

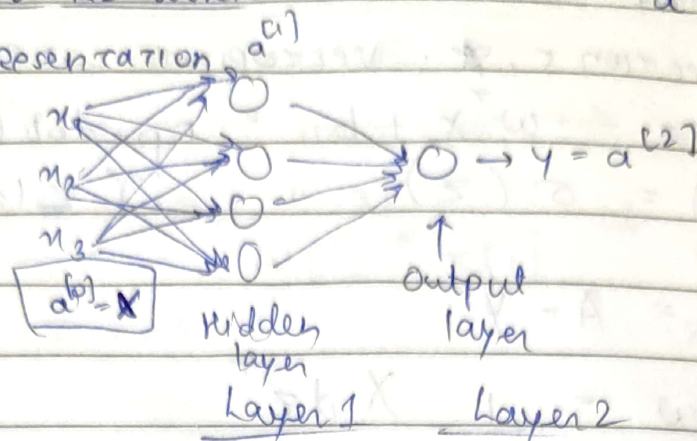
19/7/24

Page No.

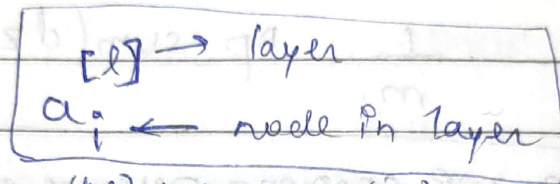
Date

II) NEURAL NETWORKS

1. Representation



Activation = $\sigma(z)$ $z = w^T x + b$



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

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

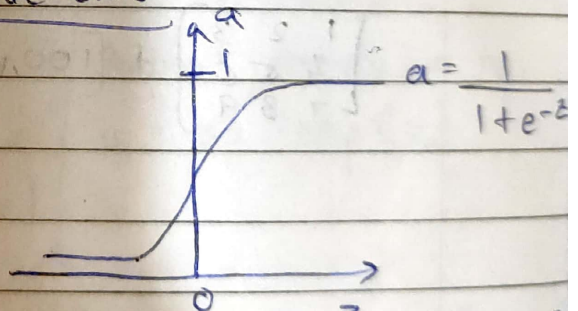
$$a^{[1]} = \sigma(z^{[1]})$$

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

2. ACTIVATION FUNCTIONS

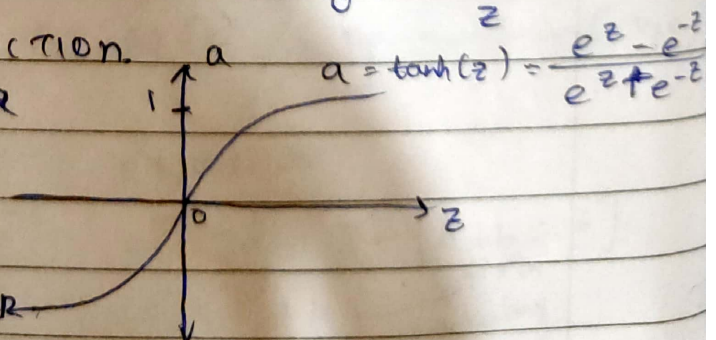
a. SIGMOID FUNCTION

[ONLY FOR BINARY CLASSIFICATION]



b. tanh function. (always better than sigmoid)

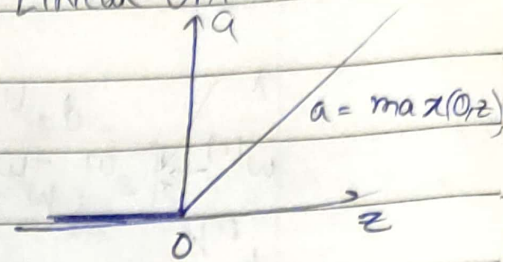
[ONLY PROBLEMATIC WHEN OUTPUT LAYER SHOULD NOT GIVE -ve values]



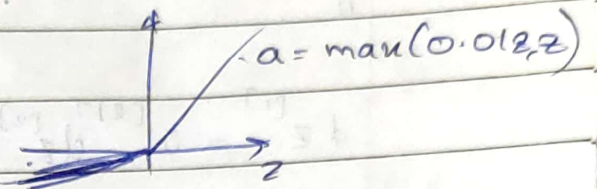
C. ReLU Function : Rectified Linear Unit

[Almost always used]

BY DEFAULT



⇒ Leaky ReLU



⊛ we DO NOT use linear activation functions
WHY? IT IS USELESS.

IF THE HIDDEN LAYERS use a linear actvn fn,
THEN IT IS AS GOOD AS NOT USING ANY HIDDEN
LAYER. Since $g(z) = z$, linear actvn fn
YIELDS THE SAME OUTPUT EVERYTIME.

THE COMPOSITION OF 2 LINEAR Fns IS A LINEAR FN.

• IF needed (say a real number output b/w $-\infty$ to ∞)
ONLY THEN can we use a linear actvn BUT
ONLY IN THE OUTPUT LAYER.

20/7/24

3. DERIVATIVES OF ACTIVATION.

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

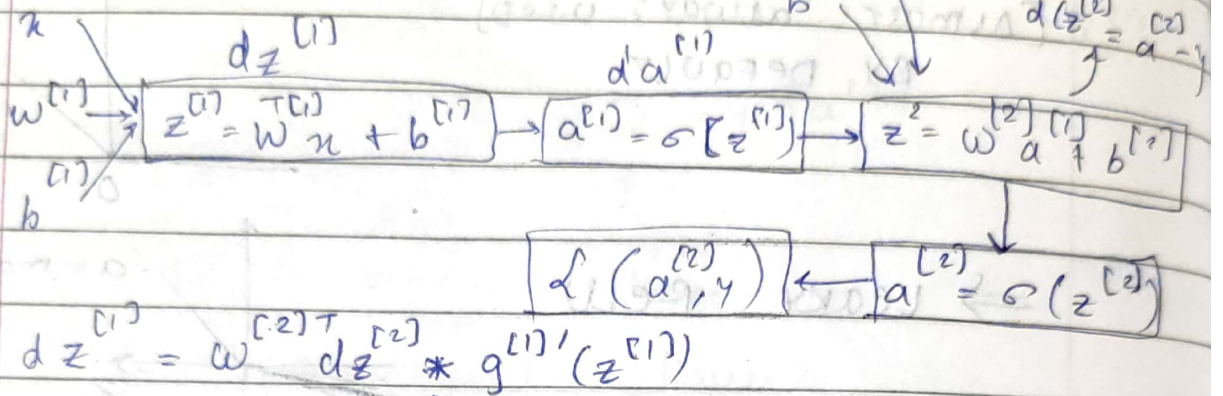
b. tanh $g(z) = \tanh z$
 $\therefore g'(z) = 1 - (\tanh z)^2$

c. ReLU $g(z) = \max(0, z)$

\therefore when $z = +ve$, $g'(z) = 1$

$z = -ve$ $g'(z) = 0.$

4. NEURAL NETWORK GRADIENTS



⇒ WHAT HAPPENS IF YOU INITIALIZE ALL WEIGHTS TO ZERO? $w^{(1)} = \text{np.random.randn}(2,2) * 0.01$

ALL THE HIDDEN UNITS ARE SYMMETRIC & EVEN IF WE CALCULATE GRADIENT DESCENT MANY TIMES, IT WILL STILL COMPUTE THE SAME FUNCTION.

BUT BIAS $b^{(1)} = \text{np.zeros}(2,1)$ CAN BE ZERO.