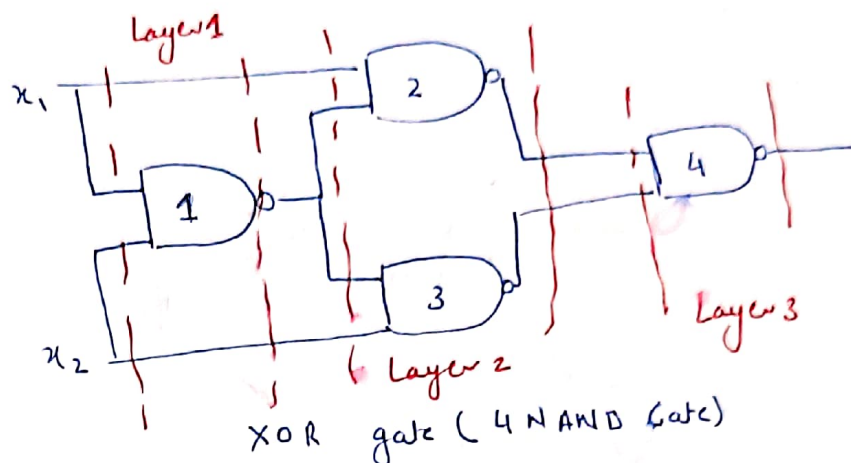**Perceptron** : Drawback : only works for linearly separable data

XOR x

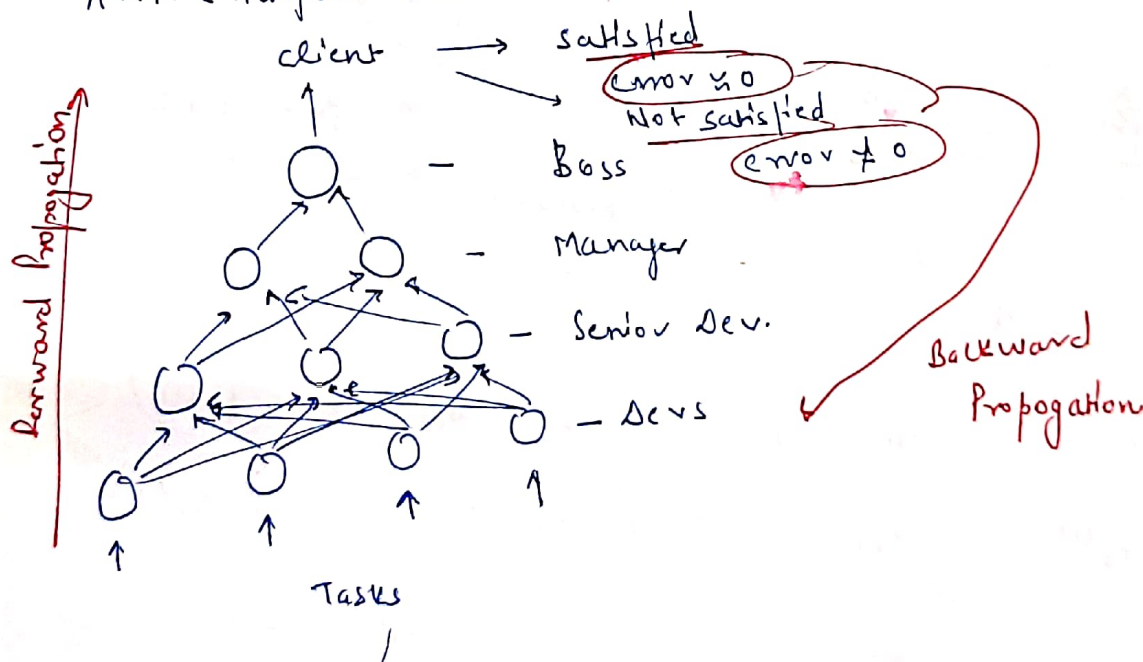① NAND : XOR can also be implemented as a NAND Gate



XOR gate ( 4 NAND Gate )

One NAND gate can be thought of as one perceptron. If we stack multiple perceptron, we can solve a non linear problem as well.

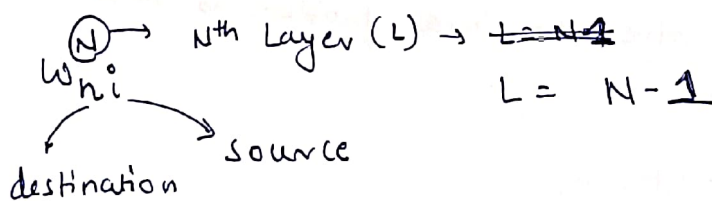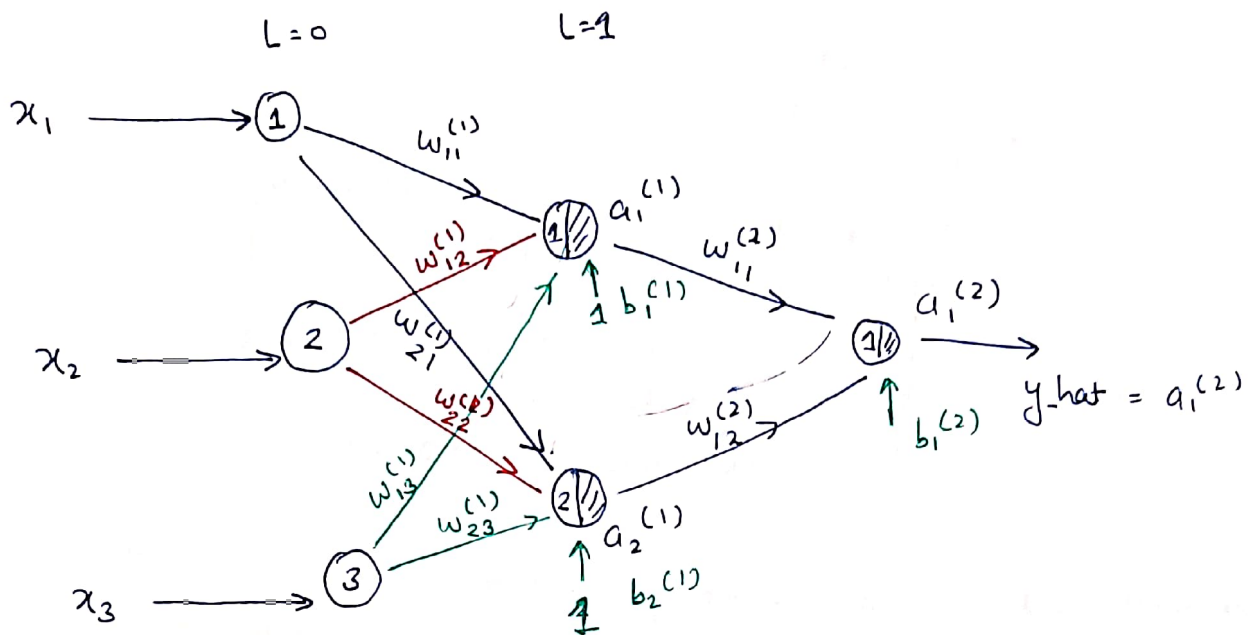(Multilayer perceptron) can solve non linear problem as well.
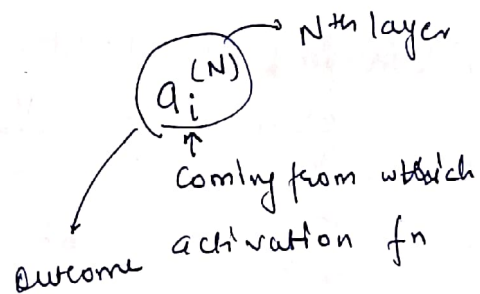↓

ANN ( Artificial Neural Network)



Activation function will filter out the unwanted information and only transmit the most important one to the further layer. There can be many types of activation functions that we can look for.

# * forward Propogation

$N^{th}$ Layer $(L) \rightarrow$ ~~L = N~~

$$L = N - 1$$

$w_{ni}$
destination → source

⊘ activation functions

$a_i^{(N)}$ → $N^{th}$ layer

↑ coming from which activation fn

outcome

## at layer 1

① — $z_1^{(1)} = w_{11}^{(1)} x_1 + w_{12}^{(1)} x_2 + w_{13}^{(1)} x_3 \; \boxed{+ \, b_1^{(1)}}$ → bias

② — $z_2^{(1)} = w_{21}^{(1)} x_1 + w_{22}^{(1)} x_2 + w_{23}^{(1)} x_3 \; \boxed{+ \, b_2^{(1)}}$

$a_1^{(1)} = \sigma(z_1^{(1)}) \quad \approx \quad \begin{bmatrix} 1 & z_1^{(1)} \geq \theta \\ 0 & z_1^{(1)} < \theta \end{bmatrix}$

↳ sigmoid fn

↳ this type can change

$a_2^{(1)} = \sigma(z_2^{(1)})$

final output layer

$$\text{③} - z_1^{(2)} = w_{11}^{(2)} a_1^{(1)} + w_{12}^{(2)} a_2^{(1)} + \boxed{b_1^{(2)}}$$

↘ adding bias.

$$a_1^{(2)} = \sigma(z_1^{(2)}) \longrightarrow \hat{y}$$

$$\text{error}/\text{loss} = y - \hat{y}$$

Cost function

$$\text{cost}(y, \hat{y})$$

## Back Propogation → weight update Rule

$$w = w + \Delta w \longrightarrow \text{change in}$$

$$\Delta w = -\eta \left(\boxed{\dfrac{\partial \text{cost}}{\partial w}}\right)$$

cost function
with respect to
small change in
weights.

learning rate

Let's introduced the bias.

$$z_1^{(1)} = w_{11}^{(1)} x_1 + w_{12}^{(1)} x_2 + w_{13}^{(1)} x_3 + b_1^{(1)}$$

$$z_2^{(1)} = w_{21}^{(1)} x_1 + w_{22}^{(1)} x_2 + w_{23}^{(1)} x_3 + b_2^{(1)}$$

$$z_1^{(2)} = w_{11}^{(2)} a_1^{(1)} + w_{12}^{(2)} a_2^{(1)} + b_1^{(2)}$$

Matrix form of same things

from ① & ② for layer ①

$$\begin{bmatrix} w_{11}^{(1)} & w_{12}^{(1)} & w_{13}^{(1)} \\ w_{21}^{(1)} & w_{22}^{(1)} & w_{23}^{(1)} \end{bmatrix}_{2 \times 3} \begin{bmatrix} x_1 \\ x_2 \\ x_3 \end{bmatrix}_{3 \times 1} + \begin{bmatrix} b_1^{(1)} \\ b_2^{(1)} \end{bmatrix}_{2 \times 1} = \begin{bmatrix} z_1^{(1)} \\ z_2^{(1)} \end{bmatrix}_{2 \times 1}$$

$$\begin{bmatrix} z_1^{(1)} \\ z_2^{(1)} \end{bmatrix}_{2\times1} \xrightarrow[\text{function}]{\text{activation}} \begin{bmatrix} \sigma(z_1^{(1)}) \\ \sigma(z_2^{(1)}) \end{bmatrix} = \begin{bmatrix} a_1^{(1)} \\ a_2^{(1)} \end{bmatrix}_{2\times1}$$

At final layer

$$\begin{bmatrix} w_{11}^{(2)} & w_{12}^{(2)} \end{bmatrix}_{1\times2} \begin{bmatrix} a_1^{(1)} \\ a_2^{(1)} \end{bmatrix}_{2\times1} + \begin{bmatrix} b_1^{(2)} \end{bmatrix}_{1\times1} = \begin{bmatrix} z_1^{(2)} \end{bmatrix}_{1\times1}$$

$$\downarrow \text{activation fn}$$

$$\hat{y} \leftarrow \begin{bmatrix} a_1^{(2)} \end{bmatrix} \leftarrow \begin{bmatrix} \sigma(z_1^{(2)}) \end{bmatrix}_{1\times1}$$

 Q Why do we need an
    Activation function?

① It helps achieve non linearity as there are variety of activation
   functions available

> 1> sigmoid fn          $\sigma(x) = \dfrac{1}{1 + e^{-x}}$

> 2) tanh(x)            $\dfrac{e^x - e^{-x}}{e^x + e^{-x}}$

> 3) ReLU(x)            $\max(x, 0)$
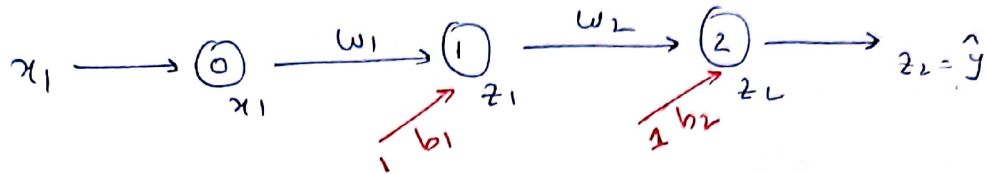
② It helps us to converge the solution

> 1> $\sigma(x)$           Range Input $(-\infty, \infty)$ : domain
>                          outcome/Range $(0, 1)$
>
>     helps reduce the solution space

2)    tanh $(x)$      domain $(-\infty, \infty)$      Range $(-1, 1)$

3)    ReLu $(x)$      domain $(-\infty, \infty)$      Range $(0, \infty)$
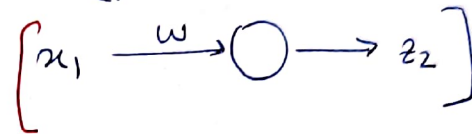
Assumption ( No activation function available)

$$x_1 \longrightarrow \underset{x_1}{\textcircled{0}} \xrightarrow{w_1} \underset{\underset{b_1}{\nearrow} z_1}{\textcircled{1}} \xrightarrow{w_2} \underset{\underset{b_2}{\nearrow} z_2}{\textcircled{2}} \longrightarrow z_2 = \hat{y}$$

$z_1 = w_1 x_1$             $z_2 = \dfrac{w x_1}{\Downarrow}$

$z_2 = w_2 z_1$

$z_2 = w_2 w_1 x_1$         $\left[ x_1 \xrightarrow{w} \bigcirc \longrightarrow z_2 \right]$

                        $\hookrightarrow$ entire network boiled down

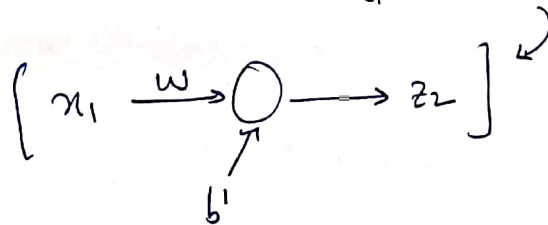                             to one single neuron

\*   No benefit of multilayer classification or multilayer effect without an activation function.

\* Lets <u>introduce the bias</u> without activation function

$z_1 = w_1 x_1 + b_1$

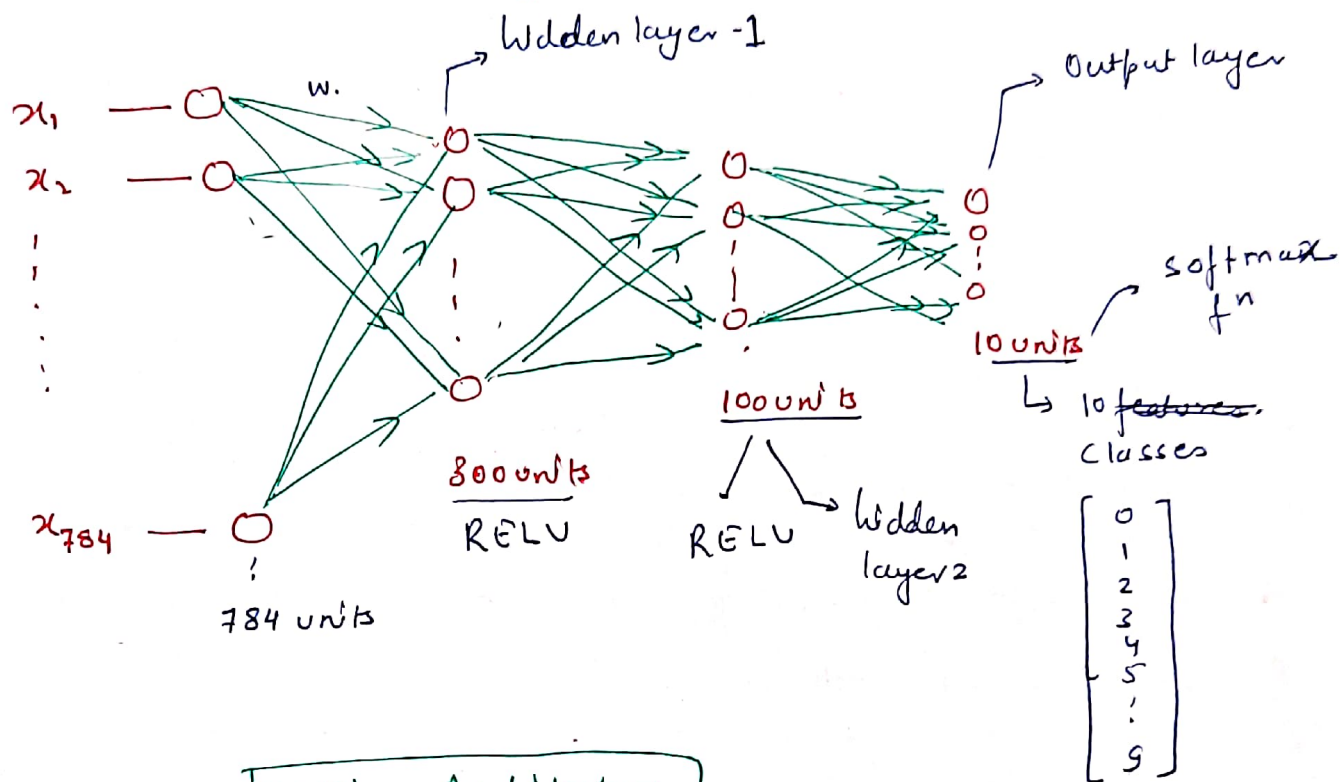$z_2 = w_2 z_1 + b_2 \quad = \quad w_2 w_1 x_1 + w_2 b_1 + b_2$

                            $\underset{\sqcup}{\quad} \quad w x_1 + b'$

$$\left[ x_1 \xrightarrow{w} \underset{\underset{b'}{\nearrow}}{\bigcirc} \longrightarrow z_2 \right]^{\curvearrowleft}$$

# MNIST data analysis

1 Data point – 28×28 $\longrightarrow$ 1D array
2D array $\downarrow$ 784

flattening
operation
1st layer flattening layer



Hidden layer-1

Output layer

$x_1$, $x_2$, ..., $x_{784}$

784 units

800 units
RELU

100 units
RELU → Hidden layer 2

10 units

softmax fn

↳ 10 features classes

$$\begin{bmatrix} 0 \\ 1 \\ 2 \\ 3 \\ 4 \\ 5 \\ \vdots \\ 9 \end{bmatrix}$$

---

ANN Architecture

fully connected ANN where every datapoint is connected to other data point.

RELU   $\max(x, 0)$

Softmax fn
↑
multi-class classification with probability distinction.

Total number of weights = $784 \times 300 + 300$ (bias units)

= 235500