

Multi-Layer Perceptrons

Nipun Batra

IIT Gandhinagar

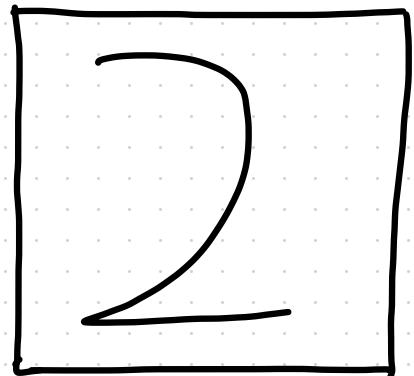
August 30, 2025

RECENT SUCCESSES OF NN

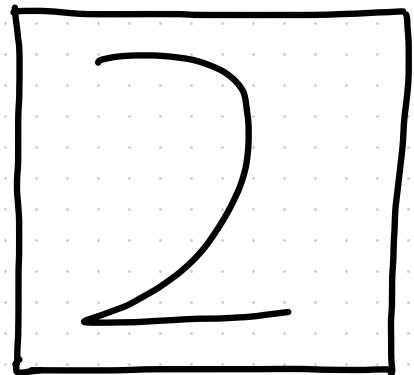
* State-of-the-art (SOTA) in most fields

PARADIGM CHANGE

PARADIGM CHANGE



PARADIGM CHANGE

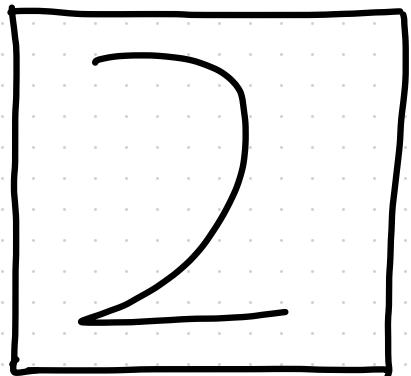


FEATURE

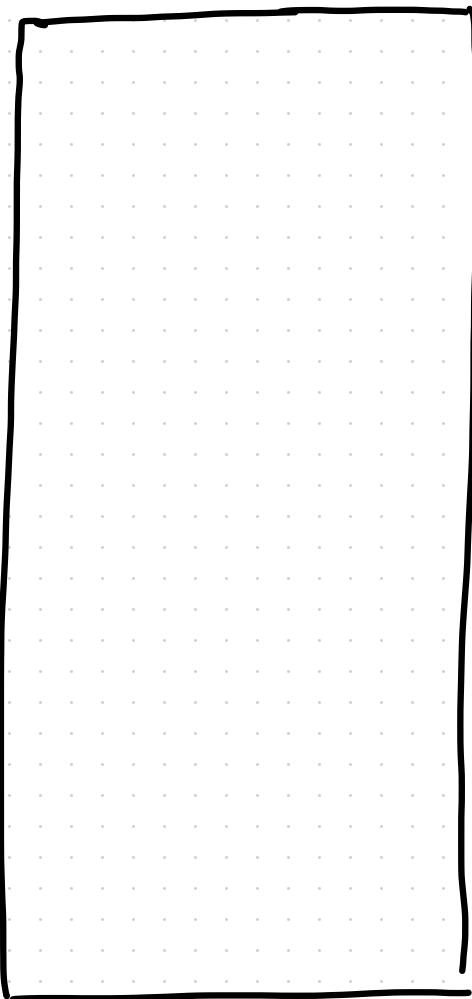


EXTRACTOR

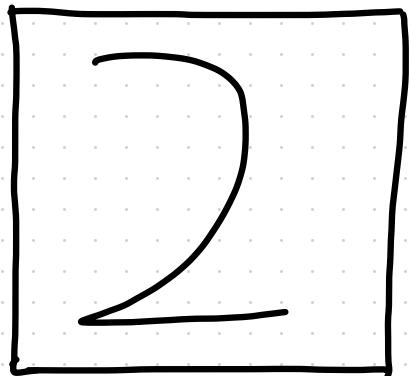
PARADIGM CHANGE



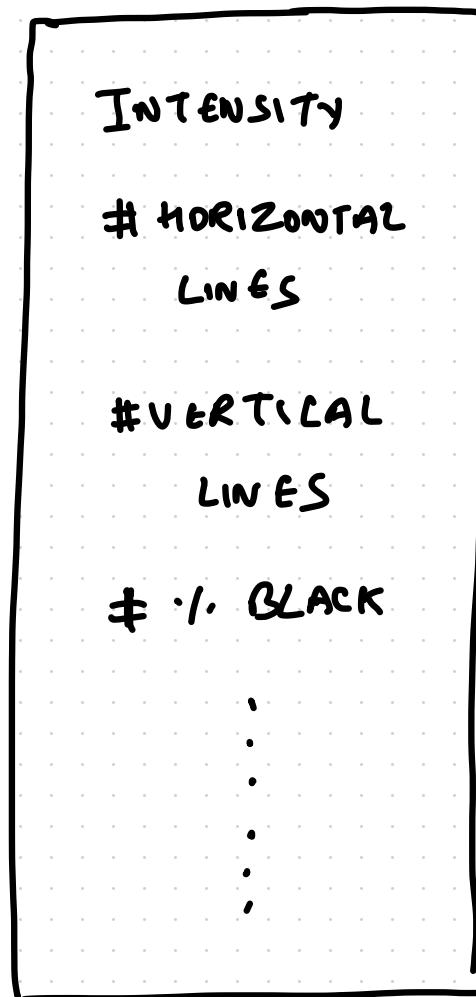
FEATURE
→
EXTRACTOR



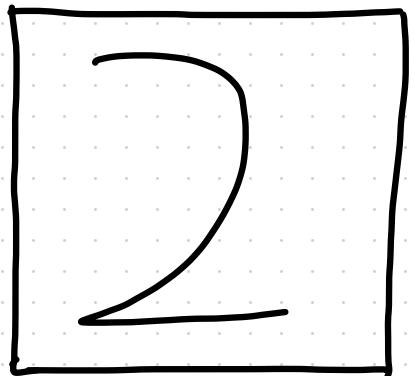
PARADIGM CHANGE



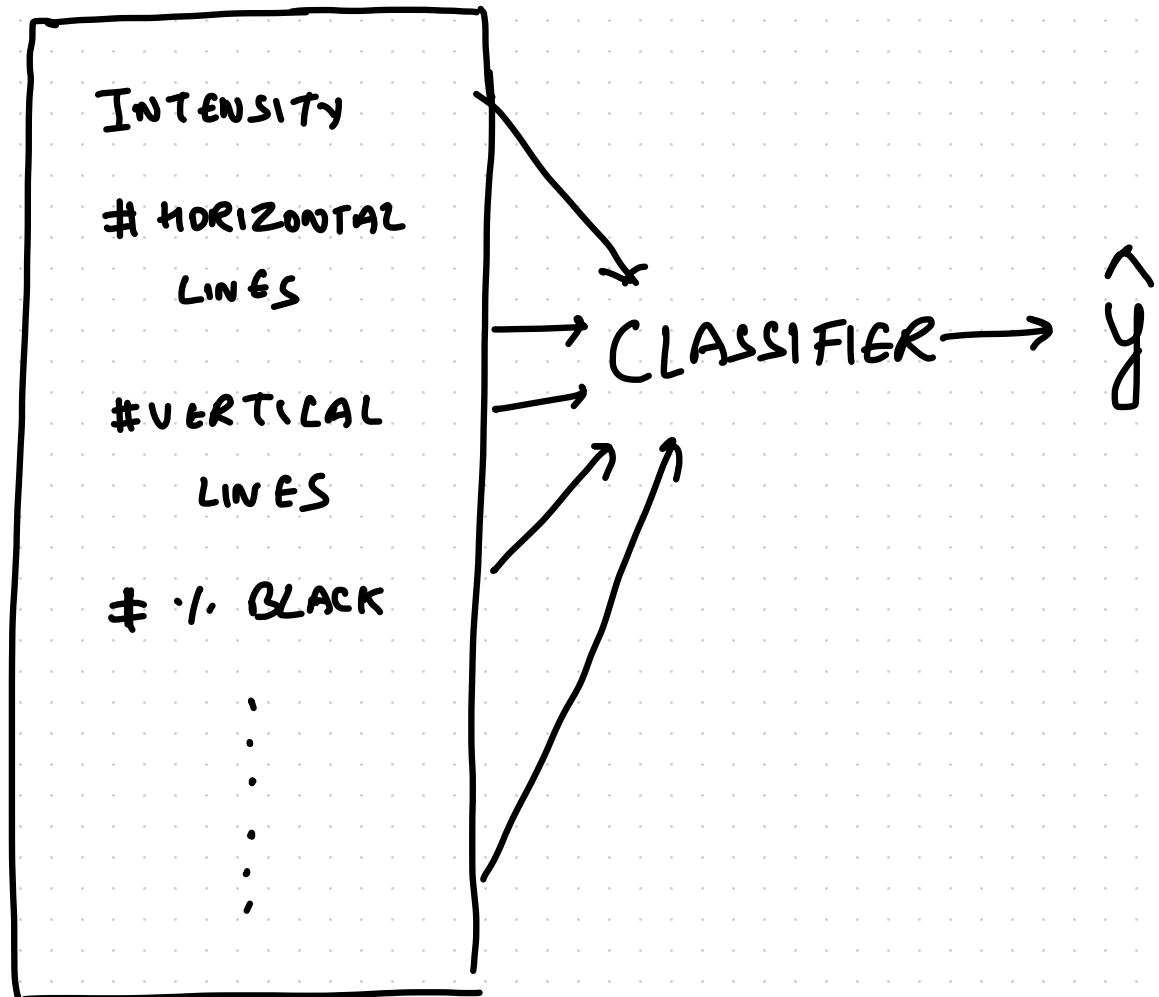
FEATURE
→
EXTRACTOR



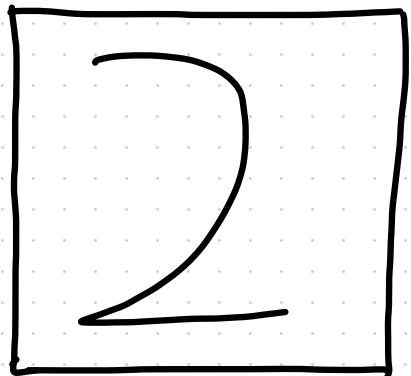
PARADIGM CHANGE



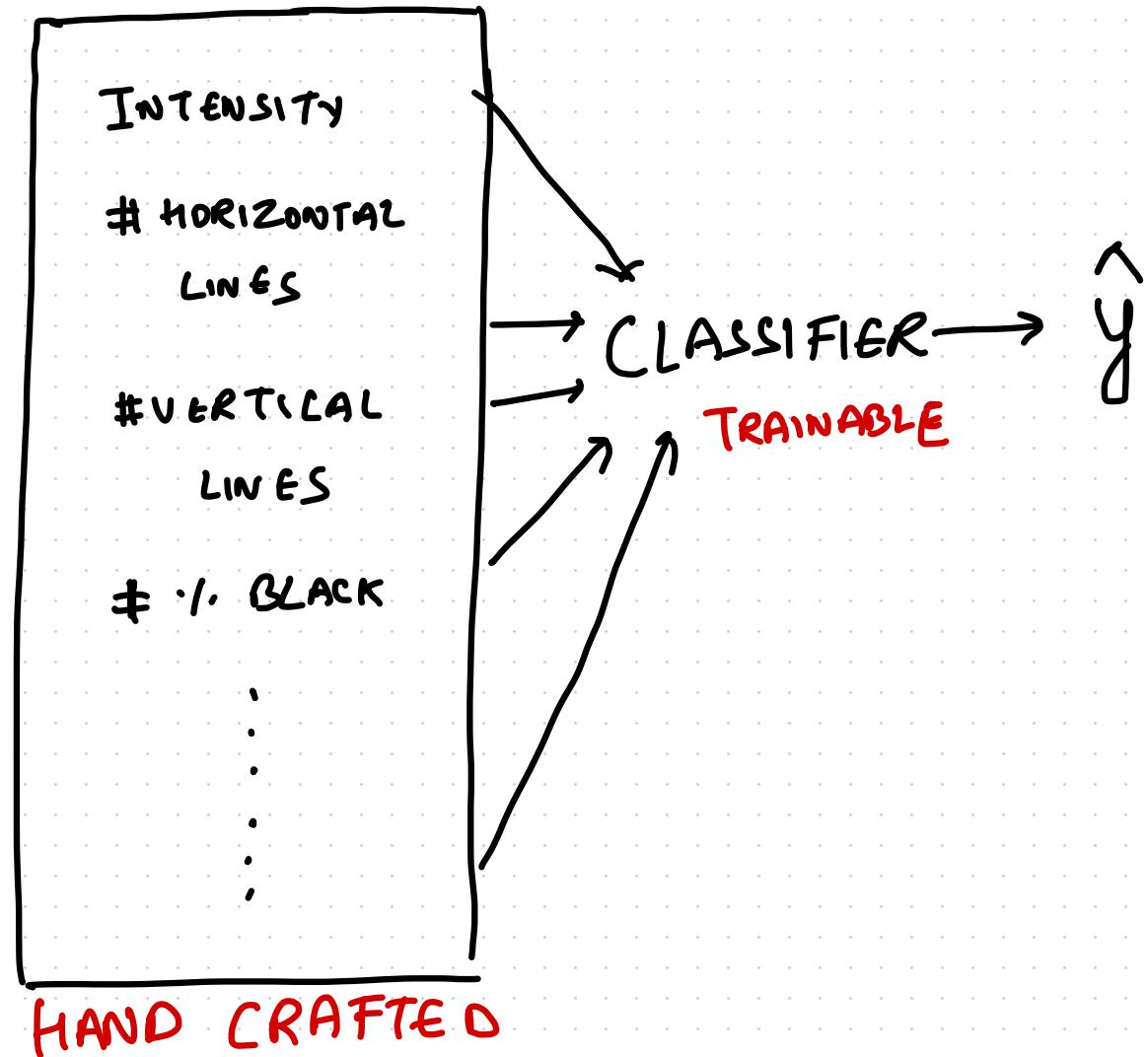
FEATURE
EXTRACTOR



PARADIGM CHANGE

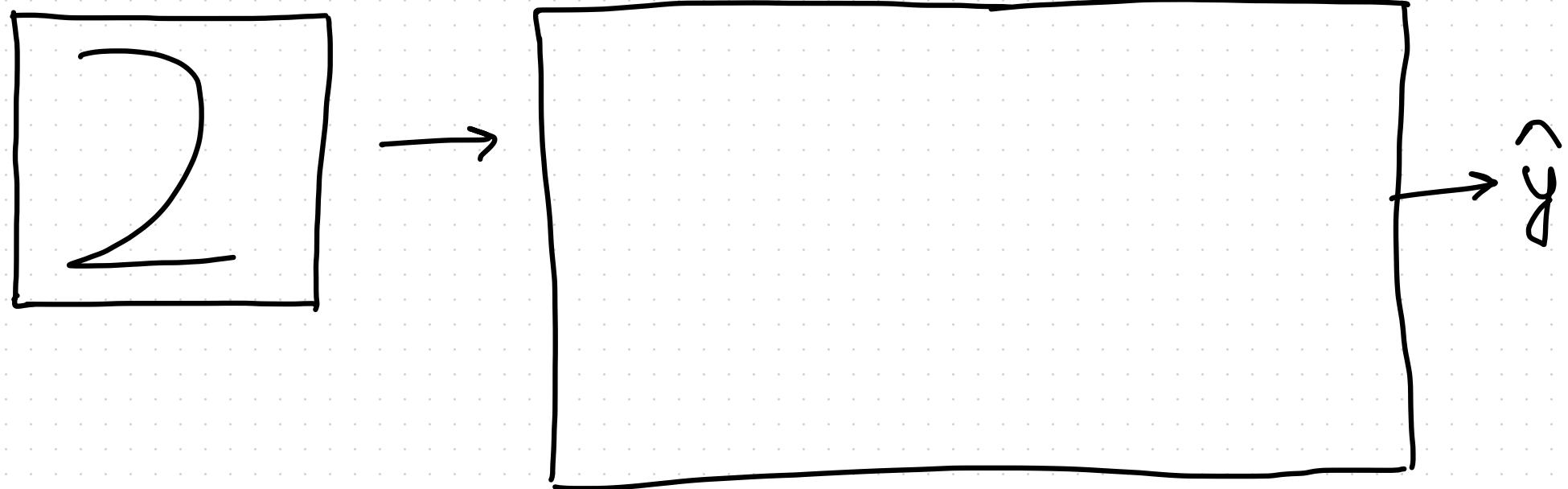


FEATURE
EXTRACTOR



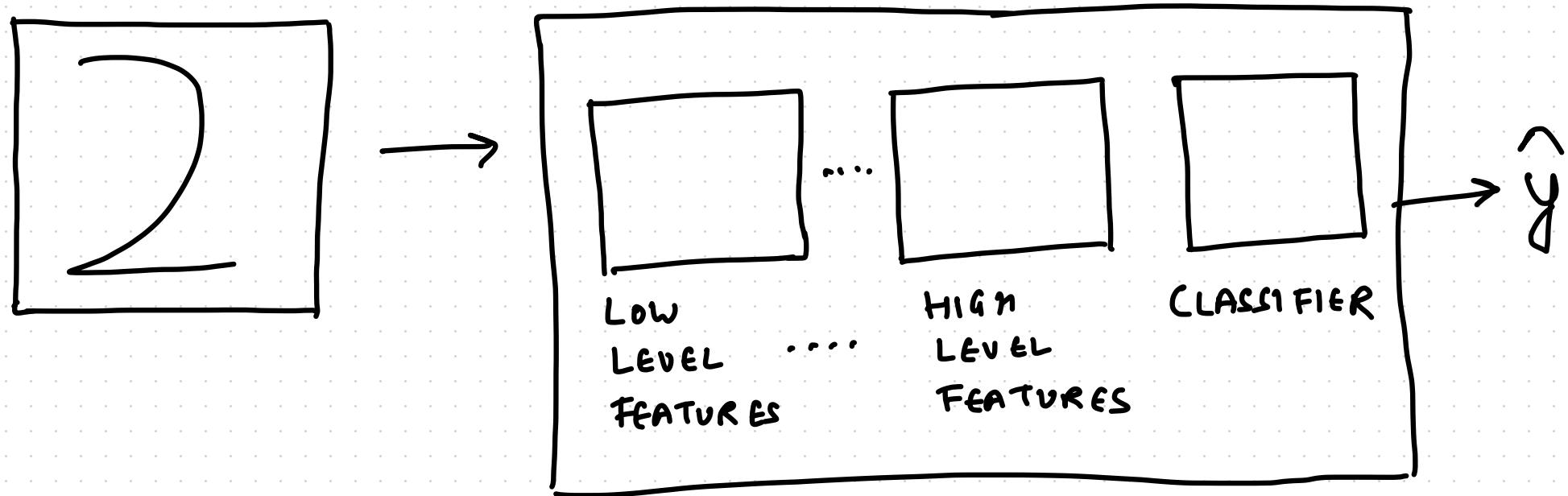
PARADIGM

CHANGE (NNS)

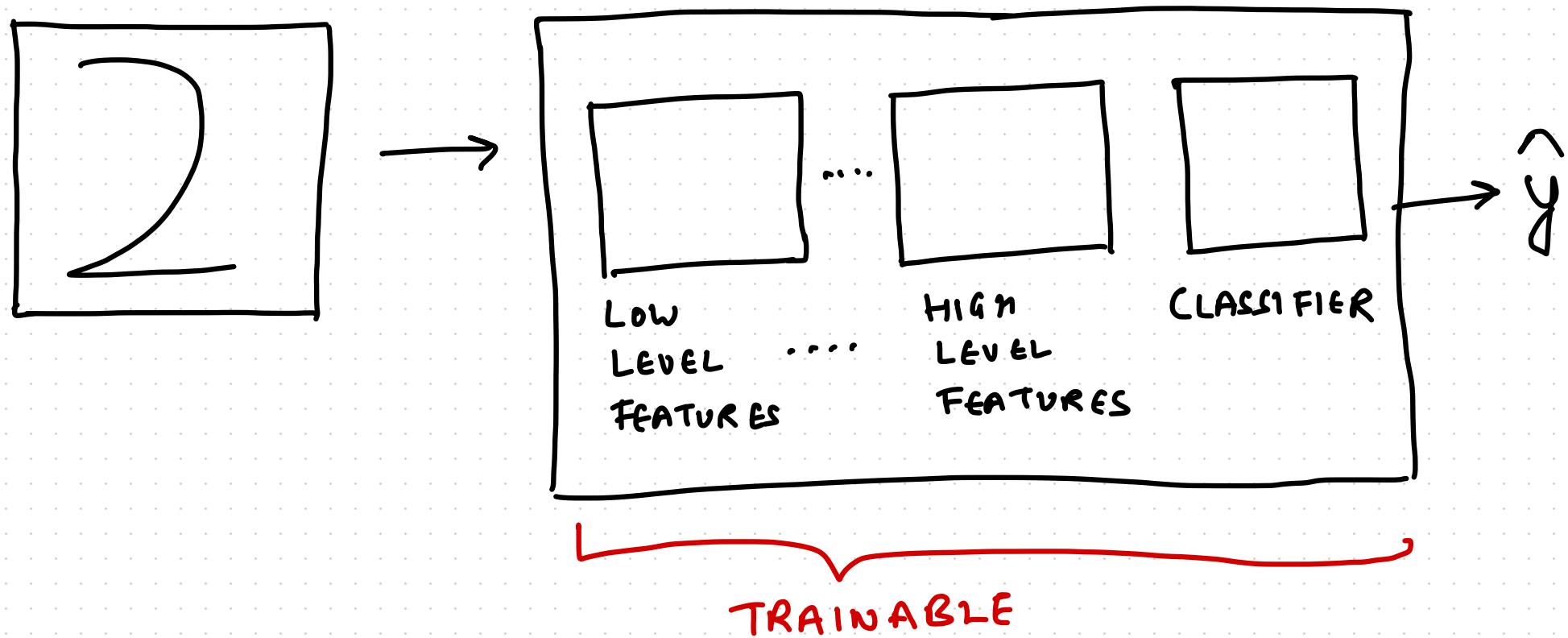


PARADIGM

CHANGE (NNs)



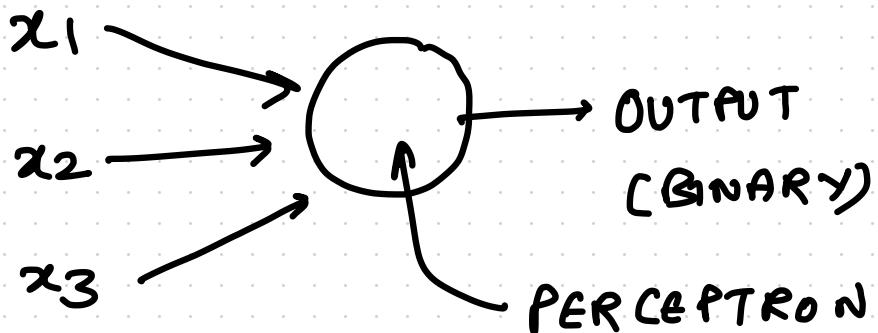
PARADIGM CHANGE (NNs)



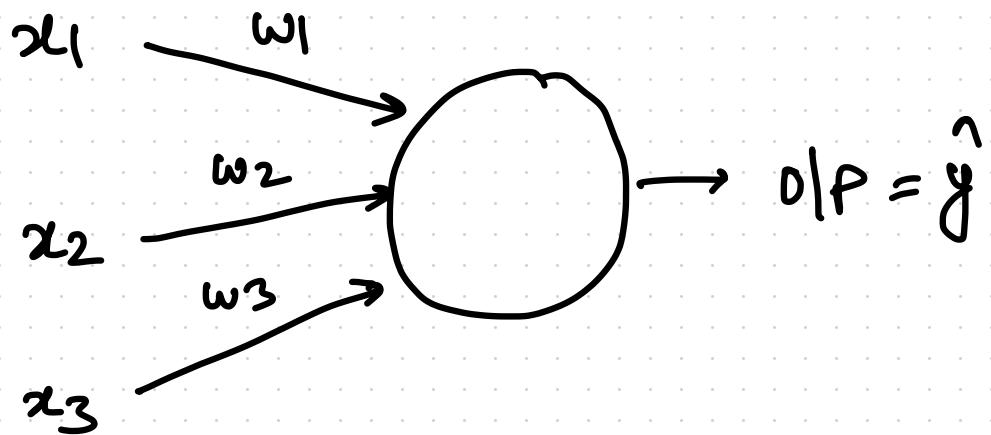
PERCEPTRON

— ARTIFICIAL NEURON DEVELOPED BY
ROSENBLATT IN 1960^s INSPIRED BY
MCCULLOCH & PITTS

BINARY IP

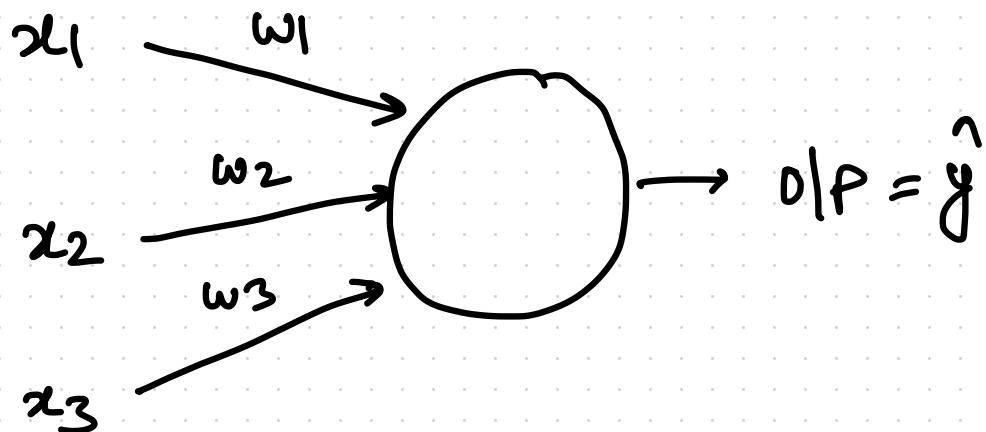


PERCEPTRON



$$O/P = \hat{y} = \begin{cases} 0 & ; \sum w_i x_i \leq \text{THRESHOLD} \\ 1 & ; \sum w_i x_i > \text{THRESHOLD} \end{cases}$$

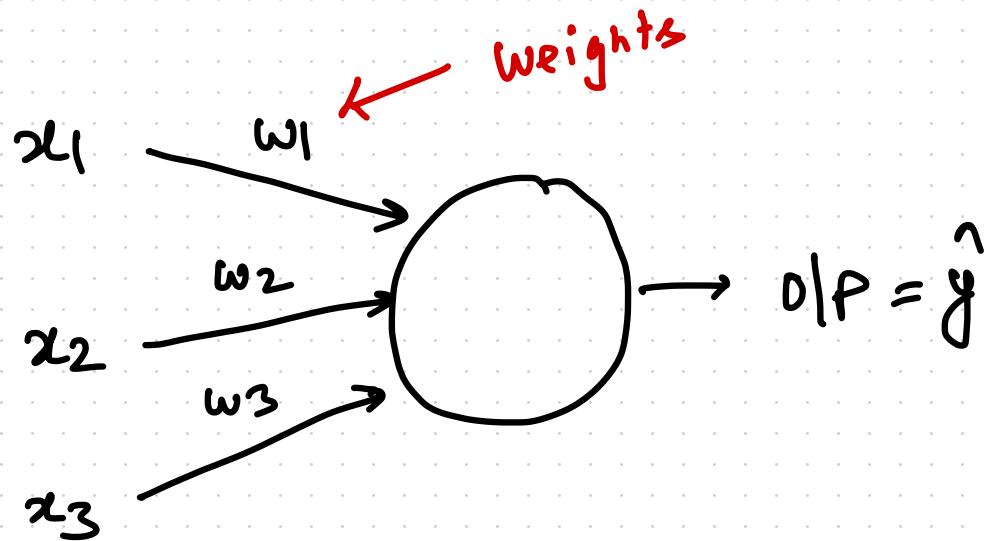
PERCEPTRON



$$O/P = \hat{y} = \begin{cases} 0 & ; \sum w_i x_i \leq \text{THRESHOLD} \\ 1 & ; \sum w_i x_i > \text{THRESHOLD} \end{cases}$$

NEURONS "FIRE" ABOVE THRESHOLD

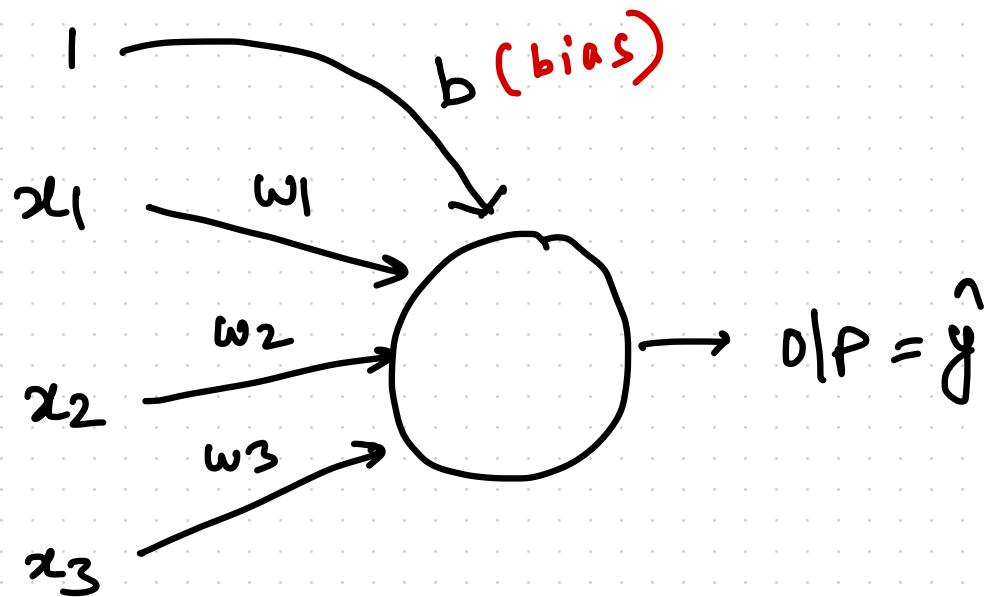
PERCEPTRON



$$O/P = \hat{y} = \begin{cases} 0 & ; \sum w_i x_i \leq \text{THRESHOLD} \\ 1 & ; \sum w_i x_i > \text{THRESHOLD} \end{cases}$$

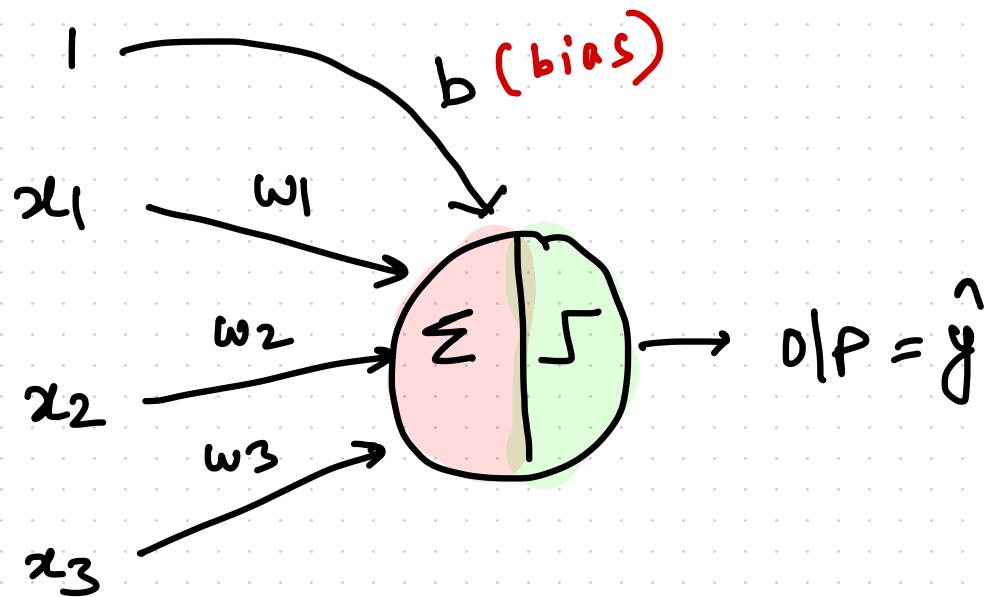
NEURONS "FIRE" ABOVE
THRESHOLD

PERCEPTRON



$$O/P = \hat{y} = \begin{cases} 0 & ; \sum w_i x_i + b \leq 0 \\ 1 & ; \sum w_i x_i + b > 0 \end{cases}$$

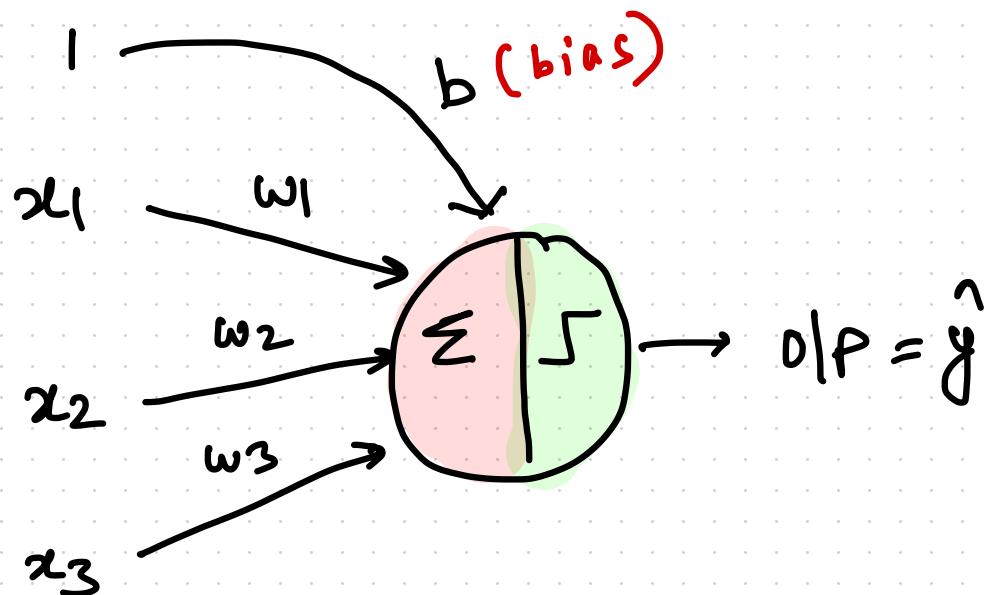
PERCEPTRON



NEURON HAS 2 COMPONENTS

- ① SUMMATION
- ② ACTIVATION : STEP FUNCTION

PERCEPTRON



NEURON HAS 2 COMPONENTS

(1) SUMMATION

(2) ACTIVATION : STEP FUNCTION

SIGMO STEP
Activation



LEARNING BINARY GATES

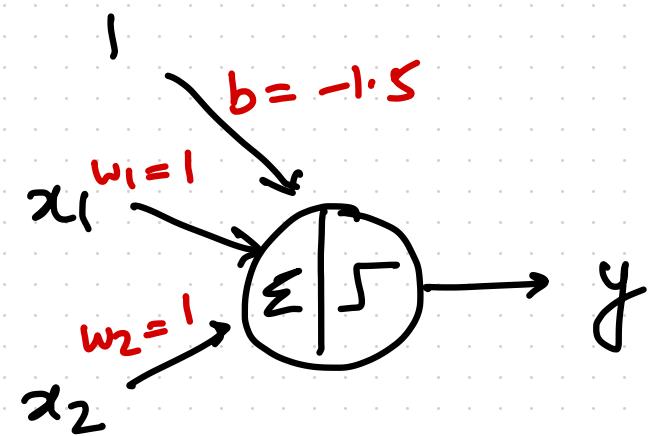
Q) FOR 2 I/P's $x_1 \& x_2$ learn w's and b for
BINARY AND

x_1	x_2	y
0	0	0
0	1	0
1	0	0
1	1	1

LEARNING BINARY GATES

Q) FOR 2 I/P's $x_1 \& x_2$ learn w's and b for
BINARY AND

x_1	x_2	y
0	0	0
0	1	0
1	0	0
1	1	1



LEARNING BINARY GATES

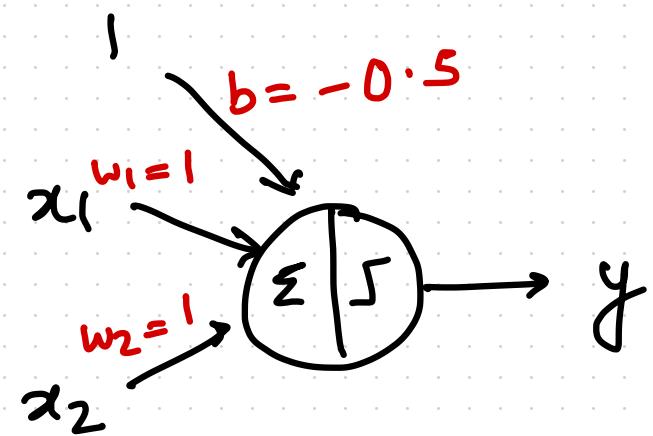
Q) FOR 2 IIPS $x_1 \neq x_2$ learn w's and b for
BINARY OR

x_1	x_2	y
0	0	0
0	1	1
1	0	1
1	1	1

LEARNING BINARY GATES

Q) FOR 2 IIPS $x_1 \neq x_2$ learn w's and b for
BINARY OR

x_1	x_2	y
0	0	0
0	1	1
1	0	1
1	1	1



LEARNING BINARY GATES

Q) FOR 1 IIPS x_1 learn w's and b for
UNARY NOT

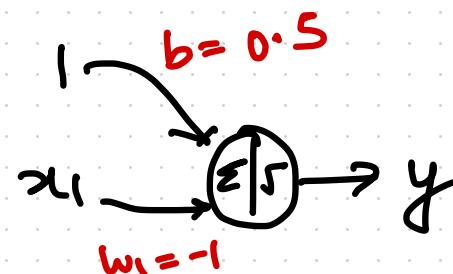
x_1	y
0	1
1	0

LEARNING BINARY GATES

Q) FOR 1 IIPS x_1
UNARY NOT

learn w's and b for

x_1	y
0	1
1	0



LEARNING BINARY GATES

Q) FOR 2 IPs $x_1 \& x_2$ learn w's and b for
NAND

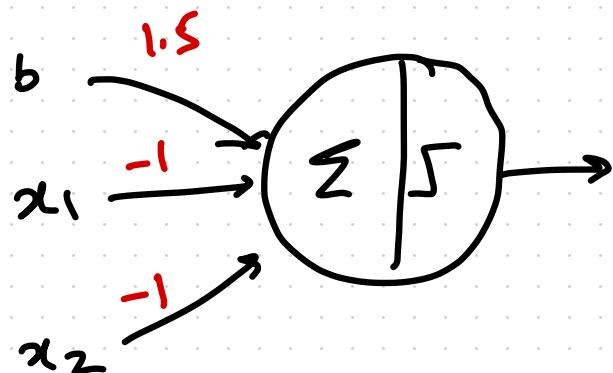
x_1	x_2	y
0	0	1
0	1	1
1	0	1
1	1	0

LEARNING BINARY GATES

Q) FOR 2 IIPS $x_1 \& x_2$ learn w's and b for
NAND

x_1	x_2	y
0	0	1
0	1	1
1	0	1
1	1	0

APPROACH #1

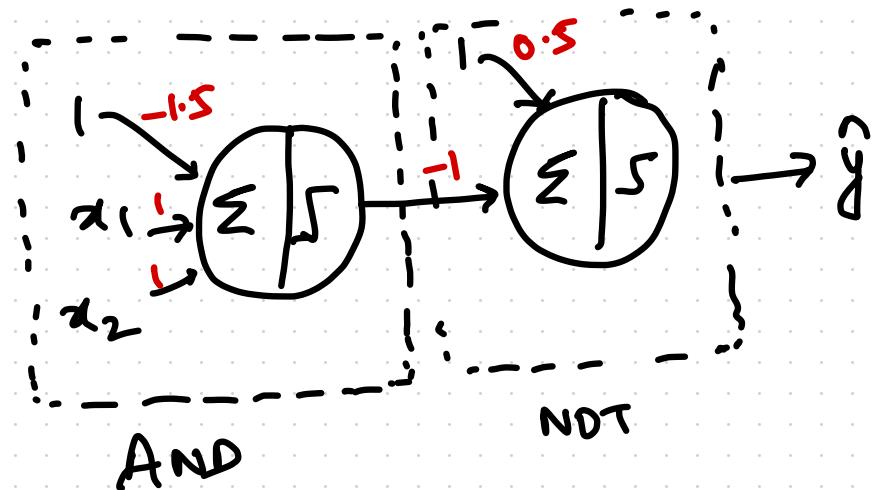
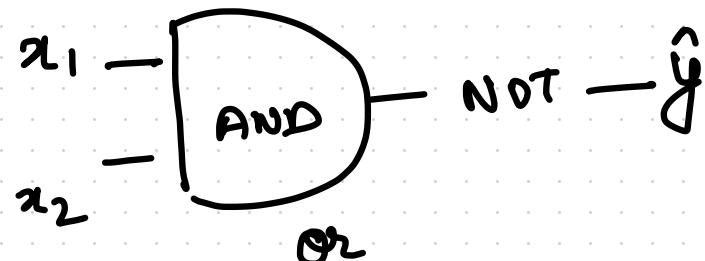


LEARNING BINARY GATES

Q) FOR 2 IIPS $x_1 \& x_2$ learn w's and b for
NAND

x_1	x_2	y
0	0	1
0	1	1
1	0	1
1	1	0

APPROACH #2



PERCEPTRON LEARNING ALGORITHM

IP: $X \xrightarrow{N \times D}$; $y \xrightarrow{N \times 1}$; $lr \xrightarrow{\text{learning rate}}$; $it \xrightarrow{\# \text{iterations}}$

PERCEPTRON LEARNING ALGORITHM

IIP: $X \xrightarrow{N \times D} ; y \xrightarrow{N \times 1} ; lr \xrightarrow{\text{learning rate}} ; it \xrightarrow{\# \text{iterations}}$

SI AUGMENT $x \rightarrow x' - s.t. x'[1:s] = x ; x'[1] = [;]$

PERCEPTRON LEARNING ALGORITHM

IIP: $X \xrightarrow{NxD} ; y \xrightarrow{Nx1} ; lr \xrightarrow{\text{learning rate}} ; it \xrightarrow{\# \text{iterations}}$

S1 AUGMENT $x \rightarrow x' - s.t. x'[1:s] = x ; x'[s] = [1; :]$

S2 INIT $W \in R^{D+1}$

PERCEPTRON LEARNING ALGORITHM

IP: $X \xrightarrow{NxD} ; y \xrightarrow{Nx1} ; lr \xrightarrow{\text{learning rate}} ; it \xrightarrow{\# \text{iterations}}$

S1 AUGMENT x to x' - s.t. $x'[1:s] = x$; $x'[1] = [1; \dots]$

S2 INIT $W \in \mathbb{R}^{D+1}$

S3 FOR i IN IT :

PERCEPTRON LEARNING ALGORITHM

IP: $x \xrightarrow{N \times D}$; $y \xrightarrow{N \times 1}$; $\eta \xrightarrow{\text{learning rate}}$; $it \xrightarrow{\# \text{iterations}}$

S1 AUGMENT x to x' - s.t. $x'[1:s] = x$; $x'[1] = [1; \dots]$

S2 INIT $w \in \mathbb{R}^{D+1}$

S3 FOR i IN IT:

S3.1 FOR d in D :

S3.1.1 $\hat{y} = \text{STEP}(x' \cdot w)$

PERCEPTRON LEARNING ALGORITHM

IP: $X \xrightarrow{N \times D}; y \xrightarrow{N \times 1}$; $\eta \xrightarrow{\text{learning rate}}$; $it \xrightarrow{\# \text{iterations}}$

S1 AUGMENT x to x' - s.t. $x'[1:s] = x$; $x'[1] = [1; \dots]$

S2 INIT $w \in \mathbb{R}^{D+1}$

S3 FOR i IN IT:

S3.1 FOR d in D :

$$\hat{y} = \text{STEP}(x' \cdot w)$$

$$\underline{\text{S3.1.2}} \quad \text{ERR} = y - \hat{y}$$

PERCEPTRON LEARNING ALGORITHM

IP: $X \xrightarrow{N \times D}; y \xrightarrow{N \times 1}; lr \xrightarrow{\text{learning rate}}; it \xrightarrow{\# \text{iterations}}$

S1 AUGMENT x to x' - s.t. $x'[1:s] = x$; $x'[s] = [1; \dots]$

S2 INIT $w \in \mathbb{R}^{D+1}$

S3 FOR i IN IT:

S3.1 FOR n in N :

$$\hat{y} = \text{STEP}(x' \cdot w)$$

$$ERR = y - \hat{y}$$

$$w \leftarrow w + lr * ERR_n * x'[n]$$

PERCEPTRON LEARNING ALGORITHM

IP: $X \xrightarrow{N \times D}; y \xrightarrow{N \times 1}; lr \xrightarrow{\text{learning rate}}; it \xrightarrow{\# \text{iterations}}$

S1 AUGMENT x to x' - s.t. $x'[1:s] = x$; $x'[s] = [1; \dots]$

S2 INIT $w \in \mathbb{R}^{D+1}$

S3 FOR i IN IT:

S3.1 FOR n in N :

$$\hat{y} = \text{STEP}(x' \cdot w)$$

$$err = y - \hat{y}$$

$$w \leftarrow w + lr * err_n * x'[n]$$

ERROR IN n^{th} sample

err_n

$x'[n]$

n^{th} sample

PERCEPTRON LEARNING ALGORITHM

IP: $X \xrightarrow{N \times D} ; y \xrightarrow{N \times 1} ; lr \xrightarrow{\text{learning rate}} ; it \xrightarrow{\# \text{iterations}}$

S1 AUGMENT x to x' → s.t. $x'[1:s] = x$; $x'[s] = [1; \dots]$

S2 INIT $w \in \mathbb{R}^{D+1}$

S3 FOR i IN IT :

S3.1 FOR n in N :

$$\hat{y} = \text{STEP}(x' \cdot w)$$

$$err = y - \hat{y}$$

S3.1.3

$$w \leftarrow w + lr * err_n * x'[n]$$

Analogous to S.G.D

Analogous to Gradient

ERROR IN n^{th} Sample

err_n

$x'[n]$

n^{th} sample

PERCEPTRON LEARNING ALGORITHM

S31.3 $w \leftarrow w + lr * \text{ERR}_n * x'[n]$

IMAGINE

$$w = \begin{bmatrix} 1 \\ -1 \\ -1 \end{bmatrix} ; x'[n] = \begin{bmatrix} 1 \\ 0 \\ 0 \end{bmatrix} ; y = 0$$

The diagram illustrates a perceptron model. It shows an input vector x' and a weight vector w . The input vector x' is represented as a column vector with elements 1, 0, and 0. The weight vector w is represented as a column vector with elements 1, -1, and -1. The output y is labeled as 0. Below the vectors, there is a circle containing the text "ε / r". An arrow points from the input vector x' to this circle, and another arrow points from the weight vector w to the same circle. From the circle, an arrow points to the output $\hat{y} = 1$.

PERCEPTRON LEARNING ALGORITHM

S3.1.3 $w \leftarrow w + lr * ERR_n * x'[n]$

IMAGINE

$$w = \begin{bmatrix} 1 \\ -1 \\ -1 \end{bmatrix}$$

$$; \quad x'[n] = \begin{bmatrix} 1 \\ 0 \\ 0 \end{bmatrix}; \quad y = 0$$

ϵ | r → $\hat{y} = 1$

$$ERR_n = 0 - 1 = -1$$

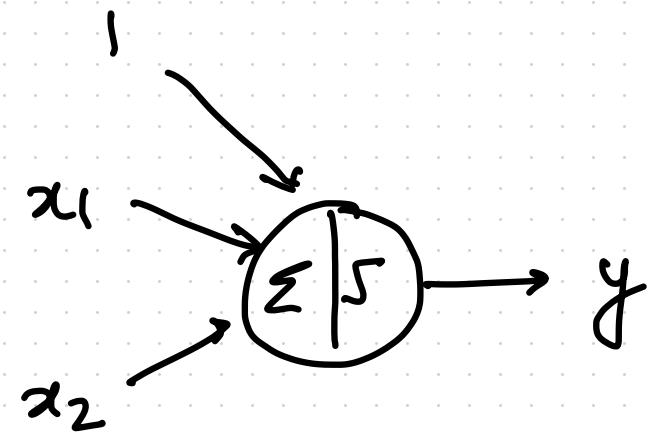
$$w \leftarrow \begin{bmatrix} 1 \\ -1 \\ -1 \end{bmatrix} + 0.01 * -1 * \begin{bmatrix} 1 \\ 0 \\ 0 \end{bmatrix}$$

$$\leftarrow \begin{bmatrix} 0.99 \\ -1 \\ -1 \end{bmatrix}$$

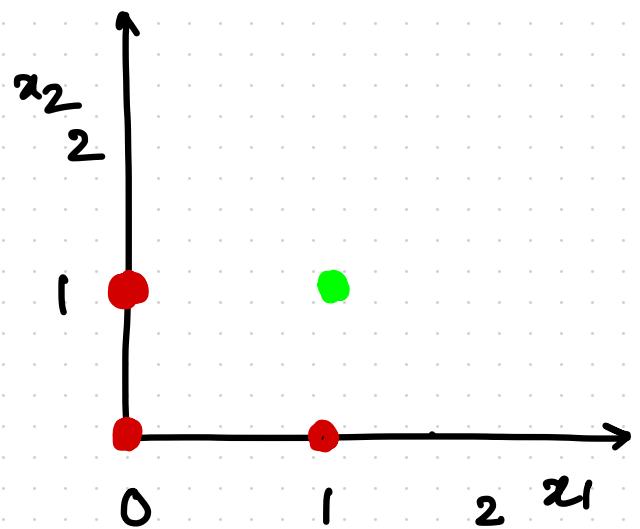
LEARNING BINARY GATES

Q) FOR 2 IIPS $x_1 \neq x_2$ learn w's and b for
BINARY XOR

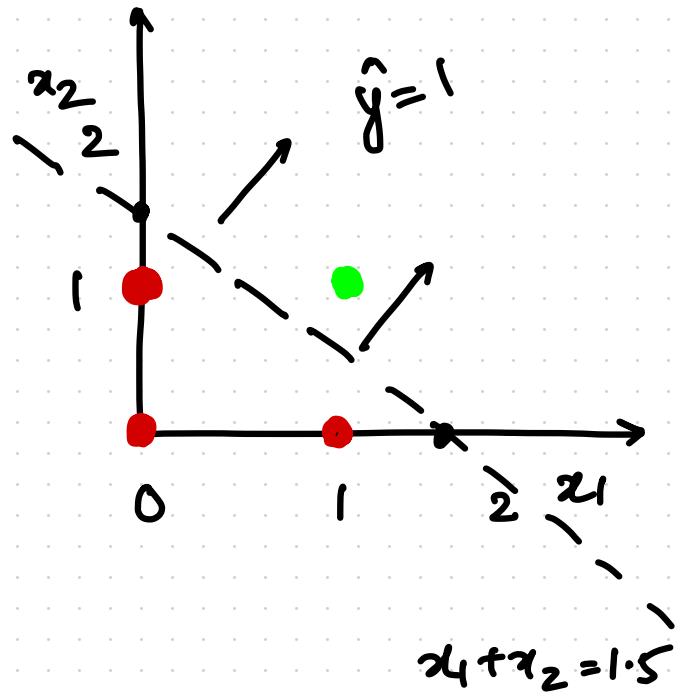
x_1	x_2	y
0	0	0
0	1	1
1	0	1
1	1	0



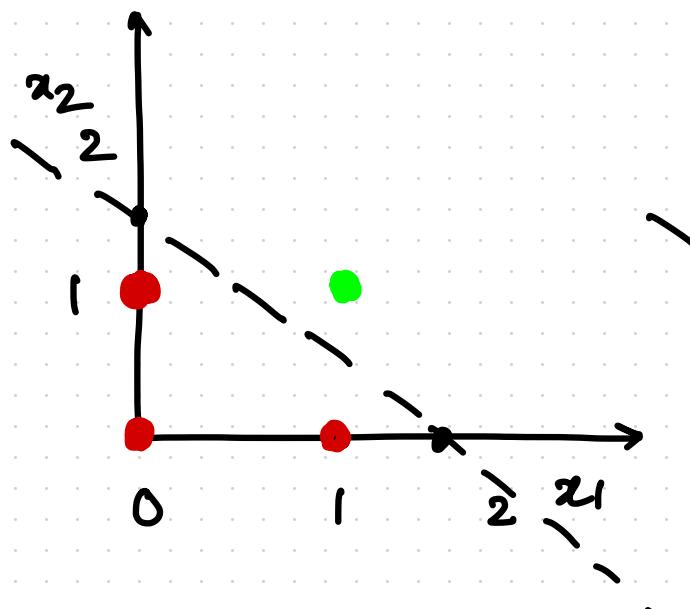
AND



AND

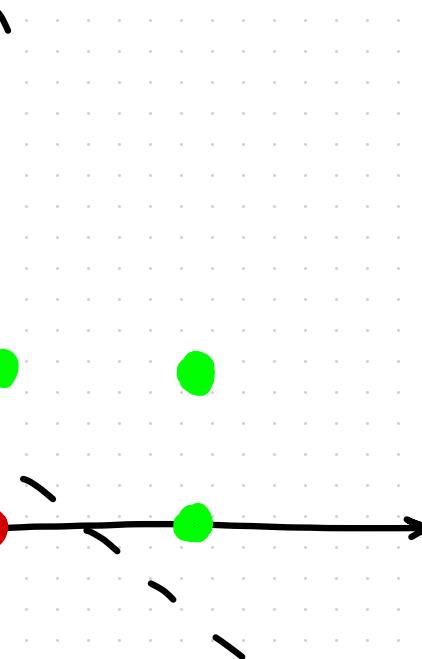


AND



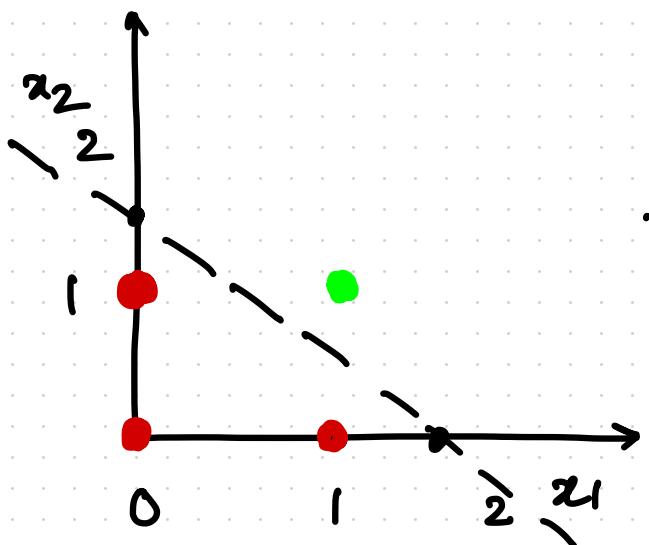
$$x_1 + x_2 = 1.5$$

OR



$$\begin{aligned}x_1 + x_2 \\= 0.5\end{aligned}$$

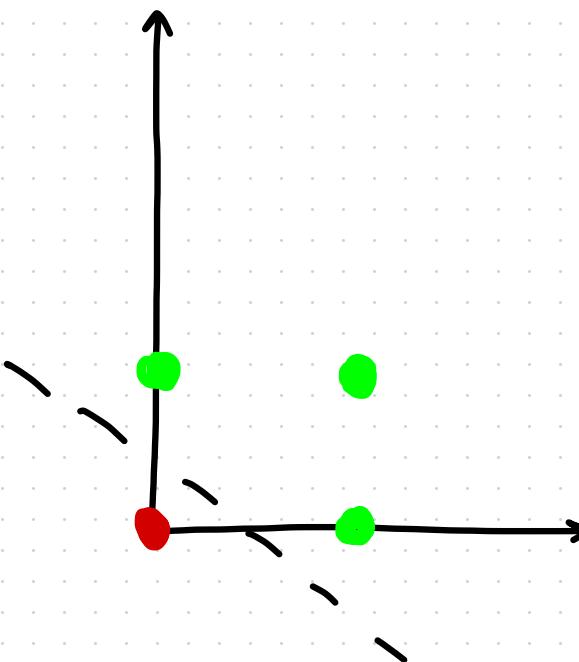
AND



$$x_1 + x_2 = 1.5$$

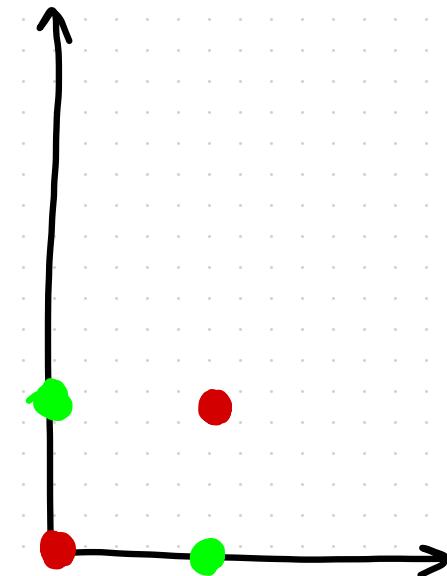
LINEARLY SEPARABLE

OR



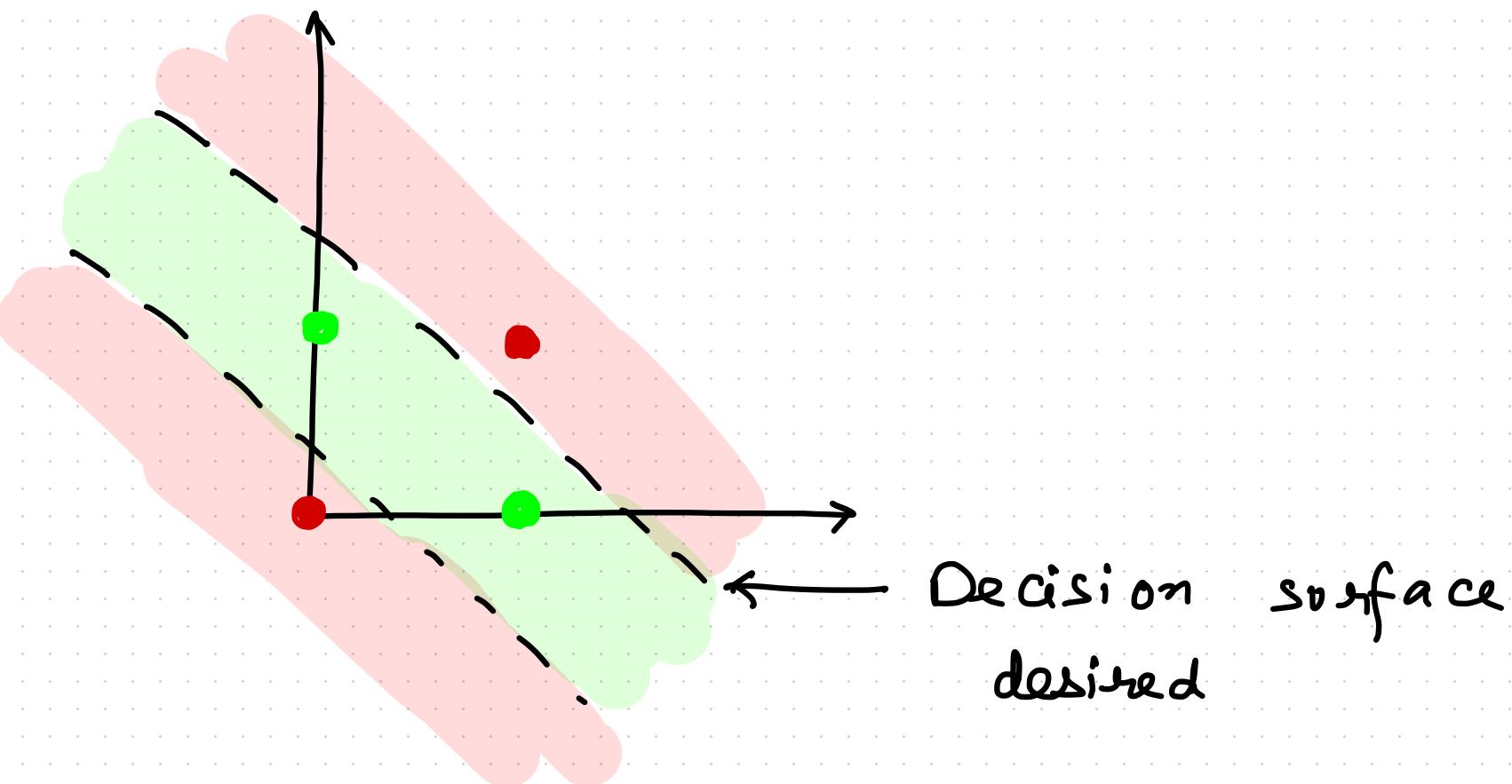
$$x_1 + x_2 = 0.5$$

XOR

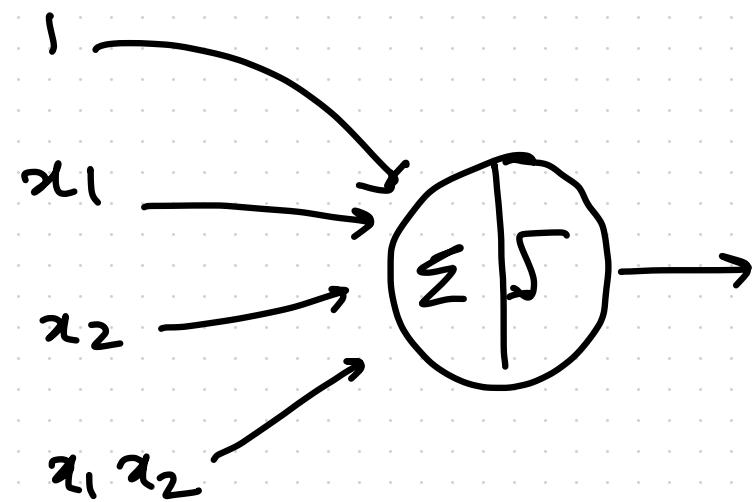
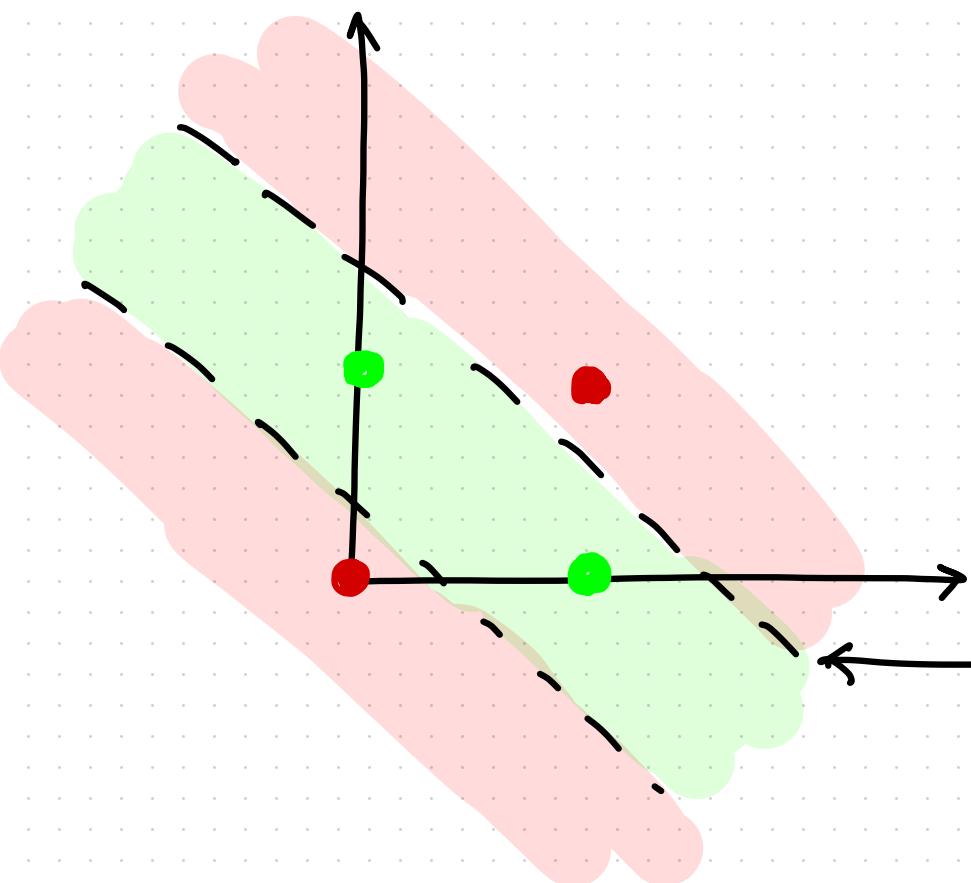


NOT
SEPARABLE
LINEARLY

XOR CLASSIFICATION

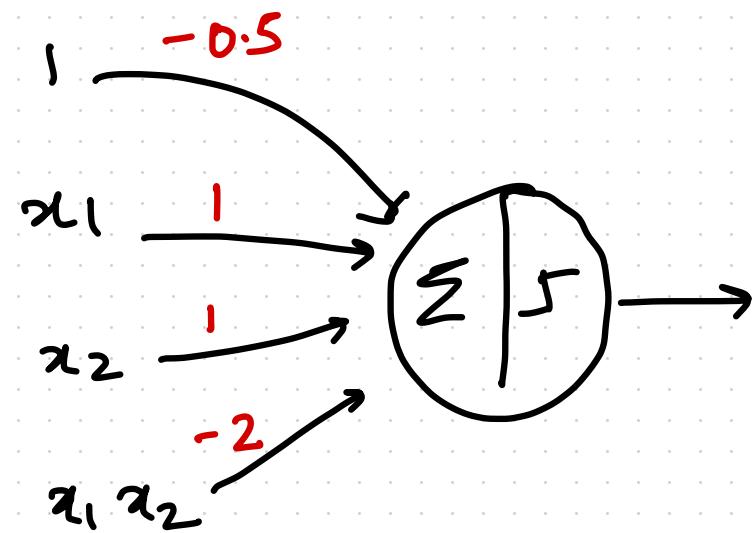
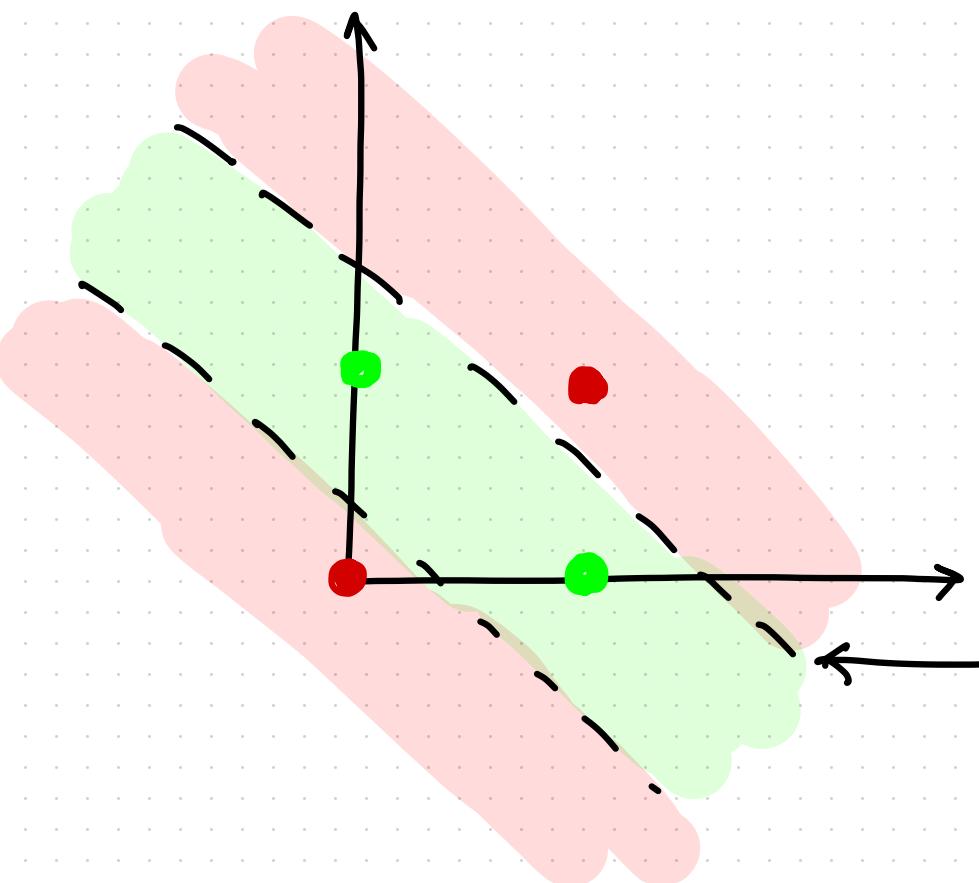


XOR CLASSIFICATION



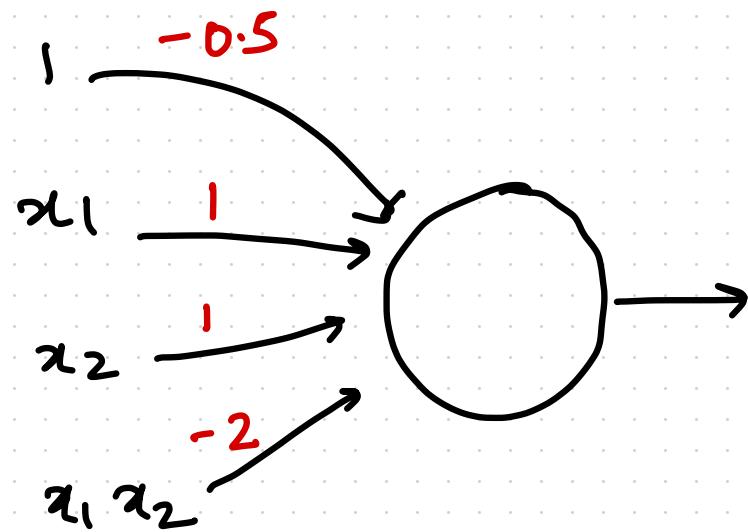
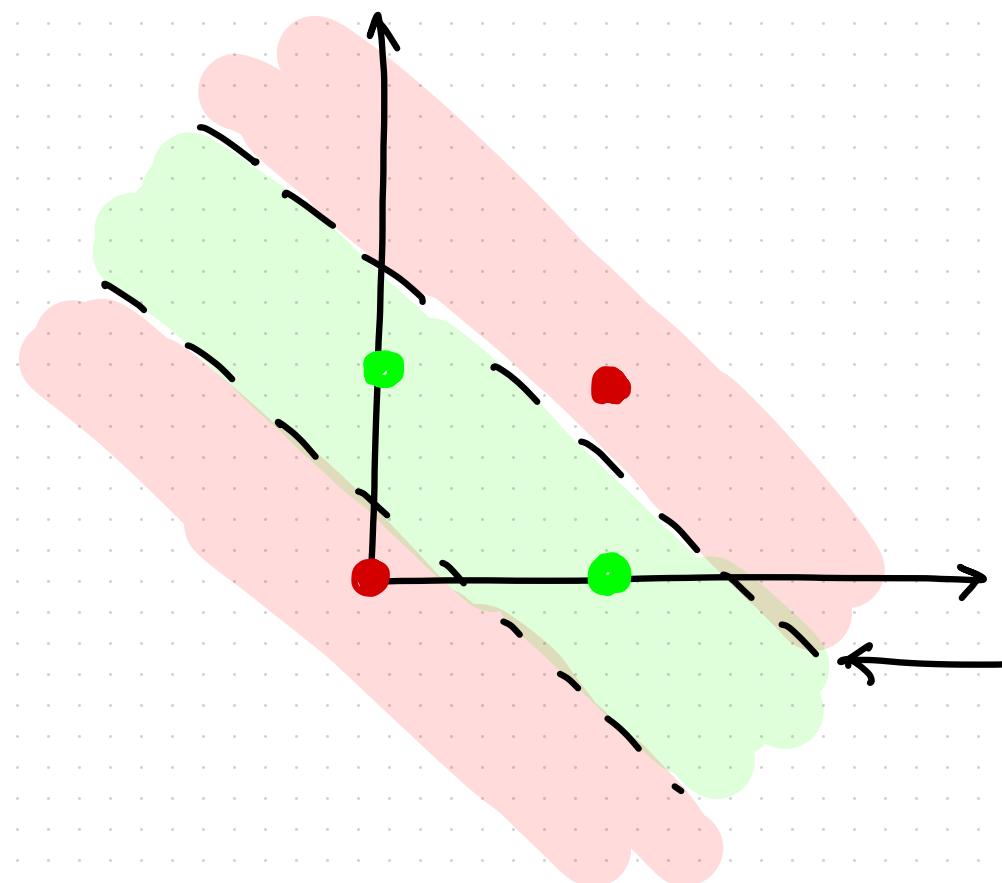
Decision surface
desired

XOR CLASSIFICATION



Decision surface
desired

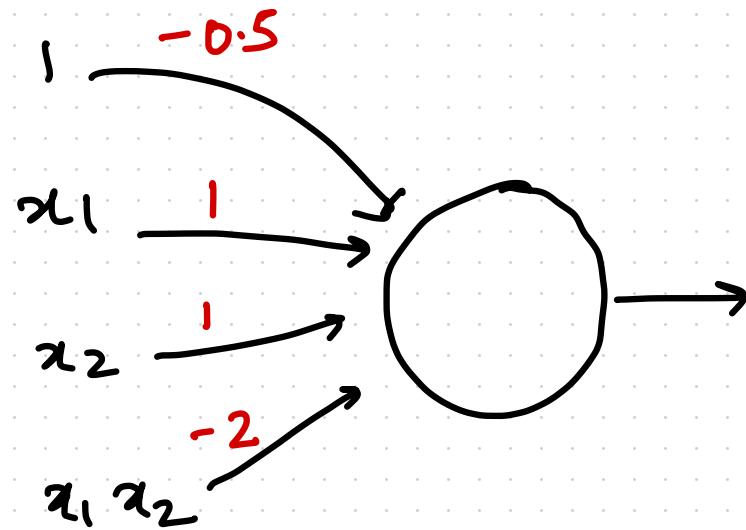
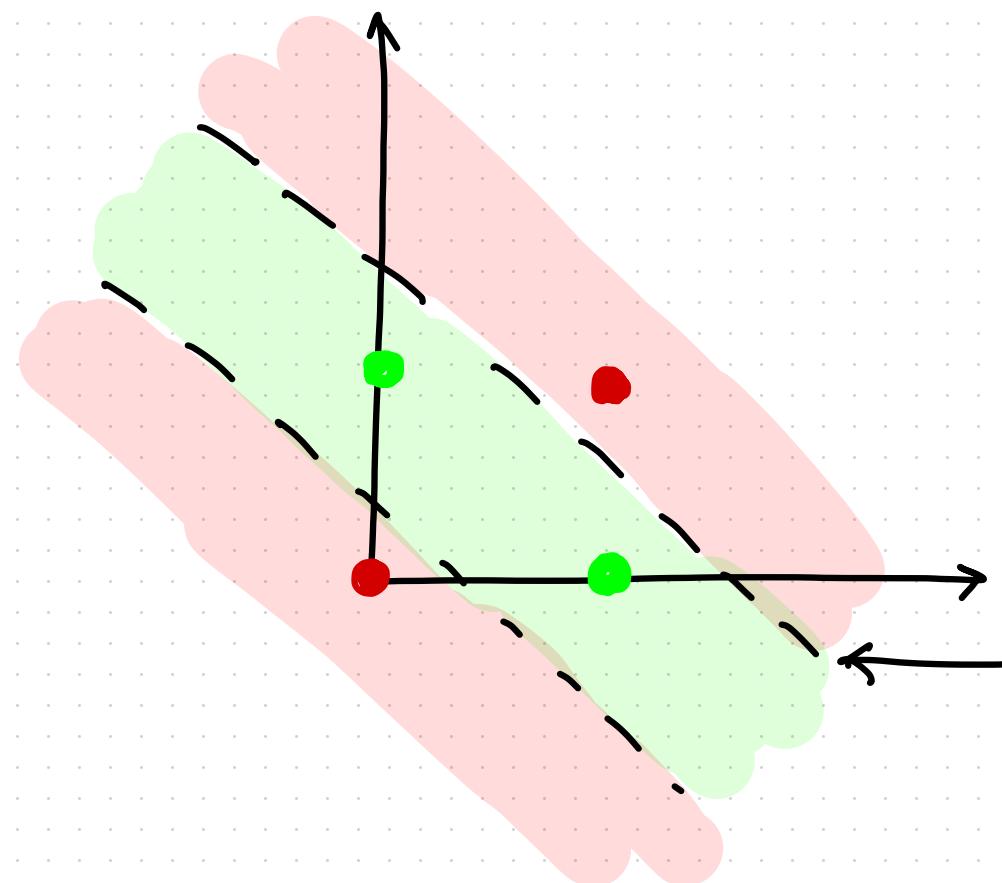
XOR CLASSIFICATION



FOR $x_1 = 0 ; x_2 = 0$; we get
 $\hat{y} = \{-0.5 \leq 0\} = \text{RED CLASS}$

Decision surface
desired

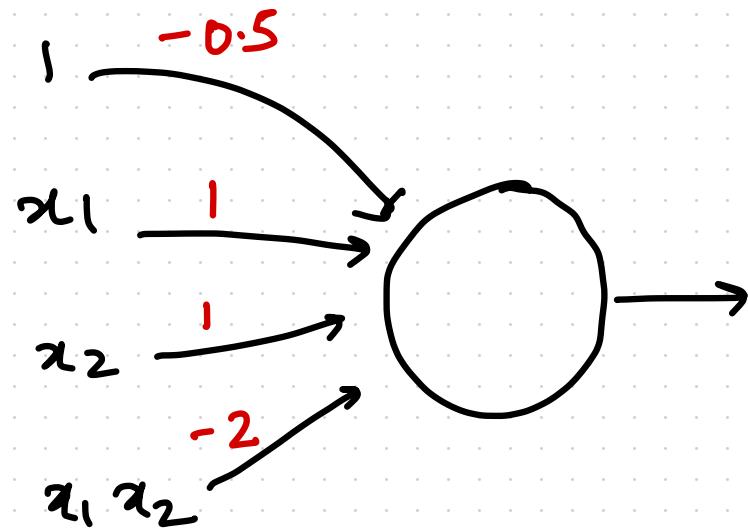
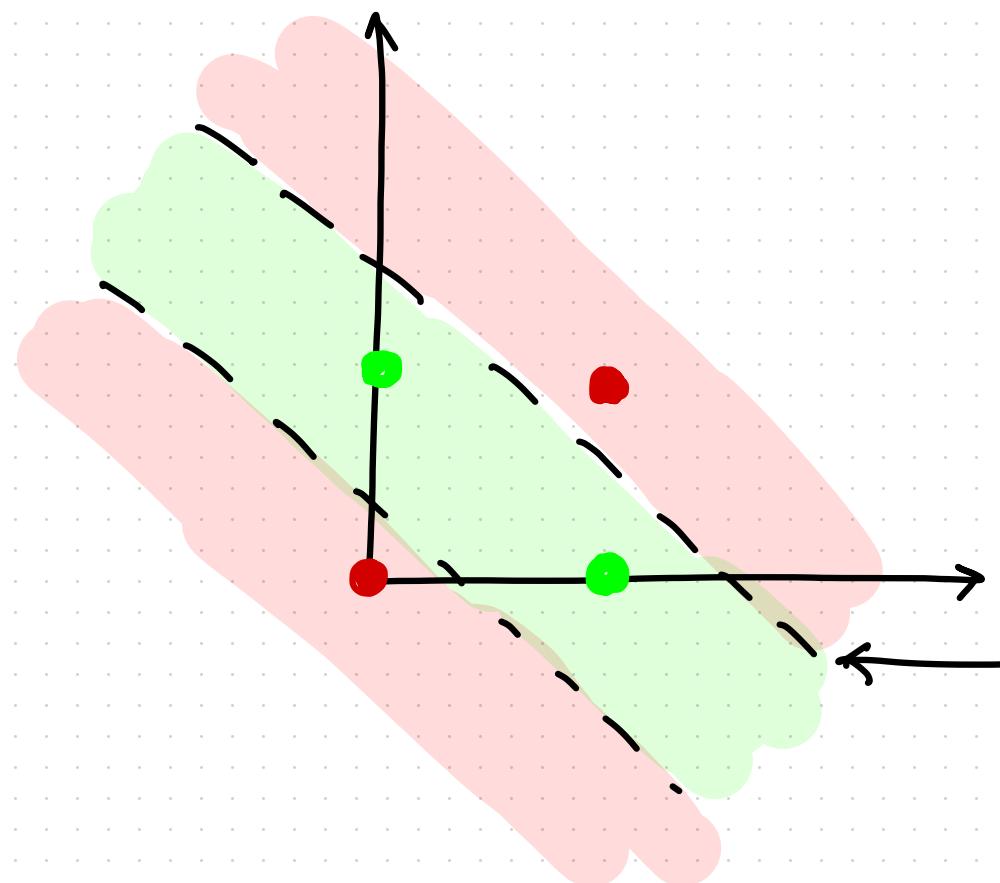
XOR CLASSIFICATION



FOR $x_1 = 0 ; x_2 = 1$; we get
 $\hat{y} = \{-0.5 + 1 \leq 0\} = \text{GREEN CLASS}$

Decision surface
desired

XOR CLASSIFICATION

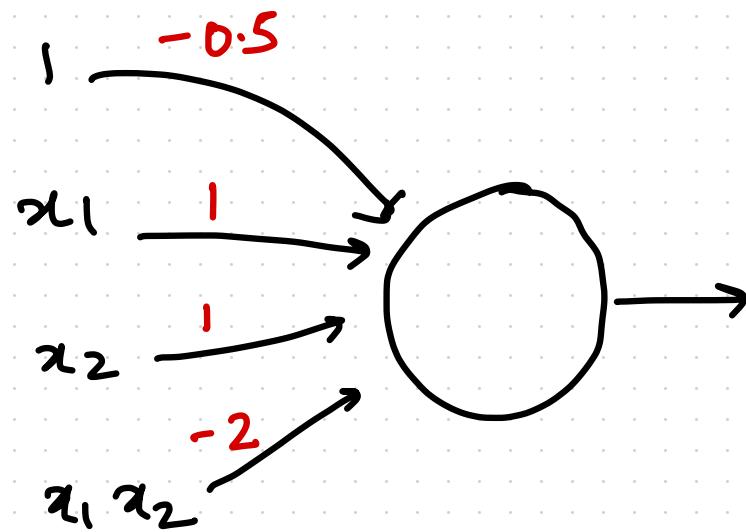
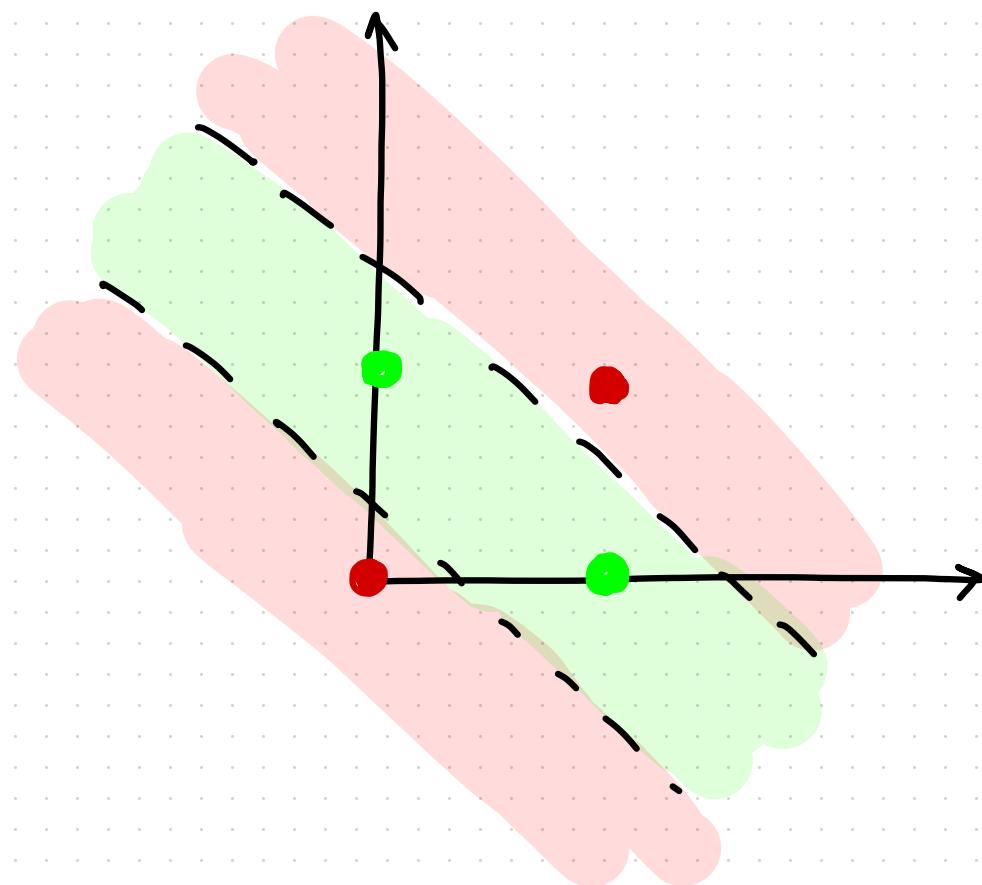


FOR $x_1 = 1 ; x_2 = 0$; we get

$$\hat{y} = \{0.5 \leq 0\} = \text{GREEN CLASS}$$

Decision surface
desired

XOR CLASSIFICATION

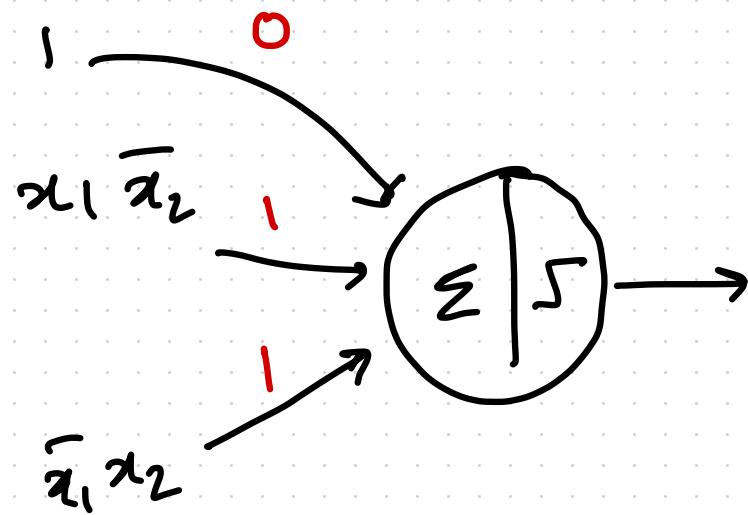
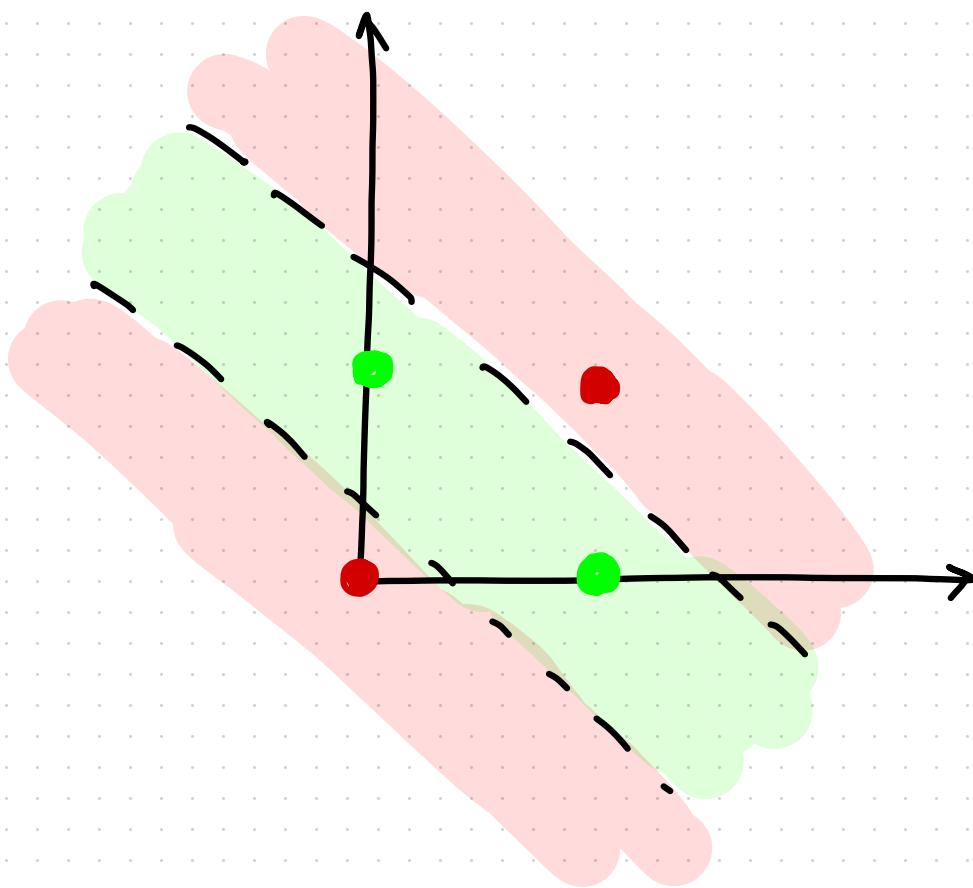


FOR $x_1 = 1 ; x_2 = 1$; we get

$$\hat{y} = \{-0.5 \leq 0\} = \text{RED CLASS}$$

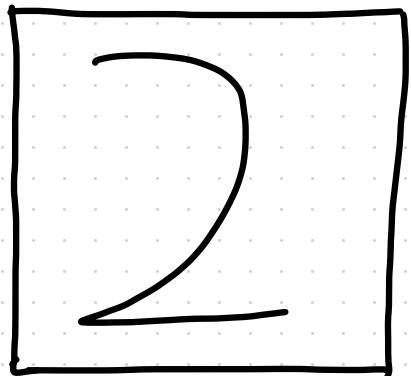
CAN ADD NON-LINEARITY
BY HAND-CRAFTING
FEATURES!

XOR CLASSIFICATION

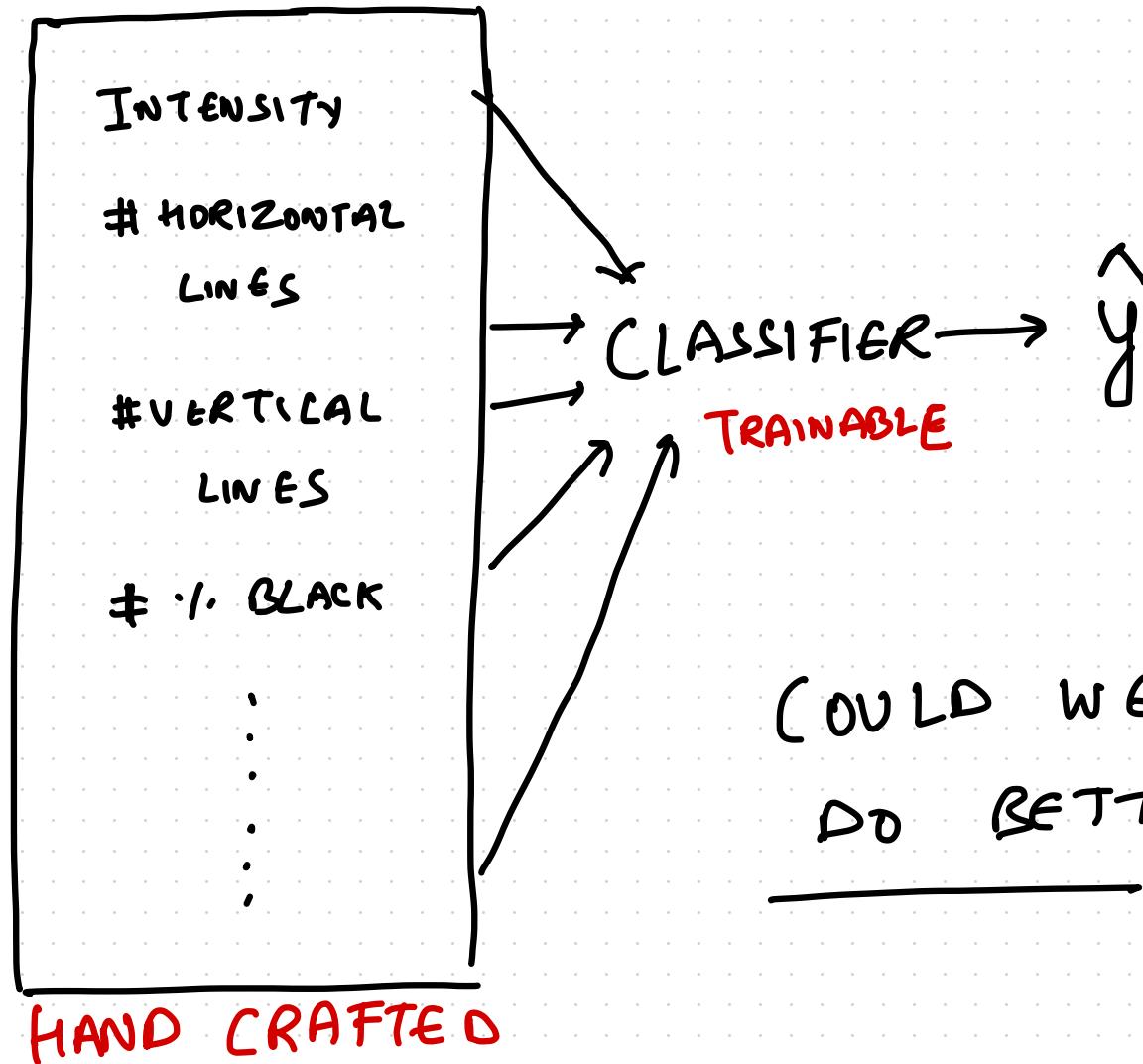


CAN ADD NON-LINEARITY
BY HAND-CRAFTING
FEATURES!

PARADIGM CHANGE



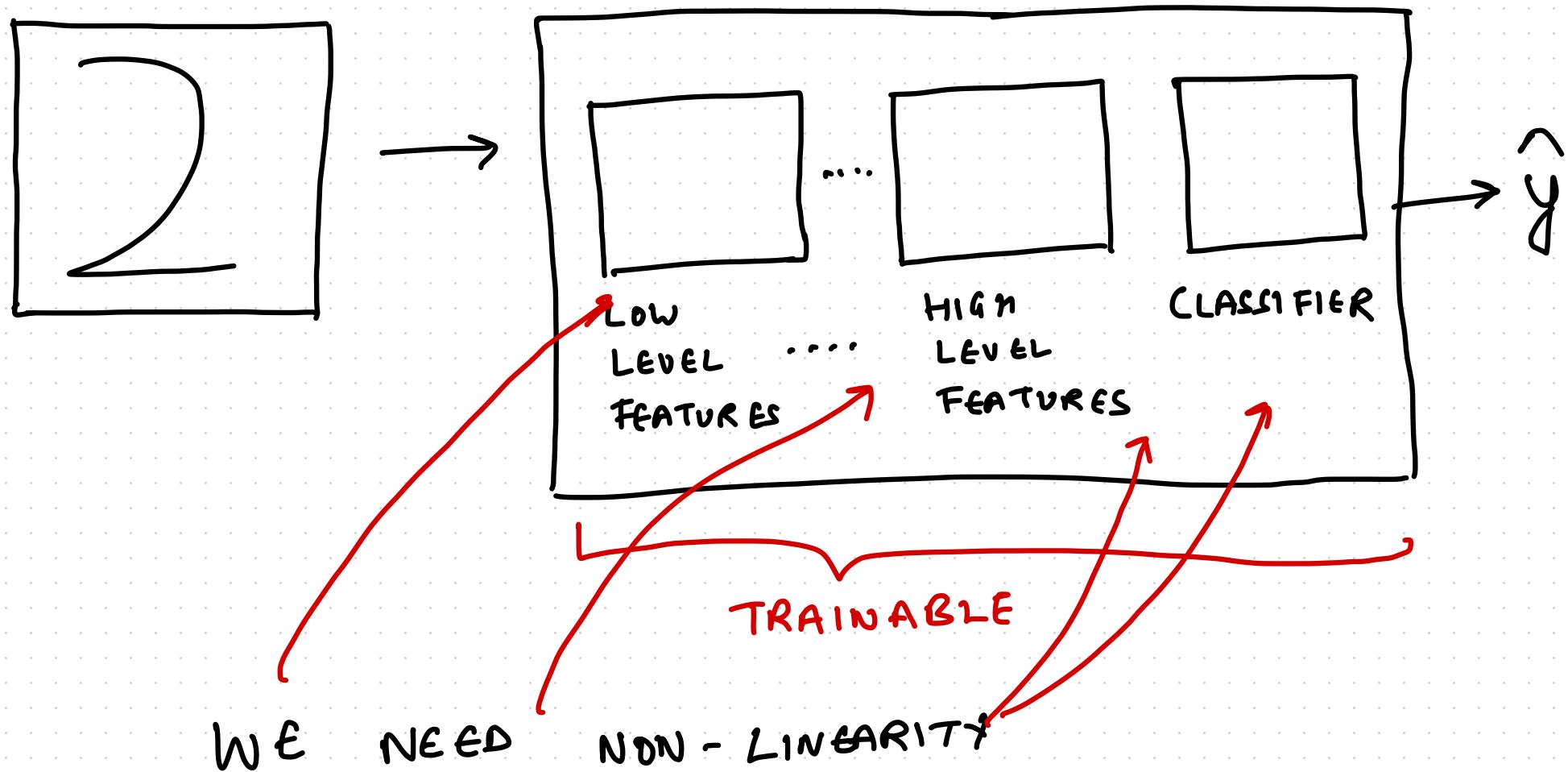
FEATURE
EXTRACTOR



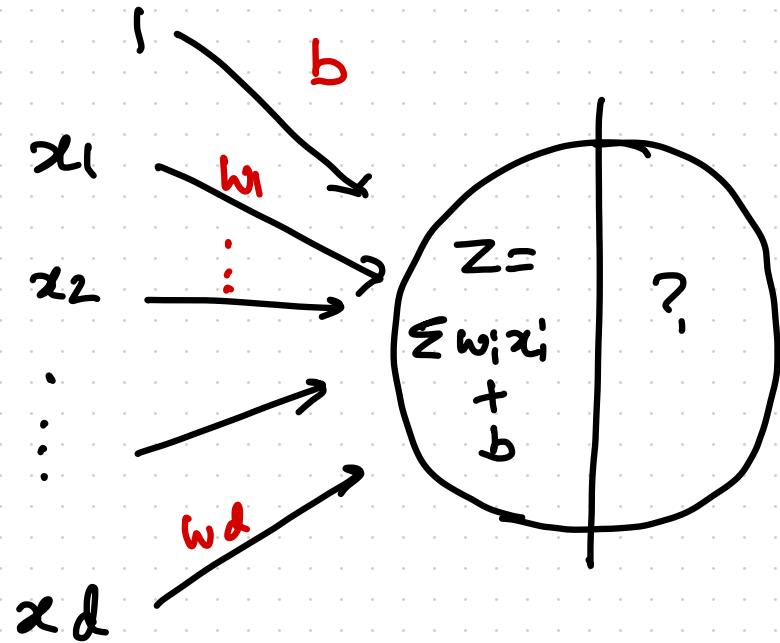
COULD WE
DO BETTER?

PARADIGM

CHANGE (NNs)

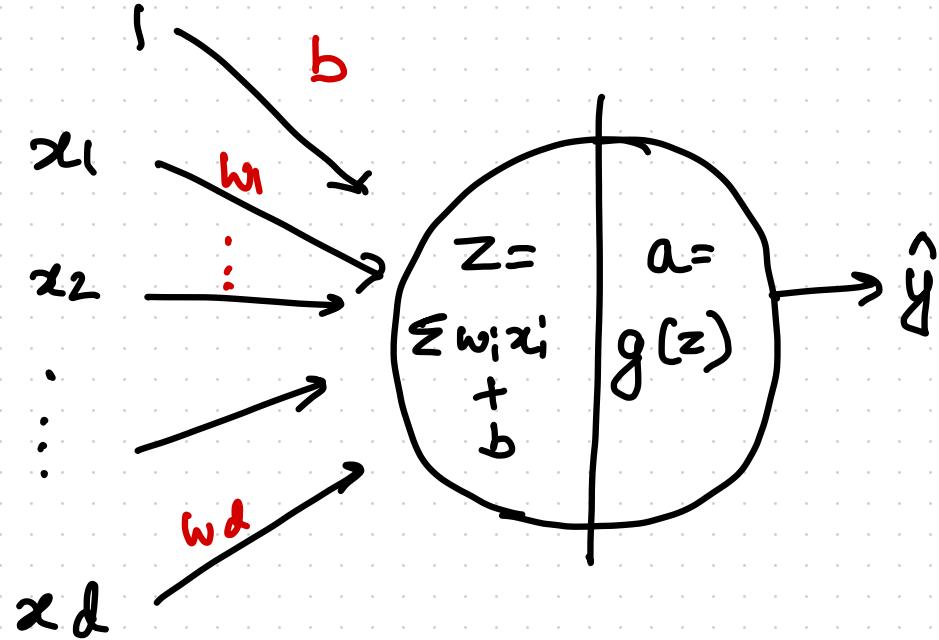


BACK PROPAGATION SUPPORTED ACTIVATIONS



key idea: use
activation
similar to
but differentiable

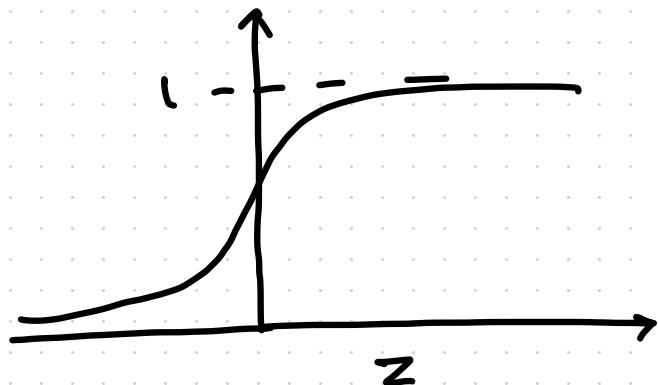
ADDING NON-LINEARITY



$g(z)$: NON-LINEAR
TRANSFORMATION

ACTIVATION FUNCTIONS

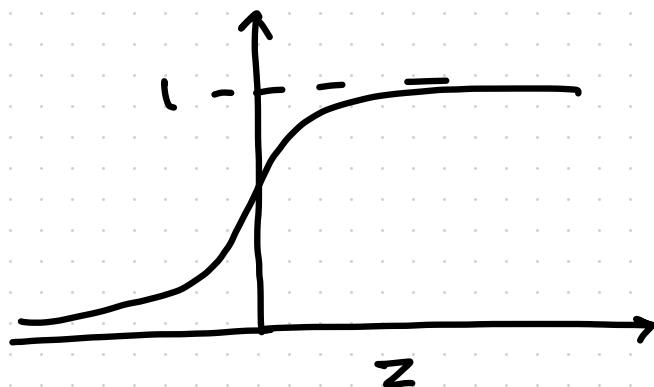
SIGMOID



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

ACTIVATION FUNCTIONS

SIGMOID



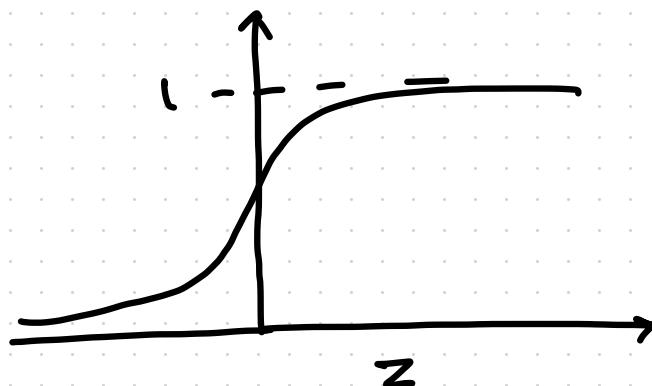
$$g(z) = \frac{1}{1+e^{-z}}$$

Q): If we have 1 neuron

$g(z) = \frac{1}{1+e^{-z}}$ what do we
get?

ACTIVATION FUNCTIONS

SIGMOID



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

Q): If we have 1 neuron

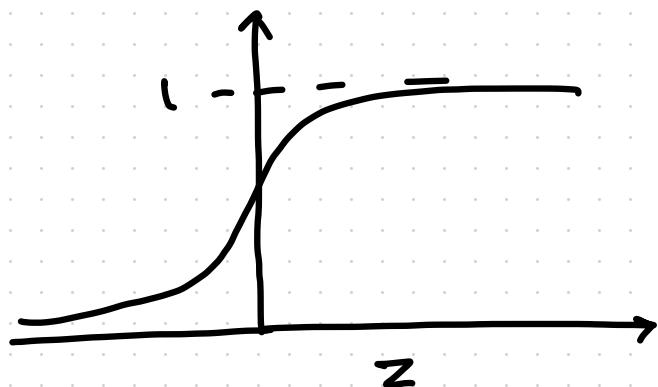
$\xrightarrow{?}$

$$g(z) = \frac{1}{1+e^{-z}} \quad \text{what do we get?}$$

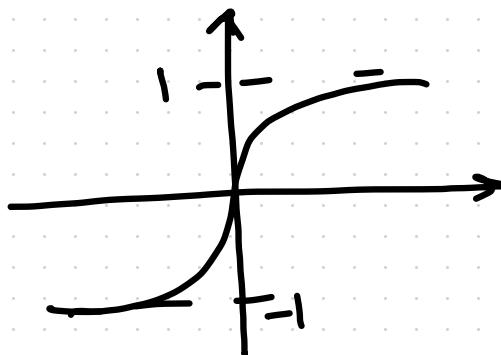
Logistic Regression

ACTIVATION FUNCTIONS

SIGMOID



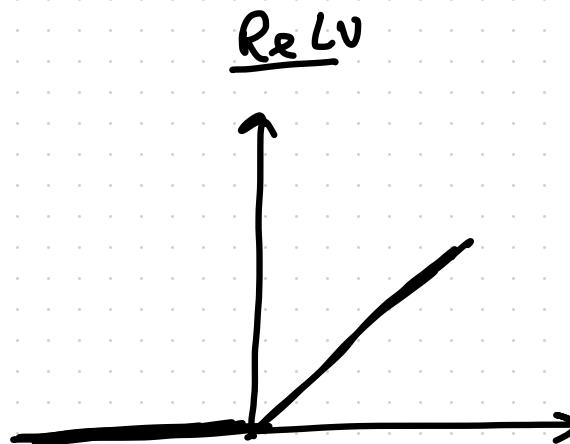
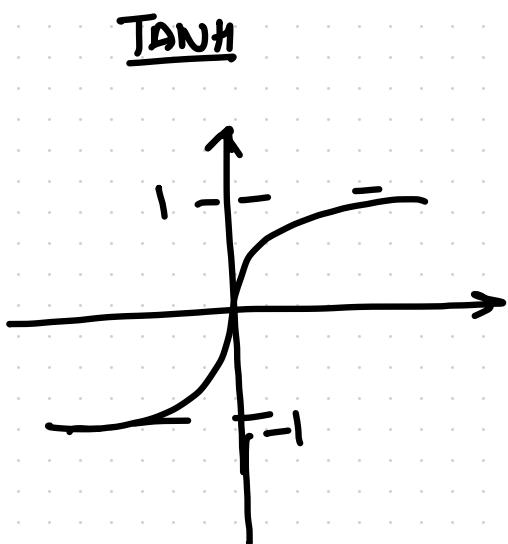
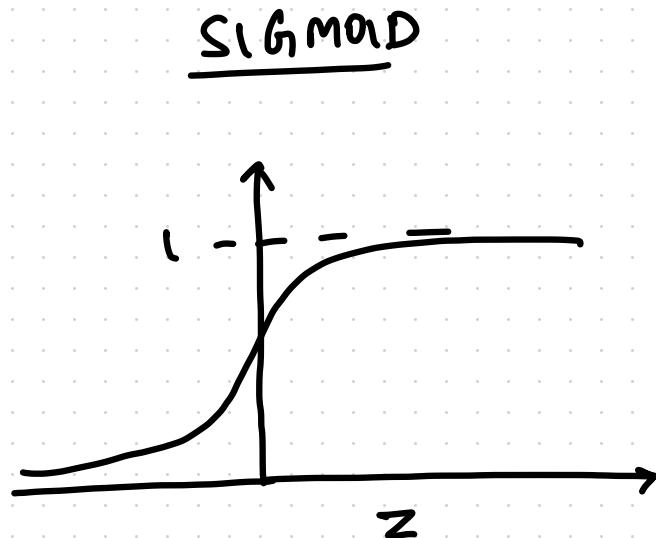
TANH



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

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

ACTIVATION FUNCTIONS



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

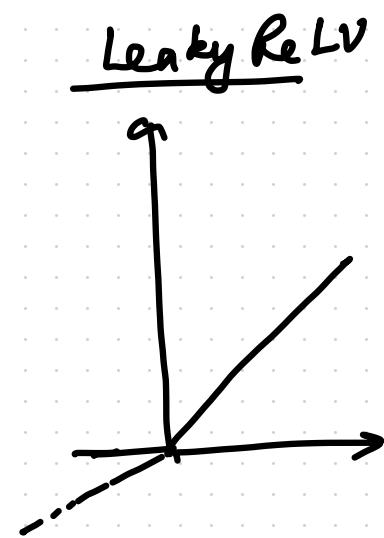
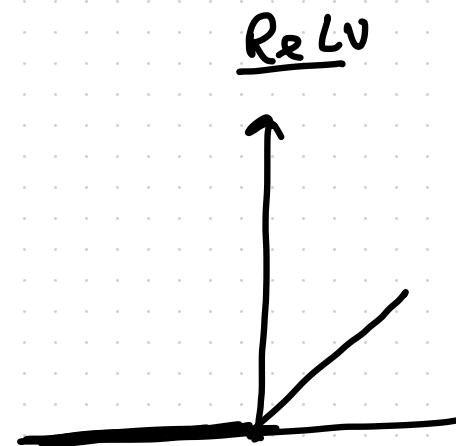
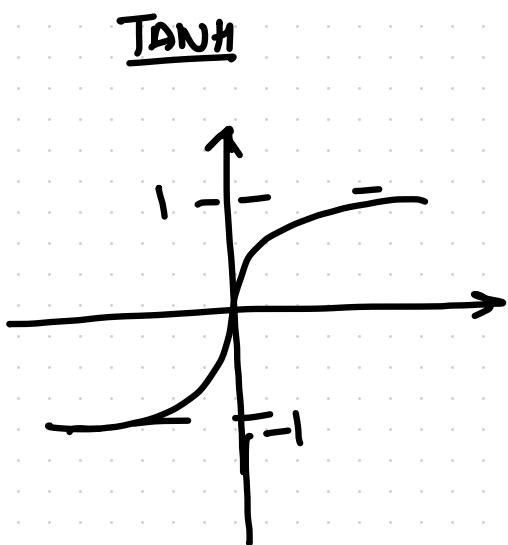
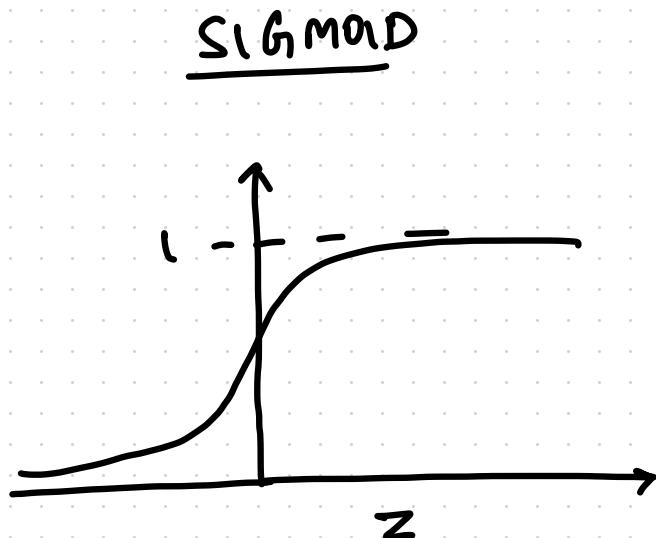
$$g(z) = \frac{e^z - e^{-z}}{e^z + e^{-z}}$$

$$g(z) = \begin{cases} z; & z \geq 0 \\ 0; & z < 0 \end{cases}$$

or

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

ACTIVATION FUNCTIONS



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

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

$$g(z) = \begin{cases} z; & z > 0 \\ 0; & z \leq 0 \end{cases}$$

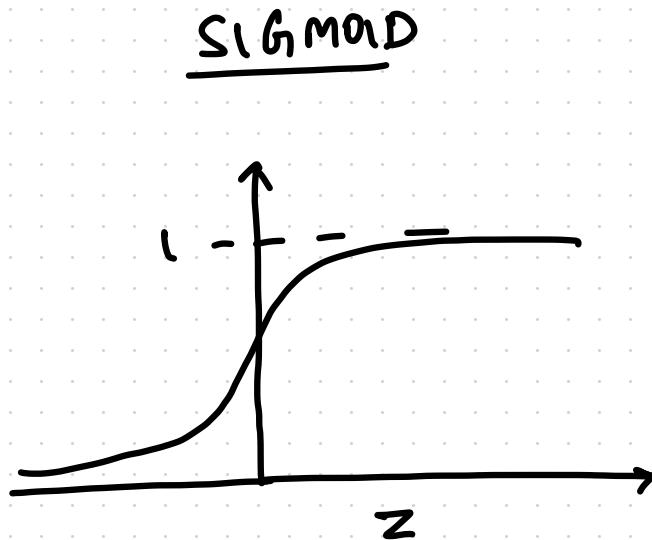
or

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

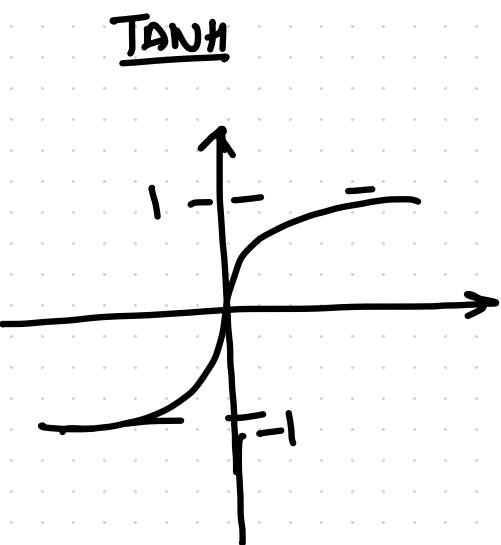
$$g(z) = \max(\alpha z, z)$$

$\alpha \rightarrow 0$

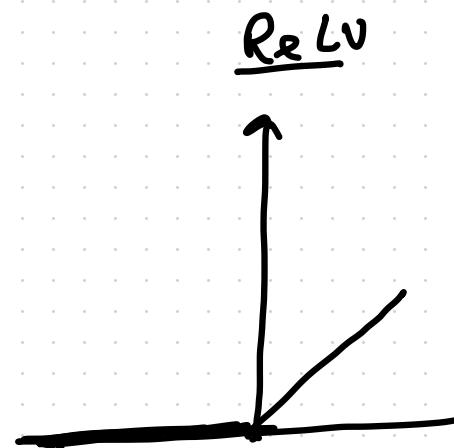
ACTIVATION FUNCTIONS



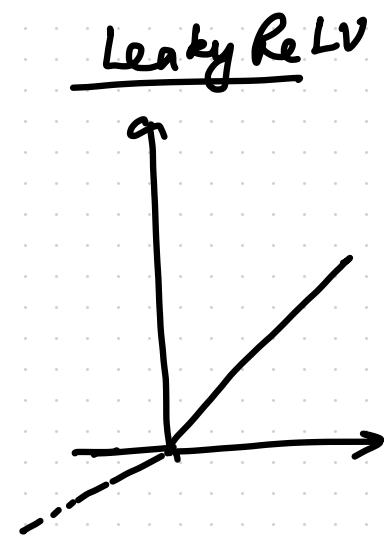
USEFUL FOR
PROBABILISTIC
ESTIMATES
 \therefore Blw 0 & 1



USEFUL IF
DATA TRANSFORMED
WITH MEAN 0



GAME
CHANGER
(Default)
Good
learning
for $|z| = \text{high}$.

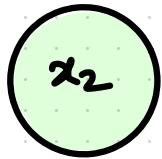


Similar
to
ReLU
Learns
for $z < 0$
also

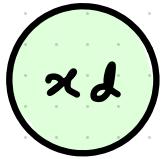
DESIRABLE ATTRIBUTES OF ACTIVATION FUNCTIONS

- 1) NON-LINEAR
- 2) (MOSTLY) SMALL CHANGE IN $I_f \Rightarrow$ SMALL CHANGE IN O_f

1 LAYER PERCEPTRON (NN)

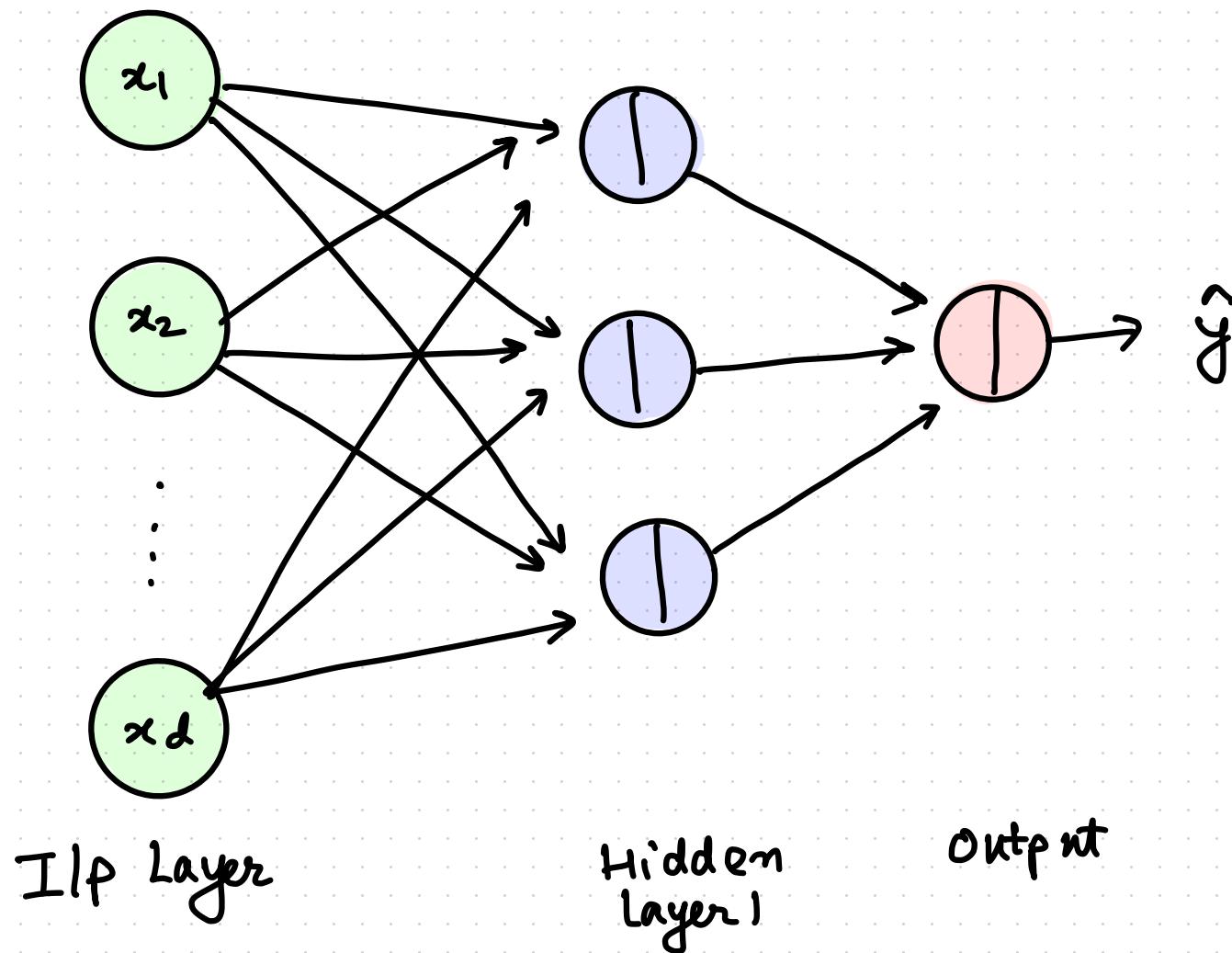


⋮
⋮

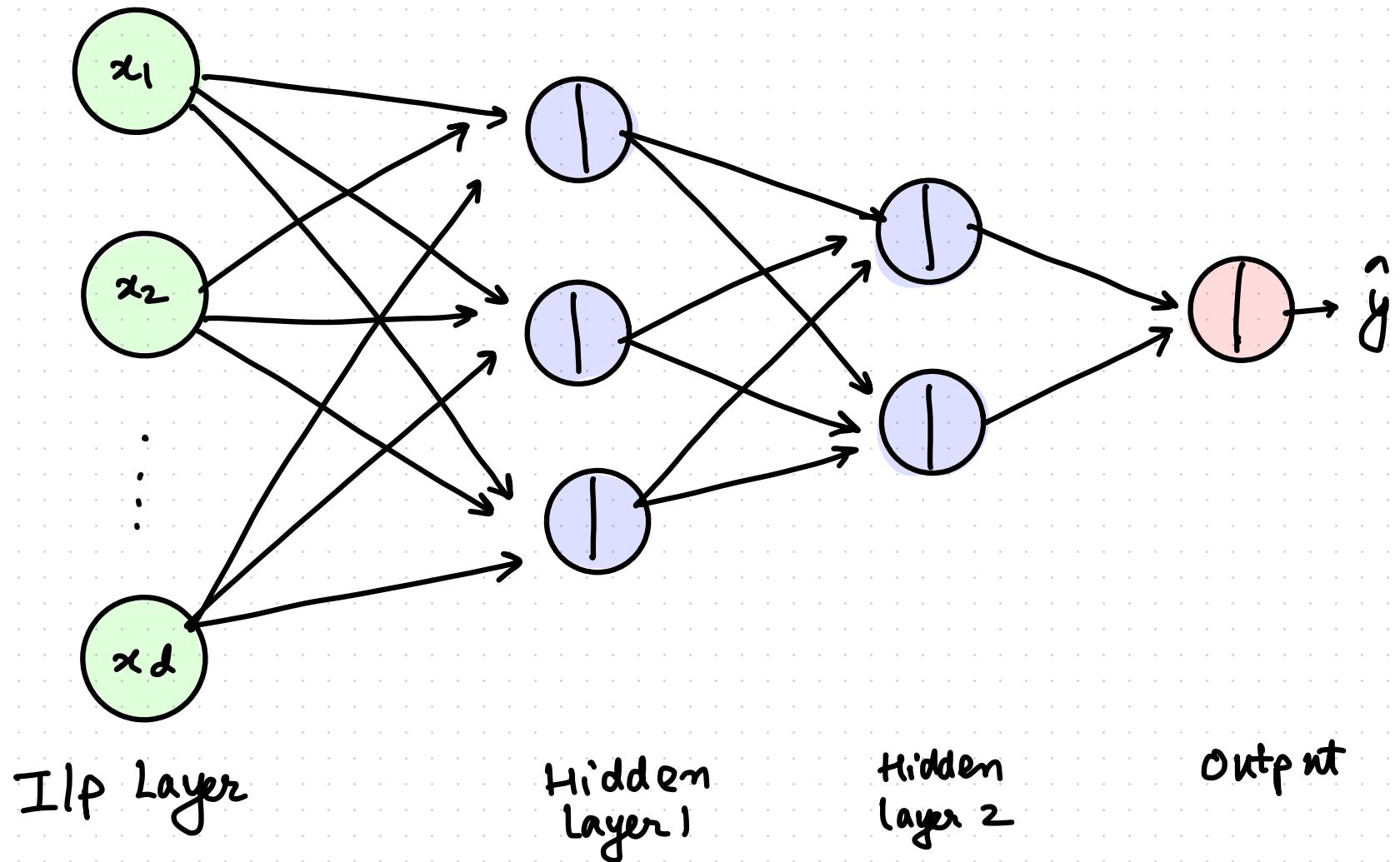


I/f Layer

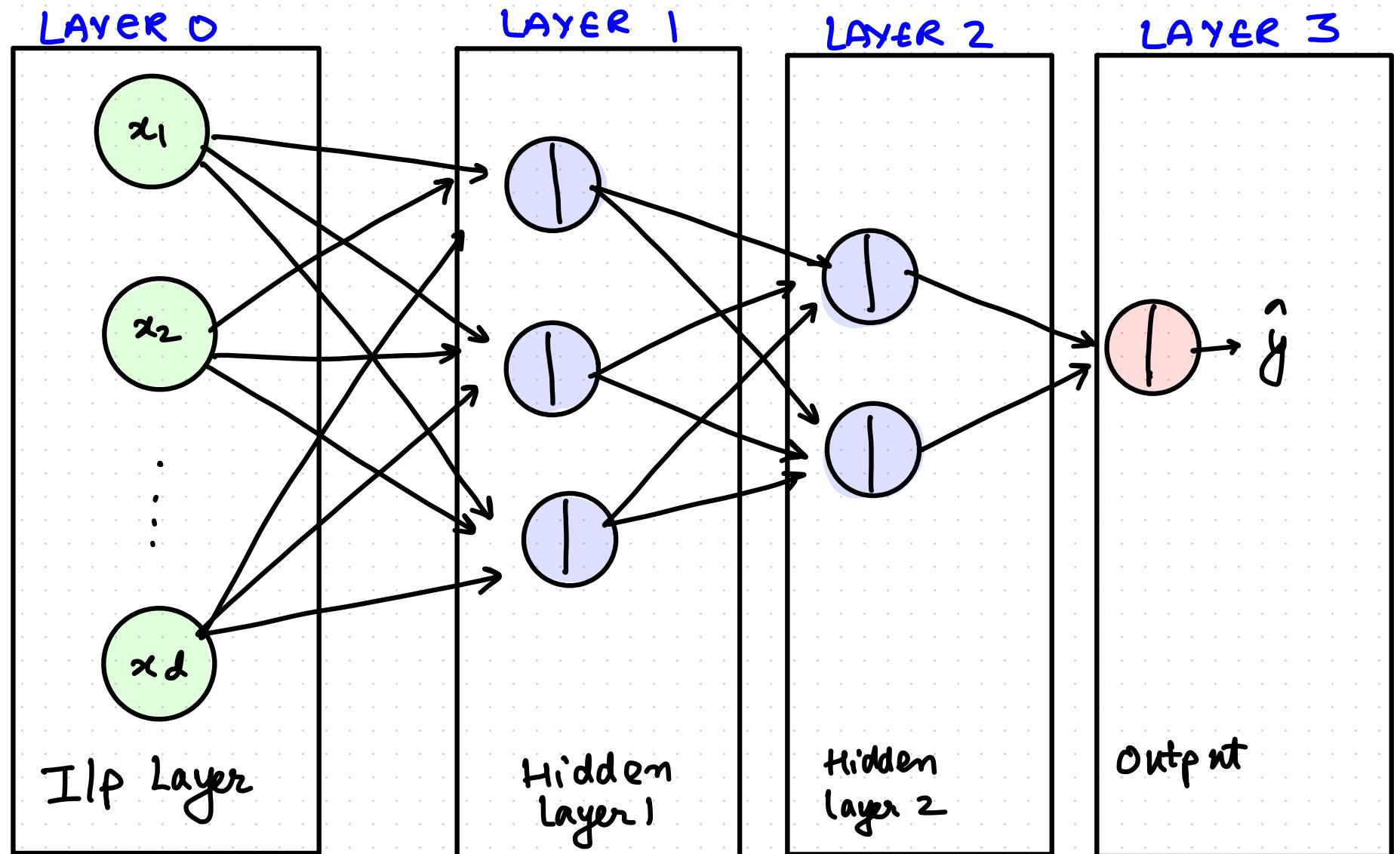
1-LAYER PERCEPTRON (NN)

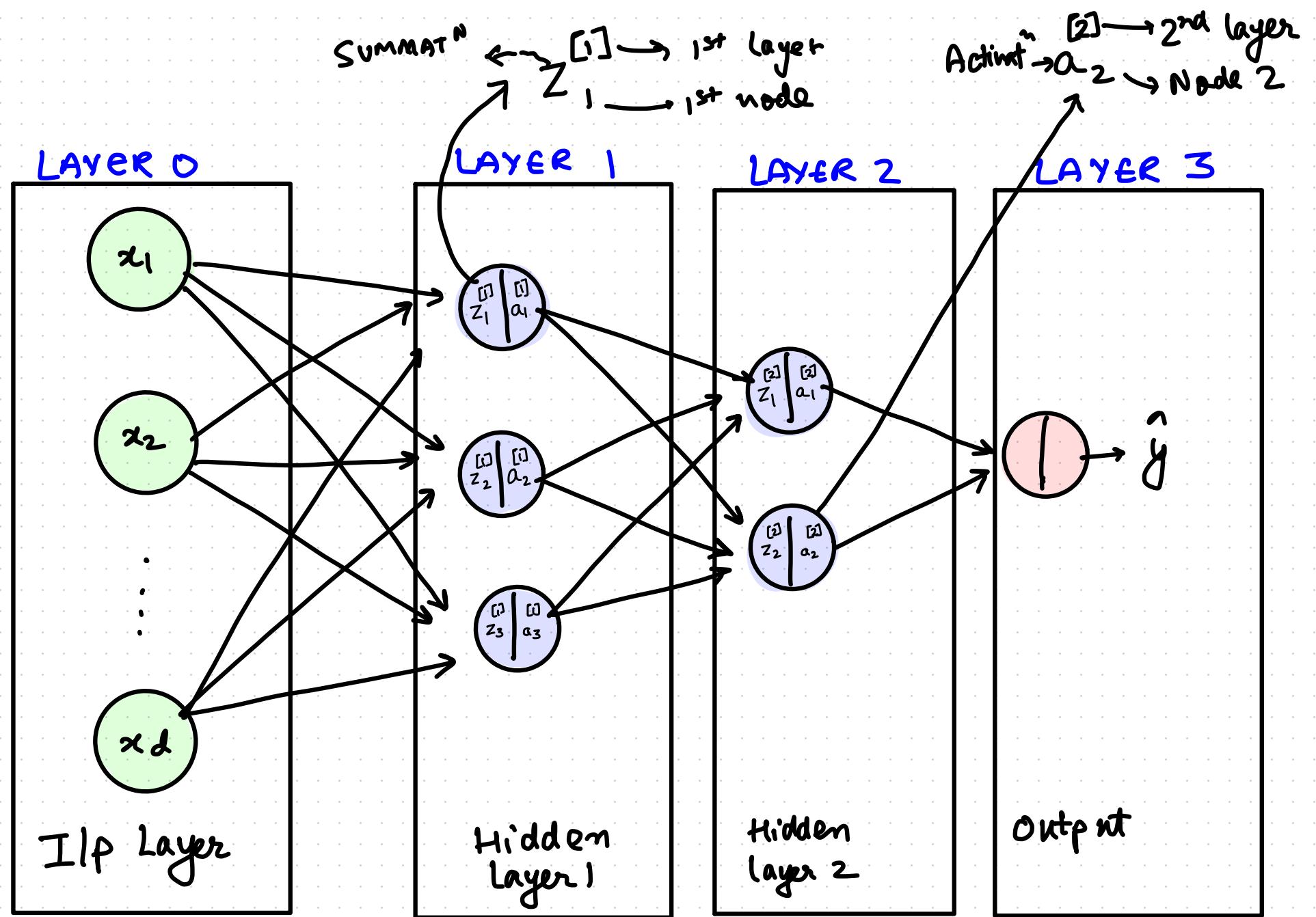


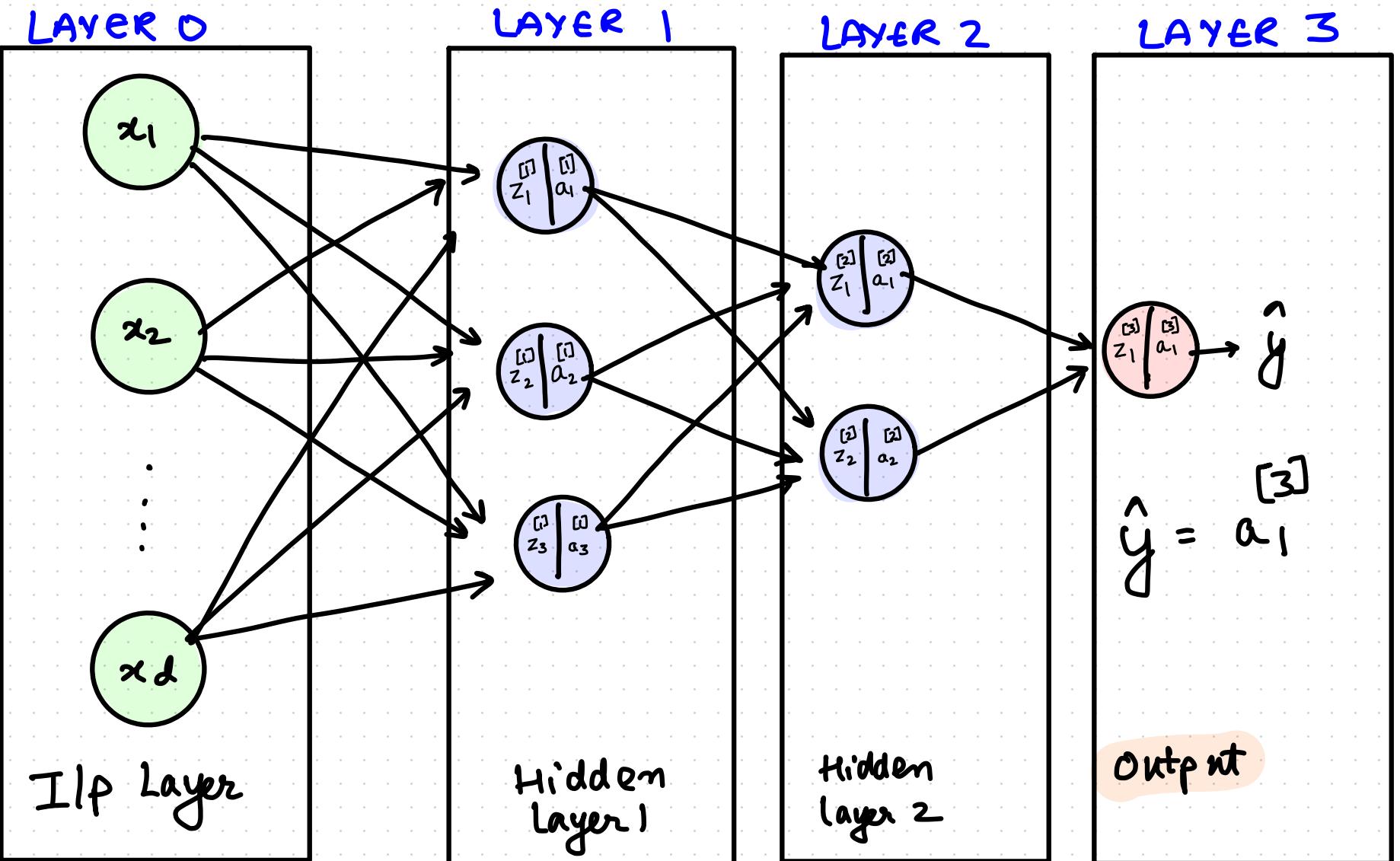
MULTI - LAYER PERCEPTRON

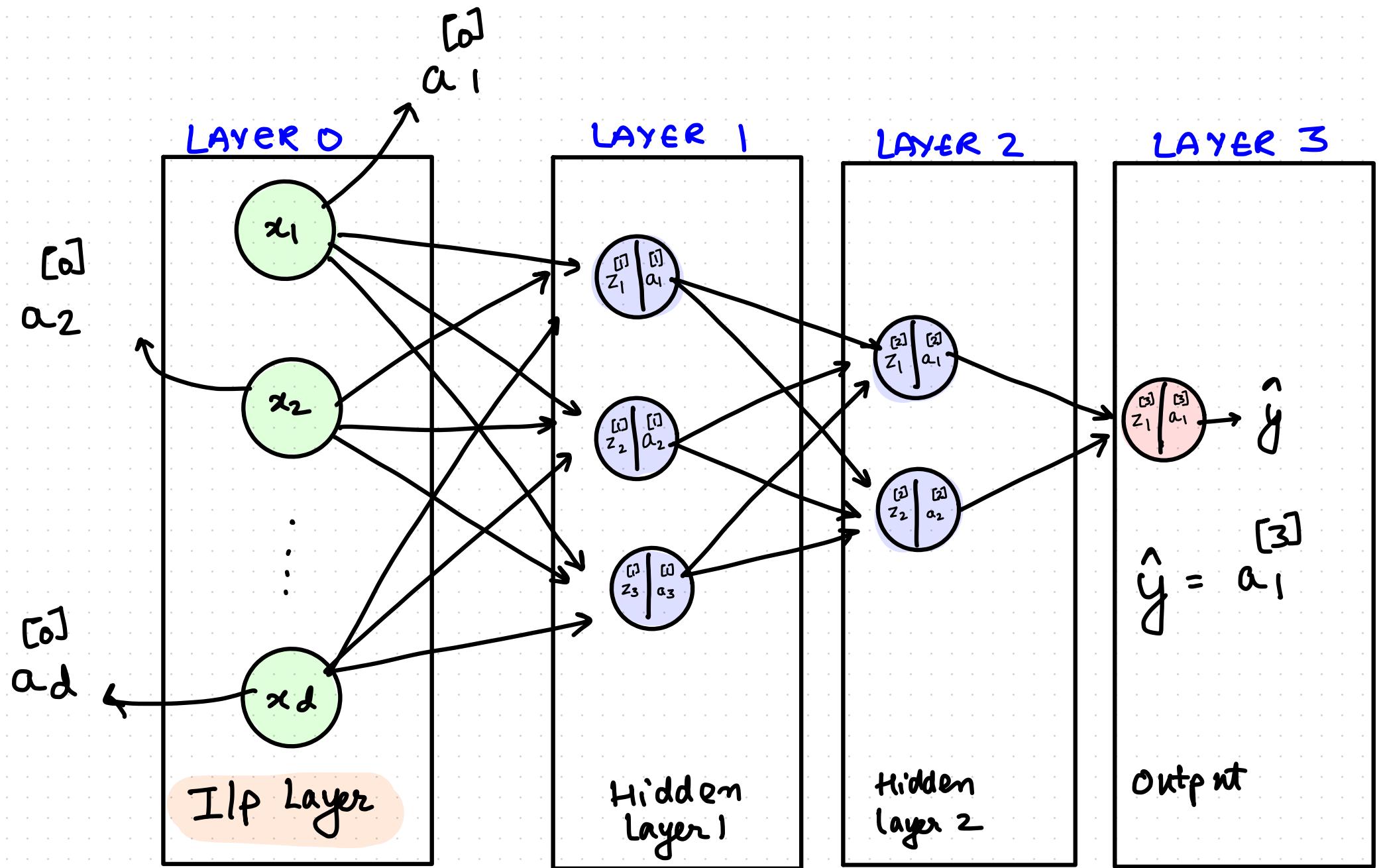


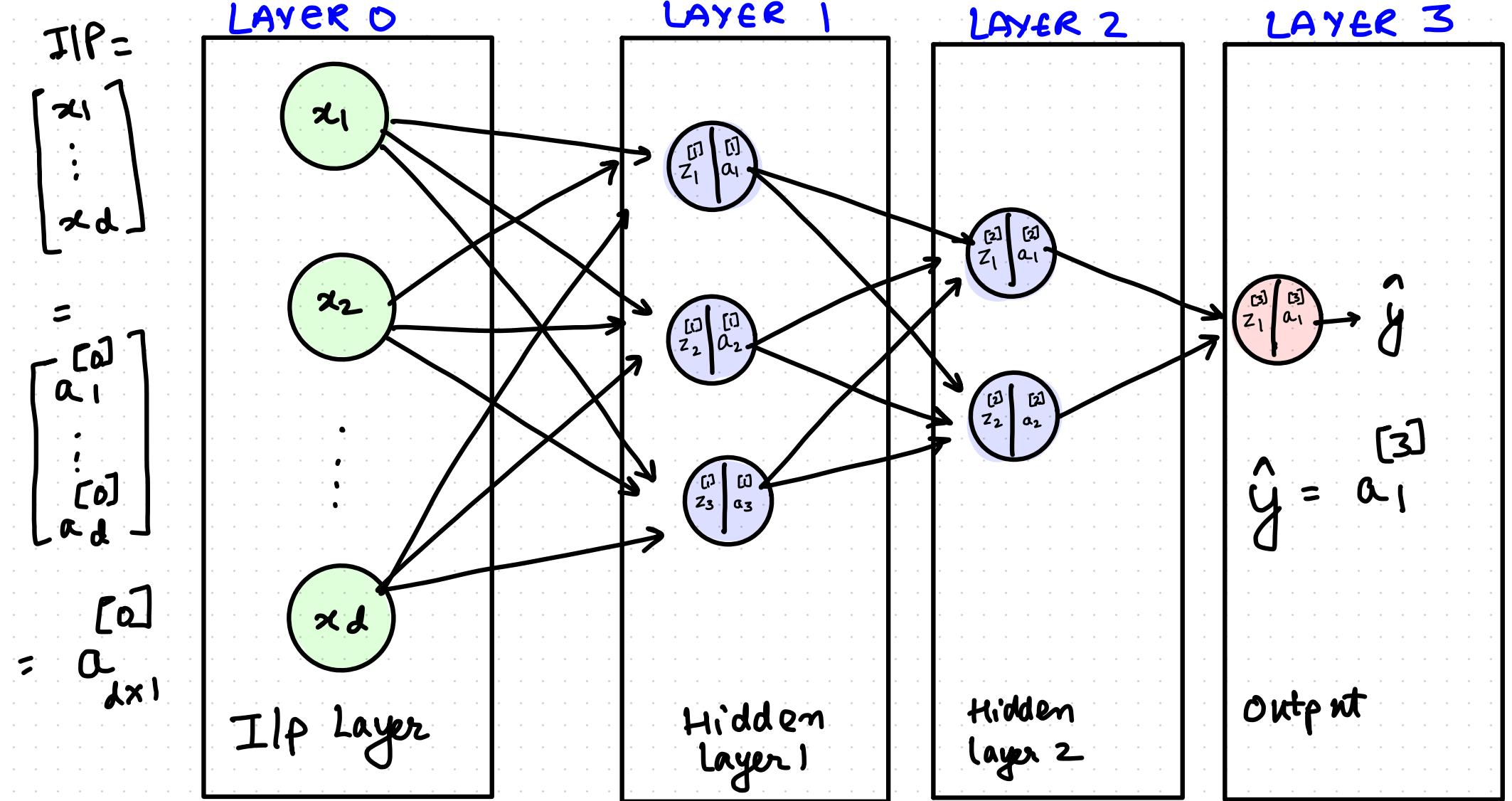
MULTI - LAYER PERCEPTRON



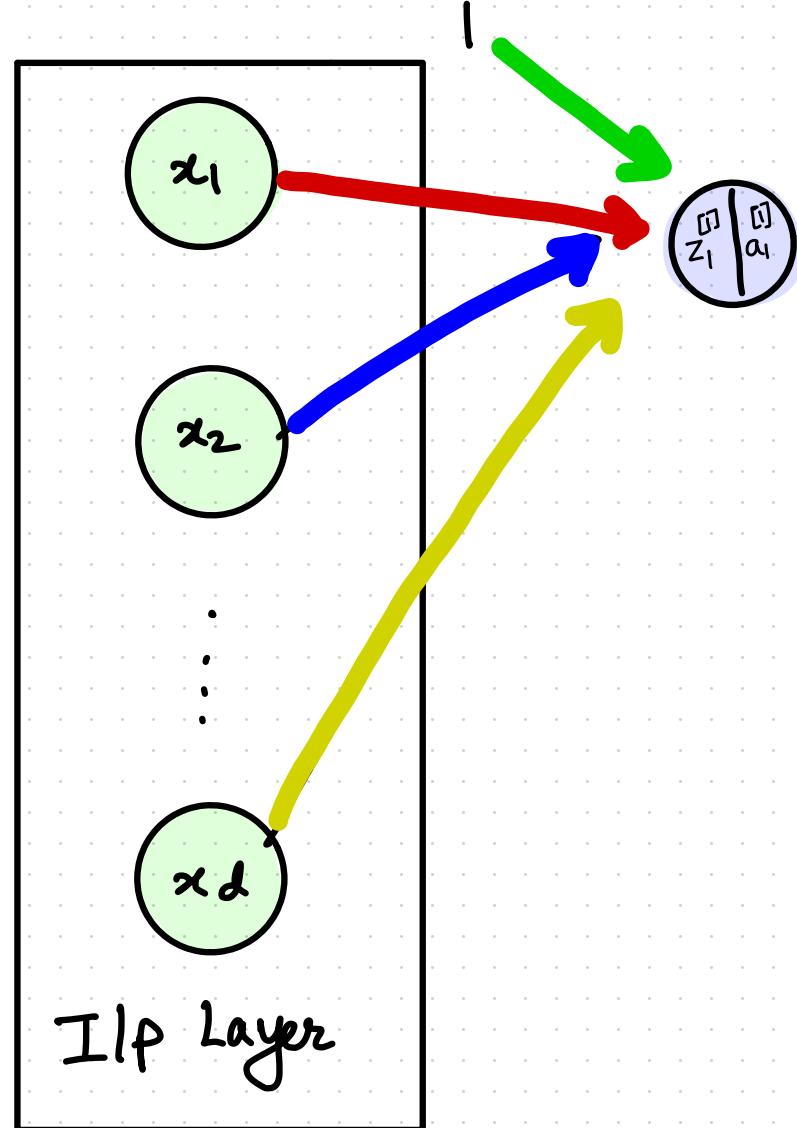








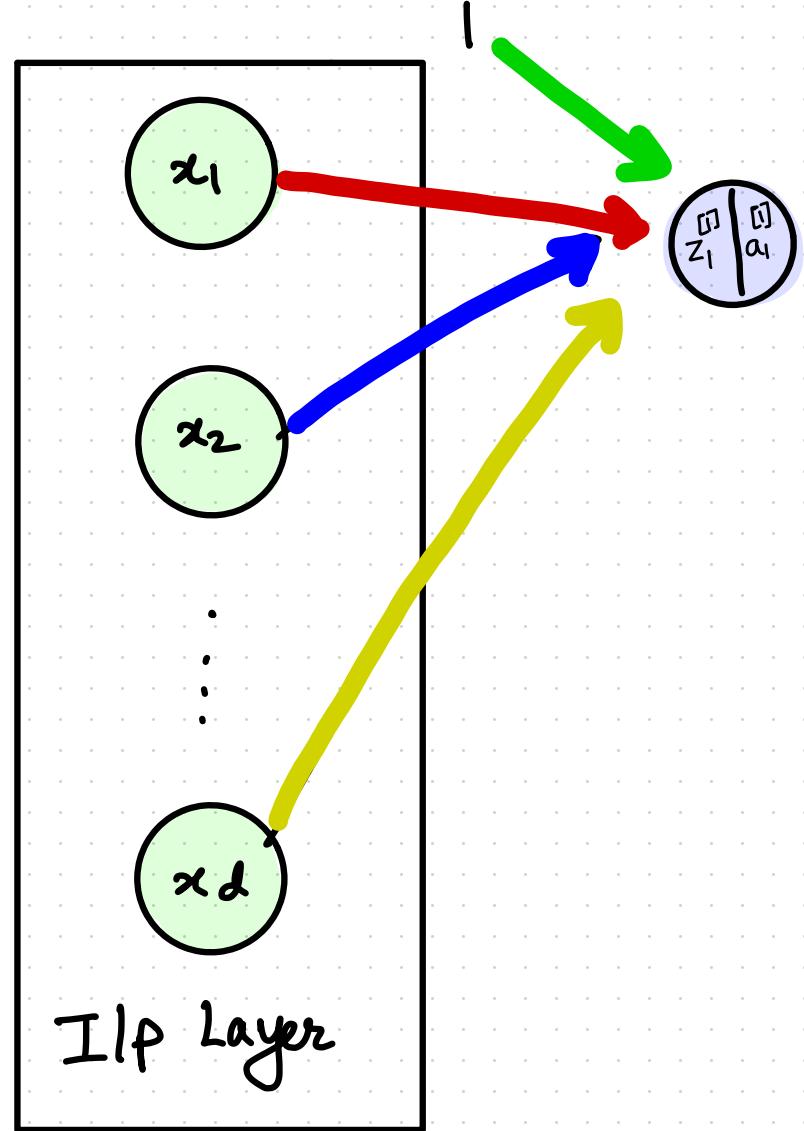
(CONSIDER SINGLE NEURON (LAYER 1, NODE 1))



$$z_1^{[1]} = b_1^{[1]} + x_1 * w_{1,1}^{[1]} + x_2 * w_{1,2}^{[1]} + \dots + x_d * w_{1,d}^{[1]}$$

bias layer
(node)

(CONSIDER SINGLE NEURON (LAYER 1, NODE 1))



$$z_1^{[1]} = l * b_1^{[1]} +$$

$$x_1 * w_{1,1}^{[1]} +$$

$$x_2 * w_{1,2}^{[1]} +$$

$$\vdots$$

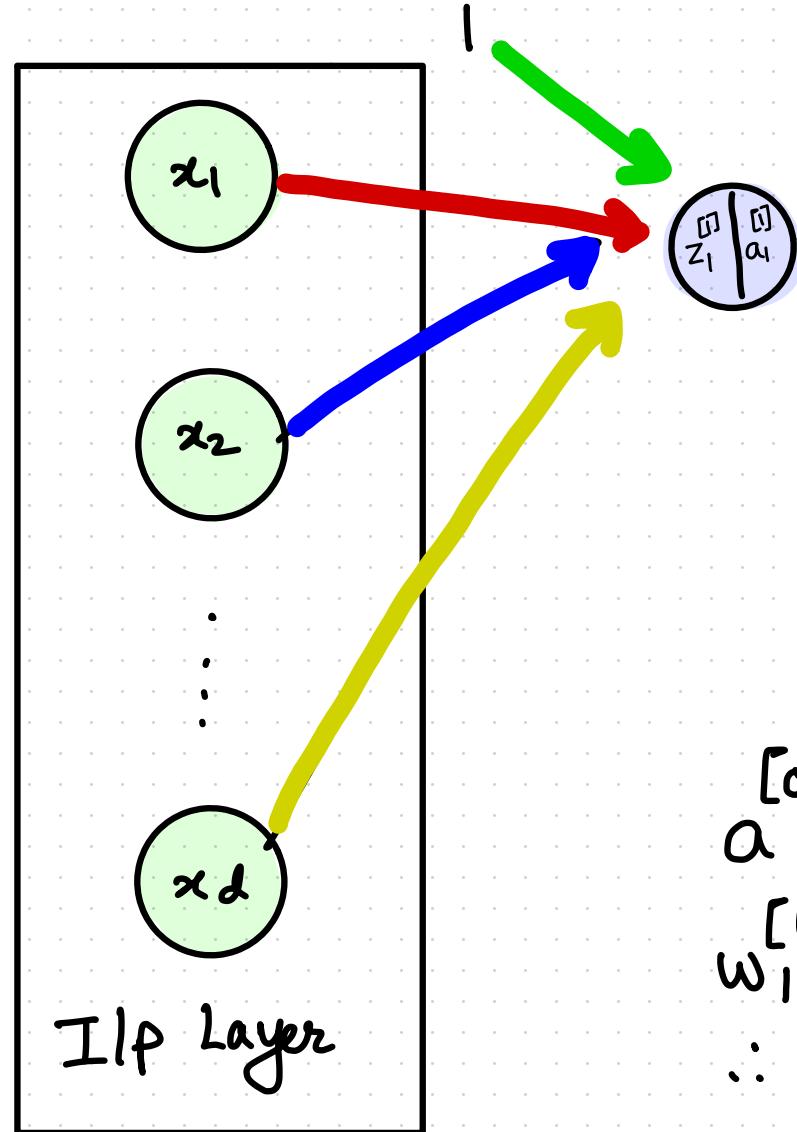
$$x_d * w_{1,d}^{[1]}$$

bias layer
node)

$w^{[l]} \leftarrow l^{\text{th}} \text{ layer}$
 $w_{a, b}^{[l]}$
 a^{th}
 node in
 l^{th}
 layer

$b^{[l]}$ component of
 prev. layer activat"

(CONSIDER SINGLE NEURON (LAYER 1, NODE 1))



$$a^{[0]} \in \mathbb{R}^D$$

$$w_1^{[i]} \in \mathbb{R}^D$$

$$\therefore z_1^{[i]} = \omega_1^{[i] T} a^{[0]} + b_1^{[i]}$$

$$z_1^{[i]} = l * b_1^{[i]} +$$

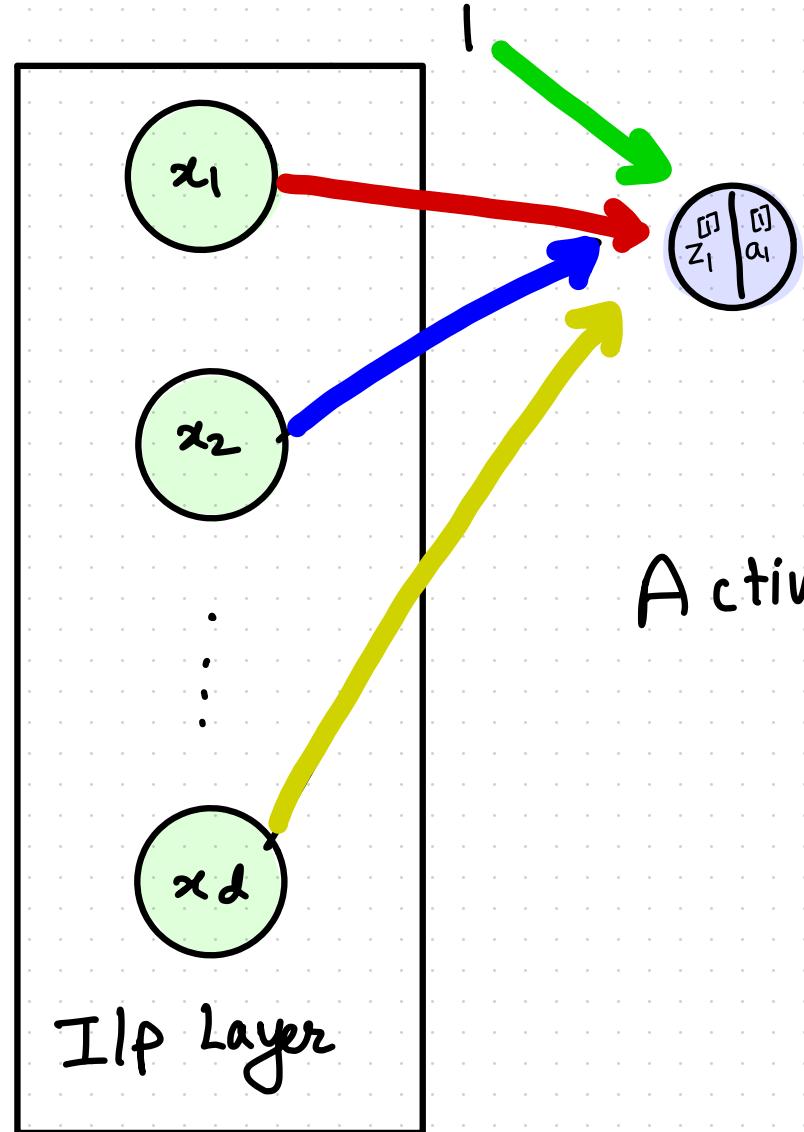
$$a_1^{[0]} * w_{1,1}^{[i]} +$$

$$a_2^{[0]} * w_{1,2}^{[i]} +$$

$$\vdots$$

$$a_d^{[0]} * w_{1,d}^{[i]}$$

CONSIDER SINGLE NEURON (LAYER 1, NODE 1)

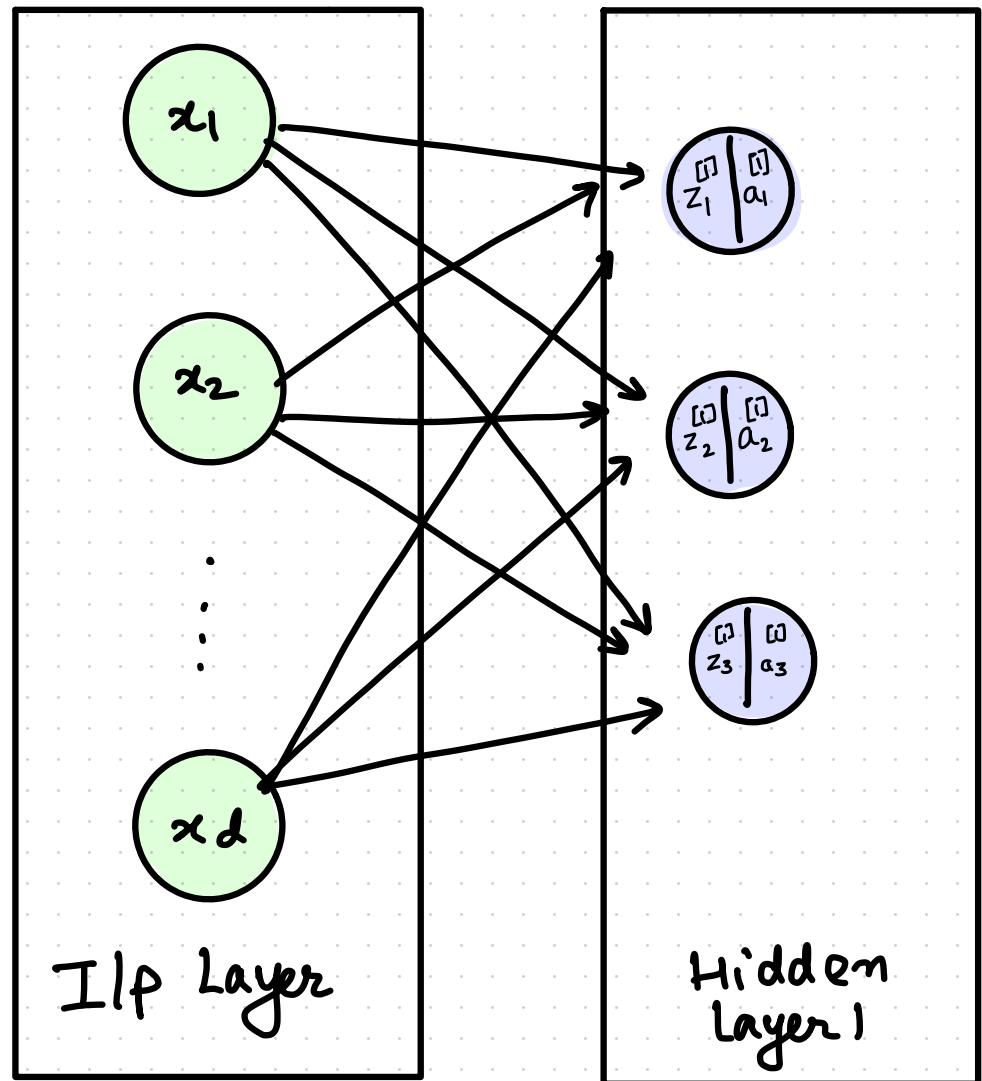


$$z_1^{[1]} = \omega_1^{[1]T} a^{[0]} + b_1^{[1]}$$

$$\text{Activation} = a_1^{[1]} = g(z_1^{[1]})$$

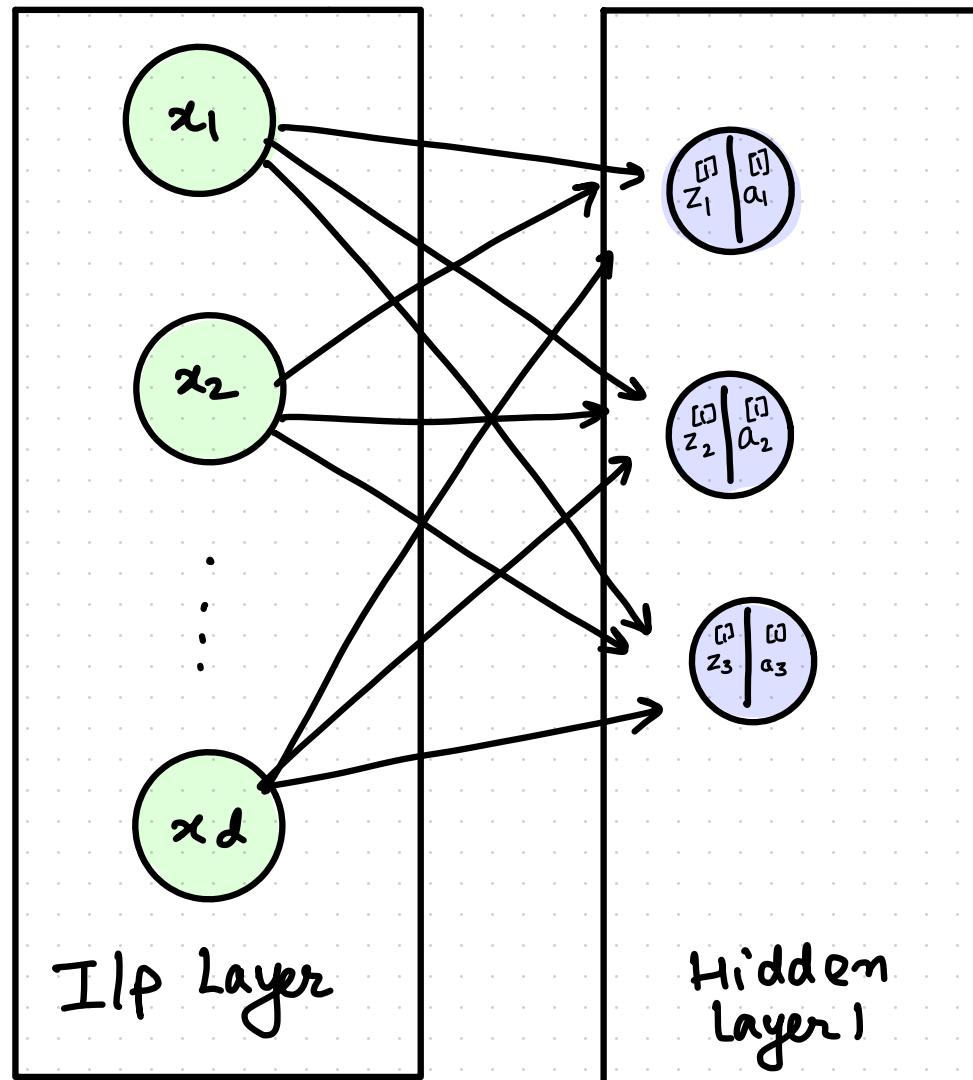
$$a_1^{[1]} \in \mathbb{R}$$

FORWARD PROPAGATION



$$a_1^{[1]} = g\left(\mathbf{w}_1^{[1] T} \mathbf{a}_{[0]}^{[0]} + b_1^{[1]}\right)$$

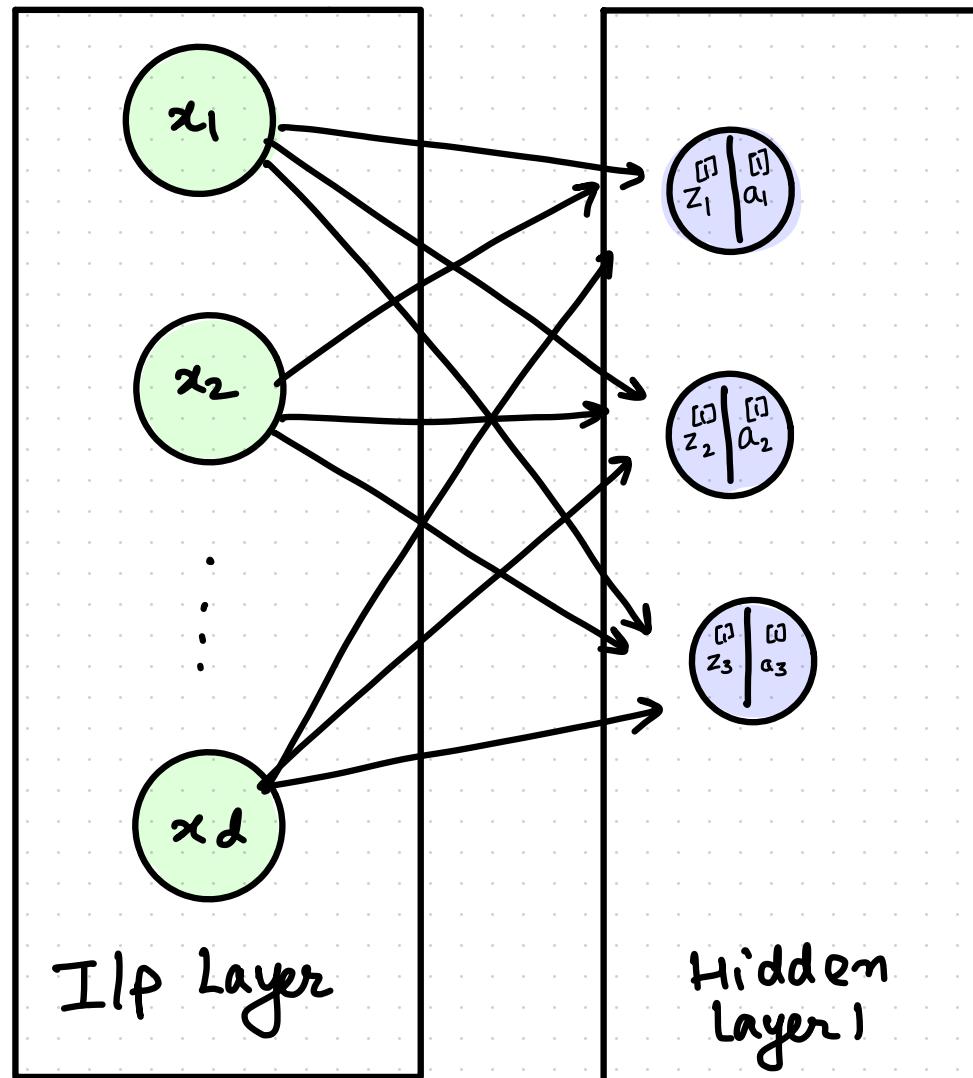
FORWARD PROPAGATION



$$a_1^{[1]} = g\left(\omega_1^{[1] T} a^{[0]} + b_1^{[1]}\right)$$

$$a_2^{[1]} = g\left(\omega_2^{[1] T} a^{[0]} + b_2^{[1]}\right)$$

FORWARD PROPAGATION

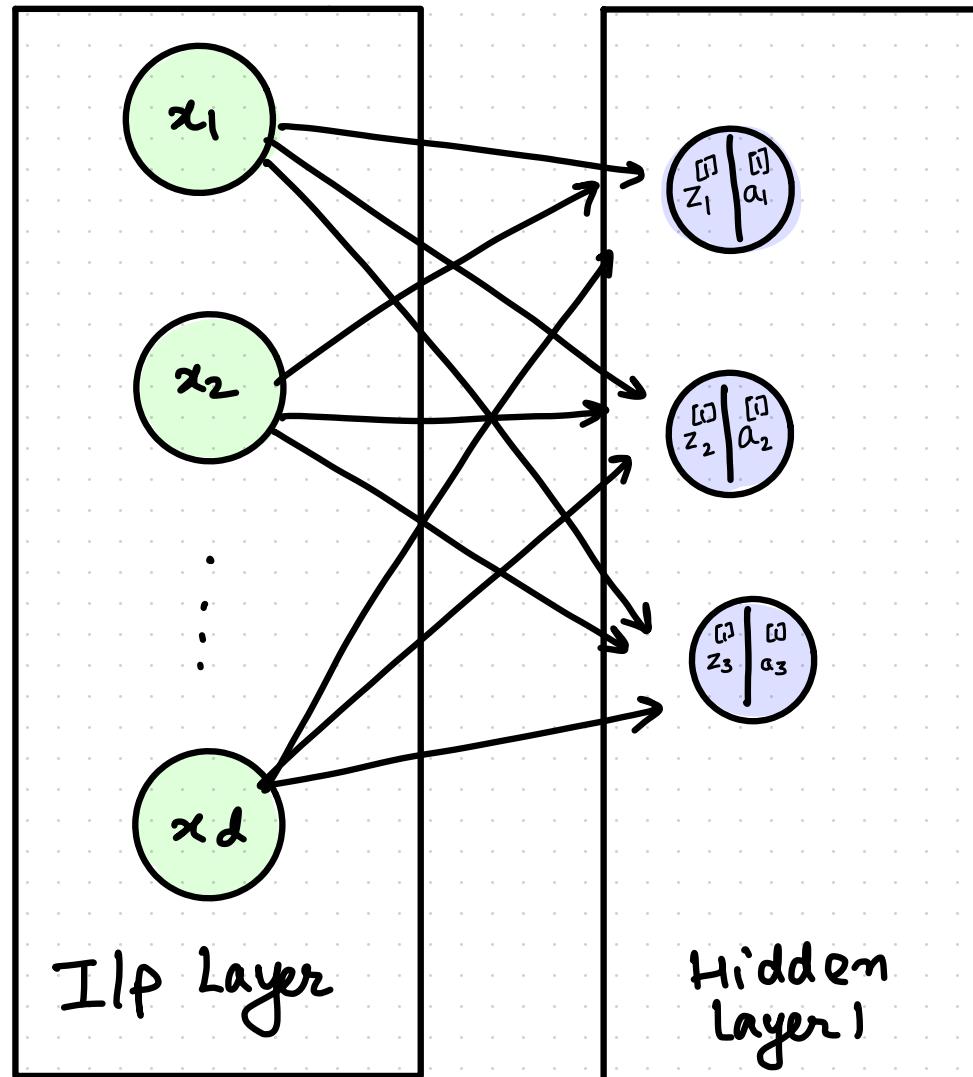


$$a_1^{[1]} = g\left(\omega_1^{[1] T} a^{[0]} + b_1^{[1]}\right)$$

$$a_2^{[1]} = g\left(\omega_2^{[1] T} a^{[0]} + b_2^{[1]}\right)$$

$$a_3^{[1]} = g\left(\omega_3^{[1] T} a^{[0]} + b_3^{[1]}\right)$$

FORWARD PROPAGATION (VECTORISATION)

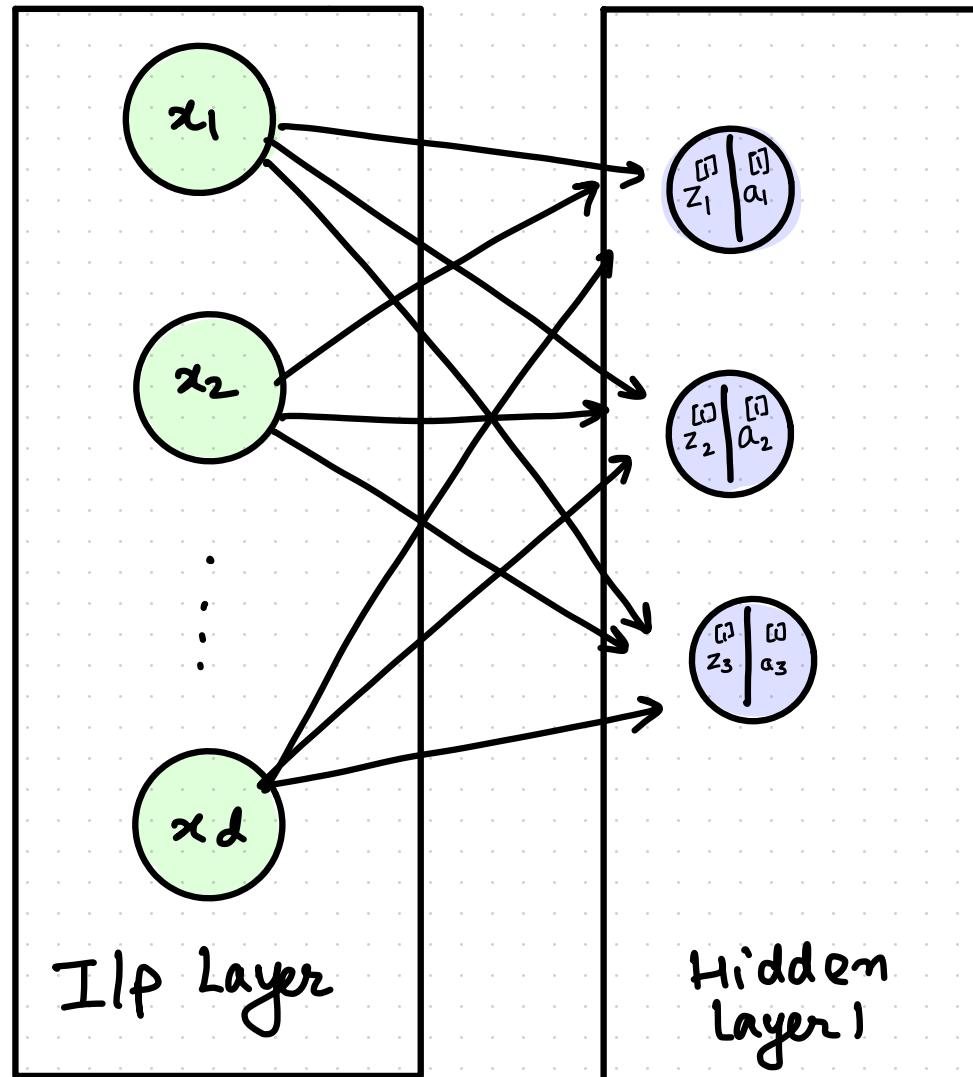


$$z_1^{[1]} = w_1^{[1]} a^{[0]} + b_1^{[1]}$$

$$z_2^{[1]} = w_2^{[1]} a^{[0]} + b_2^{[1]}$$

$$z_3^{[1]} = w_3^{[1]} a^{[0]} + b_3^{[1]}$$

FORWARD PROPAGATION (VECTORISATION)



$$z_1^{[1]} = w_1^{[1]} a^{[0]} + b_1^{[1]}$$

$1 \times 1 \quad 1 \times 3 \quad 3 \times 1 \quad 1 \times 1$

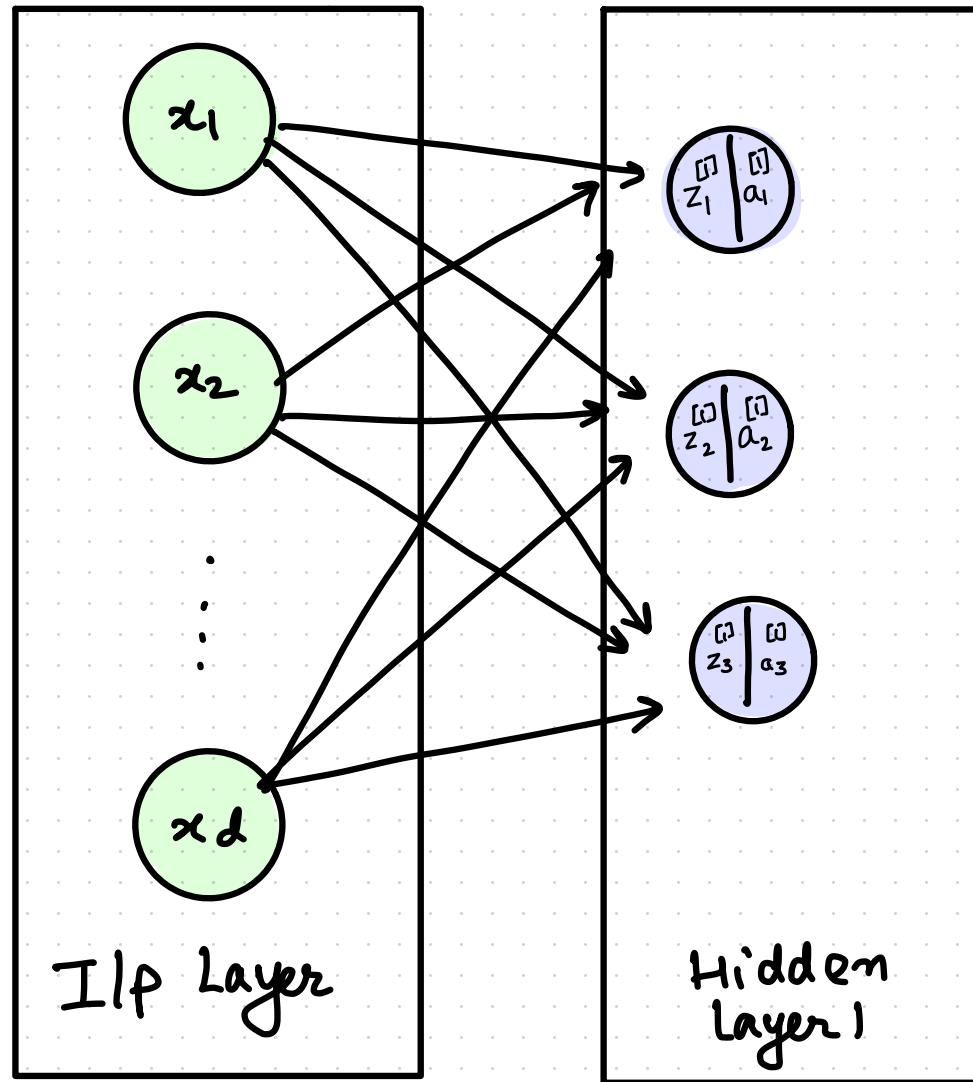
$$z_2^{[1]} = w_2^{[1]} a^{[0]} + b_2^{[1]}$$

$1 \times 1 \quad 1 \times 3 \quad 3 \times 1 \quad 1 \times 1$

$$z_3^{[1]} = w_3^{[1]} a^{[0]} + b_3^{[1]}$$

$1 \times 1 \quad 1 \times 3 \quad 3 \times 1 \quad 1 \times 1$

FORWARD PROPAGATION (VECTORISATION)



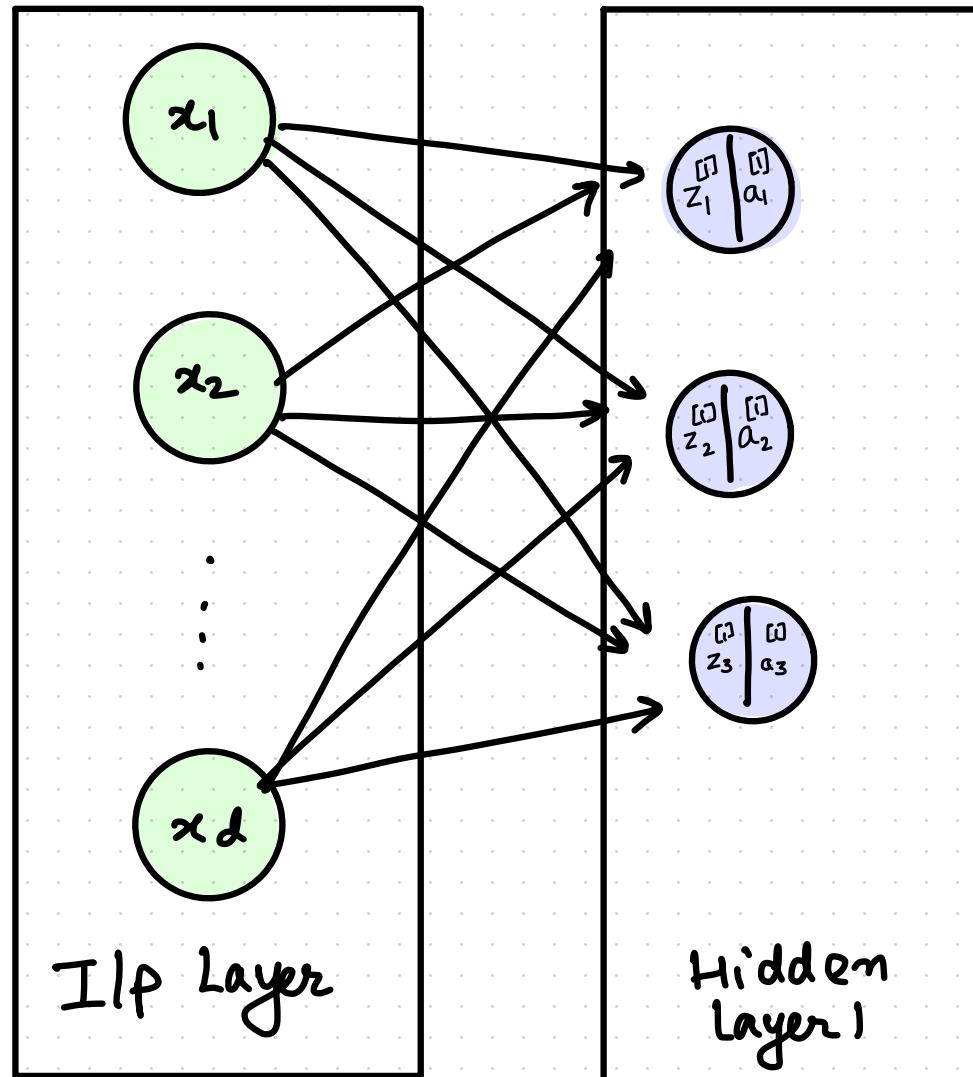
$$z_1^{[0]} = w_1^{[0] T} a^{[0]} + b_1^{[0]} \quad |x|$$

$$z_2^{[0]} = w_2^{[0] T} a^{[0]} + b_2^{[0]} \quad |n|$$

$$z_3^{[0]} = w_3^{[0] T} a^{[0]} + b_3^{[0]} \quad |x|$$

$$z^{[0]} = \begin{bmatrix} -w_1^{[0] T} \\ -w_2^{[0] T} \\ -w_3^{[0] T} \end{bmatrix} a^{[0]} + \begin{bmatrix} b_1^{[0]} \\ b_2^{[0]} \\ b_3^{[0]} \end{bmatrix}$$

FORWARD PROPAGATION (VECTORISATION)

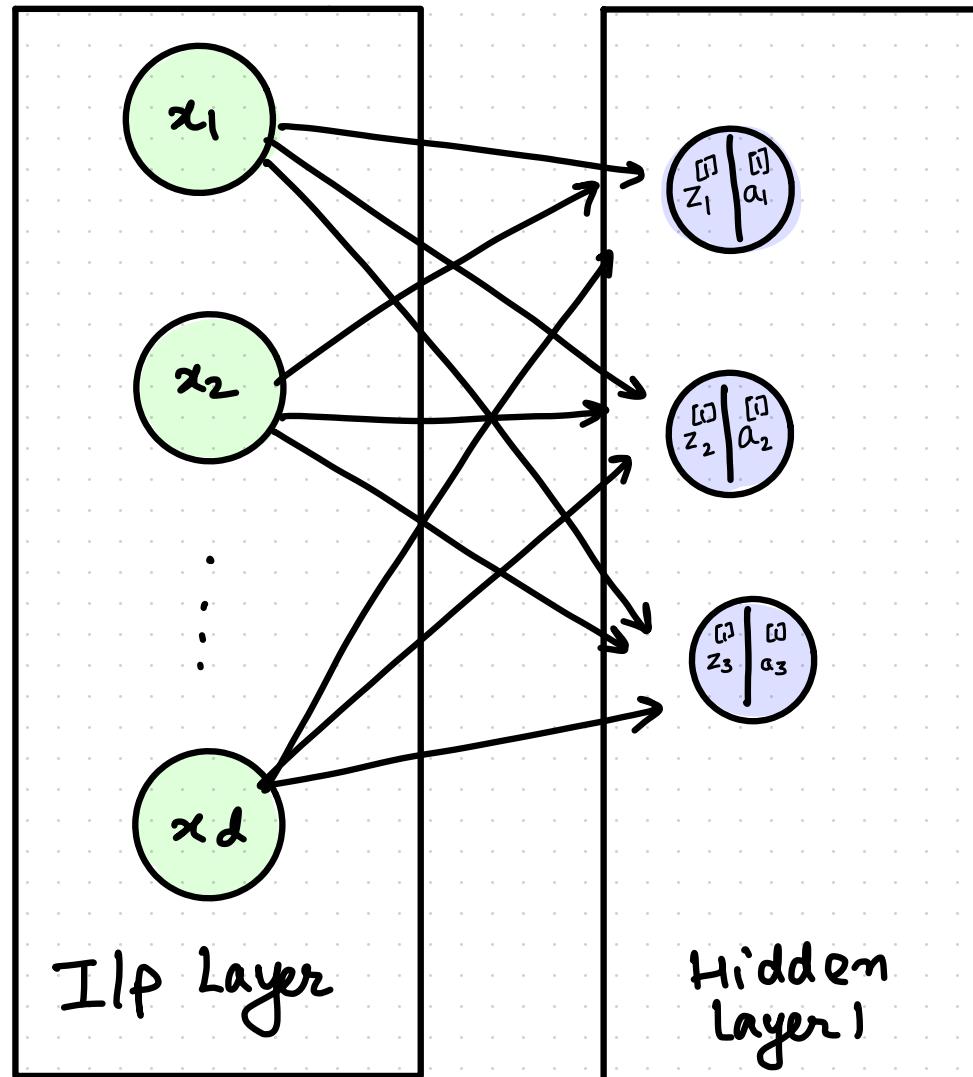


$$z^{[1]}_{3 \times 1} = \begin{bmatrix} -w_1^{[1]T} \\ -w_2^{[1]T} \\ -w_3^{[1]T} \end{bmatrix} a^{[0]}_{3 \times 1} + \begin{bmatrix} b_1 \\ b_2 \\ b_3 \end{bmatrix}_{3 \times 1}$$

$$z^{[1]} = W^{[1]} a^{[0]} + b^{[1]}$$

↑ capitals for matrices

FORWARD PROPAGATION (VECTORISATION)



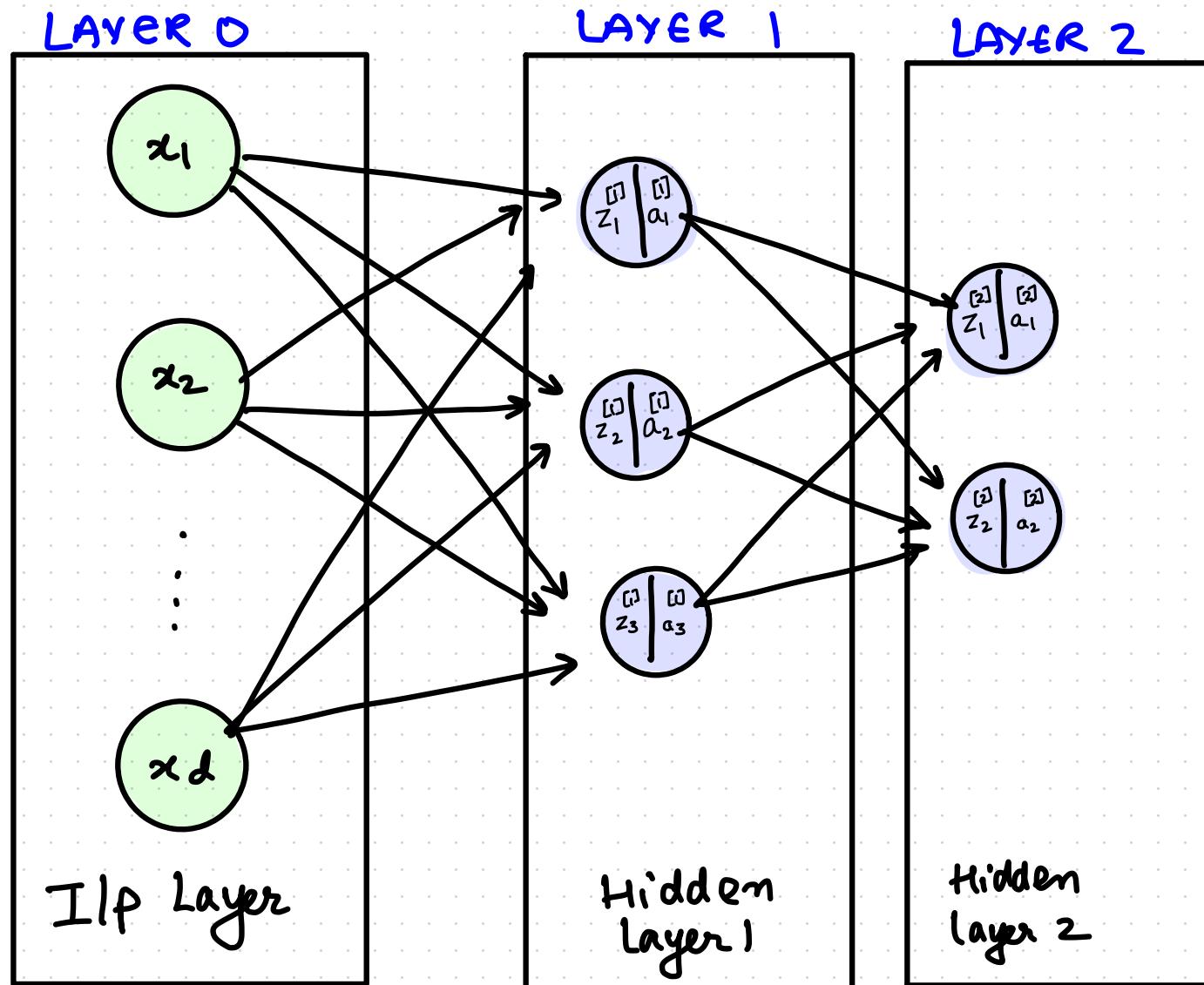
$$z^{[1]}_{3 \times 1} = \begin{bmatrix} -\omega_1^{[1]T} \\ -\omega_2^{[1]T} \\ -\omega_3^{[1]T} \end{bmatrix} a^{[0]}_{3 \times 1} + \begin{bmatrix} b_1 \\ b_2 \\ b_3 \end{bmatrix}_{3 \times 1}$$

$$z^{[1]} = W^{[1]} a^{[0]} + b^{[1]}$$

↑ capitals for matrices

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

FORWARD PROPAGATION

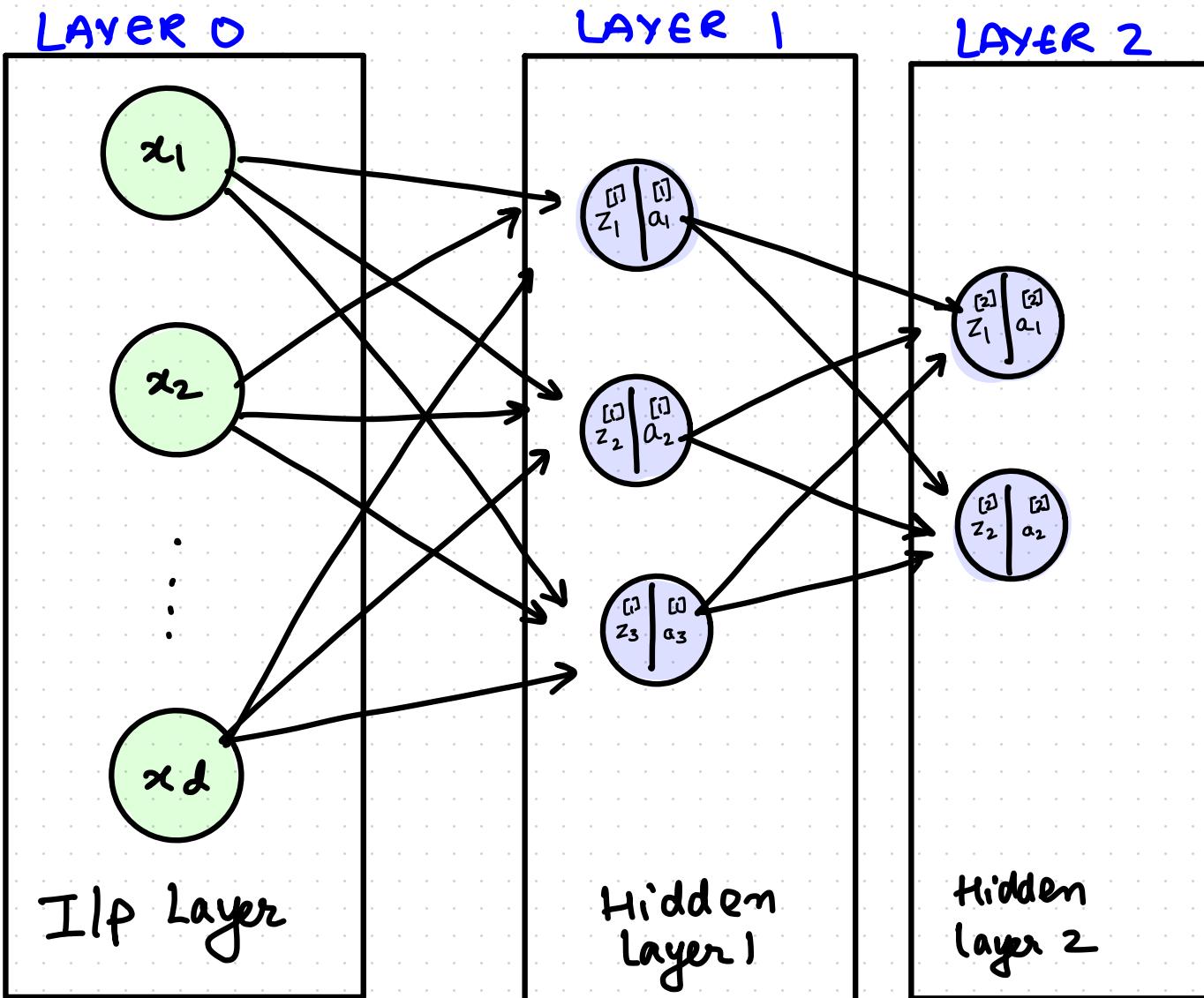


$$z^{[2]} = w^{[2]} a^{[1]} + b^{[2]}$$

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

Q. Dim. of $w^{[2]}$?

FORWARD PROPAGATION



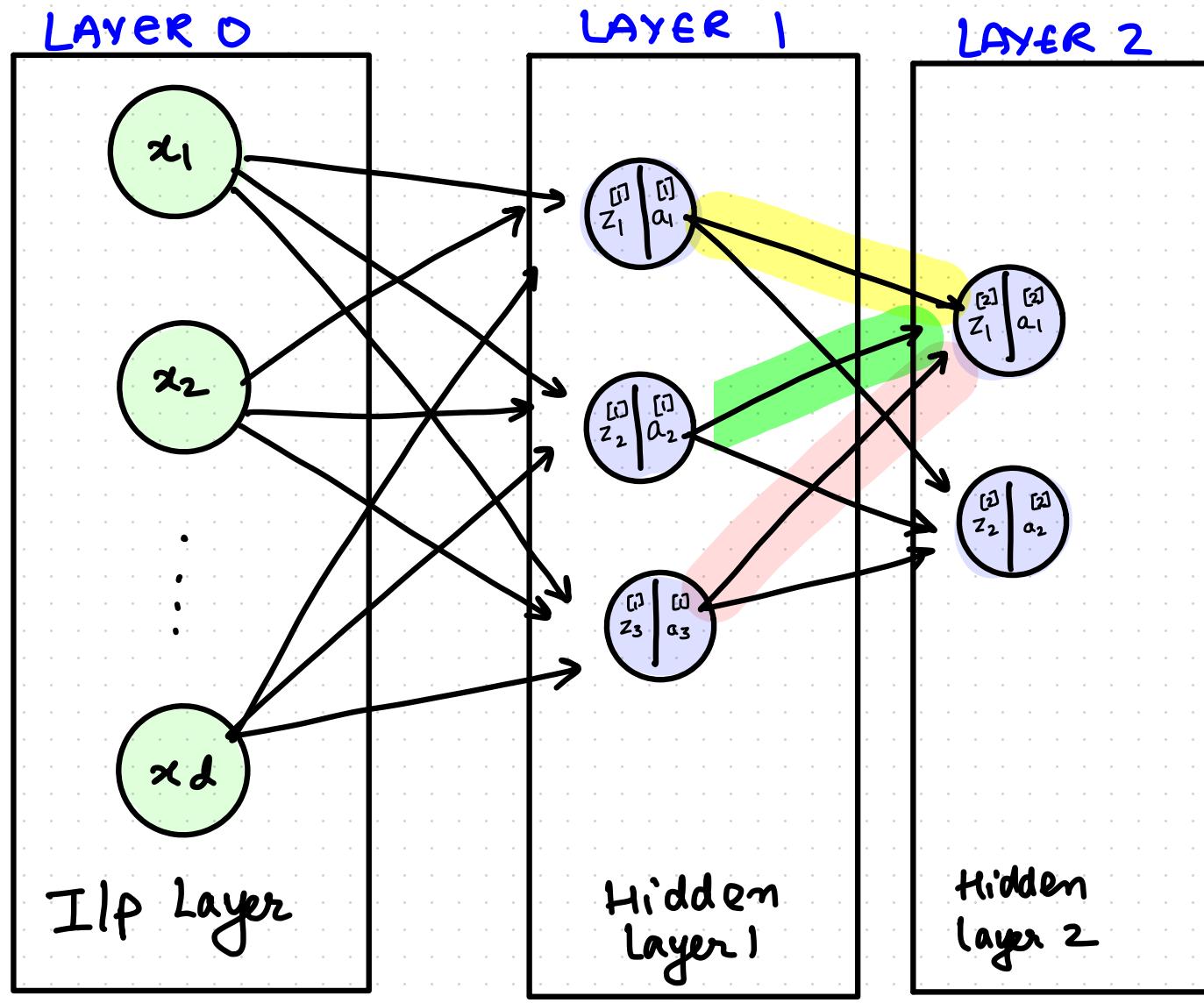
$$z^{[2]} = W^{[2]} a^{[1]} + b^{[2]}$$

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

Q. Dim. of $W^{[2]}$?

$$W^{[2]} = \begin{bmatrix} -w_1^{[2]} - \\ -w_2^{[2]} - \end{bmatrix}$$

FORWARD PROPAGATION



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

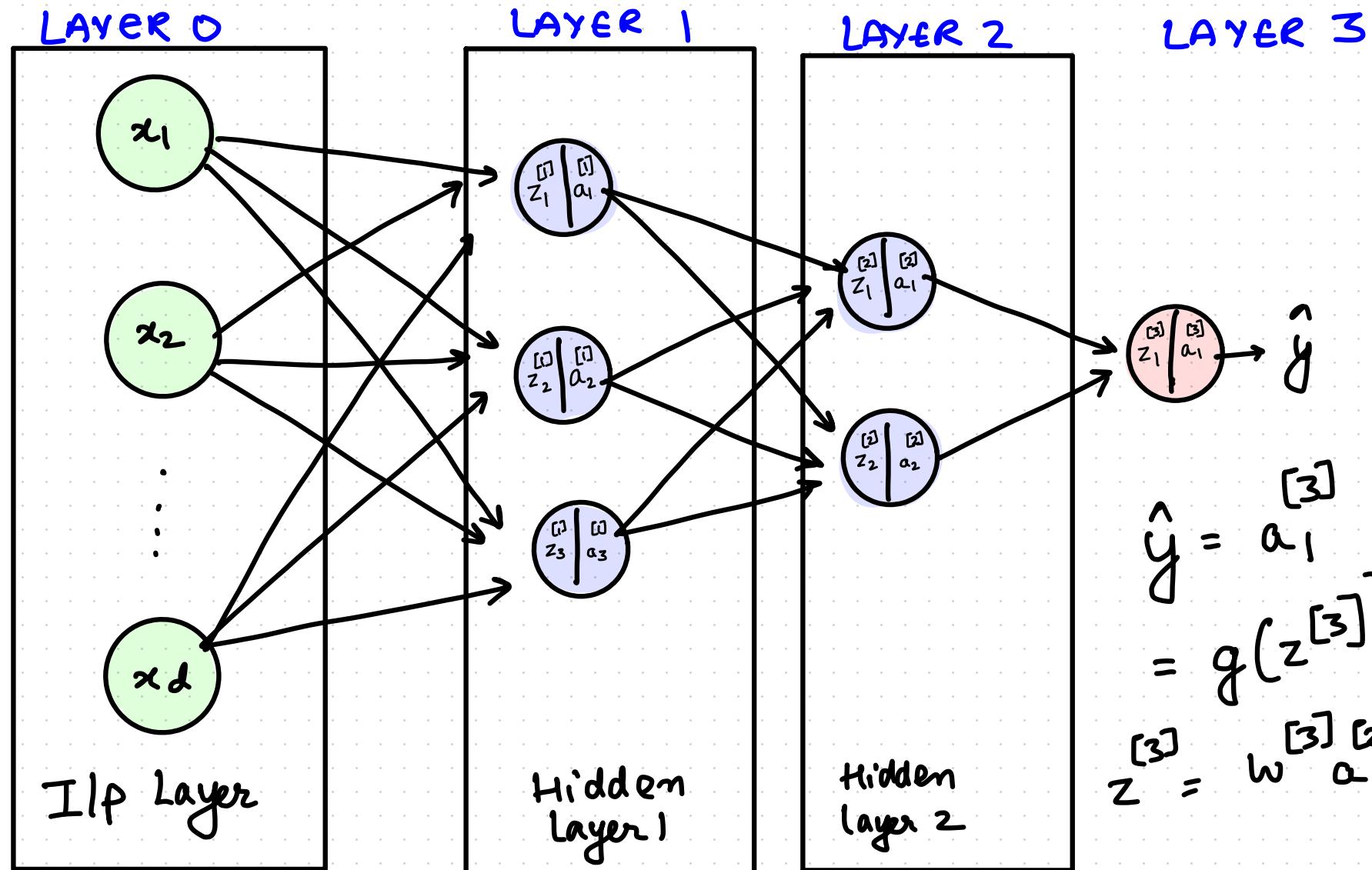
$$a^{[2]} = g(z^{[2]})$$

Q. Dim. of $w^{[2]}$?

$$w^{[2]} = \begin{bmatrix} -w_1^T \\ -w_2^T \end{bmatrix}^T$$

$$\begin{aligned} w_1 &\in \mathbb{R}^3 \\ \therefore w^{[2]} &\in \mathbb{R}^{2 \times 3} \end{aligned}$$

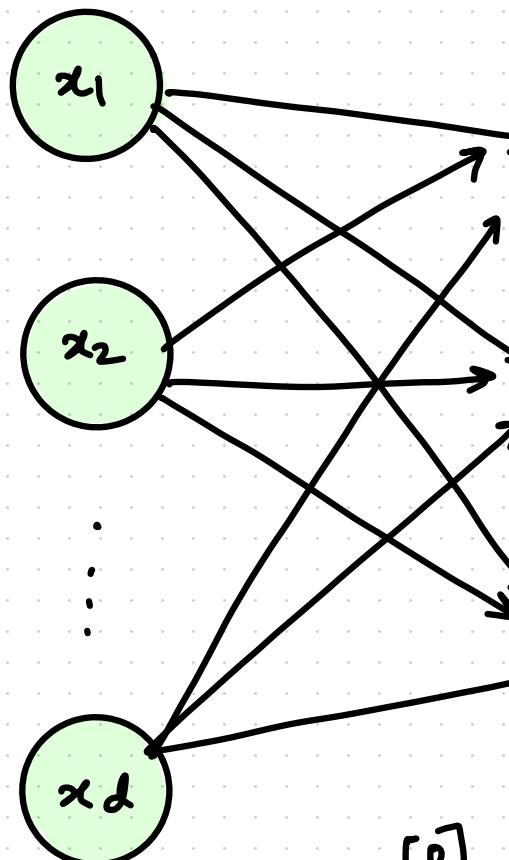
FORWARD PROPAGATION



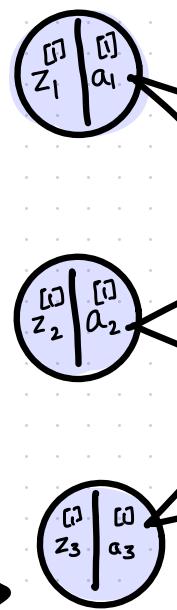
$$\begin{aligned}
 \hat{y} &= a_1^{[3]} \\
 &= g(z^{[3]}) \\
 z^{[3]} &= w^{[3]} a^{[2]} + b^{[3]}
 \end{aligned}$$

WHAT CAN WE SAY ABOUT SHAPES OF a, b, w

LAYER 0



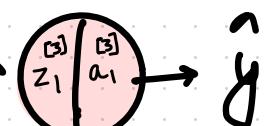
LAYER 1



LAYER 2



LAYER 3

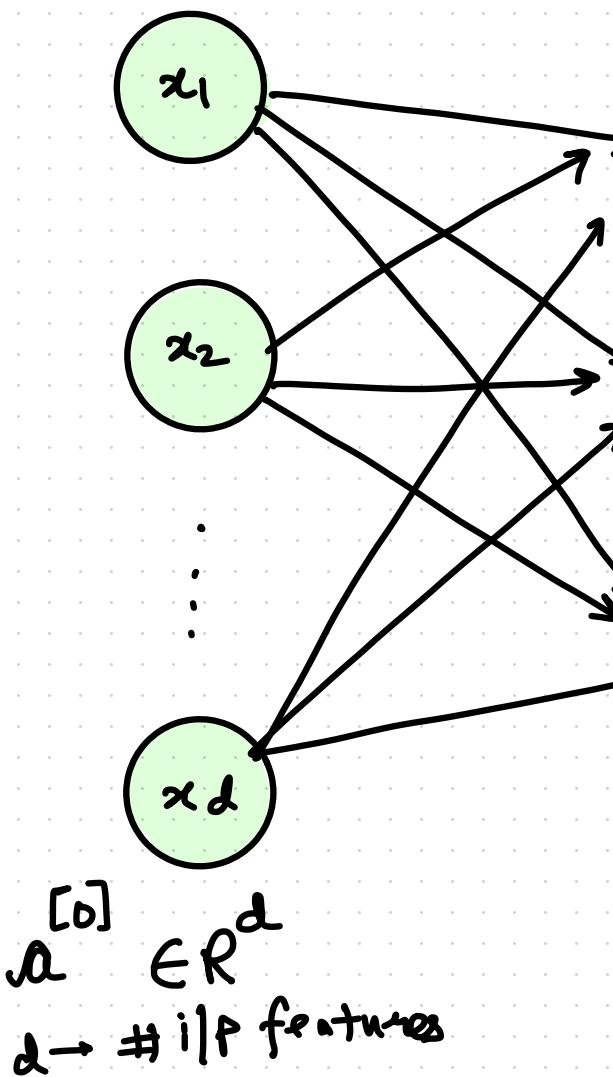


$a^{[0]} \in \mathbb{R}^d$ or $R^{N^{[0]}}$
 $d \rightarrow \# \text{ input features}$

$N^{[0]} = \# \text{ units in } 0^{\text{th}} \text{ layer}$

WHAT CAN WE SAY ABOUT SHAPES OF a, b, w

LAYER 0



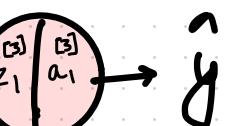
LAYER 1



LAYER 2



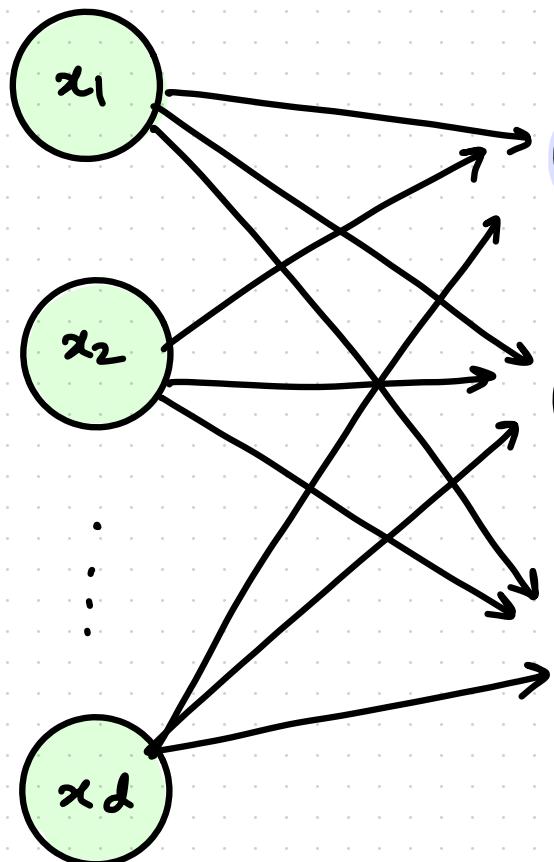
LAYER 3



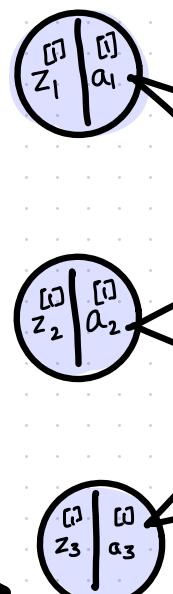
$$w^{[1]} = \begin{bmatrix} -w_1^{[0]^T} \\ -w_2^{[0]^T} \\ \vdots \\ -w_{N^0}^{[0]^T} \end{bmatrix} \Rightarrow w^{[1]} \in \mathbb{R}^{N^1 \times N^0}$$

WHAT CAN WE SAY ABOUT SHAPES OF a, b, w

LAYER 0



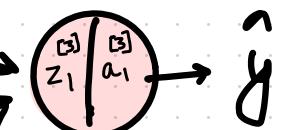
LAYER 1



LAYER 2



LAYER 3



$$b^{[0]} \in \mathbb{R}^{N^{[1]}}$$

$$\begin{aligned} a^{[0]} &\in \mathbb{R}^d \\ d &\rightarrow \# \text{ input features} \end{aligned}$$

SUMMARY OF SHAPES

$N^{[l]}$: # NODES IN l^{th} Layer

$a^{[0]} \in R^{N^{[0]}}$

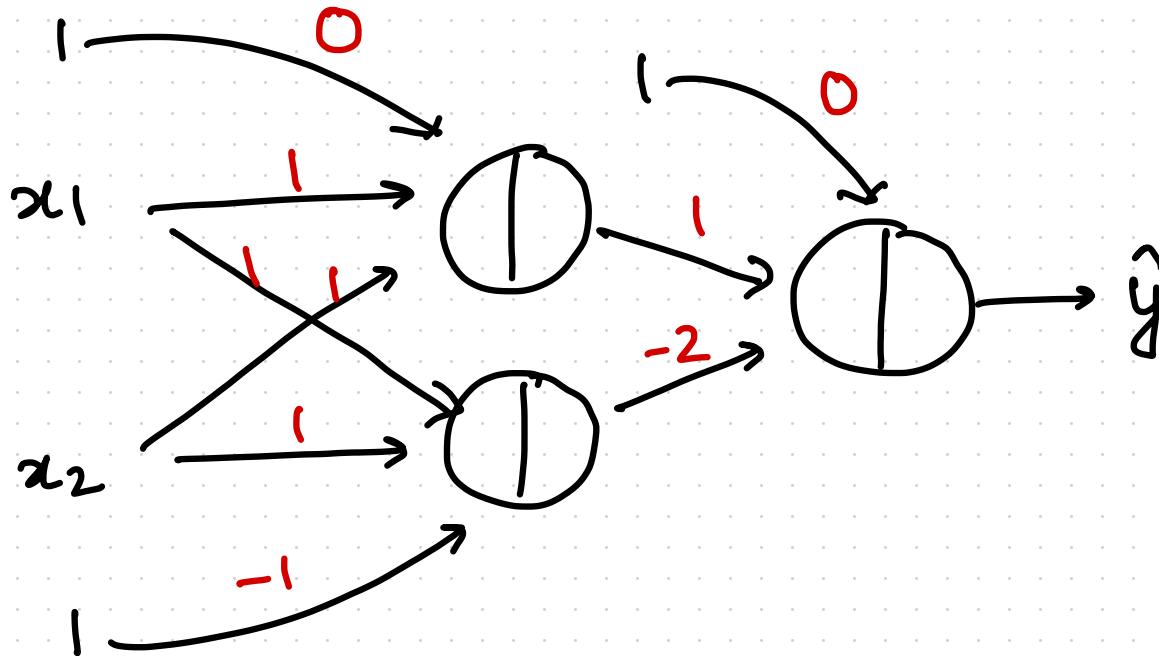
$w^{[l]} \in R^{N^{[l]} \times N^{[l-1]}}$

$b^{[l]} \in R^{N^{[l]}}$

$z^{[l]} \in R^{N^{[l]}}$

$\alpha^{[l]} \in R^{N^{[l]}}$

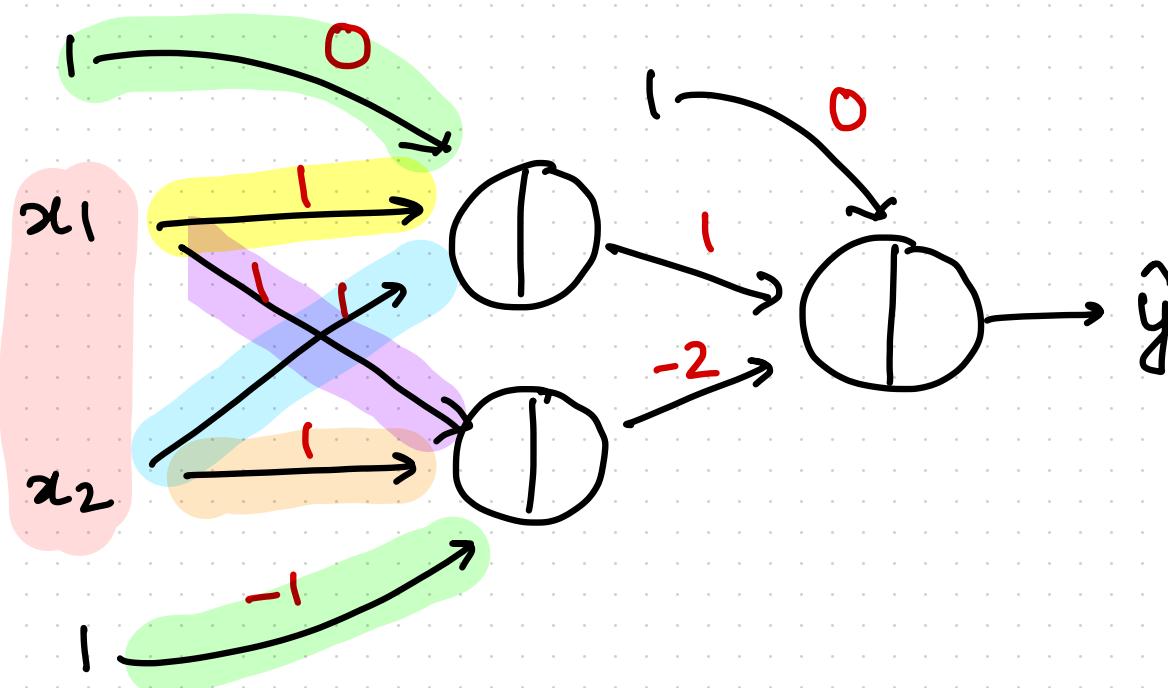
XOR USING "MLP" RELU



CONFIRM IF ABOVE N/W IS CORRECT FOR XOR.

Start with $\begin{bmatrix} x_1 \\ x_2 \end{bmatrix} = \begin{bmatrix} 0 \\ 0 \end{bmatrix}$ $y_{\text{TRUE}} = 0$

XOR USING "MLP" RELU

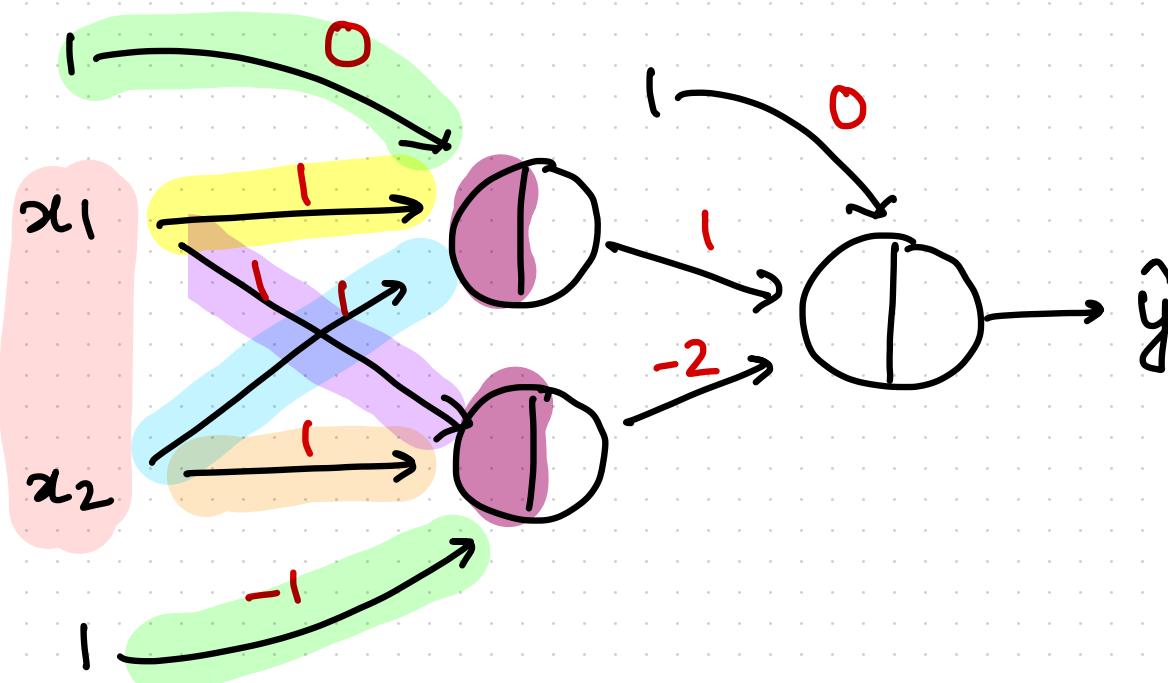


$$a^{[0]} = \begin{bmatrix} 0 \\ 0 \end{bmatrix};$$

$$b^{[1]} = \begin{bmatrix} 0 \\ -1 \end{bmatrix};$$

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

XOR USING A "MLP" RELU



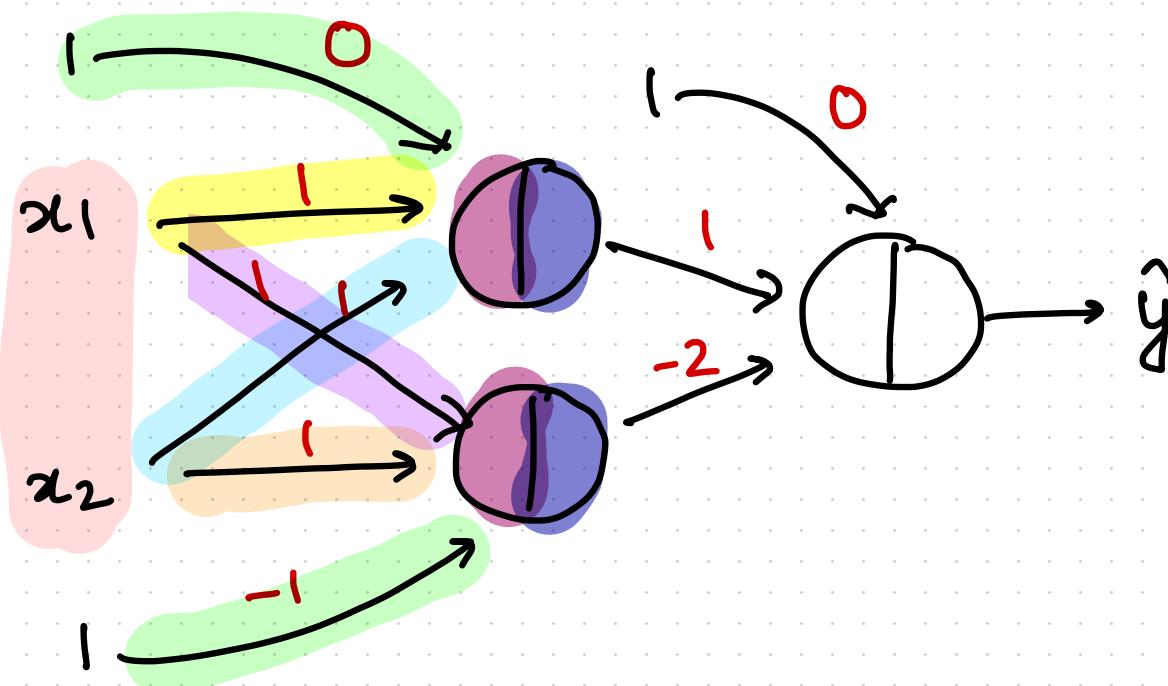
$$a^{[0]} = \begin{bmatrix} 0 \\ 0 \end{bmatrix};$$

$$b^{[1]} = \begin{bmatrix} 0 \\ -1 \end{bmatrix};$$

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

$$z^{[1]} = W^{[1]} a^{[0]} + b^{[1]} = \begin{bmatrix} 1 & 1 \end{bmatrix} \begin{bmatrix} 0 \\ 0 \end{bmatrix} + \begin{bmatrix} 0 \\ -1 \end{bmatrix} = \begin{bmatrix} 0 \\ -1 \end{bmatrix}$$

XOR USING A "MLP" RELU



$$a^{[0]} = \begin{bmatrix} 0 \\ 0 \end{bmatrix};$$

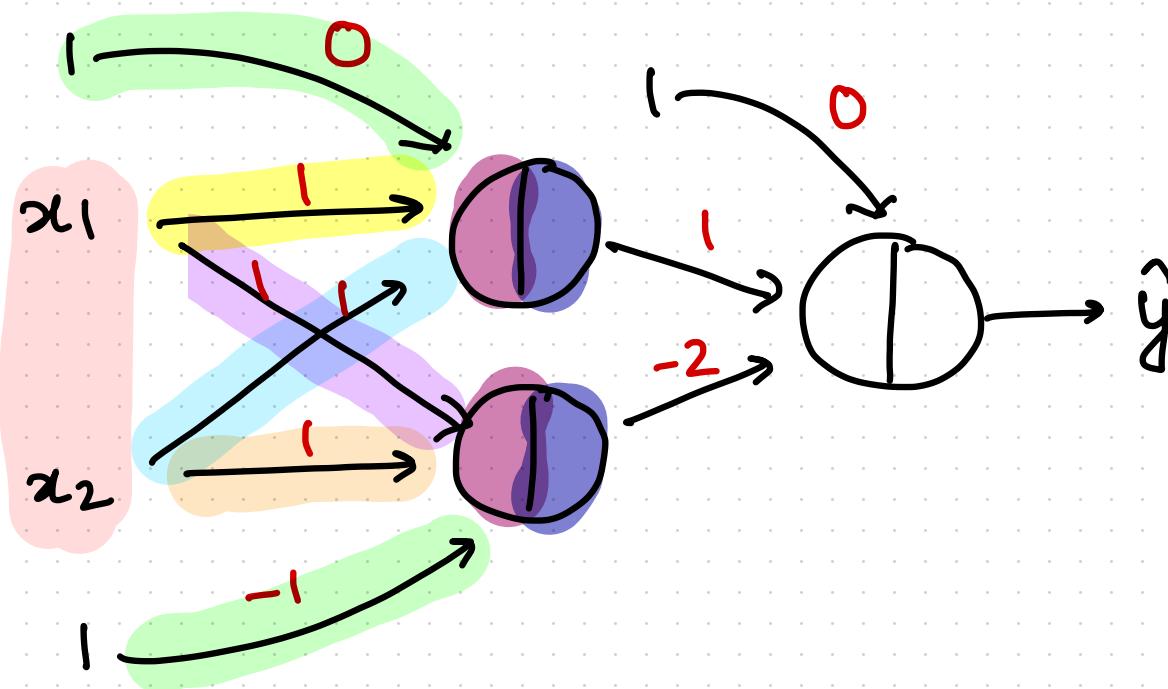
$$b^{[1]} = \begin{bmatrix} 0 \\ -1 \end{bmatrix};$$

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

$$z^{[1]} = W^{[1]} a^{[0]} + b^{[1]} = \begin{bmatrix} 1 & 1 \end{bmatrix} \begin{bmatrix} 0 \\ 0 \end{bmatrix} + \begin{bmatrix} 0 \\ -1 \end{bmatrix} = \begin{bmatrix} 0 \\ -1 \end{bmatrix}$$

$$a^{[1]} = \text{RELU} \begin{bmatrix} 0 \\ -1 \end{bmatrix} = \begin{bmatrix} 0 \\ 0 \end{bmatrix}$$

XOR USING A "MLP" RELU



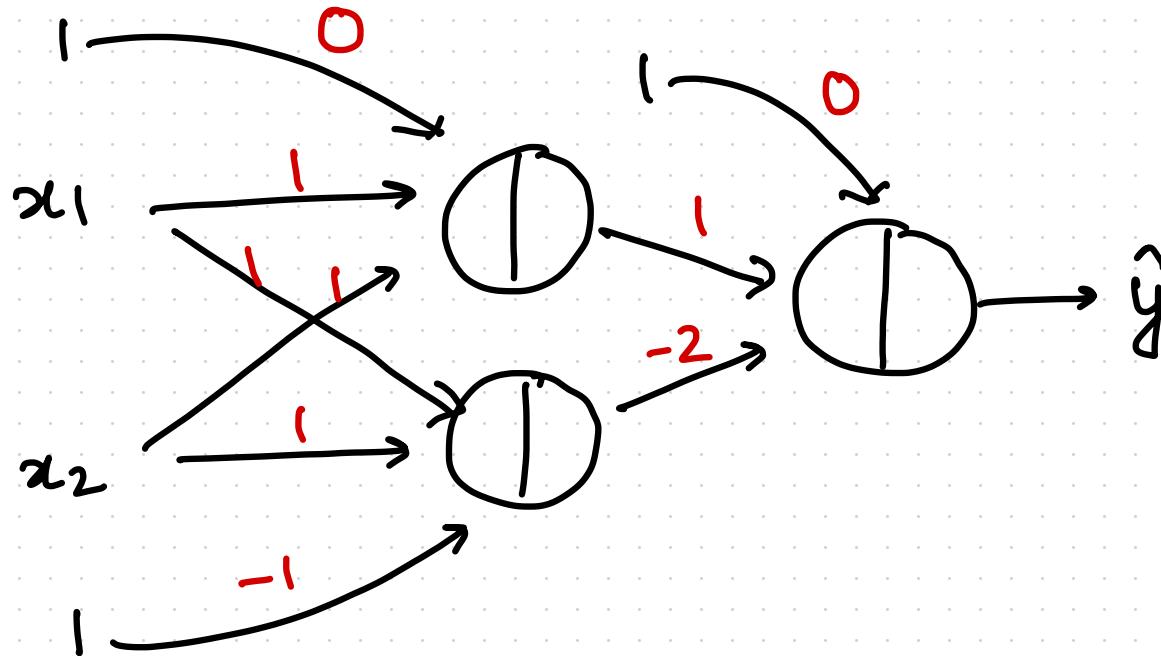
$$a^{[1]} = \begin{bmatrix} 0 \\ 0 \end{bmatrix}$$

$$; W^{[2]} = \begin{bmatrix} 1 & -2 \end{bmatrix} ; b^{[2]} = \begin{bmatrix} 0 \end{bmatrix}$$

$$z^{[2]} = \begin{bmatrix} 1 & -2 \end{bmatrix} \begin{bmatrix} 0 \\ 0 \end{bmatrix} + \begin{bmatrix} 0 \end{bmatrix} = \begin{bmatrix} 0 \end{bmatrix} ; a^{[2]} = \hat{y} = \text{ReLU}(0) = 0$$

✓

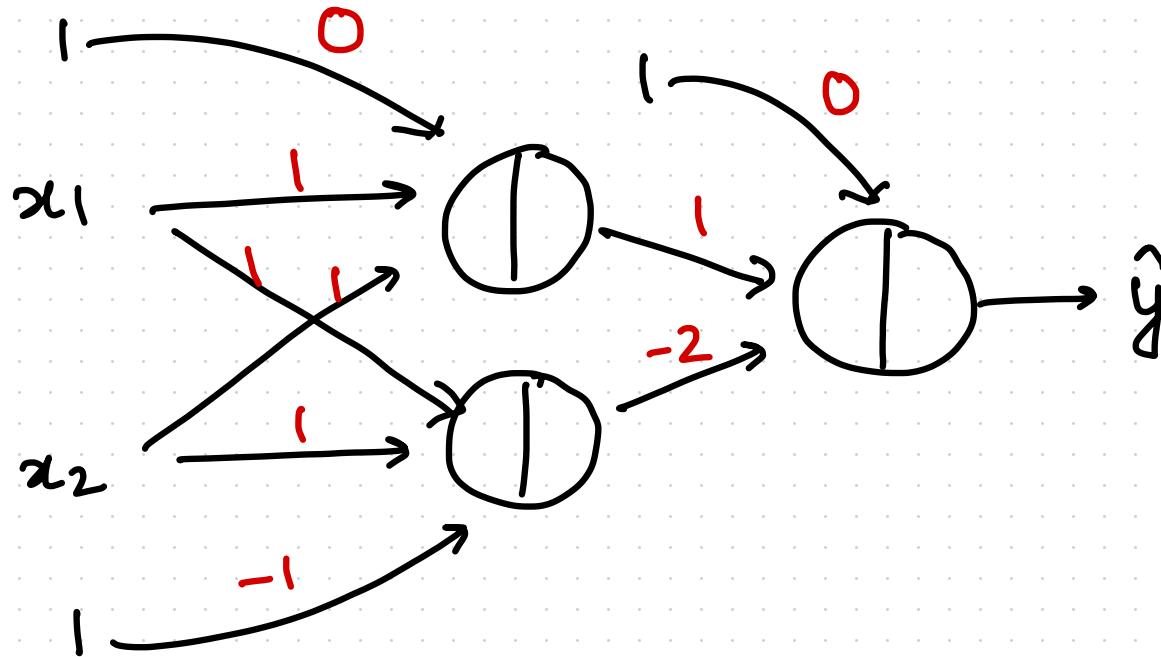
XOR USING "MLP" RELU



CONFIRM IF ABOVE N/W IS CORRECT FOR XOR.

Start with $\begin{bmatrix} x_1 \\ x_2 \end{bmatrix} = \begin{bmatrix} 0 \\ 1 \end{bmatrix}$ $y_{\text{TRUE}} = 1$

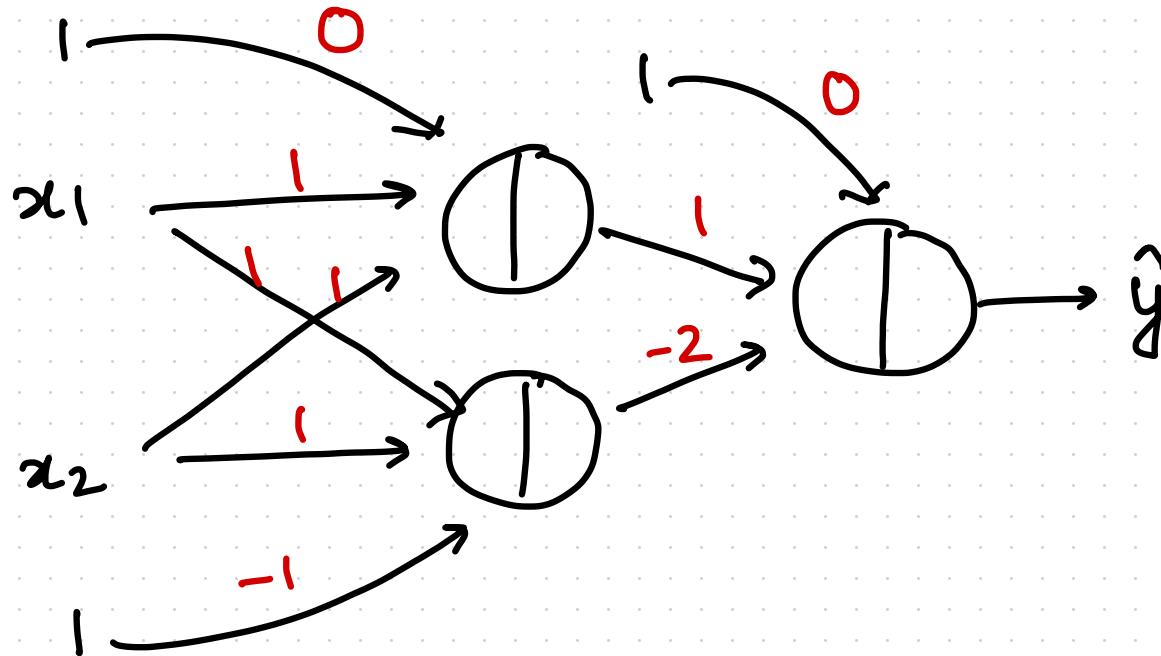
XOR USING "MLP" RELU



$$z^{[1]} = \begin{bmatrix} 1 & 1 \\ 1 & -1 \end{bmatrix} \begin{bmatrix} 0 \\ 1 \end{bmatrix} + \begin{bmatrix} 0 \\ -1 \end{bmatrix} = \begin{bmatrix} 1 \\ 0 \end{bmatrix}$$

$$a^{[1]} = \text{RELU}(z^{[1]}) = \begin{bmatrix} 1 \\ 0 \end{bmatrix}$$

XOR USING A "MLP" RELU



$$z^{[1]} = \begin{bmatrix} 1 & 1 \\ 1 & 1 \end{bmatrix} \begin{bmatrix} 0 \\ 0 \end{bmatrix} + \begin{bmatrix} 0 \\ -1 \end{bmatrix} = \begin{bmatrix} 1 \\ 0 \end{bmatrix}$$

$$a^{[1]} = \text{RELU}(z^{[1]}) = \begin{bmatrix} 1 \\ 0 \end{bmatrix}$$

$$z^{[2]} = \begin{bmatrix} 1 & -2 \\ 1 & 0 \end{bmatrix} \begin{bmatrix} 1 \\ 0 \end{bmatrix} + \begin{bmatrix} 0 \end{bmatrix} = \begin{bmatrix} 1 \end{bmatrix}$$

$$\begin{aligned} a^{[2]} &= \text{RELU}(1) = 1 \\ &= \hat{y} \end{aligned}$$

COMPUTATION FOR N EXAMPLES

$x_{(i)} \in R^d$ or $R^{N^{[0]}}$

$$X = \begin{bmatrix} -x_{(1)}^T - \\ -x_{(2)}^T - \\ -x_{(N)}^T - \end{bmatrix} = \begin{bmatrix} -a_{(1)}^{[0]} - \\ \vdots \\ -a_{(N)}^{[0]} - \end{bmatrix} = A \in R^{N \times N^{[0]}}$$

↑ matrix

COMPUTATION FOR N EXAMPLES

$x_{(i)} \in R^d$ or $R^{N^{[0]}}$

$$X = \begin{bmatrix} -x_{(1)}^T - \\ -x_{(2)}^T - \\ -x_{(N)}^T - \end{bmatrix} = \begin{bmatrix} -a_{(1)}^{[0]}^T - \\ \vdots \\ -a_{(N)}^{[0]}^T - \end{bmatrix} = A \in R^{N \times N^{[0]}}$$

$z^{[i]} \leftarrow$ layer
 $z_{(i)} \leftarrow$ Instance/
 sample

$$= w^{[i]} A^{[0]}_{(i)} + b^{[i]} \in R^{N^{[i]}}$$

Independent of "i"

$a_{(i)}^{[0]} \in R^d \equiv \begin{bmatrix} : \\ : \end{bmatrix}$

COMPUTATION FOR N EXAMPLES

$$x_{(i)} \in R^D$$

$$X = \begin{bmatrix} -x_{(1)}^T - \\ -x_{(2)}^T - \\ -x_{(N)}^T - \end{bmatrix} = \begin{bmatrix} -a_{(1)}^{[0]} - \\ \vdots \\ -a_{(N)}^{[0]} - \end{bmatrix}$$

$$z_{(1)}^{[i]} = \{W a_{(1)}^{[0]} + b^{[i]}\} \in R^{N^{[i]}}$$

$$z_{(2)}^{[i]} = \{W a_{(2)}^{[0]} + b^{[i]}\} \in R^{N^{[i]}}$$

$$\vdots$$

$$z_{(N)}^{[i]} = \{W a_{(N)}^{[0]} + b^{[i]}\} \in R^{N^{[i]}}$$

COMPUTATION FOR N EXAMPLES

$$x_{(i)} \in \mathbb{R}^D$$

$$X = \begin{bmatrix} -x_{(0)}^T - \\ -x_{(1)}^T - \\ -x_{(2)}^T - \\ \vdots \\ -x_{(N)}^T - \end{bmatrix} = \begin{bmatrix} -a_{(0)}^{[0]} - \\ \vdots \\ -a_{(N)}^{[0]} - \end{bmatrix} = A \in \mathbb{R}^{N \times N^{[0]}}$$

$$z_{(1)}^{[1]} = w^{[0]} a_{(1)}^{[0]} + b^{[1]}$$

$$\in \mathbb{R}^{N^{[1]}}$$

$$z_{(2)}^{[1]} = w^{[0]} a_{(2)}^{[0]} + b^{[1]}$$

$$\in \mathbb{R}^{N^{[1]}}$$

$$\vdots$$

$$z_{(N)}^{[1]} = w^{[0]} a_{(N)}^{[0]} + b^{[1]}$$

$$\in \mathbb{R}^{N^{[1]}}$$

COMPUTATION FOR N EXAMPLES

$$x_{(i)} \in \mathbb{R}^D$$

$$X = \begin{bmatrix} -x_{(1)}^T - \\ -x_{(2)}^T - \\ -x_{(N)}^T - \end{bmatrix} = \begin{bmatrix} -a_{(1)}^{[0]} - \\ \vdots \\ -a_{(N)}^{[0]} - \end{bmatrix} = A \in \mathbb{R}^{N \times N^{[0]}}$$

$$z_{(1)}^{[i]} = W^{[0]} a_{(1)}^{[0]} + b^{[i]}$$

$$z_{(2)}^{[i]} = W^{[0]} a_{(2)}^{[0]} + b^{[i]}$$

$$\vdots$$

$$z_{(N)}^{[i]} = W^{[0]} a_{(N)}^{[0]} + b^{[i]}$$

$$\Rightarrow Z = \begin{bmatrix} -z_{(1)}^{[i]} - \\ -z_{(2)}^{[i]} - \\ \vdots \\ -z_{(N)}^{[i]} - \end{bmatrix} \in \mathbb{R}^{N \times N^{[i]}}$$

COMPUTATION FOR N EXAMPLES

$$x_{(i)} \in \mathbb{R}^D$$

$$X = \begin{bmatrix} -x_{(1)}^T - \\ -x_{(2)}^T - \\ -x_{(N)}^T - \end{bmatrix} = \begin{bmatrix} -a_{(1)}^{[0]} - \\ \vdots \\ -a_{(N)}^{[0]} - \end{bmatrix} = A \in \mathbb{R}^{N \times N^{[0]}}$$

$$z_{(1)}^{[i]} = W^{[0]} a_{(1)}^{[0]} + b^{[i]}$$

$$z_{(2)}^{[i]} = W^{[0]} a_{(2)}^{[0]} + b^{[i]}$$

$$\vdots$$

$$z_{(N)}^{[i]} = W^{[0]} a_{(N)}^{[0]} + b^{[i]}$$

$$\Rightarrow Z = \begin{bmatrix} -z_{(1)}^{[i]} - \\ -z_{(2)}^{[i]} - \\ \vdots \\ -z_{(N)}^{[i]} - \end{bmatrix} \in \mathbb{R}^{N \times N^{[i]}}$$

\xleftarrow{N}

COMPUTATION FOR N EXAMPLES

$$A^{[0]} \in R^{N \times N^{[0]}}$$

$$W^{[i]} \in R^{N^{[i]} \times N^{[0]}}$$

$$b^{[i]} \in R^{N^{[i]}}$$

$$B^{[i]} = \begin{bmatrix} -b^{[i]T} \\ \vdots \\ -b^{[i]T} \end{bmatrix} \in R^{N \times N^{[i]}}$$

$$Z^{[i]} \in R^{N \times N^{[i]}}$$

all same entries

COMPUTATION FOR N EXAMPLES

$$A^{[0]} \in R^{N \times N^{[0]}}$$

$$W^{[1]} \in R^{N^{[1]} \times N^{[0]}}$$

$$b^{[1]} \in R^{N^{[1]}}$$

$$B^{[1]} = \begin{bmatrix} -b^{[1]T} \\ \vdots \\ -\bar{b}^{[1]T} \end{bmatrix} \in R^{N \times N^{[1]}}$$

$$Z^{[1]} \in R^{N \times N^{[1]}}$$

$$Z^{[1]} = A^{[0]} W^{[1]} + B^{[1]}$$

COMPUTATION FOR N EXAMPLES

$$A^{[0]} \in R^{N \times N^{[0]}}$$

$$W^{[1]} \in R^{N^{[1]} \times N^{[0]}}$$

$$b^{[1]} \in R^{N^{[1]}}$$

$$B^{[1]} = \begin{bmatrix} -b^{[1]T} \\ \vdots \\ -b^{[1]T} \end{bmatrix} \in R^{N \times N^{[1]}}$$

$$Z^{[1]} \in R^{N \times N^{[1]}}$$

$$Z^{[1]} = A^{[0]} W^{[1]T} + B^{[1]} \Rightarrow A^{[1]} = g(Z^{[1]})$$

$$\boxed{Z^{[l]} = A^{[l-1]} W^{[l]T} + B^{[l]} \\ A^{[l]} = g(Z^{[l]})}$$

XOR ALL EXAMPLES

① $x = A^{[0]} = \begin{bmatrix} 0 & 0 \\ 0 & -1 \\ -1 & 0 \\ 1 & 1 \end{bmatrix}_{4 \times 2}; \hat{y} = \begin{bmatrix} 0 \\ -1 \\ -1 \\ 0 \end{bmatrix}_{4 \times 1}$

XOR ALL EXAMPLES

$$\textcircled{1} \quad x = A^{[0]} = \begin{bmatrix} 0 & 0 \\ 0 & -1 \\ -1 & 0 \\ 1 & 1 \end{bmatrix}_{4 \times 2} ; \quad \hat{y} = \begin{bmatrix} 0 \\ -1 \\ -1 \\ 0 \end{bmatrix}_{4 \times 1}$$

$$w^{[1]} = \begin{bmatrix} 1 & 1 \\ -1 & 1 \end{bmatrix}_{2 \times 2} ; \quad b^{[1]} = \begin{bmatrix} 0 \\ -1 \end{bmatrix} \Rightarrow B^{[1]} = \begin{bmatrix} 0 & -1 \\ 0 & -1 \\ 0 & -1 \\ 0 & -1 \end{bmatrix}_{4 \times 2}$$

XOR ALL EXAMPLES

$$\textcircled{1} \quad X = A^{[0]} = \begin{bmatrix} 0 & 0 \\ 0 & -1 \\ -1 & 0 \\ 1 & 1 \end{bmatrix}_{4 \times 2} ; \quad \hat{y} = \begin{bmatrix} 0 \\ -1 \\ -1 \\ 0 \end{bmatrix}_{4 \times 1}$$

$$w^{[1]} = \begin{bmatrix} 1 & -1 \\ -1 & 1 \end{bmatrix}_{2 \times 2} = w^{[1]T} ; \quad b^{[1]} = \begin{bmatrix} 0 \\ -1 \end{bmatrix} \Rightarrow B^{[1]} = \begin{bmatrix} 0 & -1 \\ 0 & -1 \\ 0 & -1 \\ 0 & -1 \end{bmatrix}_{4 \times 2}$$

$$Z^{[1]} = A^{[0]} w^{[1]T} + B^{[1]} = \begin{bmatrix} 0 & 0 \\ 0 & -1 \\ -1 & 0 \\ -1 & 1 \end{bmatrix} \begin{bmatrix} 1 & -1 \\ -1 & 1 \end{bmatrix} + \begin{bmatrix} 0 & -1 \\ 0 & -1 \\ 0 & -1 \\ 0 & -1 \end{bmatrix} = \begin{bmatrix} 0 & -1 \\ -1 & 0 \\ -1 & 0 \\ 2 & -1 \end{bmatrix}$$

XOR ALL EXAMPLES

$$Z^{[1]} = A^{[0]} W^{[1]} + B^{[1]} = \begin{bmatrix} 0 & 0 \\ 0 & 1 \\ 1 & 0 \\ 1 & 1 \end{bmatrix} \begin{bmatrix} 1 & 1 \\ 1 & 1 \end{bmatrix} + \begin{bmatrix} 0 & -1 \\ 0 & -1 \\ 0 & -1 \\ 0 & -1 \end{bmatrix} = \begin{bmatrix} 0 & -1 \\ 1 & 0 \\ -1 & 0 \\ 2 & 1 \end{bmatrix}$$

$$A^{[1]} = \text{RELU}\left(\begin{bmatrix} 0 & -1 \\ 1 & 0 \\ 1 & 0 \\ 2 & 1 \end{bmatrix}\right) = \begin{bmatrix} 0 & 0 \\ 1 & 0 \\ 1 & 0 \\ 2 & 1 \end{bmatrix}_{4 \times 2}$$

$$W^{[2]} = \begin{bmatrix} 1 & -2 \end{bmatrix}$$

$$b^{[2]} = [0] \Rightarrow B = \begin{bmatrix} 0 \\ 0 \\ 0 \\ 0 \end{bmatrix}$$

XOR ALL EXAMPLES

$$A^{[1]} = \begin{bmatrix} 0 & 0 \\ 1 & 0 \\ 1 & 0 \\ 2 & 1 \end{bmatrix}_{4 \times 2}$$

$$w^{[2]} = [1 \ -2]_{1 \times 2}$$

$$b^{[2]} = [0] \Rightarrow B^{[2]} = \begin{bmatrix} 0 \\ 0 \\ 0 \\ 0 \end{bmatrix}_{4 \times 1}$$

$$z^{[2]} = A^{[1]} w^{[2]^T} + B^{[2]} = \begin{bmatrix} 0 & 0 & 0 \\ 1 & 0 & 0 \\ 1 & 0 & 0 \\ 2 & 1 & 0 \end{bmatrix} \begin{bmatrix} 1 \\ -2 \end{bmatrix} + \begin{bmatrix} 0 \\ 0 \\ 0 \\ 0 \end{bmatrix} = \begin{bmatrix} 0 \\ -1 \\ -1 \\ 0 \end{bmatrix}$$

XOR ALL EXAMPLES

$$z^{[2]} = A^{[1]} w^{[2]^\top} + b^{[2]} = \begin{bmatrix} 0 & 0 \\ 1 & 0 \\ -1 & 0 \\ 2 & 1 \end{bmatrix} \begin{bmatrix} 1 \\ -2 \end{bmatrix} + \begin{bmatrix} 0 \\ 0 \\ 0 \\ 0 \end{bmatrix} = \begin{bmatrix} 0 \\ 1 \\ 0 \\ -1 \end{bmatrix}$$

$$A^{[2]} = \hat{y} = \text{RELU}\left(\begin{bmatrix} 0 \\ 1 \\ 0 \\ -1 \end{bmatrix}\right) = y_{\text{G.T}}$$

PARAMETERS

Parameters: $w^{[l]}$; $b^{[l]}$ $\neq l$

Q: # parameters for XOR example?

PARAMETERS

Parameters: $w^{[l]}$; $b^{[l]}$ $\neq l$

Q: # parameters for XOR example?

$$N^{[0]} = 2; N^{[1]} = 2; N^{[2]} = 1$$

$$w^{[1]} \in R^{N^{[0]} \times N^{[1]}} = \begin{bmatrix} \quad \\ \quad \end{bmatrix}_{2 \times 2} = 4 \text{ params} = N^{[1]} * N^{[0]}$$

$$b^{[1]} \in R^{N^{[1]}} = \begin{bmatrix} \quad \end{bmatrix}_{2 \times 1} = 2 \text{ params} = N^{[1]}$$

$$w^{[2]} \rightarrow 2 \text{ params} = N^{[2]} * N^{[1]} \quad \& \quad b^{[2]} = N^{[2]} \text{ params}$$

PARAMETERS

Parameters: $w^{[l]}$; $b^{[l]}$ $\forall l$

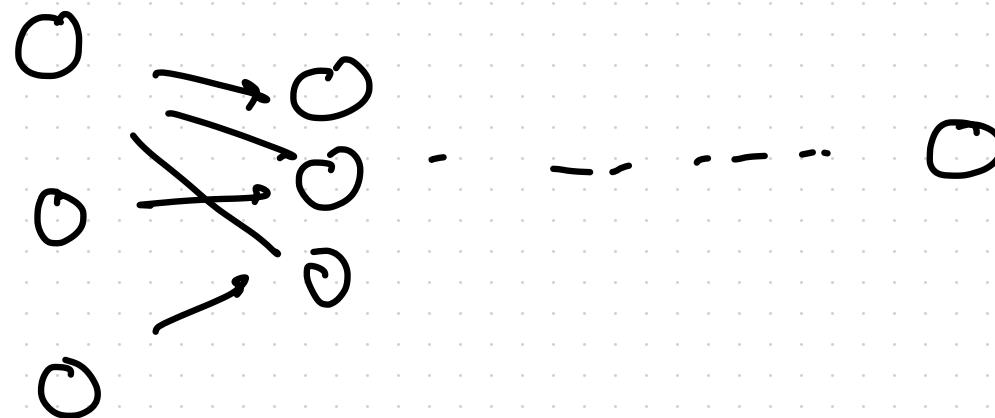
$$\text{PARAMS} = \sum_{l=1}^L N^{[l]} \cdot N^{[l-1]} + N^{[l]}$$

? XOR 4 examples, 9 parameters ??

⇒ NOTEBOOK: XOR demo

LEARNING

PARAMETERS



FORWARD

PROPAGATION (PREDICT
BASED ON
CURRENT
PARAMS)

BACKWARD

PROPAGATION (CHANGE
WEIGHTS TO
IMPROVE
OBJECTIVE)

LEARNING

PARAMETERS

Assume $\text{if } P : X \in \mathbb{R}^{N \times N^{[0]}}$

$\text{of } P : \hat{y}; \text{ G.T. : } y$

$$\text{LOSS} = \frac{1}{N} \sum_{i=1}^N L(\hat{y}^{(i)}, y) \quad \text{or} \quad \sum_{i=1}^N L(\hat{y}^{(i)}, y)$$

Params: $w^{[1]}, b^{[1]}, \dots$

LEARNING

PARAMETERS

Assume $\text{if } p : X \in \mathbb{R}^{N \times N^{[0]}}$

$\text{of } p : \hat{y}; \text{ G.T. : } y$

$$\text{LOSS} = \frac{1}{N} \sum_{i=1}^N L(\hat{y}_{(i)}, y) = J(\theta)$$

Params: $\theta = \{w^{[1]}, b^{[1]}, \dots\}$

GRADIENT DESCENT

① INIT Params randomly

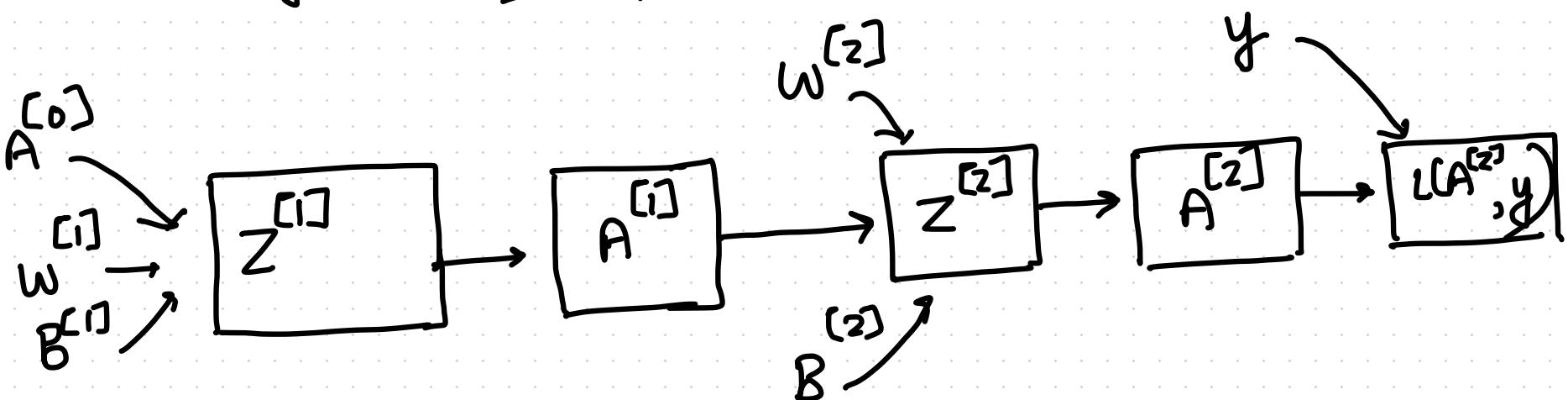
② Till convergence:

$$w^{[1]} = w^{[1]} - \alpha \frac{\partial J(\theta)}{\partial w^{[1]}} \dots$$

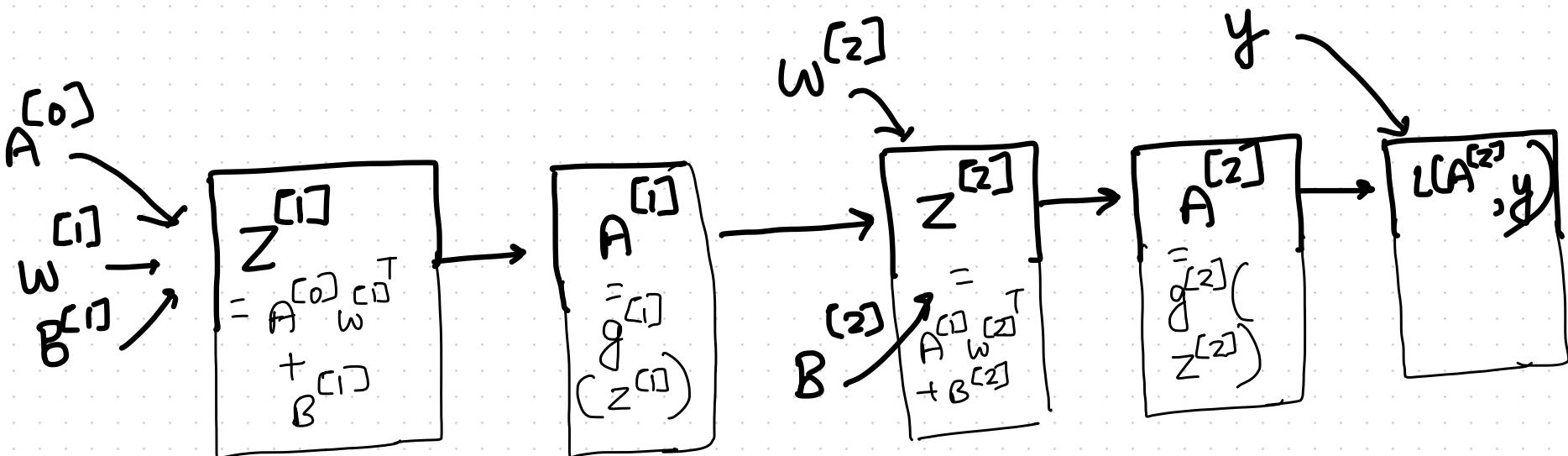
How TO COMPUTE?

COMPUTATION GRAPH (FOR XOR EXAMPLE)

$$\begin{aligned} \textcircled{1} \quad z^{[1]} &= A^{[0]} w^{[1]} + b^{[1]} \\ \textcircled{2} \quad A^{[1]} &= g^{[1]}(z^{[1]}) \\ \textcircled{3} \quad z^{[2]} &= A^{[1]} w^{[2]} + b^{[2]} \\ \textcircled{4} \quad A^{[2]} &= g^{[2]}(z^{[2]}) = \hat{y} \end{aligned}$$



COMPUTATION GRAPH (FOR XOR EXAMPLE)



Let $g^{[1]} = \text{SIGMOID}$; ASSUME CROSS ENTROPY LOSS
 $g^{[2]} = \text{SIGMOID}$

WHAT IS $\frac{\partial L(\theta)}{\partial w^{[0]}}$; WHAT IS $\frac{\partial L(\theta)}{\partial A^{[2]}}$?

DERIVATIVES OF ACTIVATION FUNCTIONS

RELU

$$g(z) = \begin{cases} z; & z > 0 \\ 0; & z \leq 0 \\ \text{undefined}; & 0/w \end{cases}$$

↓ Assume $z \neq 0$

$$g(z) = \begin{cases} z; & z > 0 \\ 0; & z \leq 0 \end{cases}$$

$$g'(z) = \begin{cases} 1; & z > 0 \\ 0; & 0/w \end{cases}$$

DERIVATIVES OF ACTIVATION FUNCTIONS

RELU (LEAKY)

$$g(z) = \begin{cases} z; & z > 0 \\ \alpha z; & z \leq 0 \end{cases}$$

unrefined; $\alpha \neq 0$



$$g(z) = \begin{cases} z; & z \geq 0 \\ \alpha z; & z < 0 \end{cases}$$

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

DERIVATIVES OF ACTIVATION FUNCTIONS

SIGMOID

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

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

$$g'(z) = \frac{1+e^{-z}}{(1+e^{-z})^2} - \frac{1}{(1+e^{-z})^2} = g(z)(1-g(z))$$

DERIVATIVES OF ACTIVATION FUNCTIONS

TANH

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

$$g'(z) = \frac{v \frac{du}{dz} - u \frac{dv}{dz}}{v^2} = \frac{\left(e^z + e^{-z} \right) \left(e^z - e^{-z} \right) - \left(e^z - e^{-z} \right) \left(e^z + e^{-z} \right)}{\left(e^z + e^{-z} \right)^2}$$
$$= 1 - \left(g(z) \right)^2$$

ACTIVATION FUNCTIONS (SUMMARY)

RELU

$$g(z) = \begin{cases} z; & z \geq 0 \\ 0; & z < 0 \end{cases}$$

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

L-RELU

$$g(z) = \begin{cases} z; & z \geq 0 \\ \alpha z; & z < 0 \end{cases}$$

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

SIGMOID

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

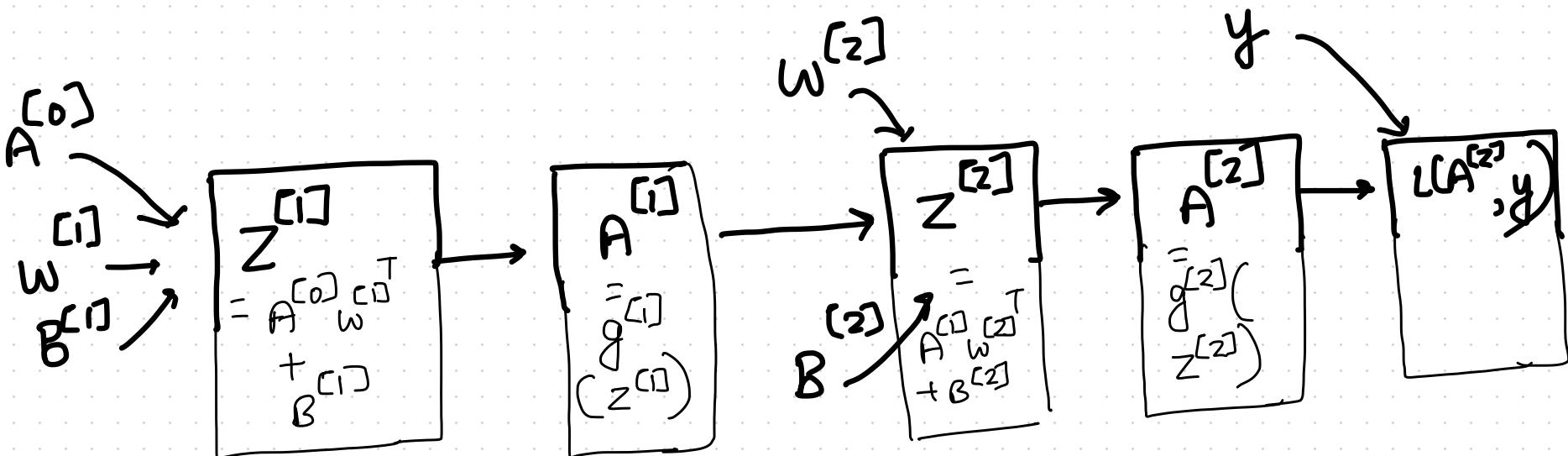
$$g'(z) = g(z) \cdot (1 - g(z))$$

TANH

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

$$g'(z) = 1 - (g(z))^2$$

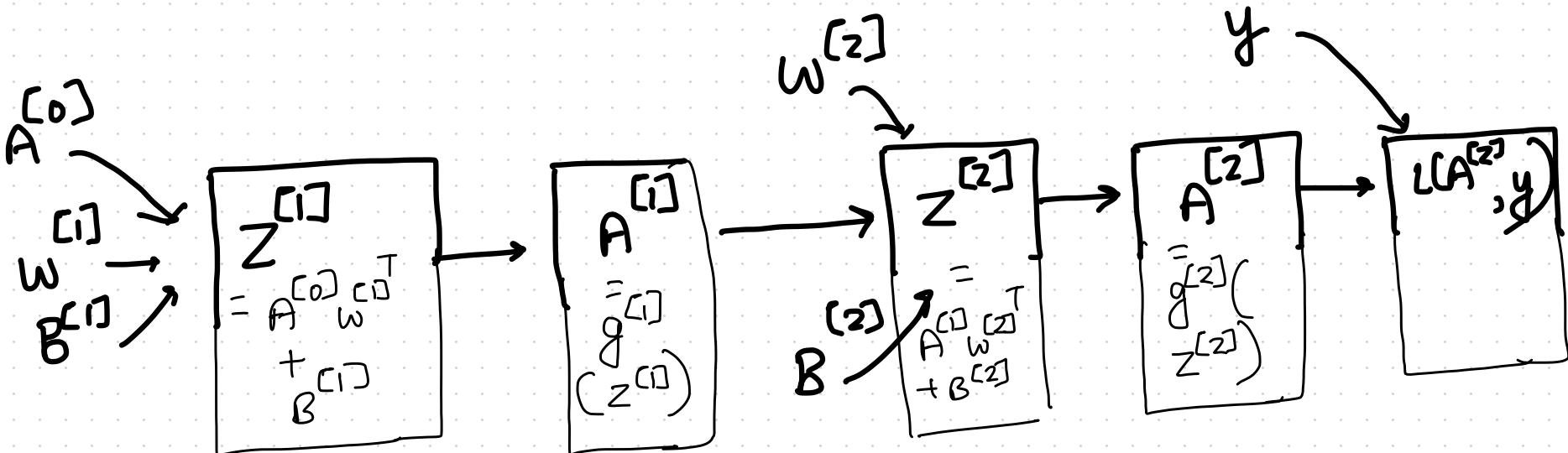
COMPUTATION GRAPH (FOR XOR EXAMPLE)



$$L(A^{[2]}, y) = \sum_{i=1}^N -y_{(i)} \log A_{(i)}^{[2]} - (1 - y_{(i)}) \log (1 - A_{(i)}^{[2]})$$

WRITE IN VECTOR FORM

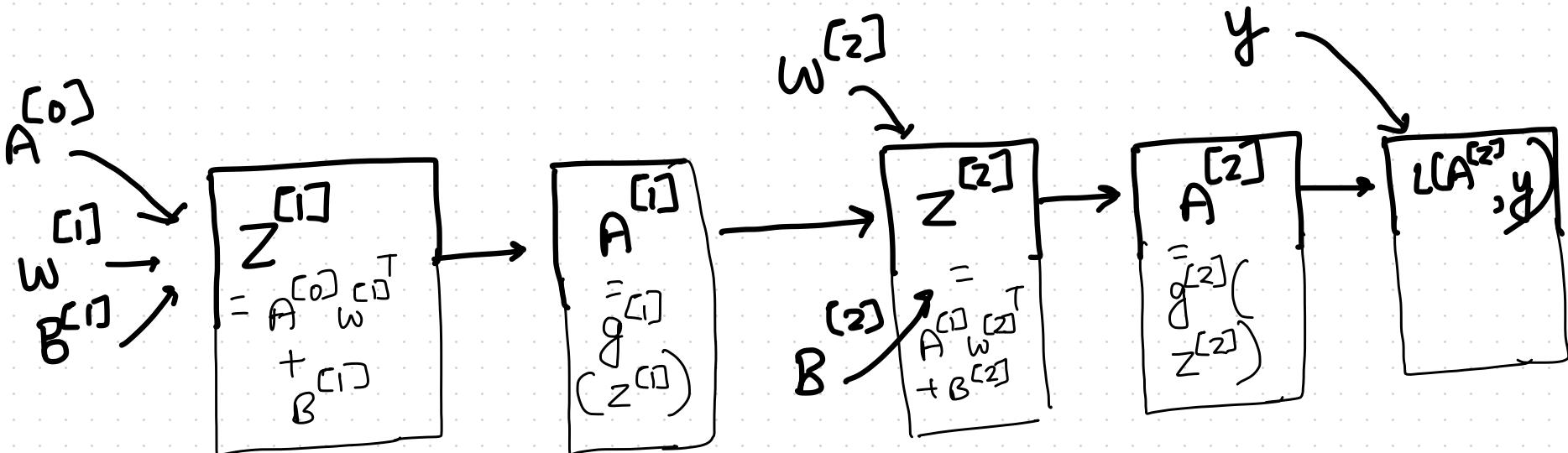
COMPUTATION GRAPH (FOR XOR EXAMPLE)



$$\begin{aligned}
 L(A^{[2]}, y) &= \sum_{i=1}^N -y_{(i)} \log A_{(i)}^{[2]} - (1 - y_{(i)}) \log (1 - A_{(i)}^{[2]}) \\
 &= -y^T \underbrace{\log(A^{[2]})}_{1 \times N} - (1 - y)^T \underbrace{\log(1 - A^{[2]})}_{N \times 1} \quad \begin{bmatrix} \vdots \\ \vdots \end{bmatrix}_{N \times 1}
 \end{aligned}$$

APPLIED ELEMENT-WISE

COMPUTATION GRAPH (FOR XOR EXAMPLE)



$$\begin{aligned}
 L(A^{[2]}, y) &= \sum_{i=1}^N -y_{(i)} \log A_{(i)}^{[2]} - (1 - y_{(i)}) \log (1 - A_{(i)}^{[2]}) \\
 &= -y^T \underbrace{\log(A^{[2]})}_{1 \times N} - (1 - y)^T \underbrace{\log(1 - A^{[2]})}_{N \times 1} \quad \begin{bmatrix} \vdots \\ \vdots \end{bmatrix}_{N \times 1}
 \end{aligned}$$

$$\frac{\partial L(A^{[2]}, y)}{\partial A^{[2]}} = ?$$

APPLIED ELEMENT-WISE

$$\text{let } q_1 = \mathbf{y}^T \log(\mathbf{A}^{[2]})$$

$$q_1 = y_{(1)} \log A_{(1)}^{[2]} + \dots + y_{(n)} A_{(n)}^{[2]}$$

$$\text{let } q_1 = \mathbf{y}^T \log(\mathbf{A}^{[2]})$$

$$q_1 = y_{(1)} \log A_{(1)}^{[2]} + \dots + y_{(N)} \log A_{(N)}^{[2]}$$

$$\frac{\partial q_1}{\partial \mathbf{A}^{[2]}} = \begin{bmatrix} \frac{\partial}{\partial A_{(1)}^{[2]}} (y_{(1)} \log A_{(1)}^{[2]} + \dots) \\ \vdots \\ \frac{\partial}{\partial A_{(N)}^{[2]}} (y_{(1)} \dots) \end{bmatrix} = \begin{bmatrix} \frac{y_{(1)}}{A_{(1)}^{[2]}} \\ \vdots \\ \frac{y_{(N)}}{A_{(N)}^{[2]}} \end{bmatrix}$$

$$\frac{\partial q_1}{\partial \mathbf{A}^{[2]}}_{N \times 1} = \mathbf{y}_{N \times 1} \odot \mathbf{A}_{N \times 1}^{[2]}$$

Element-wise division

$$\text{Let } \varphi = (1 - y)^T \log (1 - A^{[2]})$$

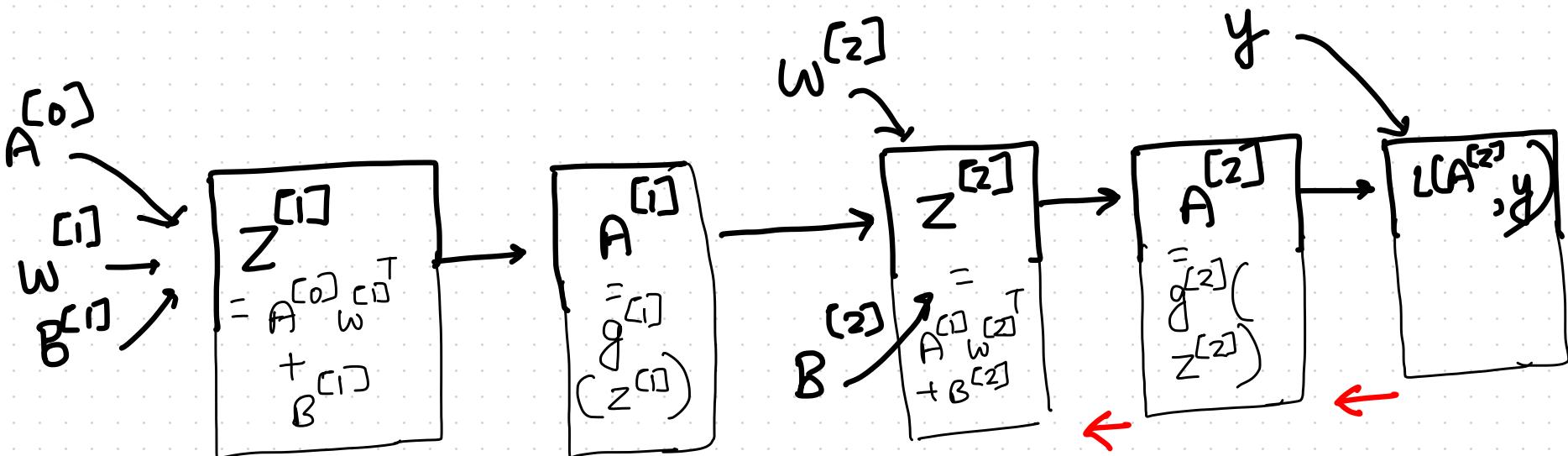
$$\varphi = (1 - y_{(1)}) \log (1 - A_{(1)}^{[2]}) + \dots$$

$$\frac{\partial \varphi}{\partial A^{[2]}} = \left[\begin{array}{c} \frac{\partial}{\partial A_{(1)}^{[2]}} \left\{ (1 - y_{(1)}) \log (1 - A_{(1)}^{[2]}) + \dots \right\} \\ \vdots \\ \frac{\partial}{\partial A_{(N)}^{[2]}} \left\{ (1 - y_{(N)}) \log (1 - A_{(N)}^{[2]}) + \dots \right\} \end{array} \right] = \begin{bmatrix} \frac{(1 - y_{(1)})}{(1 - A_{(1)}^{[2]})} (-1) \\ \vdots \\ \frac{(1 - y_{(N)})}{(1 - A_{(N)}^{[2]})} (-1) \end{bmatrix}$$

$$\frac{\partial \varphi}{\partial A^{[2]}}_{N \times 1} = -(1 - y)_{N \times 1} \odot (1 - A^{[2]})_{N \times 1}$$

$$\frac{\partial L(A^{[2]}, y)}{\partial A^{[2]}} = -y \odot A^{[2]} + (1-y) \odot (1-A^{[2]})$$

COMPUTATION GRAPH (FOR XOR EXAMPLE)

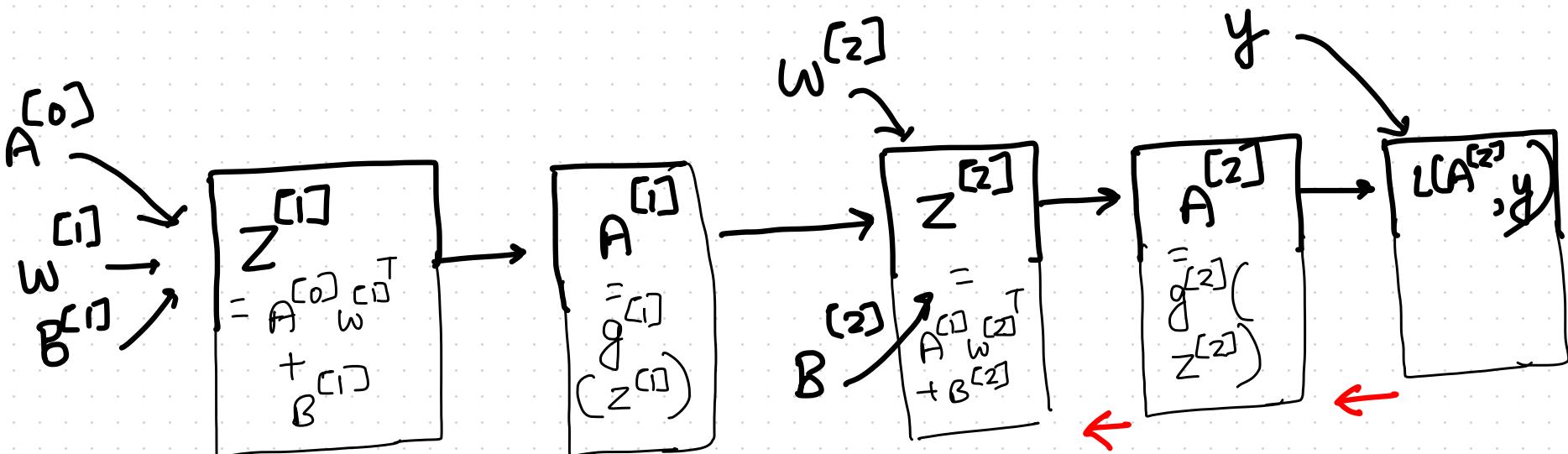


$$\frac{\partial L}{\partial Z^{[2]}} = ?$$

$$\frac{\partial L}{\partial Z^{[2]}} = \frac{\partial L}{\partial A^{[2]}} \cdot \frac{\partial A^{[2]}}{\partial Z^{[2]}}$$

(Chain Rule)

COMPUTATION GRAPH (FOR XOR EXAMPLE)



$$\frac{\partial L}{\partial Z^{[2]}} = ? \quad \frac{\partial L}{\partial Z^{[2]}} = \frac{\partial L}{\partial A^{[2]}} \cdot \frac{\partial A^{[2]}}{\partial Z^{[2]}} \quad (\text{Chain Rule})$$

$$= \frac{\partial L}{\partial A^{[2]}} \cdot \frac{\partial g^{[2]}(Z^{[2]})}{\partial Z^{[2]}} = \frac{\partial L}{\partial A^{[2]}} g(Z^{[2]}) (1 - g(Z^{[2]}))$$

$$\therefore g^{[2]} = \text{SIG MOID}$$

COMPUTATION GRAPH (FOR XOR EXAMPLE)

$$= \frac{\partial L}{\partial A^{[2]}} \cdot \frac{\partial g^{[2]}(z^{[2]})}{\partial z^{[2]}} = \frac{\partial L}{\partial A^{[2]}} g(z^{[2]}) (1 - g(z^{[2]}))$$

$$= \frac{\partial L}{\partial A^{[2]}} \circ A^{[2]} \odot (1 - A^{[2]})$$

$N \times 1$ $N \times 1$

Element-wise multiply

$$= (-y \odot A^{[2]} + (1-y) \odot (1-A^{[2]})) \left((A^{[2]} \odot (1-A^{[2]})) \right)$$

COMPUTATION GRAPH (FOR XOR EXAMPLE)

$$= \frac{\partial L}{\partial A^{[2]}} \cdot \frac{\partial g^{[2]}(z^{[2]})}{\partial z^{[2]}} = \frac{\partial L}{\partial A^{[2]}} g(z^{[2]}) (1 - g(z^{[2]}))$$

$$= \frac{\partial L}{\partial A^{[2]}} A^{[2]} \odot (1 - A^{[2]})$$

$N \times 1$ $N \times 1$

Element-wise

$$= (-y \odot A^{[2]} + (1-y) \odot (1-A^{[2]})) \left(\begin{pmatrix} A^{[2]} & (1-A^{[2]}) \end{pmatrix} \right)$$

$$= -y \odot A^{[2]} \circ A^{[2]} \odot (1-A^{[2]}) + (1-y) \odot (1-A^{[2]}) \circ A^{[2]} \odot (1-A^{[2]})$$

COMPUTATION GRAPH (FOR XOR EXAMPLE)

$$= \frac{\partial L}{\partial A^{[2]}} \cdot \frac{\partial g(z^{[2]})}{\partial z^{[2]}} = \frac{\partial L}{\partial A^{[2]}} g(z^{[2]}) (1 - g(z^{[2]}))$$

$$= \frac{\partial L}{\partial A^{[2]}} A^{[2]} \odot (1 - A^{[2]})$$

$N \times 1$ $N \times 1$

Element-wise

$$= (-y \odot A^{[2]} + (1-y) \odot (1-A^{[2]})) \left((A^{[2]} \odot (1-A^{[2]})) \right)$$

$$= -y \odot A^{[2]} \odot A^{[2]} \odot (1-A^{[2]}) + (1-y) \odot (1-A^{[2]}) \odot A^{[2]} \odot (1-A^{[2]})$$

$$= -y \odot (1-A^{[2]}) + (1-y) \odot A^{[2]} = -y_{N \times 1} + y \odot A^{[2]}_{N \times 1} + A^{[2]}_{N \times 1} - y \odot A^{[2]}_{N \times 1}$$

COMPUTATION GRAPH (FOR XOR EXAMPLE)

$$= \frac{\partial L}{\partial A^{[2]}} \cdot \frac{\partial g(z^{[2]})}{\partial z^{[2]}} = \frac{\partial L}{\partial A^{[2]}} g(z^{[2]}) (1 - g(z^{[2]}))$$

$$= \frac{\partial L}{\partial A^{[2]}} A^{[2]} \odot (1 - A^{[2]})$$

$N \times 1$ $N \times 1$

Element-wise

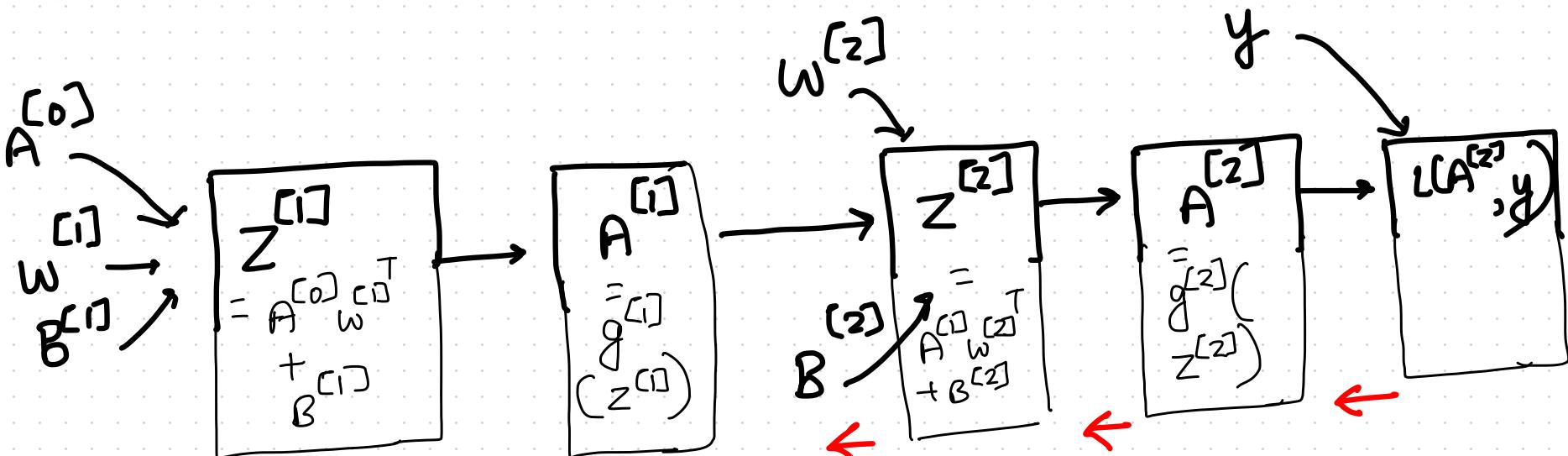
$$= (-y \odot A^{[2]} + (\neg y) \odot (1 - A^{[2]})) \left((A^{[2]} \odot (1 - A^{[2]})) \right)$$

$$= -y \odot A^{[2]} \odot A^{[2]} \odot (1 - A^{[2]}) + (\neg y) \odot (1 - A^{[2]}) \odot A^{[2]} \odot (1 - A^{[2]})$$

$$= -y \odot (1 - A^{[2]}) + (\neg y) \odot A^{[2]} = -y_{N \times 1} + y \odot A^{[2]}_{N \times 1} + A^{[2]}_{N \times 1} - y \odot A^{[2]}_{N \times 1}$$

$$\boxed{\frac{\partial L}{\partial z^{[2]}} = A^{[2]} - y}$$

COMPUTATION GRAPH (FOR XOR EXAMPLE)



$$\frac{\partial L}{\partial W^{[2]}} = \frac{\partial L}{\partial Z^{[2]}} \frac{\partial Z^{[2]}}{\partial W^{[2]}} \in R^{N^{[2]} \times N^{[1]}}$$

(Chain Rule)

SAME DIMENSION ($\because L$ is scalar)

ASIDE

- * GRADIENT: VECTOR IN, SCALAR OUT
- * JACOBIAN : VECTOR IN, VECTOR OUT

$$f: \mathbb{R}^N \rightarrow \mathbb{R}^M$$

$$\text{IIP: } \mathbb{R}^N$$

$$\text{OIP: } \mathbb{R}^M$$

$$y = f(x)$$

Derivative of 'f' at 'x' called Jacobian is $M \times N$ matrix

$$\frac{\partial y}{\partial x} = \begin{pmatrix} \frac{\partial y_1}{\partial x_1} & \dots & \frac{\partial y_1}{\partial x_N} \\ \vdots & & \vdots \\ \frac{\partial y_M}{\partial x_1} & \dots & \frac{\partial y_M}{\partial x_N} \end{pmatrix}_{M \times N}$$

GENERALISED JACOBIAN: TENSOR IN, TENSOR OUT

$$f: \mathbb{R}^{N_1 \times \dots \times N_{Dx}} \rightarrow \mathbb{R}^{M_1 \times \dots \times M_{Dy}}$$

I/P: D_x - dimensional tensor of shape $N_1 \times \dots \times N_{Dx}$

O/P: D_y - dimensional tensor of shape $M_1 \times \dots \times M_{Dy}$

$$y = f(x)$$

Then; $\frac{\partial y}{\partial x} = \text{Gen. Jacobian} = (M_1 \times \dots \times M_{Dy}) \times (N_1 \times \dots \times N_{Dx})$
shape

BACK PROP. WITH TENSORS

Let $G_1 = \alpha \beta$

$\alpha: N \times D$

$\beta: D \times M$

$G_1: N \times M$

BACK PROP. WITH TENSORS

Let $G_1 = \alpha \beta$

$\alpha: N \times D$

$\beta: D \times M$

$G_1: N \times M$

→ We are given Loss L as function of G_1

→ we also know $\frac{\partial L}{\partial G_1}$

BACK PROP. WITH TENSORS

Let $G_1 = \alpha \beta$

$$\alpha : N \times D$$

$$\beta : D \times M$$

$$G_1 : N \times M$$

→ We are given Loss L as function of G_1

→ we also know $\frac{\partial L}{\partial G_1}$

→ Q: $\frac{\partial L}{\partial \alpha} = ?$; $\frac{\partial L}{\partial \beta} = ?$

BACK PROP. WITH TENSORS

Let $G_1 = \alpha \beta$

$$\alpha : N \times D$$

$$\beta : D \times M$$

$$G_1 : N \times M$$

Let's choose: $N=1; D=2; M=3$

BACK PROP. WITH TENSORS

Let $G_1 = \alpha \beta$

$$\alpha : N \times D$$

$$\beta : D \times M$$

$$G_1 : N \times M$$

Let's choose: $N=1; D=2; M=3$

$$G_{1 \times 3} = \begin{pmatrix} G_{1,1} & G_{1,2} & G_{1,3} \end{pmatrix}$$

$$\alpha_{1 \times 2} = \begin{pmatrix} \alpha_{1,1} & \alpha_{1,2} \end{pmatrix}$$

$$\beta_{2 \times 3} = \begin{pmatrix} \beta_{1,1} & \beta_{1,2} & \beta_{1,3} \\ \beta_{2,1} & \beta_{2,2} & \beta_{2,3} \end{pmatrix}$$

BACK PROP. WITH TENSORS

$$G_{1 \times 3} = \begin{pmatrix} G_{1,1} & G_{1,2} & G_{1,3} \end{pmatrix}$$

$$\alpha_{1 \times 2} = \begin{pmatrix} \alpha_{1,1} & \alpha_{1,2} \end{pmatrix}$$

$$\beta_{2 \times 3} = \begin{pmatrix} \beta_{1,1} & \beta_{1,2} & \beta_{1,3} \\ \beta_{2,1} & \beta_{2,2} & \beta_{2,3} \end{pmatrix}$$

$$\alpha \beta = \left[\alpha_{1,1} \beta_{1,1} + \alpha_{1,2} \beta_{2,1} ; \alpha_{1,1} \beta_{1,2} + \alpha_{1,2} \beta_{2,2} ; \alpha_{1,1} \beta_{1,3} + \alpha_{1,2} \beta_{2,3} \right]$$

BACK PROP. WITH TENSORS

$$G_{1 \times 3} = \begin{pmatrix} G_{1,1} & G_{1,2} & G_{1,3} \end{pmatrix}$$

$$\alpha_{1 \times 2} = \begin{pmatrix} \alpha_{1,1} & \alpha_{1,2} \end{pmatrix}$$

$$\beta_{2 \times 3} = \begin{pmatrix} \beta_{1,1} & \beta_{1,2} & \beta_{1,3} \\ \beta_{2,1} & \beta_{2,2} & \beta_{2,3} \end{pmatrix}$$

$$\alpha\beta = \left[\alpha_{1,1} \beta_{1,1} + \alpha_{1,2} \beta_{2,1} ; \alpha_{1,1} \beta_{1,2} + \alpha_{1,2} \beta_{2,2} ; \alpha_{1,1} \beta_{1,3} + \alpha_{1,2} \beta_{2,3} \right]$$

$$\frac{\partial(\alpha\beta)}{\partial \alpha_{1,1}} = [\beta_{1,1}; \beta_{1,2}; \beta_{1,3}]$$

$$\frac{\partial(\alpha\beta)}{\partial \alpha_{1,2}} = [\beta_{2,1}; \beta_{2,2}; \beta_{2,3}]$$

BACK PROP. WITH TENSORS

$$\frac{\partial(\alpha\beta)}{\partial \alpha_{1,1}} = [\beta_{1,1} ; \beta_{1,2} ; \beta_{1,3}]$$

$$\frac{\partial(\alpha\beta)}{\partial \alpha_{1,2}} = [\beta_{2,1} ; \beta_{2,2} ; \beta_{2,3}]$$

Shape = Shape of $\alpha\beta \times$ Shape of
 $\alpha_{1,1}$

$$= (N \times M) \times 1$$

$$\text{Shape} = (N \times M) \times 1$$

BACK PROP. WITH TENSORS

$$\frac{\partial(\alpha\beta)}{\partial \alpha_{1,1}} = [\beta_{1,1} ; \beta_{1,2} ; \beta_{1,3}]$$

$$\frac{\partial(\alpha\beta)}{\partial \alpha_{1,2}} = [\beta_{2,1} ; \beta_{2,2} ; \beta_{2,3}]$$

Shape = Shape of $\alpha\beta \times$ Shape of $\alpha_{1,1}$

$$= (N \times M) \times 1$$

Shape = $(N \times M) \times 1$

Generalised
Tensor shape

BACK PROP. WITH TENSORS

$$\frac{\partial(\alpha\beta)}{\partial \alpha_{1,1}} = [\beta_{1,1}; \beta_{1,2}; \beta_{1,3}]$$

$$\frac{\partial(\alpha\beta)}{\partial \alpha_{1,2}} = [\beta_{2,1}; \beta_{2,2}; \beta_{2,3}]$$

We know $\frac{\partial L}{\partial G} = [dG_{1,1} \quad dG_{1,2} \quad dG_{1,3}]$

Generalised shape = $1 \times (N \times M)$

Why?

BACK PROP. WITH TENSORS

$$\frac{\partial(\alpha\beta)}{\partial \alpha_{1,1}} = [\beta_{1,1}; \beta_{1,2}; \beta_{1,3}]$$

$$\frac{\partial(\alpha\beta)}{\partial \alpha_{1,2}} = [\beta_{2,1}; \beta_{2,2}; \beta_{2,3}]$$

We know $\frac{\partial L}{\partial g} = [d\mathcal{L}_{1,1} \quad d\mathcal{L}_{1,2} \quad d\mathcal{L}_{1,3}]$

Generalised
shape =

$1 \times (N \times M)$



$\therefore L$ is a
scalar

BACK PROP. WITH TENSORS

$$\frac{\partial(\alpha\beta)}{\partial \alpha_{1,1}} = [\beta_{1,1}; \beta_{1,2}; \beta_{1,3}]$$

$$\frac{\partial(\alpha\beta)}{\partial \alpha_{1,2}} = [\beta_{2,1}; \beta_{2,2}; \beta_{2,3}]$$

we know $\frac{\partial L}{\partial g} = [dg_{1,1} \quad dg_{1,2} \quad dg_{1,3}]$

$$\frac{\partial L}{\partial \alpha} = \left[\begin{array}{cc} \frac{\partial L}{\partial \alpha_{1,1}} & \frac{\partial L}{\partial \alpha_{1,2}} \end{array} \right]$$

BACK PROP. WITH TENSORS

$$\frac{\partial(\alpha\beta)}{\partial \alpha_{1,1}} = [\beta_{1,1}; \beta_{1,2}; \beta_{1,3}] = \underline{\frac{\partial G}{\partial \alpha_{1,1}}}$$

$$\frac{\partial(\alpha\beta)}{\partial \alpha_{1,2}} = [\beta_{2,1}; \beta_{2,2}; \beta_{2,3}] = \underline{\frac{\partial G}{\partial \alpha_{1,2}}}$$

we know $\frac{\partial L}{\partial G} = [\underline{dG_{1,1}} \quad \underline{dG_{1,2}} \quad \underline{dG_{1,3}}]$

$$\frac{\partial L}{\partial \alpha} = \left[\frac{\partial L}{\partial \alpha_{1,1}} \quad \frac{\partial L}{\partial \alpha_{1,2}} \right]$$

$$\frac{\partial L}{\partial \alpha_{1,1}} = \frac{\partial L}{\partial G} \cdot \underline{\frac{\partial G}{\partial \alpha_{1,1}}}$$

BACK PROP. WITH TENSORS

$$\frac{\partial(\alpha\beta)}{\partial \alpha_{1,1}} = [\beta_{1,1}; \beta_{1,2}; \beta_{1,3}] = \underline{\frac{\partial G}{\partial \alpha_{1,1}}}$$

$$\frac{\partial(\alpha\beta)}{\partial \alpha_{1,2}} = [\beta_{2,1}; \beta_{2,2}; \beta_{2,3}] = \underline{\frac{\partial G}{\partial \alpha_{1,2}}}$$

we know $\frac{\partial L}{\partial G} = [\delta G_{1,1} \quad \delta G_{1,2} \quad \delta G_{1,3}]$

$$\frac{\partial L}{\partial \alpha} = \left[\frac{\partial L}{\partial \alpha_{1,1}} \quad \frac{\partial L}{\partial \alpha_{1,2}} \right]$$

$$\frac{\partial L}{\partial \alpha_{1,1}} = \frac{\partial L}{\partial G} \cdot \underline{\frac{\partial G}{\partial \alpha_{1,1}}} \xrightarrow{1 \times (N \times m)} \xrightarrow{(N \times m) \times 1}$$

BACK PROP. WITH TENSORS

$$\frac{\partial(\alpha\beta)}{\partial \alpha_{1,1}} = [\beta_{1,1}; \beta_{1,2}; \beta_{1,3}] = \underline{\frac{\partial G}{\partial \alpha_{1,1}}}$$

$$\frac{\partial(\alpha\beta)}{\partial \alpha_{1,2}} = [\beta_{2,1}; \beta_{2,2}; \beta_{2,3}] = \underline{\frac{\partial G}{\partial \alpha_{1,2}}}$$

We know $\frac{\partial L}{\partial G} = [dG_{1,1} \quad dG_{1,2} \quad dG_{1,3}]$

$$\frac{\partial L}{\partial \alpha} = \left[\frac{\partial L}{\partial \alpha_{1,1}} \quad \frac{\partial L}{\partial \alpha_{1,2}} \right]$$

$$\frac{\partial L}{\partial \alpha_{1,1}} = \frac{\partial L}{\partial G} \cdot \underline{\frac{\partial G}{\partial \alpha_{1,1}}} = [dG_{1,1} \quad \beta_{1,1} + dG_{1,2}\beta_{1,2} + dG_{1,3}\beta_{1,3}]$$

$(1 \times (N \times M)) \times ((N \times M) \times 1) \rightarrow$ Dot product of two.

BACK PROP. WITH TENSORS

$$\frac{\partial L}{\partial \alpha} = \left[\frac{\partial L}{\partial \alpha_{1,1}} ; \frac{\partial L}{\partial \alpha_{1,2}} \right] = \begin{pmatrix} \frac{\partial L}{\partial G} \end{pmatrix}^T \beta^{1 \times (N \times M)} \quad m \times D$$

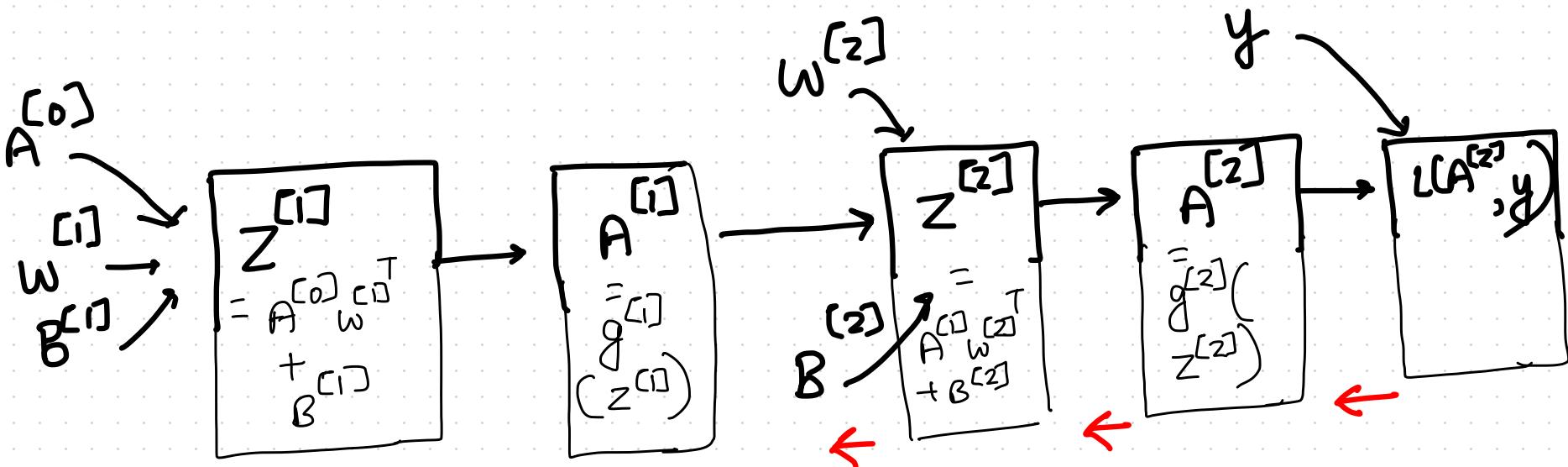
BACK PROP. WITH TENSORS

$$\frac{\partial L}{\partial \alpha} = \left[\frac{\partial L}{\partial \alpha_{1,1}} ; \frac{\partial L}{\partial \alpha_{1,2}} \right] = \begin{pmatrix} \frac{\partial L}{\partial G} \\ \vdots \end{pmatrix}^T \quad \begin{matrix} 1 \times (N \times M) \\ \beta \end{matrix} \quad \begin{matrix} M \times D \end{matrix}$$

Similarly;

$$\frac{\partial L}{\partial \beta} = \alpha^T \quad \begin{matrix} D \times N \\ 1 \times (D \times M) \end{matrix} \quad \frac{\partial L}{\partial G} \quad \begin{matrix} 1 \times (N \times M) \end{matrix}$$

COMPUTATION GRAPH (FOR XOR EXAMPLE)



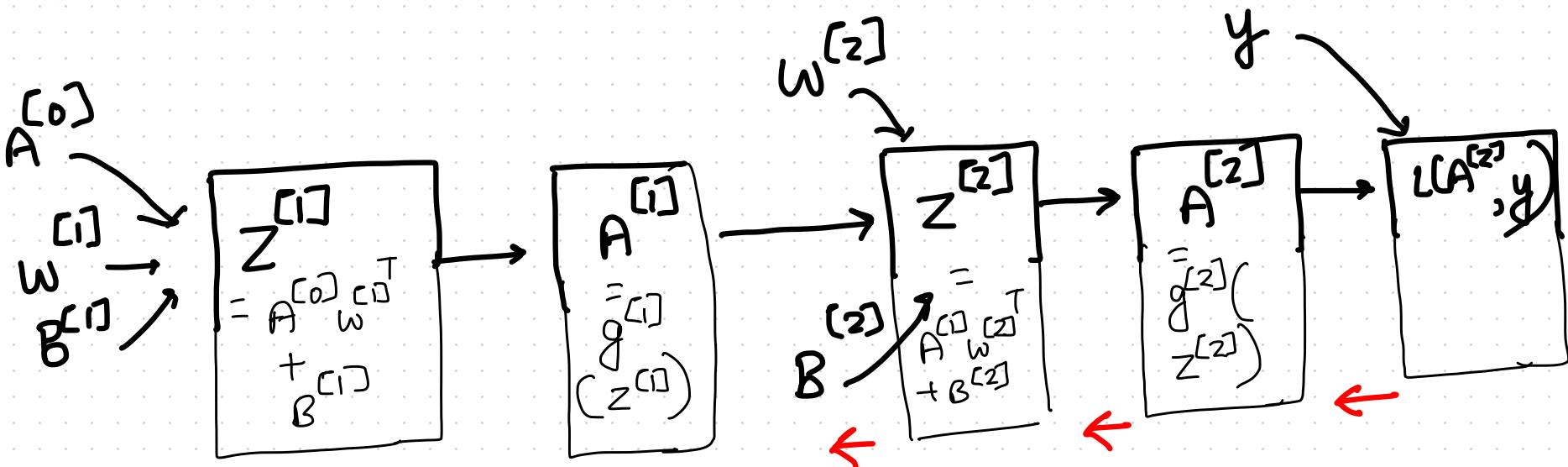
Equivalence from Aside

$$G \leftrightarrow Z^{[2]}$$

$$\alpha \leftrightarrow A^{[1]}$$

$$\beta \leftrightarrow W^{[2]T}$$

COMPUTATION GRAPH (FOR XOR EXAMPLE)

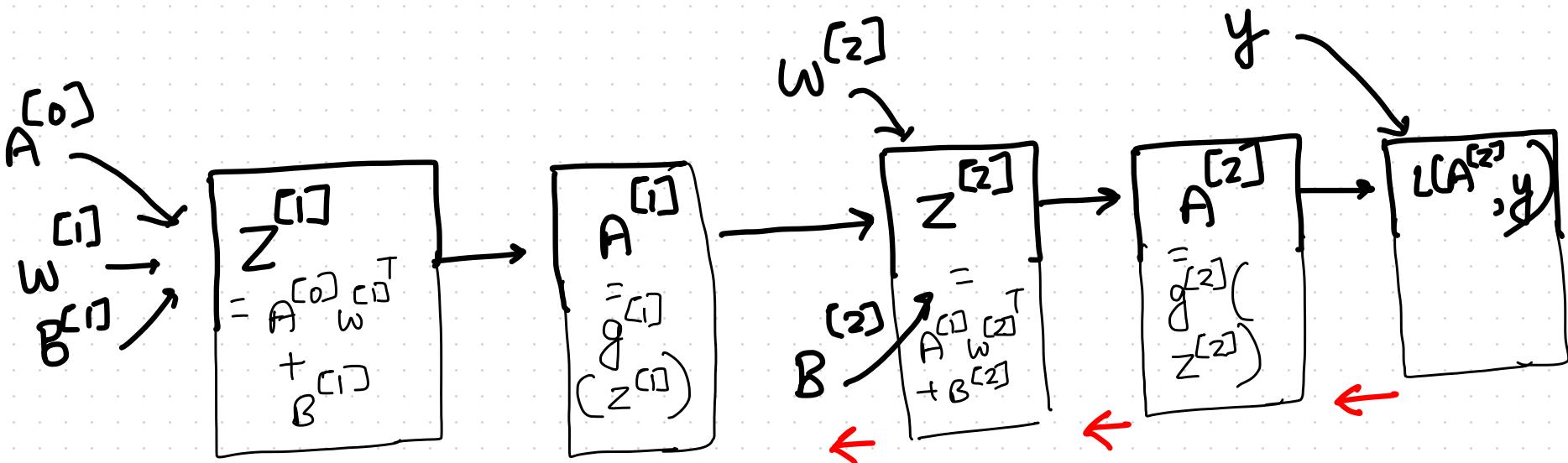


Equivalence from Aside

$$\begin{aligned} G &\leftrightarrow Z^{[2]} \\ \alpha &\leftrightarrow A^{[1]} \\ \beta &\leftrightarrow W^{[2] T} \end{aligned}$$

$$\begin{aligned} \frac{\partial L}{\partial W^{[2] T}} &= A^{[1] T} \frac{\partial L}{\partial Z^{[2]}} \\ &= A^{[1] T} (A^{[2]} - y) \end{aligned}$$

COMPUTATION GRAPH (FOR XOR EXAMPLE)



$$\frac{\partial L}{\partial W^{[2]}} = A^{[1]^T} (A^{[2]} - y)$$

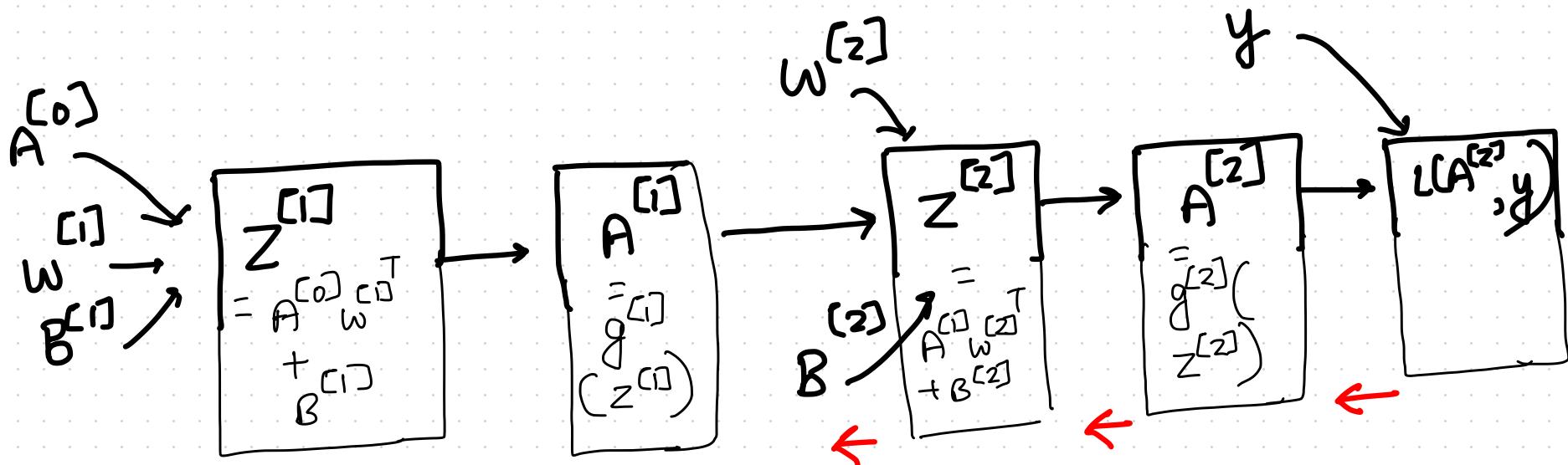
$$\Rightarrow \frac{\partial L}{\partial W^{[2]}} = (A^{[2]} - y)^T A^{[1]}$$

$N^{[2]} \times N^{[1]}$

$N \times N^{[1]}$

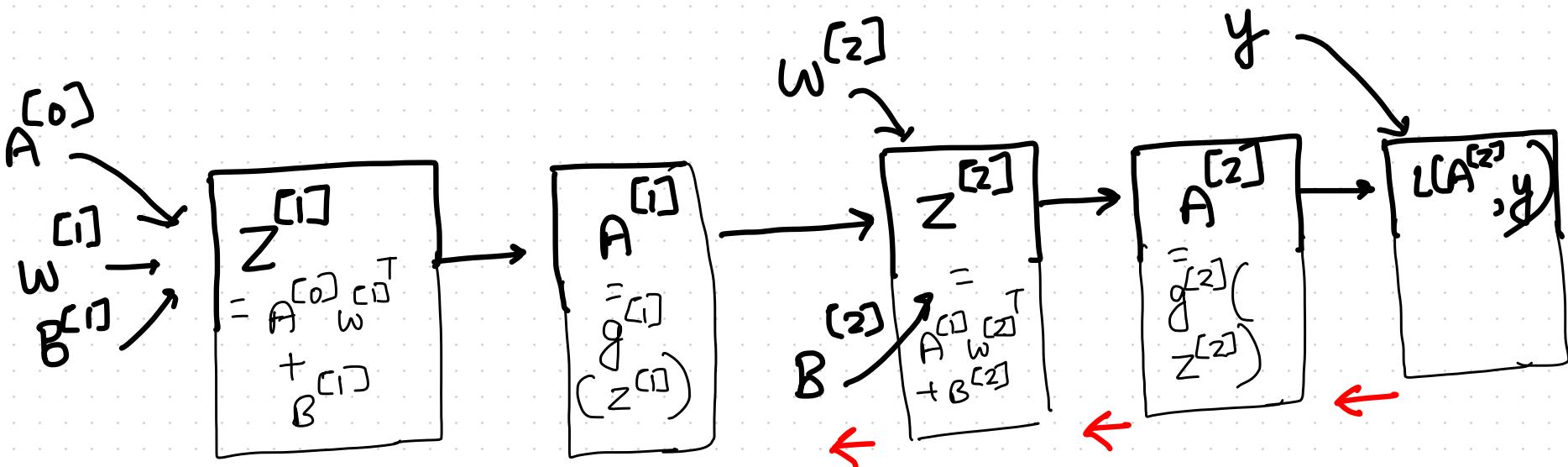
$I \times (N^{[2]} \times N^{[1]})$

COMPUTATION GRAPH (FOR XOR EXAMPLE)



$$\frac{\partial L}{\partial B^{[2]}} = ?$$

COMPUTATION GRAPH (FOR XOR EXAMPLE)

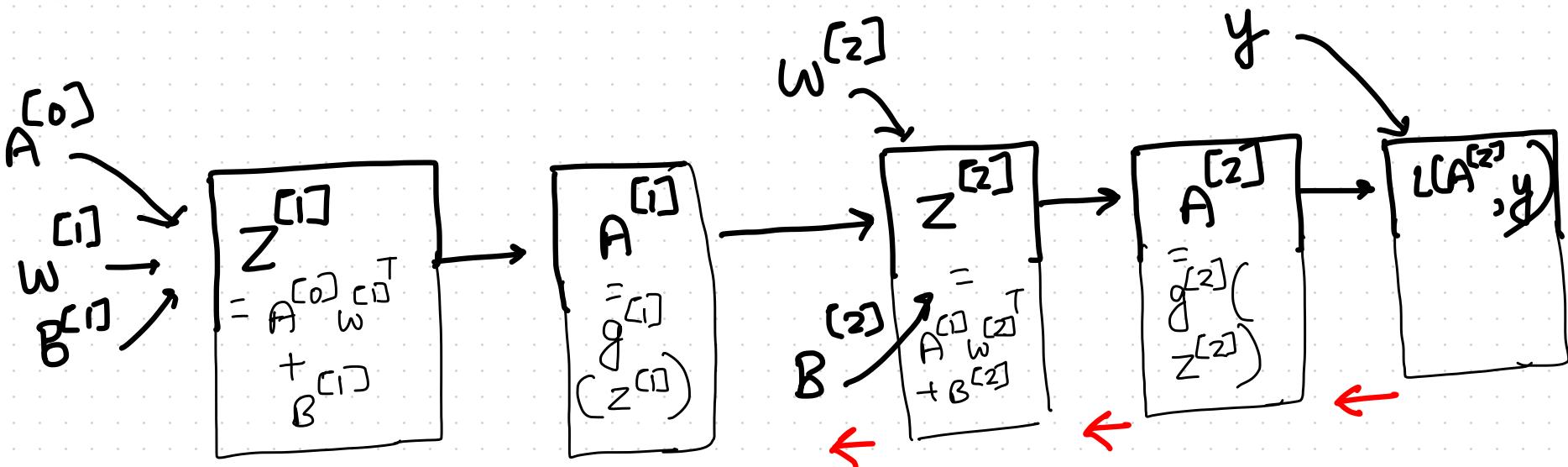


$$\frac{\partial L}{\partial B^{[2]}} = ?$$

$B^{[2]} = b^{[2]}$

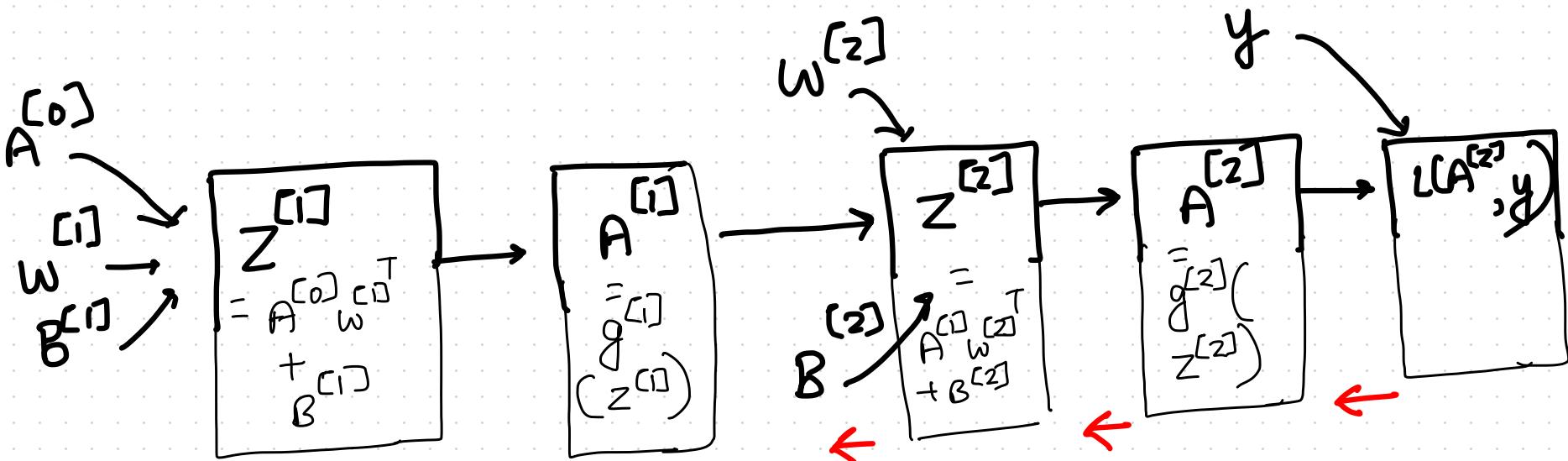
$$B^{[2]} = \begin{bmatrix} -c- \\ -c- \\ \vdots \end{bmatrix}$$

COMPUTATION GRAPH (FOR XOR EXAMPLE)



$$\frac{\partial L}{\partial c} = ?$$

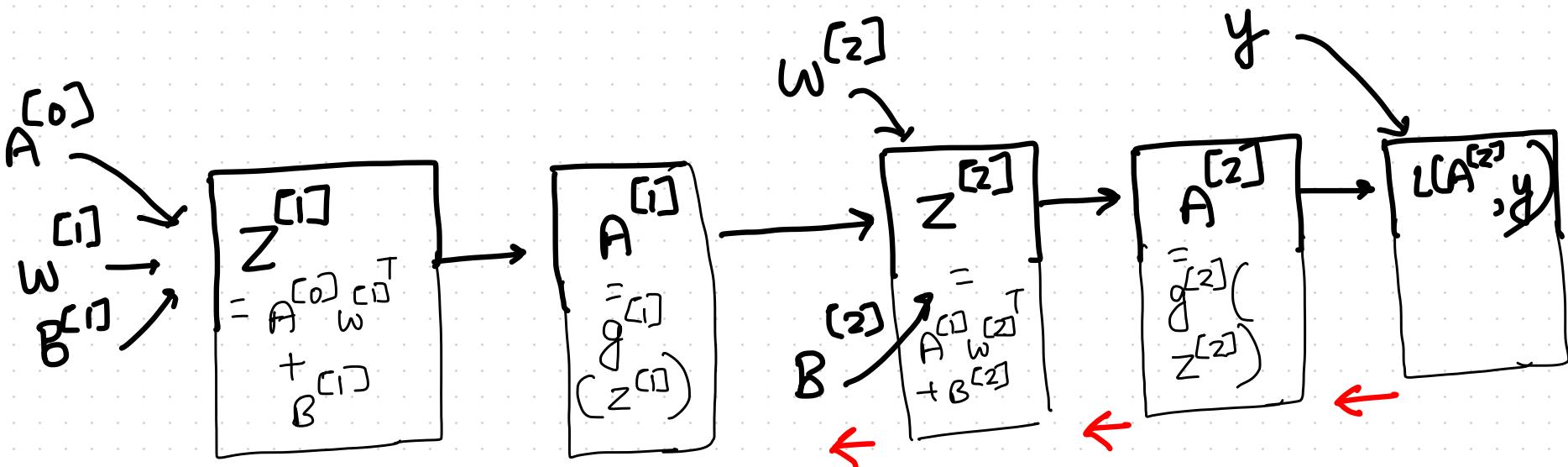
COMPUTATION GRAPH (FOR XOR EXAMPLE)



$$\frac{\partial L}{\partial c} = \begin{bmatrix} \frac{\partial L}{\partial c_1} & \dots & \frac{\partial L}{\partial c_{N^{[2]}}} \end{bmatrix}$$

$$= \left[\frac{\partial L}{\partial z^{[2]}} \cdot \frac{\partial z^{[2]}}{\partial c_1} \dots \right]$$

COMPUTATION GRAPH (FOR XOR EXAMPLE)

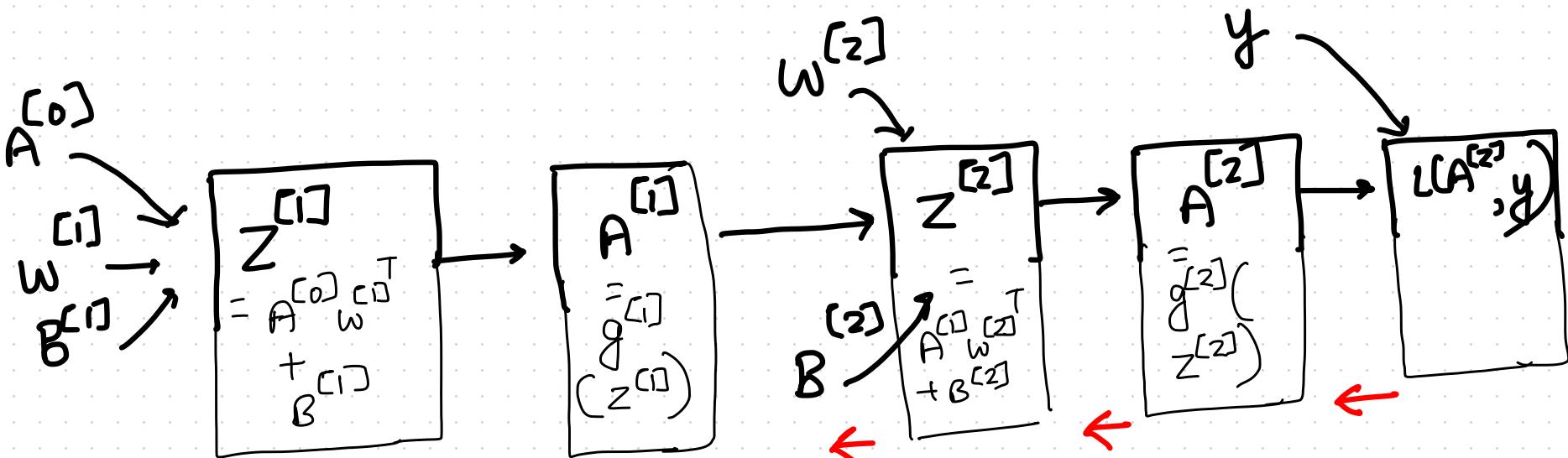


$$\frac{\partial L}{\partial c_1} = \frac{\partial L}{\partial Z^{[2]}} \cdot \frac{\partial Z^{[2]}}{\partial c_1}$$

$$Z^{[2]} = \begin{bmatrix} z_{1,1}^{[2]} & \dots & z_{1,N^{[2]}}^{[2]} \\ \vdots & & \vdots \\ z_{N,1}^{[2]} & \dots & z_{N,N^{[2]}}^{[2]} \end{bmatrix} = A^{[2]} W^{[2] \top} + \begin{bmatrix} c_1 & c_2 & \dots & c_N^{[2]} \\ c_1 & c_2 & \dots & c_N^{[2]} \\ \vdots & \vdots & \ddots & \vdots \\ c_1 & c_2 & \dots & c_N^{[2]} \end{bmatrix}$$

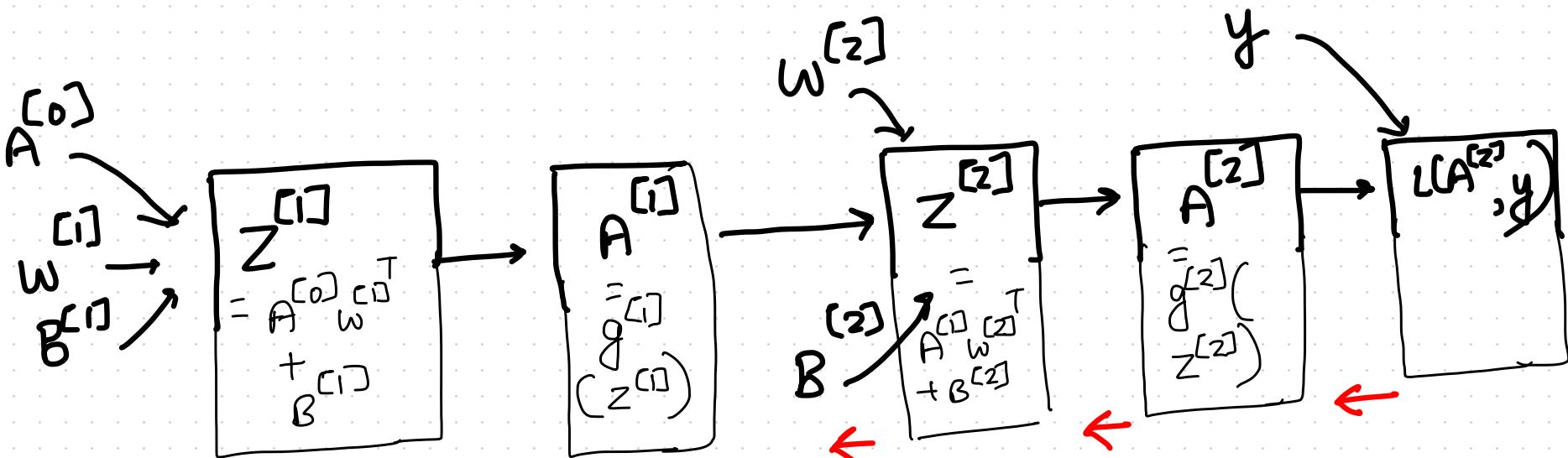
$N \times N^{[2]}$

COMPUTATION GRAPH (FOR XOR EXAMPLE)



$$\frac{\partial Z^{[2]}}{\partial c_1} = \begin{bmatrix} \frac{\partial c_1}{\partial c_1} & \frac{\partial c_2}{\partial c_1} & \dots \\ \vdots & \ddots & \dots \end{bmatrix} = \begin{bmatrix} 1 & 0 & 0 & \dots & 0 \\ \vdots & \ddots & \ddots & \ddots & \vdots \\ 1 & - & - & - & - \end{bmatrix} \quad (N \times N^{[2]}) \times (1)$$

COMPUTATION GRAPH (FOR XOR EXAMPLE)

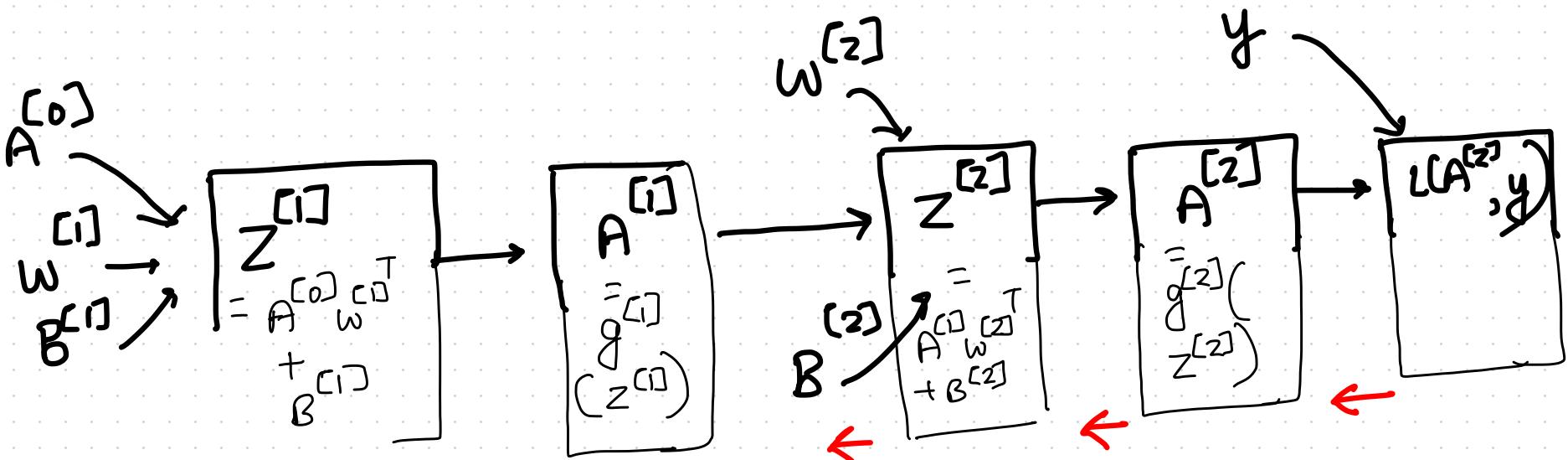


$$\frac{\partial Z}{\partial C_1} = \begin{bmatrix} \frac{\partial C_1}{\partial C_1} & \frac{\partial C_2}{\partial C_1} & \dots \\ \vdots & \ddots & \dots \end{bmatrix} = \begin{bmatrix} 1 & 0 & 0 & \dots & 0 \\ \vdots & \ddots & \ddots & \ddots & \vdots \\ 1 & - & - & - & - \end{bmatrix} \quad (N \times N^{[2]}) \times (1)$$

$$\frac{\partial L}{\partial C_1} = \frac{\partial L}{\partial Z^{[2]}} \frac{\partial Z^{[2]}}{\partial C_1}$$

$(1 \times (N \times N^{[2]})) \quad ((N \times N^{[2]}) \times 1)$

COMPUTATION GRAPH (FOR XOR EXAMPLE)



$$\text{Let } \frac{\partial L}{\partial Z^{[2]}} = \begin{bmatrix} dZ_{1,1}^{[2]} & \dots & dZ_{1,N^{[2]}}^{[2]} \\ dZ_{N,1}^{[2]} & \dots & dZ_{N,N^{[2]}}^{[2]} \end{bmatrix}$$

Then

$$\frac{\partial L}{\partial c_1} = \sum_{i=1}^N dZ_{i,1}^{[2]}$$

COMPUTATION GRAPH (FOR XOR EXAMPLE)

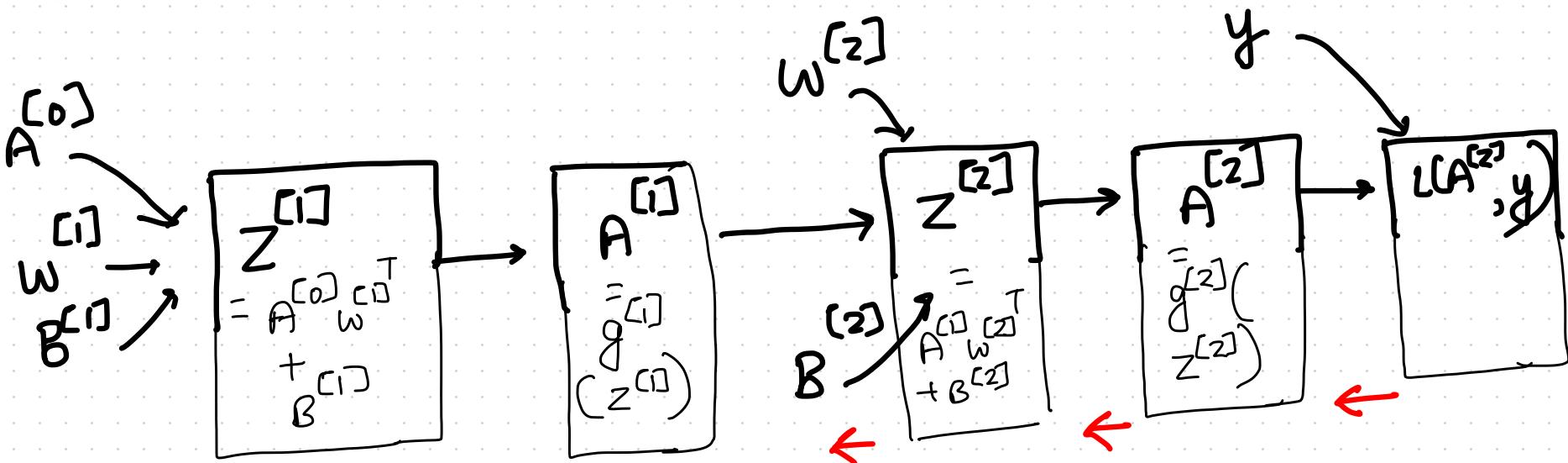
$$\frac{\partial L}{\partial c_1} = \sum_{i=1}^N dz_{i,1}^{[2]}$$

$$\Rightarrow \frac{\partial L}{\partial c} = \left[\sum_{i=1}^N dz_{i,1}^{[2]} \quad ; \quad \sum_{i=1}^N dz_{i,2}^{[2]} \dots ; \sum_{i=1}^N dz_{i,N}^{[2]} \right]$$

$$\Rightarrow \frac{\partial L}{\partial b} = \left(\frac{\partial L}{\partial c} \right)^T = \begin{bmatrix} \sum_{i=1}^N dz_{i,1}^{[2]} \\ \vdots \\ \sum_{i=1}^N dz_{i,N}^{[2]} \end{bmatrix}$$

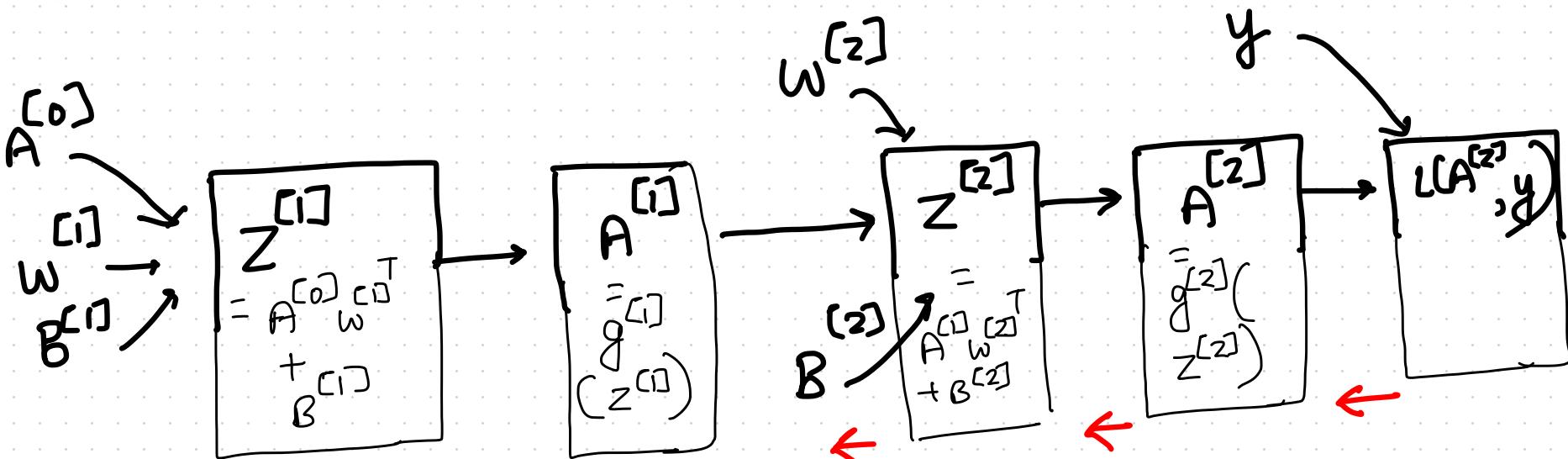
$N^{[2]} \times 1$

COMPUTATION GRAPH (FOR XOR EXAMPLE)



$$\frac{\partial L}{\partial A^{[0]}} = \frac{\partial L}{\partial Z^{[2]}} \frac{\partial Z^{[2]}}{\partial A^{[1]}} = \frac{\partial L}{\partial Z^{[2]}} \left(W^{[2]} \right)^T = \frac{\partial L}{\partial Z^{[2]}} W^{[2]}$$

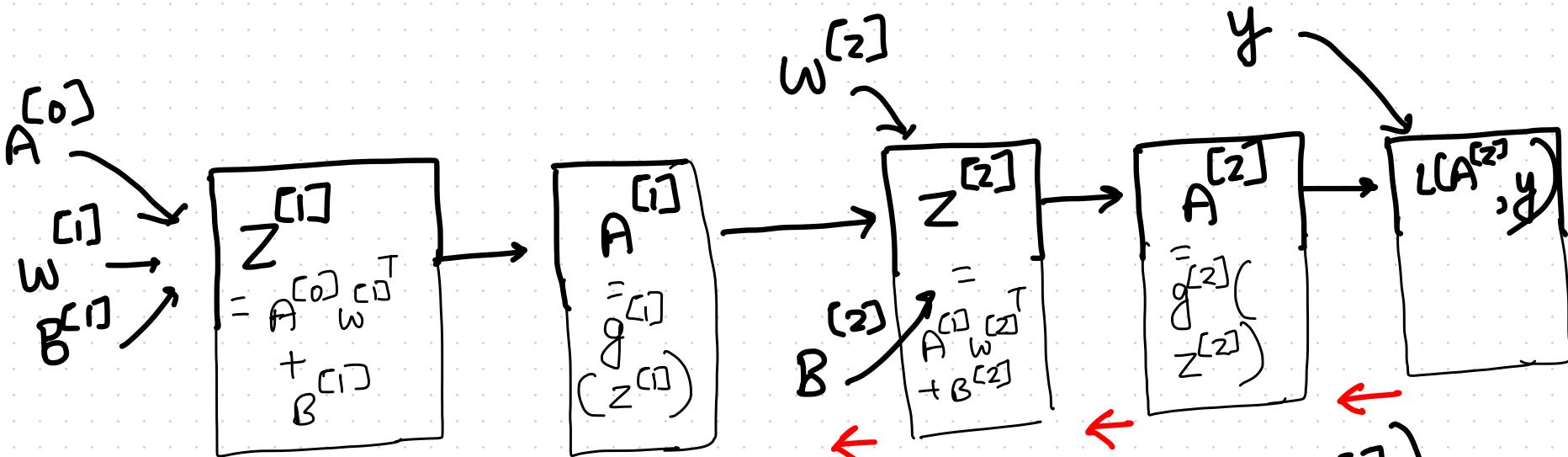
COMPUTATION GRAPH (FOR XOR EXAMPLE)



$$\frac{\partial L}{\partial A^{[1]}} = \frac{\partial L}{\partial Z^{[2]}} \frac{\partial Z^{[2]}}{\partial A^{[1]}} = \frac{\partial L}{\partial Z^{[2]}} (W^{[2]})^T = \frac{\partial L}{\partial Z^{[2]}} \frac{\partial Z^{[2]}}{\partial A^{[0]}}$$

$I \times (N \times N^{[1]})$ $I \times (N \times N^{[2]})$ $N^{[2]} \times N^{[1]}$

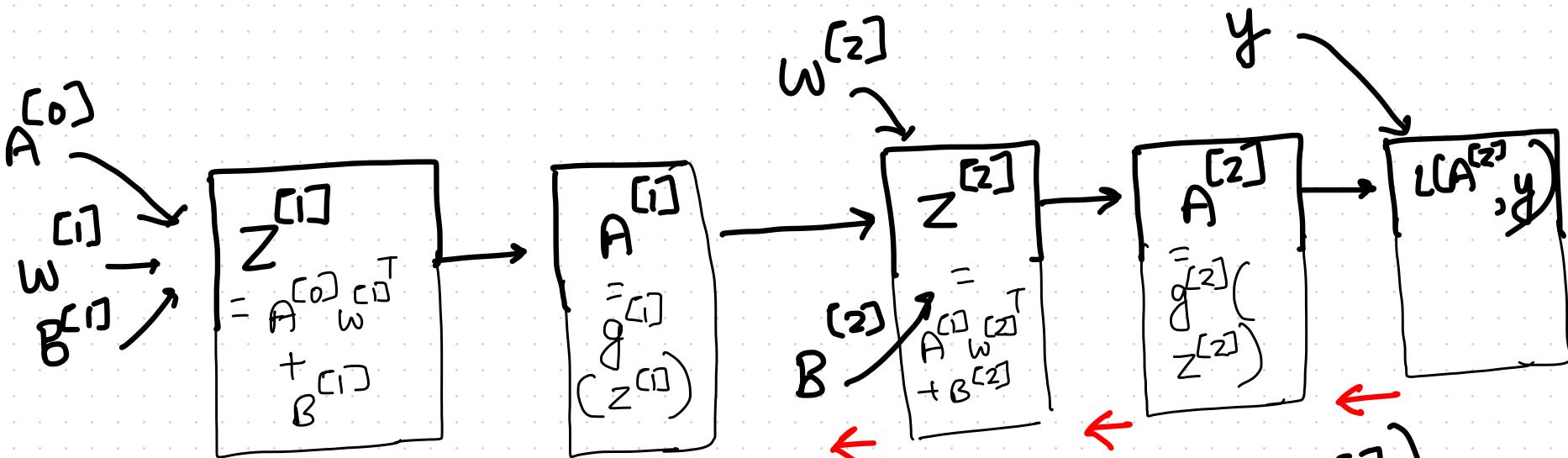
COMPUTATION GRAPH (FOR XOR EXAMPLE)



$$\frac{\partial L}{\partial Z^{[1]}} = \frac{\partial L}{\partial A^{[0]}} \frac{\partial A^{[0]}}{\partial Z^{[1]}} = \left(\frac{\partial L}{\partial Z^{[2]}} W^{[2]} \right) \left(\frac{\partial A^{[1]}}{\partial Z^{[1]}} \right)$$

\downarrow
 $1 \times (N \times N^{[1]})$
 $\times (N \times N^{[2]}) \times (N^{[2]} \times N^{[1]})$
 \times
 $(N \times N^{[1]})$
 $(N \times N)$

COMPUTATION GRAPH (FOR XOR EXAMPLE)



$$\frac{\partial L}{\partial Z^{[2]}} = \frac{\partial L}{\partial A^{[2]}} \frac{\partial A^{[2]}}{\partial Z^{[1]}} = \left(\frac{\partial L}{\partial Z^{[2]}} W^{[2]} \right) \left(\frac{\partial A^{[2]}}{\partial Z^{[1]}} \right)$$

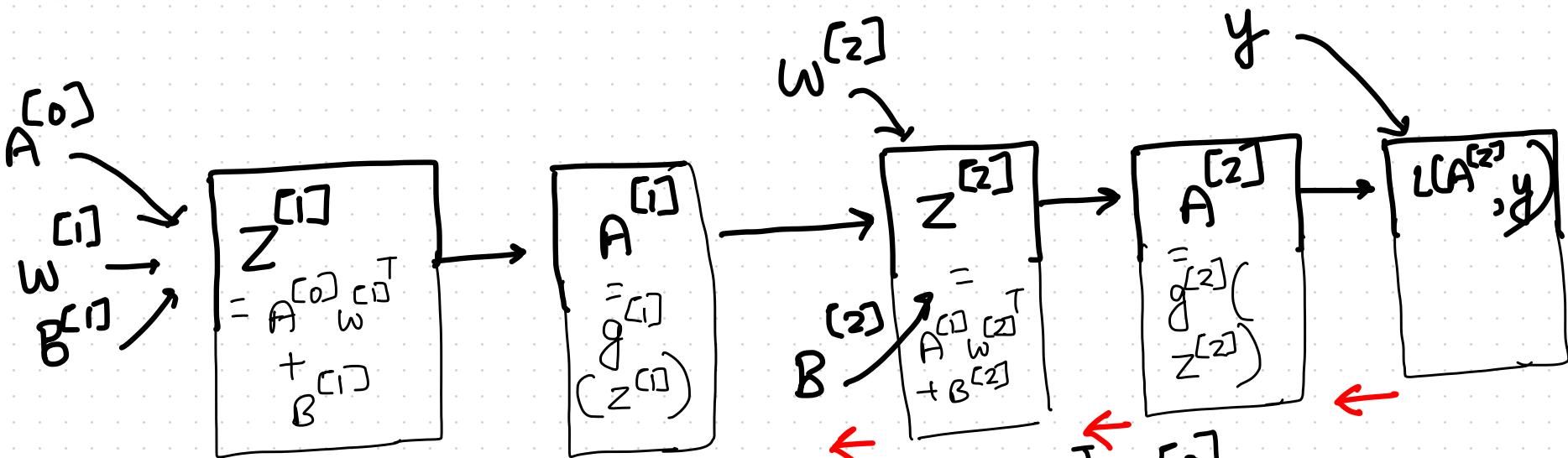
\downarrow
 $1 \times (N \times N^{[1]}) \times (N^{[2]} \times N^{[1]})$

\downarrow
 $(N \times N^{[1]})$
 \times
 $(N \times N^{[1]})$

$$= \frac{\partial L}{\partial Z^{[2]}} W^{[2]} \odot g^{[1]}'(Z^{[1]})$$

\uparrow
 Element wise
 or dot product of Jacobian

COMPUTATION GRAPH (FOR XOR EXAMPLE)



$$\frac{\partial L}{\partial W^{[1]}} = A^{[0]T} \frac{\partial L}{\partial Z^{[1]}} \quad \Rightarrow \quad \frac{\partial L}{\partial W^{[1]}} = \left(\frac{\partial L}{\partial Z^{[1]}} \right)^T A^{[0]}$$

$$\frac{\partial L}{\partial b^{[1]}} = \begin{bmatrix} \sum_{i=1}^N dz_{i,1}^{[1]} \\ \vdots \\ \sum_{i=1}^N dz_{i,N^{[2]}}^{[1]} \end{bmatrix}$$

$N^{[2]} \times 1$