# *ANNs → artificial neural networks :-



$$Z = W_0 + W_1 x_1 + W_2 x_2$$

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

activation

$$W = \begin{pmatrix} W_0 \\ W_1 \\ W_2 \end{pmatrix}$$

OVR → one vs rest

$$W^{(1)} \cdot \begin{pmatrix} W_{01} \\ W_{11} \\ W_{21} \end{pmatrix}$$

$\sigma(Z_1)$

$$W^{(2)} \cdot \begin{pmatrix} W_{02} \\ W_{12} \\ W_{22} \end{pmatrix}$$

$\sigma(Z_2)$

neurons
work independently

$\sigma(Z_3)$

Linear Classifier

$$W^{(3)} = \begin{pmatrix} W_{03} \\ W_{13} \\ W_{23} \end{pmatrix}$$

# * Connecting neurons together — 3 Layer networks
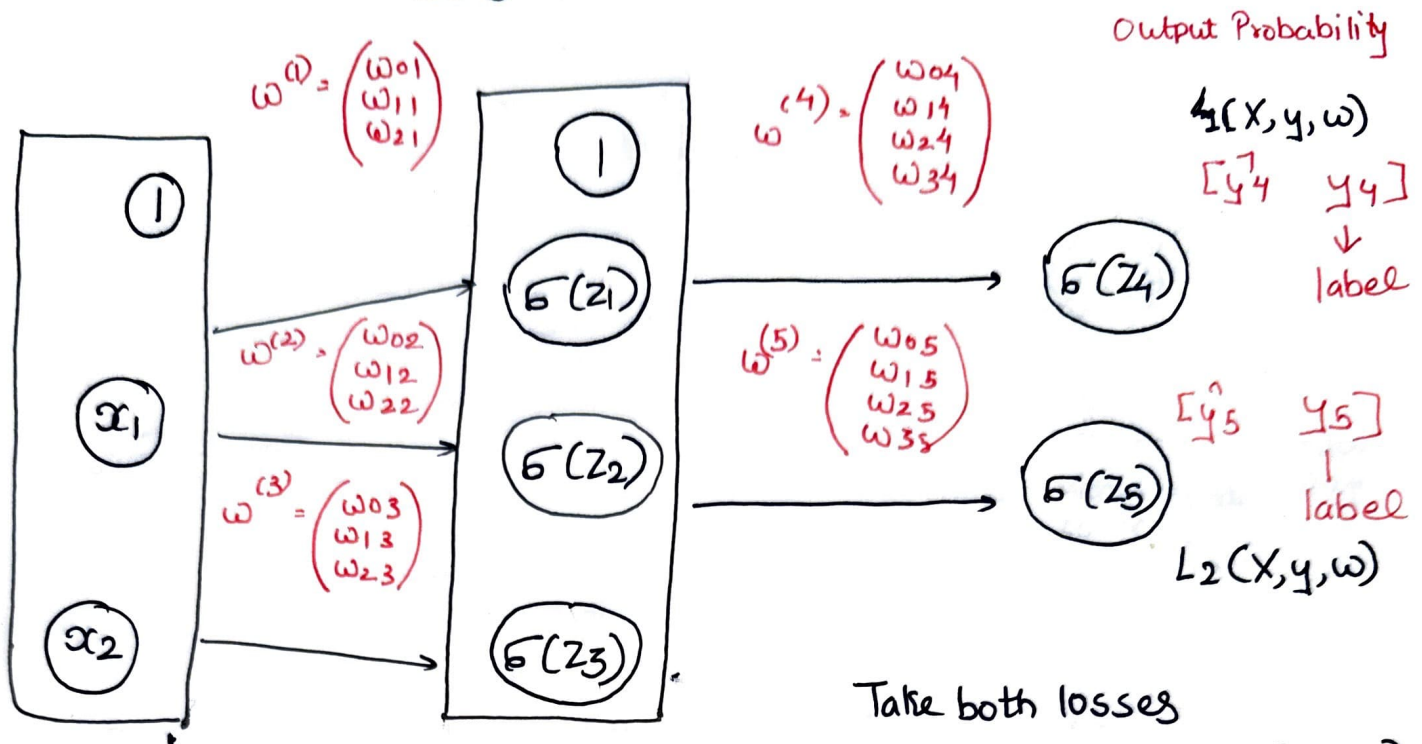
Input



Hidden Layer

→ forward computation

$$\sigma(z_4) = \sigma(w_{04} + w_{14}\,\sigma(\underbrace{z_1}_{\text{neuron 1}}) + w_{24}\,\sigma(\underbrace{z_2}_{\text{neuron 2}}) + w_{34}\,\sigma(\underbrace{z_3}_{\text{neuron 3}}))$$

$w^{(4)} = \begin{pmatrix} w_{04} \\ w_{14} \\ w_{24} \\ w_{34} \end{pmatrix}$   4 weights

$w^{(1)} = \begin{pmatrix} w_{01} \\ w_{11} \\ w_{21} \end{pmatrix}$   $w^{(2)} = \begin{pmatrix} w_{02} \\ w_{12} \\ w_{22} \end{pmatrix}$

Each neuron has its own set of independent weights.

→ apply activation over this → that gives activation value for that layer
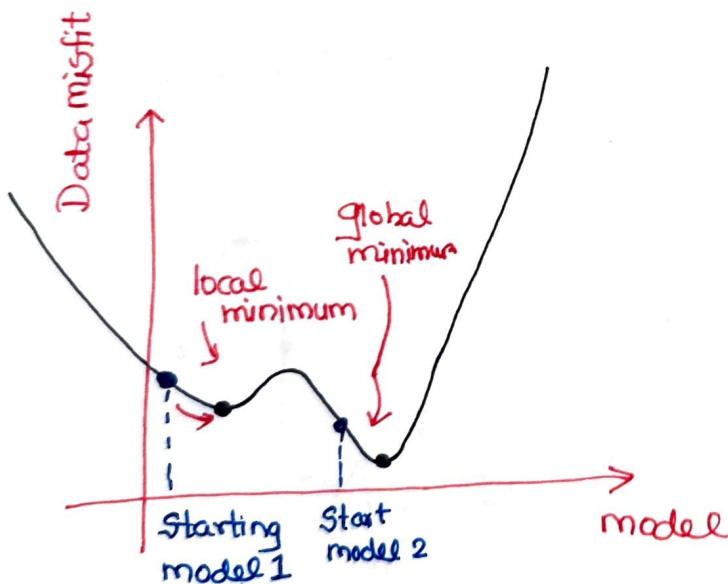
# * Training 3 Layer Networks:

$$\omega^{(1)} = \begin{pmatrix} \omega_{01} \\ \omega_{11} \\ \omega_{21} \end{pmatrix}$$

$$\omega^{(4)} = \begin{pmatrix} \omega_{04} \\ \omega_{14} \\ \omega_{24} \\ \omega_{34} \end{pmatrix}$$

$$\omega^{(2)} = \begin{pmatrix} \omega_{02} \\ \omega_{12} \\ \omega_{22} \end{pmatrix}$$

$$\omega^{(3)} = \begin{pmatrix} \omega_{03} \\ \omega_{13} \\ \omega_{23} \end{pmatrix}$$

$$\omega^{(5)} = \begin{pmatrix} \omega_{05} \\ \omega_{15} \\ \omega_{25} \\ \omega_{35} \end{pmatrix}$$

① ①
$\sigma(z_1)$ $\sigma(z_4)$
$x_1$ $\sigma(z_2)$ $\sigma(z_5)$
$x_2$ $\sigma(z_3)$

$L_1(X,y,\omega)$

$[\hat{y}_4 \quad y_4]$
↓
label

$[\hat{y}_5 \quad y_5]$
|
label

$L_2(X,y,\omega)$

Take both losses

$$L_1(X,y,\omega) + L_2(X,y,\omega) = L(X,y,\omega)$$

Updates weights over iterations via
Some gradient descent
(To Reduce Loss)

# * Complicated Landscape:-
## (1-Dimension)



Data misfit

global minimum

local minimum

Starting model 1   Start model 2

model

Local minimum → generalize

$$L(X,y,\omega)$$
↳ entire set of weights

⎛ (move forward
based on
Step size / learning
rate ⎞ → more to local minimum → not global minimum directly
↓
However if we start at model 2
(global minimum can be found)

# *Universal Approximation Theorem

**Input Layer**     **Hidden Layer**     **Output Layer**

→ output

non linear

**(4-Layer) More Powerful?**

Input Layer   Hidden Layer 1   Hidden Layer 2   Output layer
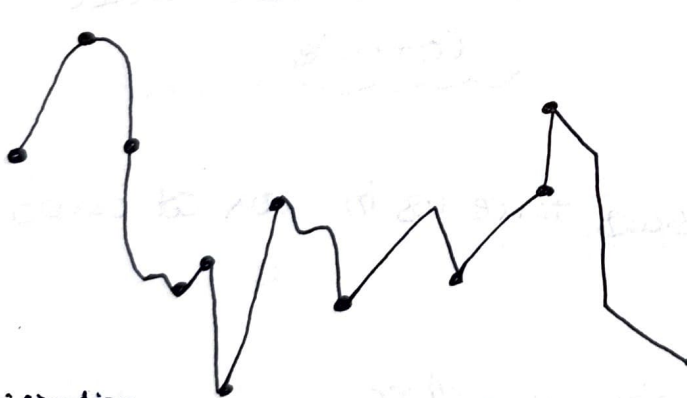
If the training dataset $(X,y)$ comes from a 'reasonable' function $y = f(X)$ (where reasonable means almost any function you can imagine), then there exists a 3 layer neural network $\underline{m}$ such that $m(X)$ approximates $f(X)$ arbitrarily well ie the loss function of $\underline{m}$ can be driven to be very small on X.

①

$(X,y)$

→ very similar
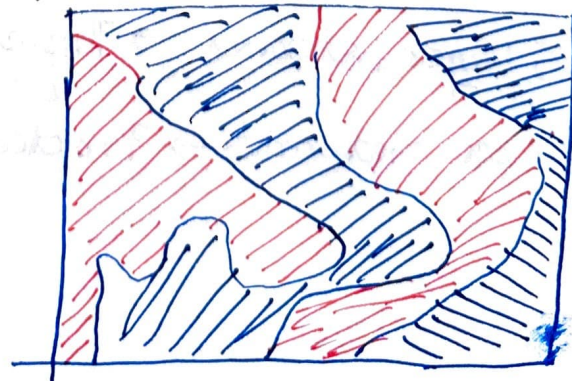
$m : m(X) \approx y$

**Classification Problem**

Classify Points →

Powerful 3 Layer Network

If you go to 4-layer
↓
you cannot do much better
↓
If you learn it on 4 layer

↳ Some weights/setting exists on a 3 layer network that produces the same result.

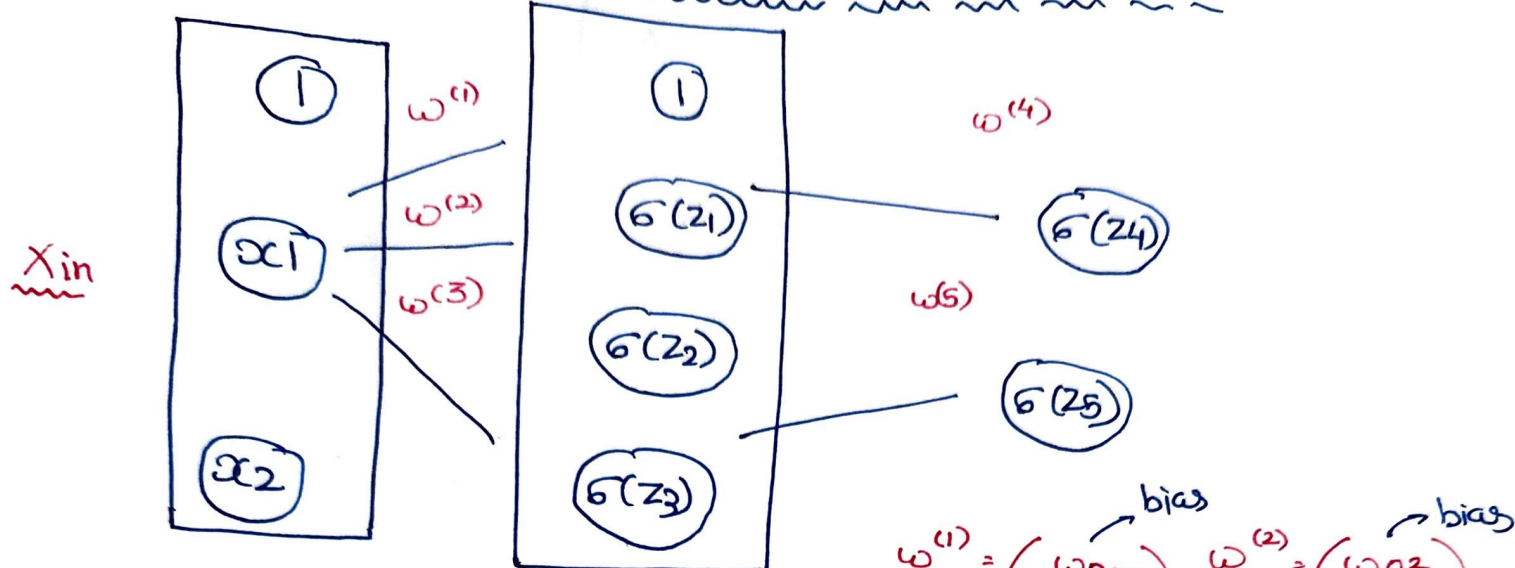So UAT is strongest for 3 Layer Networks.

It does not speak about how many nodes in hidden layer will be → So we don't know how big it will be

it cannot be processed on a local Computer

Combination of these issues force us to look at deeper networks

↓

3 Layer Networks → Theoretically better

Deeper Networks → Practically Better.

# *Algebra of Forward Computation



$$W^{(1)} = \begin{pmatrix} W_{001} \\ W_{11} \\ W_{21} \end{pmatrix} \leftarrow bias$$

$$W^{(2)} = \begin{pmatrix} W_{002} \\ W_{12} \\ W_{22} \end{pmatrix} \leftarrow bias$$

other weights →

$$W_{ih} = \begin{pmatrix} W_{11} & W_{12} & W_{13} \\ W_{21} & W_{22} & W_{23} \end{pmatrix}$$

weights ← $2 \times 3$ → neurons

$$W^{(4)} = \begin{pmatrix} W_{04} \\ W_{14} \\ W_{24} \\ W_{34} \end{pmatrix}$$

$$Z_h = x_{in} \cdot W_{ih} + [W_{01}, W_{02}, W_{03}]$$

↓ input Point

biases

---

## The forward function (3 Layer)

nets  $Z_h = x_{in} \cdot W_{ih} + [W_{01}, W_{02}, W_{03}] \rightarrow ih \rightarrow$ input to hidden layer

$a_h = \sigma(Z_h)$  activations (hidden neuron) → become input for ne

nets  $Z_0 = a_h \cdot W_{ho} + [W_{04}, W_{05}] \rightarrow ho \rightarrow$ hidden to output layer

$a_0 = \sigma(Z_0)$  activations (output activation) ↓

Final Output

*Take batch of points → $X_1, X_2, X_3, \ldots m(X_1), m(X_2), m(Y_3)$
Send one by one to model ↗ / works on all parallel by making a matrix

$$X_{in} = \begin{pmatrix} X_1 \\ X_2 \\ X_3 \end{pmatrix}$$

$$3 \times 2$$

Parallism is important instead of single points (for GPU usage) efficient use of computation power.

$$Z_h = X_{in} \cdot W_{ih} + [W_{01}, W_{02}, W_{03}]$$

# * Training with BackPropagation

$y(u(x))$

**chain rule**

$$\frac{dy}{dx} = \frac{dy}{du} \cdot \frac{du}{dx}$$

**Example**

TESLA $\overset{u}{\phantom{a}}$  Bmw $\overset{x}{\phantom{a}}$

Speed of Tesla w.r.t Bmw $2x$

$$\frac{dy}{du} = 2 \quad \left.\right] 200\,m/h$$

Bmw
(on ground) $\dfrac{du}{dx} = 100$

$\left.\right]$ chain rule

Tesla $\rightarrow$ find that
(ground)

loss

$Z = \omega_0 + \omega_1 x_1 + \omega_2 x_2$



$$\frac{1}{2}\left(y - \sigma(z)\right)^2$$

derivative of
loss fn
w.r.t $\omega_1$

$$\frac{dL}{d\omega_1}$$

---



$$\frac{du_2}{du_1} \quad \frac{du_3}{du_2} \quad \frac{du_1}{dw_{ij}}$$

③  ②  ①

Lets focus on $\dfrac{dL}{\partial \sigma(z)}$

**Loss w.r.t activation**

① $\dfrac{dL}{\partial \sigma(z)} = \dfrac{(-y - \sigma(z))}{\phantom{a}}$ ⟶ single variable

Back Propagation → go to ②$-$⑥②

② $\dfrac{d\,\sigma(z)}{dz} = \sigma(z)(1 - \sigma(z))$

then more to 1st & 2nd Layer

③ $\dfrac{dz}{d\omega_1} = \dfrac{\partial(\omega_0 + \omega_1 x_1 + \omega_2 x_2)}{d\omega_1} = X_1$

Algorithm (Automatic Derivations)

Similar to Tesla, Bmw, ground

(2) layer  (1) layer

Now apply chain rule $= \dfrac{dL}{d\omega_1} = \dfrac{dL}{d\sigma}\overset{(3)}{\phantom{a}} \cdot \dfrac{d\sigma}{dz} \cdot \dfrac{dz}{d\omega_1}$