→ AI { Application which can do its own task without any human intervention }

→ DL

ML

It is a subset an AI, which can make pattern on historical data often for making predictions.

→ Netflix App
→ Self Driving Cars
→ ANPR
→ Alexa, Sofia
→ Chat GPT

DL → ("Neural Networks" (Geovffrey Hinton)

CNN
RNN   ML

Mimic the human brain) {2006}

Structured

. csv
. excel
. html

Unstructured Data

Videos
Text
Audio
Images

Q.) Why Deep Learning is becoming so popular?

⇒ 2005 → ORKUT, facebook, Insta, Twitter
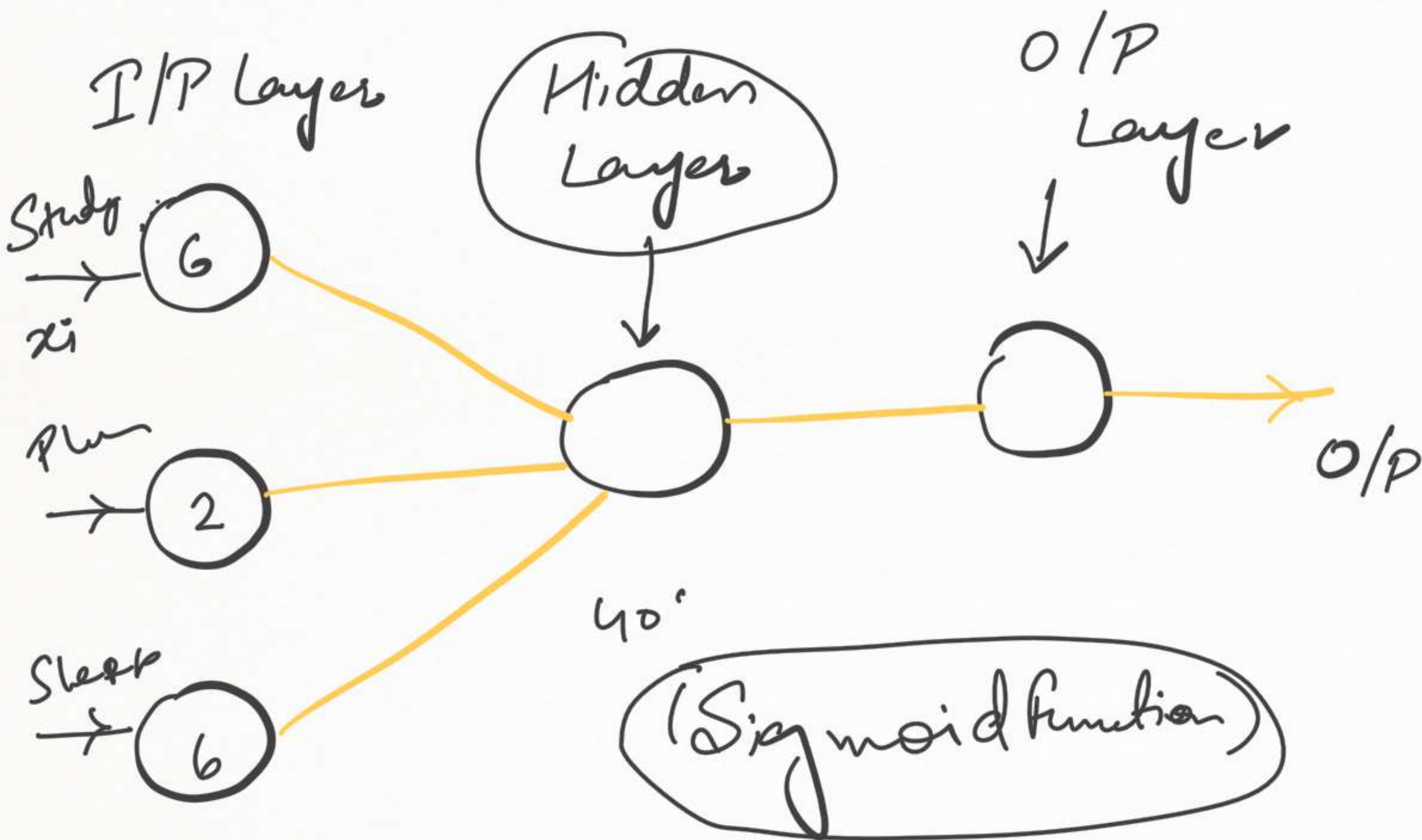
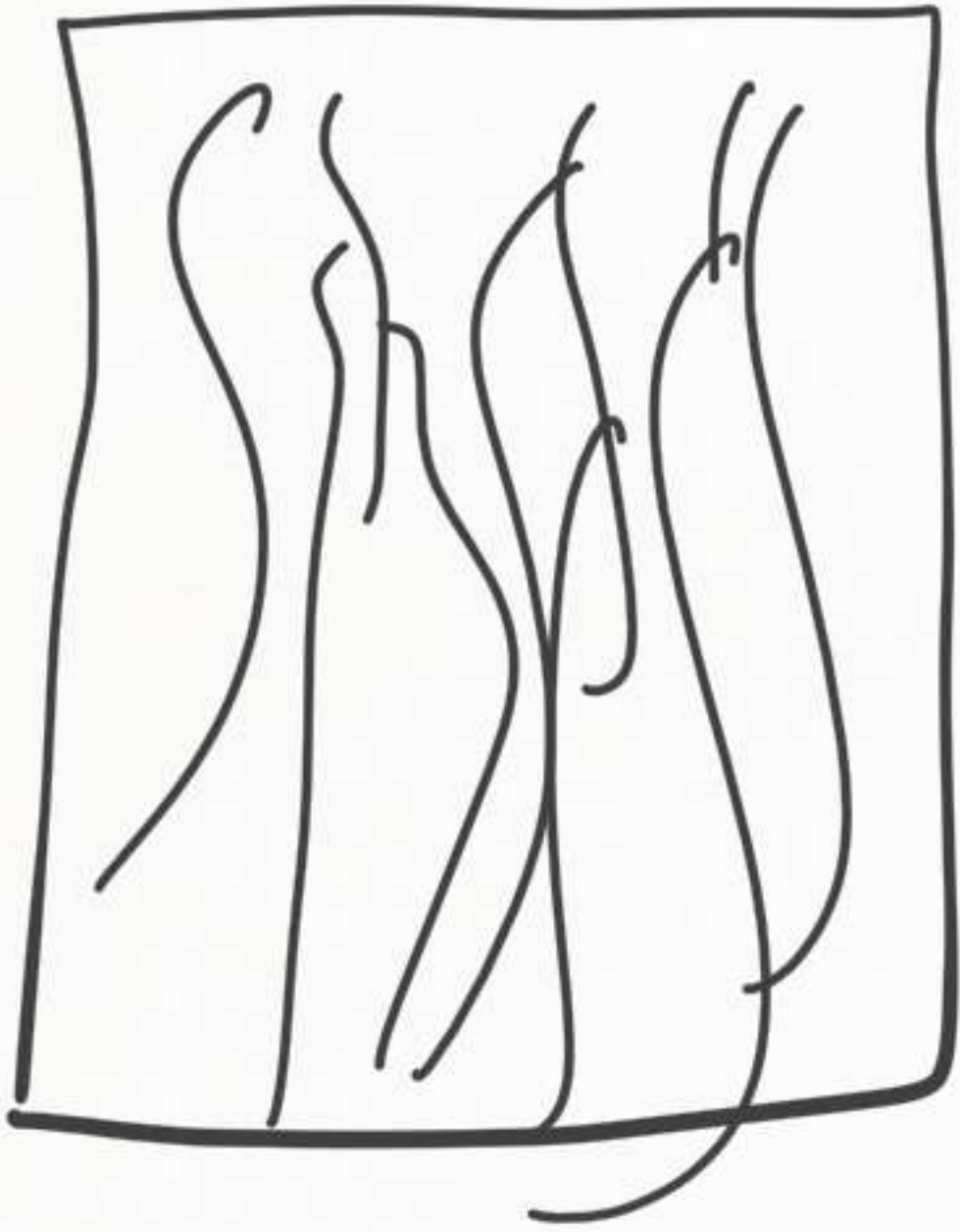Data ↑↑↑↑ exponentially

(2016) "PetaBytes"

So

GPU's
TPU
$\}$ Hardware Ac

# Perceptron {Single Layer Neural N/w}

**I/P Layer**

Hidden Layers

O/P Layer

Example



| Feature | | | Target |
|---------|------|-------|-----|
| Study | Play | Sleep | P/F |
| 6 | 2 | 6 | 1 |
| 2 | 5 | 8 | 0 |
| 5 | 3 | 7 | 1 |

Study $x_i$

6

Plw

2

Sleep

6

40°

O/P

(Sigmoid function)

1- year

Bias

I/P Layer

Hidden Layer

Forward Propagation

$x_1 = 6$

$w_1$

$w_2$

$w_3$

O/P Layer

$y = 1$

$0.3$

① ②

$\sum x_i w_i + b$  $Act(y)$

$O_1$  $w_4$

$\sum$ $Act()$

$\hat{y} = 0$

$x_2 = 2$

$x_3 = 6$

$Loss = (\dot{y} - \hat{y})$

$0.2$  $=$

BACK PROPAGATION

$\sum x_i w_i = x_1 w_1 + x_2 w_2 + x_3 w_3 + b \Rightarrow y = 1_0$

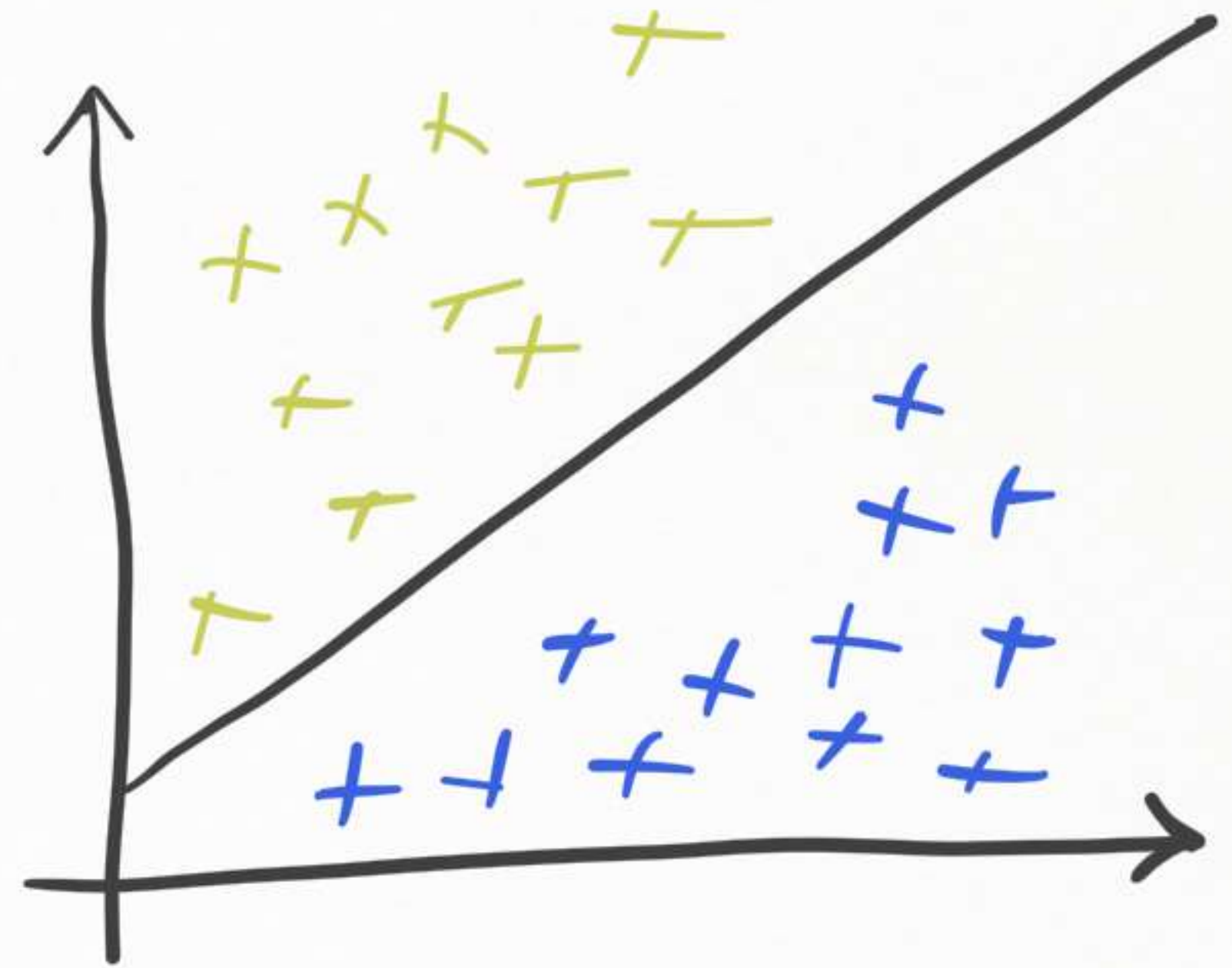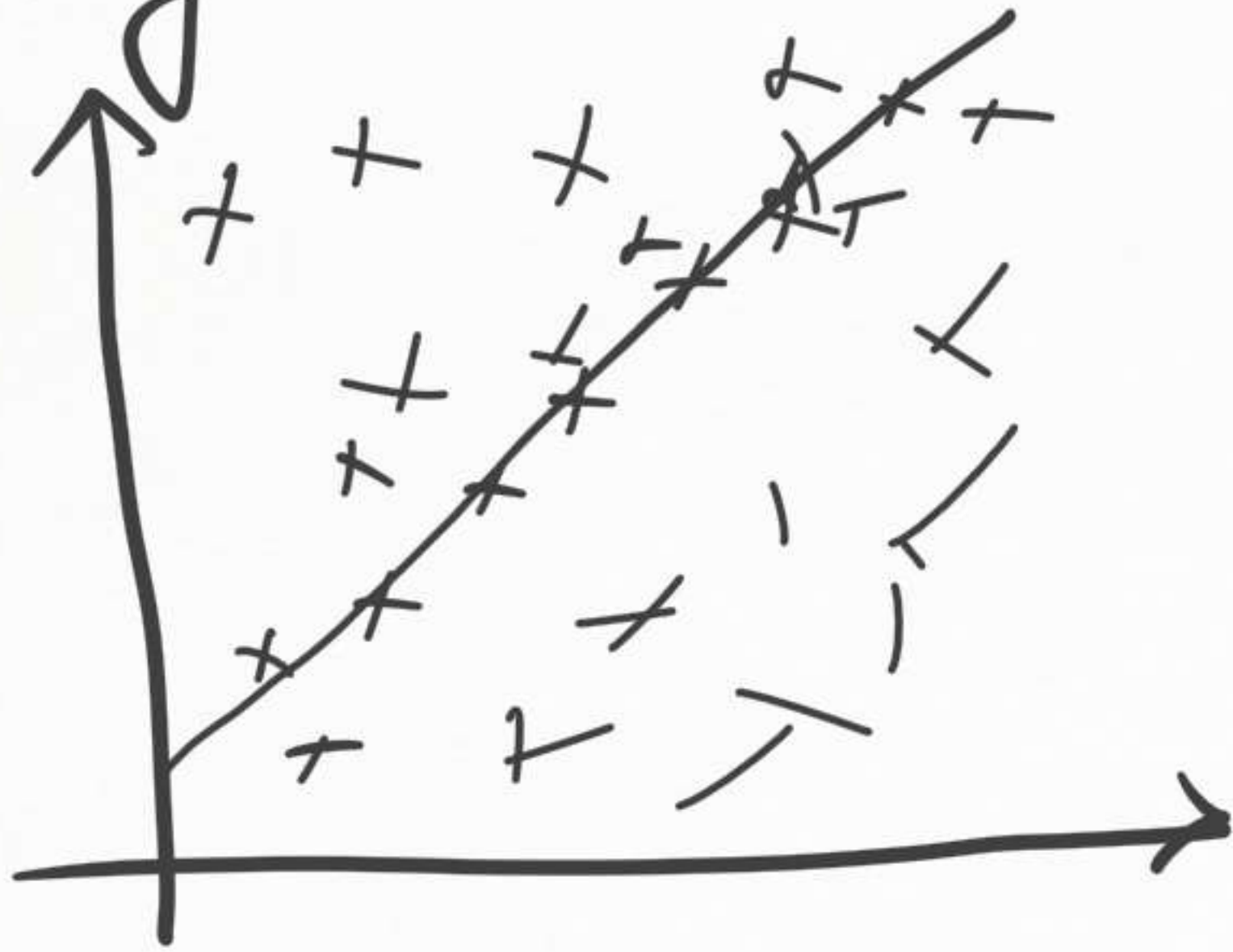$y = mx + c \quad , \quad \boxed{\hat{y} \sim \beta_0 + \beta_1 x}$  $= 0.7$

# Activation functions

→ Support non-linear properties.

→ Bring all the data into same scale.

# Sigmoid A f

$$z = \frac{1}{1 + e^{-y}} = \frac{1}{1 + e^{-(\Sigma x_i \omega_i + b)}}$$

$$\downarrow$$

$$= 0 \text{ to } 1$$

$$\geqslant 0.5 \Rightarrow 1$$

$$< 0.5 \Rightarrow 0$$

# Forward Propagation

1.) I/P Layer

2.) Weights

3.) Activation functions
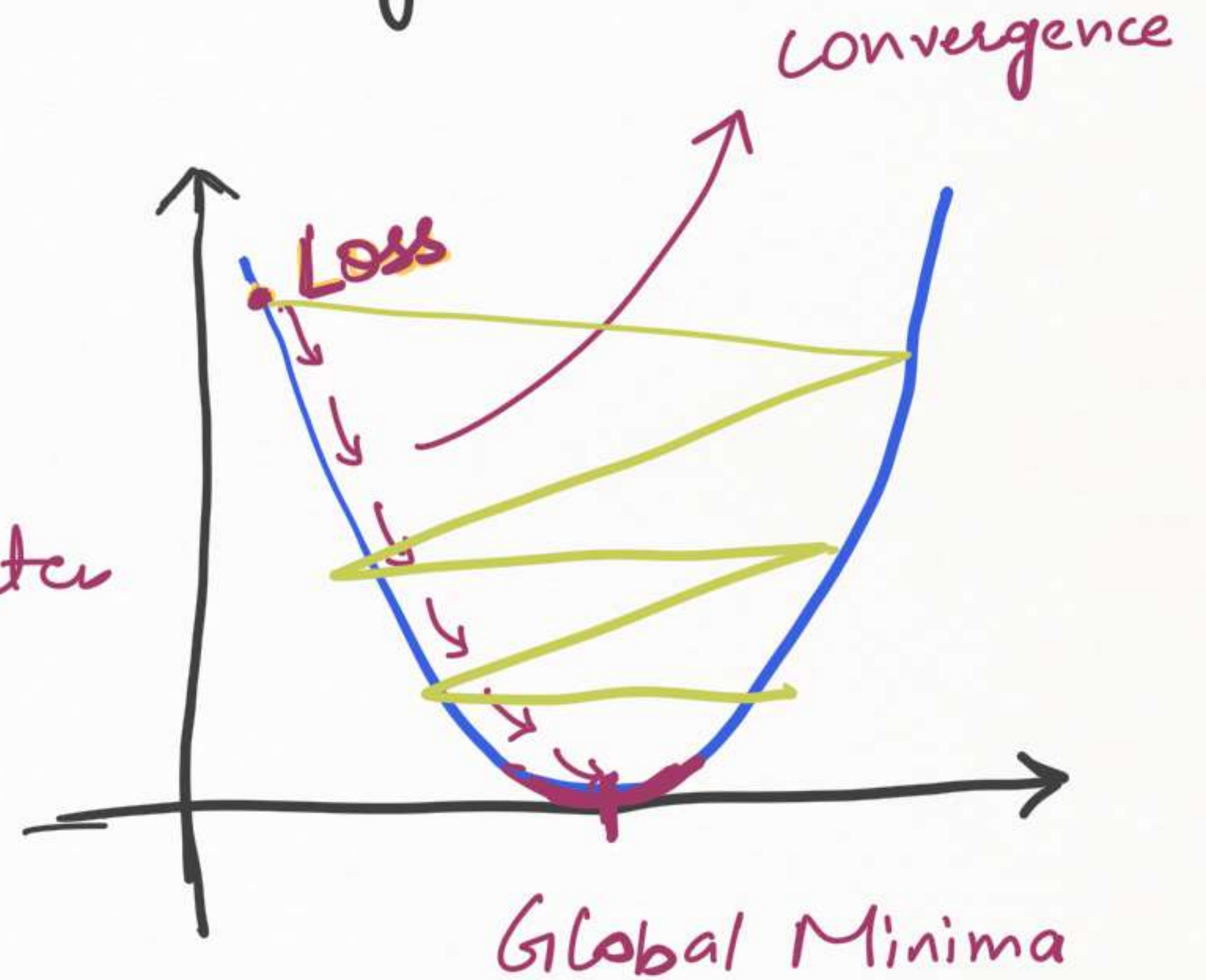
4.) O/P Layer

5.) Loss Calculation
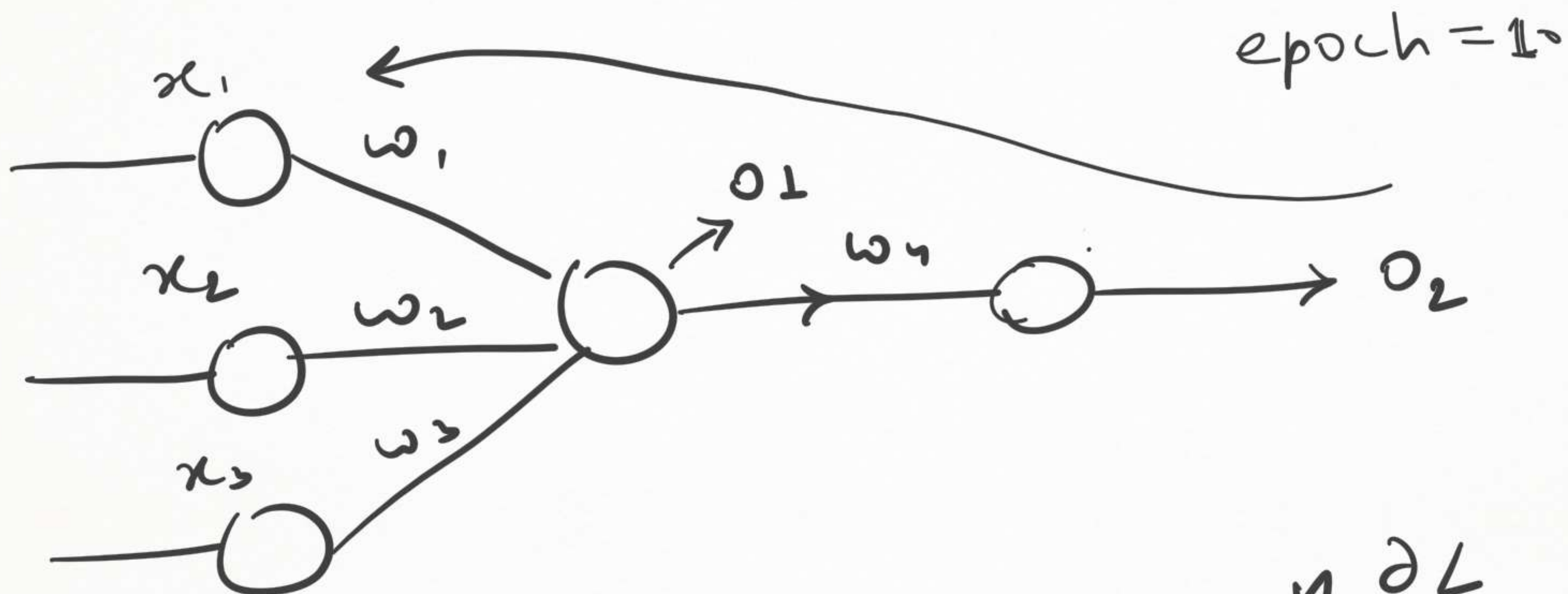
# Back Propagation

① Update the weights

② Optimizers

## Weight Updation formula

$$W_{new} = W_{old} - \eta \boxed{\dfrac{\partial L}{\partial W_{old}}}$$

Learning Rate $= 10^{-1}$ to $10^{-3}$

convergence

Loss

Global Minima

* Learning Rate is a hyperparameter which controls the speed of convergence.

$$\text{epoch} = 10$$
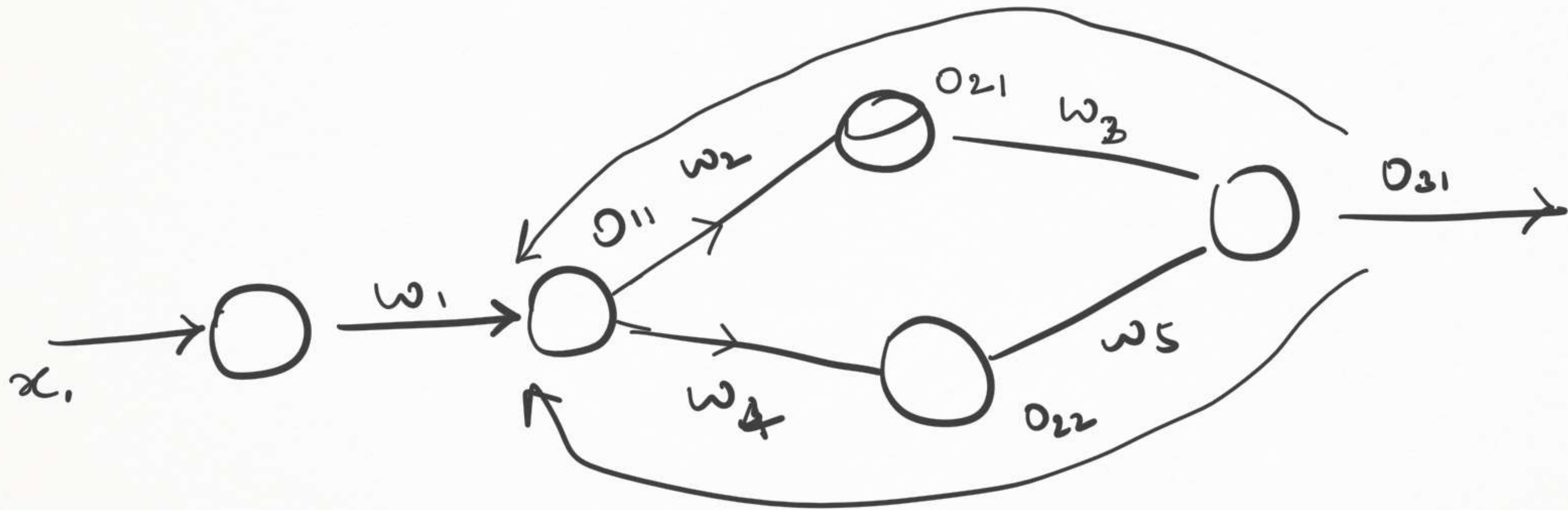
$$w_{4\,new} = w_{4\,old} - \eta \frac{\partial L}{\partial w_{4\,old}}$$

$$\frac{\partial L}{\partial w_{4\,old}} = \frac{\partial L}{\partial O_2} \times \frac{\partial O_2}{\partial w_{4\,old}} \quad \Big\} \quad \text{chain Rule} \atop \text{differentiation}$$
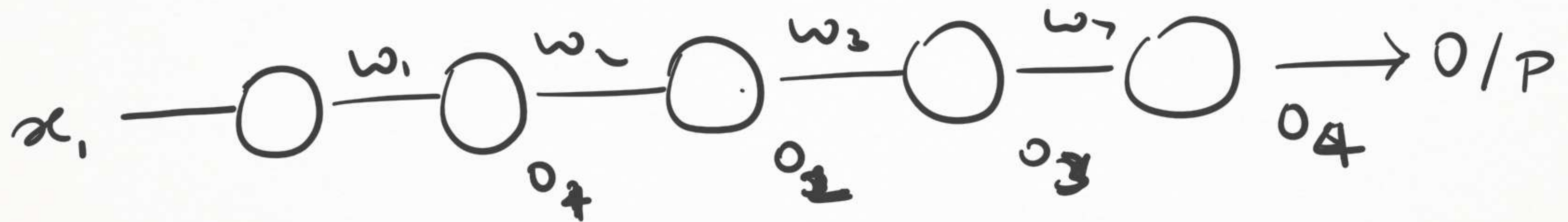
$$\boxed{\frac{\partial L}{\partial w_{1new}} = \frac{\partial L}{\partial O_2} \times \frac{\partial O_2}{\partial O_1} \times \frac{\partial O_1}{\partial w_{1old}}}$$

$$w_{1new} = w_{1old} - \eta \frac{\partial L}{\partial w_{1old}}$$

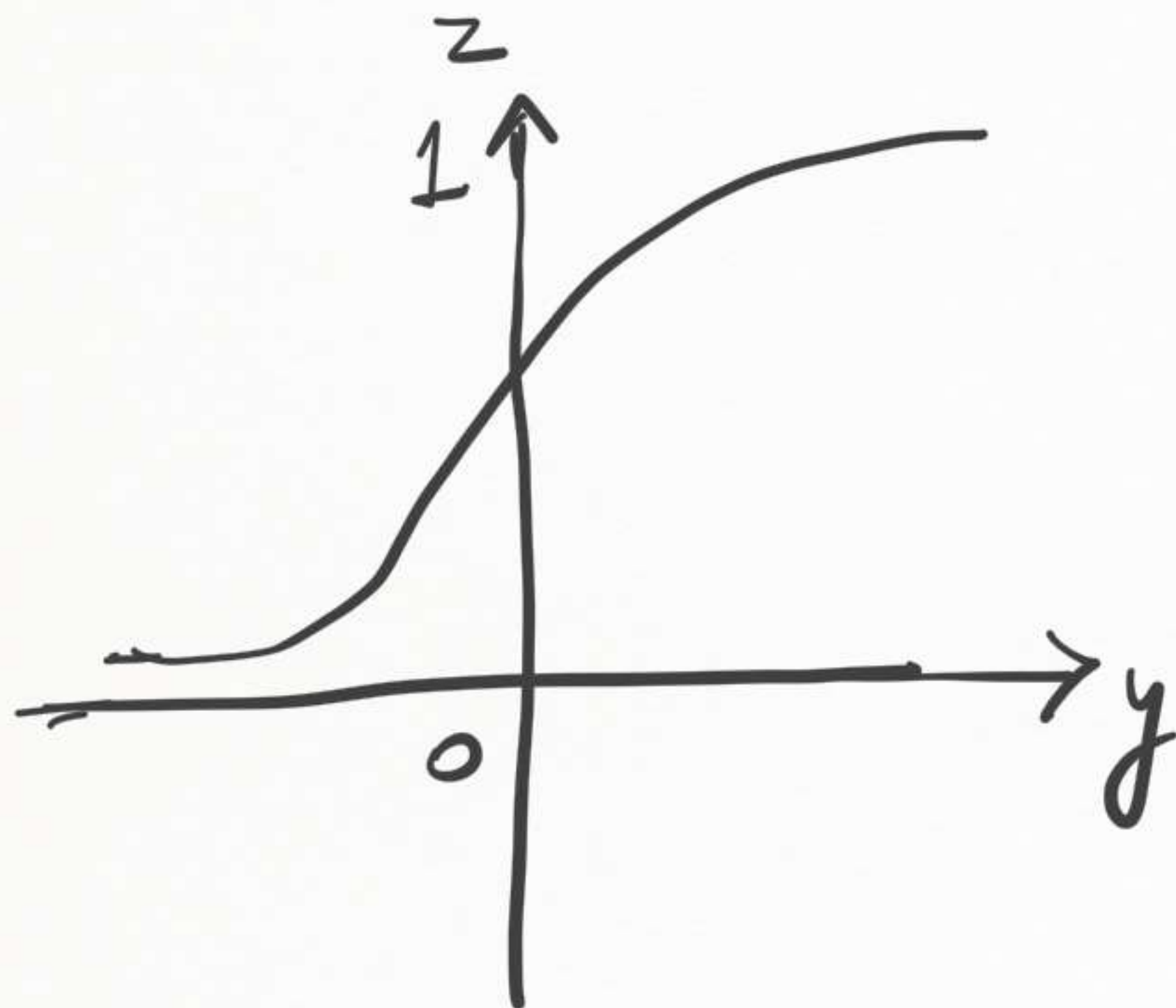$$W_{1new} = W_{1old} - \eta \frac{\partial L}{\partial W_{1old}}$$

$$\frac{\partial L}{\partial W_{1old}} = \frac{\partial L}{\partial O_{31}} \times \frac{\partial O_{31}}{\partial O_{21}} \times \frac{\partial O_{21}}{\partial O_{11}} \times \frac{\partial O_{11}}{\partial W_{1old}}$$

$$x_1 \longrightarrow \bigcirc \xrightarrow{w_1} \bigcirc \xrightarrow{w_2} \bigcirc \xrightarrow{w_3} \bigcirc \xrightarrow{w_4} \bigcirc \longrightarrow O/P$$
$$O_1 \qquad O_2 \qquad O_3 \qquad O_4$$

$$\frac{\partial L}{\partial w_1} = \frac{\partial L}{\partial O_4} \times \frac{\partial O_4}{\partial O_3} \times \frac{\partial O_3}{\partial O_2} \times \frac{\partial O_2}{\partial O_1} \times \frac{\partial O_1}{\partial w_1}$$

$$= 0.25 \times 0.25 \times 0.25 \times 0.25 \times 0.25$$

$$= 0.0000625 \approx 0$$

$( \text{small value} ) \times ( \text{very small value} )$

$0.01 \times 0.0000625$

$=$

Very very very small

$z = 0 \text{ to } 1$

$z = \dfrac{1}{1+e^{-x}}$

Derivative $(z) = 0 \text{ to } 0.25$

$W_{new} = W_{old} - \left( \eta \dfrac{\partial L}{\partial w} \right)$

$$W_{new} \approx W_{old}$$

$\rightarrow$ { "Vanishing Gradient Problem" }

0.25

0.25

0

"Sigmoid Af"

$$\Sigma \rightarrow (AF) \xrightarrow{FP} O/P \longrightarrow (Loss)$$

Update Weight $\longrightarrow$ Optimizers

BACK PROPAHATION

$$\frac{1}{1+e^{-x}} = \frac{1}{1+e^{-\xi}}$$

$$3$$

$$\sum w_i x_i \cancel{\leqq} 17$$

$$\sigma(\cdot) = \frac{1}{1+e^{-17}} = 0.7$$

$$\sum w_i x_i = -17$$

$$\max(0, -17) = 0$$
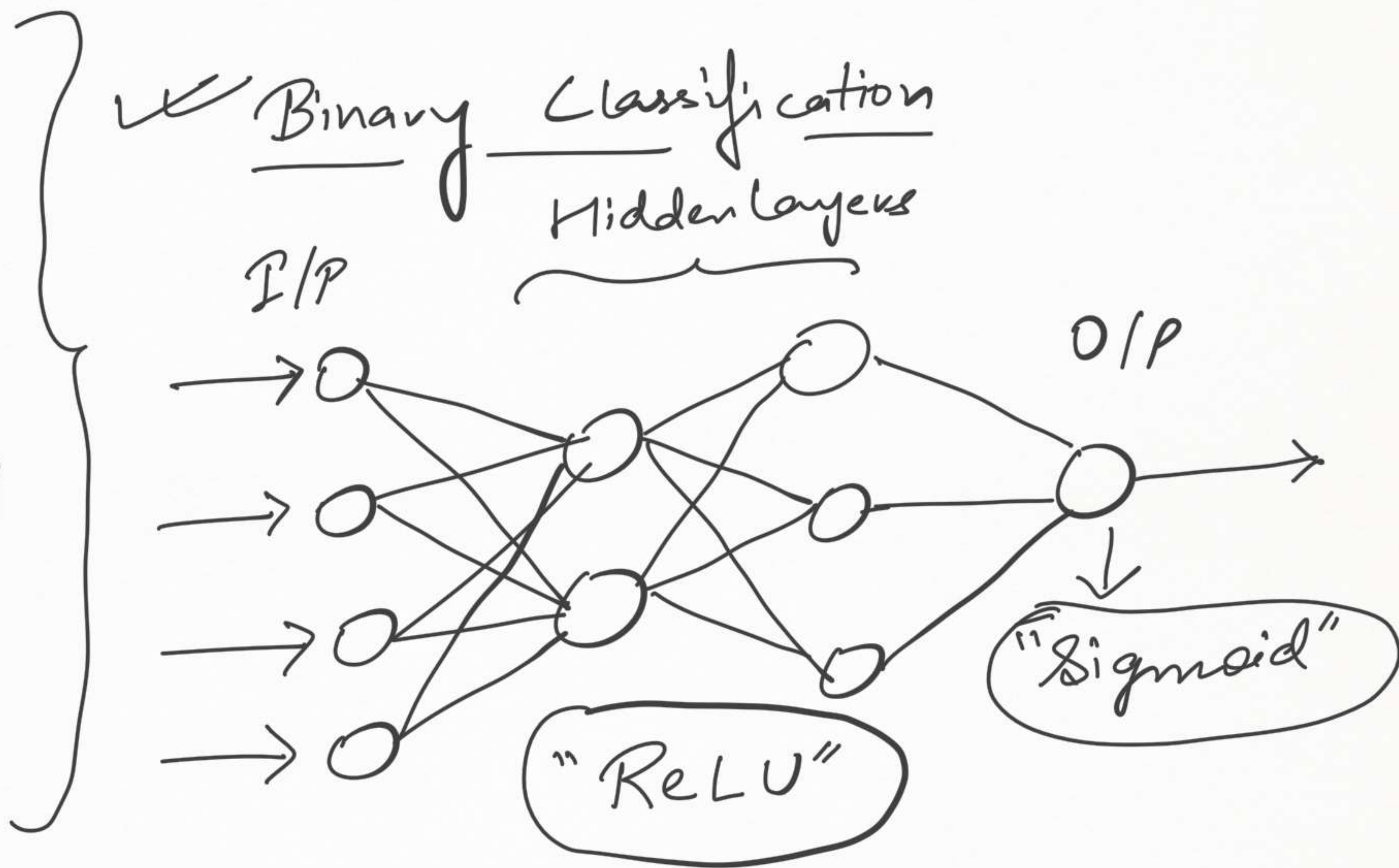
$$ReLU(x) = \max(0, 17)$$
$$= 17$$

## Softmax AF

$$S(i) = \frac{e^i}{\sum e^i}$$

$$1, 2, 3, 4, 5$$
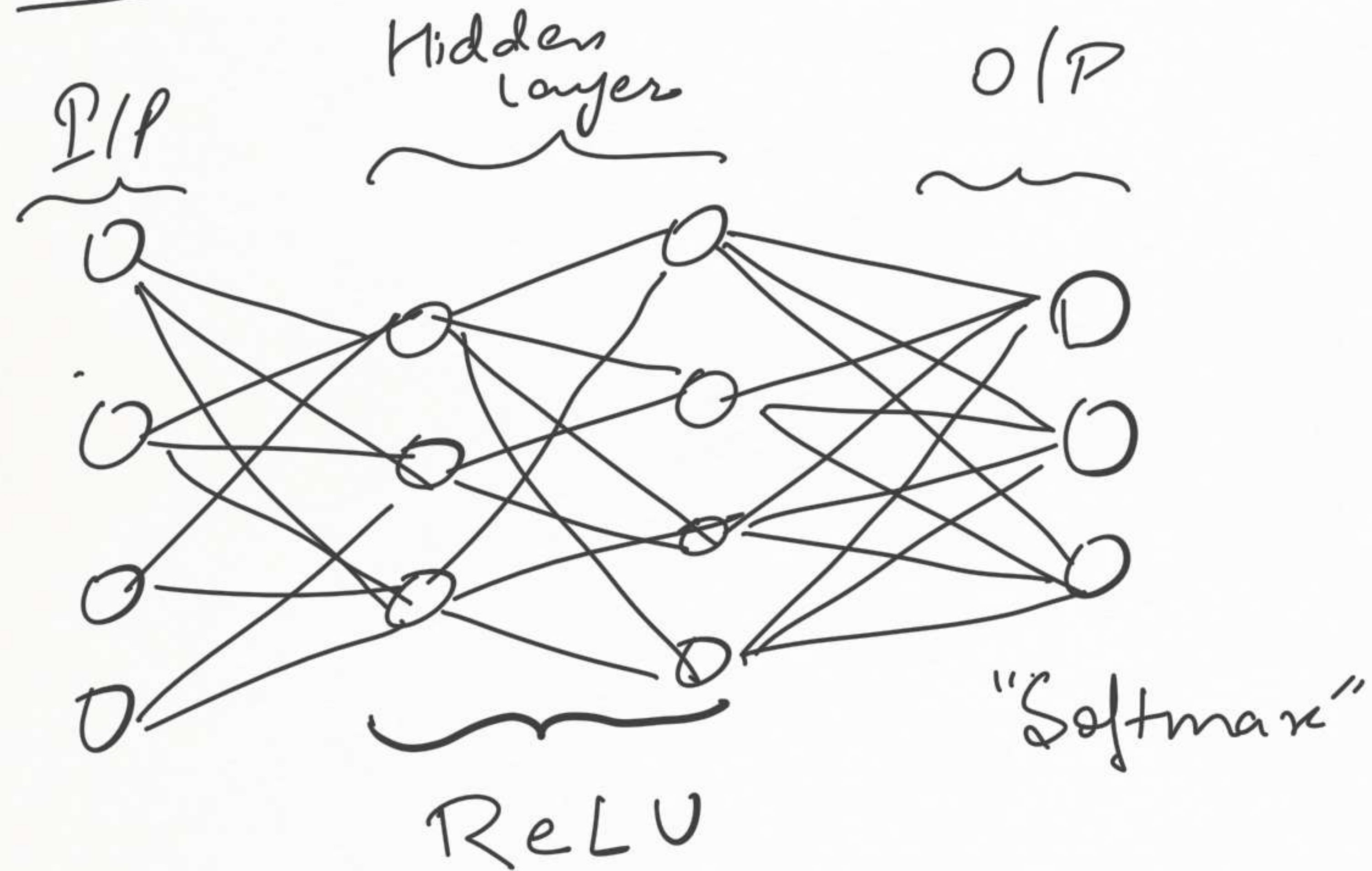
$$S(i=2) = \frac{e^2}{e^1 + e^2 + \dots e^5}$$

$$S(i=1) = \frac{e^1}{e^1 + e^2 + e^3 + e^4 + e^5} =$$

1.) Sigmoid

2.) Tanh

3.) ReLU

4.) Leaky ReLU

5.) ELU

6.) PReLU

7.) Softmax

Binary Classification

Hidden layers

I/P

O/P

"ReLU"

"Sigmoid"

# Multi Class Classification

Hidden Layers

I/P

O/P

"Softmax"

ReLU

# Regression

"ReLU"

ReLU
or
Linear
AF

Loss or cost

$$\downarrow$$

Error $\Bigg\{$

MSE

$$Loss = \frac{1}{2}(y-\hat{y})^2$$

$$Cost = \frac{1}{2n}\sum_{i=1}^{n}(y-\hat{y})^2$$

# Regression Task

1.) Mean Squared Error (MSE)

$$\frac{1}{2n} \sum_{i=1}^{n} (y-\hat{y})^2$$

Adv.

1.) Diffentiable

2.) Convergence is fast

3.) It has only 1 local or global minima.

## Disadvantage

1.) Not robust to an outlier.

## Mean Absolute Error

$$\frac{1}{2n} \sum_{i=1}^{n} |y - \hat{y}|$$

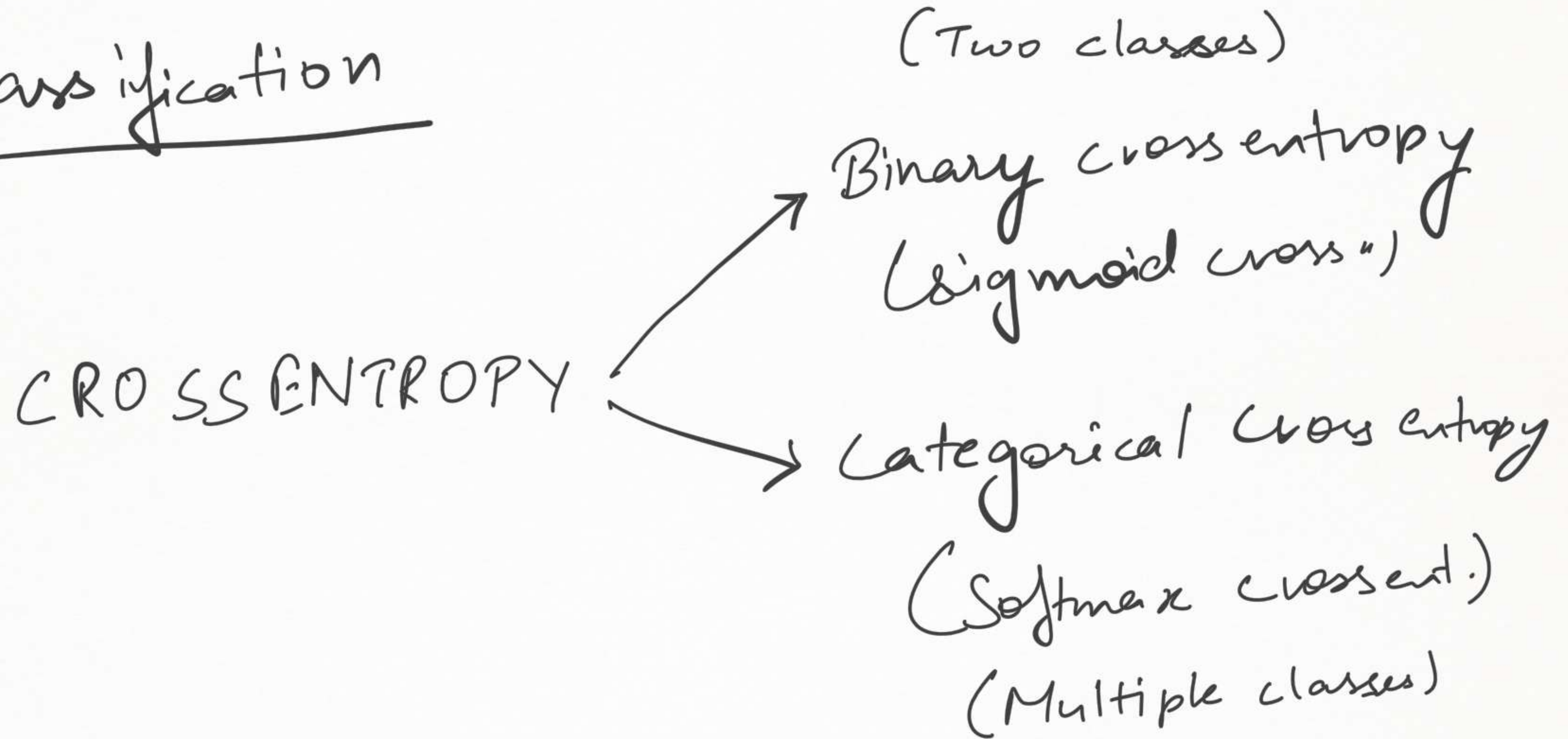① Robust to an outlier.

③ Huber Loss

$$
\text{Loss} = \begin{cases} \frac{1}{2}(y-\hat{y})^2, & \text{if } |y-\hat{y}| \leq \delta \\ \delta|y-\hat{y}| - \frac{1}{2}\delta^2, & \text{otherwise} \end{cases}
$$

$\delta$ → $10$ → hyper parameter

Classification

CROSS ENTROPY

(Two classes)
Binary cross entropy
(sigmoid cross ")

Categorical cross entropy

(Softmax crossent.)
(Multiple classes)

## Binary Cross Entropy —

$$Loss = -y * \log(\hat{y}) - (1-y) * \log(1-\hat{y})$$

$$Loss = \begin{cases} -\log(1-\hat{y}) & \text{if } y = 0 \\ -\log(\hat{y}) & \text{if } y = 1 \end{cases}$$

$$\hat{y} = \frac{1}{1+e^{-z}}$$

$$\text{cost}(J) = -\left[\sum_{i=1}^{n} y * \log(\hat{y}) + (1-y) * \log(1-\hat{y})\right]$$

$$= 7.3 \approx 0$$

$$\hat{y} = \frac{1}{1+e^{-z}}$$

# Categorical Cross Entropy (Multi-class Classification)

| $f_1$ | $f_2$ | $f_3$ | O/P | fever | Malaria | Jaundice |
|-------|-------|-------|-----|-------|---------|----------|
| 2 | 6 | 4 | fever | 1 | 0 | 0 |
| 5 | 1 | 8 | Malaria | 0 | 1 | 0 |
| 7 | 3 | 9 | Jaundice | 0 | 0 | 1 |

$$L(x_i, y_i) = -\sum_{i=1}^{c} y_{ij} * \ln(\hat{y}_{ij})$$

$$y_{ij} = \begin{cases} 1, & \text{if the element in the class} \\ 0, & \text{otherwise} \end{cases}$$

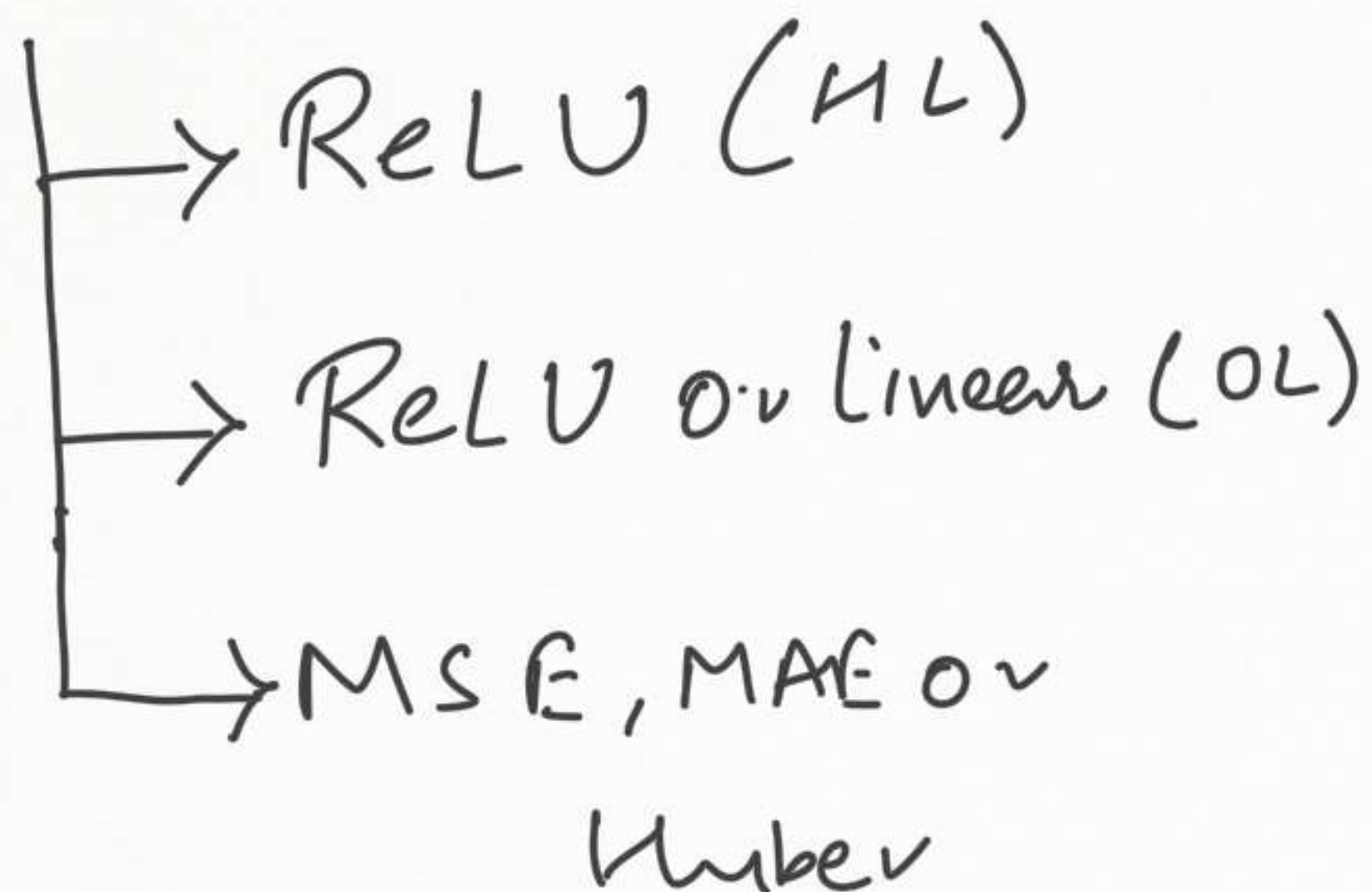$$\sigma(z) = \frac{e^{z_i}}{\sum_{i=1}^{c} e^{z_i}}$$
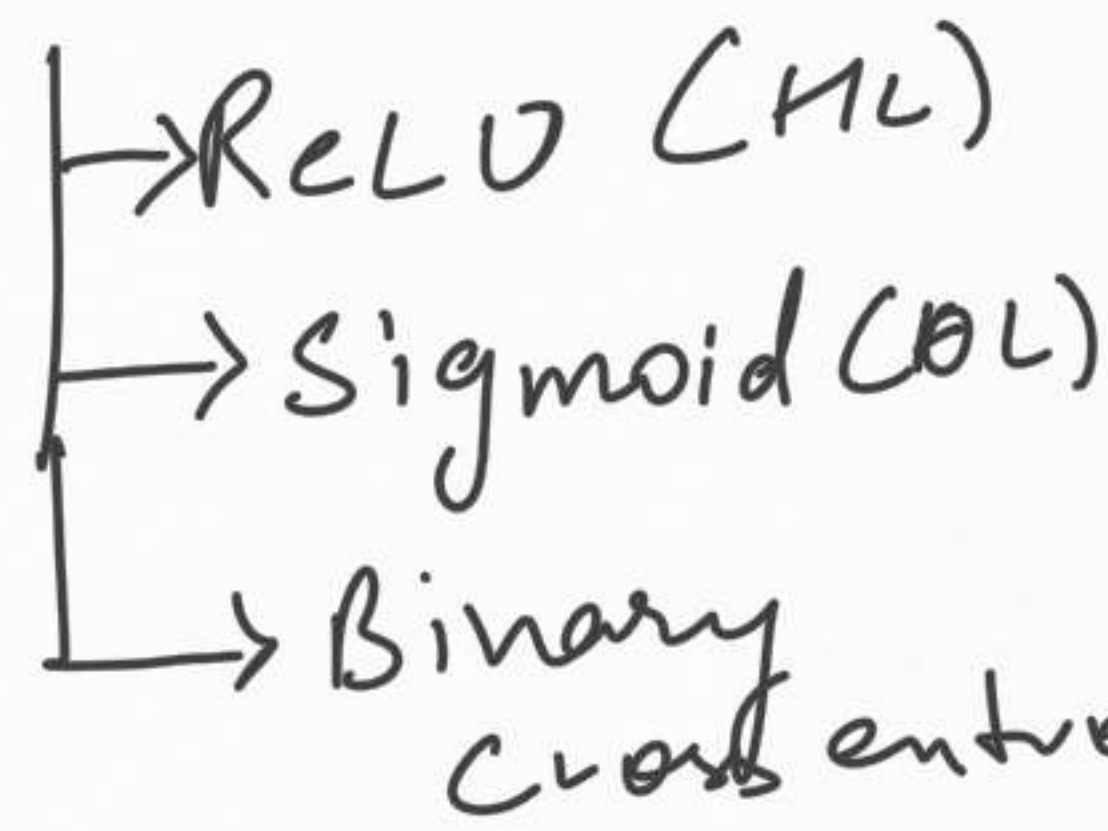$(\hat{y})$

$$\sigma(f) = 0.2$$

$$\sigma(m) = 0.1$$

$$\sigma(J) = 0.7$$

# Conclusion

## Regression

- → ReLU (HL)
- → ReLU or linear (OL)
- → MSE, MAE or Huber

## Classification

- → Binary
  - → ReLU (HL)
  - → Sigmoid (OL)
  - → Binary Crossentropy
- → Multiclass
  - → ReLU (HL)
  - → Softmax (OL)
  - → Categorical Crossentropy

# Optimizers

1. Gradient Descent (Batch)

2.) Stochastic GD (SGD)

3.) Mini-Batch GD (MBGD)
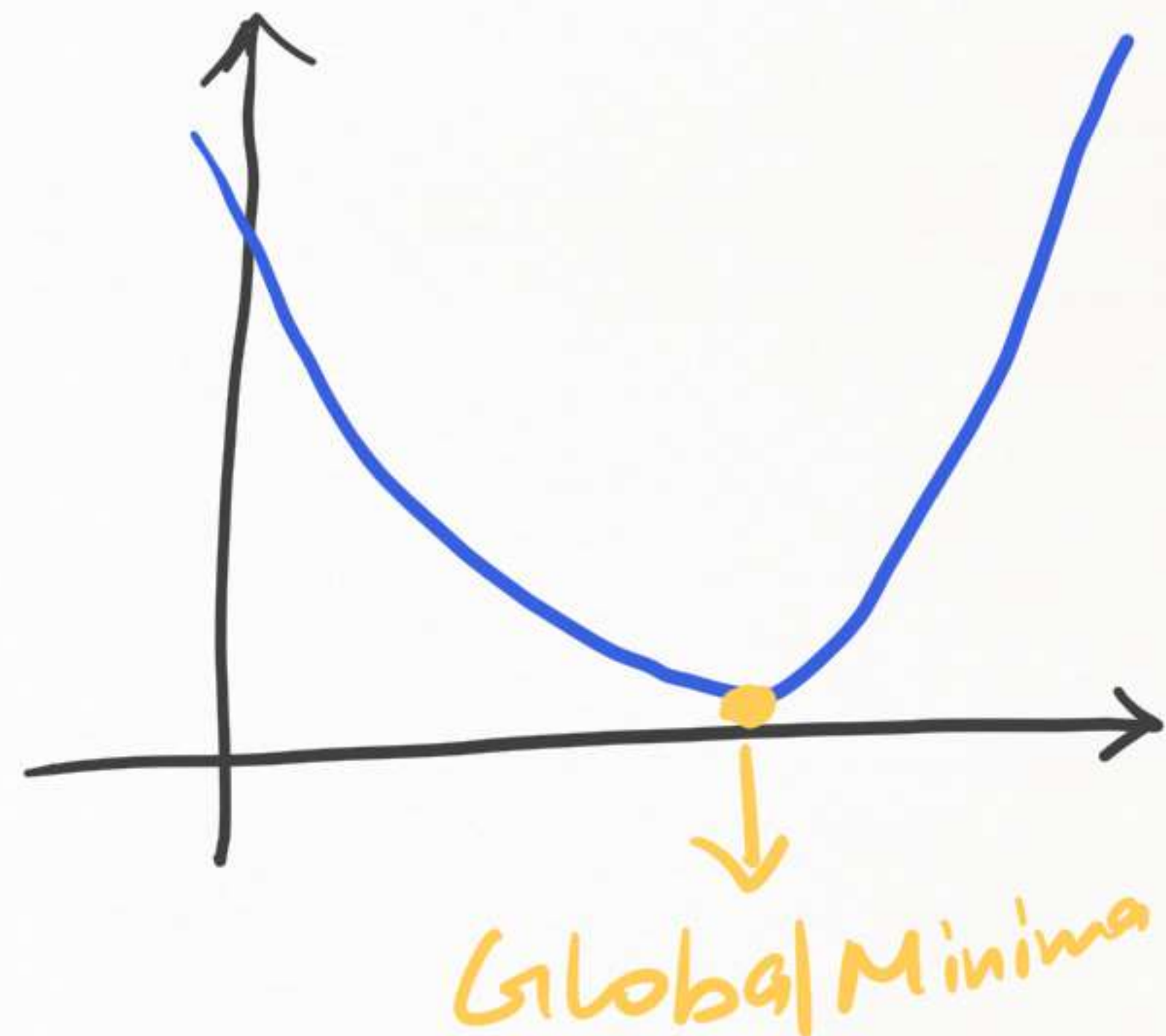
4.) SGD with momentum

5.) Adagrad

6.) RMSProp

7.) Adam

# ① Gradient Descent (Batch)

## Weight Updation formula

$$W_{new} = W_{old} - \eta \frac{\partial L}{\partial W_{old}}$$

Learning Rate

$$\eta = 0.1 \text{ to } 0.001$$

Global Minima

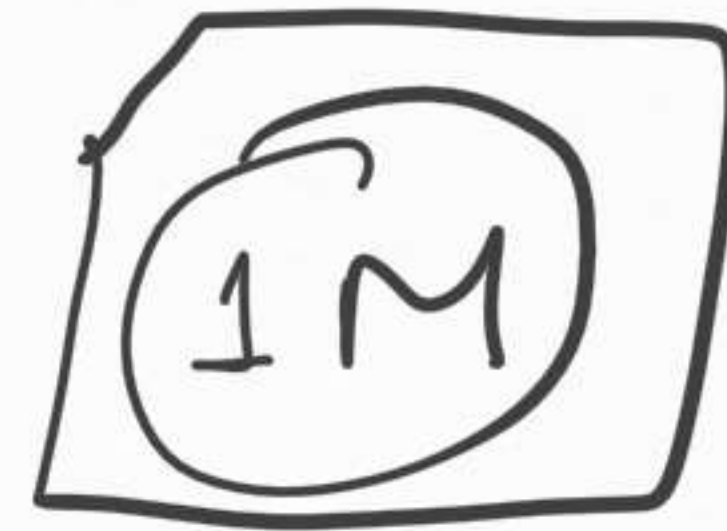# Batch vs. Epoch vs. Iterations

$n = 1000$

Forward Propagation $\longrightarrow$

$\left.\begin{array}{c} \\ \end{array}\right\}$ 1 Epoch

$\longleftarrow$ Back Propagation

1M

Dist.

(1) Resource Extensive {Huge RAM}

1000 rows

$\boxed{100}$

$\boxed{100}$

$\left.\begin{array}{c} \longrightarrow \\ \longleftarrow \end{array}\right\}$ Iter 1 (10)

$\left.\begin{array}{c} \longrightarrow \\ \longleftarrow \end{array}\right\}$ Iter 1 (10)

$\left.\begin{array}{c} \longrightarrow \\ \longleftarrow \end{array}\right\}$ 2 (10)

$\dfrac{1000}{100} = 10$

$\vdots$

$\left.\begin{array}{c} \longrightarrow \\ \longleftarrow \end{array}\right\}$ 10 (10)

② (Stochastic) Gradient Descent (SGD)  (1 M)

Epochs

RAM ↓↓

1 record

$\}$ Iterations

Update weight

Dist.

⟶ Time consuming.

⟶ convergence will be

very slow.

1 record

$\}$ It ~ 2

③ (Mini Batch) SGD

→ Resource Intensive

→ Convergence
  will be better

→ Time complexity
  will improve.

1 M rows

$(Batch\_size) = (1000)$

$$No. \ of \ ituation = \frac{1000\cancel{000}}{1\cancel{000}}$$
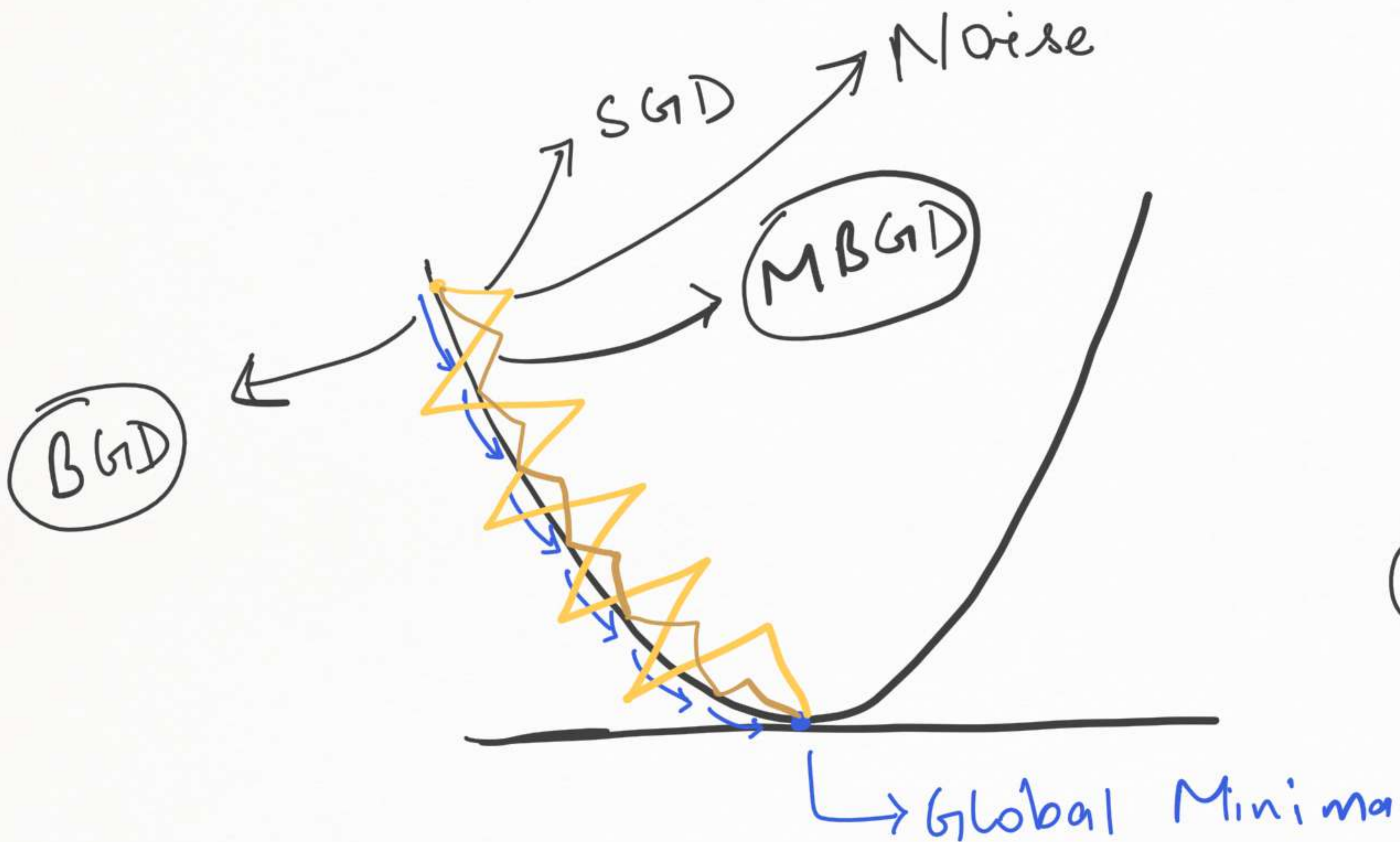
Epoch = 2

SGD → Noise

MBGD

BGD

Global Minima

Reduce Noise
⇓
"Momentum"

④ SGD with Momentum $\underline{\underline{\phantom{=}}}$

$$W_{new} = W_{old} - \eta \frac{\partial L}{\partial W_{old}}$$

$$W_t = W_{t-1} - \eta \frac{\partial L}{\partial W_{t-1}}$$

$W_1 \quad , W_2, W_3$

$$W_4 = W_3 - \eta \frac{\partial L}{\partial W_3}$$

$W_5$

"Exponential Weighted Average"

$t = $ current time.

$t-1 = $ previous time stamp

# Exponential Weight Average

$$t_1 \quad t_2 \quad t_3 \quad t_4 \quad ----\cdot- t_n$$

$$a_1 \quad a_2 \quad a_3 \quad a_7 \quad --- \quad \cdot-a_n$$

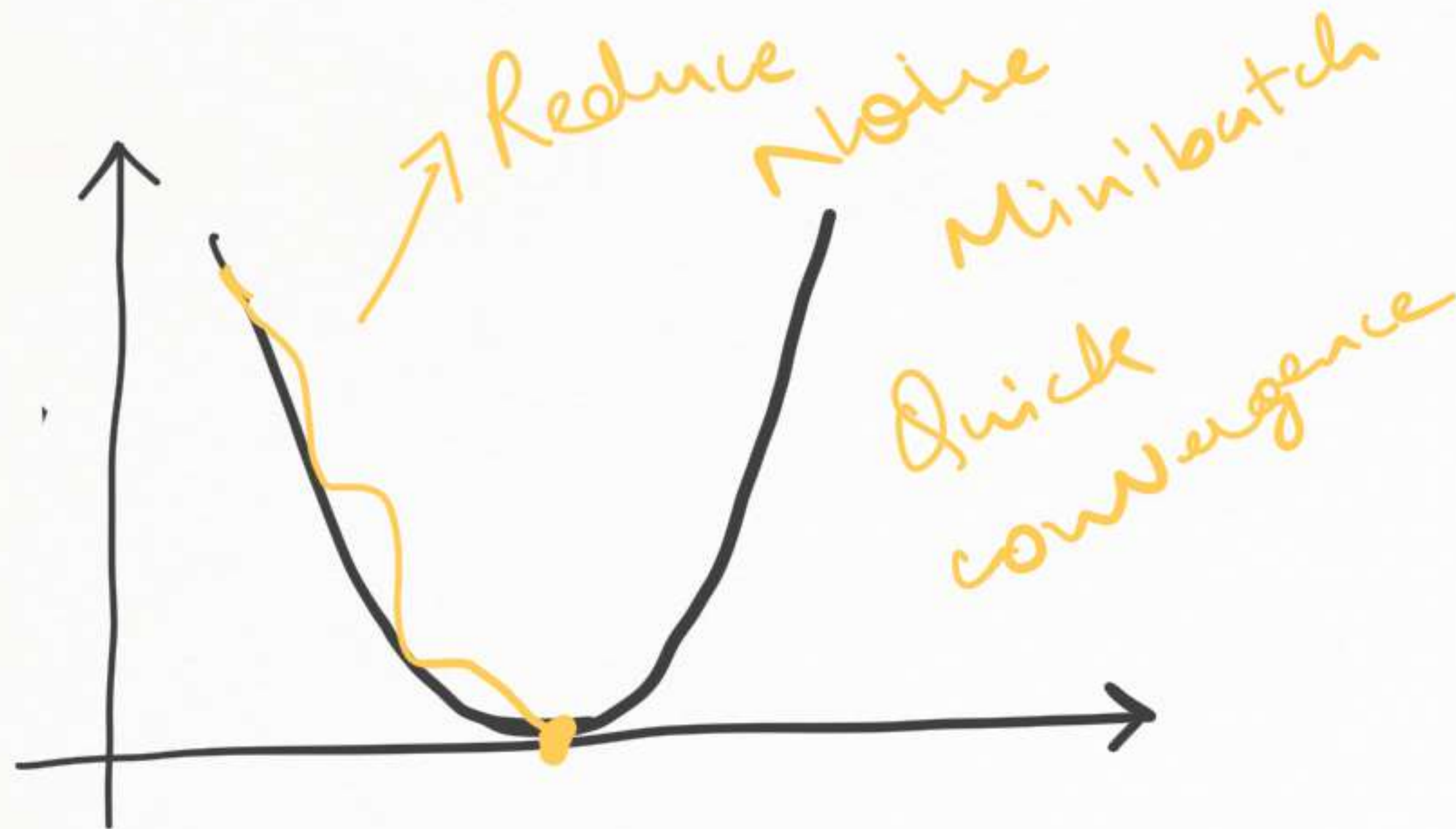$$\beta = 0 \text{ to } 1$$

$$\beta = 0.95$$

$$V_{t_1} = a_1$$

$$V_{t_2} = \beta * V_{t_1} + (1-\beta) * a_2$$

$$= 0.95 \times a_1 + 0.05 \times a_2$$

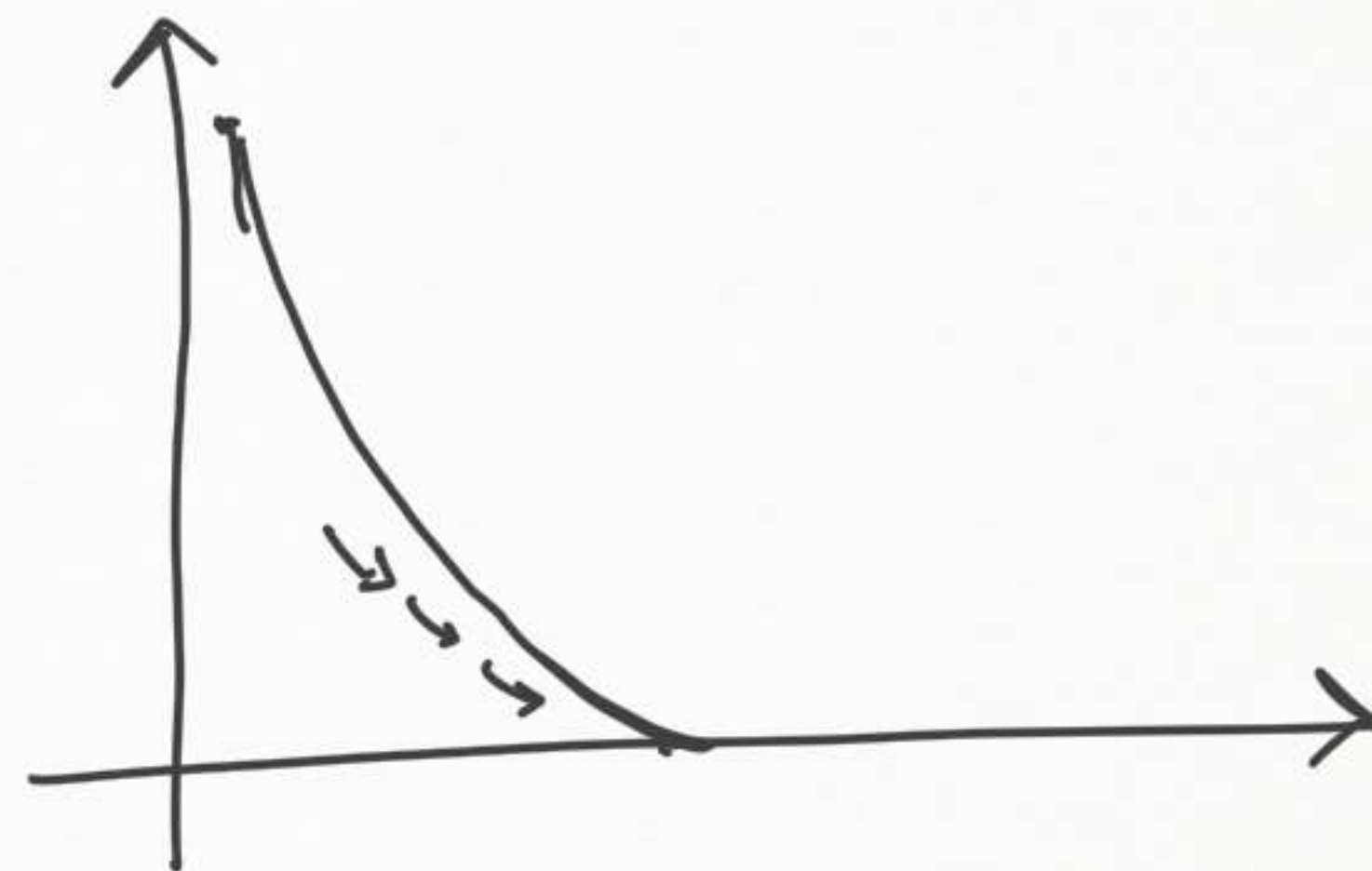$$W_t = W_{t-1} - \eta \, V_{dw}$$

"SGD with Momentum"

$$V_{dw} = \beta \times V_{dw_{t-1}} + (1-\beta) * \frac{\partial L}{\partial w_{t-1}}$$

Reduce Noise

Minibatch

Quick convergence

$\boxed{\eta} \longrightarrow$ fixed value

$\boxed{0.1} = 0.1$

⑤ Adagrad $\longrightarrow$ Adaptive Gradient Descent

$$\boxed{W_t = \omega_{t-1} - \eta' \frac{\partial L}{\partial \omega_{t-1}}}$$

$$\eta' = \frac{\eta}{\sqrt{\alpha_t + \epsilon}} \quad \xrightarrow{\text{small value}}$$

$$\alpha_t = \sum_{i=1}^{t} \left( \frac{\partial L}{\partial \omega_t} \right)^2$$

⑥ Adadelta & $\underline{\underline{RMSProp}}$

$$\boxed{\beta = 0.95}$$

$$S_{dw} = \textcircled{$\beta$} * S_{dw_{t-1}} +$$

$$\eta^1 = \frac{\eta}{\sqrt{\hat{S}_{dw} + \epsilon}}$$

$$(1-\beta)\left(\frac{\partial L}{\partial w_{t-1}}\right)^2$$

$\left\{ \begin{array}{c} \downarrow \\ \downarrow \\ \downarrow \end{array} \right\}$

$$\boxed{w_t = w_{t-1} - \eta^1 \frac{\partial w}{\partial w_{t-1}}}$$

$$0.95 + S_{dw_{t-1}}^L + 0.05 \times$$

(7) Adam Optimizer (Best Optimizer)

Adaptive + Momentum

(RMSProp + SGD with Mom.)

$$W_t = W_{t-1} - \eta' V_{dw}$$

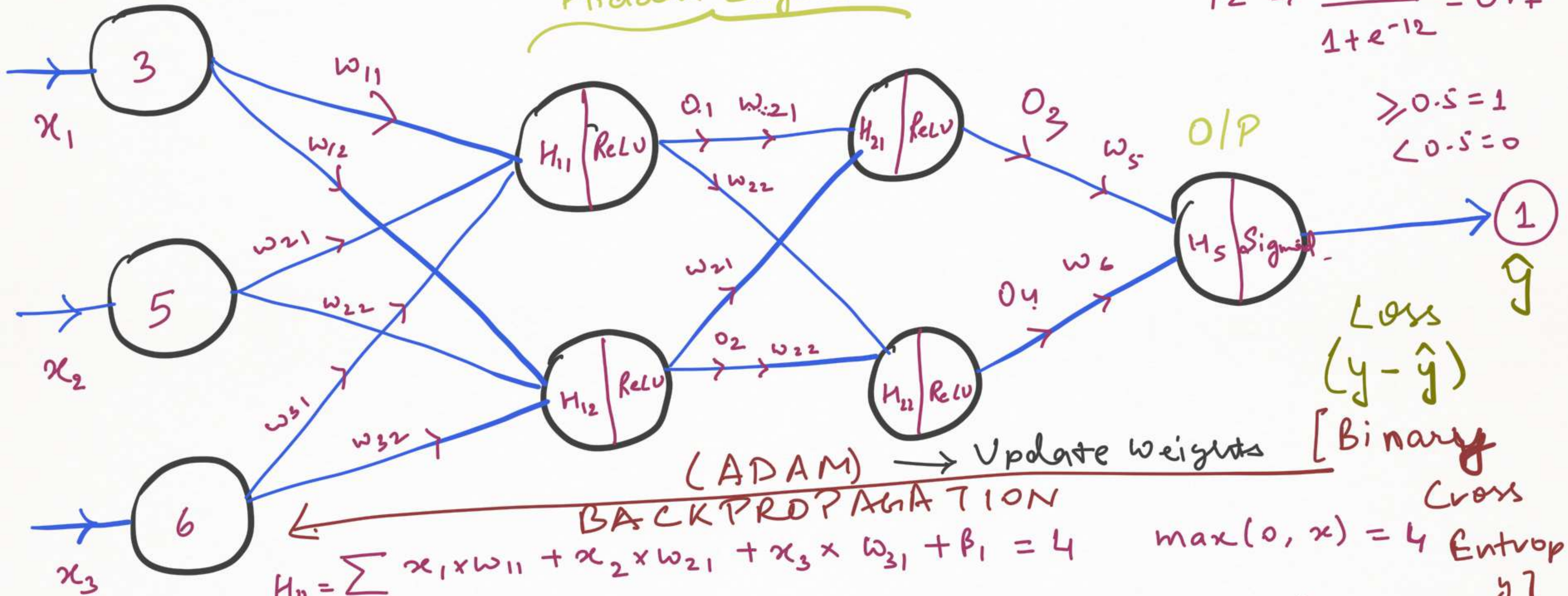$$\eta' = \frac{\eta}{\sqrt{S_{dw} + \epsilon}}$$

Epoch = 10

forward Propagation

Hidden Layers

I/P

3    $x_1$

5    $x_2$

6    $x_3$

$W_{11}$    $W_{12}$    $W_{21}$    $W_{22}$    $W_{31}$    $W_{32}$

$H_{11}$ ReLU    $H_{12}$ ReLU

$O_1$  $W_{21}$    $O_2$  $W_{22}$

$H_{21}$ ReLU    $H_{22}$ ReLU

$O_3$    $W_5$    $O_4$    $W_6$

$12 \Rightarrow \dfrac{1}{1+e^{-12}} = 0.7$

$\geqslant 0.5 = 1$

$< 0.5 = 0$

O/P

$H_5$ Sigmoid

1    $\hat{y}$

Loss

$(y - \hat{y})$

[Binary Cross Entropy]

(ADAM) → Update weights

BACKPROPAGATION

$H_{11} = \sum x_1 \times W_{11} + x_2 \times W_{21} + x_3 \times W_{31} + \beta_1 = 4$

$H_{12} = \sum x_1 W_{12} + x_2 \times W_{22} + x_3 \times W_{32} + \beta_2 = 3$

$max(0, x) = 4$

$max(0, x) = 3$

"Nested Loop"

Epoch
.
Batch Size