db^{[3]} \end{array}$$

The eqn't of backward prop can be generalised — for multiple layers.
So learn them

# Hyperparameters    vs    parameters
$$\qquad\qquad \lower{L}\qquad\qquad\qquad\qquad\qquad \lower{L} \to w, b$$

      $\to$ eg :- $\Big\{$ learning rate $\alpha$
            no. of iterations
            no. of hidden layers $L$
            no. of hidden units
            choice of activation fn$^{\text{ctn}}$

they control the actual parameters $w$ and $b$

other hyper parameters :-  momentum.
                         mini batch size.
                         regularizations --

1) L layered network — Get the dimensions right.

2) Initialize random parameters

forward {
3) Compute linear part $z = wx + b$

4) Apply activation on z.
↓
ReLU.
( Sigmoid is lost

repeat L-1 times

5) compute loss
↳ you have the function.

6) compute linear part of backward prop

7) get combine linear and activation parts
descent

8) Stack backward prop combinations.