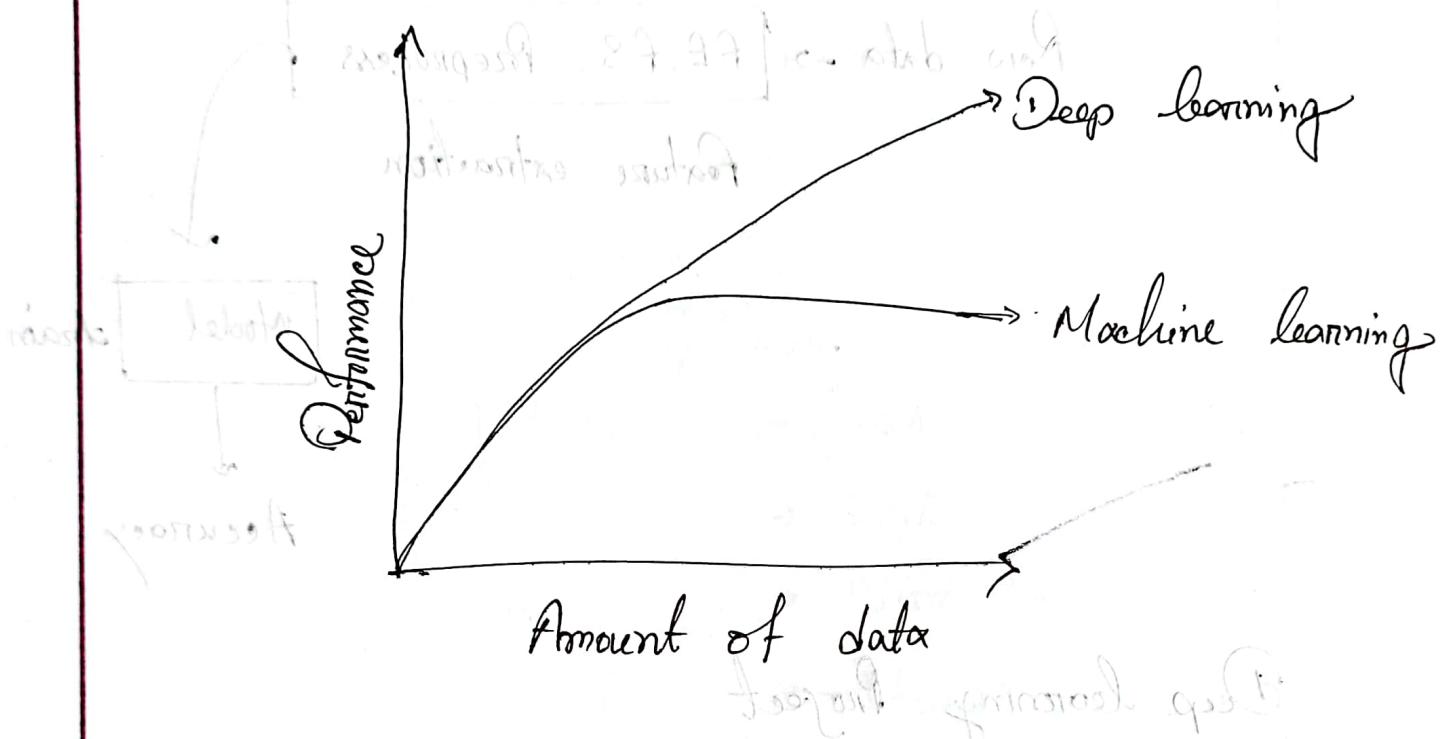


why Deep learning



Exponential growth of data

Deep learning model.

Accuracy ↑

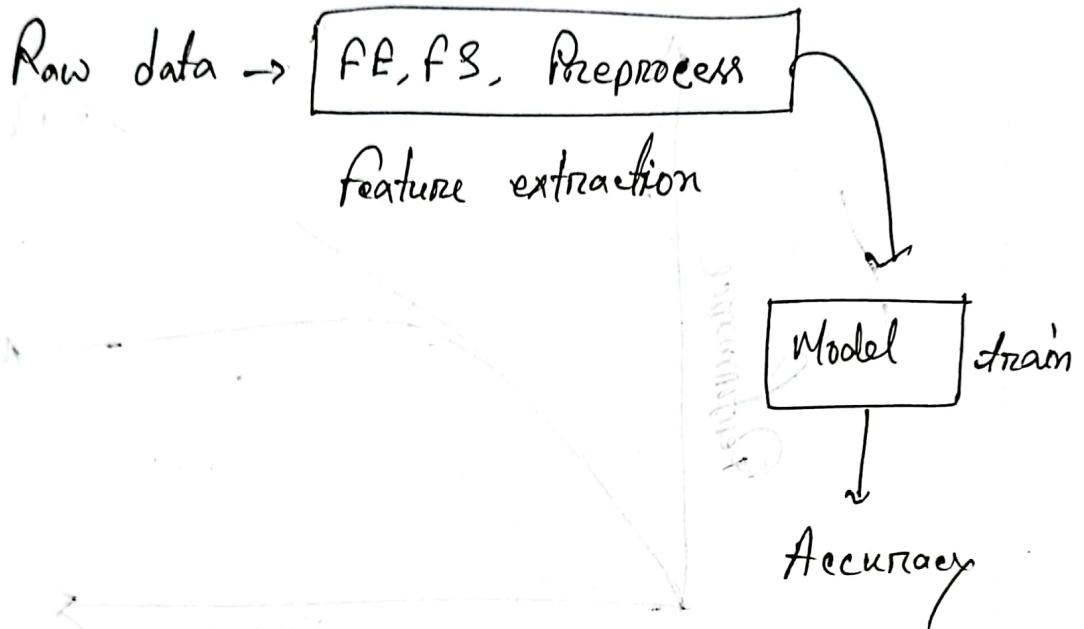
→ Technology Upgradation.

cheap hardware

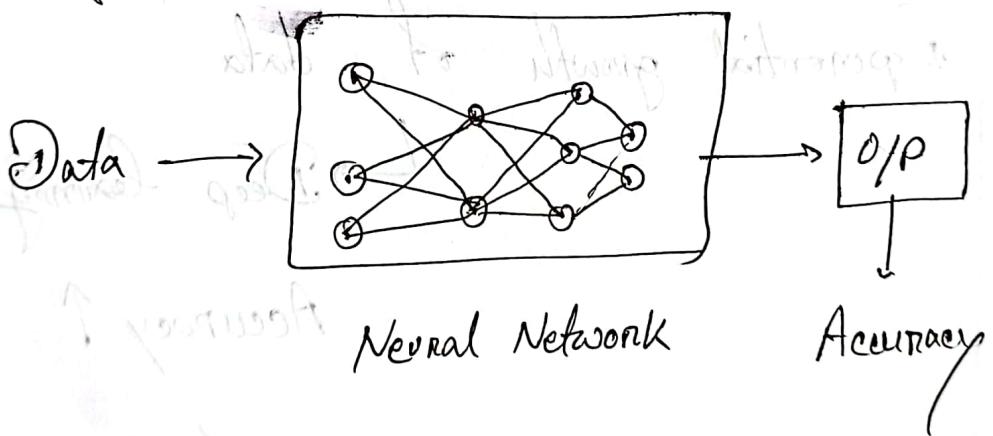
NVIDIA → GPU Price,

↓
used in cloud.

ML Project



Deep learning Project



→ Deep learning can solve complex problems
statement

- NLP
- Image classification
- Speech Recognition.

Deep learning Roadmap

① Introduction to Neural Network.

Learning Rule
Loss Function

Optimizer

Gradient descent

SGD

Mini-batch GD

Adagrad

RMSProp

Adam

Activation function

Sigmoid

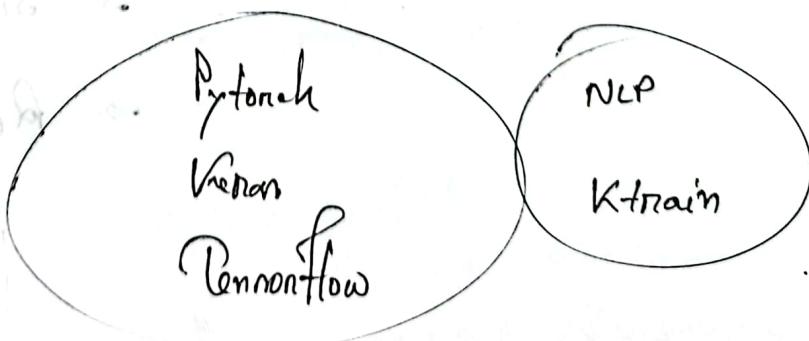
ReLU

- Global minimal
- local minimal
- Convergence point.

② Artificial Neural Network (An like ML)

↳ Deployment project.

- Heroku
- M Azure
- Flank
- weight initialize
- Parameters
- hyperparameter tuning
- Hidden layers
- Neuron



③ CNN - Convolutional Neural Network.

⇒ 1 in dm done here.

⇒ How does the convolution actually work.

⇒ Input: Images, video, sound.

⇒ Transfer learning, ResNet,

Alexnet, vgg16.

⇒ Object Detection.

⇒ R-CNN, Masked R-CNN, SSD,

YOLO

⇒ Deploy Project.

④ RNN → Recurrent Neural Network.

⇒ when our input data Sequence to sequence

⇒ LSTM , Bidirectional LSTM

⇒ GRU

⇒ Word Embedding

⇒ Word2Vec

⇒ Attention Model, Encoder, decoder.

Autoencoder.

⇒ Transfer. (NLP)

⇒ BERT (NLP)

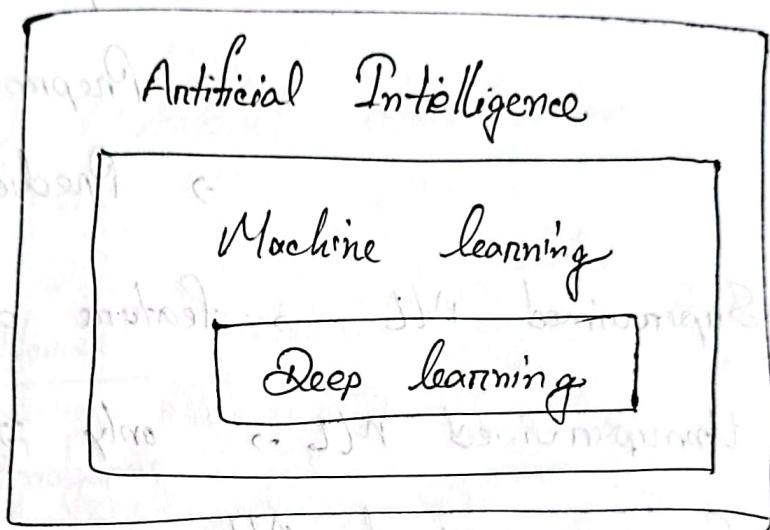
⇒ Hugging face, Ktrain (library)

AI vs ML vs DL vs DS ?

Computer Vision, NLP

AI Application

- Netflix
- Self driving
- Amazon
- YouTube.



→ Netflix Recommendation is a part of ML, DL.

→ In AI Application not everything is AI.

Paying to creating an application, and that application perform any human intervention.

Machine learning ⊂ Subset of AI

- Use Stat tool to analyze the data
 - explore the data
 - Preprocess the data.
 - Prediction, forecasting
- Supervised ML → feature and label
- Unsupervised ML → only feature. Dependant feature.
- Semisupervised ML.
- Reinforcement learning
- Supervised ML → dependent AI
 - Regression: output will be continuous.
 - Classification: classify or category something. generate value.

Classification.

① Binary classification

② Multiclass classification

③ Multilabel classification

Movies Action Romance Thriller Suspense

Endgame

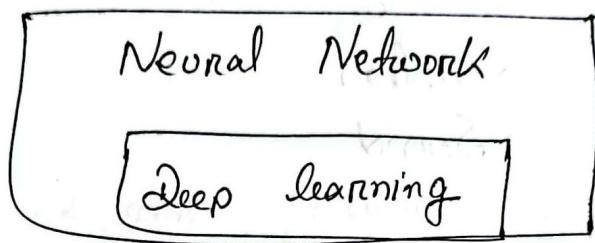
Drishyam

Regression Problem.

Time series analysis.

- Unsupervised ML
 - Clustering
 - Segmentation
 - Reduce dimension.
- Semisupervised ML
 - initially → labeled data - Model (train) → Rec.
 - Based on the user interaction one
 - the user data → train model
 - Continuously → New Recommendation.
 - Reinforcement learning
 - New born baby → Young man.

Deep learning is Subnet of machine learning.



- How our brain works?
- Main aim of Deep learning is Minimize the human brain.
- ANN → Artificial neural network.

Input : Tabular kind of dataset

- All the task done by NL can be done by ANN.

CNN → Convolutional Neural Network

Input : → Image (matrix format)
→ Video
→ Sound

Application :

1. Image classification & Recognize
2. Object detection.
3. Object segmentation.
4. Tracking

Recognition → detection → Segmentation.

- Image classification technique. — CNN
- Parameter learning learning technique.
CNN (Image Detection)

Another branch of AI is computer vision.

ii. Object Detection

RCNN, FASTER RCNN, FASTER RCNN

SSD, YOLO, Detectron.

iii. Segment, Detect, Recognize and classify.

iv. Tracking : Capture the object and try to track where it's moving.

⇒ RNN : Recurrent Neural Network.

Input data : Text data

Time Series

Sequential data

→ RNN

→ Transformer

→ Bidirectional RNN

→ Encoder, Decoder.

→ LSTM

→ B-LSTM

→ GAN

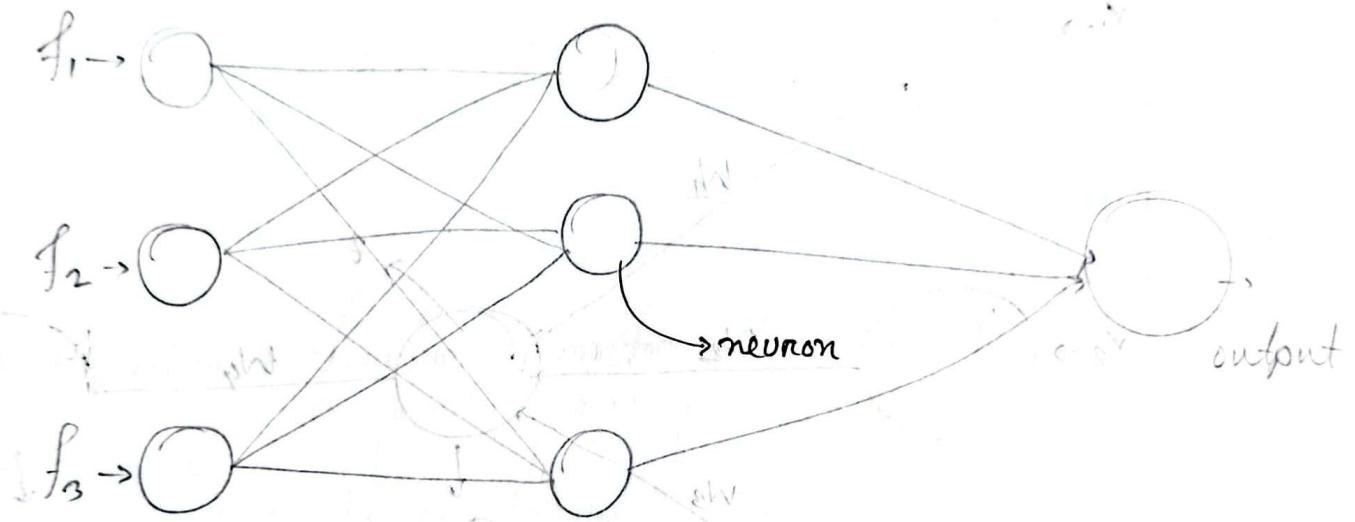
→ BERT

→ GRPT1, GRPT2, GRPT3

Introduction to Neural Network.

- Neural Network : Combination of many neuron build a network mean neural Network.
- ⇒ Perception : Single layer neural Network.
- ⇒ Multi-layered Neural Network
- Perception }
- ⇒ Not able to learn very properly
- ⇒ In 1986 Geoffrey hinton invented "back-propagation", Able to learn efficiently.
- ⇒ Not invented actually, what he claimed is that he was the person to clearly demonstrate that backpropagation could learn interesting

internal representations and that this is what made it popular.

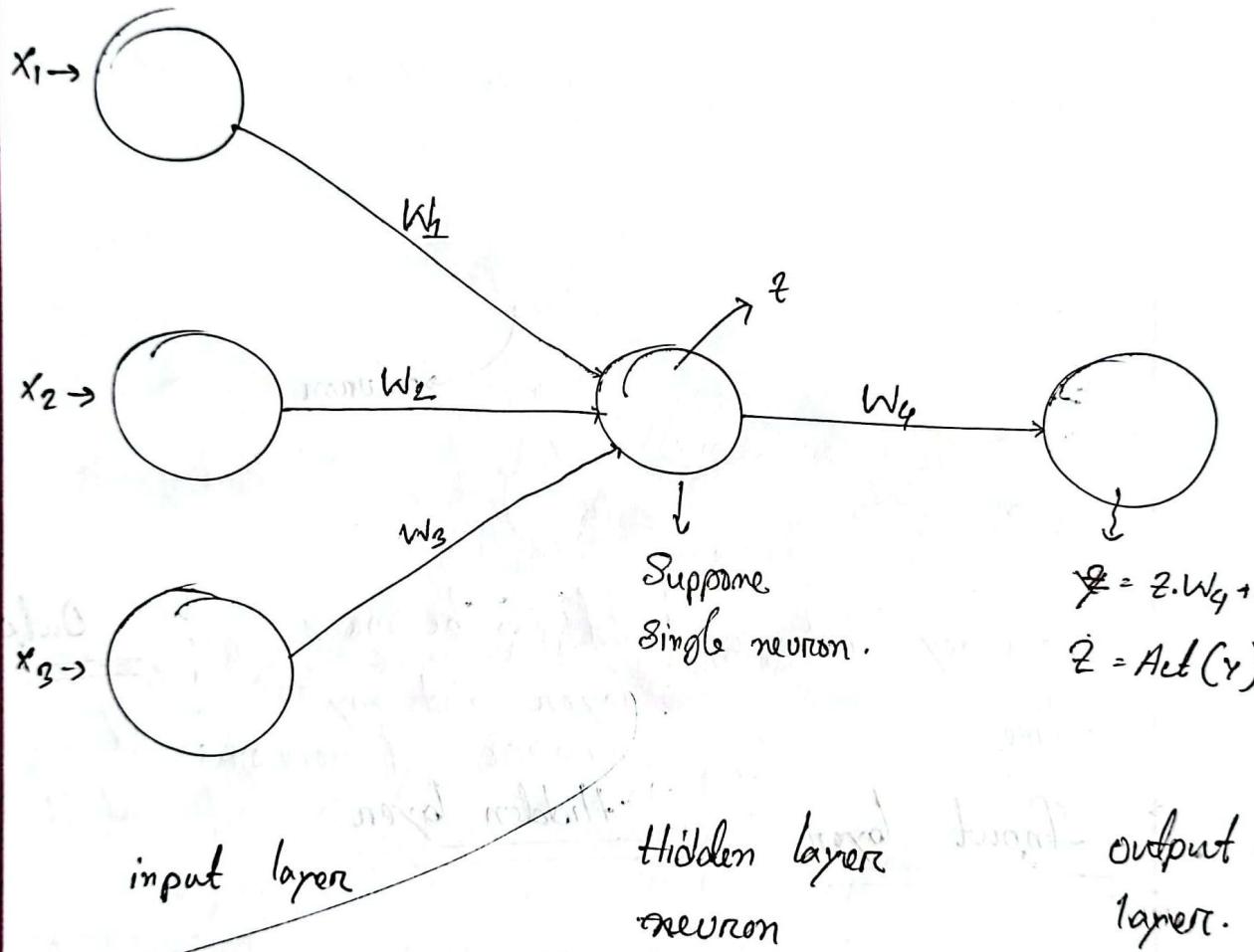


It can have any number of feature. It can be many layers and any number of neurons. It may be my Input layer, Hidden layer, Output layer.

→ Hidden layer can have any number of neurons.

Basic NN Architecture

How Neural Network works.



Process

Step 1 : $\sum_{i=1}^n w_i x_i$

Step 2 : Activation function $(\sum_{i=1}^n w_i x_i)$

$$y = w_1x_1 + w_2x_2 + w_3x_3 + \text{bias}.$$

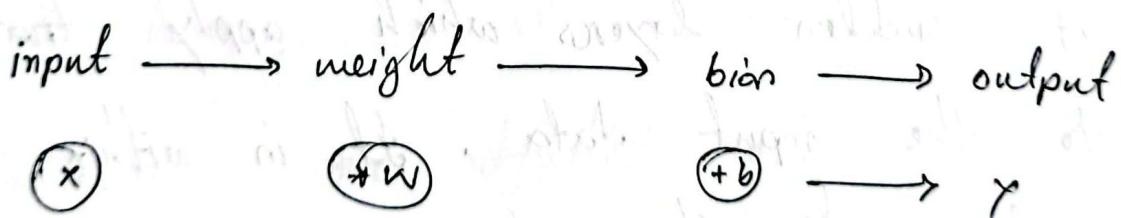
$$g = \text{Act}(y)$$

→ Sigmoid function like Activation function.

$$\frac{1}{1+e^{-y}}$$

→ Activation function transform the result between 0 to 1.

Question: why we need or what are the importance of weight and Activation function.



→ Weight help to determine the importance of any given variable.

- Determine the importance of any given variable
 - ↳ every input.
- longer weights, greater importance to the decision or outcome.
- DeepAI.org.
 - Within a neural network there is an input layer, that takes the input signal and passes them to the next layer.
 - Next, the neural network contains a series of hidden layers which apply transformation to the input data. If it is within the nodes of the hidden layer the weight are applied.

→ for example, a neural node may take the input data and multiply it by an assigned weight value, then add bias before passing the data to the next layer. The final layer of the neural network also known as the output layer.

- Weight and bias both learnable parameters.
- Initially randomly or any way assign weight and bias.
- An training neural network continues, both parameters are adjusted towards the desire value and the correct output
- "Bias" represent how far off the prediction are from the intended value.

- The weight is connection between two neurons.
- The strength of the connection between two neurons.
- Weight decides how much influence the input will have on the output.
- Weight affects the amount of influence a change in the input will have upon the output.
- A low weight value will have no change on the input, and alternatively a larger weight value will more significantly change the output.

→ example : If we connect a neuron to another with a weight, the higher the weight the higher the activation is of the second neuron.

Activation Function

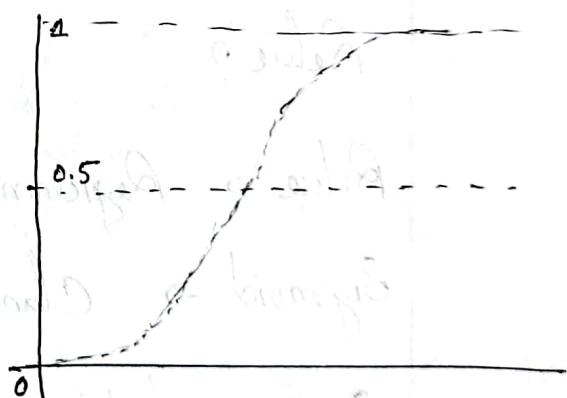
① Sigmoid AF

$$\frac{1}{1 + e^{-y}}$$

$$y = \sum_{i=1}^n w_i x_i + b_{bias}$$

Act (y)

Transform the result between 0 to 1.



Used in Logistic Regression.

↓
Decide your neuron Active or not.

Active means transfer the signal to the next layer.

② ReLU Activation Function

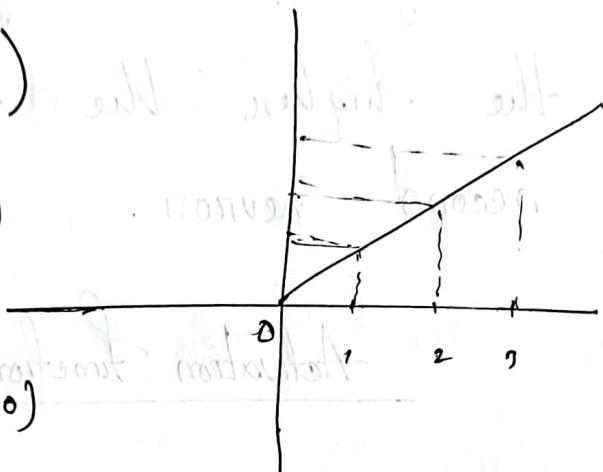
$$\max(y, 0)$$

$$-ve \rightarrow \max(y, 0)$$

$$= 0$$

$$+ve \rightarrow \max(y, 0)$$

$$\Rightarrow +ve$$



→ when we use sigmoid? when we use ReLU?

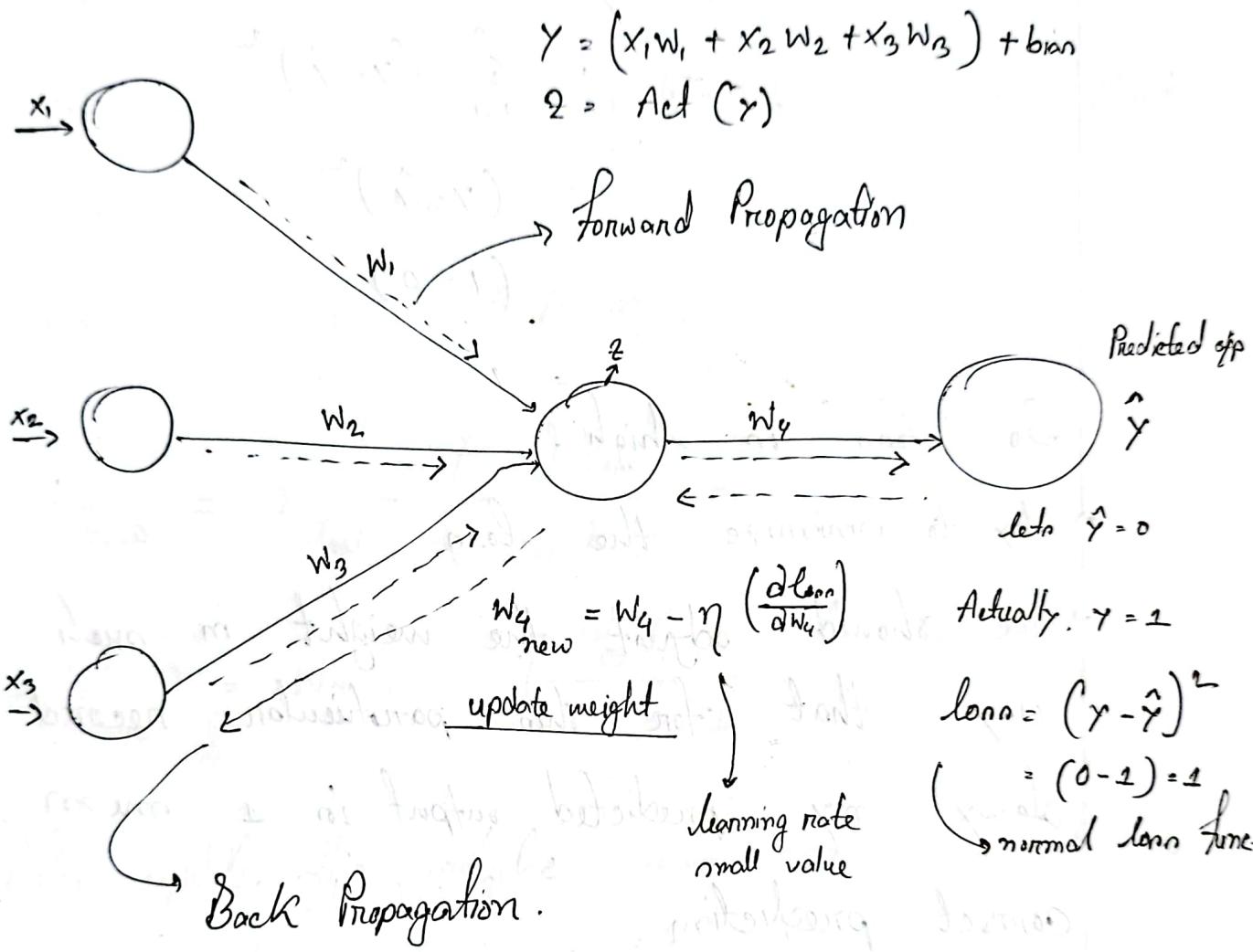
ReLU → Regression problem.

Sigmoid → Classification problem.

→ In hidden layer we can use ReLU.

→ But the final output layer, should always use sigmoid.

Neural network training with backpropagation.



One Case:

Play	Study	Sleep	Op
2h	4h	8h	$\begin{cases} 1 & \rightarrow \text{Pass the exam} \\ 0 & \rightarrow \text{Fail the exam} \end{cases}$

loss function normally

$$\text{loss} = \sum_{i=1}^n (\gamma_i - \hat{\gamma}_i)^2$$

$$= (\gamma - \hat{\gamma})^2$$
$$= (1 - o)^2$$

So loss is high.

- try to minimize the loss.
- we should adjust the weight in such a way that for this particular record always my predicted output is a mean, correct prediction
- loss value decreasing (minimize the loss)
- Accuracy increasing

It can be done by using optimizers.

"Forward Propagation and Backpropagation"

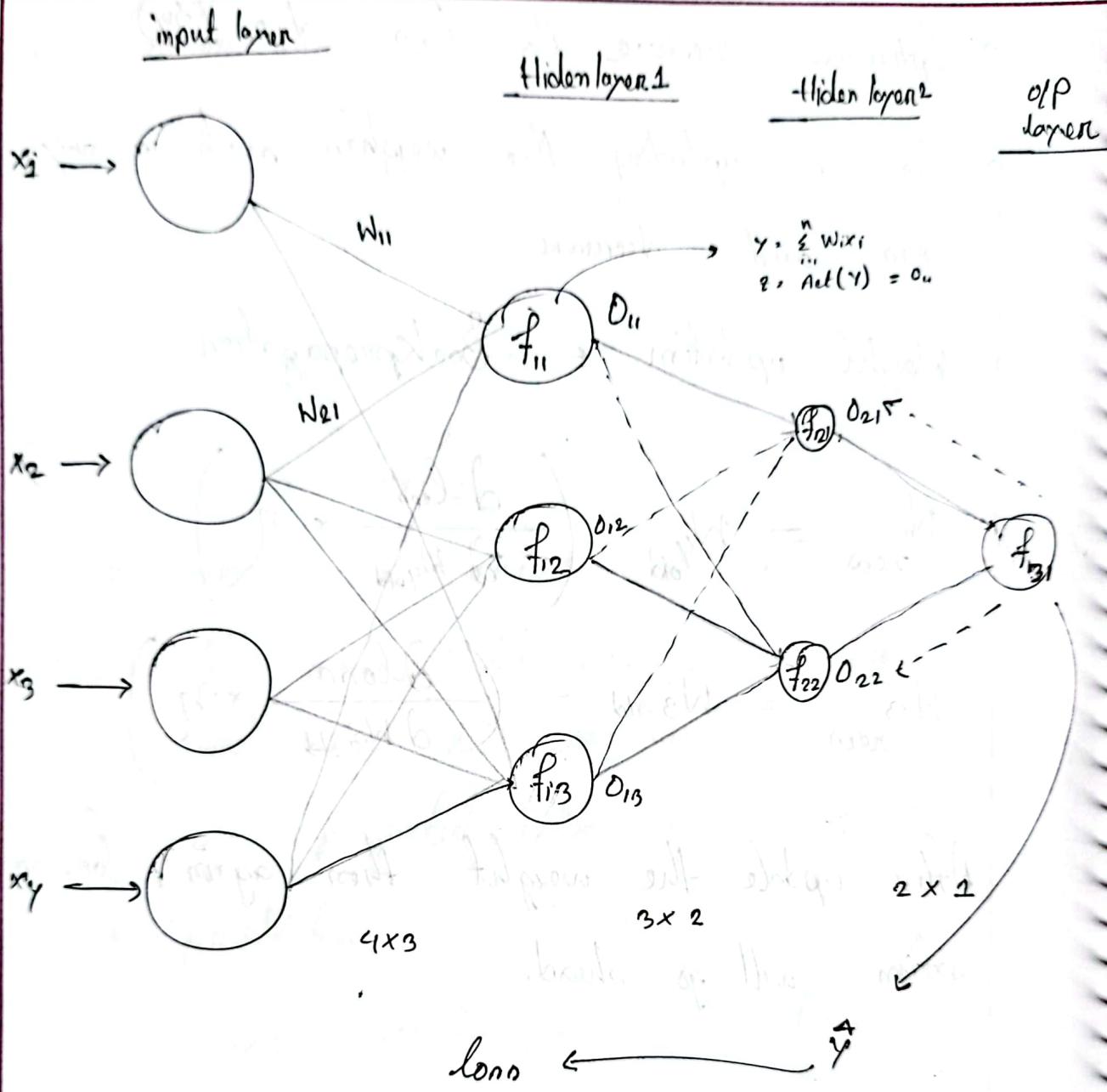
- Optimizer reduce the loss value (for forward pass)
- So we updating the weights such a way that loss will decrease.
- Weight updation \leftarrow Backpropagation.

$$w_4^{\text{new}} = w_4^{\text{old}} - \left(\frac{\partial \text{loss}}{\partial w_4^{\text{old}}} \times \eta \right)$$

$$w_3^{\text{new}} = w_3^{\text{old}} - \left(\frac{\partial \text{loss}}{\partial w_3^{\text{old}}} \times \eta \right)$$

After update the weight then again forward propagation will go ahead.

"Multi layer Neural Network"



$$\text{loss function: } (\gamma - \hat{\gamma})^2 \downarrow$$

- = If we want to reduce loss, we use **optimizer**.
- = Optimizer update the weight such a way that loss will decrease mean $\gamma, \hat{\gamma}$ almost same

- Each and every layer updating the weight on efficient way. How it is done?

"Backpropagation"

Optimizer: Gradient descent, -----

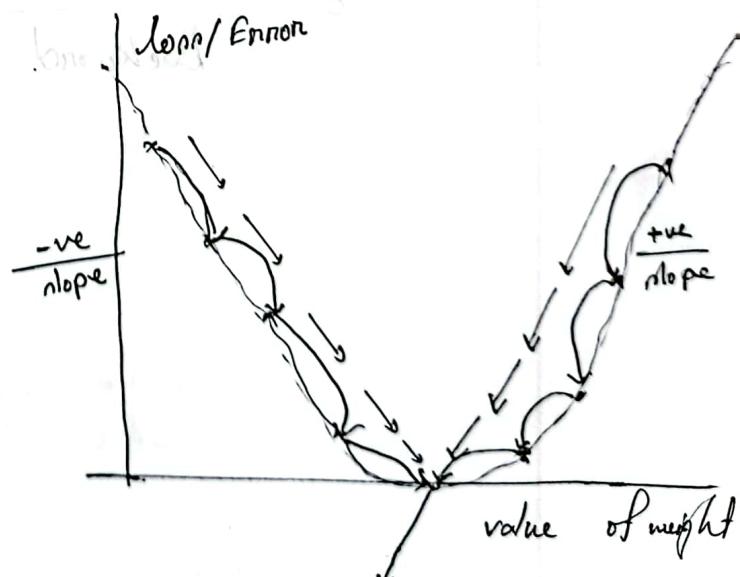
- first weight and biases initialize randomly and then update in Backpropagation.

- Weight update,

$$W_{\text{new}} = W_{\text{old}} - \left(\frac{\partial \text{Loss}}{\partial W_{\text{old}}} \times n \right)$$

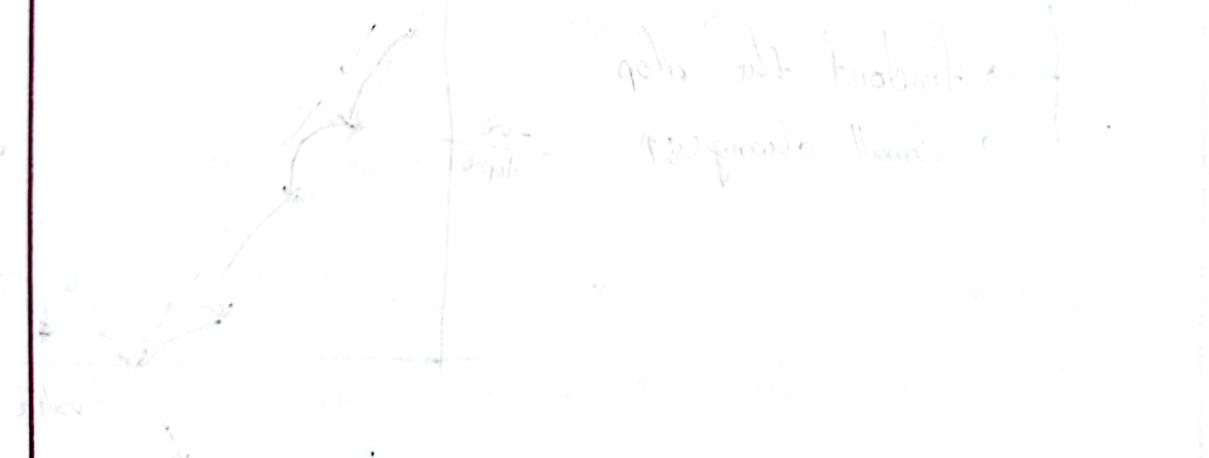
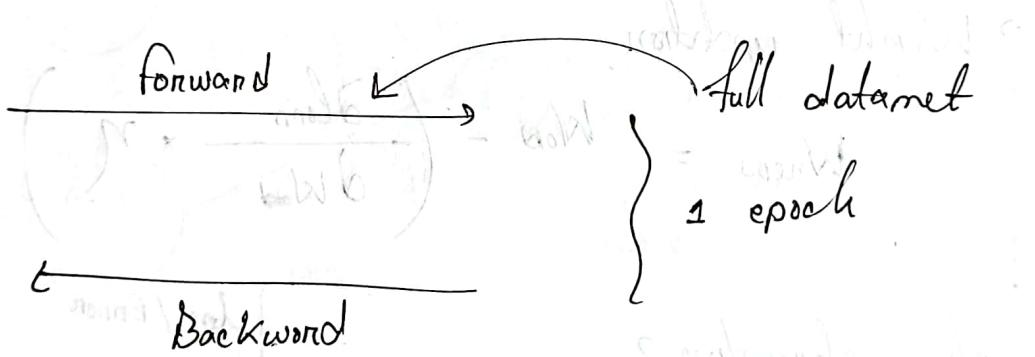
why derivative?

- Findout the slope
- Small changes ↑



Convergence point - - - Global minima

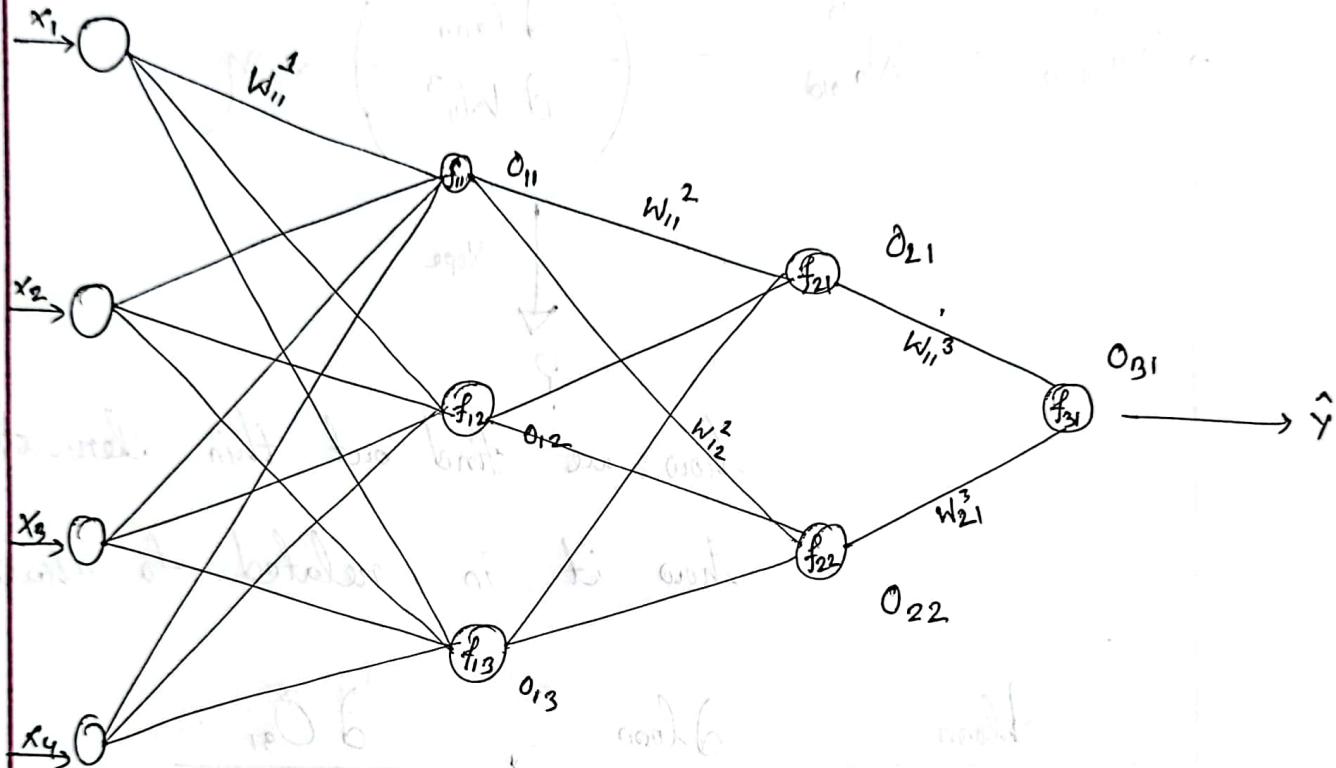
- $\eta \rightarrow$ learning rate is important.
- \rightarrow should not be very small or should not be very large.
- \rightarrow How we know that we are on the right way by selecting correct learning rate?
- \rightarrow After some epochs loss will be decreasing



Final loss \rightarrow final estimation

"Chain Rule of Differentiation with Backpropagation"

Input layer H_l1 H_l2 opp layer



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

Backpropagation : weight update process . How ?

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

↓
learning rate (how to choose it?)

learning rate choose by hyperparameter optimization.

$$W_{11}^3 \text{ new} = W_{11}^3 \text{ old} - \left(\frac{\partial \text{loss}}{\partial W_{11}^3} \right) \times \eta$$

↓
Slope
?

how we find out this derivative

how it is related to chain rule.

$$\frac{\partial \text{loss}}{\partial W_{11}^3} = \frac{\partial \text{loss}}{\partial O_{31}} * \frac{\partial O_{31}}{\partial W_{11}^3}$$

$$W_{21}^3 \text{ new} = W_{21}^3 \text{ old} - \left(\frac{\partial \text{loss}}{\partial W_{21}^3} \right) \times \eta$$

$$\frac{\partial \text{loss}}{\partial W_{21}^3} = \frac{\partial \text{loss}}{\partial O_{31}} * \frac{\partial O_{31}}{\partial W_{21}^3}$$

(i) some time back ground

$$W_{11}^2 = W_{11\text{old}}^2 - \left(\frac{\partial L_{\text{err}}}{\partial W_{11}^2} \right) \times n$$

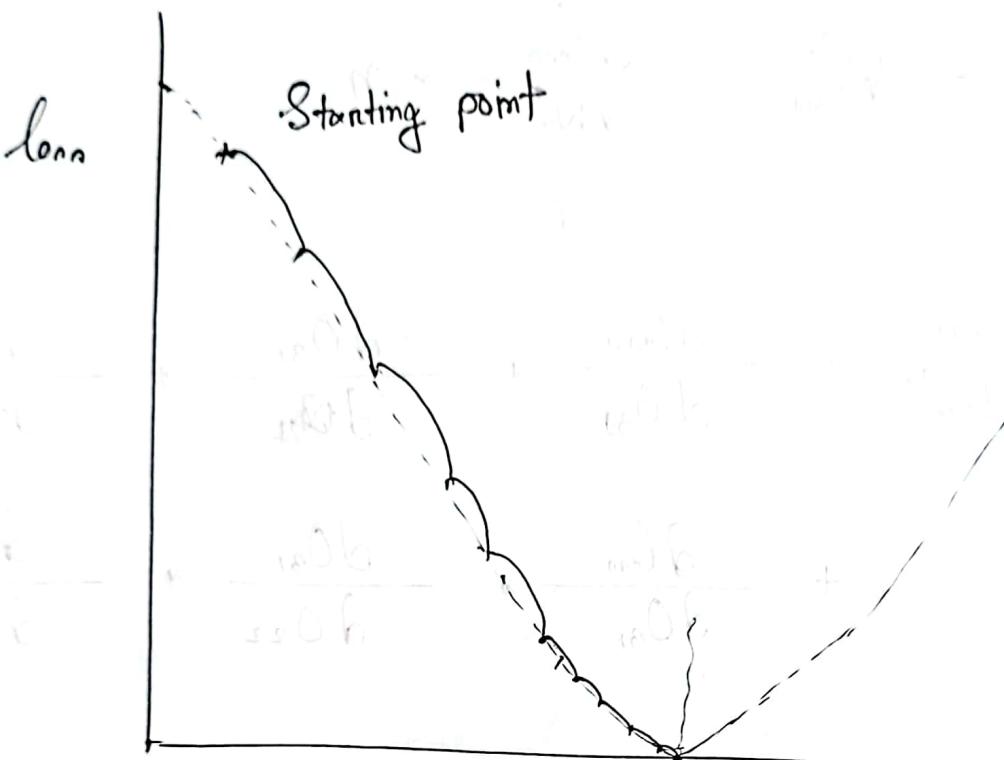
?

$$\begin{aligned} \frac{\partial L_{\text{err}}}{\partial W_{11}^2} &= \frac{\partial L_{\text{err}}}{\partial O_{31}} \times \frac{\partial O_{31}}{\partial O_{21}} \times \frac{\partial O_{21}}{\partial W_{11}^2} \\ &+ \frac{\partial L_{\text{err}}}{\partial O_{31}} \times \frac{\partial O_{31}}{\partial O_{22}} \times \frac{\partial O_{22}}{\partial W_{12}^2} \end{aligned}$$

→ Neural Network for better understanding.

→ 3 blue 1 brown

→ 1 Bn

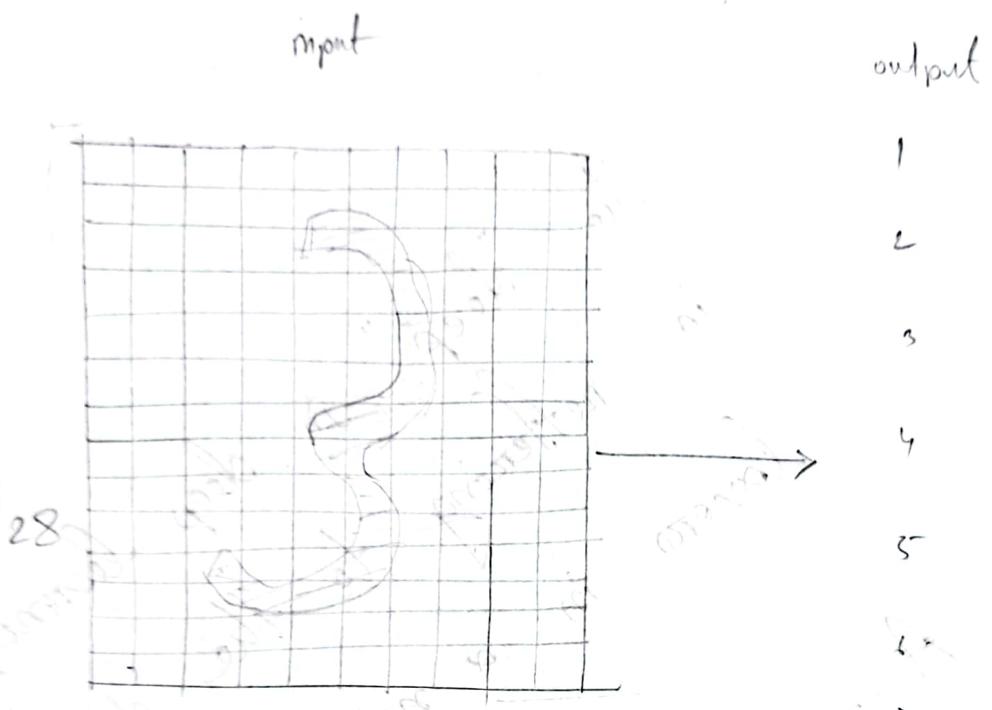


Point of convergence where the
 Cost function is at its minimum

The "deep" in deep learning
is referring to the depth of
layers in a neural network.

A neural network that consists of
more than three layers - which would
be inclusive of the "inputs",
and output can be considered
a deep learning algorithm.

① Digit Recognize



LSTM - for speech recognize
CNN → Image recognize

} Good Performance

$$28 \times 28 = 784 \text{ neuron}$$

flatten

$\Rightarrow 784$
neuron.

$$\begin{bmatrix} 0 \\ 0 \\ 0 \\ 0 \\ 0 \\ 0 \\ 0 \\ 0 \\ 0 \\ 0 \end{bmatrix}$$

input layer				hidden layer	hidden layer	hidden layer	output layer
				0	0	0	0
<u>28x28</u>	0	0	0				0 1
<u>284</u> neuron	0	0	0				0 2
	0	0	0				0 3
	0	0	0				0 4
	0	0	0				0 5
	0	0	0				0 6
	0	0	0				0 7
	0	0	0				0 8
	0	0	0				

\Rightarrow Suppose 2 hidden layer

Hidden layer 1 = 16 neuron

Hidden layer 2 = 10 neuron

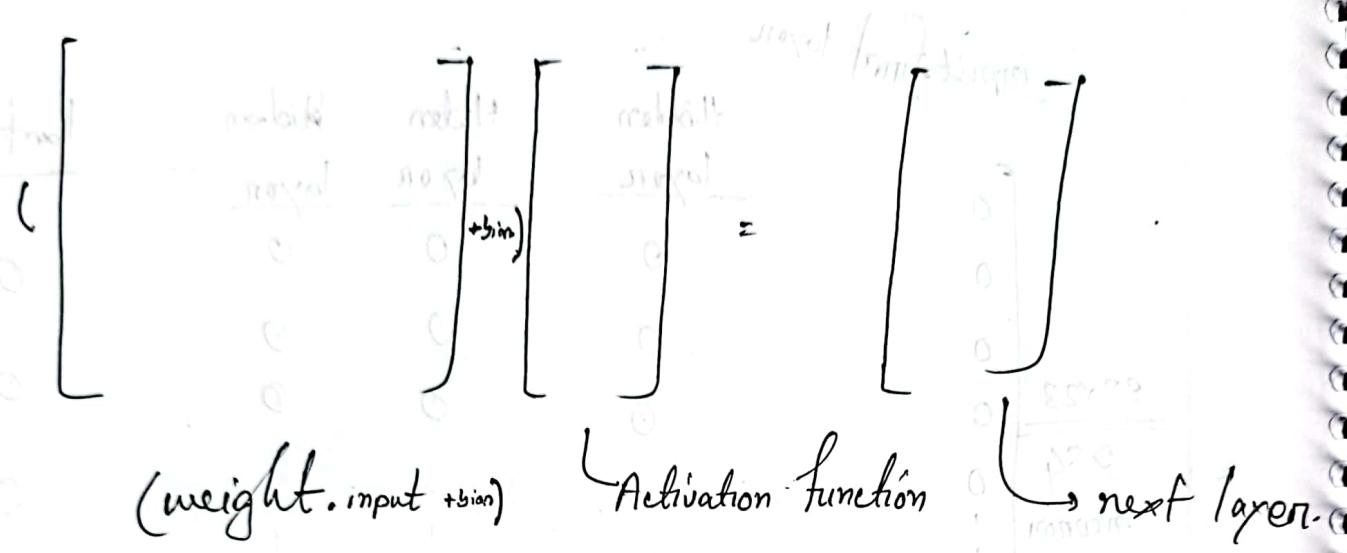
$$\underline{284 \times 16 + 16}$$

$$\underline{16 \times 16 + 16}$$

$$\underline{16 \times 10 + 10}$$

weight and biases

3 blue 4 Brown



$$\text{Act} \left(\sum_{i=1}^n w_i x_i + \text{bias} \right) = z$$

Matrix - function (Composition) \Rightarrow linear Algebra.

sigmoid function is composed of

cosine at $\pi/2$ sigmoid which

sigmoid at $\pi/2$ sigmoid addition

$$w_1 x_1 + b_1$$

$$w_2 x_2 + b_2$$

$$w_3 x_3 + b_3$$

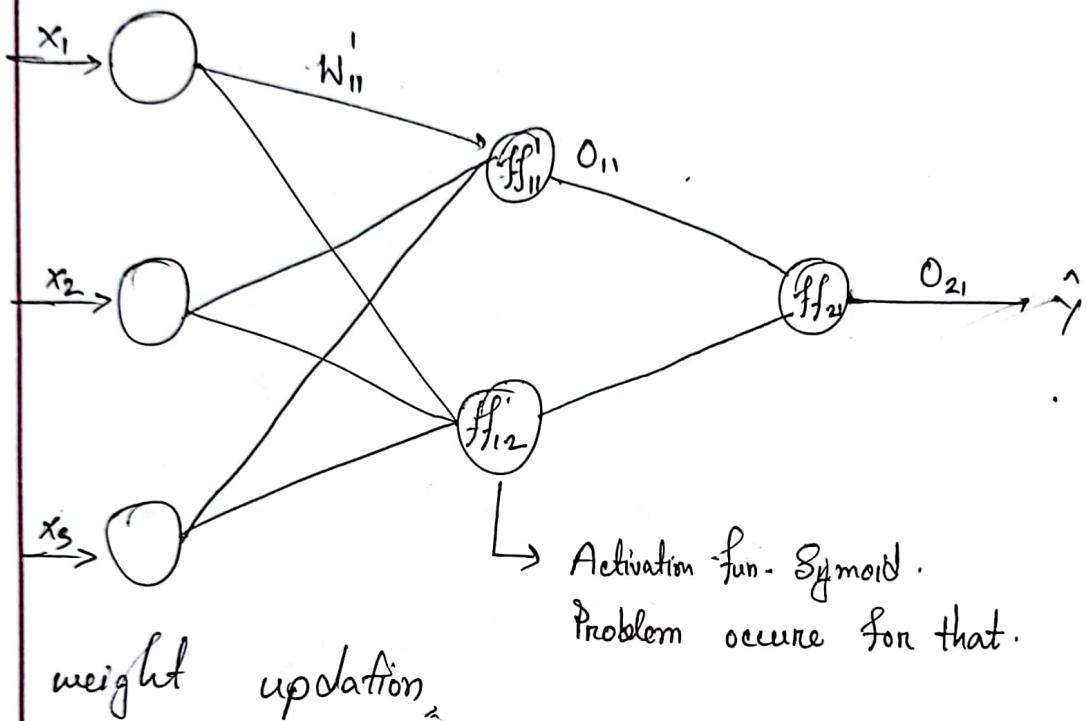
addition linear algebra

Gradient Descent \rightarrow BNL Article.

Atlanta - 1985 - 1986 - 1987

Waldschmidtia (Lindb.) Lindb.

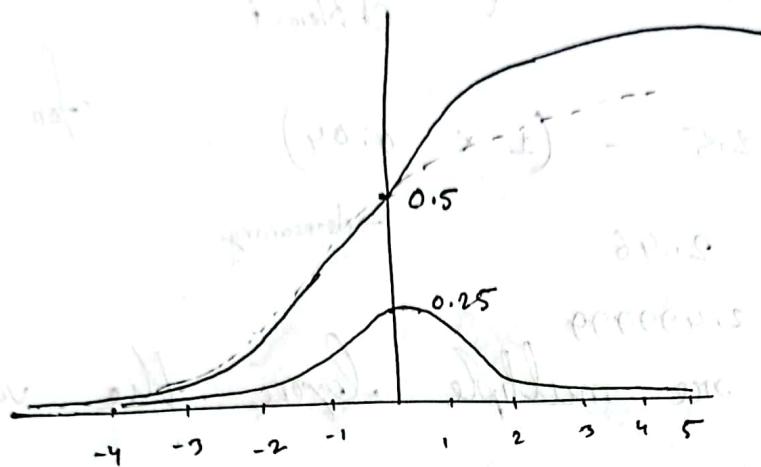
Vanishing Gradient Problem.



$$w_{11}'_{\text{new}} = w_{11}'_{\text{old}} - \eta \frac{\partial L_{\text{out}}}{\partial w_{11}'}$$

$$\frac{\partial L_{\text{out}}}{\partial w_{11}'} = \frac{\partial O_{21}}{\partial O_{11}} \cdot \frac{\partial O_{11}}{\partial w_{11}'}$$

Sigmoid function convert output 0 - 1
 derivative of sigmoid : 0 - 0.25



$$z = \sum w_i x_i + b_{bias}$$

$$\text{Sigmoid } f = \frac{1}{1+e^{-z}}$$

$$\frac{d\sigma}{dz} = 0 - 0.25$$

$$0 \leq \sigma(z) \leq 0.25$$

$$\frac{dl}{dw_{ii}} = \frac{\partial O_{2i}}{\partial O_{ii}} \cdot \frac{\partial O_{ii}}{\partial w_{ii}} \quad \text{[Chain Rule]}$$

$$(0 - 0.25) \quad (0 - 0.25)$$

$$0.2 \quad 0.02 = 0.04$$

$$\begin{aligned} \text{new weight} &= 2.5 - 1 \times 0.04 \\ &= 2.46 \end{aligned}$$

$$W_{i,new} = W_{i,old} - \eta \frac{\delta_{loss}}{\delta W_i}$$

$$= 2.5 - (1 \times 0.04)$$

$$= 2.46 \quad \text{↳ decreasing}$$

If there are multiple layers the value of derivative is decreasing.

So once a time old weight and new weight have very small change (may be unchanged).

Approximately old and new weight will be matching.

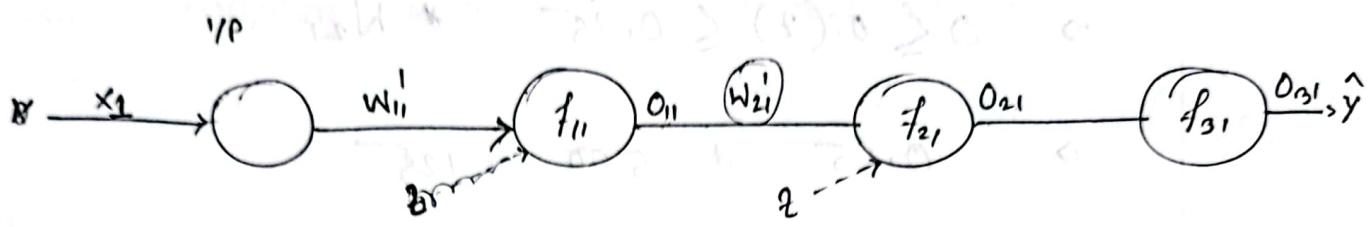
That is the reason the Sigmoid in bad

most popular activation function
is ReLU. And "ReLU" is not used in
any layer. But "ReLU" is used in "each"
and every layer. ~~layer~~
"ReLU" is used in every layer.
~~ReLU~~ Activation function solve this.

As a result the gradient descent never converge to the optimum. This is known as vanishing gradient problem.

- Sigmoid function. $\text{weight}_{\text{old}} \approx \text{weight}_{\text{new}}$
- tanh also. almost same.
- It occurs when gradient is too small
[derivative]
$$\frac{\partial \text{Loss}}{\partial \text{Weight}}$$

Exploding gradient Problem.



$$W'_{11, \text{new}} = W'_{11, \text{old}} - \frac{\partial \text{loss}}{\partial W'_{11}} * \eta$$

$$\frac{\partial \text{loss}}{\partial W'_{11}} = \frac{\partial \text{loss}}{\partial O_{31}} * \frac{\partial O_{31}}{\partial O_{21}} * \frac{\partial O_{21}}{\partial O_{11}} * \frac{\partial O_{11}}{\partial W'_{11}}$$

Main Reason,

Sigmoid and "Weights"

$$O_{21} = \phi(z) \rightarrow \text{Activation function } \phi$$

$$z = W_{21} \theta_{11} + b_2$$

$$\frac{\partial O_{21}}{\partial O_{11}} = \frac{\partial \phi(z)}{\partial z} * \frac{\partial z}{\partial O_{11}}$$

$0 \leq \phi(z) < 0.25$ $\frac{\partial \phi(z)}{\partial z}$

$$= 0 \leq \phi(z) < 0.25 * W_{21}$$

$$\Rightarrow 0 \leq \rho(z) \leq 0.25 \quad \# W_{21}$$

$$\Rightarrow 0.25 \times 500 = 125$$

$$\frac{d\text{loss}}{dW_{11}} = \frac{\partial O_{31}}{\partial O_{21}} \cdot \frac{\partial O_{21}}{\partial O_{11}} \cdot \frac{\partial O_{11}}{\partial W_{11}}$$

200 × 125 × 100

[large value.]

$$\text{Weight new} = W_{11}^{\text{old}} - \eta \frac{d\text{loss}}{dW_{11}}$$

[- large value]

New weight and old weight difference will be very high.

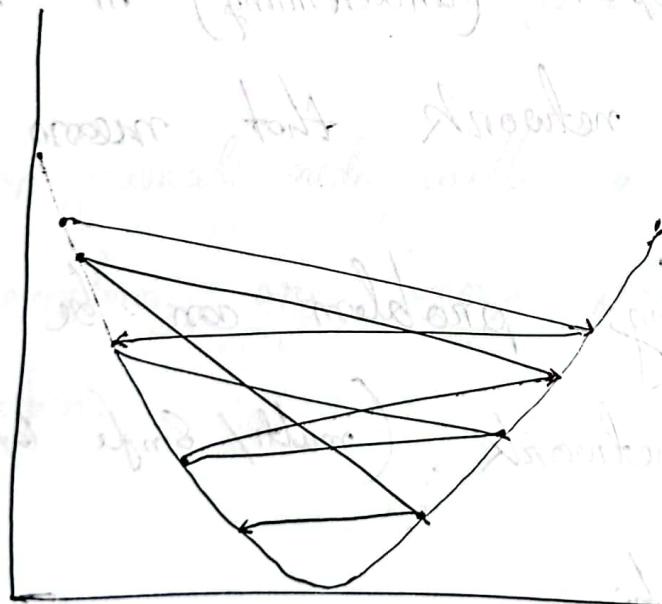
Gradient descent will never converge.
it will jumping here and there.

This happens when the gradient is too large.

gradient from layer to layer derivative is large

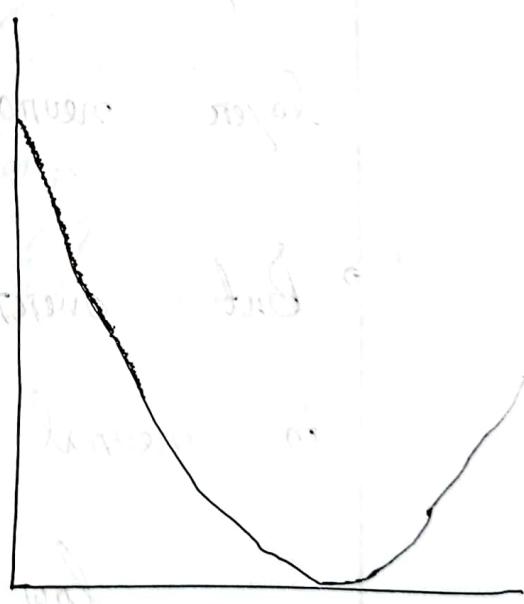
fitting at about min \rightarrow weight difference is high.
good news: multiplying with a small number

- The challenges of vanishing/exploding gradients in deep neural network - Analytics Vidhya.



convex

exploding gradient problem



vanishing gradient
problem.

but for non-convex?

Dropout layers.

- Huge data.
huge amount of weight and biases
Parameters, nn tends to overfit
to the particular use-case.
- We will never face underfitting problem
in multilayer neural network.
- May be happens (underfitting) in single
layer neural network that means perceptron.
- But overfitting problem can be happens
in neural network (multi/ single layer).

low bias

high variance.

→ large neural nets trained on relatively small datasets can overfit the training data.

Approach to reduce overfitting.

① Regularization.

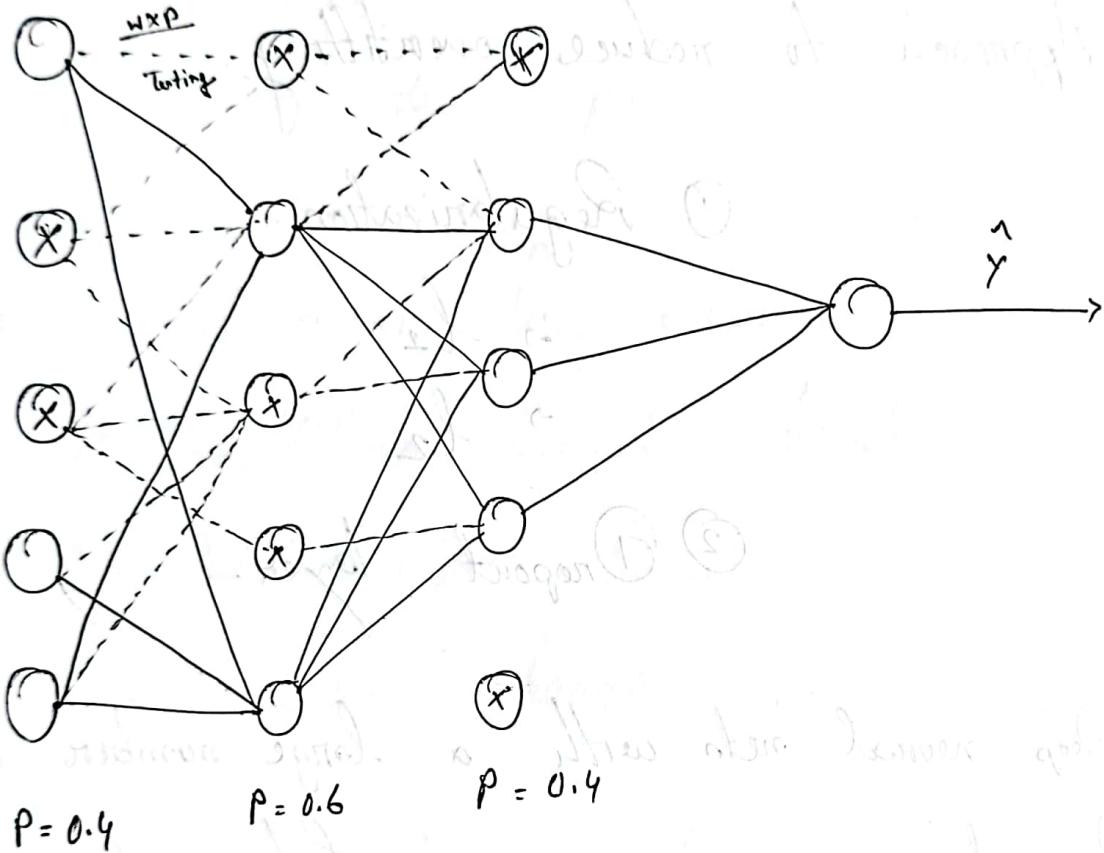
$$\rightarrow l_1$$
$$\rightarrow l_2$$

② Dropout layer.

Deep neural nets with a large number of parameters are very powerful machine learning systems.

Dropout Ratio

$$0 \leq p \leq 1$$



↳ "Randomly selected" some neuron and rest of them will be "deactivate".

training data $\xrightarrow{\text{activating}}$
 $\xrightarrow{\text{deactivating}}$

what about my test data? ?

for my [but data] every neuron will be connected

Additional work.

All these weight during training which are fixed. in multiply by probability

$w \times p$ weight and probability
→ dropout ratio.

Q: How do we select the $w \times p$ value & mean ratio?

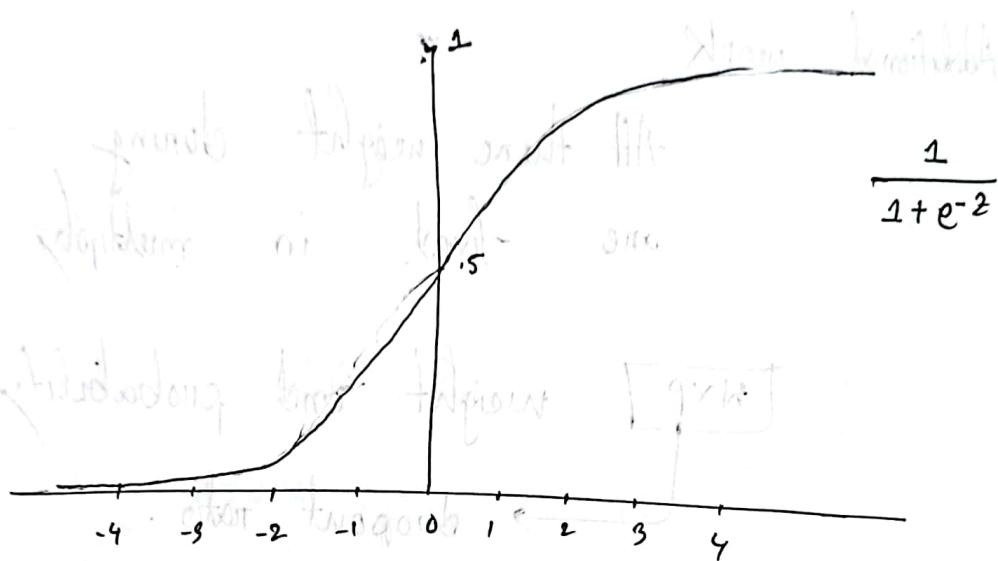
- Hyperparameter optimization (S) soft
- if high overfitting problem p should be higher then 0.5.

Question : During forward propagation are these dropped out neuron also zero (deactive / turned off)

during back-prop ?

Activation Function / Transfer Function.

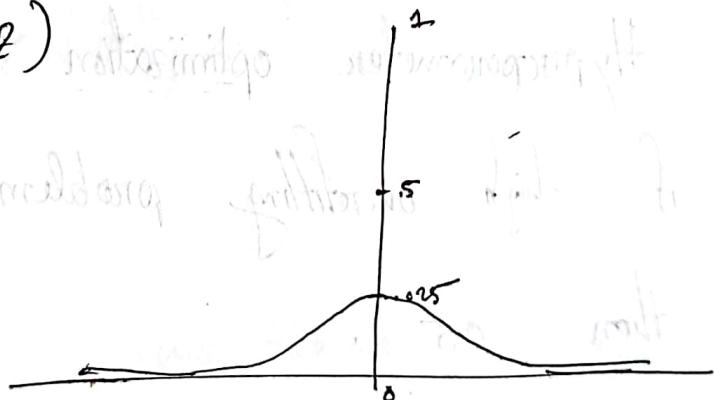
① Sigmoid Activation Function



$$z = \sum_{i=1}^n w_i x_i + b \quad 0 \leq \sigma(z) \leq 1$$

$$z = \text{Act}(z)$$

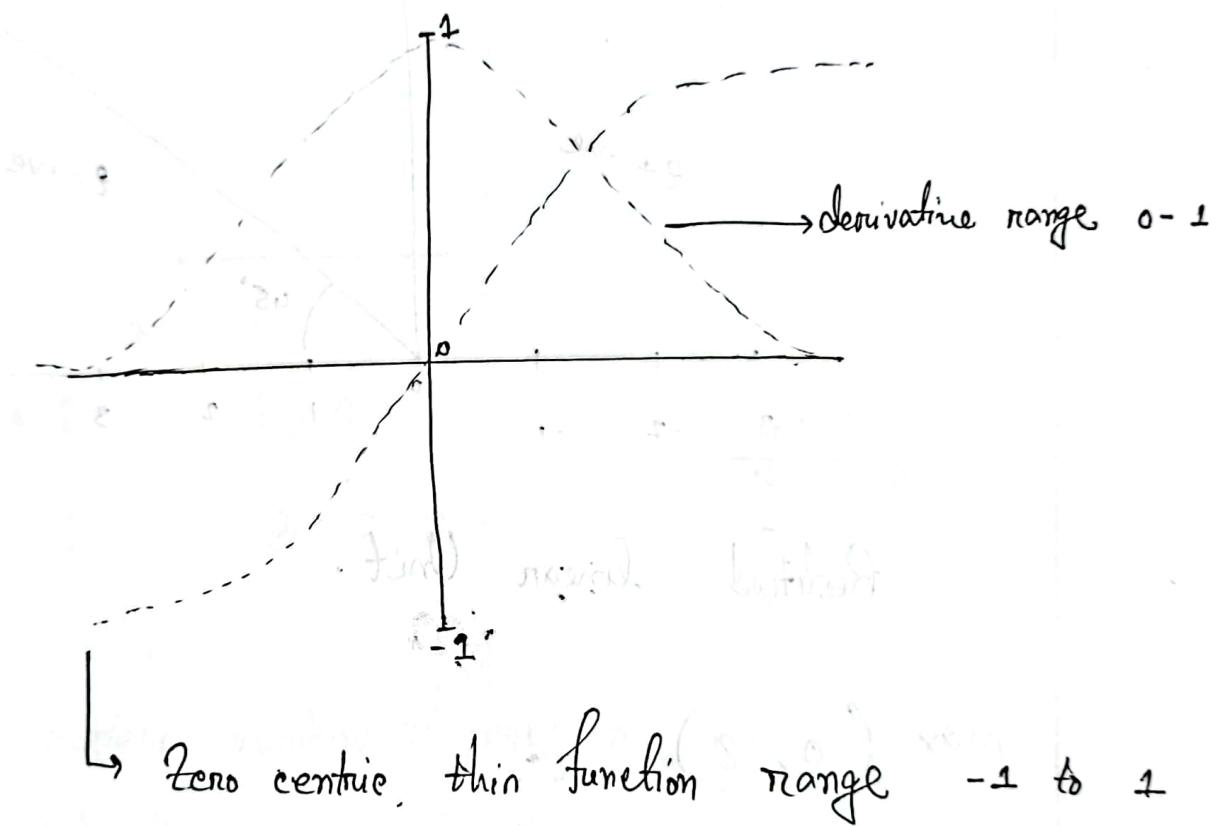
$$\frac{d(\sigma)}{dz}$$



derivative of Sigmoid $0 \leq \frac{d\sigma(z)}{dz} \leq 0.25$

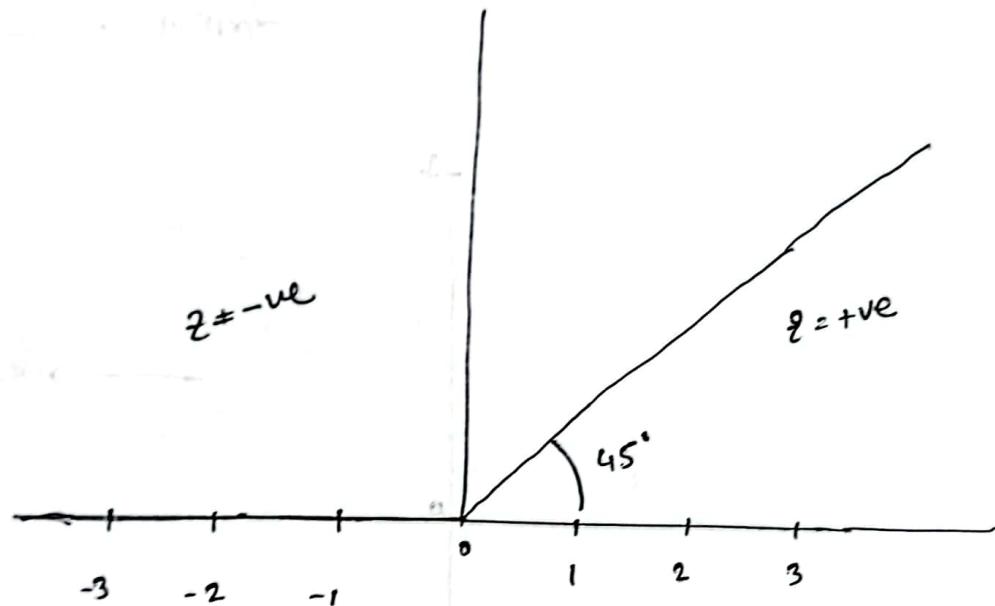
② Threshold Activation function,

$$\tanh = \frac{1 - e^{-2z}}{1 + e^{-2z}}$$



Both Sigmoid and \tanh face vanishing gradient problem.

③ ReLU activation function.



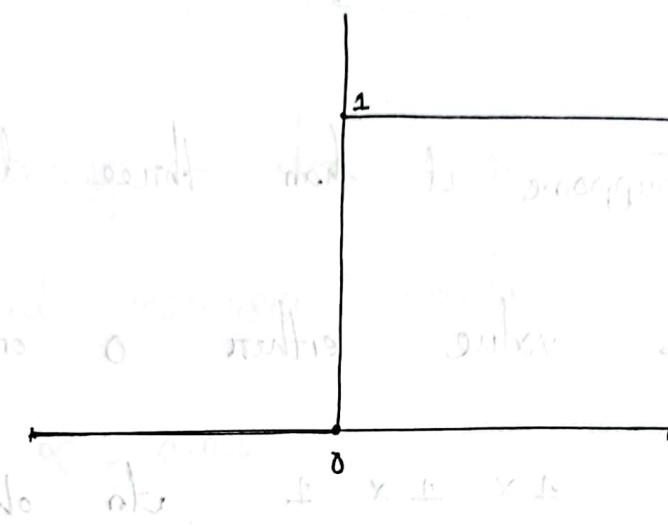
Rectified linear Unit.

$$\max(0, z)$$

$$R = \begin{cases} z & z > 0 \\ 0 & z \leq 0 \end{cases}$$

- The main catch here is that the ReLU function doesn't activate all the neurons at the same time.

Derivative of ReLU



$$R = \begin{cases} 1 & z > 0 \\ 0 & z \leq 0 \end{cases}$$

$$\frac{dR}{dz} ?$$

** Zero has no derivatives.

→ Only certain numbers of neurons are activated.

→ Computationally more efficient than sigmoid, tanh.

→ limitation
dying relu.

$$\text{Weight}_{\text{new}} = \text{Weight} - \frac{\text{dLoss}}{\text{dW}} * n$$

Suppose it has three derivative.

derivative value either 0 or 1.

$$1 \times 1 \times 1 \text{ it's ok}$$

$$0 \times 1 \times 1 = 0 \quad \text{but } 0 < 0.5$$

weight are not updated and neuron is

in not activate

⇒ Dead neuron

This is called
Dying ReLU

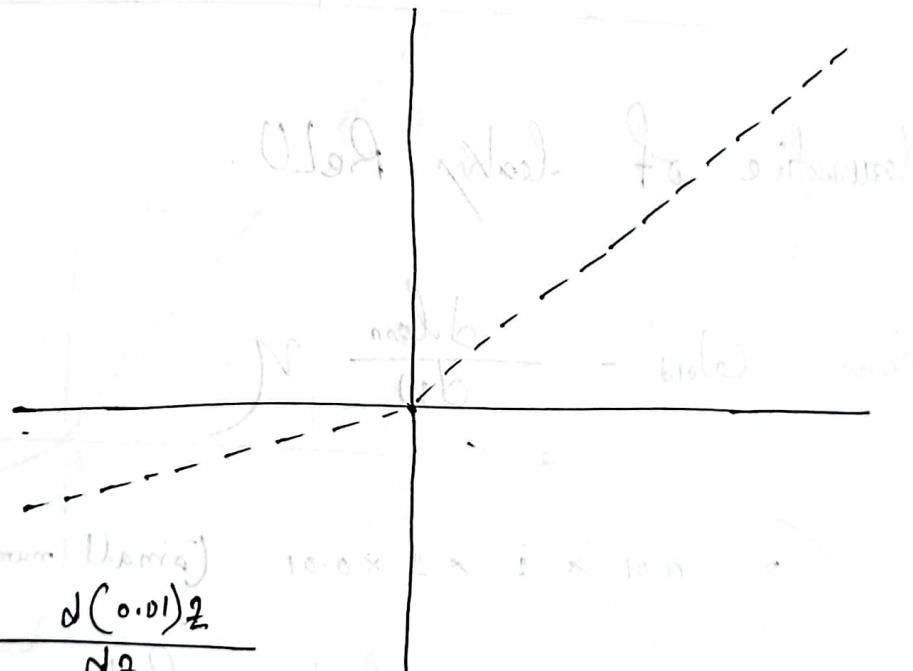
To fix this problem we use leaky ReLU.

$$R = \begin{cases} z & z > 0 \\ \alpha z & z \leq 0 \end{cases}$$

$f(z) = \max(0.01z, z)$

Instead of being 0, when $z < 0$, a leaky ReLU "allows" small non-zero constant gradient α .

Normally $\alpha = 0.01$



[actual theory]

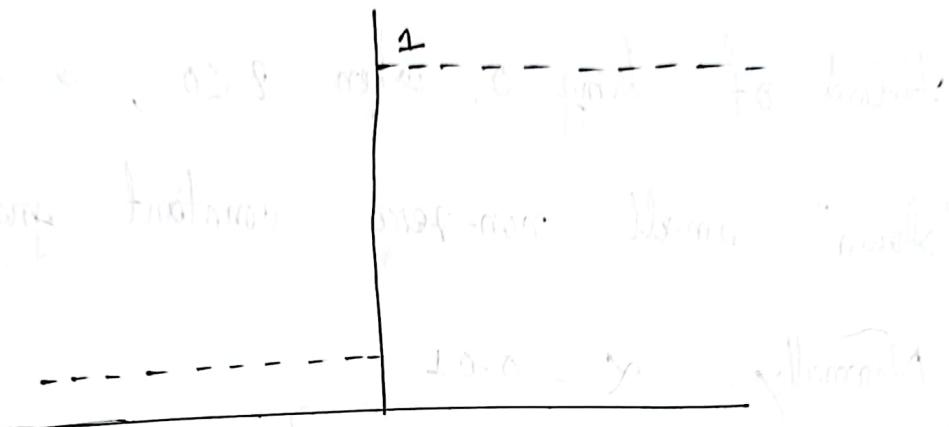
$$\frac{\partial f(0.01)z}{\partial z}$$

= some value,
not zero

$$\approx 0.01$$

so we can ignore the derivative if

$$R' = \begin{cases} 1 & z > 0 \\ \alpha & z < 0 \end{cases}$$



derivative of leaky ReLU.

$$W_{new} = W_{old} - \frac{d \text{loss}}{dW} \eta$$

→ $0.01 \times 1 \times 1 \times 0.01$ (small number)

Whenever we use leaky ReLU there will be a chance of vanishing gradient problem.

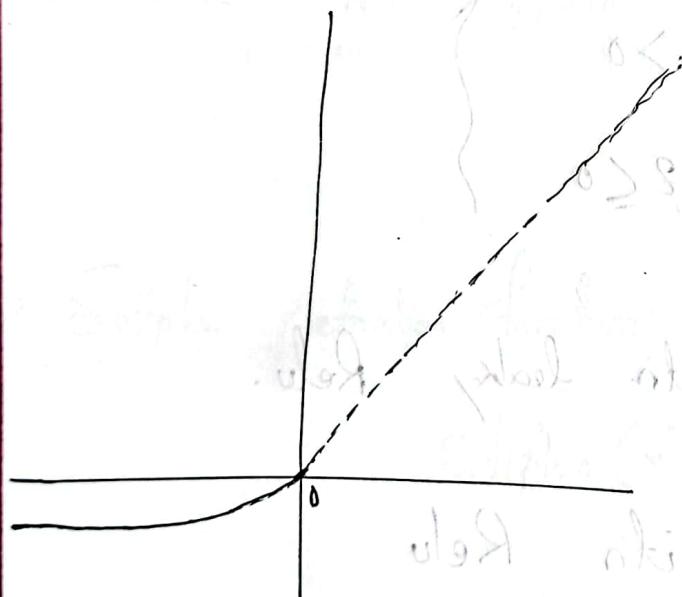
→ if maximum number of weight is -ve.

Exponential Linear Unit.

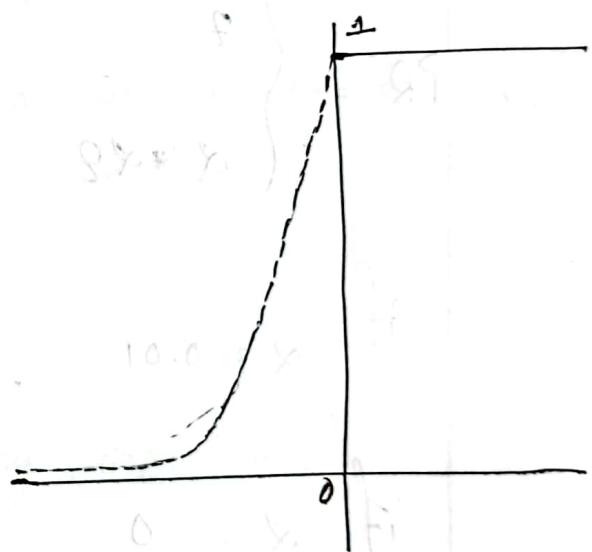
ELU is very similar to RELU except negative inputs.

$$f(x) = \begin{cases} x & \text{if } x > 0 \\ \alpha(e^x - 1) & \text{otherwise.} \end{cases}$$

↳ learnable parameter α can be fixed.



zero centered



derivative

$$\frac{d\alpha(e^x - 1)}{dx}$$

derivative is not possible
for 0.
it's computationally
expensive.

$$F'(z) = \begin{cases} 1 & z > 0 \\ \alpha \cdot e^z & z \leq 0 \end{cases}$$

derivative of
ReLU function.

PRelu (Parametric ReLU)

$$PR = \begin{cases} z & z > 0 \\ \alpha * z & z \leq 0 \end{cases}$$

If $\alpha = 0.01$, it's leaky ReLU.

If $\alpha = 0$, it's ReLU

If $\alpha = \text{any value}$, can be solve any kind of problem. it's learnable parameter.

→ it become PReLU.

→ Swish activation function.

$$\rightarrow \text{Swish}(x) = x * \text{sigmoid}(x)$$

→ Used in LSTM

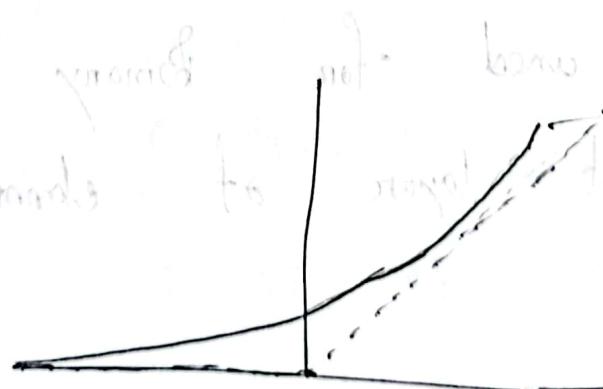
→ Also called Self-gating function.

whenever NN layers > 40

then we should use swish.

→ Softplus activation function

$$\text{Softplus}(x) = \log(\exp(x) + 1)$$



⇒ Softmax Activation Function.

→ Used as the activation for the last layer of a classification network.

→ Used for multiclass classification.

→ Converts a vector value to a probability distribution.

→ Returns tensor.

→ Used in a dense layer.

⇒ Sigmoid used for Binary classification in the last layer of classification.

Softmax

Suppose output

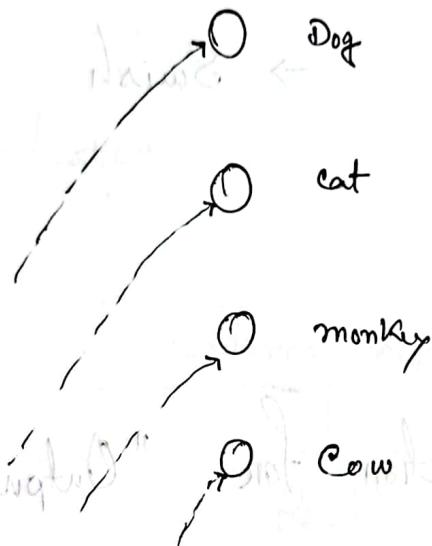
$$[40, 30, 10, 5]$$

$$\frac{e^{40}}{e^{40} + e^{30} + e^{10} + e^5} =$$

$$\frac{e^{30}}{e^{40} + e^{30} + e^{10} + e^5} =$$

$$\frac{e^{10}}{e^{40} + e^{30} + e^{10} + e^5} =$$

$$\frac{e^5}{e^{40} + e^{30} + e^{10} + e^5} =$$



Lemon

$$[0.6, 0.2, 0.1, 0.1] = 1$$

it's dog

highest probability
that in the final
output

Sigmoid

$$\begin{array}{ccc} 0 & \text{Dog} & 0.6 \\ 0 & \text{Cat} & 0.4 \end{array}$$

it's dog

→ Activation function for hidden layer.

→ ReLU

→ Swish if, for neural network

having a depth greater
than 40 layers.

→ Activation function for "Output layer".

→ Regression : linear Activation function.

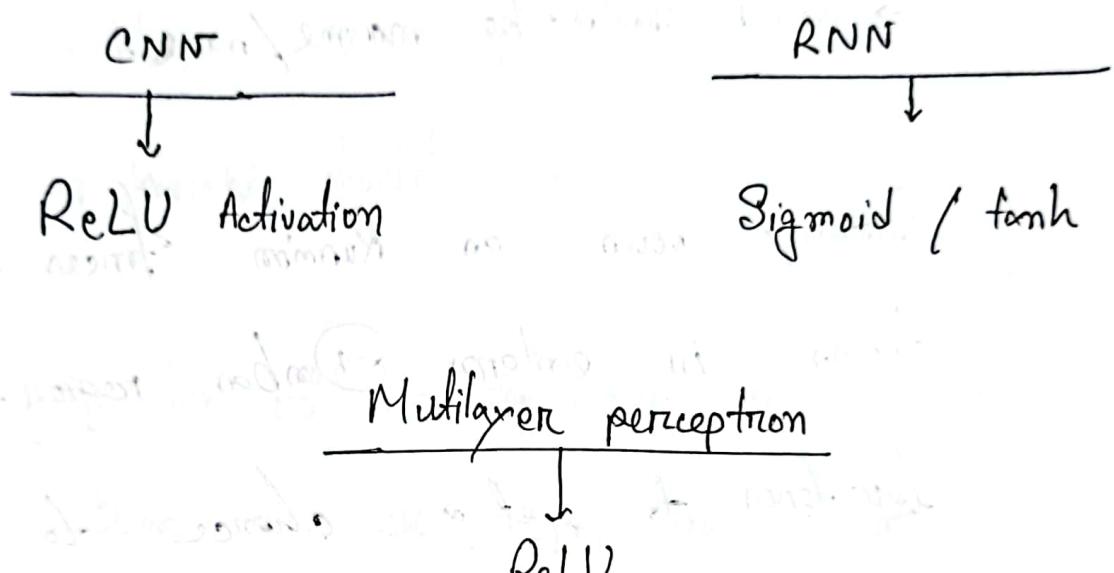
→ Binary classification : Sigmoid

→ Multiclass classification : Softmax.

→ Multilabel classification : Sigmoid.

Denne layer

Activation function in "hidden layers"



All activation function

- ① Sigmoid
- ② Threshold / tanh
- ③ Relu
- ④ leaky Relu
- ⑤ PRelu
- ⑥ ELU
- ⑦ Softplus
- ⑧ Swish
- ⑨ Softmax
- ⑩ linear

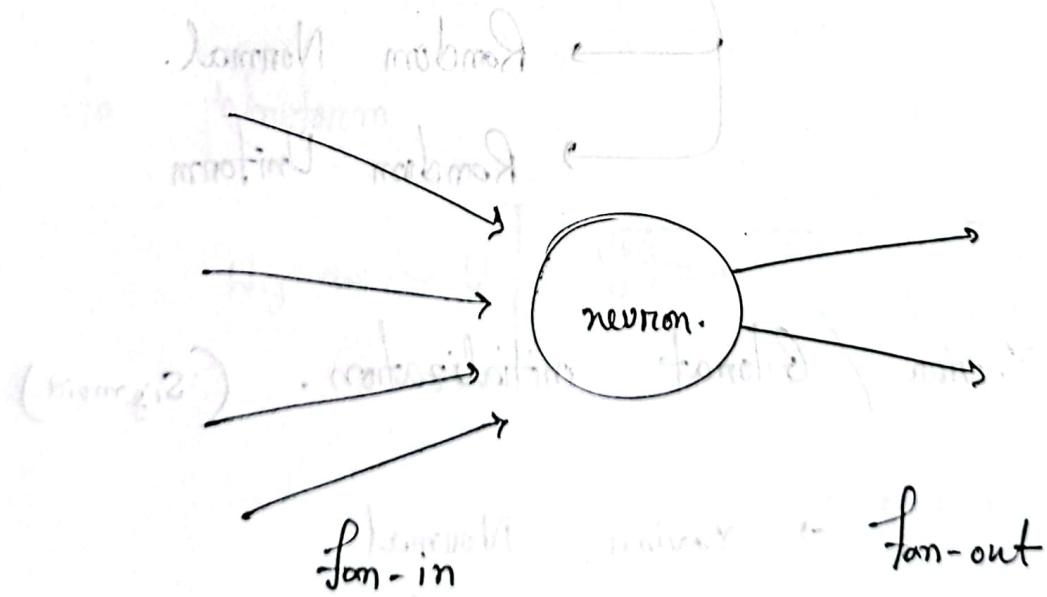
Weight Initialization.

- \rightarrow Weight should be small. (not very small or big)
- \rightarrow Weight should not be same.
- \rightarrow Weight should have good variance.

Technique: None of the technique can say best.

Some work for some dataset.

Some of the work for some of activation function properly



① Uniform distribution.

$W_{ij} \sim \text{Uniform}$

$$\left[\begin{array}{c} \text{lower bound} \\ -1 \\ \frac{1}{\sqrt{\text{fan-in}}} \end{array} \right] \quad \left[\begin{array}{c} \text{upper bound} \\ 1 \\ \frac{1}{\sqrt{\text{fan-out}}} \end{array} \right]$$

→ Not zero initialize

② Random initialize

- Random Normal.
- Random Uniform.

③ Xavier / Glorot initialization. (Sigmoid)

→ Xavier Normal

→ Xavier Uniform

→ Xavier Normal

(Suitable for Sigmoid.)

$$w_{i,j} \sim N(0, \sigma)$$

$$\sigma = \sqrt{\frac{2}{(\text{fan-in} + \text{fan-out})}}$$

→ Xavier Uniform

$$w_{i,j} \sim U\left[\frac{-\sqrt{6}}{\sqrt{\text{fan-in} + \text{fan-out}}}, \frac{\sqrt{6}}{\sqrt{\text{fan-in} + \text{fan-out}}}\right]$$

→ He - Uniform

$$w_{i,j} \sim U\left[-\sqrt{\frac{6}{\text{fan-in}}}, \sqrt{\frac{6}{\text{fan-out}}}\right]$$

Activation function = ReLU.

→ the Normal.

$$W_{ij} \approx N(0, \sigma)$$
$$\sigma = \sqrt{\frac{\sigma^2}{\text{fan-in}}}$$

Suitable for Relu - leaky Relu

→ Gradient descent,

Stochastic GrD.

Mini-batch GrD

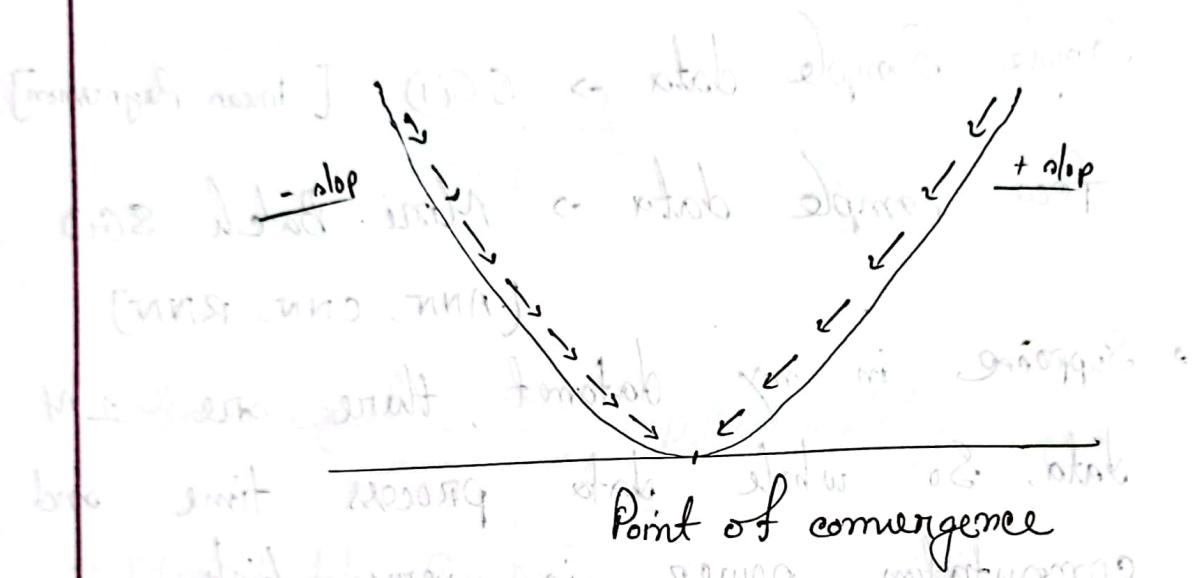
Batch GrD

Stochastic GrD with momentum.

local minima

Global minima.

→ Gradient Descent - IBNI Article.



$$w_{\text{new}} = w_{\text{old}} - \frac{\partial J(w)}{\partial w} n$$

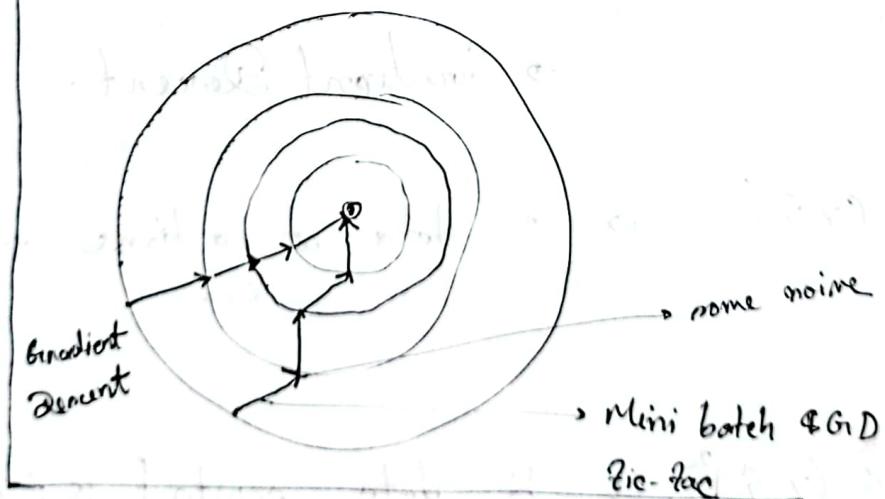
$\sum_{i=1}^n (y_i - \hat{y}_i)^2$ ⇒ All the data points(n) of the dataset
Gradient Descent.

$(y_i - \hat{y}_i)^2$ ⇒ 1 data at a time → Stochastic gradient Descent.

$\sum_{i=1}^k (y_i - \hat{y}_i)^2$ ⇒ k data points [$k < n$] → Mini-batch SGD

- ① All data points \rightarrow Gradient Descent.
 - ② Single Sample data \rightarrow SGD [linear Regression]
 - ③ few sample data \rightarrow Mini-Batch SGD
[ANN, CNN, RNN]
- Suppose in my dataset, there are $1M$ data. So whole data process time and computation power is very high.
 - \Rightarrow for that case we use Mini-Batch SGD

SGD



Sample

$$\left[\frac{\partial \text{Loss}}{\partial \text{Word}} \right] \approx \begin{cases} \text{Minibatch SGD} \\ \text{SGD} \end{cases}$$

Population

$$\left[\frac{\partial \text{Loss}}{\partial \text{Word}} \right] \text{GD}$$

Problem higher

Computational memory

→ Never equal but approximately equal.

→ Mini-batch SGD naive solution : Stochastic Gradient Descent with momentum.

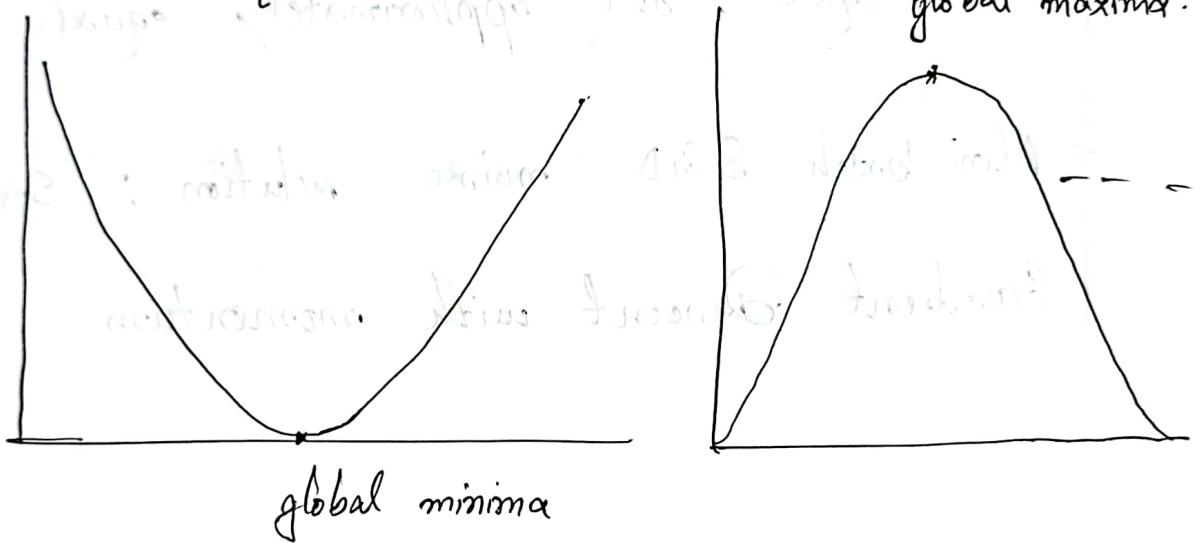
Global minima , local minima

④ MB- SGD

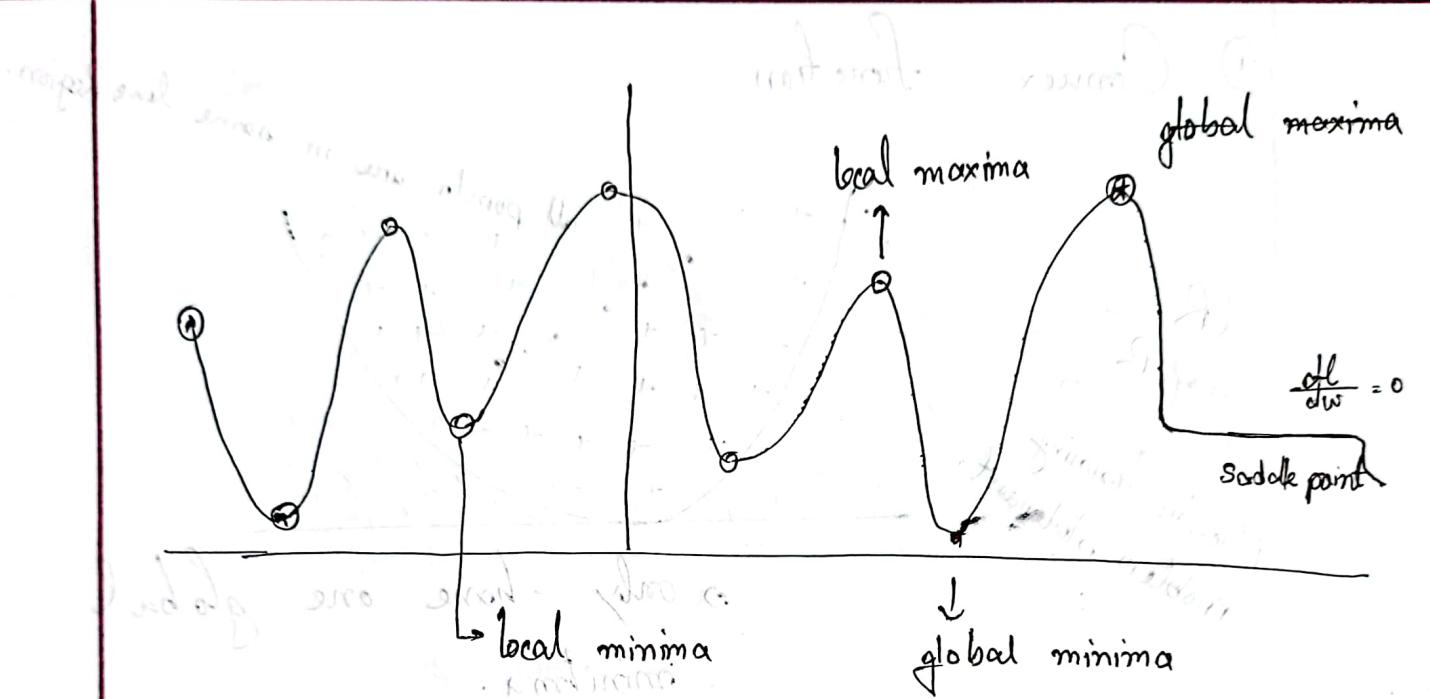
$$\text{Loss}(w) = \sum_{i=1}^K (y - \hat{y})^2$$

Mean square error

global maxima



But in deep learning there are different loss function and higher dimension value.
Curve are not like that.



derivative of points = 0

$$w_{\text{new}} = w_{\text{old}} - 0$$

$w_{\text{new}} = w_{\text{old}}$ ← Convergence point

Gradient Descent



Problem

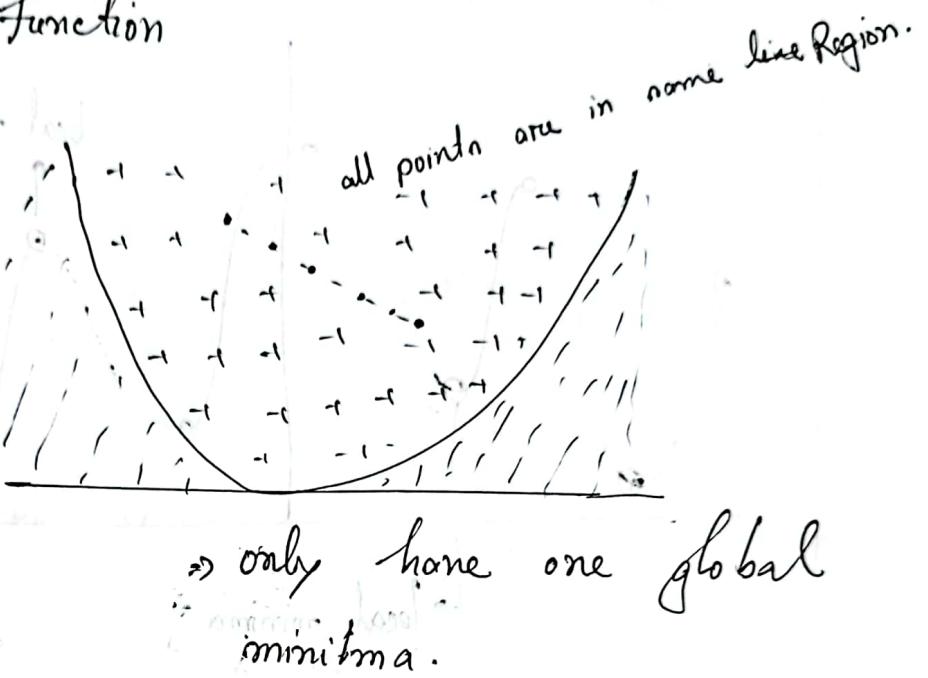
not reach the convergence point for non-convex function.

gradient descent gets stuck in local

① Convex function

LR
log-R

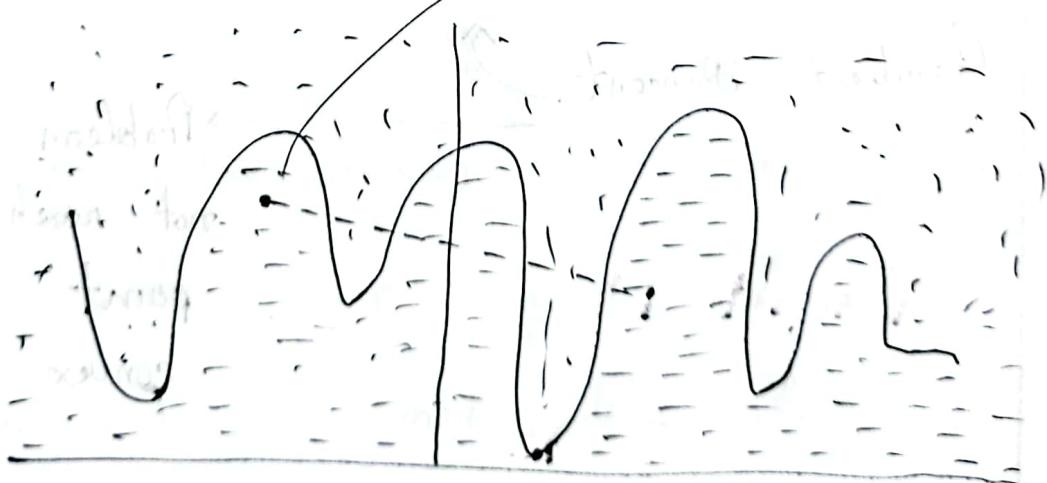
Machine learning
Problem statement.



② non-convex function

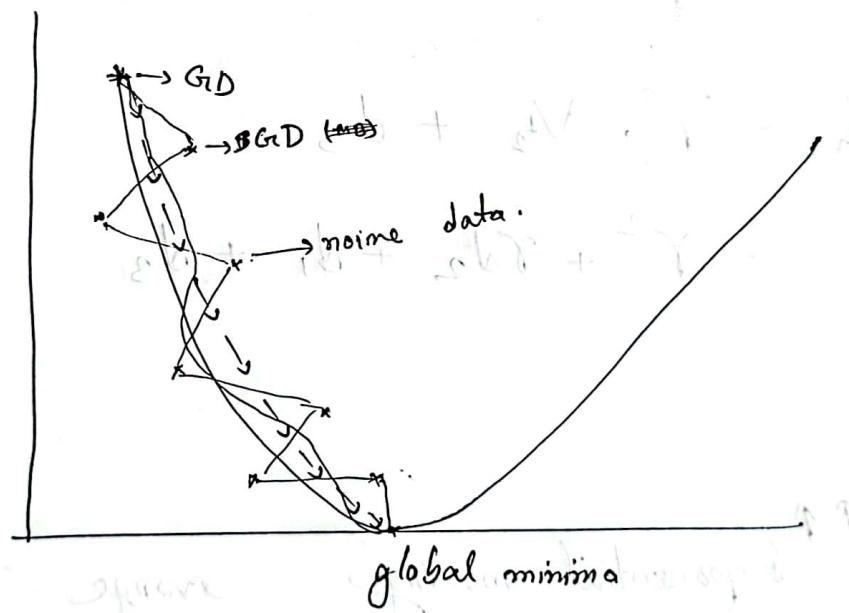
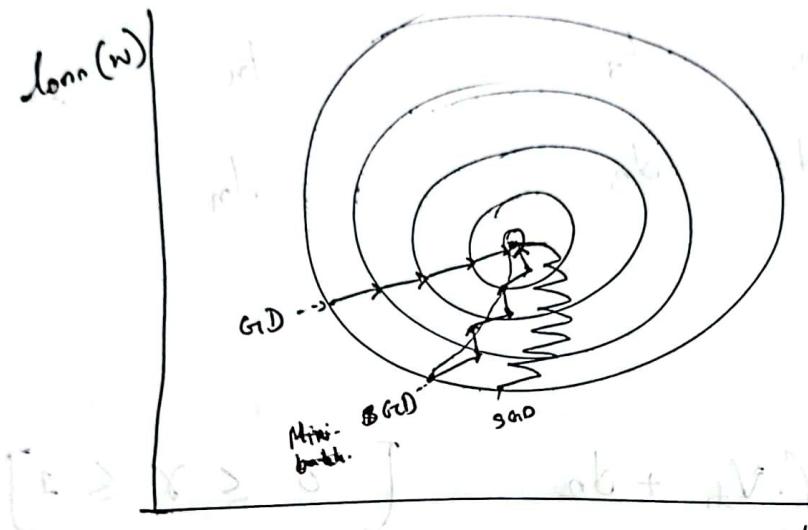
Deep learning

more than two region.



Occurred in mostly deep learning technique.

SGD with Momentum.



$S.G.D \Rightarrow$ Reach the global minima take large amount of time. So, we try to remove those noise.

How? : Exponential moving average.

$$t_1 \quad t_2 \quad t_3 \quad \dots \quad t_n$$

$$d_1 \quad d_2 \quad d_3 \quad \dots \quad d_n$$

$$V_{t_1} = d_1$$

$$V_{t_2} = \gamma \cdot V_{t_1} + d_2 \quad [0 \leq \gamma \leq 1]$$

$$V_{t_3} = \gamma \cdot V_{t_2} + d_3$$

$$= \gamma^2 + \gamma d_2 + d_1 + d_3$$

↑↑ Exponential moving average formula.

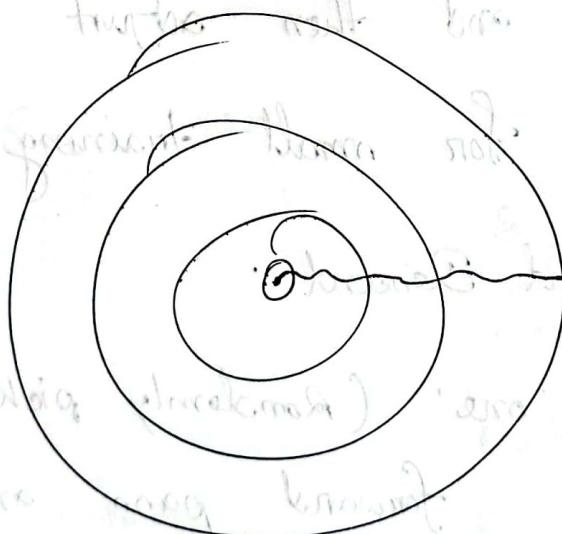
Batch gradient descent and gradient descent are same!

GD / SGD

$$w_{\text{new}} = w_{\text{old}} - \eta \frac{\partial \text{loss}}{\partial w_{\text{old}}}$$

By applying moving average (momentum)

$$\Rightarrow w_{\text{old}} \left[v_{t-1} + \eta \frac{\partial \text{loss}}{\partial w_{\text{old}}} \right]$$



$$v_t = \gamma \left[\frac{\partial l}{\partial w_{\text{old}}} \right]_t + \gamma^2 \left[\frac{\partial l}{\partial w_{\text{old}}} \right]_{t-1} + \dots$$

Analytic vidya →

Batch Gradient Descent

Stochastic "

Mini Batch Gradient Descent.

Batch Gradient Descent

- Use all training sample for one forward pass and then adjust weights.
- Good for small training set.

Stochastic Gradient Descent.

- Use 'one' (Randomly picked) sample for a forward pass and then adjust weight.
- Good when training set is very big and we don't want too much computation.

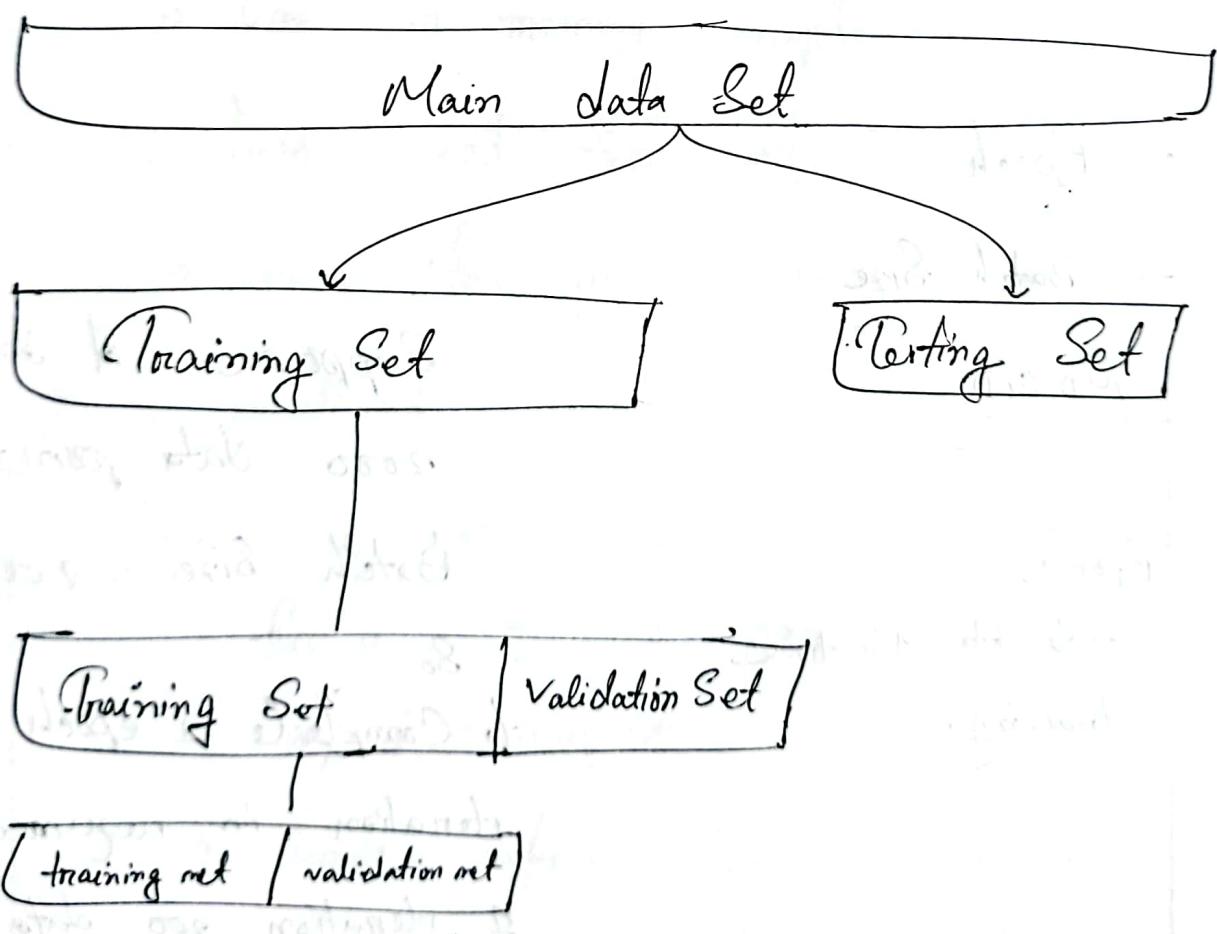
- Batch Gradient Descent
Batch Size = Size of training set.
- Stochastic Gradient Descent
Batch Size = 1
- Mini-Batch Gradient Descent
 $1 < \text{Batch size} < \text{Size of training set}$.
Sample of dataset.

- Epoch
 - Batch Size
- MD-GD

Epoch,
whole data in time
training.

Suppose, A dataset
2000 data point.
Batch Size = 200
So,
to Complete 1 epoch = 10
iteration is required.
1 iteration 200 data in
Pass.

- Training Set : Used to train the model.
- Validation Set : Used to optimize model parameters
- Testing Set : Used to get an "unbiased" estimate of the "final model" "Performance".



Adagrad Optimizer

Adaptive gradient descent.

$$W_{\text{new}} = W_{\text{old}} - \eta \left(\frac{\partial L_{\text{nn}}}{\partial W_{\text{old}}} \right)$$

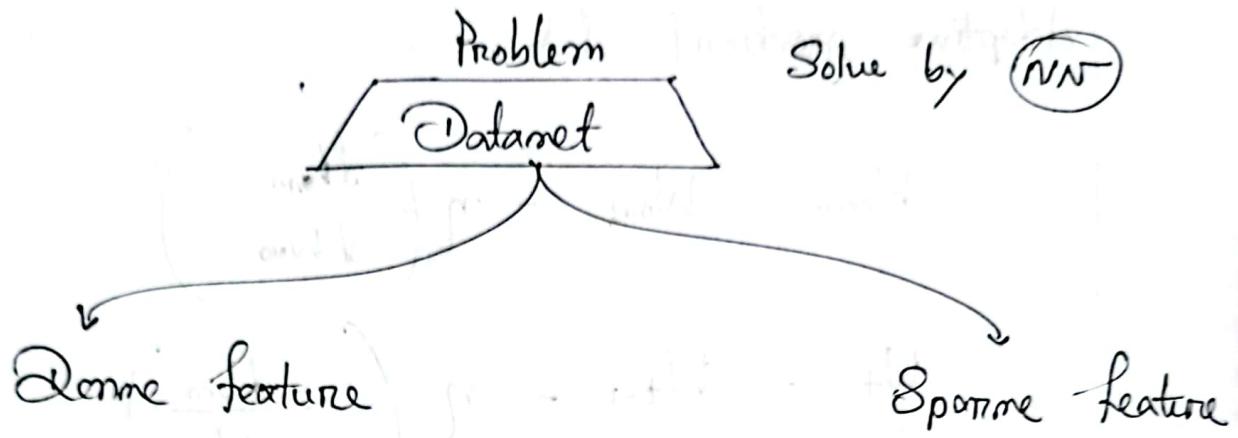
$$W_t = W_{t-1} - \eta \left(\frac{\partial L_{\text{nn}}}{\partial W_{t-1}} \right)$$

→ This learning rate used to be same in GD, SGD, MB-GD, for all the neurons for all the hidden layers.

In adaptive GD, "learning rate" should be different

- ① for each and every neuron.
- ② for each hidden layer.
- ③ based on different iteration.

already learned



→ Most of the feature
are non zero.

→ Most of the feature/value
are $\rightarrow 0$ (zeroes).

Initialize $\eta = 0.01$

↓
here we can't
use same learn-
ing rate?

Adaptive optimizer main concept
is that,

Can we use different learning
rate with respect to.

Different feature,
Weight,
layer,
Neuron,
iteration.

Dense/Sparse feature \leftarrow use different learning rate (Adaptive)

Adagrad Optimizer formula,

$$W_t = W_{t-1} - \eta_t^+ \left(\frac{\partial \text{loss}}{\partial W_{t-1}} \right)$$

$$\eta_t^+ = \frac{\eta}{\sqrt{\alpha_t + \epsilon}}$$

Constant learning rate.
initial 0.01

Epsilon always small positive number.

$$\alpha_t = \sum_{i=1}^t \left(\frac{\partial \text{loss}}{\partial w_i} \right)^2$$

t = 3 [3rd iteration] First calculate α_t .

t = 2

t = 1

α_t will high number.

So, η_t^+ will be small

when,

$$\alpha_t \uparrow$$

γ_+^1 will be low during
updation.

So, Weight will be decreasing slowly. and

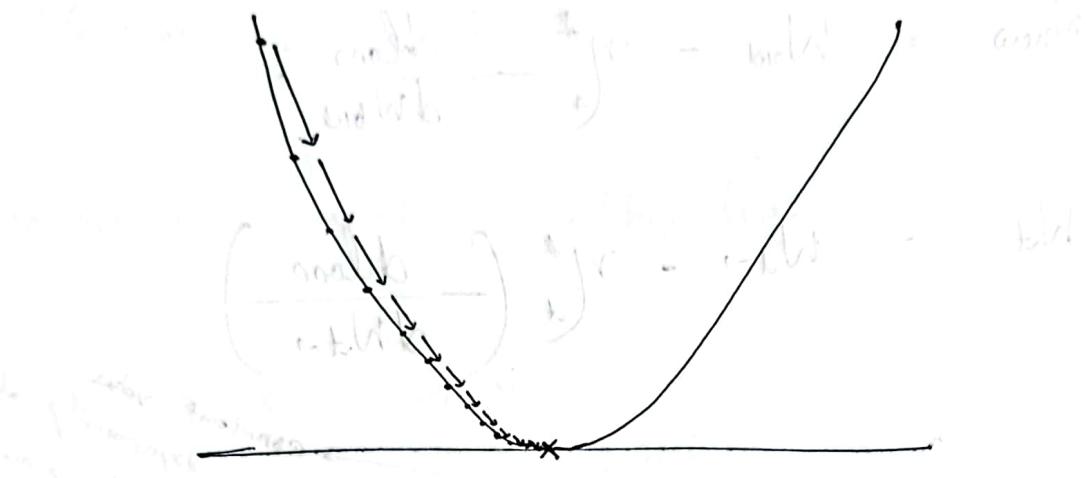
it will convergence very nicely.

first, little bigger $\rightarrow \beta$ (learning rate)

then, when it is small \rightarrow Smaller \rightarrow Smaller \rightarrow ...

So, Adagrad optimizer changing the learning
rate slowly and reach the conver-
gence point very nicely.

gradient Descent



One Disadvantage:

As the number of iteration increases

Some point of time α_t is
very high. $[w_{new} \approx w_{old}]$

The Solution : AdaDelta and RMSProp optimizer.

Prevent of higher value.

"define iteration step"

Adadelta and RMSProp

$$W_{\text{new}} = W_{\text{old}} - \eta^* \frac{\text{d} \text{loss}}{\partial W_{\text{old}}}$$

$$W_t = W_{t-1} - \eta^* \left(\frac{\text{d} \text{loss}}{\partial W_{t-1}} \right)$$

$$\eta^* = \frac{\eta}{\sqrt{W_{\text{avg}}(t) + \epsilon}}$$

→ exponential moving average.

$$W_{\text{avg}}(t) = \gamma \cdot W_{\text{avg}}(t-1) + (1-\gamma) \left(\frac{\text{d} \text{loss}}{\partial W_t} \right)^2$$

$$\gamma = 0.95 \text{ (most cases)}$$

→ There is no summation, so we restrict it
to restrict the number (t).

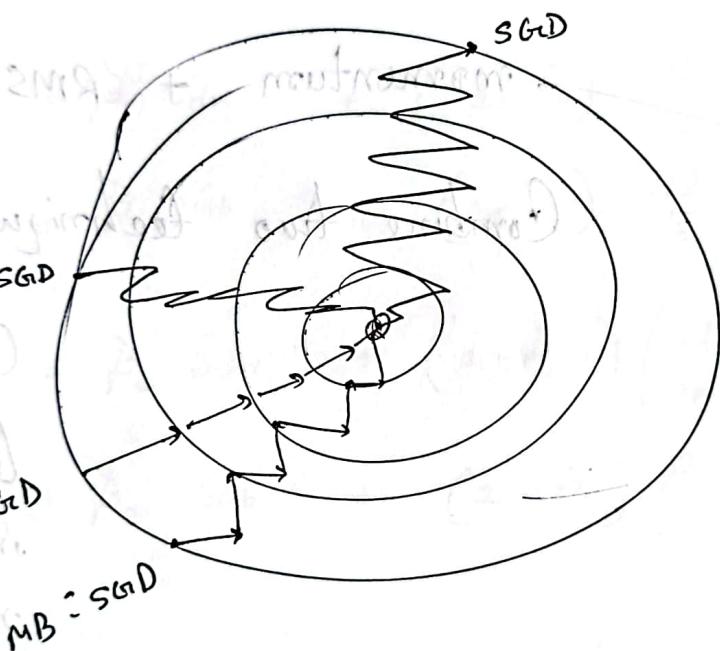
$$\eta' = \frac{\eta}{\sqrt{SdW_t + \epsilon}}$$

→ exponential moving avg.

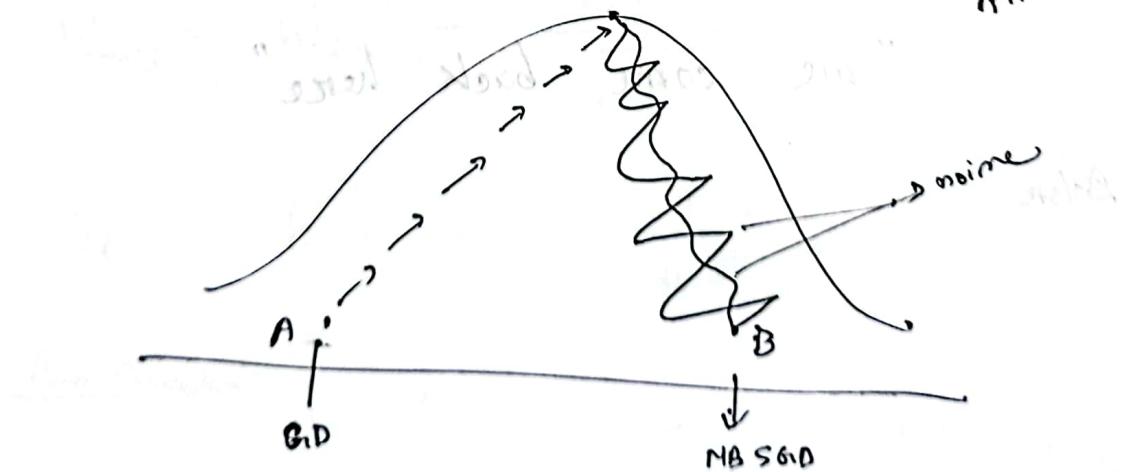
$$\left. \left(\beta SdW_{t-1} + (1-\beta) \left(\frac{\text{d} \text{loss}}{\partial W_t} \right)^2 \right) \right\} \text{Same}$$

$$\rightarrow \text{Iteration} = \frac{\text{Dataset size}}{\text{Batch Size}}$$

Gradient Descent \leftarrow Wikipedia.



Reduce noise with momentum.



check it in hmt
Adam optimizer.

Best optimizer algorithm \rightarrow "Adaptive Moment Estimation"
and digital \rightarrow from the previous.

Adam optimizer = Gradient Descent with

momentum + RMSProp

Smoothing

Combine two technique.

Change the learning rate in a efficient manner.

"we come back here"

After

Momentum,

RMS Prop,

$$V_{dw}, V_{db} =$$

$$S_{dw}, S_{db} \sim \text{[initially]}$$

On iteration t , compute $\frac{\partial l}{\partial w}$, $\frac{\partial l}{\partial b}$ using
the minibatch.

$$V_{dw} = \beta_1 V_{dw} + (1 - \beta_1) \frac{\partial l}{\partial w}$$

$$V_{db} = \beta_1 V_{db} + (1 - \beta_1) \frac{\partial l}{\partial b}$$

$$S_{dw} = \beta_2 S_{dw} + (1 - \beta_2) \left(\frac{\partial l}{\partial w} \right)^2$$

$$S_{db} = \beta_2 S_{db} + (1 - \beta_2) \left(\frac{\partial l}{\partial b} \right)^2$$

$$w_{newt} = w_{t-1} -$$

$$b_t = b_{t-1} -$$

$$\frac{\eta + V_{dw} \text{ correction}}{\sqrt{S_{dw} + \epsilon}}$$

$$\frac{\eta + V_{db} \text{ correction}}{\sqrt{S_{db} + \epsilon}}$$

Bias Correction

$$V_{dw} \text{ correction} =$$

$$\frac{V_{dw}}{(1 - \beta_1)}$$

$$S_{dw} = \frac{S_{dw}}{(1 - \beta_2)}$$

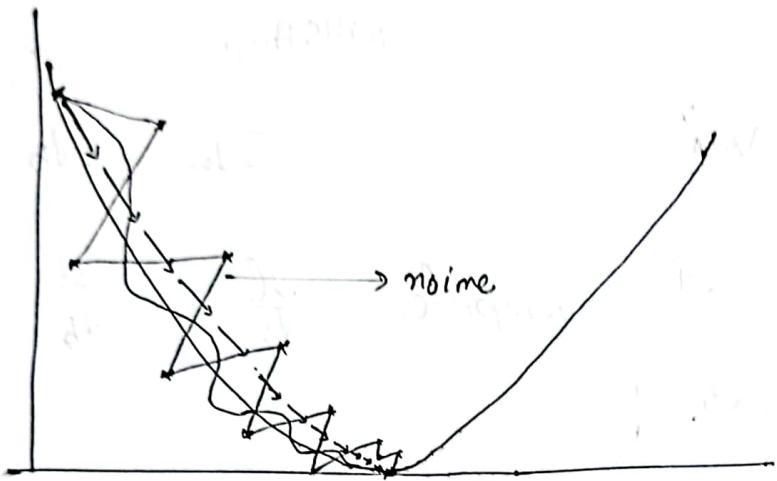
$$V_{db} =$$

$$\frac{V_{db}}{(1 - \beta_1)}$$

$$S_{db} = \frac{S_{db}}{(1 - \beta_2)}$$

Adam

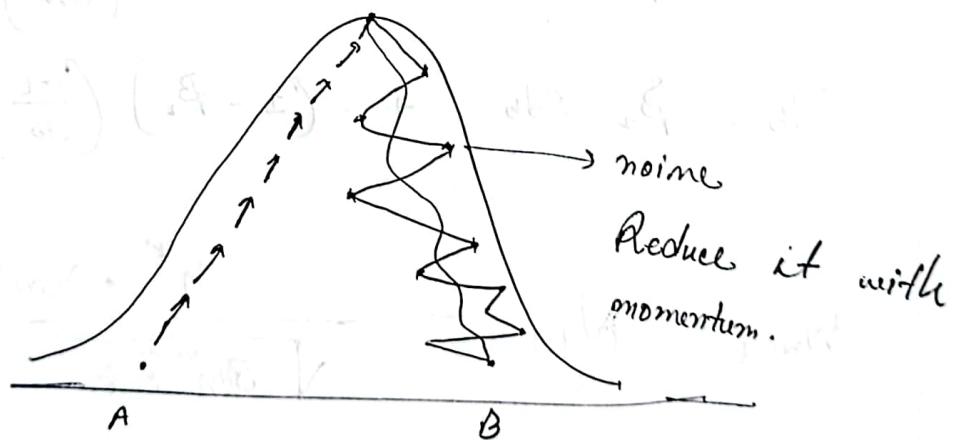
Back in SGD with momentum.



→ GD

→ MB SGD

apply momentum.



GD

SGD (MB)

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

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

Instead of computing gradient descent, we calculate Exponential weighted average

Let's consider time series problem.

Time, $t_1, t_2, t_3, \dots, t_n$

Data, $a_1, a_2, a_3, \dots, a_n$

Smoothing concept:

initially
 $V_{t_1} = a_1$

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

$$= \beta a_1 + (1-\beta) * a_2$$

$$= 0.95 a_1 + (1-0.95) * a_2$$



Learnable parameters / hyperparameters tuning.

$$\text{if } \beta = 0.1$$

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

mont carlo.

$$V_{d_2} = 0.1 a_1 + (0.9) * a_2$$

$$V_{d_3} = \beta V_{d_2} + (1-\beta) * a_3$$

$$= \beta (\beta * V_{d_1} + (1-\beta) * a_2) + (1-\beta) * a_3$$

Exponential Weighted average; weight updation in.

$$w_t = w_{t-1} - \eta \frac{\partial \text{loss}}{\partial w_{t-1}}$$

$$= w_{t-1} - \eta \cdot \boxed{Vdw} \rightarrow ?$$

for bias

$$b_t = b_{t-1} - \eta \boxed{Vdb} \rightarrow ?$$

$$Vdw = \beta Vdw_{t-1} + (1-\beta) \frac{\partial l}{\partial w_{t-1}}$$

$$Vdb = \beta Vdb_{t-1} + (1-\beta) \frac{\partial l}{\partial b_{t-1}}$$

Initially, V_{dw} and V_{db} will be zero (0).

Implementation details.

Initially $V_{dw} = 0$, $V_{db} = 0$.

On iteration t , in middle epoch

Compute Δw , Δb on the current mini-batch.

$$V_{dw,t} = \beta V_{dw,t-1} + (1-\beta) \frac{\partial \text{loss}}{\partial w_{t-1}}$$

$$V_{db,t} = \beta V_{db,t-1} + (1-\beta) \frac{\partial \text{loss}}{\partial b_{t-1}}$$

$$w = w_{t-1} - \eta V_{dw,t}$$

$$b = b_{t-1} - \eta V_{db,t}$$

Something of convergence.

SGD with momentum.

Adadelta implementation details (optional notes)

- On iteration t in minibatch Compute

$$\frac{\partial l_{\text{loop}}}{\partial w}, \frac{\partial l_{\text{loop}}}{\partial b} \text{ on current}$$

mini-batch, keeps states for next

$$S_{\partial w} = \beta S_{\partial w} + (1-\beta) \left(\frac{\partial l}{\partial w} \right)^2$$

$$S_{\partial b} = \beta S_{\partial b} + (1-\beta) \left(\frac{\partial l}{\partial b} \right)^2$$

$$w_t = w_{t-1} - \eta_t \frac{\partial l}{\partial w_{t-1}}$$

$$b_t = b_{t-1} - \eta_t \frac{\partial l}{\partial b_{t-1}}$$

$$\eta_t = \frac{\eta}{\sqrt{S_{\partial w} + \epsilon}}$$

$$= \frac{\eta}{\sqrt{S_{\partial b} + \epsilon}}$$

Loss Function in Deep Learning.

→ Loss function / Cost function.

Error for 1/Single record, like in SGD first pass single record and calculate the error then backpropagate.

$$\text{Loss function} = (y - \hat{y})^2$$

→ In Mini-Batch gradient descent, we pass batch of record, or in gradient descent pass whole record, then we calculate loss function called.

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

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

Check Konan losses

→ Regression Problem losses.

① Squared Error Loss

$$l_{\text{loss}} = (\gamma - \hat{\gamma})^2$$

② Mean Square Error

$$l_{\text{loss}} = \sum_{i=1}^n (\gamma - \hat{\gamma})^2 * \frac{1}{n}$$

$$\text{/ Actually cost} = \sum_{i=1}^n (\gamma - \hat{\gamma})^2 * \frac{1}{n}$$

③ Mean Absolute Error Loss

$$l_{\text{loss}} = |\gamma - \hat{\gamma}|$$

$$= abn(\gamma - \hat{\gamma}) + \frac{1}{n}$$

$$= \sum_{i=1}^n abn(\gamma - \hat{\gamma}) + \frac{1}{n}$$

④ Huber Loss

Combine $[MSE + MAE]$

→ Quadratic equation.

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

→ linear equation.

δ → delta value, hyperparameter.

→ Classification problem losses,

• Cross Entropy

$$\text{loss} = -\gamma \log(\hat{\gamma}) - (1-\gamma) \log(1-\hat{\gamma})$$

→ for Binary cross entropy.

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

→ Solve only
Binary classification

$$\hat{\gamma} = \text{Sigmoid function} \cdot \frac{1}{1+e^{-x}}$$

⇒ Multiclass Cross entropy loss: \rightarrow category

$$L_{\text{ce}}(x_i, y_i) = - \sum_{j=1}^c y_{ij} \times \log(\hat{y}_{ij})$$

I/P \rightarrow Multiclass.

Good, bad, Neutral.

f_1	f_2	f_3	category	y_{11}	y_{12}	y_{13}	
2	4	2	good	[1, 0, 0]	y_1	y_2	y_3
5	6	9	bad	[0, 1, 0]	y_1	y_2	y_3
3	2	5	neutral	[0, 0, 1]	y_1	y_2	y_3

3 Bits

→ This we called "Categorical Cross Entropy".

y_i is one hot encoded target vector.

$y_i = [y_{i1}, y_{i2}, y_{i3}]$ i.e. Row number.

$y_{i1} \rightarrow$ category.

y_{ij} = $\begin{cases} 1 & \text{if } i\text{-th element is in } j\text{-th class} \\ 0 & \text{otherwise} \end{cases}$

\hat{y}_{ij} → use Softmax activation function.

$$\text{softmax}(z) = \frac{e^{z_i}}{\sum_{j=1}^k e^{z_j}}$$

→ Categorical cross entropy vs Sparse categorical cross entropy.

→ both have same loss function,

$$J(w) = -\frac{1}{N} \sum_{i=1}^N [\gamma_i \log(\hat{y}_i) + (1-\gamma_i) \log(1-\hat{y}_i)]$$

γ_i = true value

\hat{y}_i = predicted value

w = weight.

→ if your y_i are one hot encoded, we use categorical crossentropy. Example,

for 3 class classification:

$$\frac{e^{y_1}}{e^{y_1} + e^{y_2} + e^{y_3}} = \left(\begin{matrix} 1 \\ 0 \\ 0 \end{matrix}\right) \rightarrow [1, 0, 0], [0, 1, 0], [0, 0, 1]$$

⇒ if your y_i are integers, we use sparse categorical crossentropy. Example for 3 class classification,

$$[1], [2], [3]$$

Most Common Loss Function.

① Probabilistic / Classification losses.

- Binary Crossentropy class.
- Categorical Crossentropy class.
- SparseCategoricalCrossentropy class.

② Regression losses.

- Mean Squared error.
- Mean Absolute error.
- Huber class.

Hyperparameter tuning

1. How many numbers of hidden layers we should have?
2. How many numbers of neurons each hidden layers.
3. What should be the learning - Rate.

Keras - tuner

→ RandomSearch

→ GridSearchCV

- Go to the Keras tuner site.
- Embedding dimension.
- Dense function → create neuron.

Dense layer \hookrightarrow Dense function

[Keras/api/layers/Core-layers/Dense]

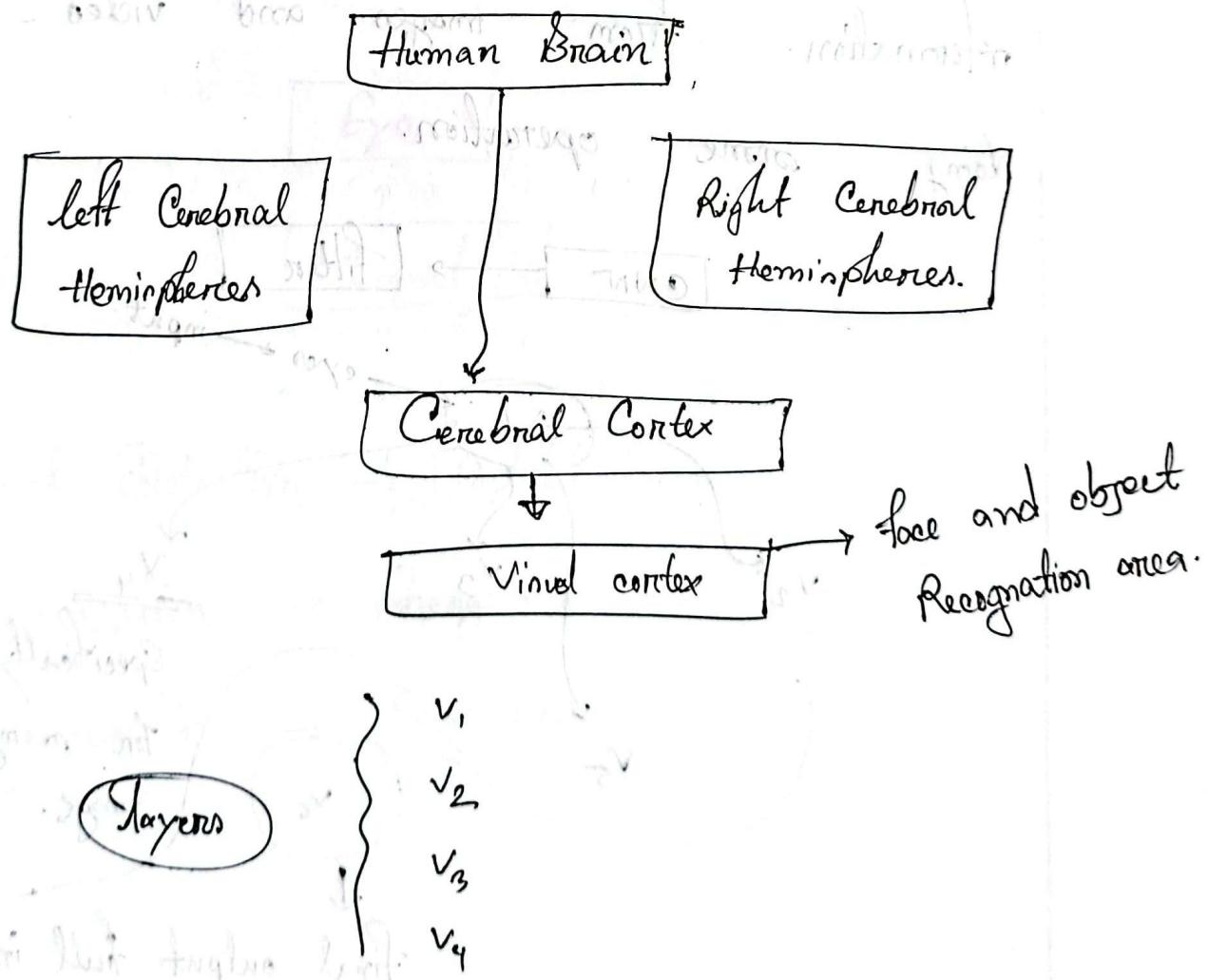
Brief understanding.

- Perplexity is a measurement of how well a probability distribution or probability model predicts a sample.
- It may be used to compare probability models.
- A low perplexity indicates the probability distribution is good at predicting the sample.

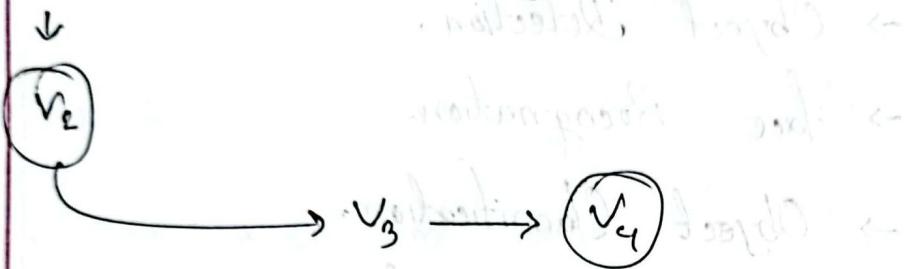
Convolution Neural Network.

Input format : Image, video

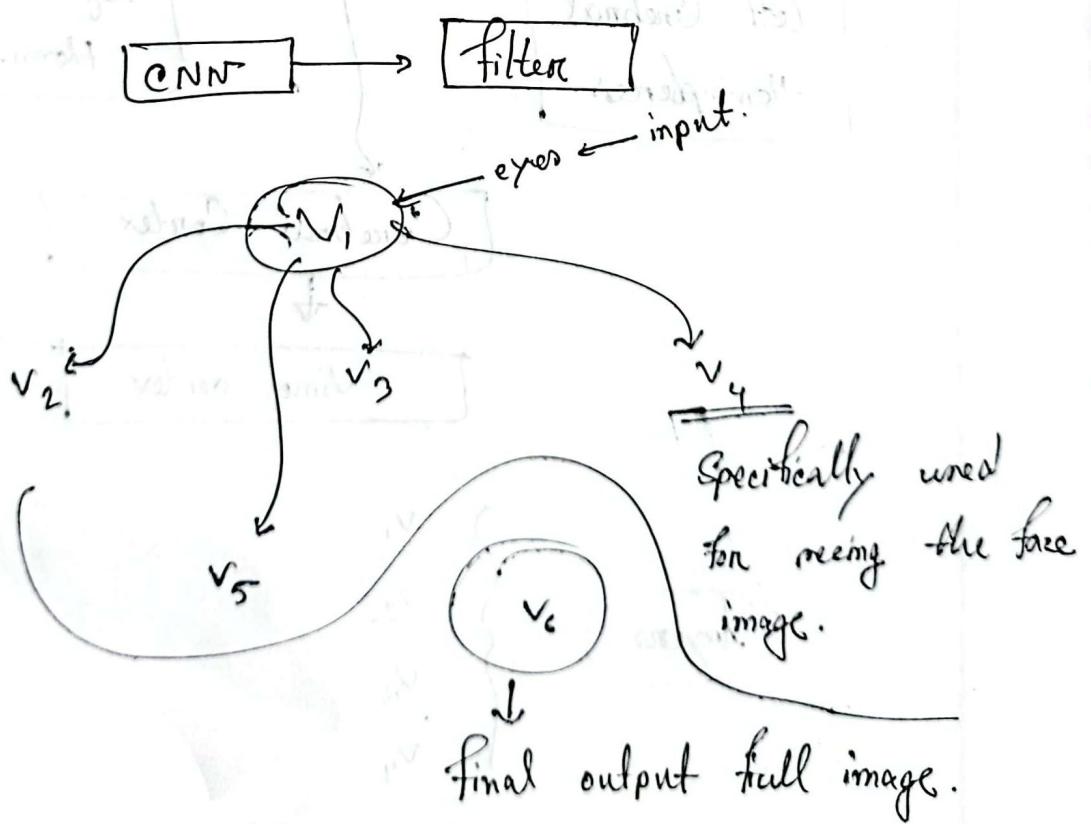
- Object Detection.
- Face Recognition.
- Object Classification.
- Image Recognition.

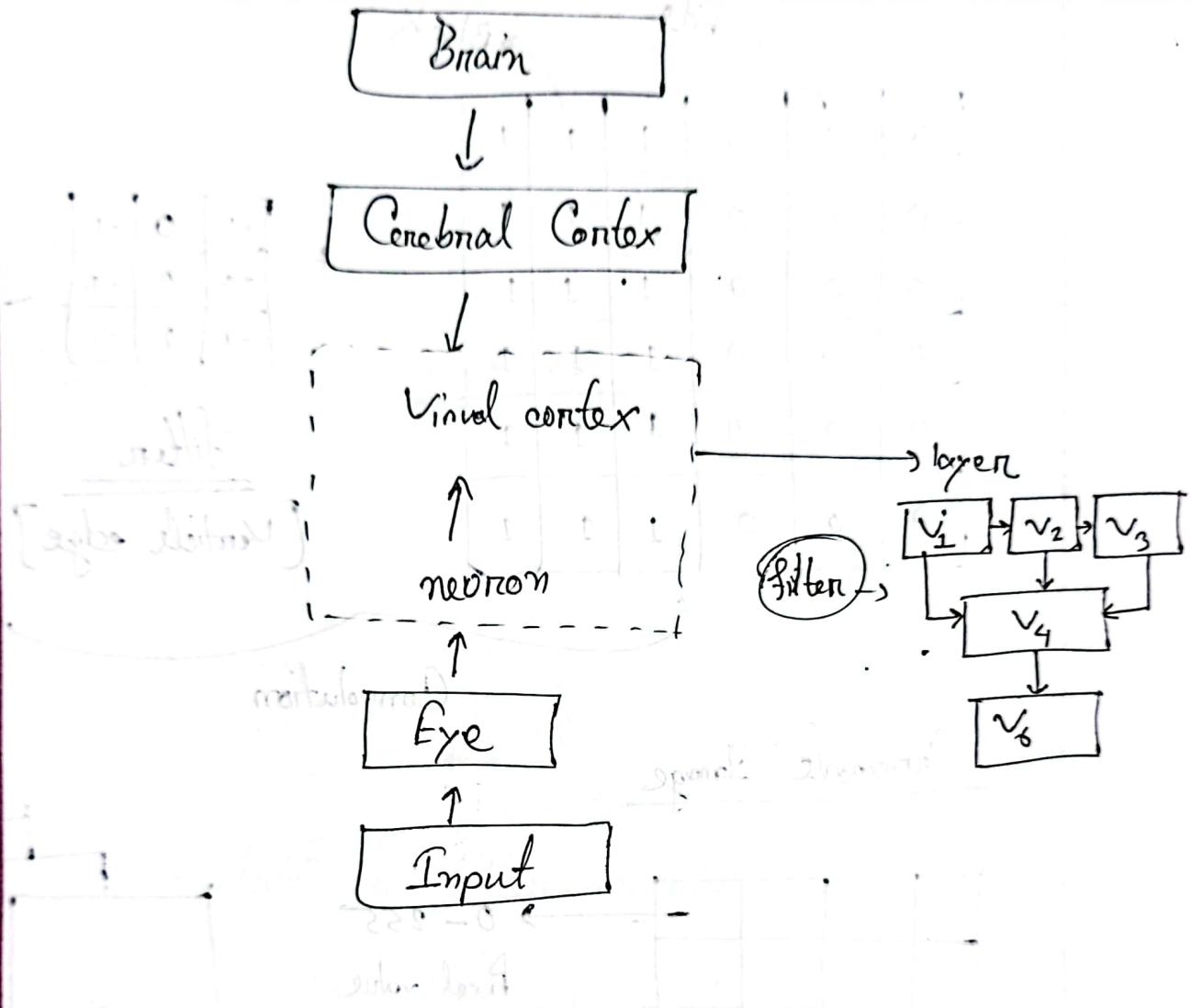


v_1 layer may be responsible for finding the edges of images (like cat, dog)



Each layer is responsible for Extract the information from images and video after doing some operation.

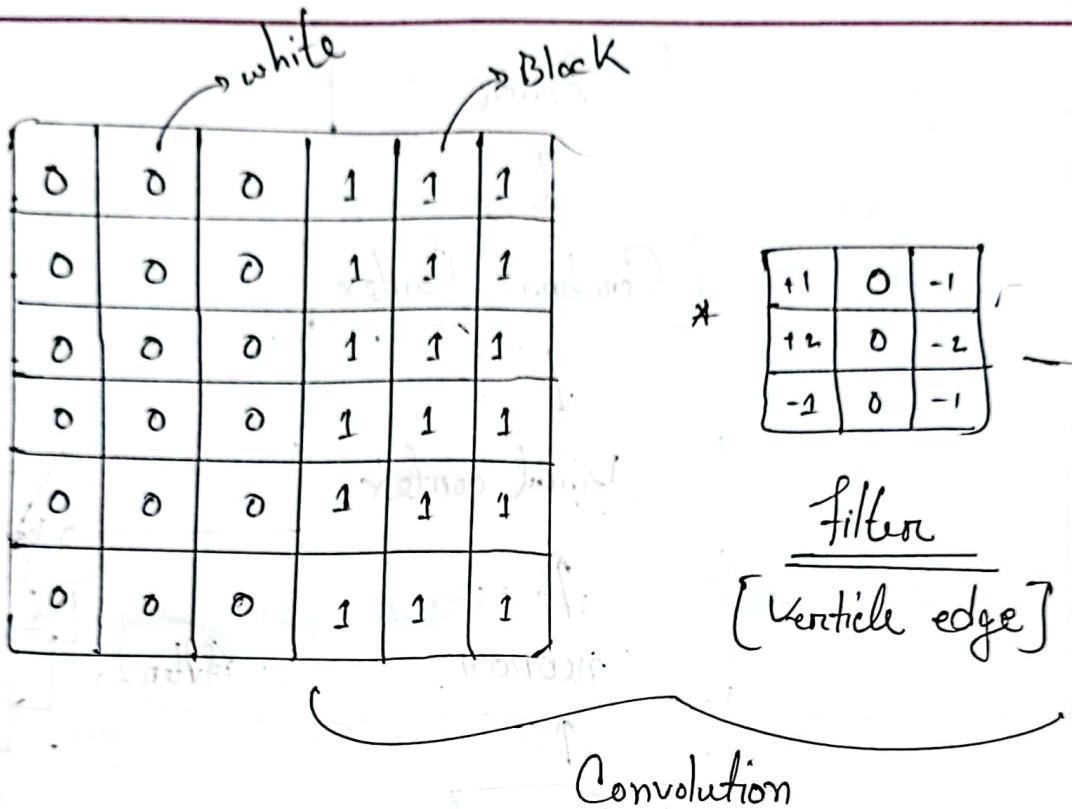




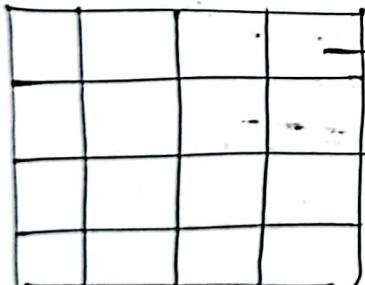
Object Detection → Detect

Recognition → Recognize

Segmentation → Segment.



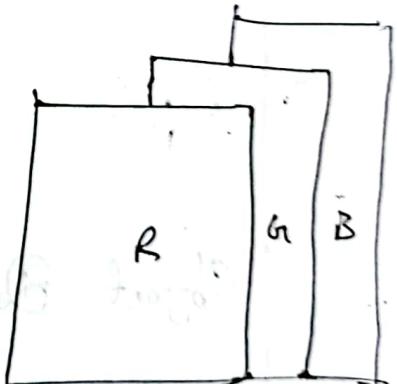
Grayscale Image



→ 0 - 255
Pixel value

channel → 1

4x4 x 1



channel → 3

4x4 x 3

0 - 255

normalize (0-1) if we use
minmax.

Vertical edge filter or detector

+1	0	-1
+1	0	-1
+1	0	-1

+ve

Horizontal edge detection

+1	+1	+1
0	0	0
-1	-1	-1

+ve

2	0	1
2	0	1
2	0	1

input matrix * filter matrix = ?

$$n_{\text{out}} = \left\lceil \frac{n_{\text{in}} + 2P - K}{S} \right\rceil + 1$$

n_{in} = number of input feature / input size

n_{out} = number of output feature / output size

P = Padding size

K = filter size / kernel size

S = Stride size.

CNN

- Convolution layer.
- Pooling layer.
- fully connected layer.

I		C	
0	0	0	1
0	0	0	1
0	0	0	1
0	0	0	1
0	0	0	1
0	0	0	1

+1	0	-1
+2	0	-2
+1	0	-1

Kernel size 3x3

Vertical edge detector.

$6 \times 6 \rightarrow$ thin edge

if stride = 1

1 step forward ----- ratio $\frac{6+0-3}{2} + 1$

1 step down ----- $\frac{3+0-3}{2} + 1$

thin kernel sum width = 3

256 channels

0	-4	-4	0
0	-4	-4	0
0	-4	-4	0
0	-4	-4	0

normalize
softmax

4×4

2×2 pool step

2×2 pool step

255	0	0	255
255	0	0	255
255	0	0	255
255	0	0	255

255	0	0	255
255	0	0	255
255	0	0	255
255	0	0	255

Block
in Dark
mole

White mole

Dark mole

Different filter detect the different edge /
different feature.

Padding in CNN

input image 6×6 * filter 3×3 = output 4×4

(yielding same amount of data/info.)
here losing some amount of data/info.

$$6 \times 6 * 3 \times 3 = 4 \times 4$$

but it should be

6×6 , because
of input size 6×6

$$n_{\text{out}} = \left(\frac{n_{\text{in}} + 2P - K}{S} \right) + 1$$

$$\Theta = \left[\frac{n_{\text{in}} + 2P - K}{S} \right] + 1 \quad [\text{if want output } 6 \times 6]$$

$$n_{\text{in}} = \Theta + (K - 1)$$

$$n_{\text{in}} + 2 \times 0 - 3 = 5$$

$$n_{\text{in}} = 5 + 3$$

$n_{\text{in}} = 8$ [input size should be $(6 \times 6) \rightarrow 8 \times 8$ do it]

by padding]

0	0	0	0	0	0	0
0						0
0						0
0						0
0						0
0						0
0	view	6x6	subgr.	0		
0	0	0	0	0	0	0

8×8

→ 0 padding.

0		

$= 6 \times 6$

$$m_{out} = \frac{6 + 2 - 3}{4} + 1$$

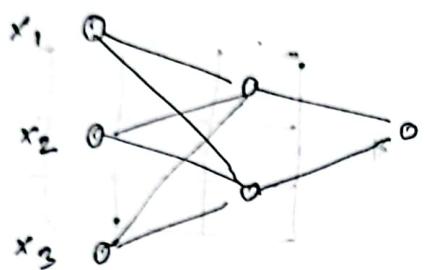


input $\rightarrow 6 \times 6 \rightarrow$ by padding ($p=1$) $\rightarrow 8 \times 8$

output
 6×6

0 padding means fill every pixel with 0.

ANN - CNN

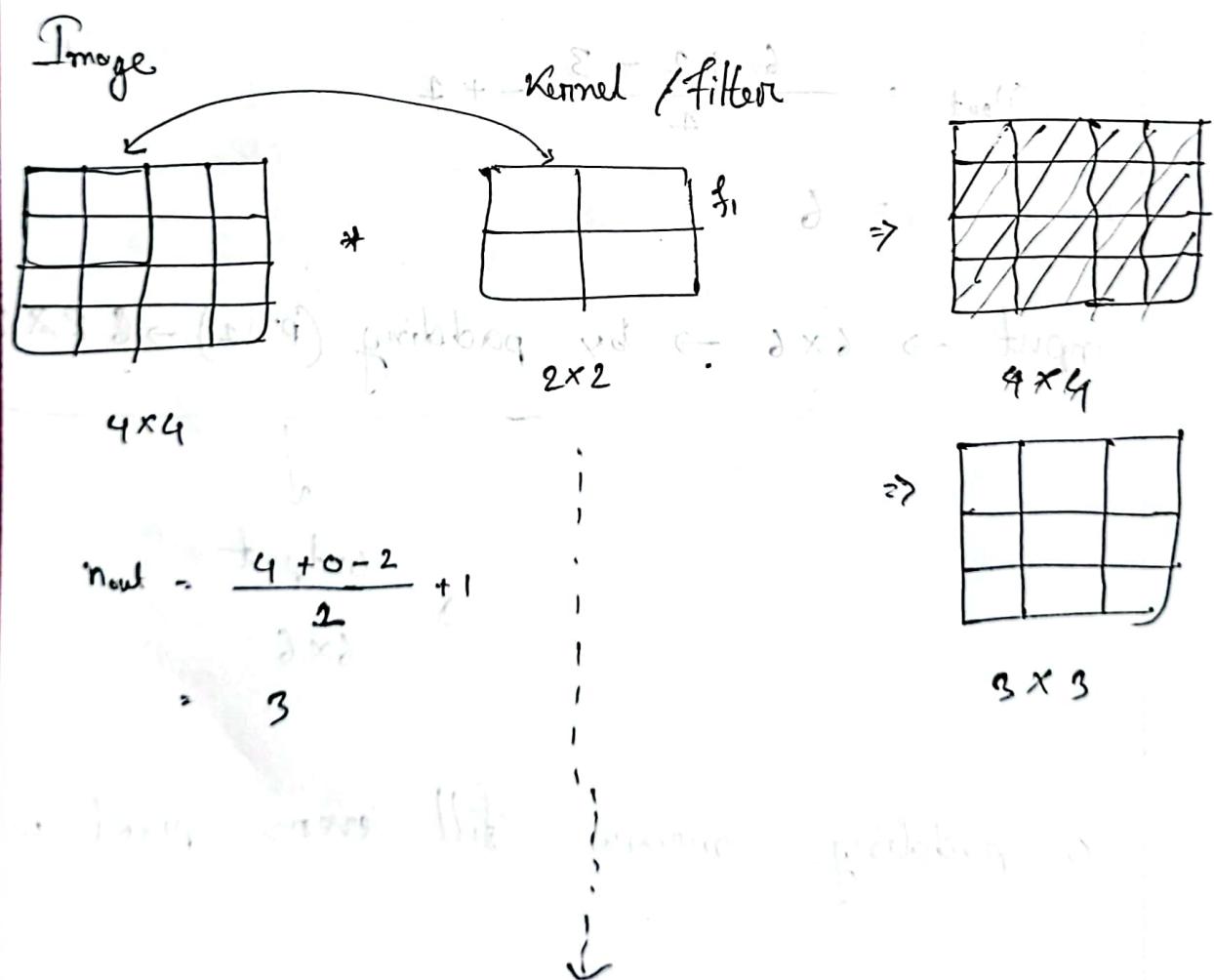


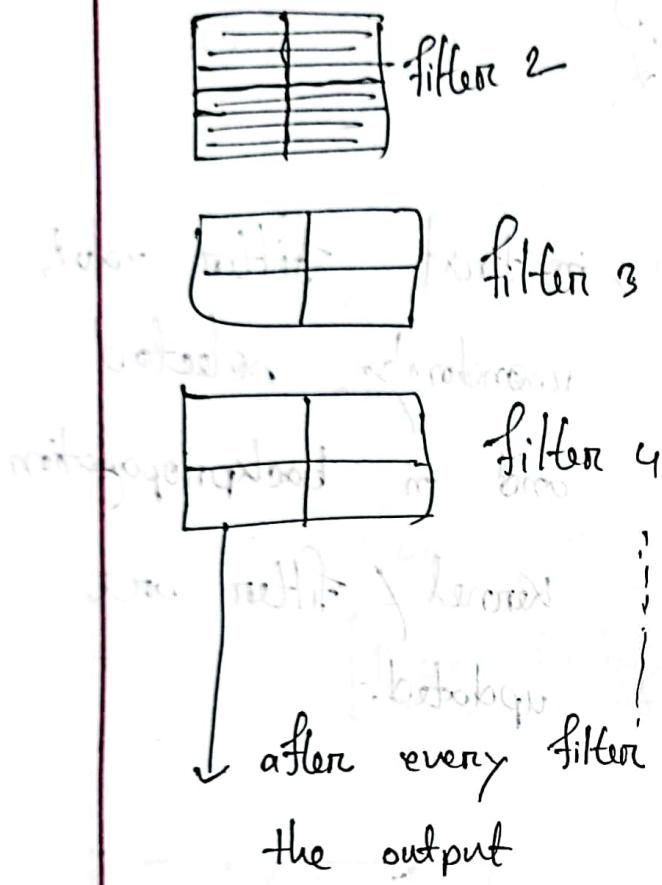
ANN

$$z = \mathbf{w}^T \mathbf{x} + \text{bias}$$

$$\hat{y} = \text{Act}(z)$$

In Backpropagation update the weight.





In Backpropagation

try to update
the filter value.

And use different
filters detect different
information.

apply activation function.
(relu) field.



ANIN-CNN → operation almost same (process)

Just calculation is diff.

v_1



v_2



v_3



v_4



v_5



v_6



↓
output

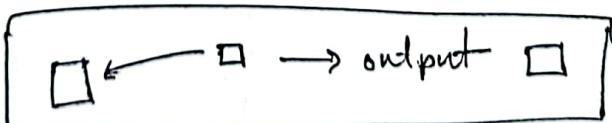
→ initial filter value

randomly selected

and in backpropagation

Kernel / filter are
updated.

↑ Convolution operations



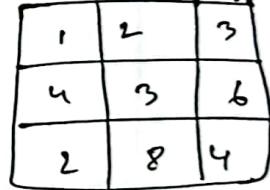
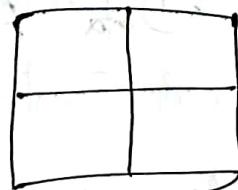
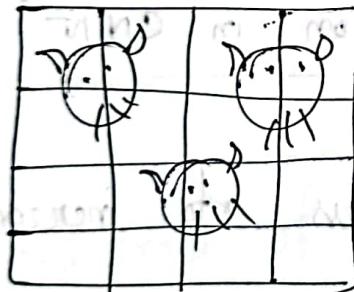
(Input) Input blocks consisting of 2×2 blocks

Filter consisting of 2×2 blocks

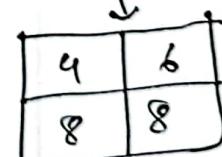
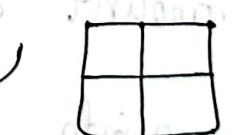
Term: location invariant

When human Brain see some faces
some of the neuron automatically
triggered

→ Convolution NN do that / done by
Max Pooling layer



Convolution operation



another filter
max pooling

Max pooling → clearly / more precisely

detect. (maximum intensity value)

(maximum intensity value)

→ Properly detected.

There is also

Average pooling

Original grid of values

Data Augmentation in CNN

Data augmentation is a process to increase the amount of data by generating new data points from existing data.

(Slightly modified)

Scale, crop, flip, rotate, transform
different formats.

"Synthetic data": When data is generated artificially without using real world images.
Produced by (GAN)

"Augmentation data": Derived from "original images" with some sort of minor geometric transformation like (flipping, translation, rotation, or the addition of noise) in order to increase the diversity of the training image.

"Importance of Data Augmentation"

- Improve the performance.
- Reduce overfitting.