

1

#Deep learning

Artificial Intelligence

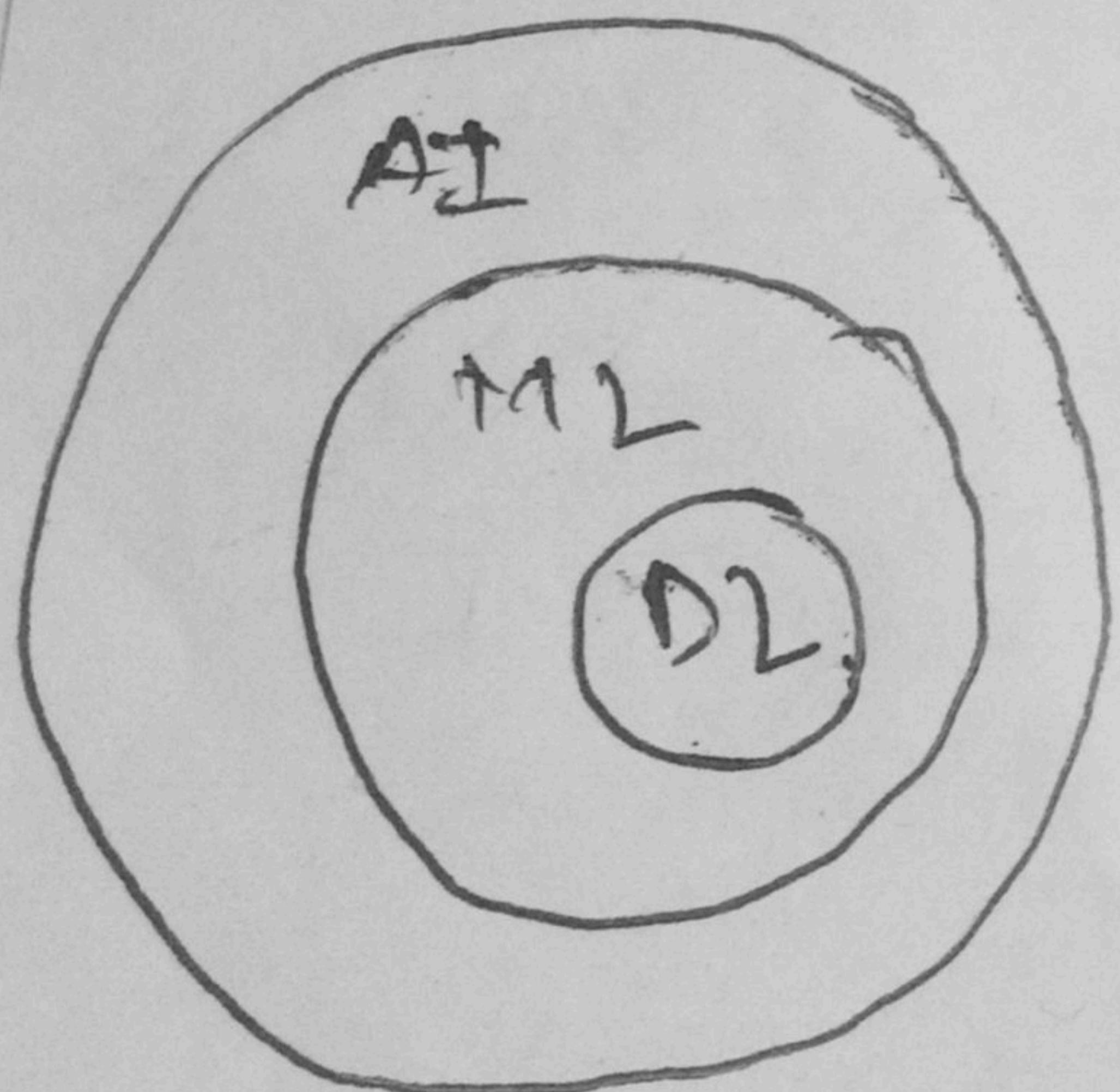
- Theory and development of computer systems able to perform tasks normally require human intelligence
- Anything that includes in ML, DL

Machine Learning

- Gives computer "the ability to learn without being explicitly programmed"
- SVM, KNN, RF, LinR, LogR, NB etc. rule based

Deep learning

- ML algorithms with brain-like logical structure of algorithms called artificial neural networks
- NN, CNN, RNN, LSTM, RL, Transformer etc. learning based



Deep learning Timeline (History of DL)

- 1940 → Dark Era (Until 1940)
- 1943 → Neutral Nets (McCulloch & Pitt) (study about biological neurons)
- 1950 → Computing Machinery and Intelligence (Alan Turing)
- 1958 → Perceptron (~~Krebspeblatt~~) (Rosenblatt)

2

(adaptive linear elements)

- 1960 → ADALINE (~~Widrow & Hoff~~)
- 1969 → XOR problem (Minsky & Papert) [Improved]
- 1974 → Backpropagation (Werbos and Mone)
- 1980 → ~~Neocitation~~ Neocognition (Fukushima)
- ↳ Self Organizing Map (Kohonen)
- 1982 → Hopfield Networks (John Hopfield)
- 1985 → Boltzmann Machine (Hinton ~~& Sejnowski~~)
- 1986 → Restricted Boltzmann Machine (Smolensky)
- ↳ Multilayer Perception (Rumelhart, Hinton & Williams)
- ↳ RNNs (Jordan) (Recurrent Neural Network)
- 1990 → LeNet (LeCun) meta(AI scientist)
- 1997 → LSTMs (Hochreiter & Schmidhuber) [long short term memory]
 ↳ Bidirectional RNN (Schuster & Paliwal)
- 2006 → Deep Boltzmann Machines (Salakhutdinov & Hinton)
- ↳ Deep Belief Networks - pretraining (Hinton)
- 2012 → Dropout (Hinton)
- 2014 → GRNs (Goodfellow)
- 2017 → Capsule Networks (Sabour, Frosst, Hinton)

AI Winter

ILSVRC - 2012 (ImageNet)

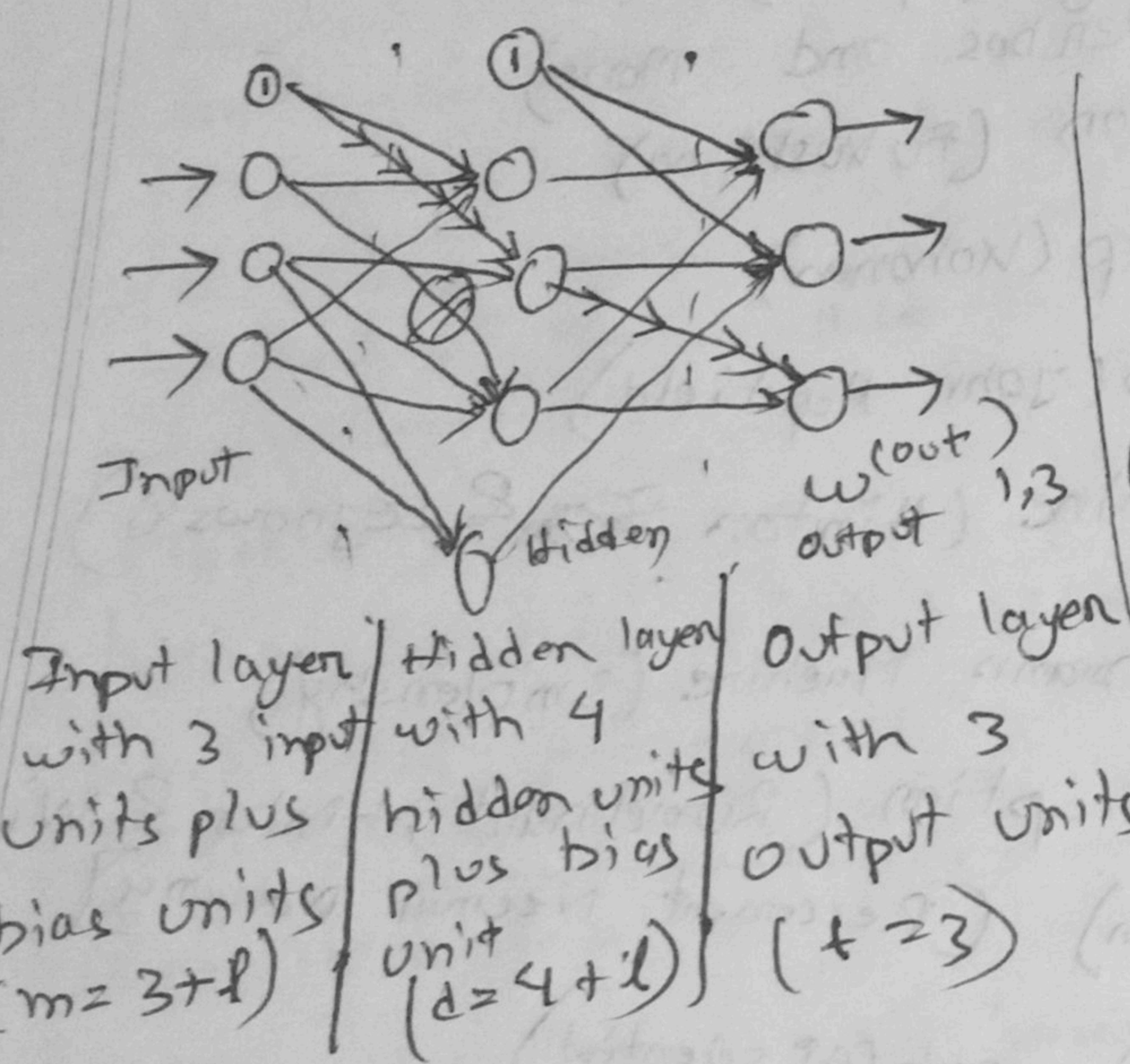
AlexNet - 2012

Hinton

GPU (2011-12)

3

Multi layer perceptron, 1986 (David Rumelhart)
Multiple Neuron



[Neural Networks - PESI]
Engine = Backpropagation

(1st input to 2nd hidden layer to the 3rd layer of output)

Number of layers $L = 3$

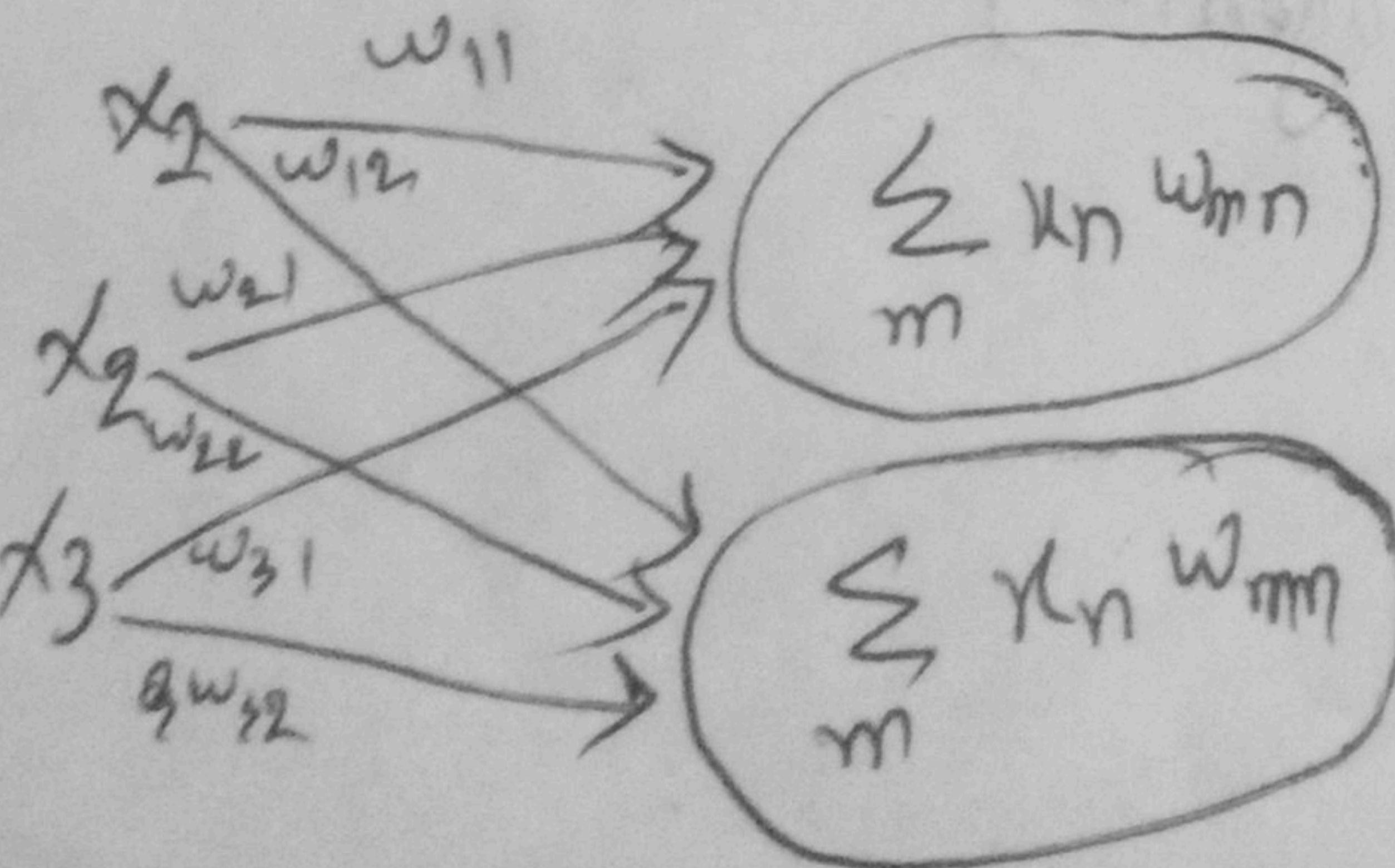
Non-linear units (how to train?)
Hinton picked it up \rightarrow back propagation (1974) chapter
 \rightarrow to train MLP

He published Nature paper

simple stacks of perceptrons?

MLP sounds good?

Components of MLP, Linear function



function output

$$z_m = f(\underbrace{u_m}_{\text{function inputs}}, \underbrace{w_{mn}}_{\text{function name}}) = b + \sum_m u_m w_{mn}$$

bias term

index for each neuron and matrix now

element n of feature vector(u)

~~Cost function~~
~~sum of squares errors for the entire dataset~~

element (m, n) for feature matrix(W)

$$x = \begin{bmatrix} u_1 \\ u_2 \\ u_3 \end{bmatrix} \Rightarrow w^T \cdot x = \begin{bmatrix} w_1 & w_2 & w_3 \end{bmatrix} \cdot \begin{bmatrix} u_1 \\ u_2 \\ u_3 \end{bmatrix}$$

$$z = w^T \cdot x = \begin{bmatrix} w_1 & w_2 & w_3 \end{bmatrix} \cdot \begin{bmatrix} u_1 \\ u_2 \\ u_3 \end{bmatrix}$$

Sigmoid function

\rightarrow An Activation function (a)

function output $\rightarrow a_m = \sigma$

(Activation function)

Non-linearit

Output functi

$$\hat{y} = f(a_m)$$

$$\hat{y} = f(a_m)$$

$$\hat{y} = \sigma$$

$$X = \begin{bmatrix} u_1 \\ u_2 \\ u_3 \end{bmatrix} \quad w = \begin{bmatrix} w_{11} & w_{12} & w_{13} \\ w_{21} & w_{22} & w_{23} \\ w_{31} & w_{32} \end{bmatrix}$$

$$\begin{array}{c} w_{11} & w_{12} & w_{13} \\ w_{21} & w_{22} & w_{23} \\ w_{31} & w_{32} \end{array} \quad \begin{array}{c} 3 \times 3 \\ 3 \times 1 \end{array}$$

4

$$Z = w^T \times X = \begin{bmatrix} w_{11} & w_{12} & w_{13} \\ w_{21} & w_{22} & w_{23} \\ w_{31} & w_{32} \end{bmatrix} \begin{bmatrix} u_1 \\ u_2 \\ u_3 \end{bmatrix} = \begin{bmatrix} u_1 w_{11} + u_2 w_{12} + u_3 w_{13} \\ u_1 w_{21} + u_2 w_{22} + u_3 w_{23} \\ u_1 w_{31} + u_2 w_{32} \end{bmatrix}$$

Sigmoid function

An Activation function (a)

function output

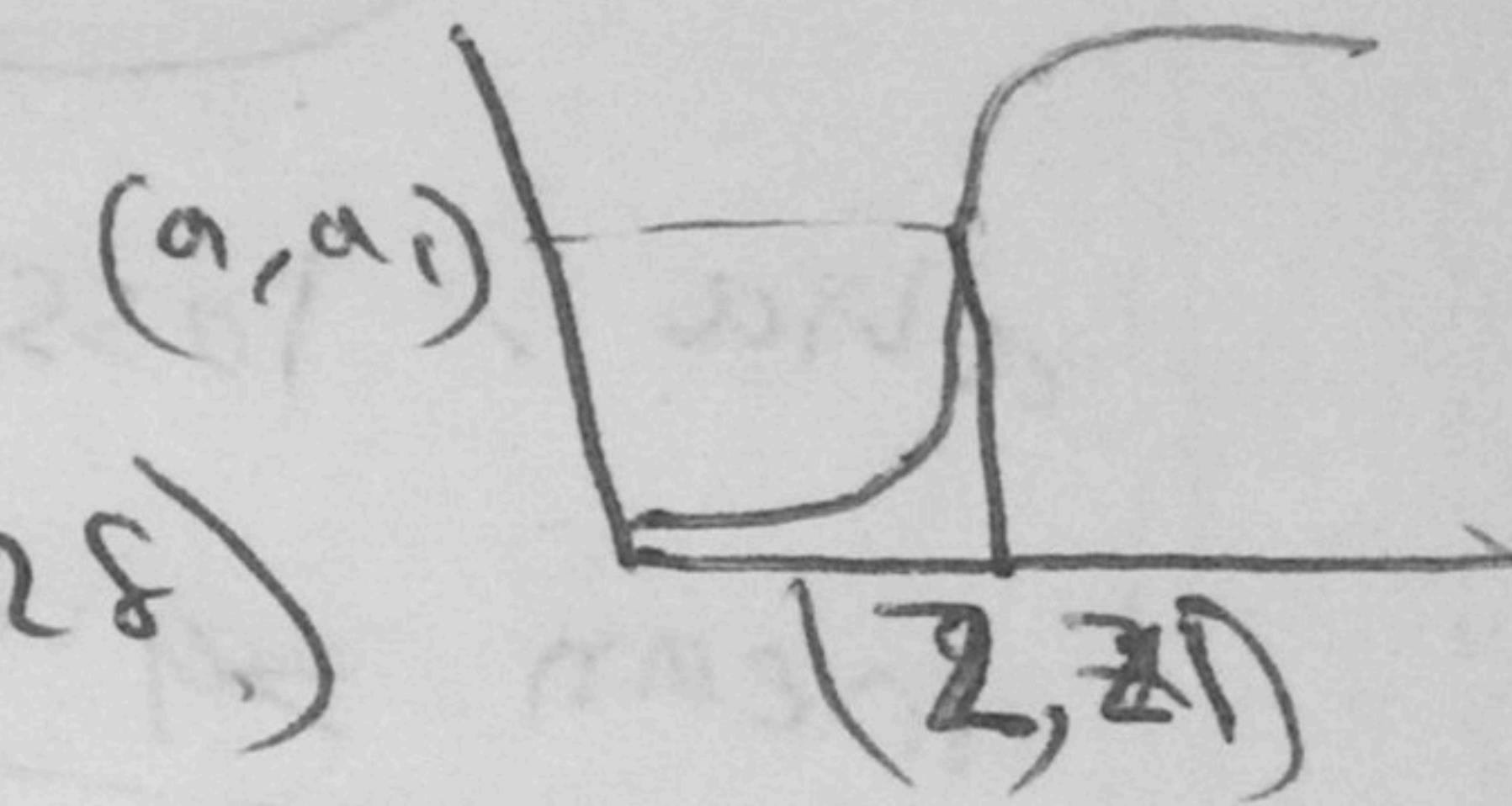
(Activation) function

linear function output becomes sigmoid input

sigma (sigmoid function)

Euler's number (≈ 2.71828)

$$a = \begin{bmatrix} a_1 \\ a_2 \end{bmatrix}$$



Non-linearity \rightarrow (to separate Non-linear MLP)

Output function

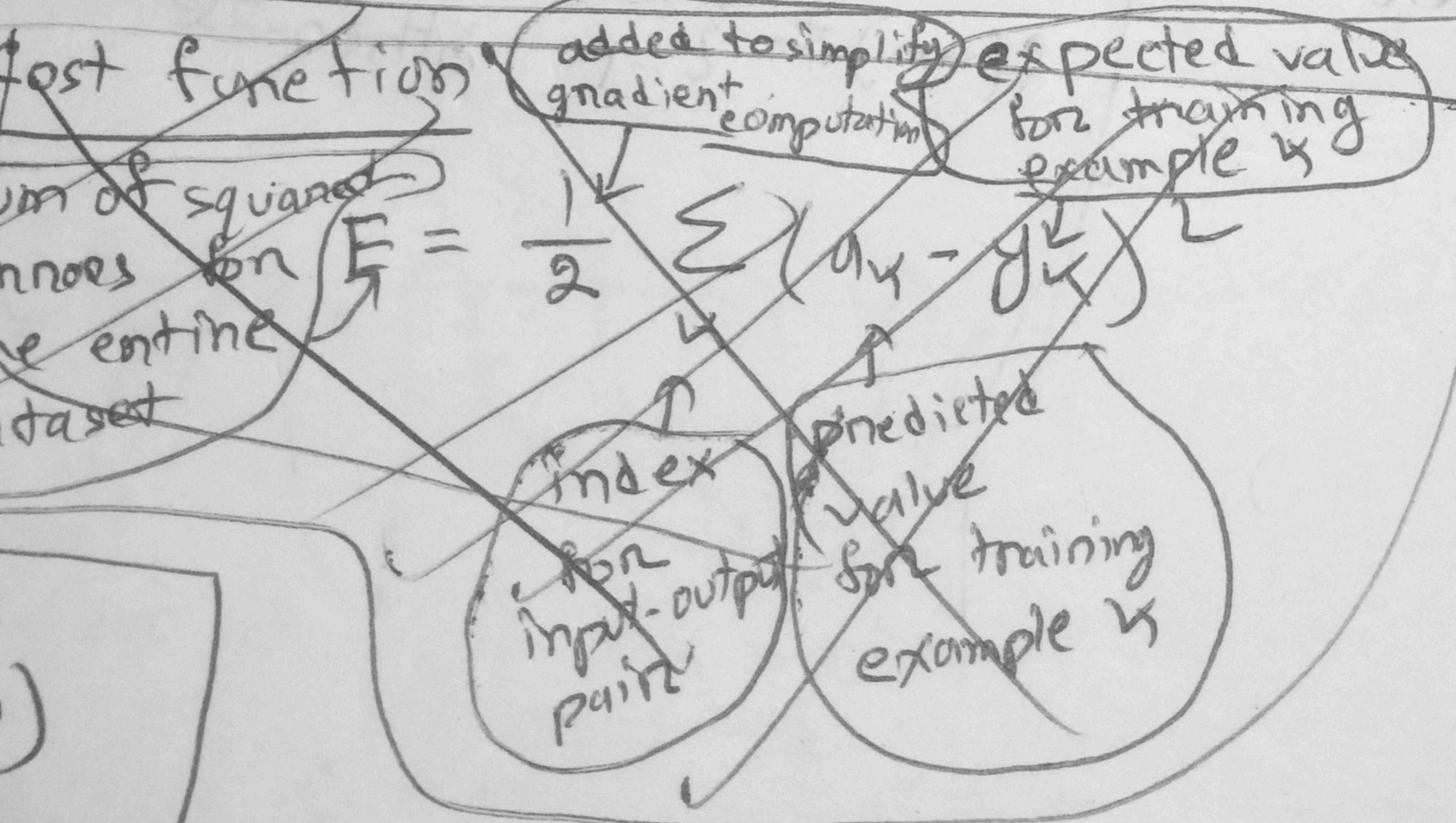
$$\hat{y} = f(a_m) \begin{cases} +1, & \text{if } a > 0.5 \\ -1, & \text{otherwise} \end{cases}$$

- Threshold function
- Binary classification

- Identity function
- Regression problem

$$\hat{y} = \sigma(a); = \frac{e^{Ba_i}}{\sum_{j=1}^k e^{Bz_j}}$$

- Softmax function
- Multiclass classification



Cost function

$$E = \frac{1}{2} \sum_k (a_{ik} - y_{ik})^2$$

For training example k
 ↓
 Predicted value
 ↓
 expected value

↓
 index for input-output pair

added to simplify gradient computation.

forward function to optimize parameters

also: loss function, cost function, objective function

mean squared error, sum of squared errors, binary cross-entropy

Forward propagation (Forward pass)

$$\begin{bmatrix} w_{11} & w_{12} & w_{13} \\ w_{21} & w_{22} & w_{23} \end{bmatrix}^T \begin{bmatrix} u_1 \\ u_2 \end{bmatrix}$$

matrix multiplication
input to hidden layer

matrix multiplication
hidden to output layer

$$\begin{bmatrix} w_1 & w_2 & w_3 \end{bmatrix} \begin{bmatrix} a_1 \\ a_2 \\ a_3 \end{bmatrix}$$

linear function

$$z_m = b + \sum_n u_n w_{mn}$$

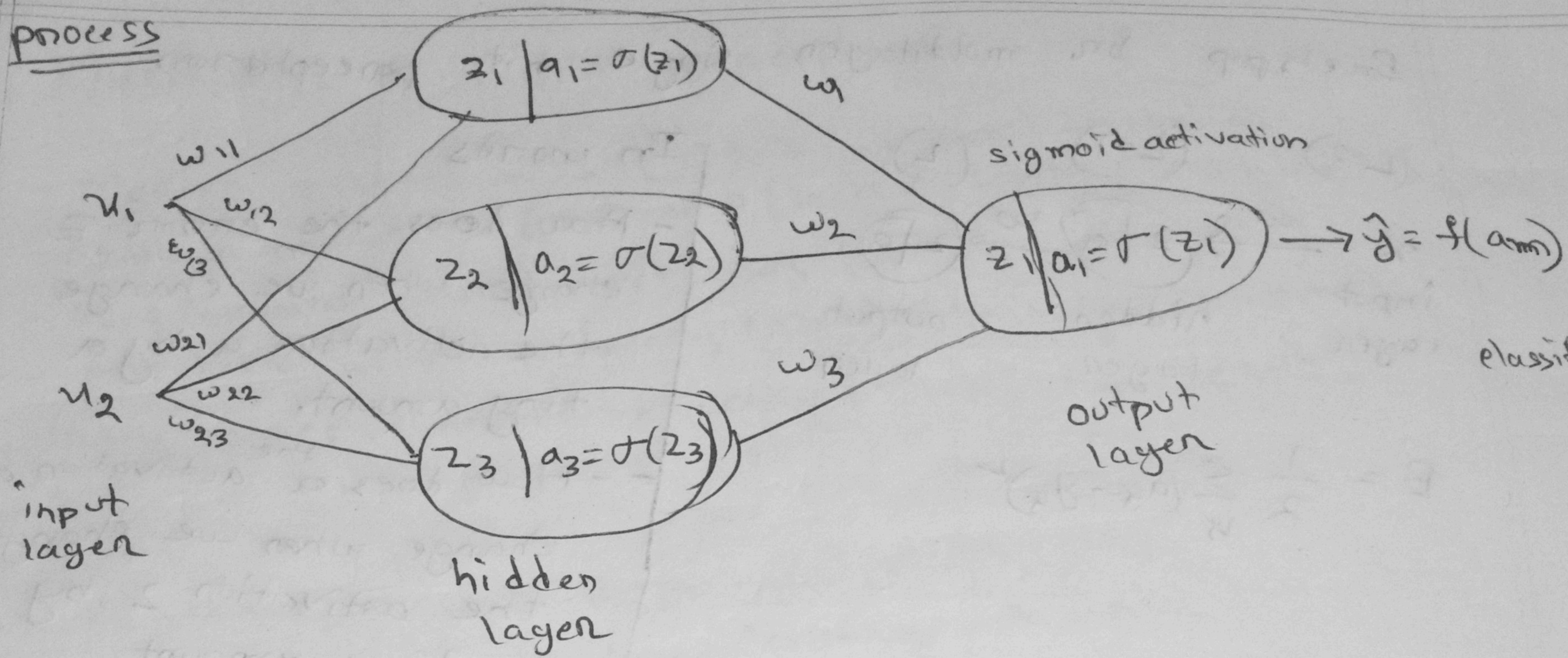
sigmoid function

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

decision threshold

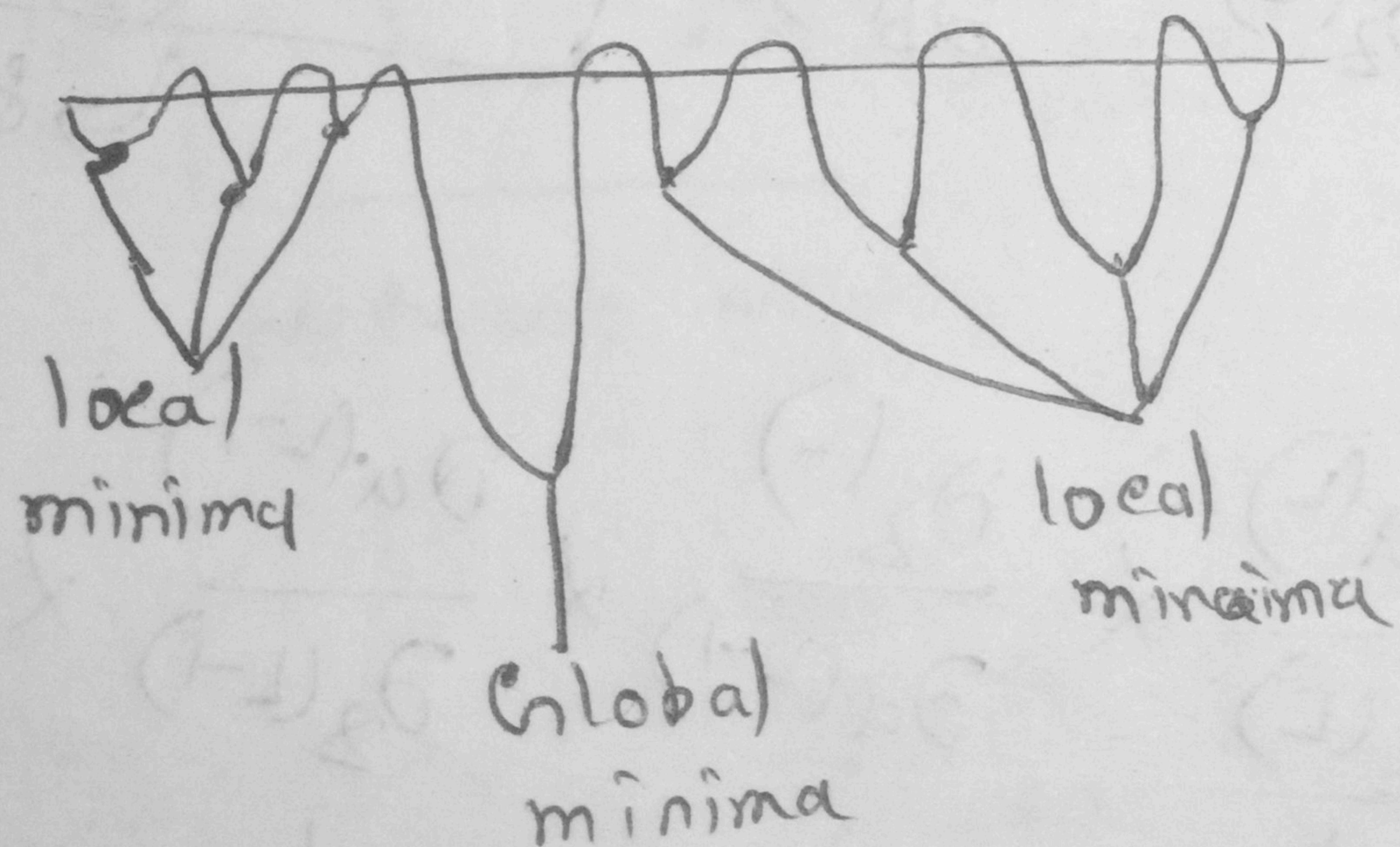
$$f(a_m) = \begin{cases} +1 & \text{if } a_m > 0.5 \\ -1 & \text{otherwise} \end{cases}$$

6



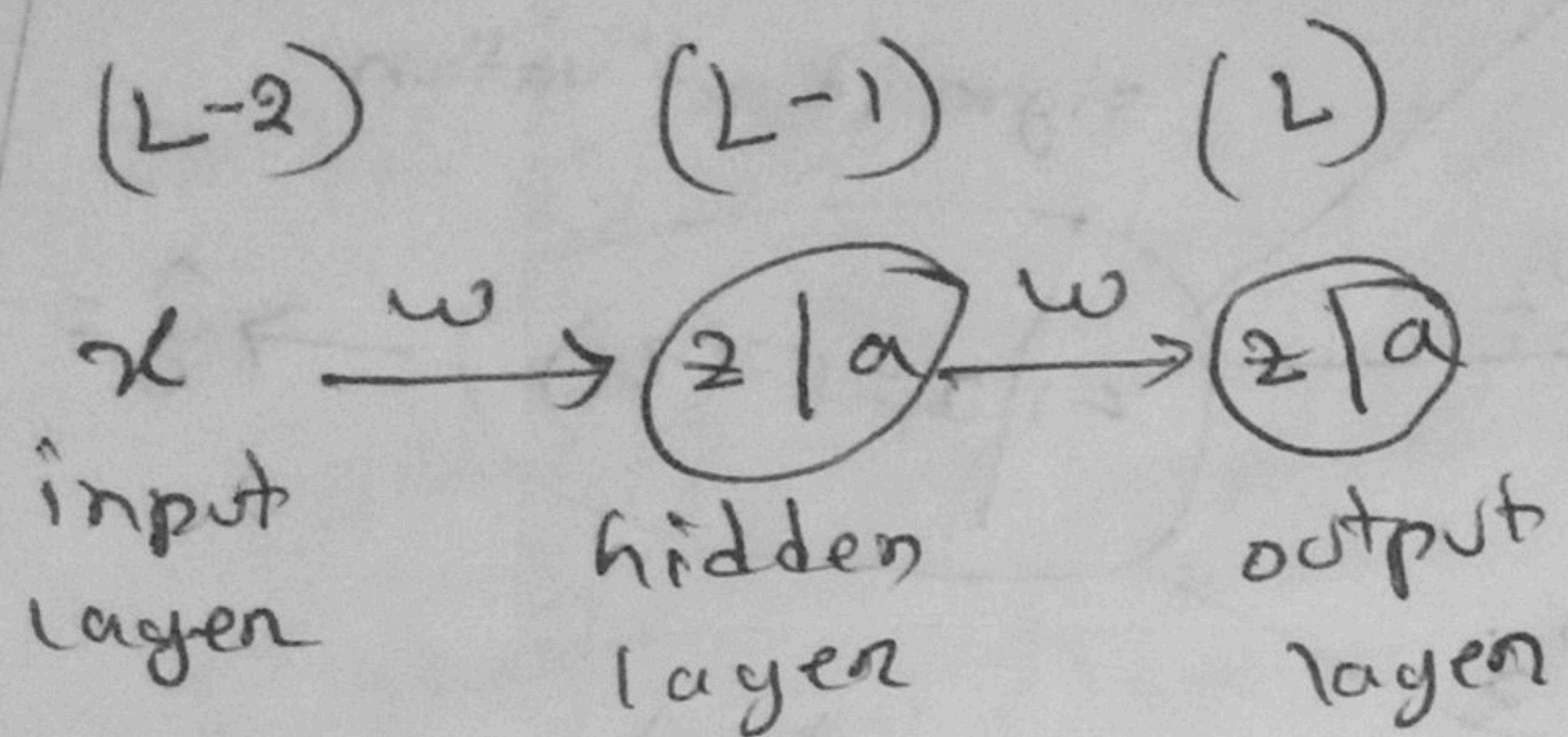
Backpropagation algorithm

- minimize the error gradient descent
- convex and non-convex optimization
- introducing nonlinearities
- multiple "valleys" with "local minima"



✓

Backprop for multilayer single unit perceptron



$$E = \frac{1}{2} \sum_{i=1}^n (a_i - y_i)^2$$

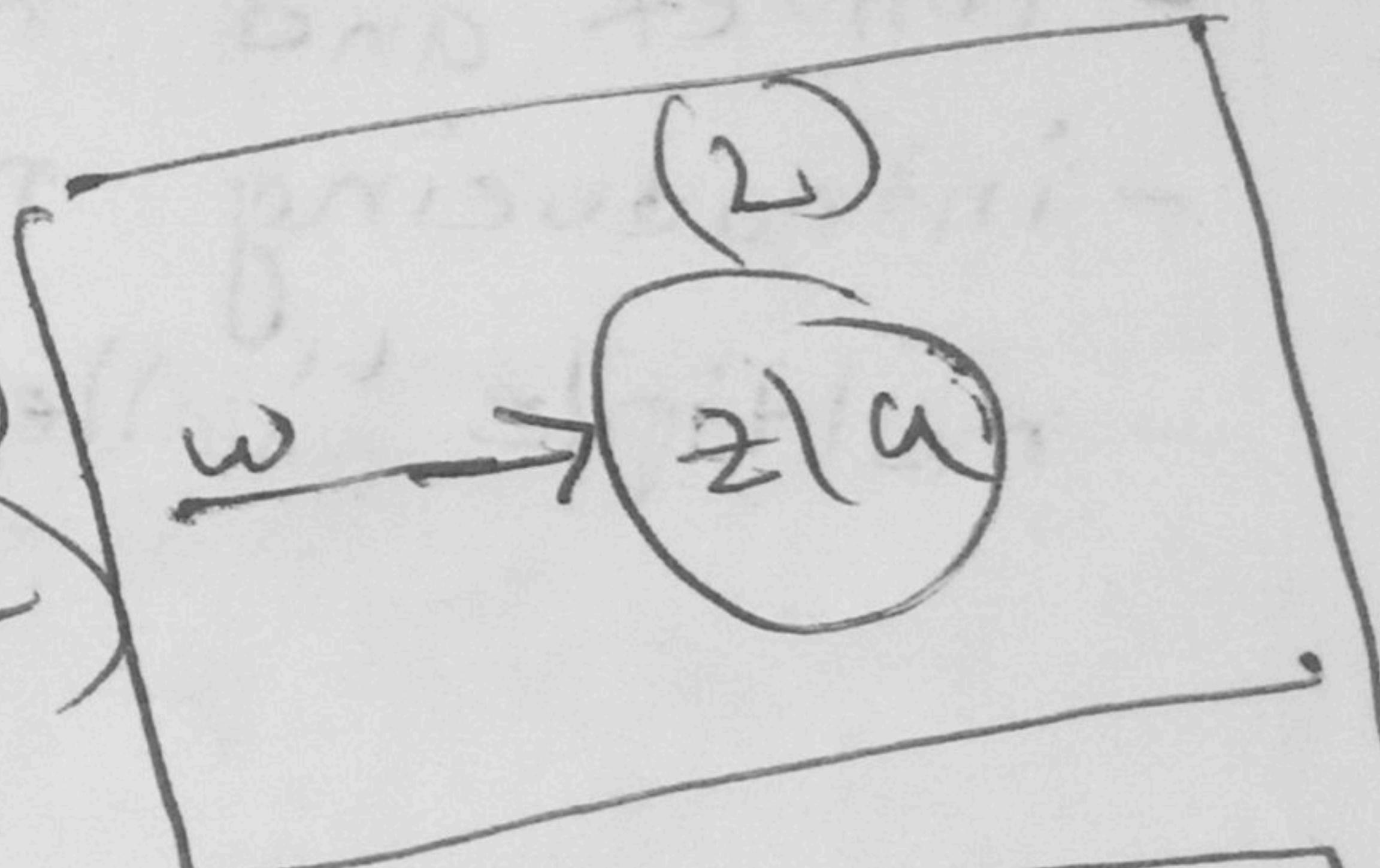
In words

- How does the error E change when we change the activation a by a tiny amount.
- How does the activation a change when we change the activation z by a tiny amount
- How does z change when we change the weights w by a tiny amount

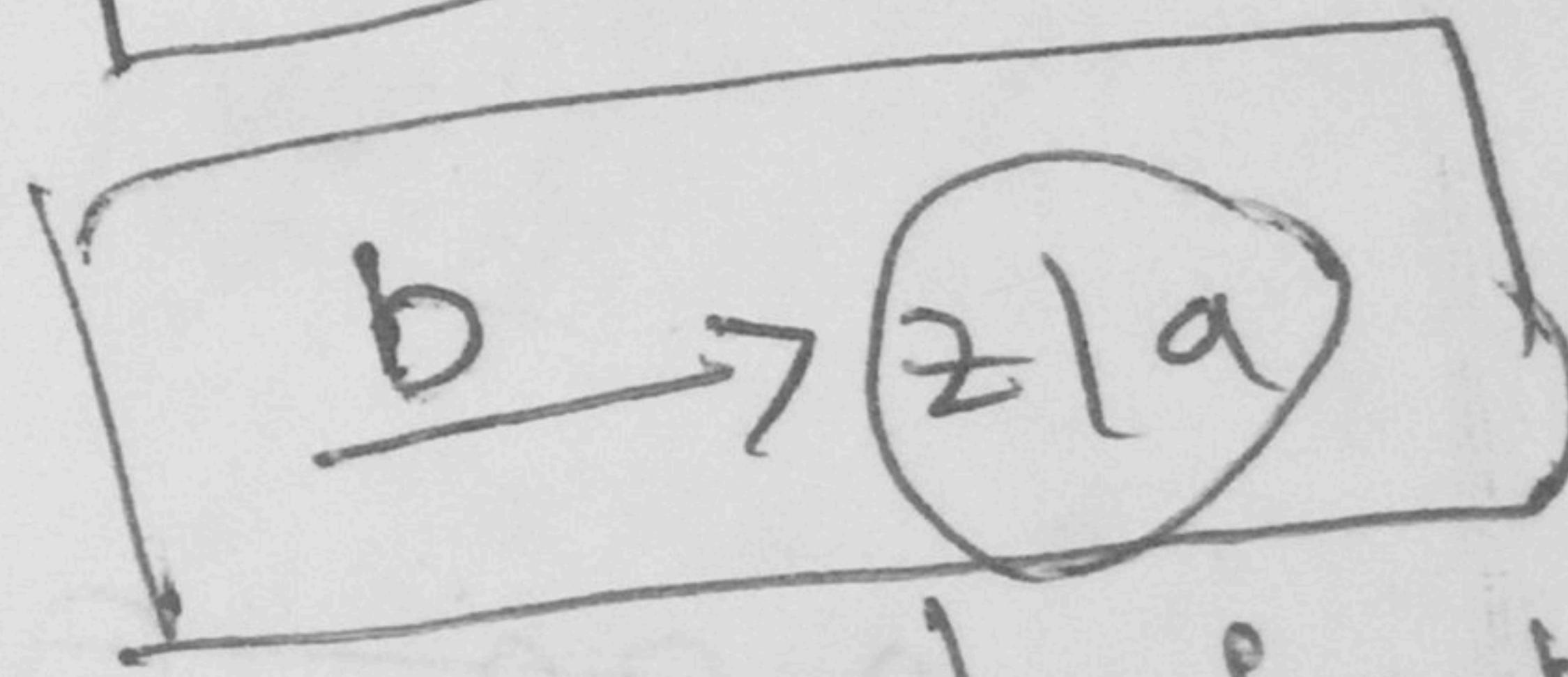
Partial derivative

single unit

$$\frac{\partial E}{\partial w^{(L)}} = \frac{\partial E}{\partial a^{(L)}} \times \frac{\partial a^{(L)}}{\partial z^{(L)}} \times \frac{\partial z^{(L)}}{\partial w^{(L)}}$$



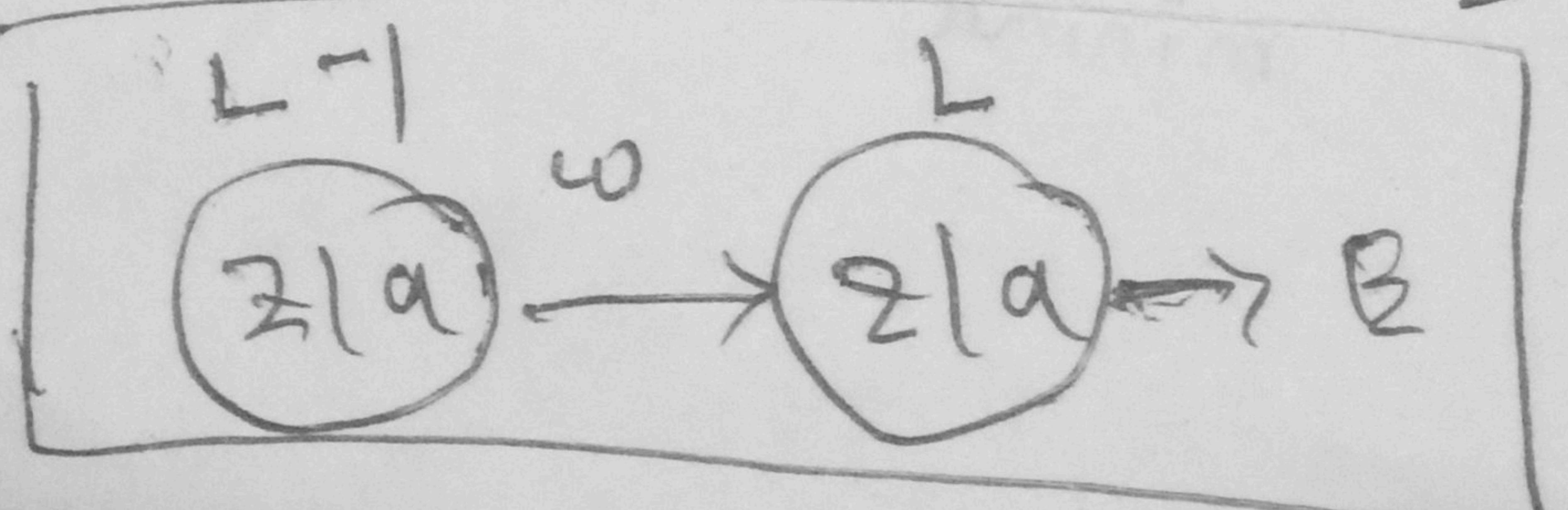
$$\frac{\partial E}{\partial b^{(L)}} = \frac{\partial E}{\partial a^{(L)}} \times \frac{\partial a^{(L)}}{\partial z^{(L)}} \times \frac{\partial z^{(L)}}{\partial b^{(L)}}$$



for bias
inaccurate

Double unit

$$\frac{\partial E}{\partial w^{(L-1)}} = \frac{\partial E}{\partial a^{(L)}} \times \frac{\partial a^{(L)}}{\partial z^{(L)}} \times \frac{\partial z^{(L)}}{\partial a^{(L-1)}} \times \frac{\partial a^{(L-1)}}{\partial z^{(L-1)}} \times \frac{\partial z^{(L-1)}}{\partial w^{(L-1)}}$$



[17]

[18]

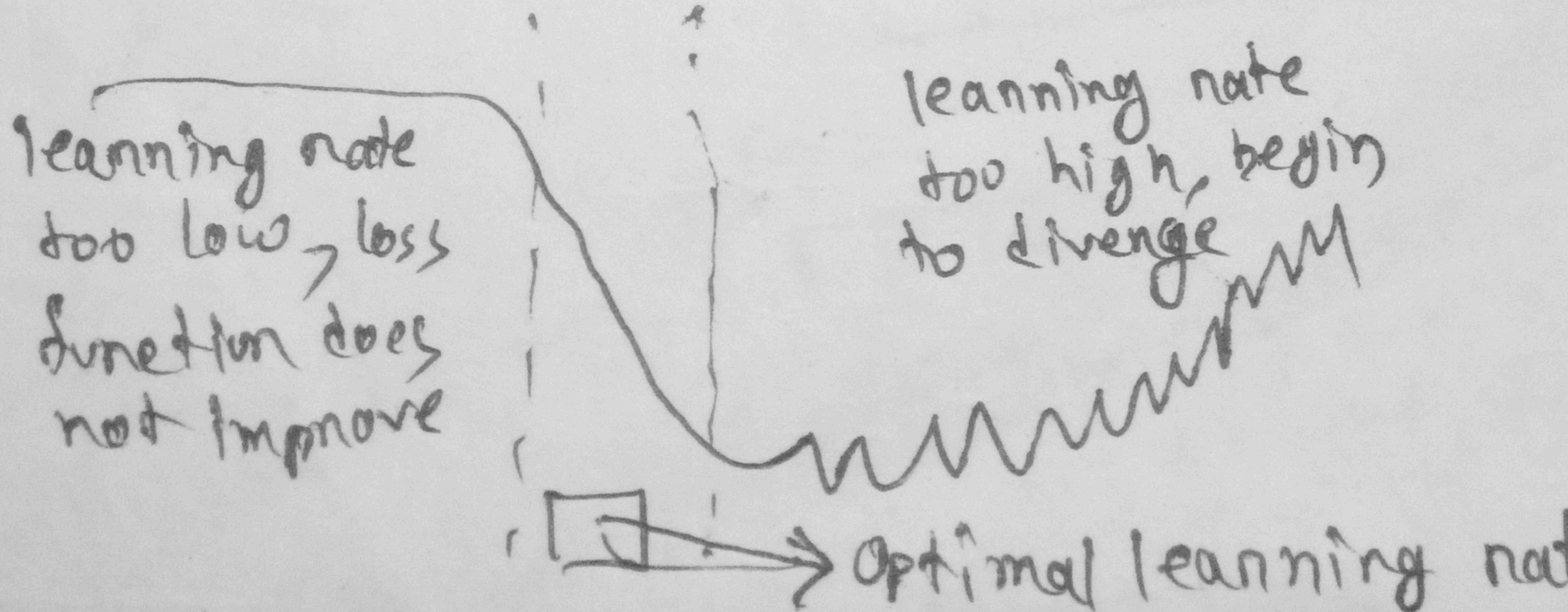
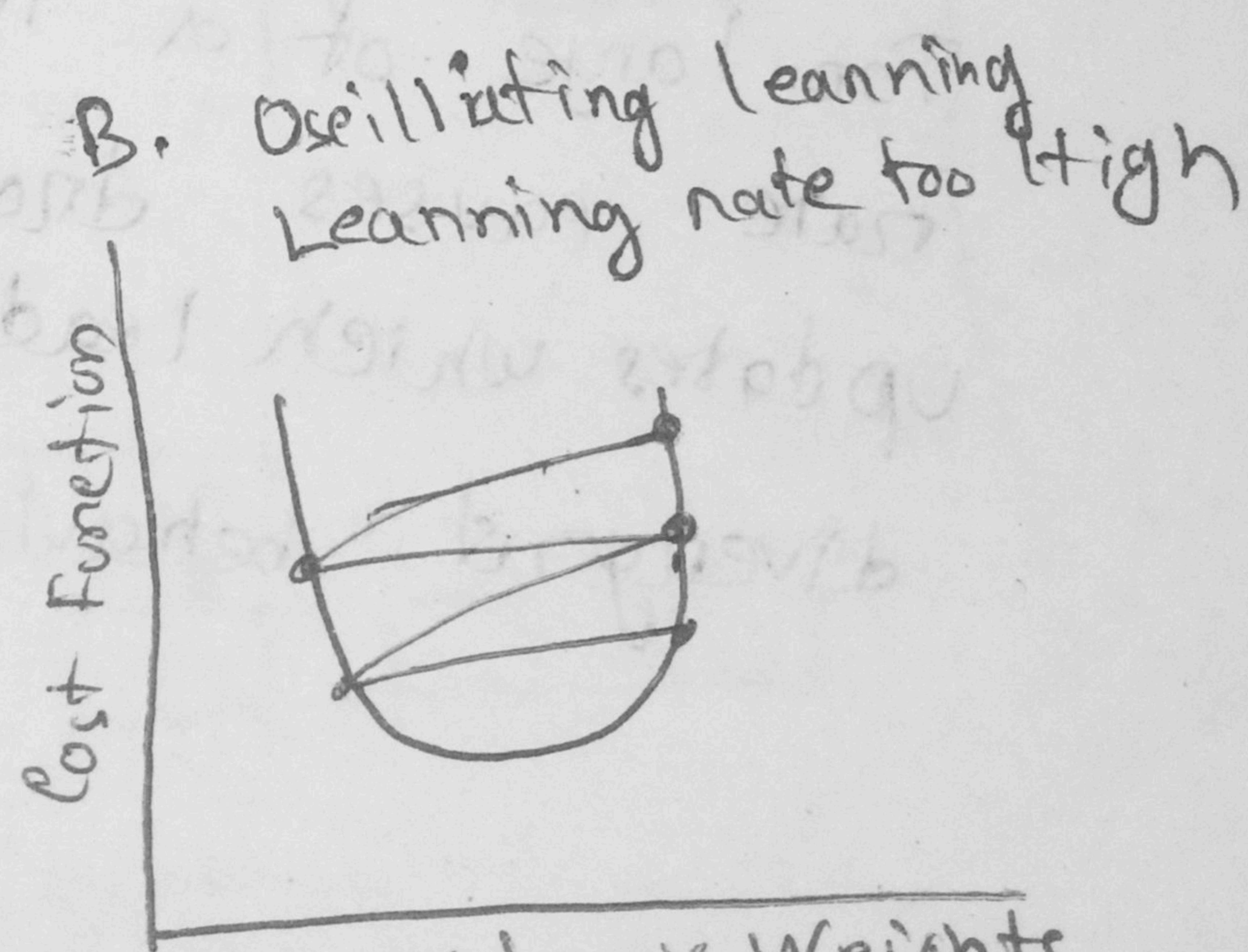
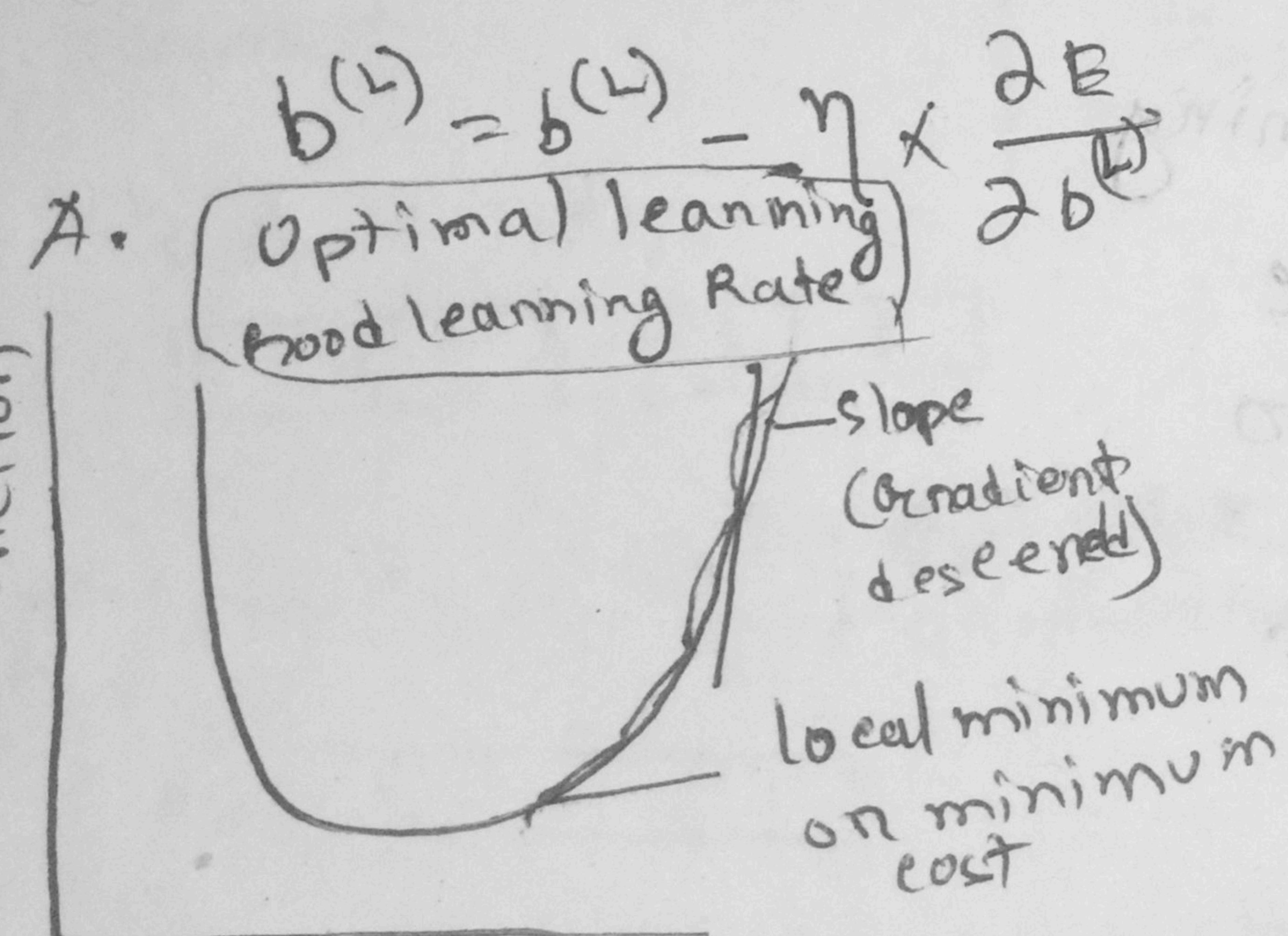
Multilayer Multi-unit Perception

- Modern neural nets (1980's - present) → similar
- shallow networks (not deep | lack of depth) (less hidden layers)
- Feed forward neural networks (Forward propagation)
- Artificial neural networks (ANN) [brain like structure]
- Fully connected neural networks (FCNN) (All connected)
- FC
- 1980's
- why not deep?
 - Resource limitation
 - Unavailability of Data

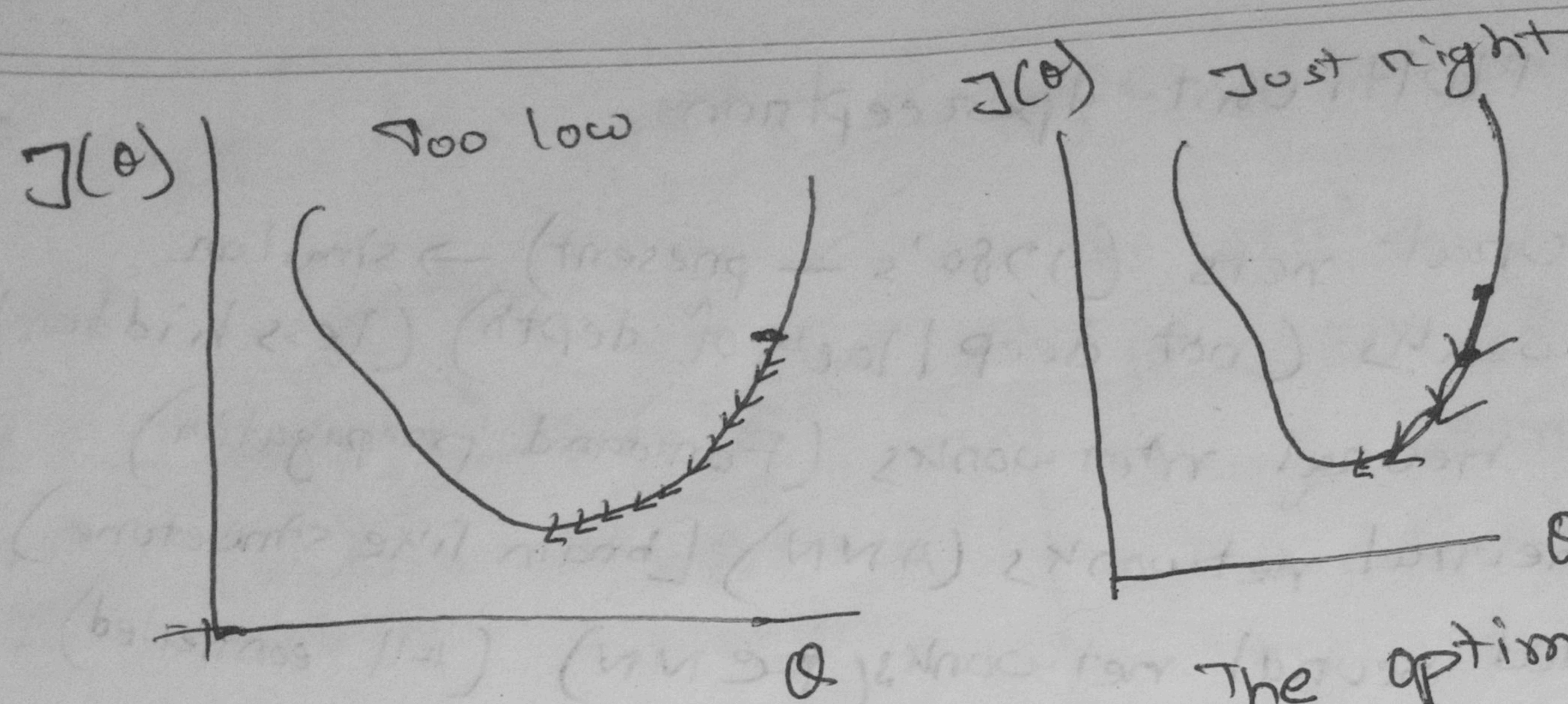
MLP Learning rule

$$w_{jk}^L = w_{jk}^L - \eta \times \frac{\partial E}{\partial w_{jk}^L}$$

η = learning rate
↓ partial derivative

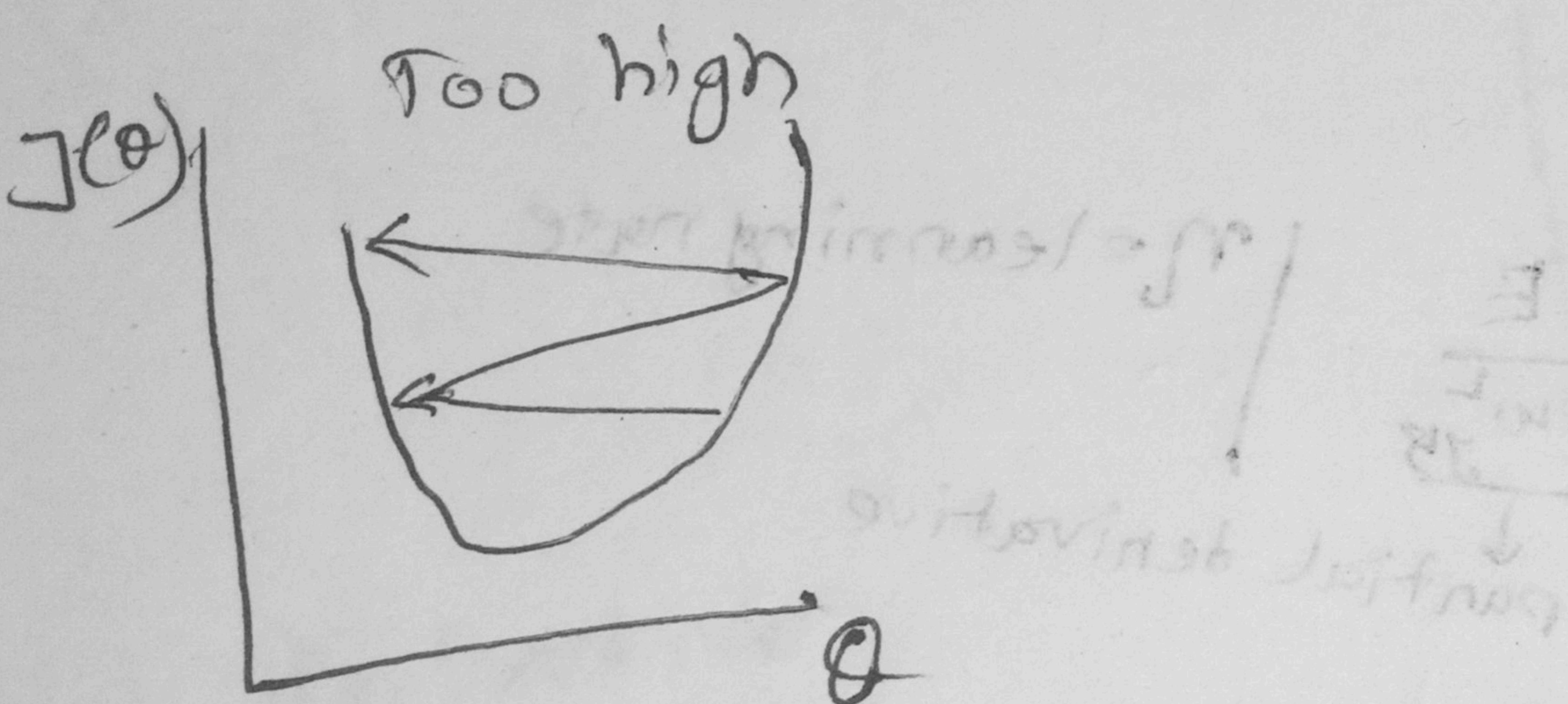


9



A small learning rate requires many updates before reaching the minimum point.

The optimal learning rate swiftly reaches the minimum point.



Too large of a learning rate causes drastic updates which lead to divergent behaviors.

10

MP neuron	Peneeptron	ADALINE	MLP
1943	1958	1960	1986
Linear aggregation, f	Linear aggregation, f	Linear aggregation, f	Linear aggregation, f
No weight	Weighted sum	Weighted sum	Weighted sum
No bias	No bias	No bias	Bias
Manual threshold, f	Heaviside step,f	Heaviside step,f	Sigmoid,f
Linear units	Linear units	Linear units	Non-linear units
can not solve XOR	can not solve NOR	sub-optimal solution	can solve
No Learning rule	Simple learning rule	Learning rule on GD(Gradient descent)	Learning rule on BP(Back-prop- agation)
No LR (Learning rate)	LR	LR	LR
Slowest	Faster	More Faster	Fastest

11

convolution (in 2D)

Polynomials

$$2 + 3u + 4u^2 \rightarrow 2 \text{ Degree}$$

$$2 + 3u + 4u^2 + 5u^3 + 6u^4 + 7u^5 \rightarrow 5 \text{ Degree}$$

$$4 + 12u + 25u^2 + 34u^3 + 43u^4 + 54u^5 + 45u^6 + 28u^7 \rightarrow 7 \text{ Degree}$$

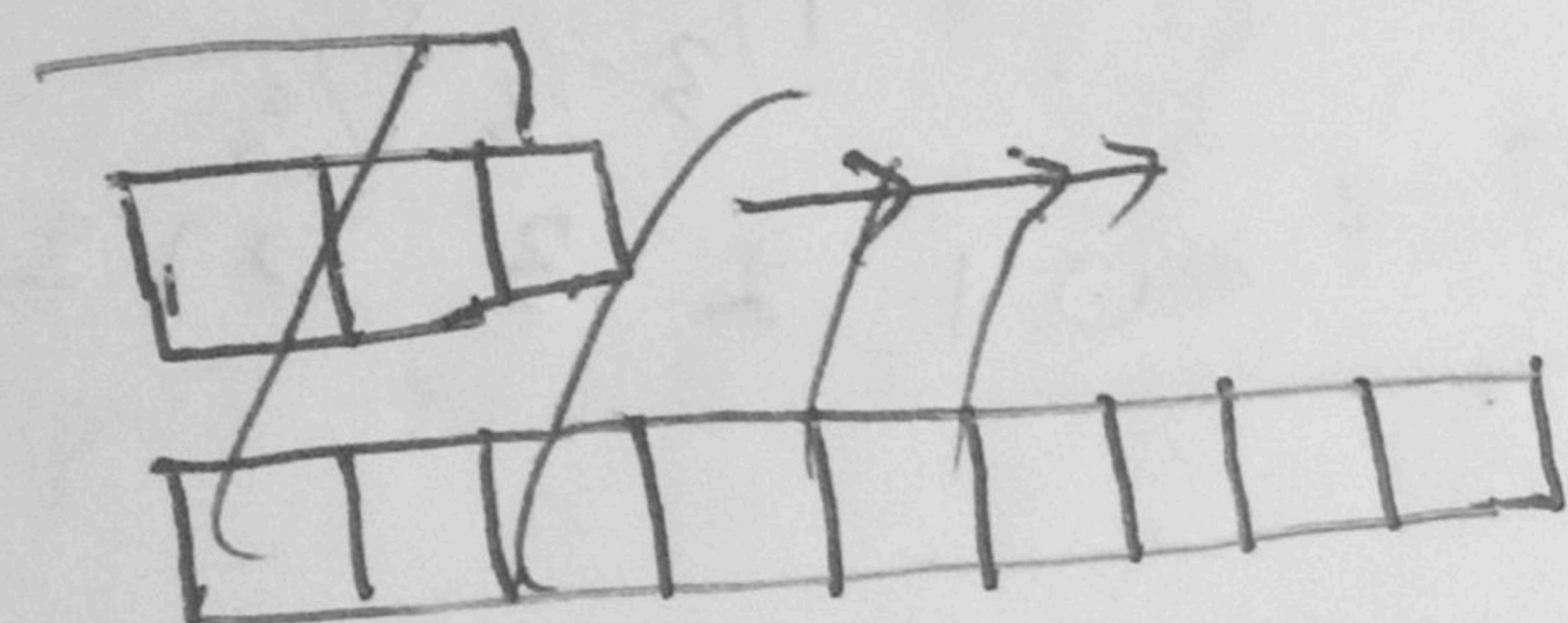
$$(u_0 + u_1 u + u_2 u^2)$$

$$(a_0 + a_1 u + a_2 u^2 + a_3 u^3 + a_4 u^4 + a_5 u^5)$$

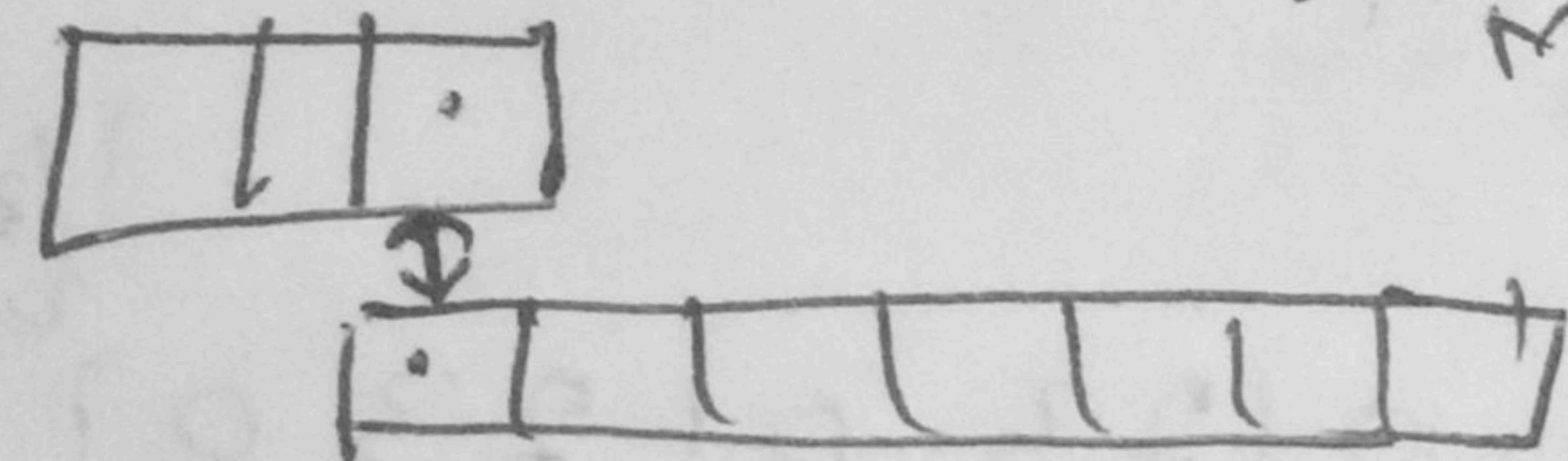
$$(u_2 a_2 + u_1 a_3 + u_0 a_4) u^4$$

→ Multiplying Both ⇒

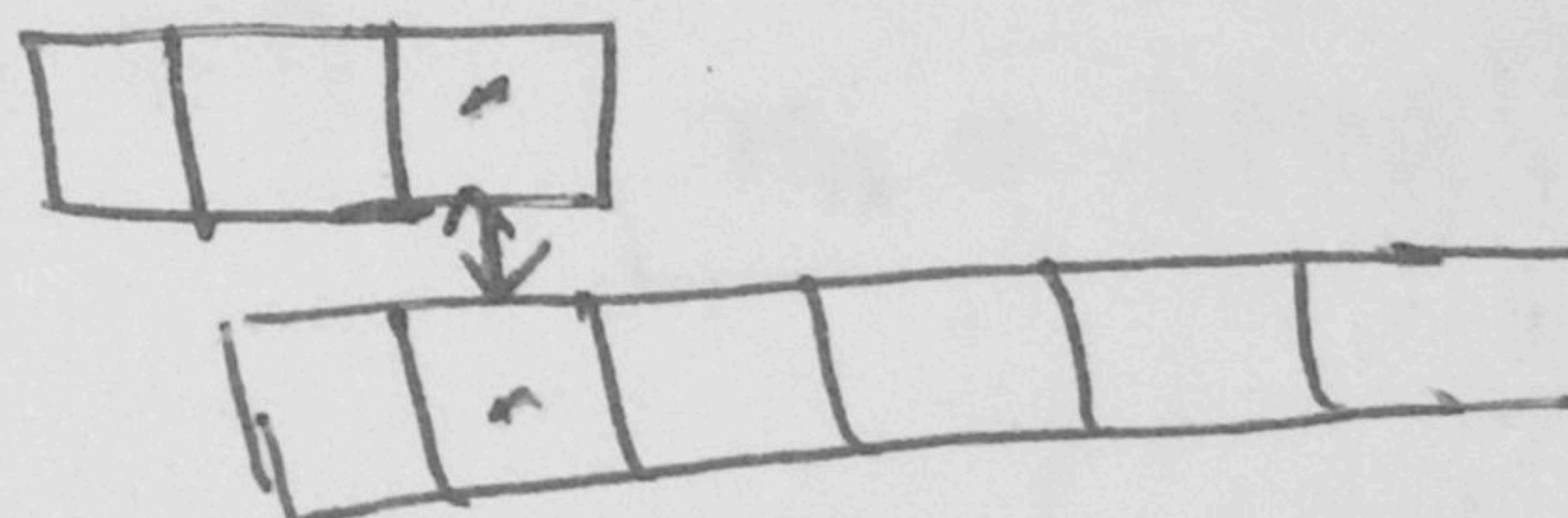
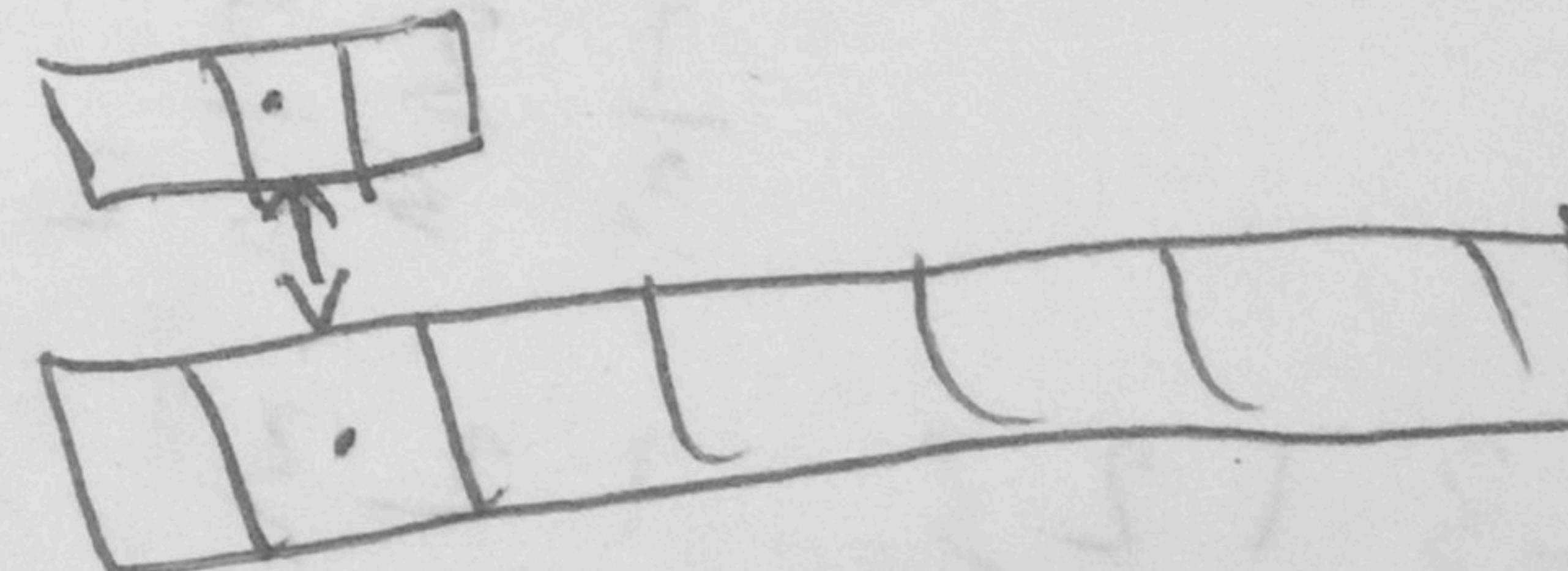
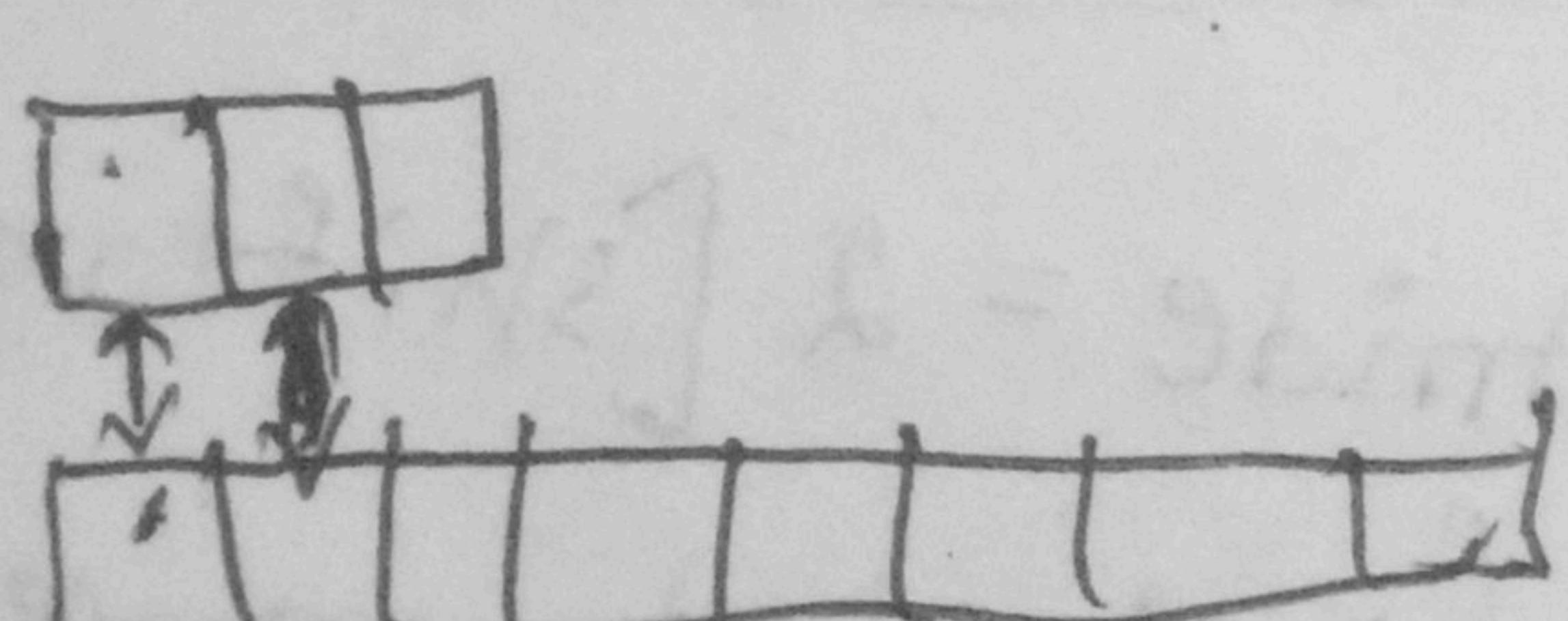
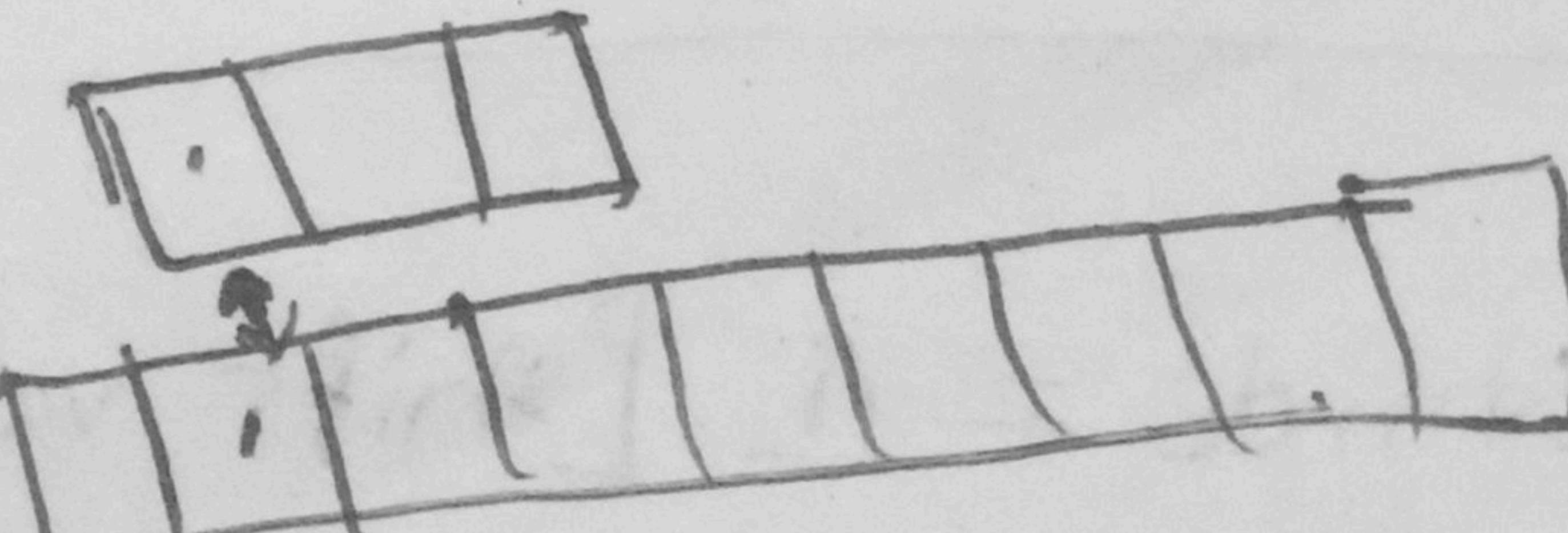
$$(K_0 a_0) + (u_1 a_0 + u_0 a_1) u + (u_2 a_0 + u_1 a_1 + u_0 a_2) u^2 + (u_2 a_1 + u_1 a_2 + u_0 a_3) u^3 + \\ (u_2 a_2 + u_1 a_3 + u_0 a_4) u^4 + (u_2 a_3 + u_1 a_4 + u_0 a_5) u^5 + (u_2 a_4 + u_1 a_5) u^6 + (u_2 a_5) u^7$$



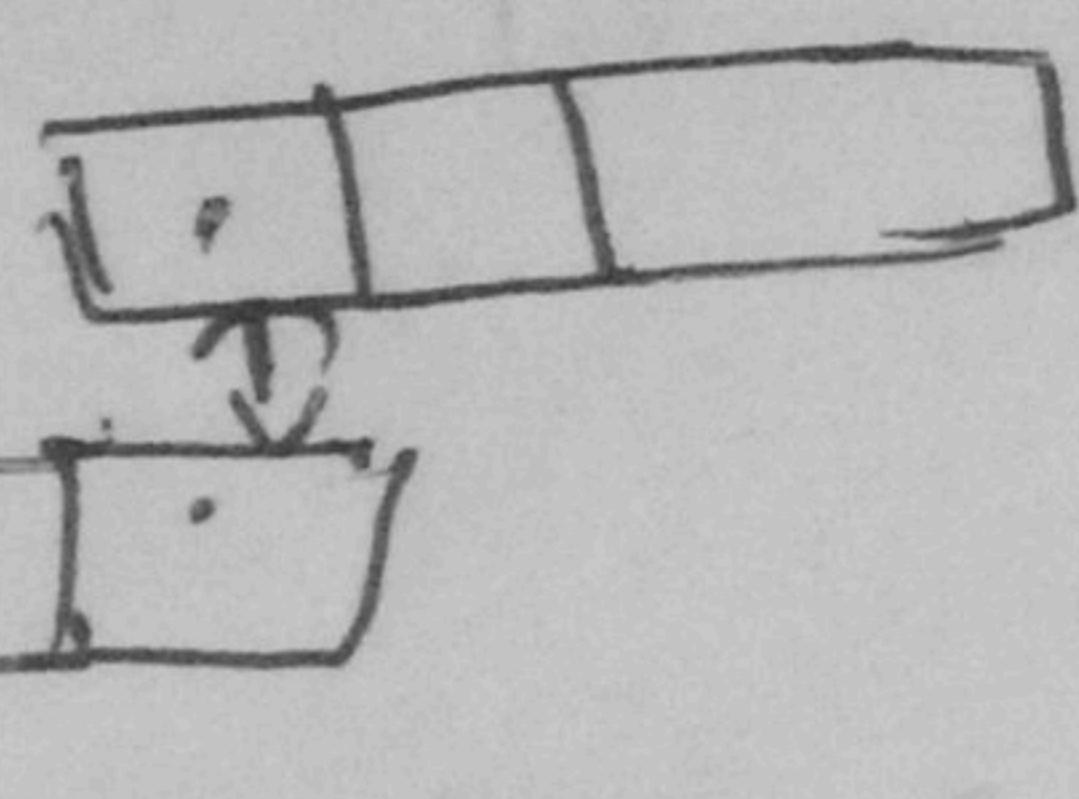
continues



→ ~~4th~~ multiplication

5th6th

nth (last)



It is convolution simulation with steps

convolution in one dimension

$$a = [1/3, 1/3, \cancel{1/3}]$$

$$u = [1, 2, 3, 1, 2, 3, 1, 2, 3]$$

$$\begin{matrix} 1^{\text{st}} & \frac{1}{3} & \frac{1}{3} & \frac{1}{3} \\ & 1 & 2 & 3 & 1 & 2 & 3 & 1 & 2 & 3 \end{matrix}$$

$$= 2 \quad \begin{matrix} \frac{1}{3} & \frac{1}{3} & \frac{1}{3} \\ 0 & 1 & 2 & 3 \end{matrix}$$

$$\begin{matrix} 2^{\text{nd}} & \\ & 1 & 2 & 3 \\ & \frac{1}{3} & \frac{1}{3} & \frac{1}{3} \\ & 1 & 2 & 3 & 1 & 2 & 3 & 1 & 2 & 3 \end{matrix} = 2$$

$$\begin{matrix} 2^{\text{nd}} & \\ & \frac{1}{3} & \frac{1}{3} & \frac{1}{3} \\ & 0 & 1 & 2 & 3 & 0 & 1 & 2 & 3 & 0 \end{matrix}$$

$$\begin{matrix} 3^{\text{rd}} & \\ & : & : & : \\ & \vdots & \vdots & \vdots \\ & 1 & 2 & 3 & 1 & 2 & 3 & 1 & 2 & 3 \end{matrix} = 2$$

$$\begin{matrix} 3^{\text{rd}} & \\ & \frac{1}{3} & \frac{1}{3} & \frac{1}{3} \\ & 0 & 1 & 2 & 3 & 0 & 1 & 2 & 3 & 0 \end{matrix}$$

$$\begin{matrix} 7^{\text{th}} & \\ & 1 & 2 & 3 \\ & \frac{1}{3} & \frac{1}{3} & \frac{1}{3} \\ & 1 & 2 & 3 & 1 & 2 & 3 \end{matrix} = 2$$

$$1 \ 2 \ 3 \ 1 \ 2 \ 3 \ 1 \ 2 \ 3$$

$\text{stride} = 1$ [shift value]
↳ steps to move

$\text{stride} = 2$ [shift value]
↳ steps to move

12

13

$$y(1) = \left\{ \frac{1}{3} \frac{1}{3} \frac{1}{3} \right. \\ \left. 0012312312300 \right\} = \frac{1}{3}$$

$$y(2) = \left\{ \frac{1}{3} \frac{1}{3} \frac{1}{3} \right. \\ \left. 0012312312300 \right\} = -1$$

$$y(1) = \left\{ 0012312312300 \right\} = 1$$

Total number of padding, Extra (rightmost) leftmost zero's (counted)
 we apply padding so that we do not loss any information.

$$\text{Output size} = \frac{n_x + 2P - n_h}{s} + 1$$

$$= \frac{9 + 2 \cdot 2 - 3}{1} + 1 = 11$$

$$= \frac{9 + 2 \cdot 0 - 3}{1} + 1 = 7$$

$$= \frac{9 + 2 \cdot 1 - 3}{2} + 1 = 5$$

[
 P = Padding
 S = stride]

[
 n_x = length of Input length
 n_h = length of filter]

Convolution 'n 2D (Gray scale Image)

Kernel = filter

	m	0	1	2
0	n	-1	-2	-1
1		0	0	0
2		1	2	1

2D image
Gray scale
image

Kernel
filter

kernel size = Matrix shape

↳ normally square shape

→ Average value always in center

Fill up extra 0's

↳ (Apply padding)

$\Rightarrow -13$

$\boxed{111} \rightarrow \text{Kernel}$

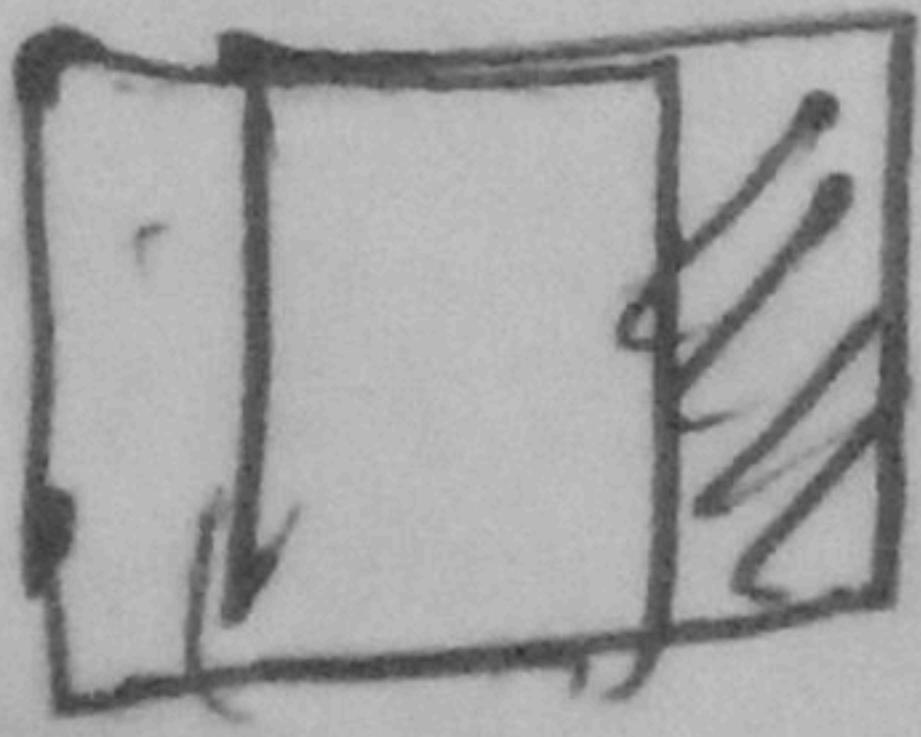
$\Rightarrow -20$

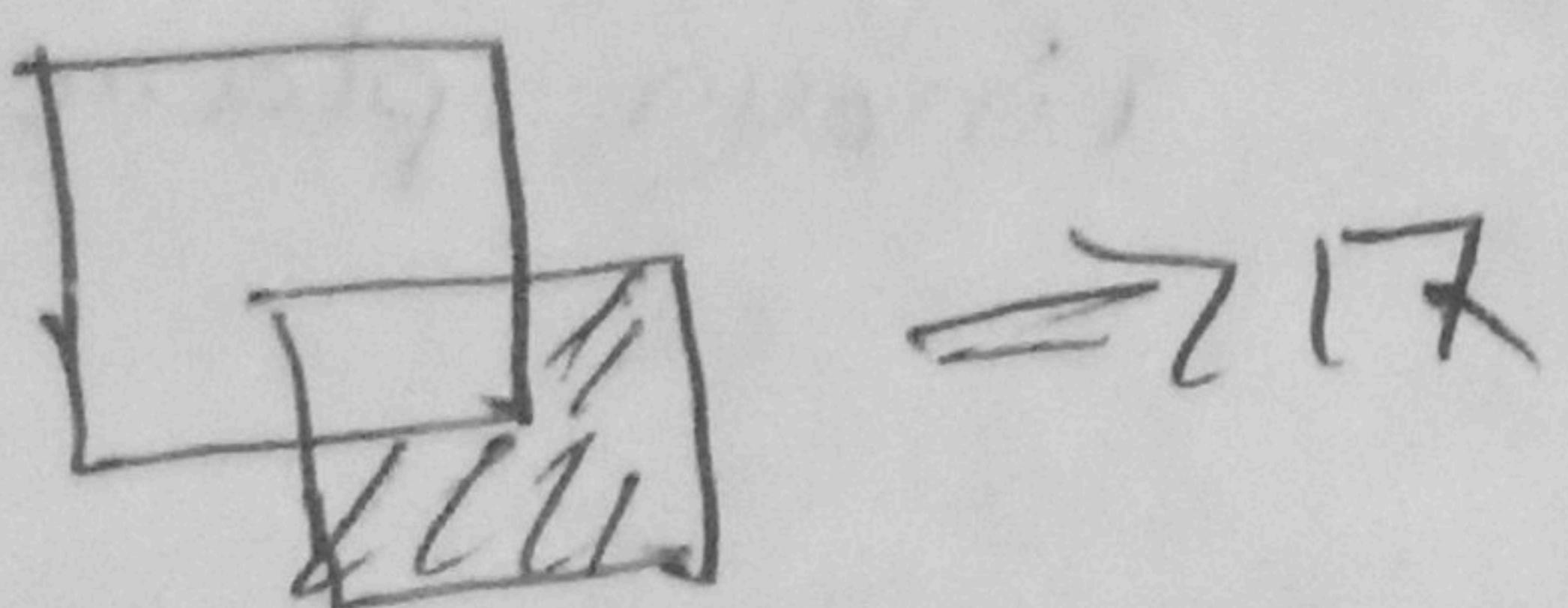
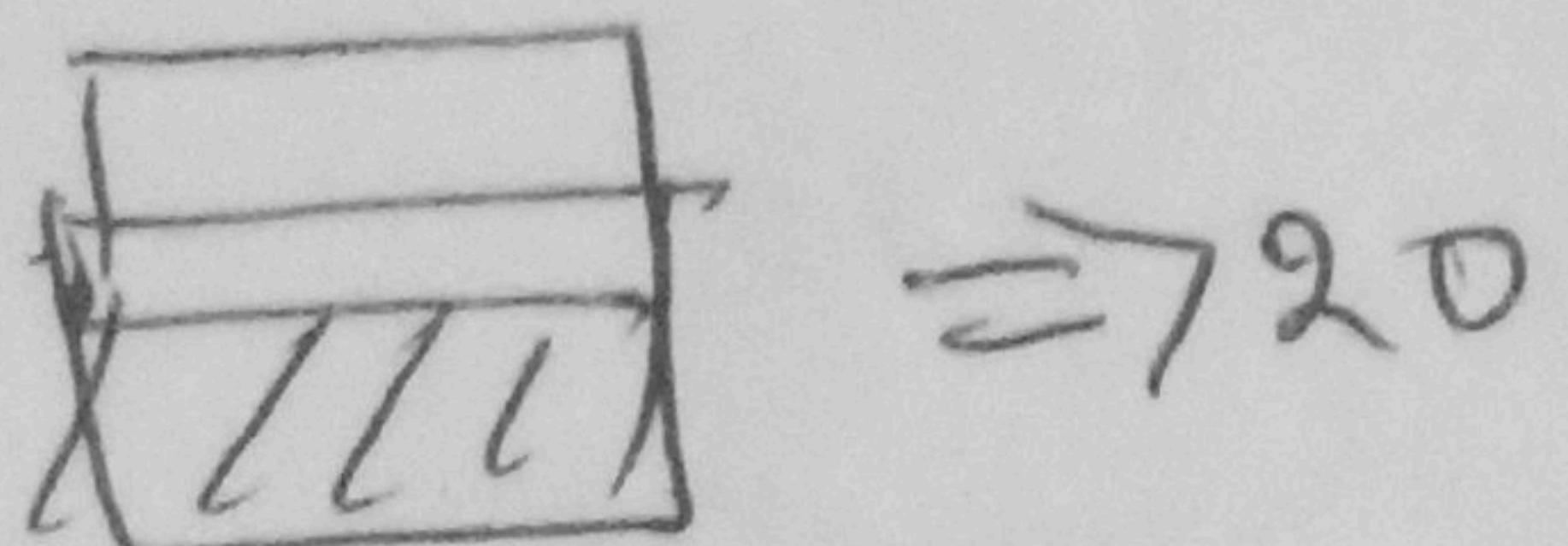
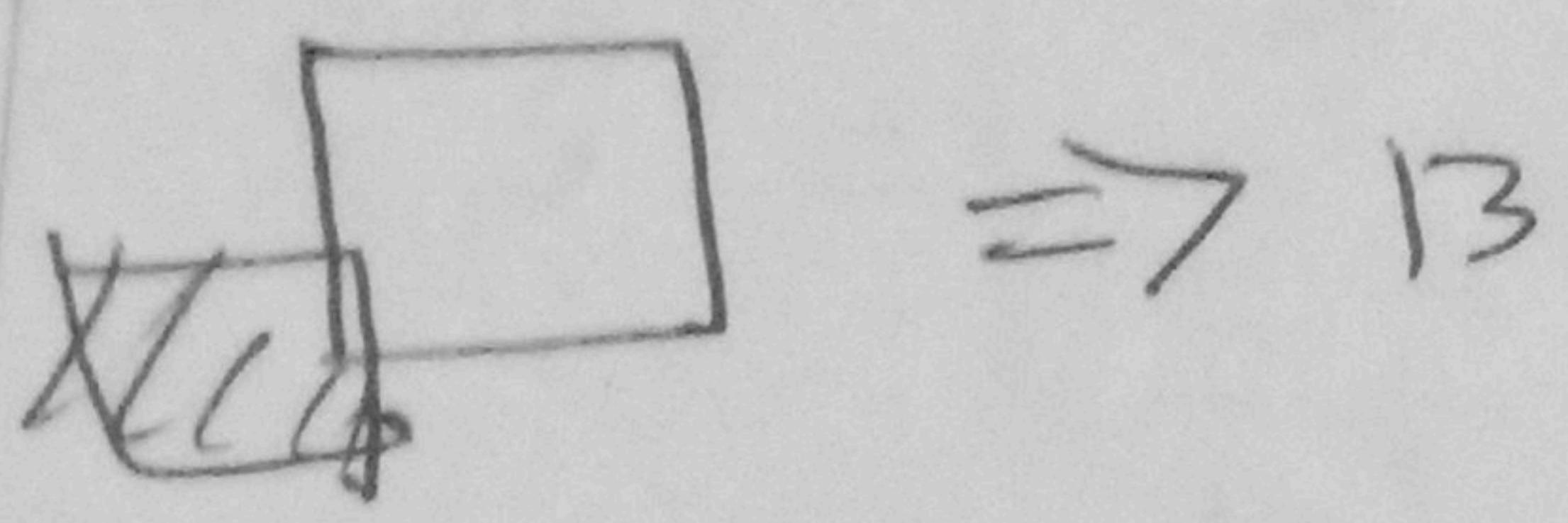
$\boxed{} \rightarrow \text{Image}(2D)$
(Applied padding)

$\Rightarrow -17$

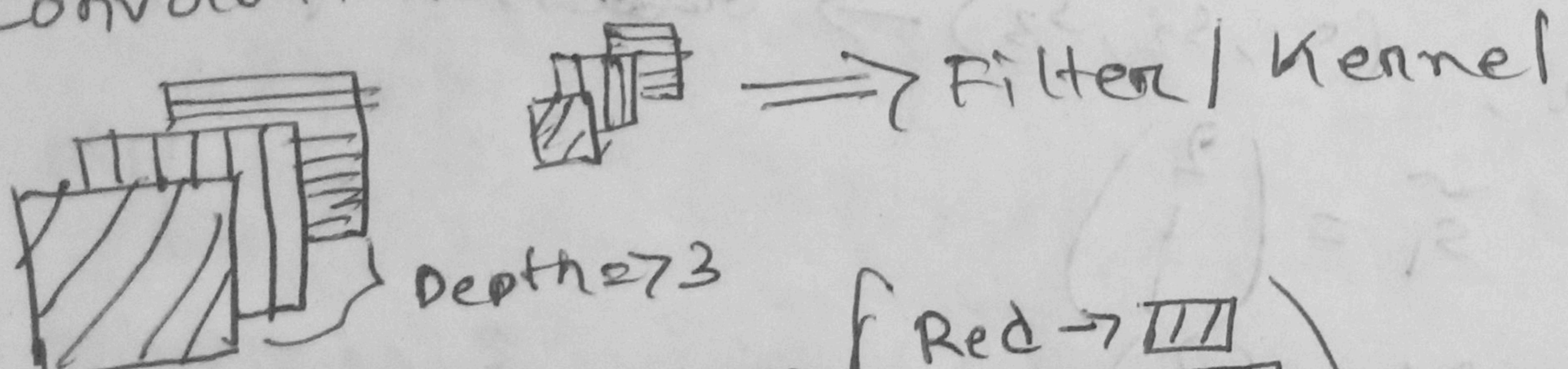
$\Rightarrow -18$


(Inside 2D image) $\Rightarrow -24$


 $\Rightarrow -18$



convolution in 3 dimension



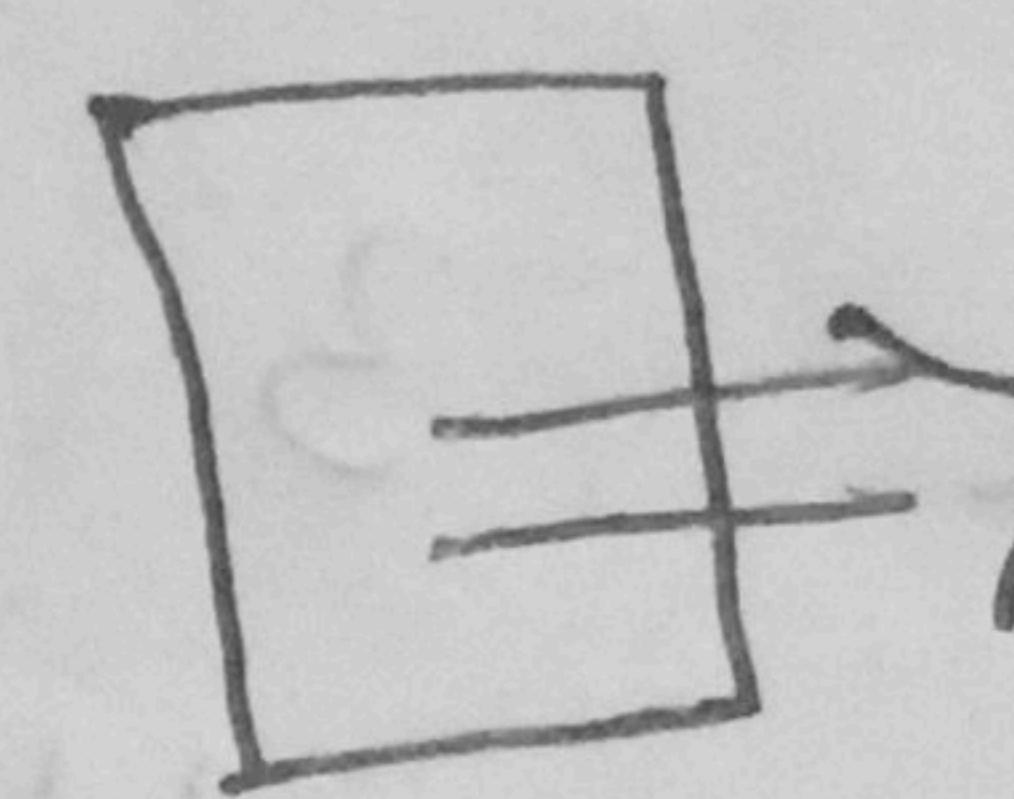
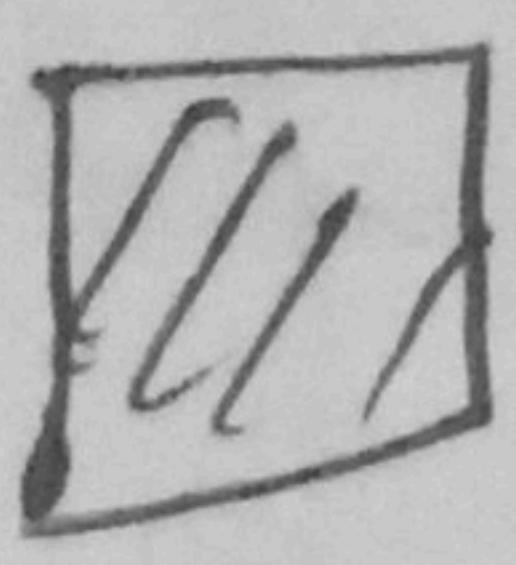
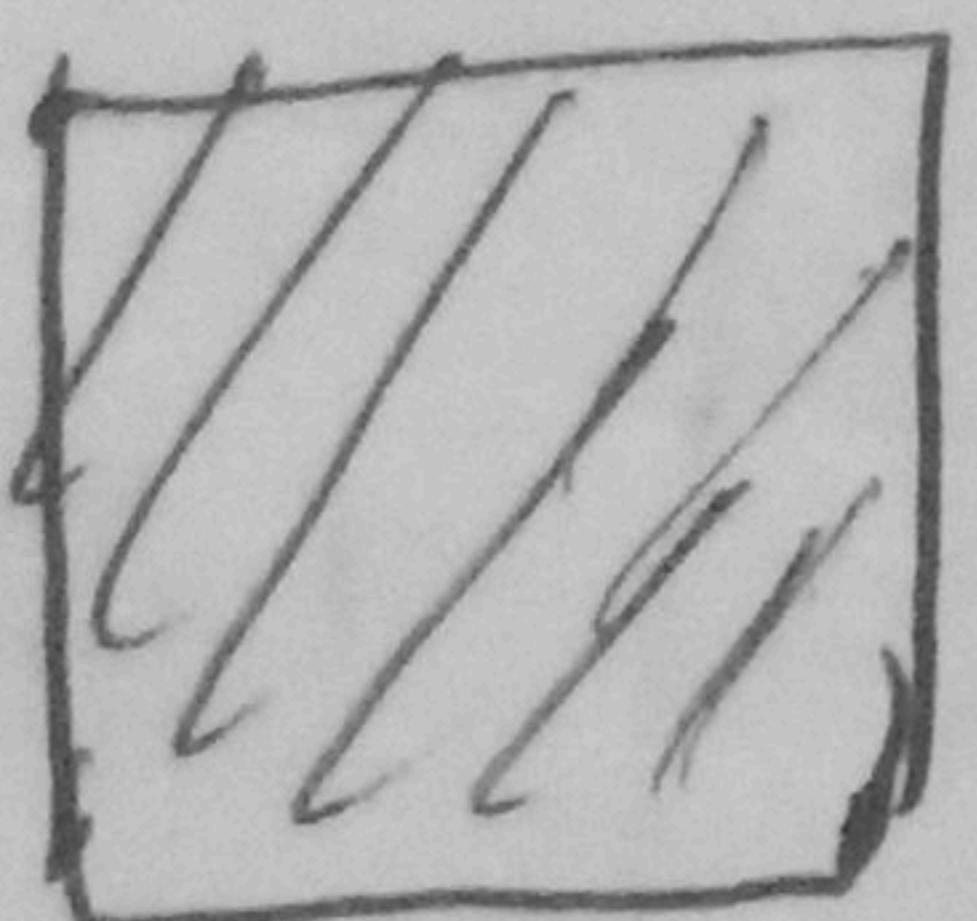
Depth = 3

RGB images

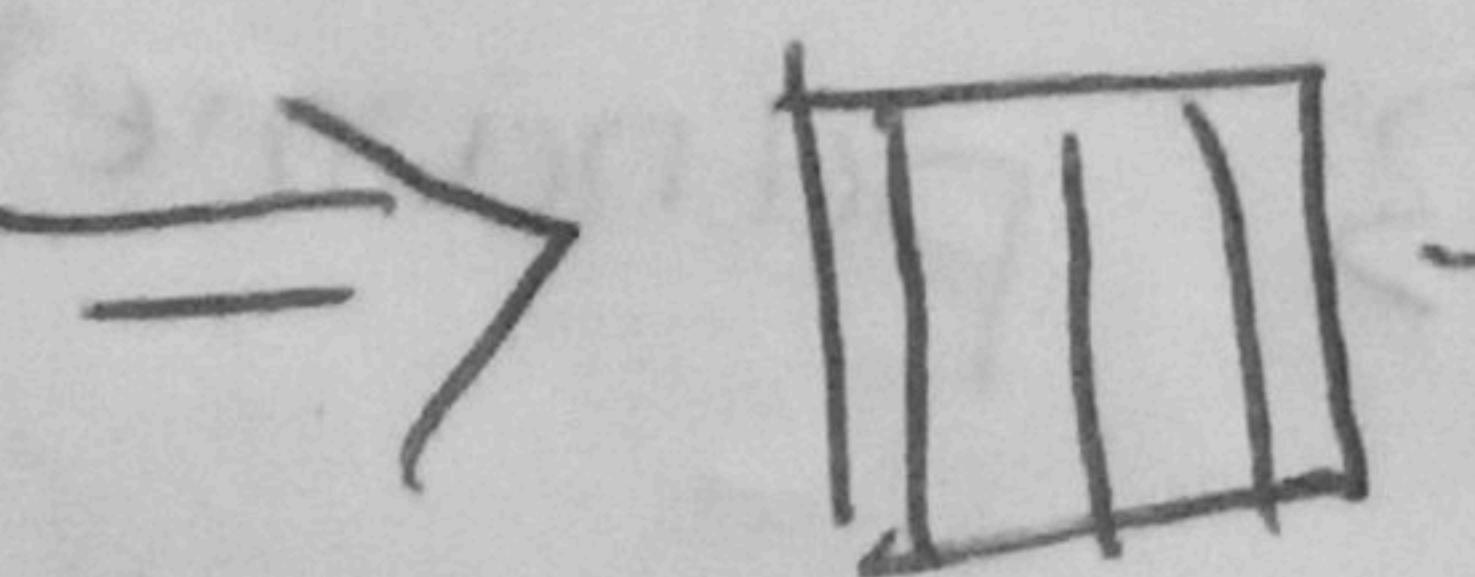
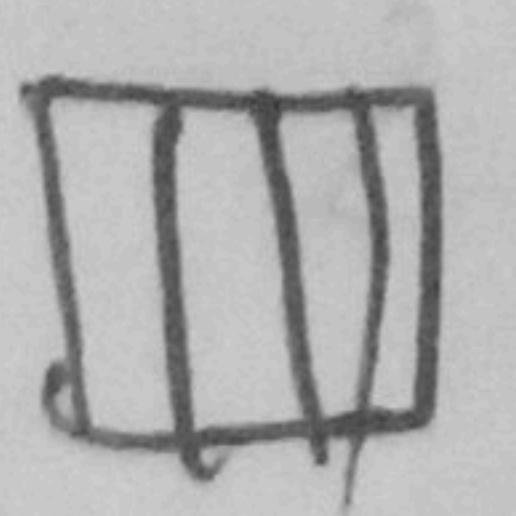
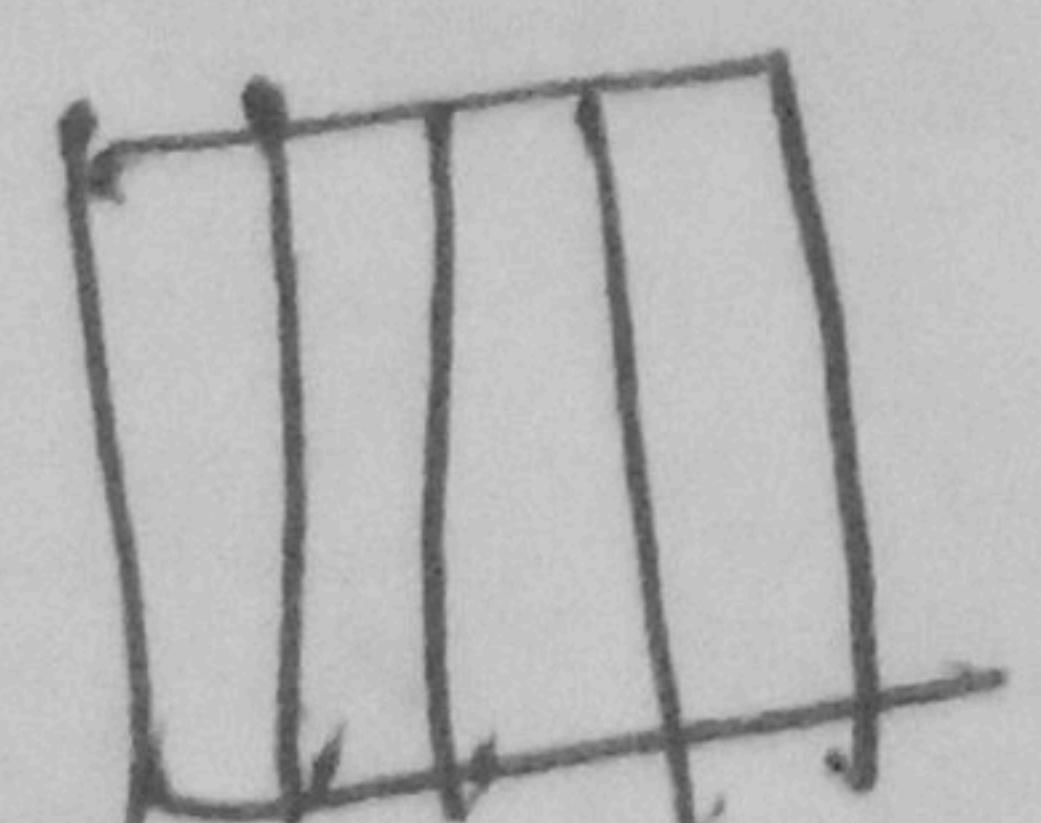
{ Red \rightarrow Green \rightarrow Blue \rightarrow

[Their Depth must be same]

Intermediate Output



Final output
(summation of these output)



Feature map

