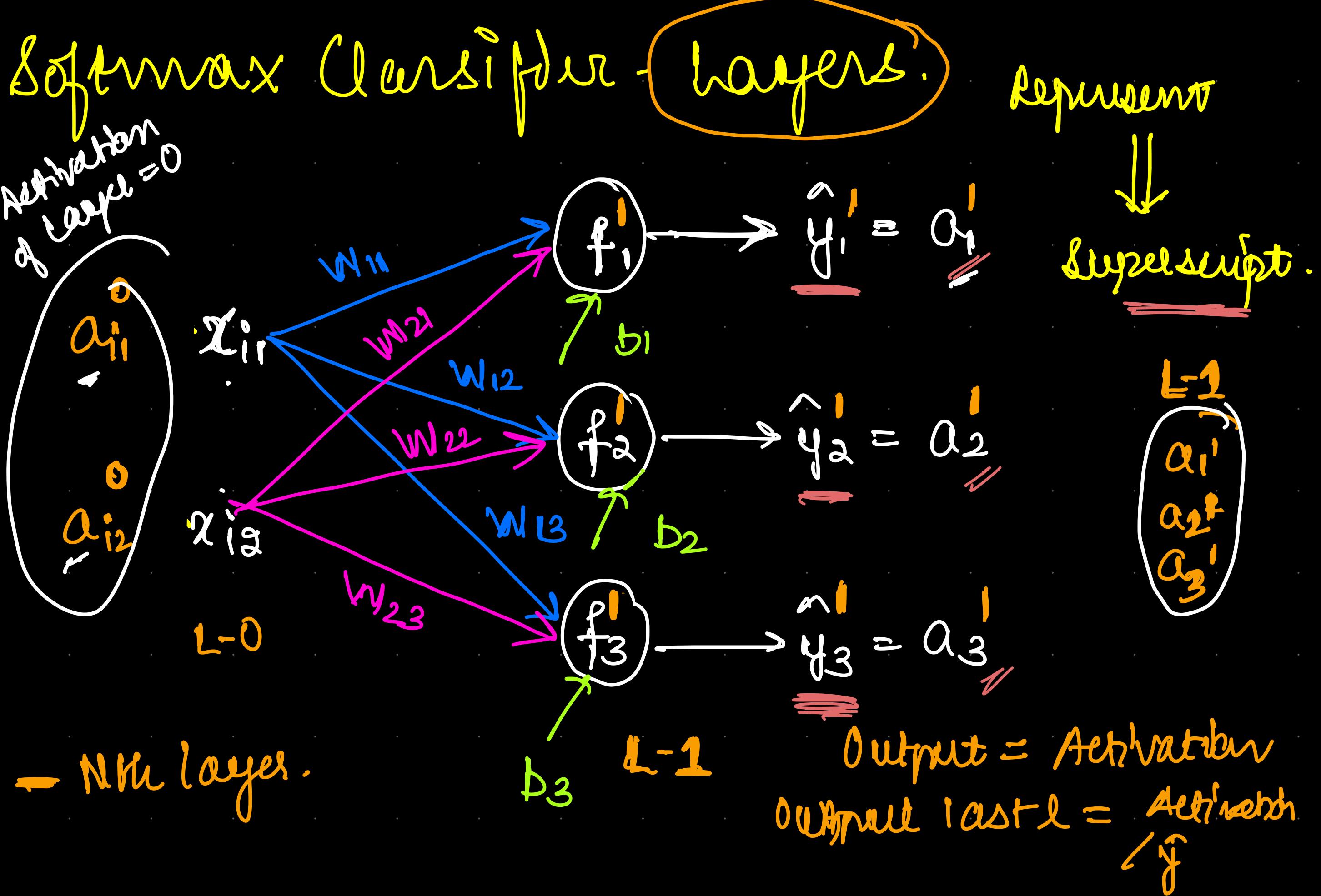


NN- week 4

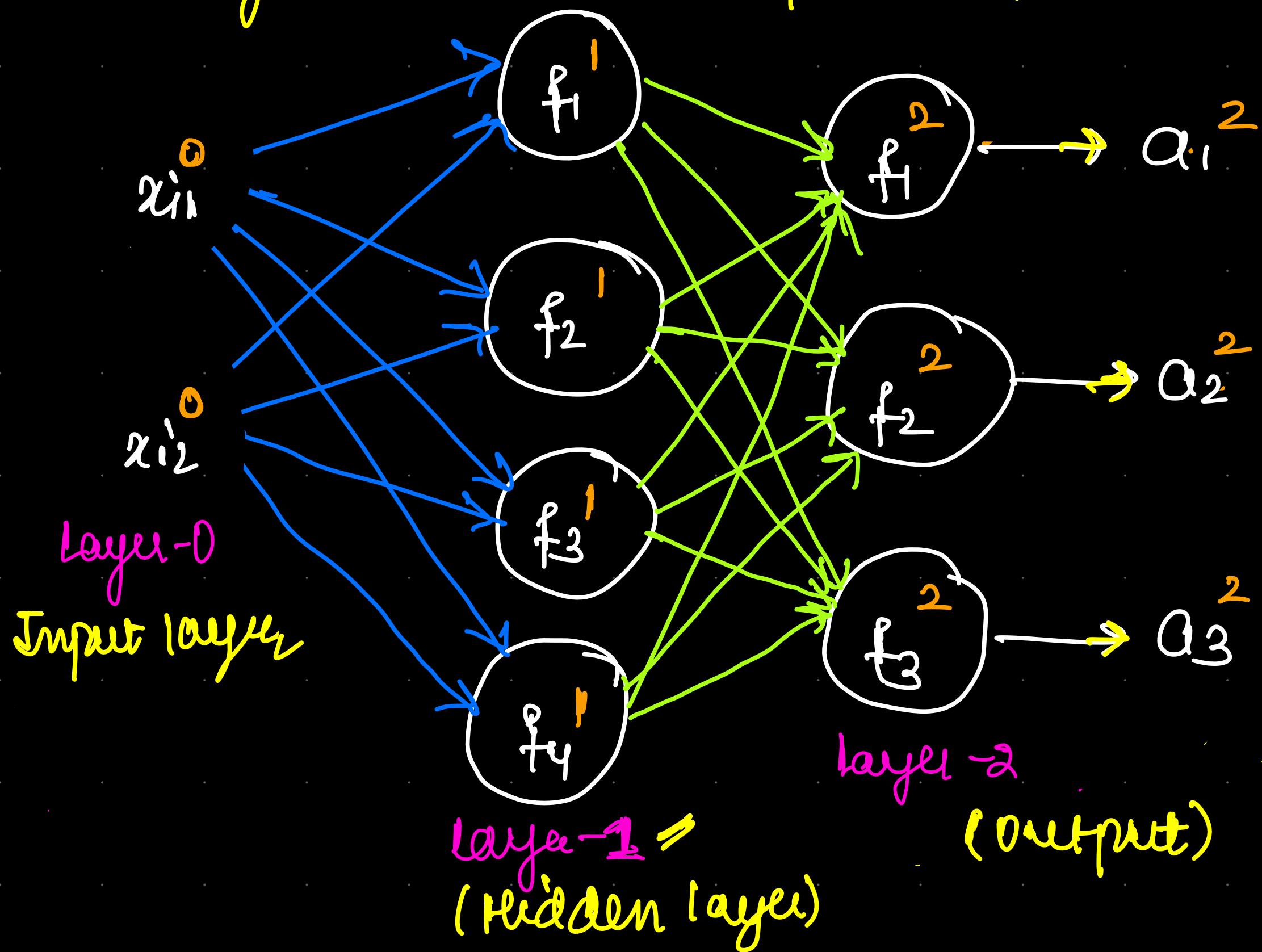
N-layer Neural Networks

- Non-Linear decision boundary
 - Activation Function
 - Forward Prop. N-layer.
- (Optional)
- Chain rule for M-layer
 - Backpropagation rule for N-layer.

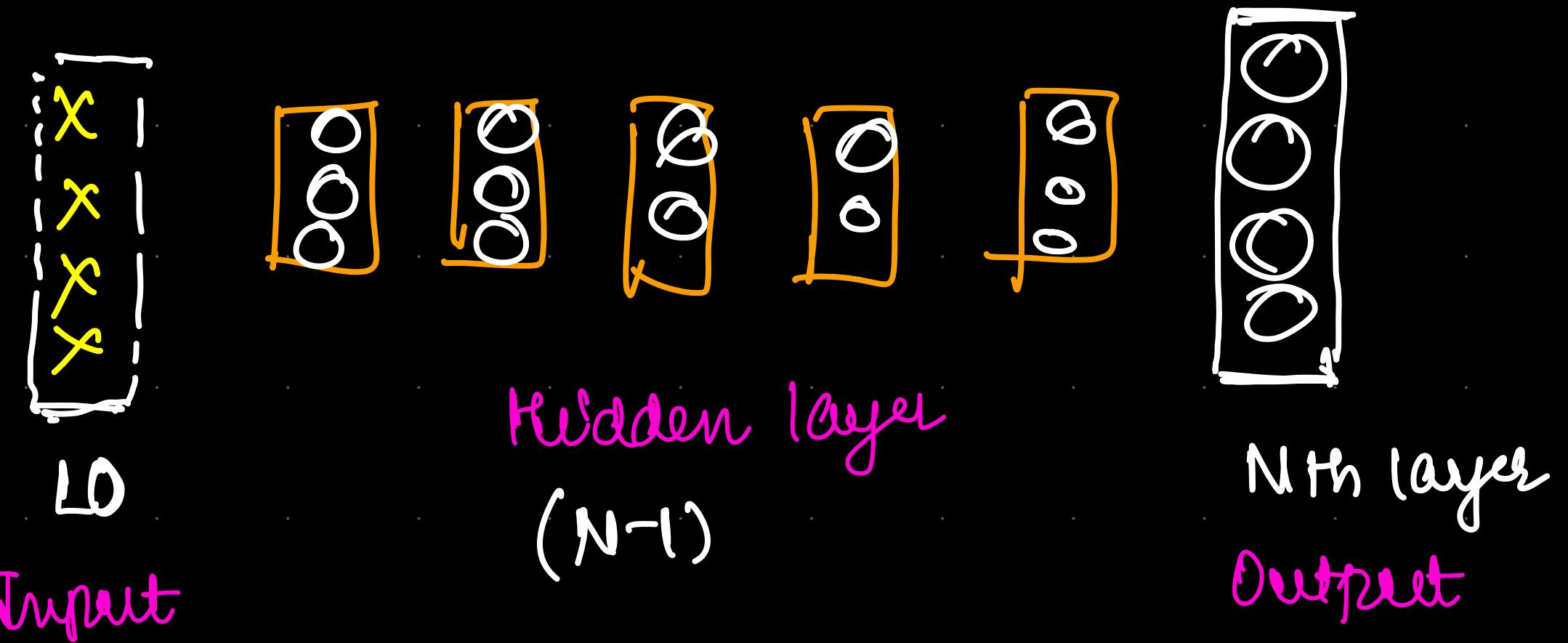


Mult-layer NN (Perceptron)

N+layer



N-layer NN.



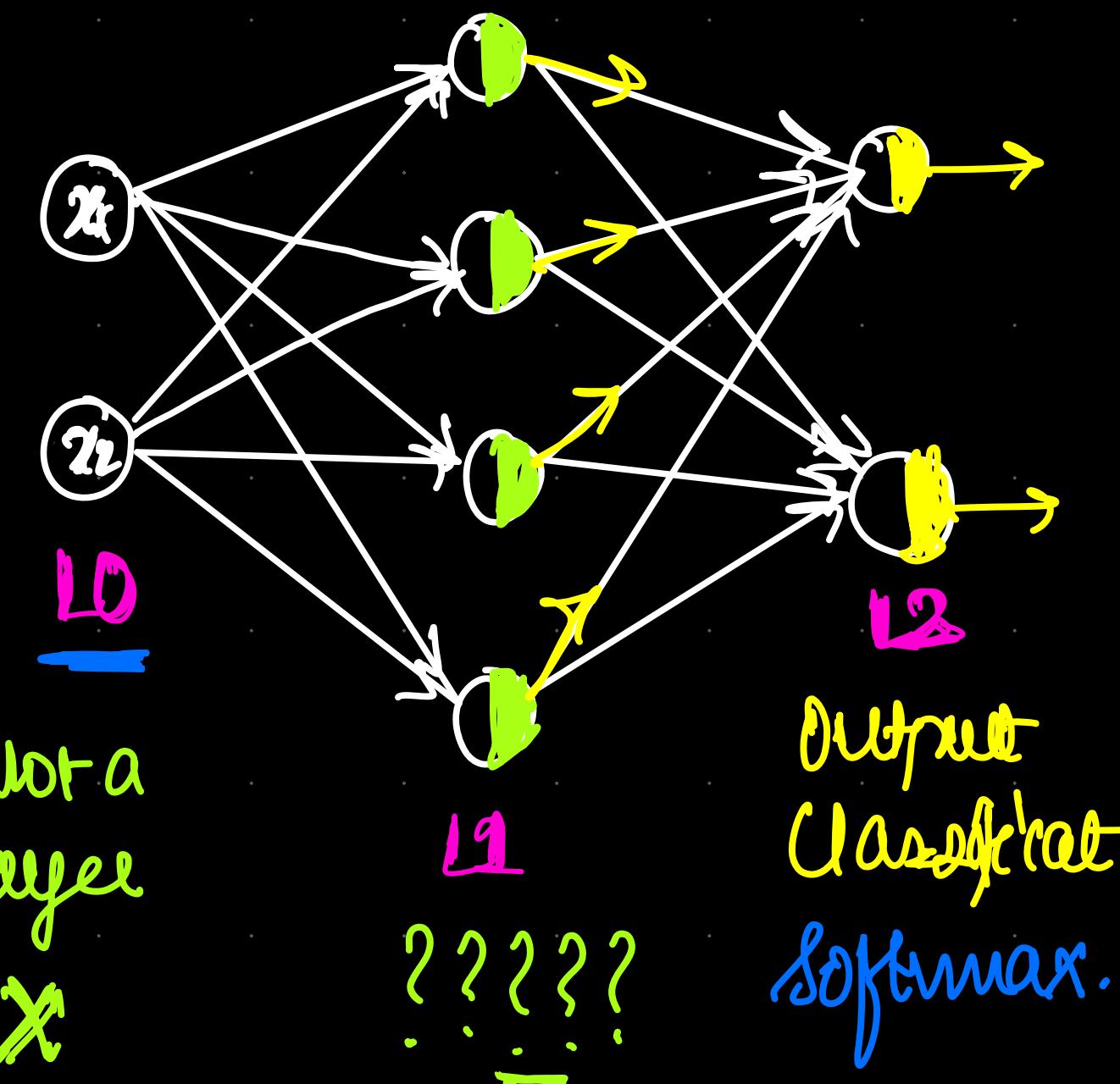
$$0 + (N-1) + 1 = N \text{-layer NN}$$

Hidden Output

What will be the activation of L-1?

Activation Function
Activation

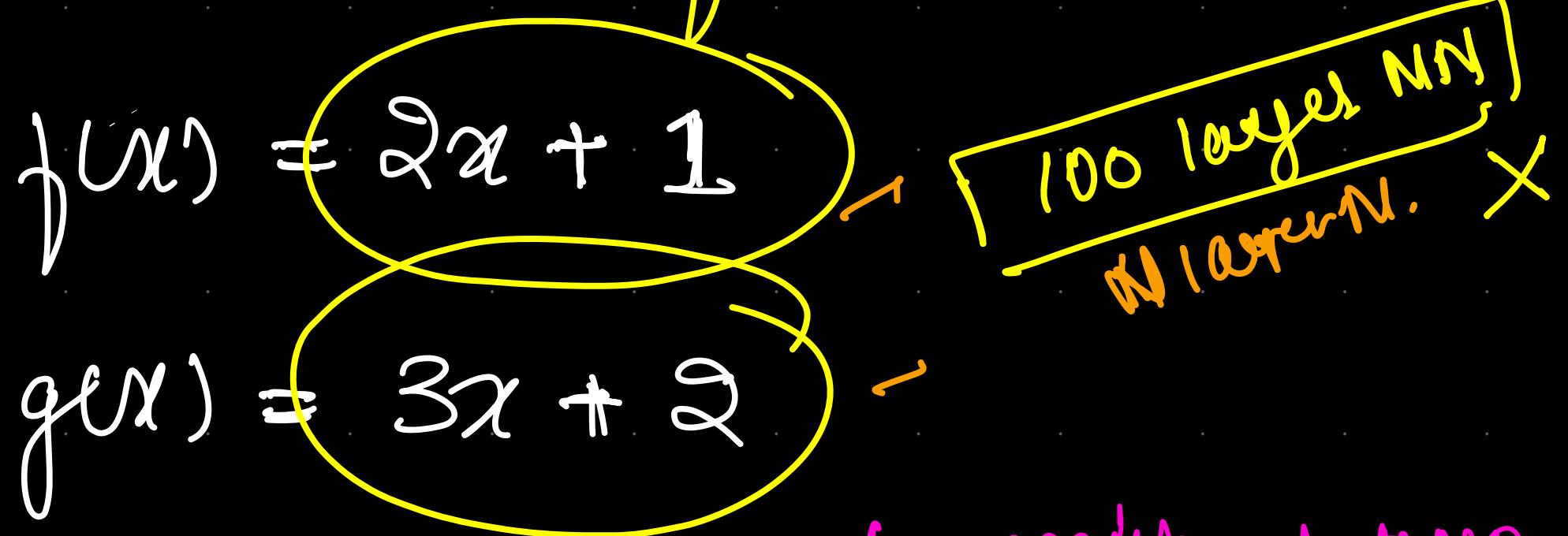
softmax



diff.
Non-
linear
function.

suggestion = step function. Sigmoid. ReLU, L.ReLU, tanh.

Why not a linear function?



$$\begin{aligned} g(f(x)) &= 3(2x+1) + 2 \\ &= \underline{6x+3+2} \\ &= \underline{\underline{6x+5}} \end{aligned}$$

Composition of two
linear
functions is a
linear function
but using 1-N.

Why a non-linear function?

$$f(x) = 2x + 1$$

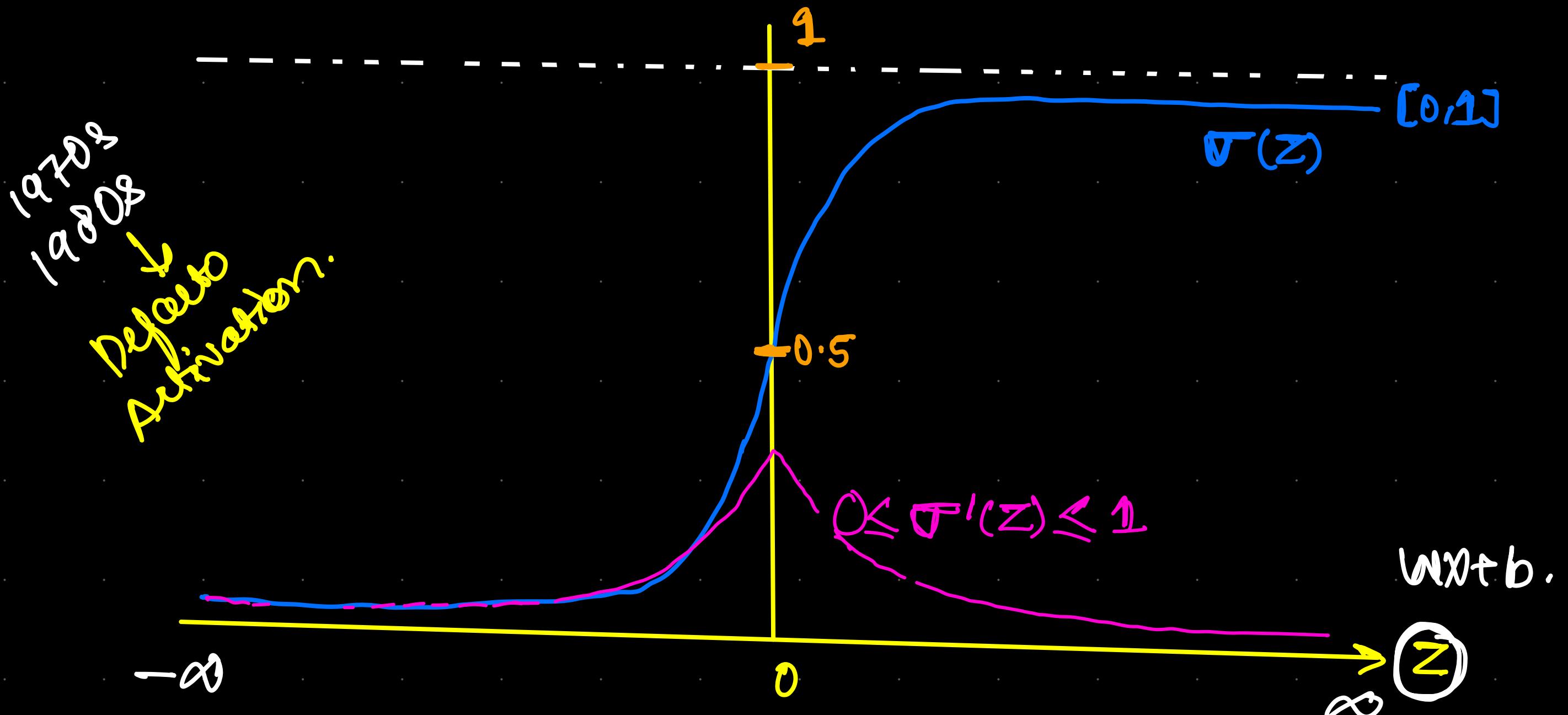
$$g(x) = x^2 + 1$$

$$\begin{aligned}gf(u) &= 2(2x+1)^2 + 1 \\&= 2(4x^2 + 4x + 1) + 1\end{aligned}$$

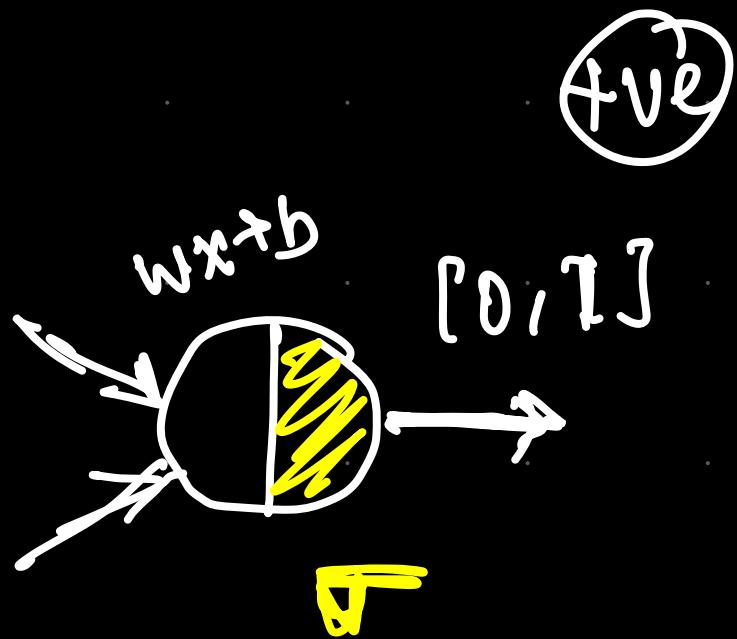
$$= 8x^2 + 8x + 2 + 1$$

$$= \boxed{8x^2 + 8x + 3} \rightarrow \text{Non-linear func.}$$

Activation function - Sigmoid ~~***~~
Add non-linearity



Sigmoid:



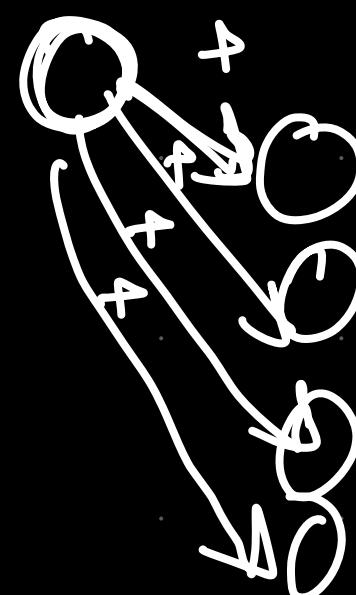
logistic Regression
Output of Binary classifier

$z \rightarrow p$
for converting $z \in \mathbb{R}$ into prop.

$$[-\infty, \infty] \rightarrow [0, 1]$$



"Non-Linear"
to the
system



Activation function - tanh

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

Why? similar behavior?

$(-\infty, \infty) \rightarrow [0, 1]$ sigmoid

$(-\infty, \infty) \rightarrow [-1, 1]$ tanh.

Centering the data

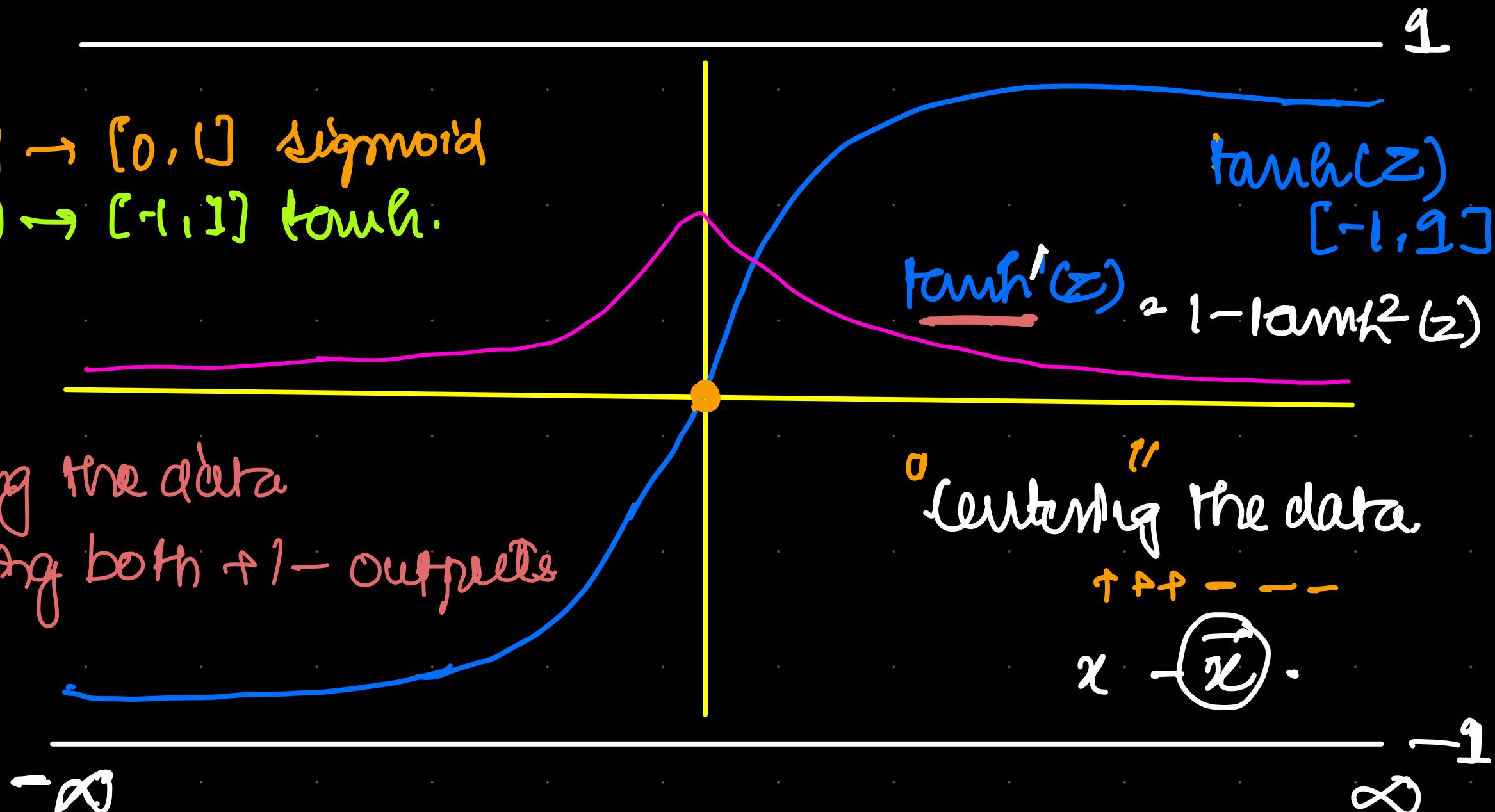
Generating both +1 - outputs

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

Centering " "

+ --- - ---

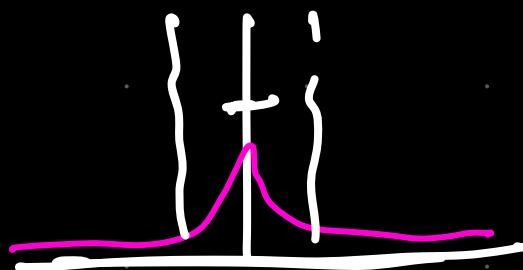
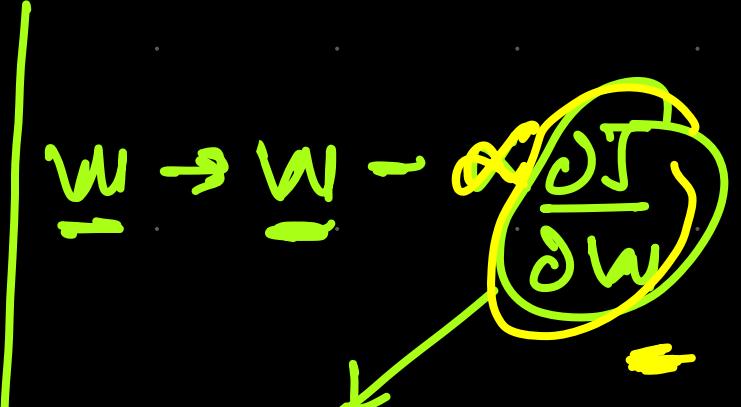
$x - \bar{x}$.



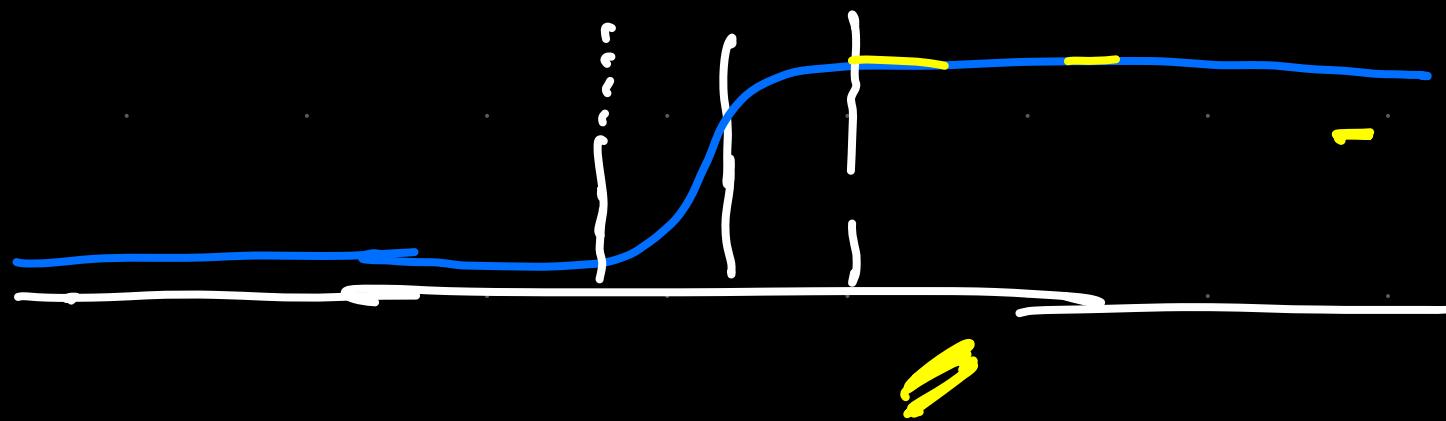
Issue with sigmoid / tanh

$$\frac{\partial J}{\partial w} = \frac{\partial J}{\partial p} \frac{\partial p}{\partial z} \frac{\partial z}{\partial w}$$

$[0,1] \quad [0,1] \quad [0,1]$



- All derivatives $0 < \partial < 1$



- Dynamic in a very small range $g \rightarrow 0$
for most values z

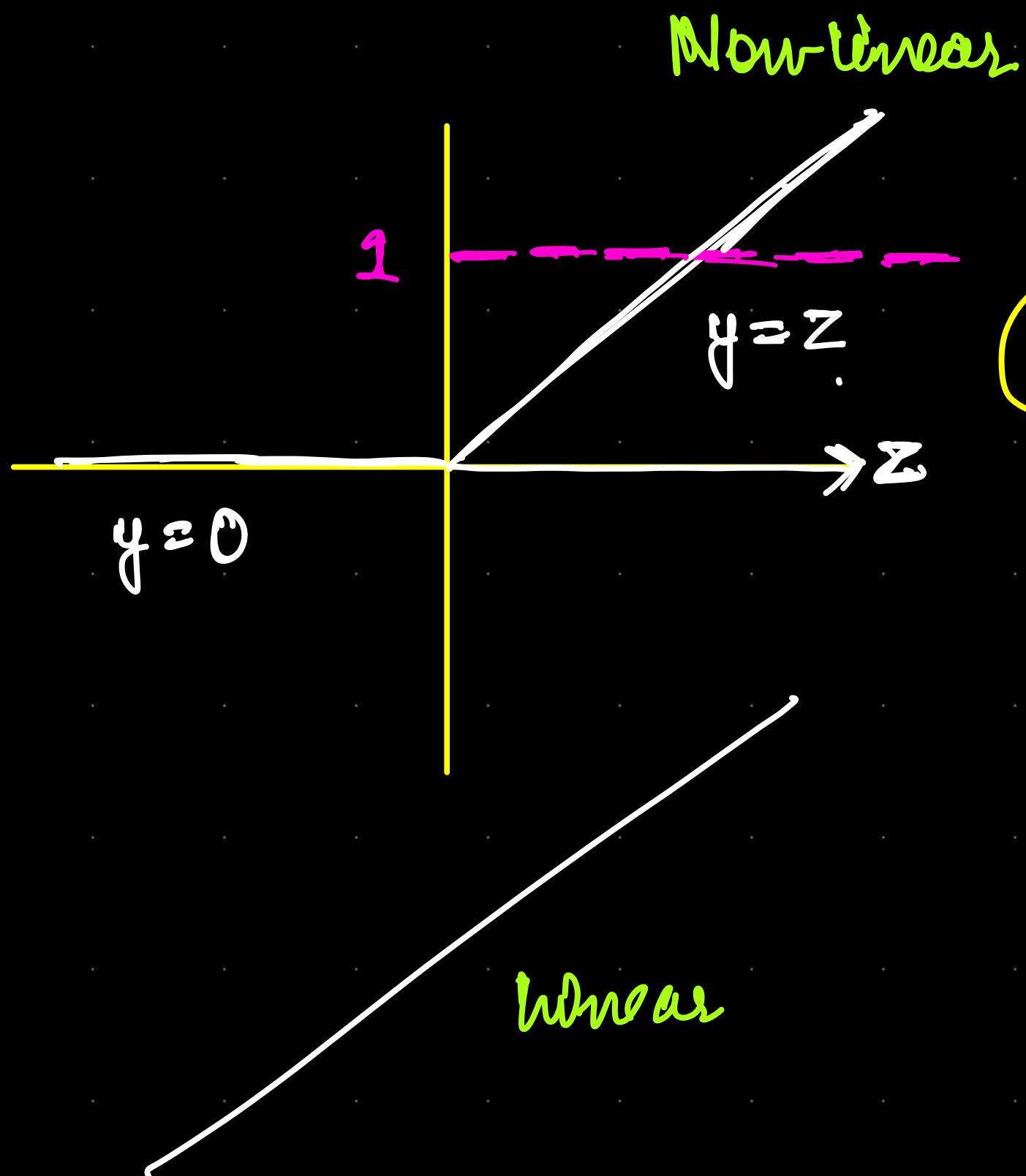
for most $z \rightarrow \partial$ will $\rightarrow 0$

$$0.0001 \times 0.0001 \sim 0$$

Ideal properties of activation function:

- Non-linear
- Differentiable
- Gradient > 1 for almost a big range Σ
- Easy to calculate. → Not strict requirement
- $[0, 1]$, $[-1, 1]$ → At least in that range

How shall we deal with this?



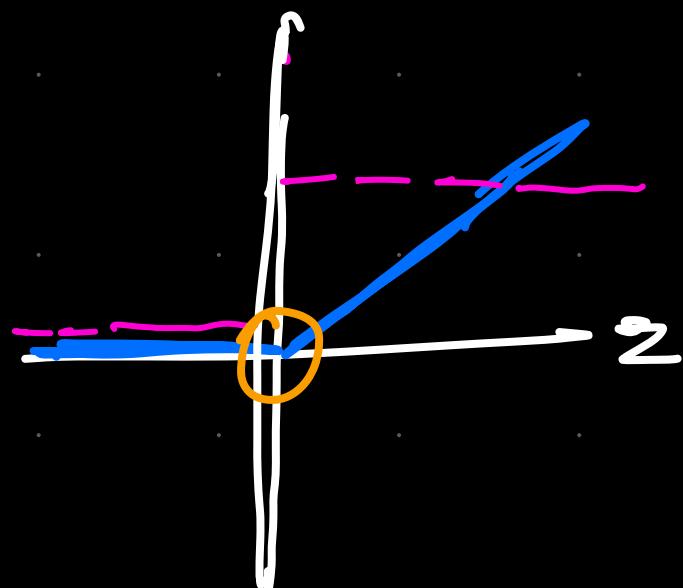
$$\text{ReLU} \begin{cases} 1 & z \geq 0 \\ 0 & z < 0 \end{cases}$$

Rectified Linear Unit

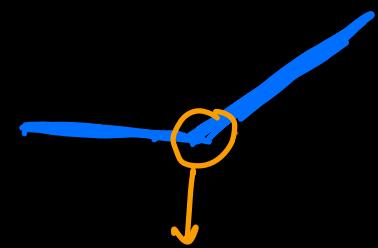
Simple and fast

$$\text{ReLU} \rightarrow \max(z, 0)$$

Problems with ReLU activation



→ Not differentiable at $z=0$



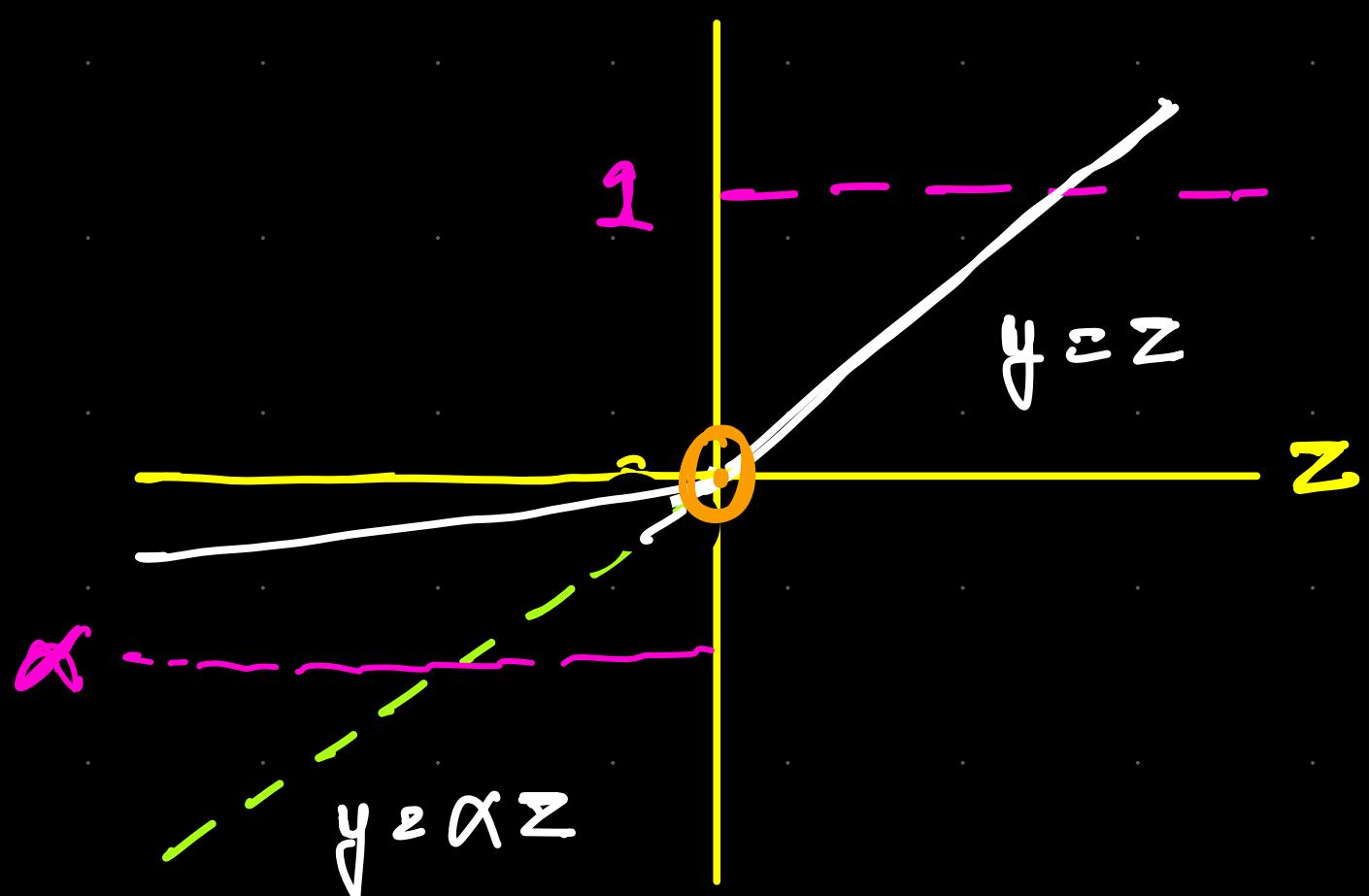
Not differentiable.

1 → still gradient is zero for half of the range.

$$0 \xrightarrow{x_1} 1 \times 1 \times 1 \times 1 = \underline{1}$$

$$1 \times 1 \times 1 \times 0 = \underline{0} \rightarrow \text{No update}$$

Another simple activation.



- - - linear.

leaky ReLU -

$$\begin{cases} z & \text{if } z \geq 0 \\ \alpha z & \text{if } z \leq 0 \end{cases}$$

Hyperparameter.

L'ReLU =

$$\begin{cases} 1 & \text{if } z \geq 0 \\ \alpha & \text{if } z \leq 0 \end{cases}$$

Not differentiable.

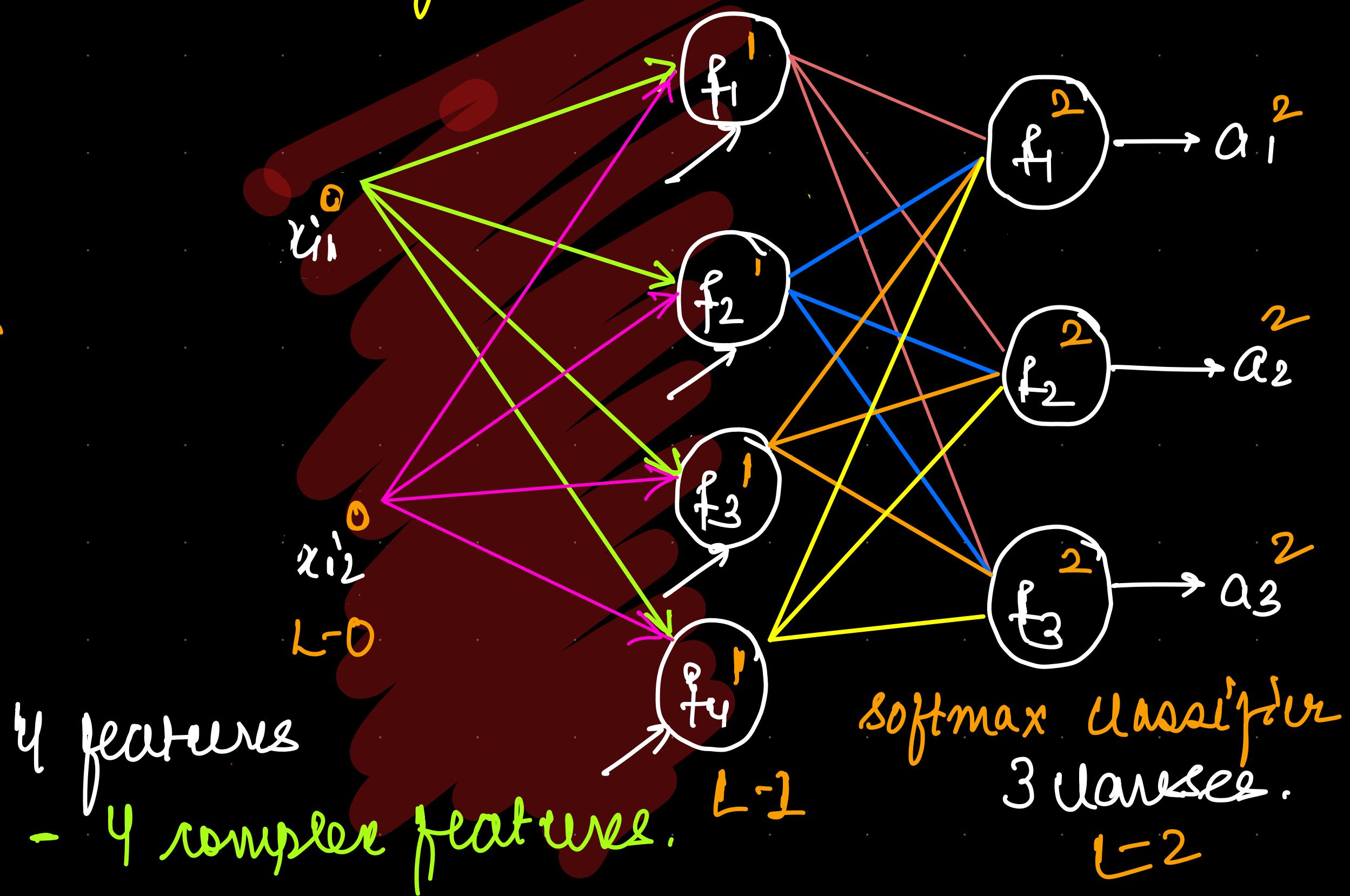
$$z \rightarrow z + \Delta$$

$\frac{z \rightarrow z + \Delta}{\Delta}$

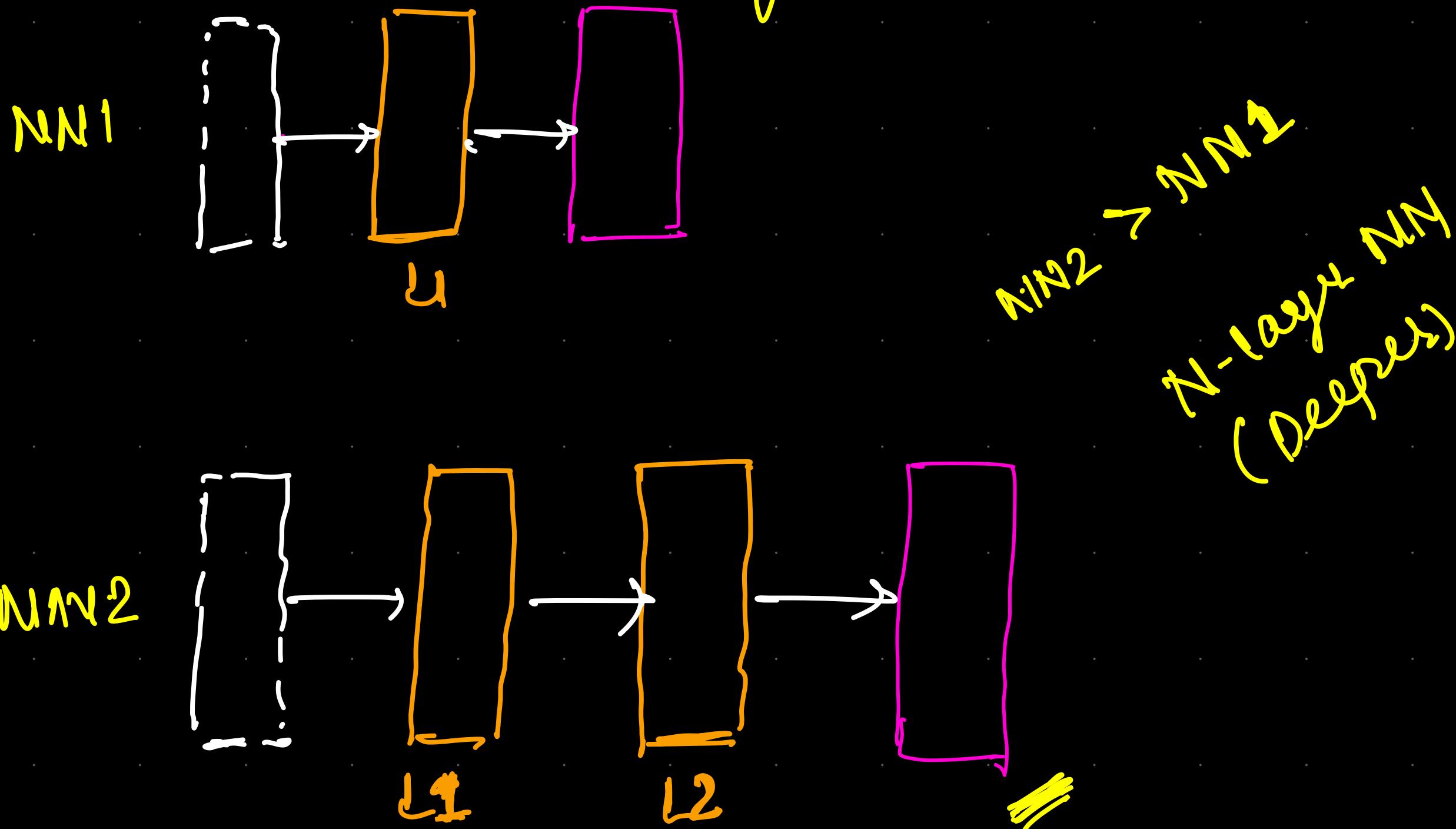
$z = 0 \neq \infty$

$z = 0.0000001$

Multi-layer NN (Notations)



creation of complex features.



Notations

w_{ij}^L

b_j^L

Neuron_i^{L-1} → Neuron_j^L

w_{ij}^L

Neuron_j^L

Notations

w^1

$$W = \begin{bmatrix} N_1^1 & N_2^1 & N_3^1 & N_4^1 \\ W_{11} & W_{12} & W_{13} & W_{14} \\ W_{21} & W_{22} & W_{23} & W_{24} \end{bmatrix}$$

2×4

inputs # neurons

b^1

$$B = \begin{bmatrix} N_1^1 & N_2^1 & N_3^1 & N_4^1 \\ b_1 & b_2 & b_3 & b_4 \end{bmatrix}$$

1×4

neurons.

Notations

W^2

$$\begin{matrix} & N_1^2 & N_2^2 & N_3^2 \\ I_1^2 & W_{11}^2 & W_{12}^2 & W_{13}^2 \\ J_2^2 & W_{21}^2 & W_{22}^2 & W_{23}^2 \\ I_3^2 & W_{31}^2 & W_{32}^2 & W_{33}^2 \\ J_4^2 & W_{41}^2 & W_{42}^2 & W_{43}^2 \end{matrix} \quad 4 \times 3$$

b^2

$$\begin{bmatrix} b_1^2 & b_2^2 & b_3^2 \end{bmatrix} \quad 1 \times 3$$

Pending.

- Fwd. Prop.
 - Backpropagation
 - Why minimize weight with random value?
-
- Keras and Tensorflow
 - EDA of Business case.