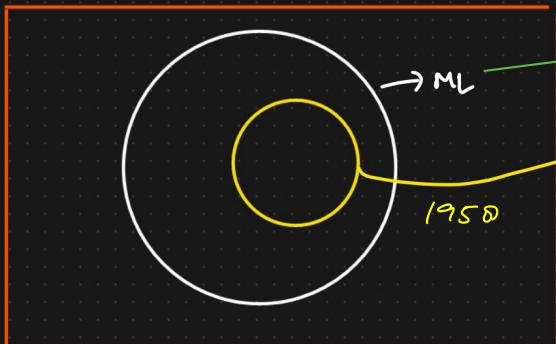


Deep learning → recognizes patterns in the data and then makes predictions based on Multilayered Neural Network

Imp. (mimics human brain)

AI ⇒ Artificial Intelligence



→ Stats tools to Analyse Data

Multilayered
⇒ Neural Networks

Mimic the Human
Brain

ML - CPU
DL - GPU

Deep learning → huge amounts of Data

- ① ANN → Artificial Neural N/w ↗ Classification }
 - ② CNN → Convolutional Neural N/w → I/P : Image, Video → RCNN, MASKED RCNN, frames
Detection, YOLO vs V6, V7
 - ③ RNN → Recurrent Neural N/w → NLP ↗ NLP, Time Series
- Sequential Data ↗ I/P: Text, Time Series

Computer Vision
Object Detection
↑

FRAMEWORK

TENSORFLOW Keras

End to End Project

{ Word Embedding, LSTM RNN, GRU RNN,
Bidirectional LSTM RNN, Encoder Decoder,
Transformers, BERT }

LSTM = Long Short Term Memory → precise & slow

GRU = Gated Recurrent Unit → simple & fast

Tensorflow → Deep Learning Library
by Google

↳ Keras → API on top of Tensorflow

PyTorch → Deep Learning Library
by Meta

② Why Deep Learning Is Becoming popular?

2005 → Facebook, YouTube, WhatsApp, LinkedIn, Twitter } ⇒ Social Media platform
↓

2011-2012

DATA WAS GENERATED

Image, Text, Documents,

Exponentially

Videos

Big DATA → HADOOP → unstructured

DATA

↓
[2011] ⇒ Cloudera, Hortonworks

Netflix → movie streaming

2011-2012 :

① Hardware Requirement → GPU's cost ↓ ↓ → Nvidia GPU

TRAIN MODELS

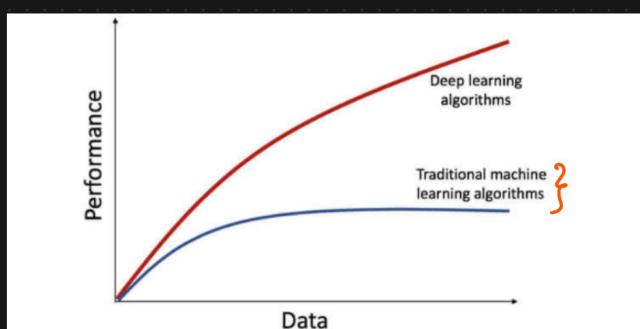
Price ↓

Recommendation System

↓
Increase Revenue

② Huge amount of data is getting generated → Deep learning

Models Perform well



③ Deep learning is been used in Many Domains

① Medical

② E-commerce

③ Retail

④ Marketing . . .

⑤ Frameworks Open Source

Community size ↑

Tensorflow
↓
Google

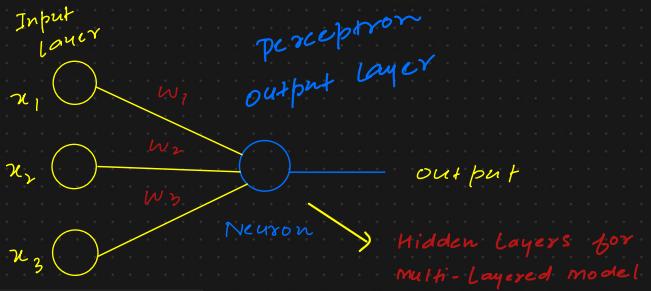
Pytorch
↓
Facebook

More Research

Perceptron; Simplest Model that Mimics Brain
Simplest form of neural network Unit.

🧠 How It Works:

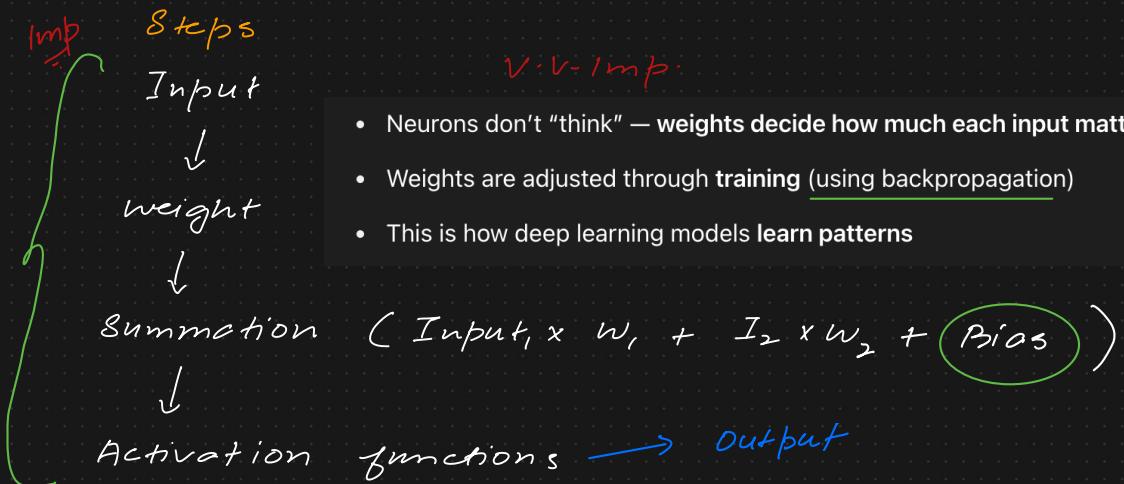
1. Inputs → Like features (e.g., height, weight).
 2. Weights → Importance of each input.
 3. Summation → Adds the weighted inputs.
 4. Activation Function → Decides the output (e.g., 0 or 1).
- 💡 Output = Activation((Input₁ × Weight₁) + (Input₂ × Weight₂) + ... + Bias)



Imp.

🧠 Quick Recap:

Type	Hidden Layer?	Example
Single-layer NN	✗ No	Perceptron
Multi-layer NN (MLP)	✓ Yes	Deep Learning models



V: V-1mp:

- Neurons don't "think" — weights decide how much each input matters
- Weights are adjusted through training (using backpropagation)
- This is how deep learning models learn patterns

So, the entire process is called Forward Propagation

Actual Equation

$$y = \sum x_i w_i + \text{bias} \Leftrightarrow mx + c \quad \text{think of it as bias.}$$

Important so we start from some default value and not just zero

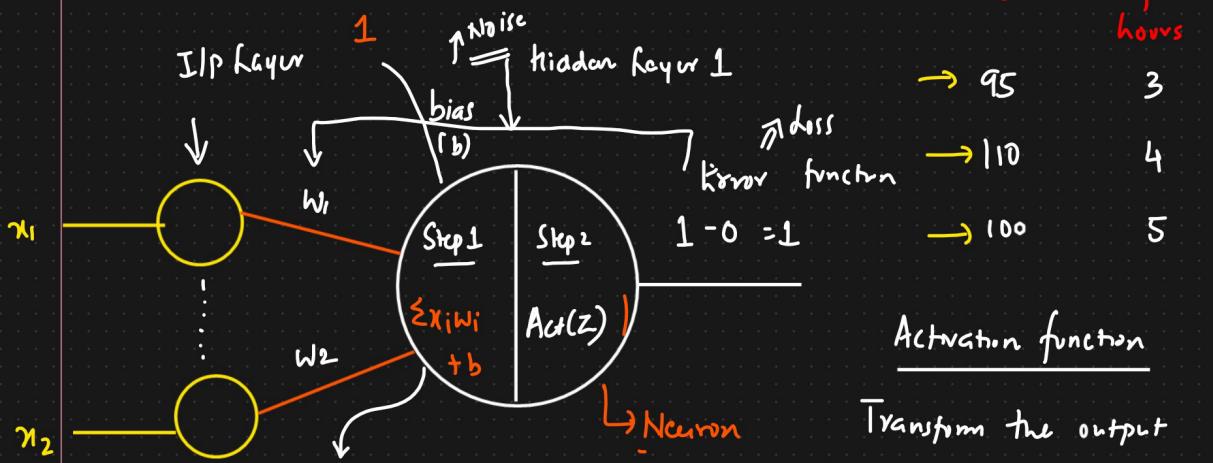
③ Perceptron [Artificial Neuron or Neural Network Unit]

① Input layer ✓
 → no hidden layer [Single layered NN]

② Hidden layer ✓ in perceptron ↓

③ Weights ✓ Binary classifier

④ Activation function ✓

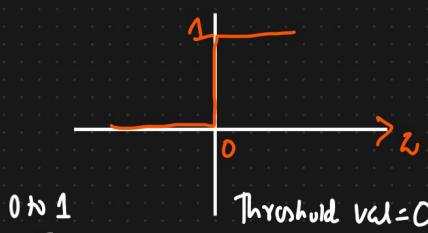


Imp. $z = w_1 x_1 + w_2 x_2 + b \rightarrow \beta_0 x_0 + \beta_1 x_1$

$$z = \sum_{i=1}^n x_i w_i + b$$

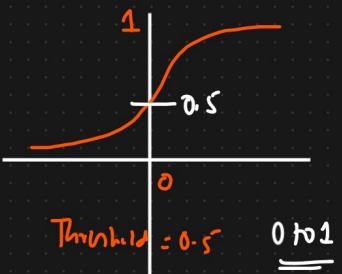
8 similar

Step Function

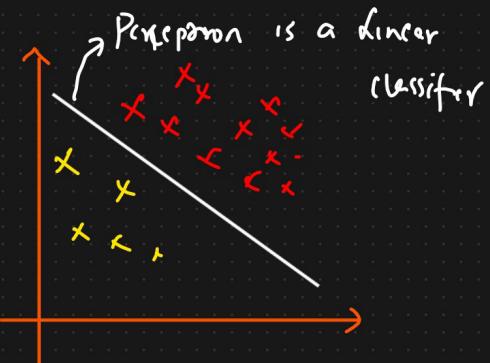


$$\begin{cases} 0 & z \leq 0 \\ 1 & z > 0 \end{cases}$$

Sigmoid Function



$$\begin{cases} 1 & z > 0.5 \\ 0 & z \leq 0.5 \end{cases}$$



Step 1

$$z = \sum_{i=1}^n w_i x_i + b \rightarrow \text{Bias}$$

$$z = b + w_1 x_1 + w_2 x_2 + w_3 x_3 + \dots + w_n x_n$$

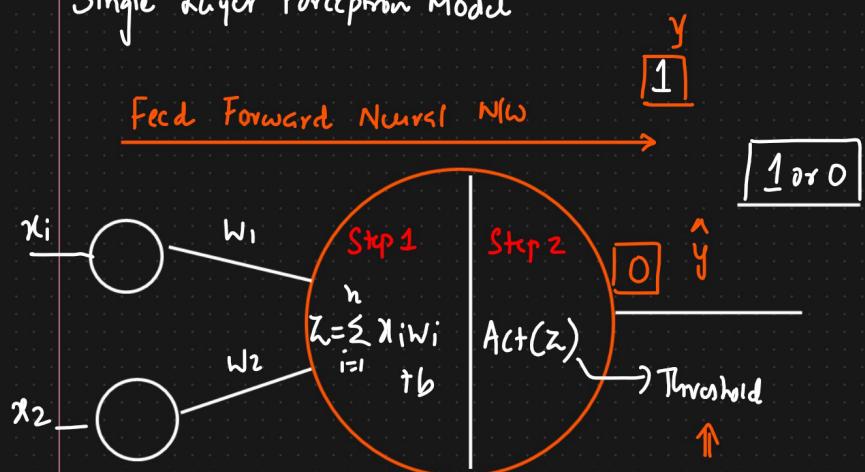
$$y = mx + c$$

$$y = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \beta_3 x_3 + \dots + \beta_n x_n$$

↓
linear
problem
statement

Perceptron Models

Single layer Perceptron Model

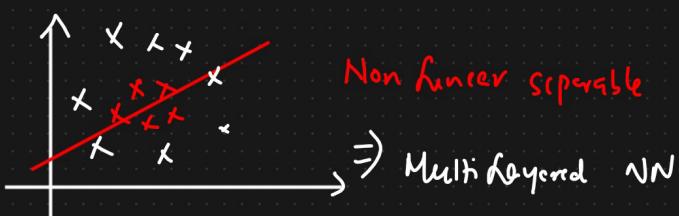
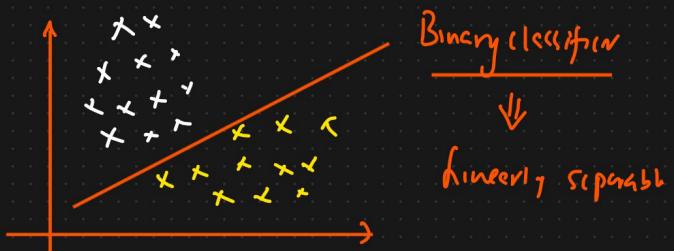


↓
[ANN]

Multi Layered Perceptron Model

- ① Forward Propagation
- ② Backward "
- ③ Loss functions
- ④ Activation functions
- ⑤ Optimizers

Linear Separable problem



Activation function

→ $\frac{1}{e^{-y}} = \frac{1}{e^{-(w_0 + w_1 x_1 + w_2 x_2 + \dots + w_n x_n)}}$ → gives output in 0 to 1.

$\geq 0.5 = 1$
$< 0.5 = 0$

Summary

Deep learning → mimic human brain with the help of neural networks

Working ;

Input Data → calculate weights → Add weights + Bias (Sum)

Output ← Activation functions
to get probability

$$\sum x_i w_i + b$$

- Pros / Cons →
 - More Accurate
 - Solves complex problems
 - work on huge Data
 - Requires GPU's
 - More Time
 - Blackbox (Hidden layers)

Neural Networks

Linear classification

perceptron



LSTM

Single layered
Neural Network

- single layer

- no hidden layers

- only linear tasks

- And / OR Gate

- Layer = 1

multi. layered

Neural Network

- multiple layers

- hidden layers

- complex & non linear

- Image classifier / NLP

- Layers ≥ 2

Qn. What is non-linearity concept? Usually Activation functions ability of a model to learn complex patterns; like curves, shapes etc. not just simple linear relationships.

(imp)

In short:

Non-linearity gives neural networks the power to learn complex, real-world relationships.

Without it, even a deep network would just behave like a linear regression model.

Qn. What do you understand by activation functions?

decides whether neuron should be activated.

- introduces probability & non-linearity

- helps model to learn beyond linear data

Activation Function Comparison Cheat Sheet

Activation	Used For	Non-Linearity	Probability Output	Common Use Case
ReLU	Hidden Layers (default)	Yes	No	Enables deep networks to learn patterns
Leaky ReLU	Hidden Layers (ReLU alternative)	Yes	No	Solves "dead ReLU" problem
Sigmoid	Output Layer (Binary Classification)	Yes	Yes	Converts output to probability (0-1)
Softmax	Output Layer (Multi-class Classification)	Yes	Yes	Distributes output as class probabilities
Tanh	Sometimes in Hidden Layers	Yes	No	Like Sigmoid but zero-centered

Forward & Backward Propagation

Forward Propagation (Prediction Phase)

♦ Purpose To compute the output (prediction) of the model.

♦ Direction From input → hidden layers → output

♦ What happens?

- Input data is passed through the network.
- At each layer:
 - Weighted sum is calculated: $Z = W \cdot X + b$
 - Activation function is applied: $A = \text{activation}(Z)$
- Final output (\hat{y}) is produced.

| ♦ Example | Predicting whether an email is spam or not. |

Backward Propagation (Learning Phase)

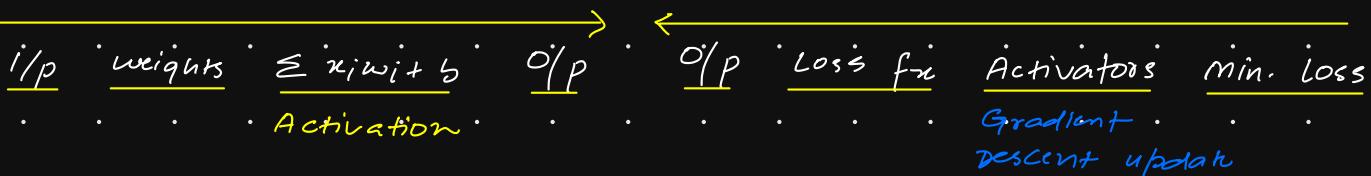
♦ Purpose To update the weights to reduce error/loss.

♦ Direction From output → hidden layers → input (reverse)

♦ What happens?

- Compute the loss between predicted output and actual output.
- Apply chain rule of calculus (using gradients):
 - Compute how much each weight contributed to the error.
- Calculate gradients (dW , db) and update weights using gradient descent.

| ♦ Example | Adjusting model so it better predicts spam in future. |



→ What is ReLU Rectified Linear Unit

✓ ReLU Meaning (In Simple Terms)

- If the input is positive, keep it as it is.
- If the input is negative, change it to zero.

👉 In short:

$$\text{ReLU}(x) = x \text{ if } x > 0, \text{ else } 0$$

⚠ Limitation: → Important

- "Dead Neurons": If a neuron always receives negative inputs, it outputs 0 forever and stops learning.
- To solve this, we use Leaky ReLU sometimes.

Summary:

ReLU = Let the positive values pass, block the negative ones.
It's simple, fast, and the default choice for hidden layers in deep learning.

Q: what is leaky ReLU? Leaks negative values

Keeps neurons alive

● What is Leaky ReLU?

Leaky ReLU is a modified version of ReLU that fixes its main problem:

It allows a small, non-zero output when the input is negative.

⚠ Why do we need it?

Problem with ReLU:

- ReLU outputs 0 for all negative inputs.
- Some neurons may get stuck and never activate or learn — called the "dead ReLU" problem.

✓ What Leaky ReLU does:

- For positive inputs: same as ReLU \rightarrow output = input.
- For negative inputs: output a small negative value (like $0.01 \times$ input) instead of 0.

🧠 Intuition:

Input (x)	ReLU(x)	Leaky ReLU(x)
-3	0	-0.03
0	0	0
2	2	2

✓ Benefits of Leaky ReLU:

Feature

Keeps neurons alive

Allows small gradients

Works better for deep networks

Q: what is Tanh hyperbolic Tangent

● What is tanh in Deep Learning?

tanh, short for hyperbolic tangent, is an activation function used in neural networks — especially in hidden layers.

✓ Key Characteristics:

Property	Description
Output Range	-1 to 1
Shape	S-shaped (like sigmoid, but zero-centered)
Non-linearity	✓ Yes — allows the network to learn complex patterns
Used In	Hidden layers (especially in RNNs or older networks)

🧠 Why use tanh?

- Unlike sigmoid (which outputs between 0 and 1), tanh is centered around 0.
- This makes optimization easier and faster, especially for data with negative values.

-1 to 1

unlike softmax

0 to 1

Conclusion;

1. I/P layers
 2. weight
 3. $\sum x_i w_i + \text{Bias}$
 4. Activation function
 5. Loss function ($\hat{y} - y$)
 6. minimize loss with optimizers
 7. update weights
- } Forward propagation
- } Backward propagation

Imp.

Machine learning = white Box
(except Random Forest)

Deep learning = Black Box models

5 Multi layered Perceptron Model [Artificial Neural N/w]

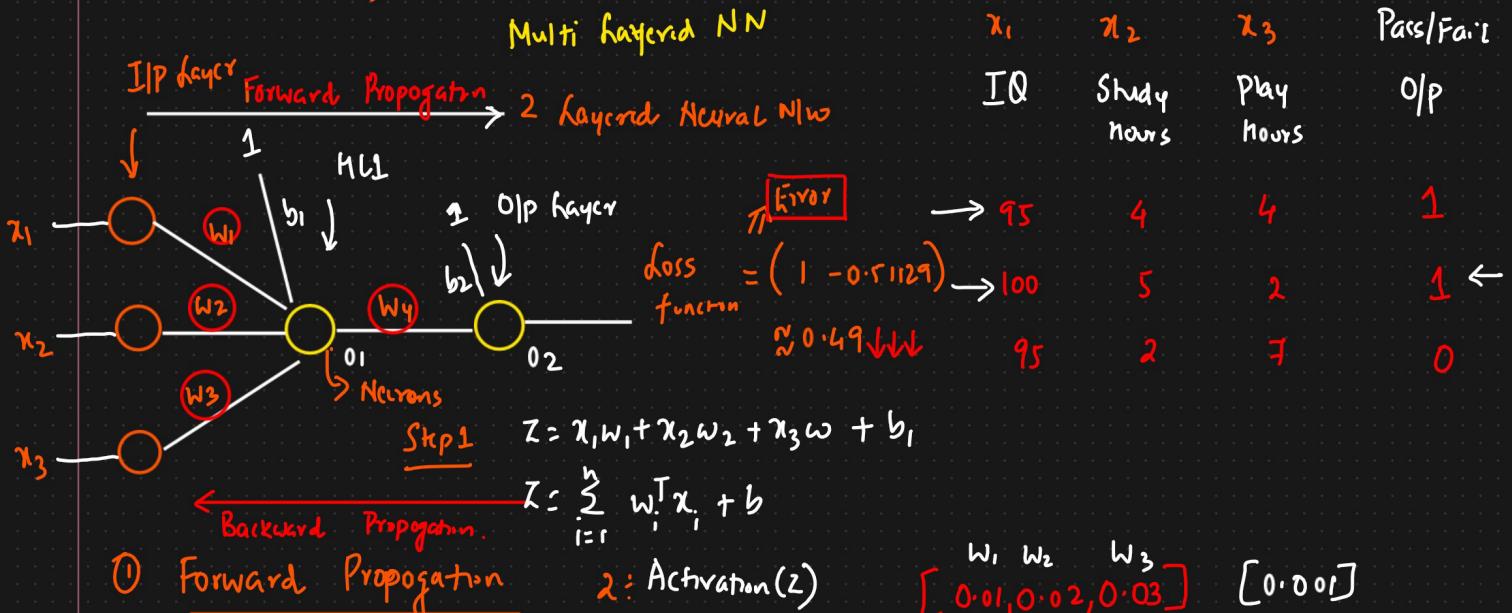
① Forward Propogation

② Backward Propogation → Geoffrey Hinton →

③ Loss function

④ Optimizers ✓ Gradient Descent.

⑤ Activation function.



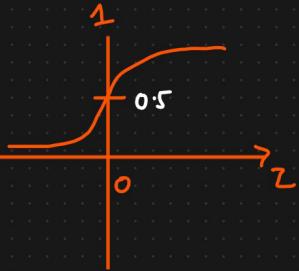
Hidden Layer 1

$$\begin{aligned} \text{Step 1: } Z &= 95 \times 0.01 + 4 \times 0.02 + 4 \times 0.03 + 1 \times 0.01 \\ &= 1.151 \end{aligned}$$

Sigmoid. $= \frac{1}{1+e^{-Z}}$

Step 2: Activation (Z)

$$f(Z) = \frac{1}{1+e^{-1.151}} = 0.759$$



$$0_1 = 0.759$$

Hidden Layer 2

$$w_4 = 0.02$$

$$b_2 = 0.03$$

$$\text{Step 1: } Z = 0_1 * w_4 + b_2$$

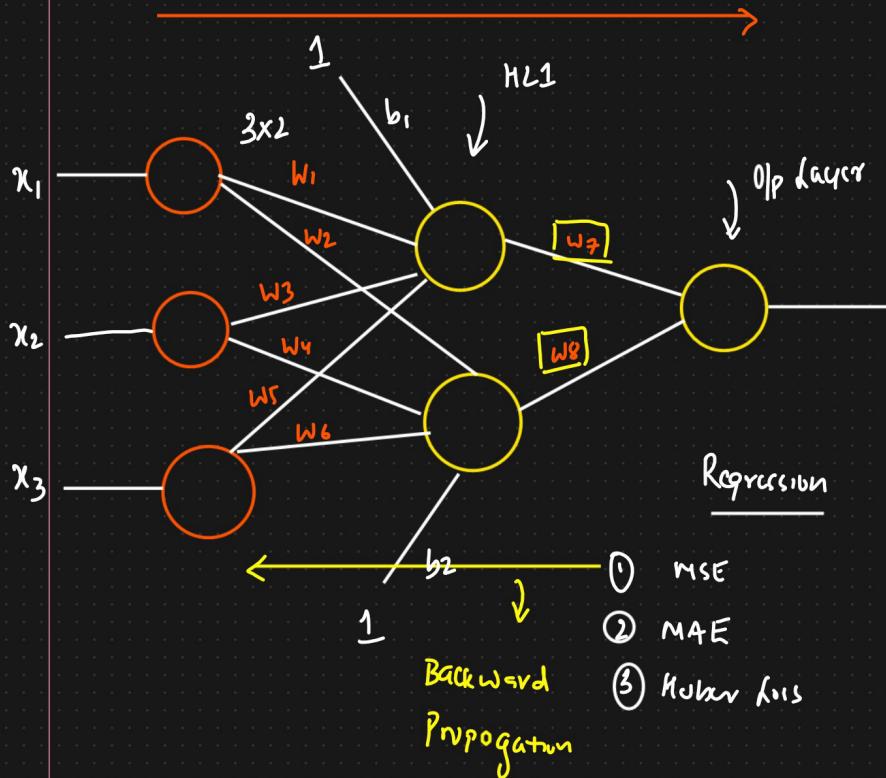
$$= 0.759 * 0.02 + 0.03$$

$$= 0.04518$$

Step 2 : Activation(z) $\frac{1}{1+e^{-(0.04518)}} = 0.51129$

$$O_2 = 0.51129 \Rightarrow \hat{y}$$

⑥ Back Propagation And Weight Updation Formula



Weight update Formula

$$w_{7\text{new}} = w_{7\text{old}} - \eta \left[\frac{\partial L}{\partial w_{7\text{old}}} \right]$$

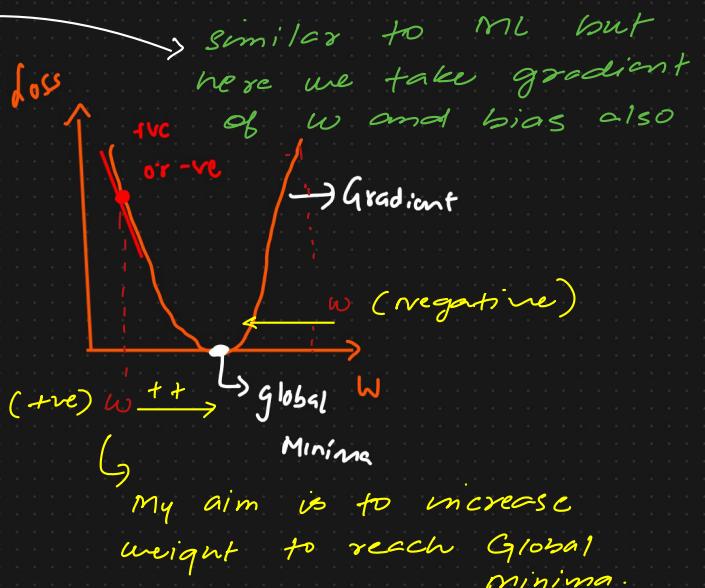
slope

$$w_{8\text{new}} = w_{8\text{old}} - \eta \left[\frac{\partial L}{\partial w_{8\text{old}}} \right]$$

learning Rate

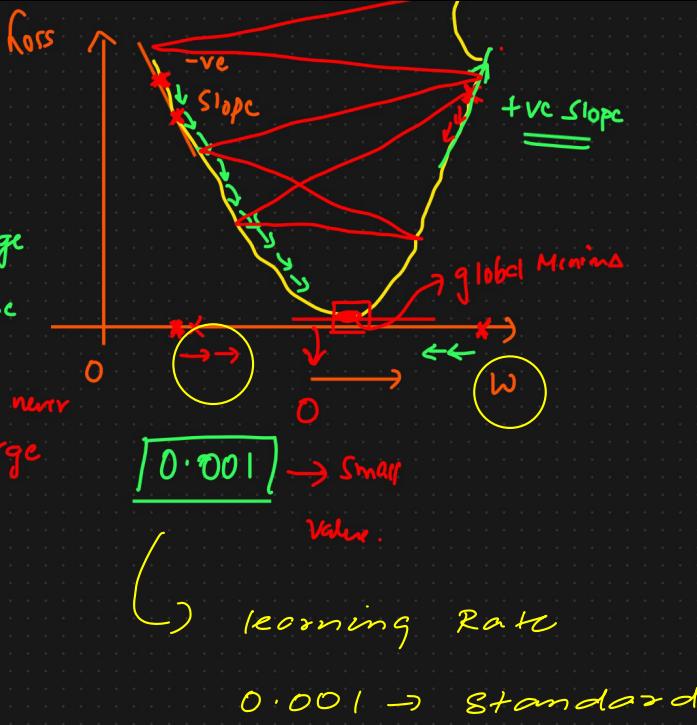
V.V. Imp

$$w_{\text{new}} = w_{\text{old}} - \eta \left[\frac{\partial L}{\partial w_{\text{old}}} \right] \Rightarrow \text{Weight Updation Formula.}$$



Gradient Descent
Optimizers

Optimizers : To reduce the loss value



$$W_{\text{new}} = W_{\text{old}} - \eta (-\text{ve})$$

$$= W_{\text{old}} + \eta (+\text{ve})$$

$$\boxed{W_{\text{new}} > W_{\text{old}}}$$

η = large value
↓
It may never converge

Learning Rate

$$\boxed{0.001} \rightarrow \text{Small value.}$$

) learning Rate

$0.001 \rightarrow \text{standard.}$

$$W_{\text{new}} = W_{\text{old}} - \eta (+\text{ve})$$

$$= W_{\text{old}} - \eta (+\text{ve})$$

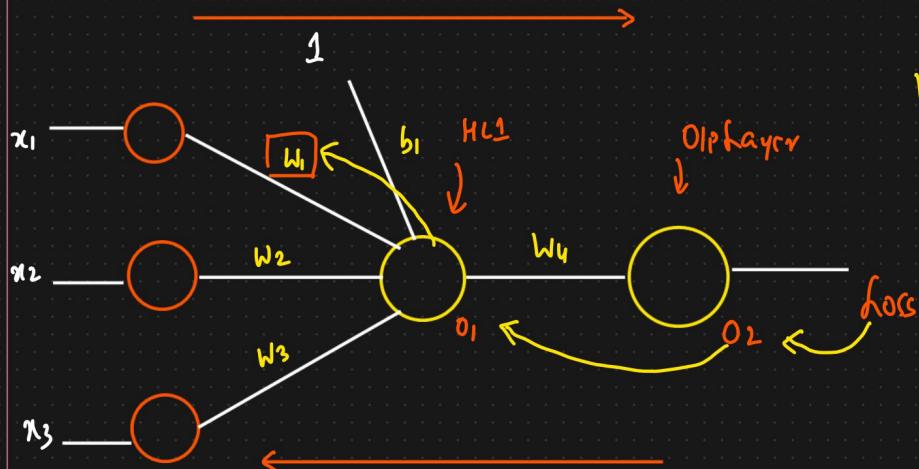
$$\boxed{W_{\text{new}} < W_{\text{old}}}$$

When w reaches Global Minima

$$\boxed{W_{\text{new}} = W_{\text{old}}}$$

→ used to find derivatives of composite functions (function inside a function).

⑦ Chain Rule of Derivative



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

$$W_{4\text{new}} = W_{4\text{old}} - \eta \boxed{\frac{\partial h}{\partial W_{4\text{old}}}}$$

$$\frac{\partial h}{\partial w_{i,old}} = \frac{\partial h}{\partial o_2} * \frac{\partial o_2}{\partial w_{i,old}} \Rightarrow \text{Chain Rule of Derivation}$$

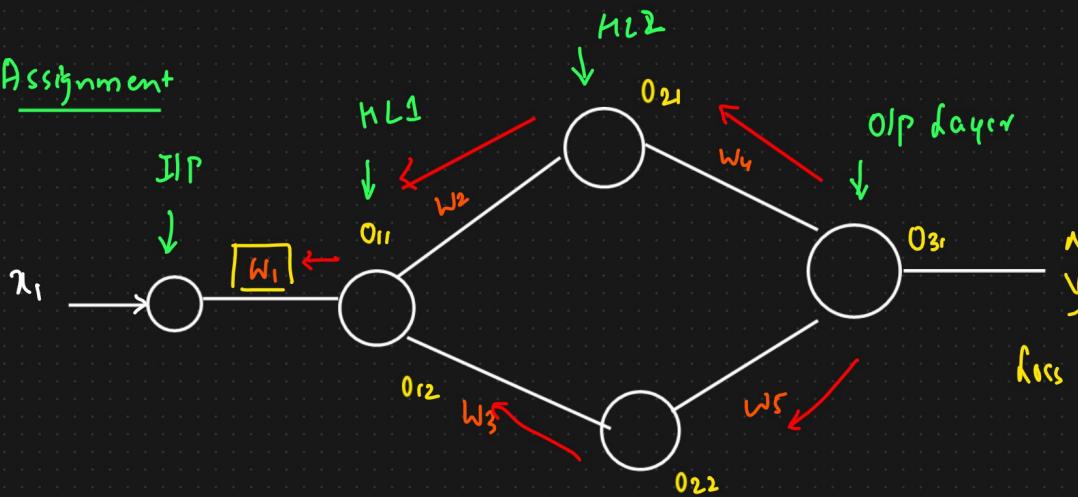
$$w_{i,new} = w_{i,old} - \eta \left[\frac{\partial h}{\partial w_{i,old}} \right]$$

$$\frac{\partial h}{\partial w_{i,old}} = \frac{\partial h}{\partial o_2} * \frac{\partial o_2}{\partial o_1} * \frac{\partial o_1}{\partial w_{i,old}}$$

$w_{2,new}$

$w_{3,new}$

Assignment



$$w_{i,new} = w_{i,old} - \eta \left[\frac{\partial h}{\partial w_{i,old}} \right]$$



$$\frac{\partial h}{\partial w_{i,old}} = \left[\frac{\partial h}{\partial o_{31}} * \frac{\partial o_{31}}{\partial o_{21}} * \frac{\partial o_{21}}{\partial o_{11}} * \frac{\partial o_{11}}{\partial w_{i,old}} \right]$$

+

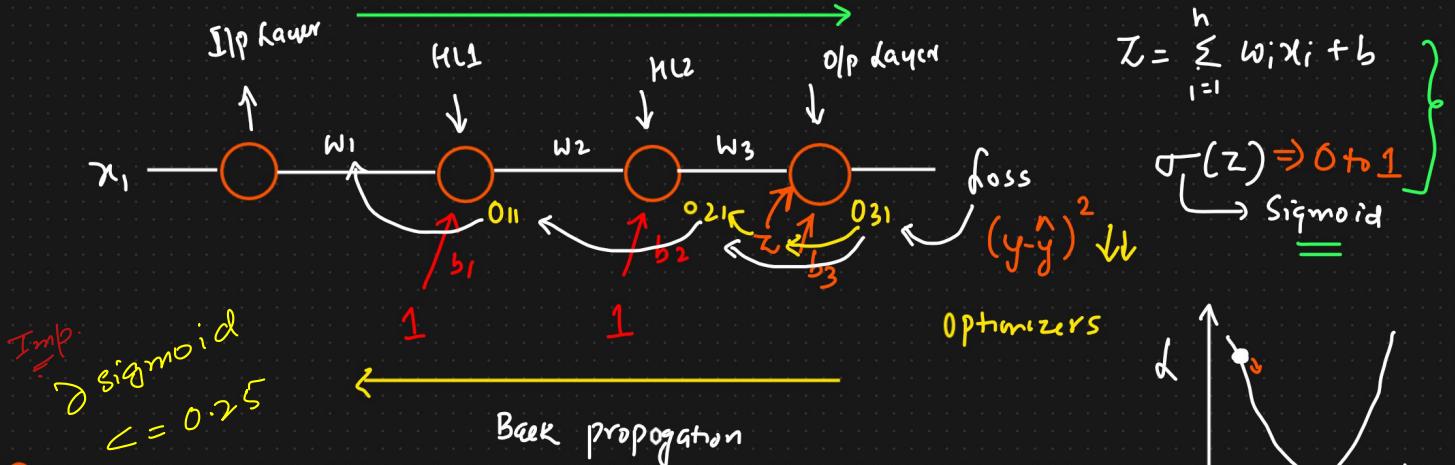
$$\left[\frac{\partial h}{\partial o_{31}} * \frac{\partial o_{31}}{\partial o_{22}} * \frac{\partial o_{22}}{\partial o_{12}} * \frac{\partial o_{12}}{\partial w_{i,old}} \right] \quad \begin{matrix} \text{dh} \rightarrow \frac{\partial h}{\partial w_{i,old}} \\ (\text{Chain Rule}) \end{matrix}$$

vv imp
==

Backpropagation = Chain Rule Repeated Backward

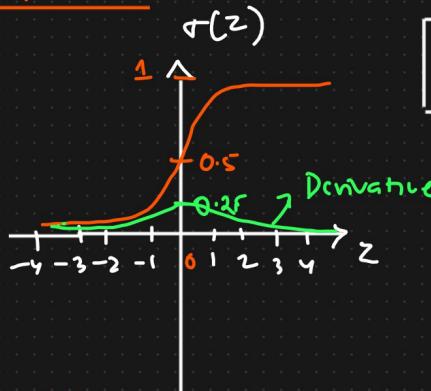
- Output Layer: Compute how loss changes with output (easy)
- Then apply chain rule to pass gradient back through each layer
- Eventually, you get the gradient of loss with respect to every weight

(Solution = Change Activation Function) ⑧ Vanishing Gradient Problem And Activation functions



⑨ Sigmoid Activation function

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



$$0 \leq \sigma(z) \leq 1$$

Vanishing
Gradient
Problem.

$$\text{Derivative } (\sigma(z))$$

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

$$w_{1,\text{new}} = w_{1,\text{old}} - \eta \left[\frac{\partial h}{\partial w_{1,\text{old}}} \right] \Rightarrow \text{small value} \Rightarrow \boxed{w_{1,\text{new}} \approx w_{1,\text{old}}}$$

$$\frac{\partial h}{\partial w_{1,\text{old}}} = \frac{\partial h}{\partial o_{31}} * \frac{\partial o_{31}}{\partial o_{21}} * \frac{\partial o_{21}}{\partial o_{11}} * \frac{\partial o_{11}}{\partial w_{11}}$$

$$o_{31} = \sigma \left(\underbrace{w_3 * o_{21}}_z + b_3 \right) \quad z = w_3 * o_{21} + b_3$$

$$o_{31} = \sigma(z)$$

$$\frac{\partial o_{31}}{\partial o_{21}} = \frac{\partial (\sigma(z))}{\partial (z)} * \frac{\partial z}{\partial o_{21}} \quad \left\{ \text{Chain Rule} \right\}$$

$$\boxed{0 \leq \sigma(z) \leq 0.25} \quad * \quad \frac{\partial((w_3 + b_3) + b_3)}{\partial(w_2)}$$

↓

Derivative of
Sigmoid

$$\boxed{\frac{\partial o_3}{\partial o_2} = 0 \leq \sigma(z) \leq 0.25 * w_3_{\text{old}}}$$

This causes
loss of Gradient
which means, we can't
have updated values

- ① To fix this problem Researchers started exploring other
Activation function
 In this case, we will use Tanh instead of Sigmoid function.
- ④ Tanh ④ ReLU ④ PReLU ④ Swish

Conclusion;

- Perceptron → No chain Rule
- multi Layered Neural Network ; we do use chain Rule of derivatives to update all the weights (composite fns).

Steps ;

1. Inputs
2. Weights
3. $\sum w_i x_i + \text{Bias}$
4. Hidden Layers
5. Activation functions
6. Loss functions ($y - \hat{y}$)
7. calculate derivatives for w/B
8. Update $w_{\text{new}} = w - \eta \frac{\partial h}{\partial w_{\text{old}}}$
 $B_{\text{new}} = B - \eta \frac{\partial h}{\partial B_{\text{old}}}$

Derivative		
Activation	Derivative Range	Max Derivative
ReLU	0 or 1	✓ 1
Leaky ReLU	Small positive to 1	✓ 1
Sigmoid	0 to 0.25	✓ 0.25
Tanh	0 to 1	✓ 1
Softmax	0 to 0.25 (per class)	✓ 0.25

Imp: Sigmoid value derivative ≤ 0.25 always

Q11 Difference b/w Sigmoid and Tanh

Tanh = 0 centered, -1 to 1, hyperbolic

Sigmoid = Not 0 centered, 0 to 1, sigmoid

V.V Imp Derivative value for different Gradients.

Activation	Derivative Range	Max Derivative
ReLU	0 or 1	✓ 1
Leaky ReLU	Small positive to 1	✓ 1
Sigmoid	0 to 0.25	✓ 0.25
Tanh	0 to 1	✓ 1
Softmax	0 to 0.25 (per class)	✓ 0.25

Gradient Vanishing problem since low numbers.

Remember → Gradient vanishing can still occur in Tanh

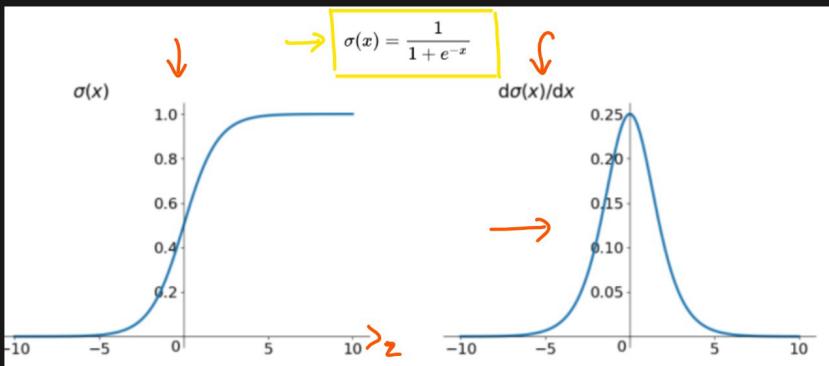
Evolution of ACTIVATION FUNCTIONS

1. Sigmoid / softmax
2. Tanh
- Best 3. ReLU → one problem
↳ Dead Neuron when Negative
4. Leaky ReLU → solves problem of Dead Neurons, by leaking Negative values.

Activation Functions

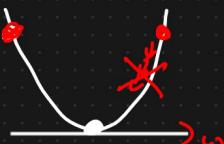
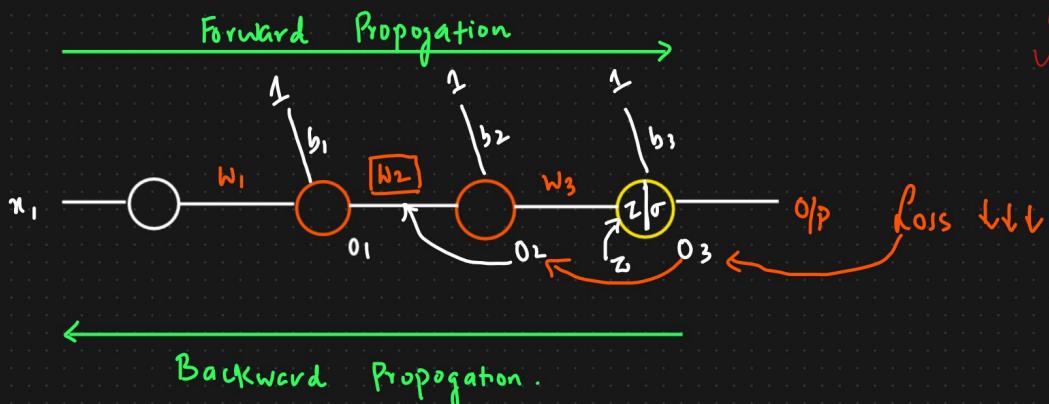
① Sigmoid Activation function [0 to 1]

$$z = \sum_{i=1}^n w^T x + b$$



$$\sigma(z) \Rightarrow 0 \leftrightarrow 1$$
$$\phi(z) \Rightarrow$$

$$\frac{\partial \sigma(z)}{\partial z} = 0 + 0.25$$



$$w_{\text{new}} = w_{\text{old}} - \eta \left[\frac{\partial h}{\partial w_{\text{old}}} \right] \quad \Rightarrow \quad w_{\text{new}} \approx w_{\text{old}}$$

$$\frac{\partial h}{\partial w_{201d}} = \frac{\partial h}{\partial o_3} \neq \boxed{\frac{\partial o_3}{\partial o_2}} \neq \frac{\partial o_2}{\partial w_2}$$

0.20  $* 0.01$ \times $\det z = (o_2 * w_3) + b_3$

$$\frac{\partial \theta_3}{\partial \theta_2} = \frac{\partial (\sigma(z))}{\partial z} * \frac{\partial z}{\partial \theta_2}$$

$$= [0 - 0.25] * \frac{\partial [(0_2 * \omega_3) + b_3]}{\partial 0_2}$$

$$= \left[0 - 0.25 \right] \neq w_3 \Rightarrow \text{small value} \Rightarrow$$

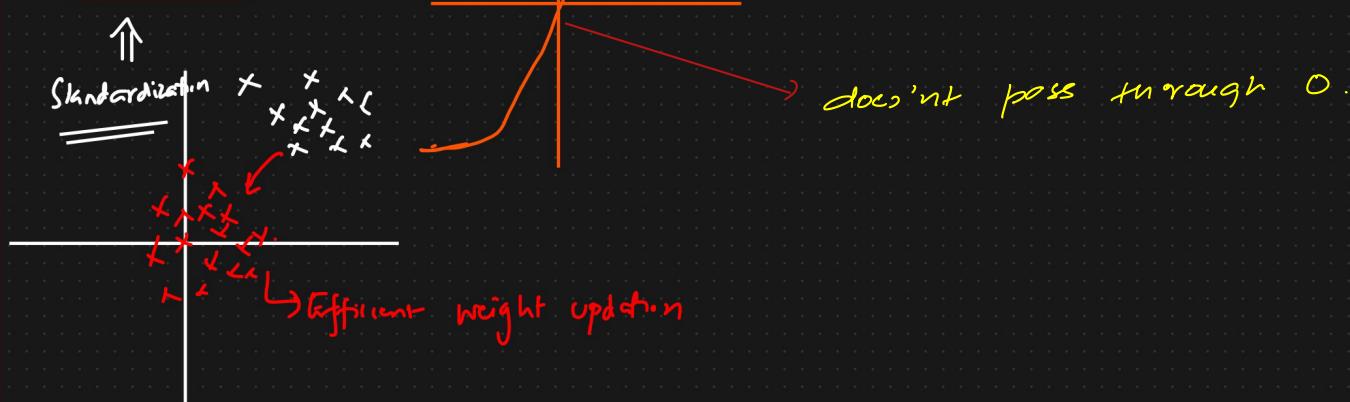
Advantages

- ① Binary classification Suitable.
- ② Clear prediction i.e. very close 1 or 0

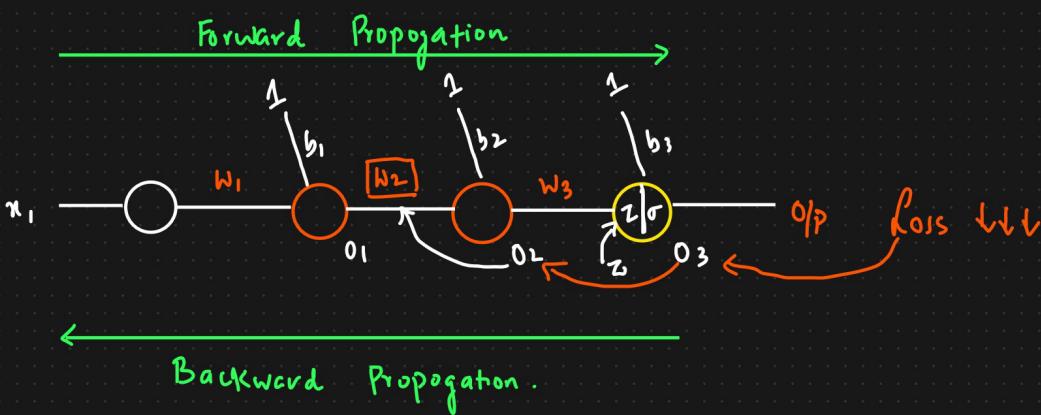
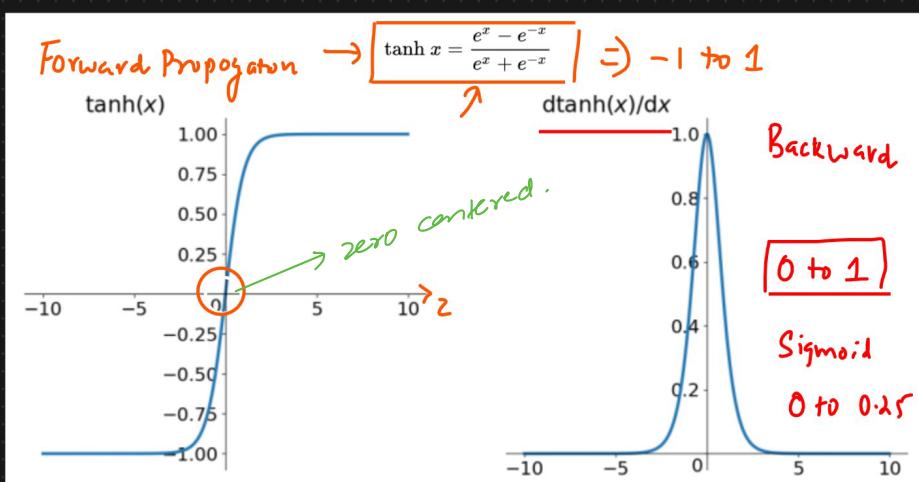
Disadvantages

- ① Prone to vanishing Gradient Problem.
- ② Function output is not Zero centered \Rightarrow Efficient weight update
- ③ Mathematical operation are relatively time consuming

Zero centered



② Tanh Activation Function



$$\frac{\partial h}{\partial w_{2,0,1,0}} = \frac{\partial h}{\partial o_3} \neq \boxed{\frac{\partial o_3}{\partial o_2}} \neq \frac{\partial o_2}{\partial w_2}$$

$\downarrow \downarrow \downarrow$

$$0 \cdot 20_{11} * 0 \cdot 01 \times \text{det } z = (o_2 * w_3) + b_3$$

$$\frac{\partial o_3}{\partial o_2} = \boxed{\frac{\partial (\tanh(z))}{\partial z} * \frac{\partial z}{\partial o_2}} [0 \rightarrow 1]$$

$$= [0 - 1] * \frac{\partial [(o_2 * w_3) + b_3]}{\partial o_2}$$

$$= [0 - 1] * w_3 \Rightarrow \text{Small value} \Rightarrow$$

Advantages

① Zero Centric \Rightarrow Weight Updation is Efficient

① Prone to Vanishing Gradient Problem

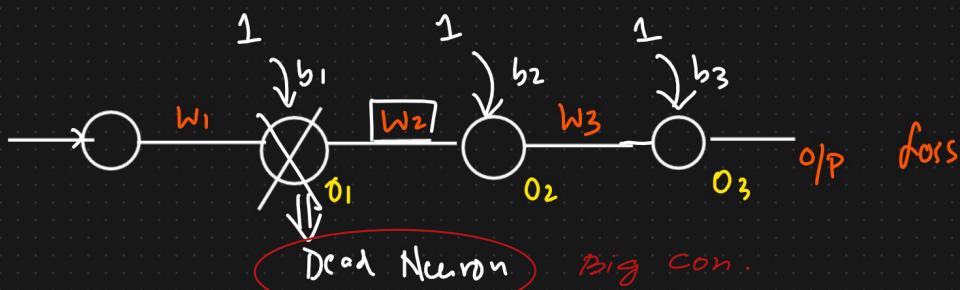
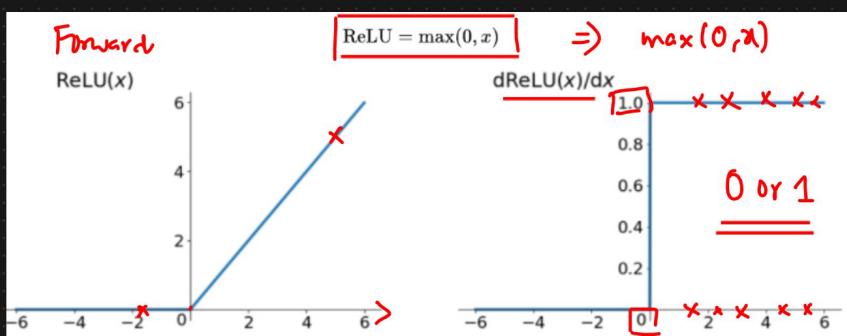
② Time Complexity

[Rectified Linear Unit].

③ ReLU Activation Function

Tanh $\Rightarrow 0 \rightarrow 1$

Sigmoid $\Rightarrow 0 \rightarrow 0.25$



$$\frac{\partial h}{\partial w_{201d}} = \frac{\partial h}{\partial o_3} * \boxed{\frac{\partial o_3}{\partial o_2}} * \frac{\partial o_2}{\partial w_2}$$

↓ ↓ ↓ ↓

$$0.20_{11} * 0.01 * \text{det } z = (o_2 * w_3) + b_3$$

$$\frac{\partial o_3}{\partial o_2} = \boxed{\frac{\partial (\text{relu}(z))}{\partial z}} * \frac{\partial z}{\partial o_2}$$

[0 to 1]

$$= [0 \text{ or } 1] * \frac{\partial [(o_2 * w_3) + b_3]}{\partial o_2}$$

$$= \begin{bmatrix} -ve & +ve \\ 0 \text{ or } 1 \end{bmatrix} * w_3 \Rightarrow \text{Small value} \Rightarrow$$

If Derivative of ReLU output is $\boxed{1}$ \Rightarrow Weight updation will happen

If ReLU output is $\boxed{0}$ \Rightarrow Dead Neuron

$$w_{2\text{new}} = w_{201d} - \eta \boxed{\frac{\partial h}{\partial w_{201d}}} \Rightarrow 0$$

If Derivative of $\text{ReLU}(z)$ is 0

$$\boxed{w_{2\text{new}} \approx w_{201d}} \Rightarrow \text{Dead Neuron}$$

If $\underline{z = +ve}$ $\frac{\partial \text{ReLU}(z)}{\partial z} = 1$

If $\underline{z = -ve}$ $\frac{\partial \text{ReLU}(z)}{\partial z} = 0$

Advantages

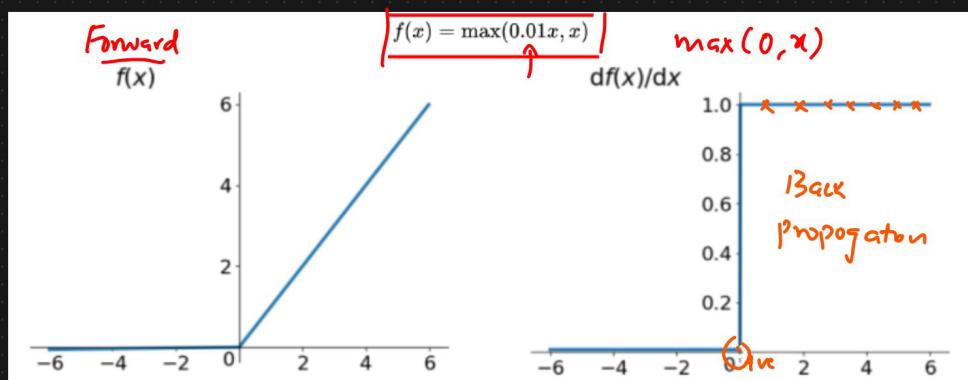
- ① Solving Vanishing Gradient Problem
- ② $\text{Max}(0, x) \rightarrow$ Calculation is Superfast. The ReLU function has a linear relationship.
- ③ It is much faster than Sigmoid or Tanh.

Disadvantages

- ① Dead Neuron
- ② ReLU function is $0/P$
 $(0, x) \Rightarrow 0$ or zero number
 \downarrow
 It is not zero centric

④ Leaky ReLU And Parametric ReLU

$\max(\lambda x, x)$ \rightarrow hyperparameter $\lambda = \alpha = 0.01, 0.02, 0.03$



ReLU \rightarrow Dead Neuron \rightarrow Dead ReLU Problem

Advantages

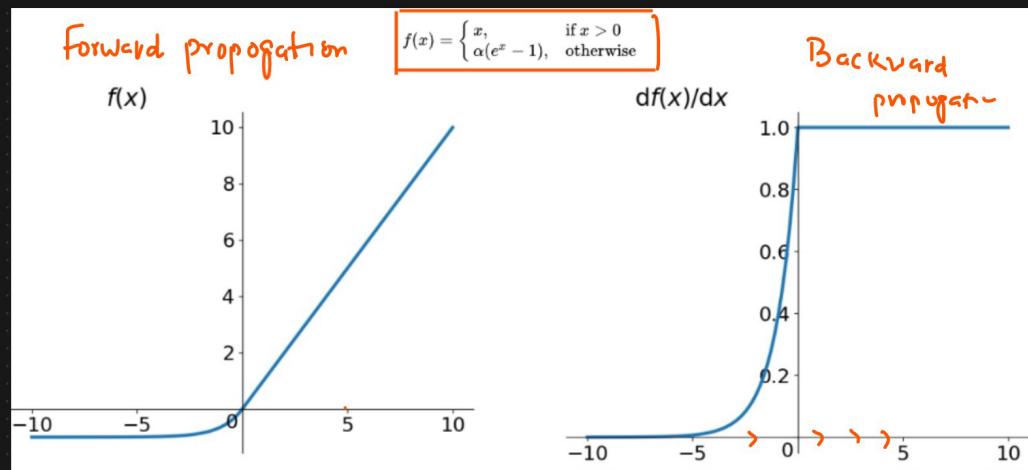
- ① Leaky ReLU has all the advantages of ReLU.
- ② It removes the Dead ReLU Problem.

Disadvantage

- ① It is not zero centric

⑤ ELU (Exponential Linear Units)

$$\alpha(e^x - 1).$$



Advantages

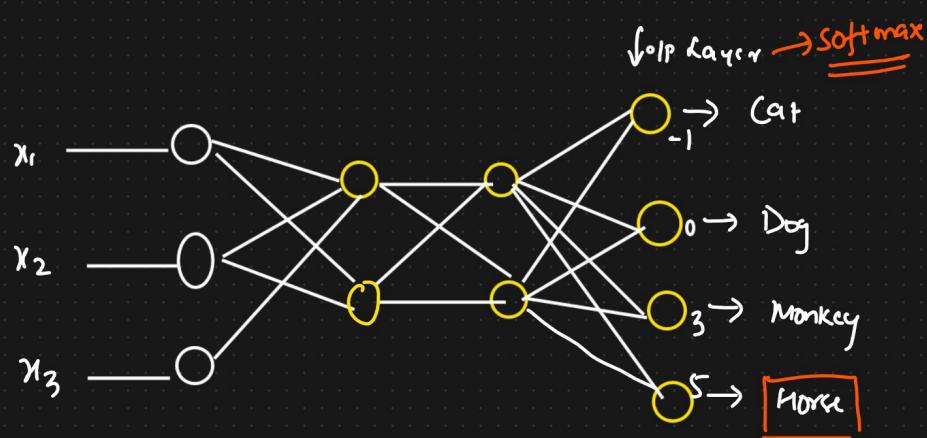
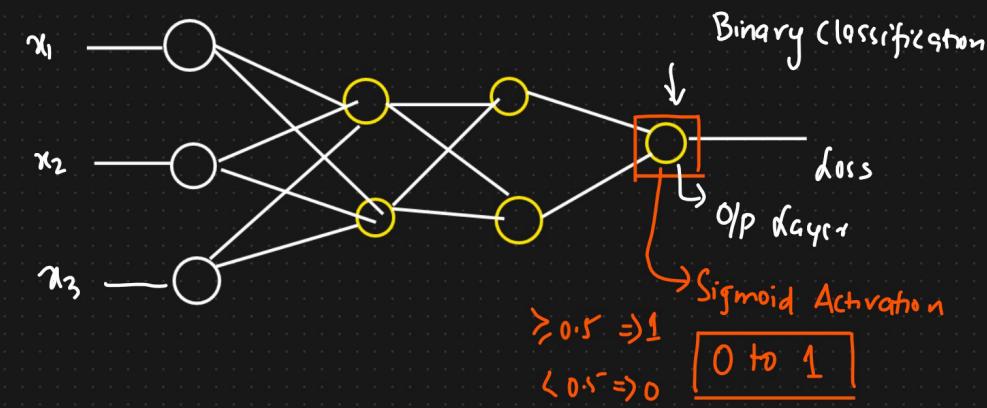
[It is used to solve
ReLU problems]

Disadvantage

- ① No Dead ReLU Issues
- ② Zero centered

i) Slightly more computationally intensive.

⑥ Softmax Activation function [Multiclass classification problem]



Softmax Activation

$$\text{Softmax} = \frac{e^{y_i}}{\sum_{k=0}^n e^{y_k}}$$

$$y_i = \mathbf{w} \cdot \mathbf{x} + b$$

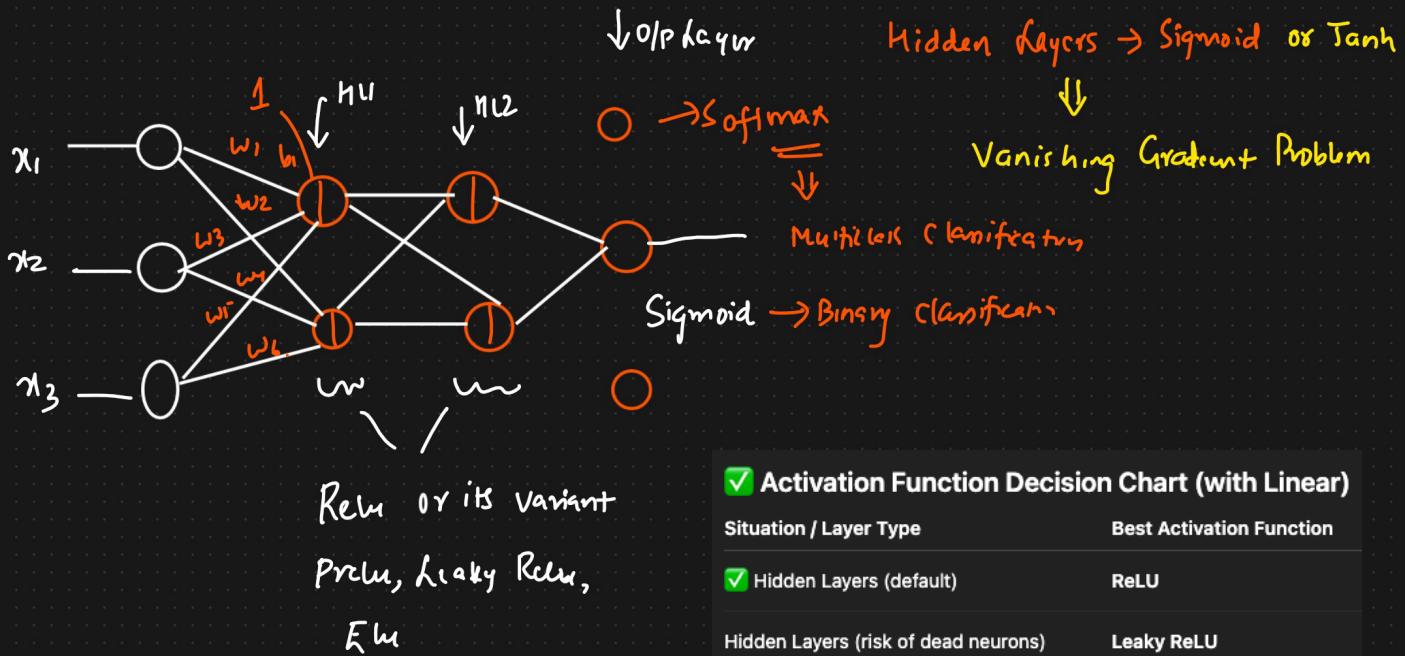
$$\text{Softmax} \Rightarrow \text{Cat} = \frac{e^{-1}}{e^{-1+0+3+5}} = 0.00033 \quad \Pr(\text{Horse}) = \frac{0.1353}{0.00033 + 0.0024 + 0.0183 + 0.1353}$$

$$\text{Dog} = \frac{e^0}{e^{-1+0+3+5}} = 0.0024 \approx 86\%$$

$$\text{Monkey} = \frac{e^3}{e^{-1+0+3+5}} = 0.0183$$

$$\text{Horse} = \frac{e^5}{e^{-1+0+3+5}} = 0.1353$$

7 Which Activation Function To Use When?



In short;

Hidden Layer — ReLU

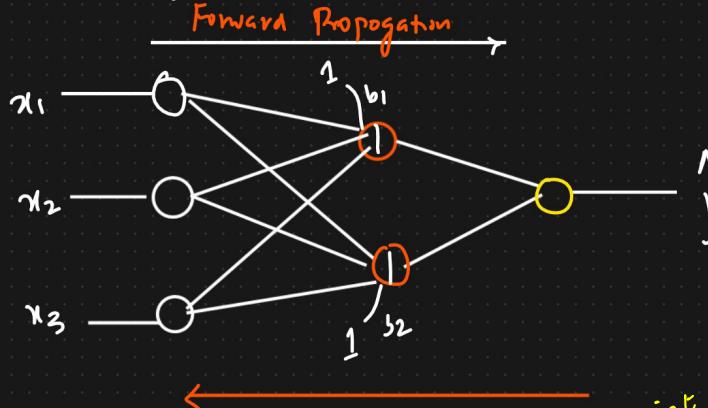
Output Layer (Classification) — Sigmoid
Softmax

Output Layer (Regression) — Regression Activator

Activation Function Decision Chart (with Linear)

Situation / Layer Type	Best Activation Function
Hidden Layers (default)	ReLU
Hidden Layers (risk of dead neurons)	Leaky ReLU
Hidden Layers (RNNs / time series)	tanh
Output – Binary Classification	Sigmoid
Output – Multi-class Classification	Softmax
Output – Regression (real number prediction)	Linear
Advanced models (deep / experimental)	Swish / ELU

Loss function And Cost Function



loss function → single datapoint Error

$$MSE = (y - \hat{y})^2$$

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

\uparrow Error

$x_1 \quad x_2 \quad x_3 \quad \text{O/P}$

\hat{y}

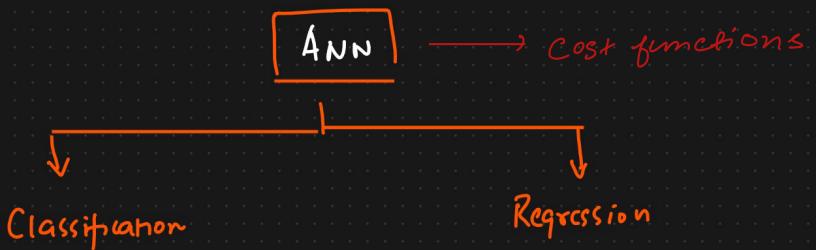
$\boxed{\text{loss} = (y - \hat{y})^2}$

$\left\{ \begin{array}{l} - - - 0 \\ - - - 1 \\ - - - 0 \\ - - - 1 \end{array} \right.$

Gradient Descent

Costfunction → Entire Data Error.

Costfunction → loss function for batches of data points

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


- ① Mean Squared Error (MSE)
- ② Mean Absolute Error (MAE)
- ③ Huber loss
- ④ RMSE

① Mean Squared Error (MSE)

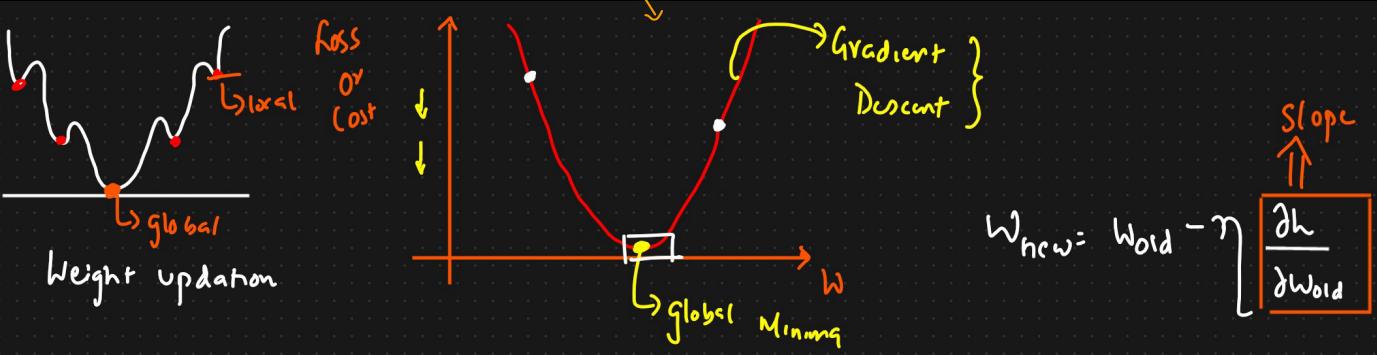
$$\text{loss function} = (y - \hat{y})^2$$

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

↳ Quadratic Equations

If we plot Loss function
we get Gradient Descent





Advantages

① MSE is Differentiable

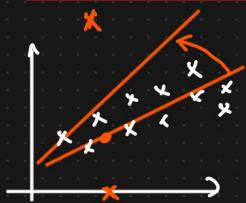
Imp: ② It has 1 local OR global Minima

③ It converges faster

Disadvantages

Imp:

① Not Robust to outliers



penalizing
the error

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

② Mean Absolute Error (MAE)

$$\text{Loss fn} = |y - \hat{y}|$$

$$\text{Cost fn} = \frac{1}{n} \sum_{i=1}^n |y_i - \hat{y}_i|$$

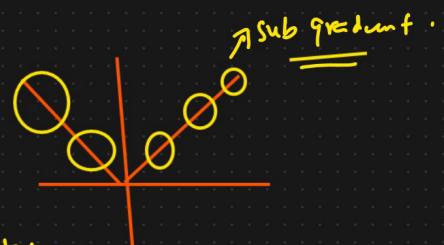
Advantages

① Robust to outliers



Disadvantages

① Convergence usually takes
time in MAE



③ Huber loss

Combination of MSE and MAE where we use MSE for no outliers and vice versa.

① MSE

② MAE

MSE

No outlier

Hypoparameter

$$\text{Cost fn} = \begin{cases} \frac{1}{n} \sum_{i=1}^n (y_i - \hat{y}_i)^2 & \text{if } |y_i - \hat{y}_i| \leq \boxed{s} \\ \boxed{s} & \text{if } |y_i - \hat{y}_i| > \boxed{s} \end{cases}$$

$$\begin{cases} \delta |\hat{y} - y| - \frac{1}{2} \delta^2, & \text{otherwise} \\ \downarrow \\ \text{MAE} \end{cases}$$

④ RMSE (Root Mean Squared Error)

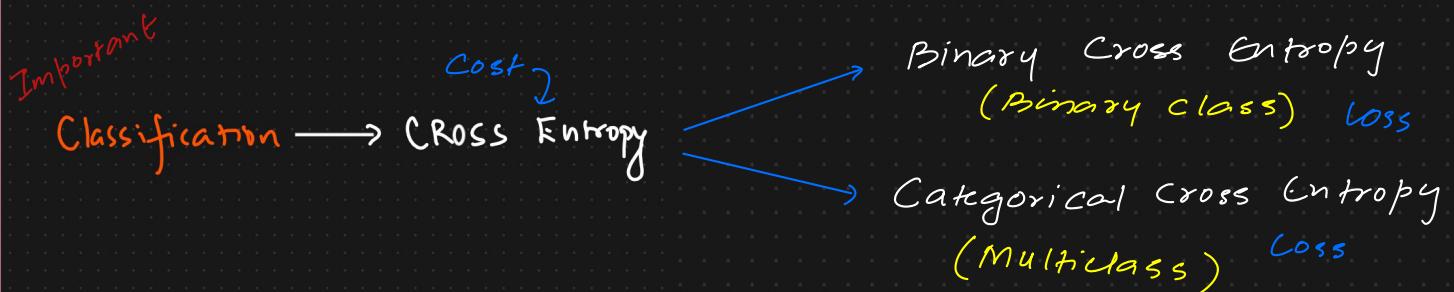
=

Cost function = $\sqrt{\frac{\sum_{i=1}^N (\hat{y}_i - y_i)^2}{N}}$

Advantages

Disadvantages

② Loss or Cost function For Classification Problems



① Binary Cross Entropy

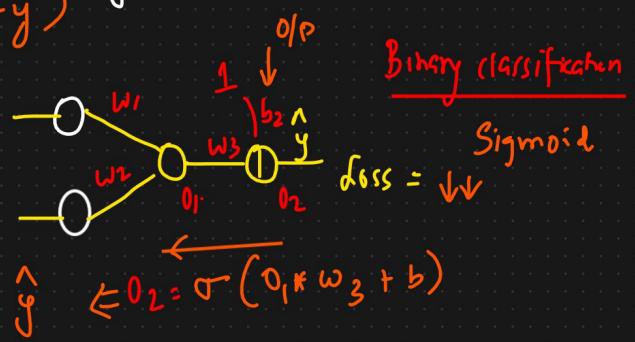
log loss

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

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

y = Actual Value

\hat{y} = Predicted Value



$$\hat{y} = \frac{1}{1+e^{-z}} \Rightarrow \text{Sigmoid Activation function}$$

② Categorical Cross Entropy (Multiclass classification) ONE is used

$\xrightarrow{\text{ONE} \rightarrow \text{One Hot Encoding}}$

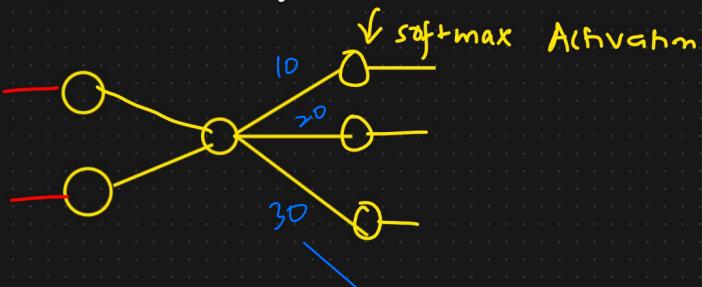
f_1	f_2	f_3	O/P	$j=1$ Good	$j=2$ Bad	$j=3$ Neutral	$C = \text{No. of categories}$
$\rightarrow 2$	3	4	Good	1	0	0	i = 1 to n
$\rightarrow 5$	6	7	Bad	0	1	0	
$\rightarrow 8$	9	10	Neutral	0	0	1	

$$d(x_i, y_i) = - \sum_{j=1}^C y_{ij} * \ln(\hat{y}_{ij})$$

Actual value $\leftarrow y_{ij} = \begin{bmatrix} y_{11} & y_{12} & y_{13} & \dots & y_{1c} \\ y_{21} & y_{22} & y_{23} & \dots & y_{2c} \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ y_{n1} & y_{n2} & y_{n3} & \dots & y_{nc} \end{bmatrix}$

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

Prediction $\leftarrow \hat{y}_{ij} \Rightarrow \text{Softmax Activation} = \text{Soft}(z) =$



$$\frac{e^{z_i}}{\sum_{j=1}^k e^{z_j}}$$

$$\begin{aligned}
 & \text{O/p } y_j = \text{Probabilities} \\
 & \text{Categorical} \Rightarrow [0.2, 0.3, 0.5] \\
 & \text{Cross Entropy} \\
 & \quad \downarrow \\
 & \quad \boxed{\text{This also gives the probability of other categories}}
 \end{aligned}$$

$$\begin{aligned}
 & [0.1, 0.2, 0.3, 0.2, 0.2] \leftarrow 1 \\
 & \quad \downarrow \\
 & = \frac{e^{z_i}}{e^{z_j}} \\
 & = \frac{e^{10}}{e^{10+20+30}}
 \end{aligned}$$

③ Sparse Categorical Cross Entropy

$$\begin{array}{c}
 \xrightarrow{\text{0 } 1^{\text{st}} \text{ } 2^{\text{nd}}} \text{Categories} \\
 \boxed{[0.2, 0.3, 0.5]} \\
 \downarrow
 \end{array}$$

Disadvantage

- ① losing info about the probability of other category.

$$\boxed{2^{\text{nd}} \text{ Index}} \Rightarrow \underline{\text{O/p}}$$

$$\begin{array}{c}
 \boxed{0 \ 1 \ 2^{\text{nd}} \ 3^{\text{rd}} \ 4^{\text{th}}} \\
 \underbrace{[0.2, 0.3, 0.1, 0.2, 0.2]}_{\text{Index}} \Rightarrow 1^{\text{st}} \text{ Index} \\
 \downarrow \\
 \text{Category} \Rightarrow \underline{\text{O/p}}
 \end{array}$$

④ Right Combination

	Activation applied	Hidden Layers	O/p Layer	Problem Statement	Loss function
①	Sigmoid	ReLU or its Variants		Binary Classification	Binary Cross Entropy
②	Softmax	ReLU or its Variants		Multi Class	Categorical or Sparse CE
③	Linear	ReLU or its Variants		Regression	MSE, MAE, Huber loss, RMSE

Conclusion

✓ Types of Neural Networks & Their Use Cases

Neural Network Type	Best Suited For	Data Type
ANN (Artificial Neural Network)	Regression & classification on tabular data	Numeric/tabular data
CNN (Convolutional Neural Network)	Image recognition, video analysis, object detection	Images, video, spatial data
RNN (Recurrent Neural Network)	Sequence modeling, text, time-series forecasting	Text, speech, time-series

- Loss function → error for single data point
- Cost function → error for entire Data (epoch or batch)

Regression



Loss → MSE → outlier sensitive
 MAE → outlier insensitive
 Huber loss → comb MSE/MAE
 RMSE

Q: Why do we need Regression Activator func in LR?
 Unlike softmax, it is used to output raw unbounded values, not non-linearity.

Q: "No Activation" = Linear Activation

That means:

The output layer just returns the raw value from the final neuron, with no transformation applied.

✓ Summary:

Term	Meaning
No Activation	Return raw output (linear function)
Used in	Regression (real-valued outputs)
Why	Keeps values unbounded ($-\infty$ to $+\infty$)

Classification

Classification Cheat Sheet

Classification Type	Loss Function	Activation Function
Binary Classification	Binary Cross Entropy (Log Loss) $L = -[y \cdot \log(p) + (1-y) \cdot \log(1-p)]$	Sigmoid $\sigma(x) = 1 / (1 + e^{-x})$
Multi-Class Classification	Categorical Cross Entropy $L = -\sum y_i \cdot \log(p_i)$ (for i = 1 to C)	Softmax $p_i = e^{z_i} / \sum e^{z_j}$ (for j = 1 to C)

Important

✓ Activation Functions by Problem Type

Problem Type	Hidden Layer Activation	Output Layer Activation
Binary Classification	ReLU	Sigmoid
Multi-Class Classification	ReLU	Softmax
Regression (Single Output)	ReLU	Linear

\$1M Question ?

"Why should I use **Artificial Neural Networks (ANNs)** for regression or classification when I can just use **Linear Regression** or **Logistic Regression**?"

Short Answer:

You don't always need ANN. But you might want to use ANN when **traditional ML models aren't enough**.

Deeper Comparison:

Aspect	Linear/Logistic Regression	ANN (Deep Learning)
Complexity of patterns	Handles linear relationships only	Captures non-linear relationships
Feature engineering	Often requires manual feature crafting	Learns features automatically
Scalability	Fast and interpretable	Scales better with large data
Overfitting risk	Low on small data	High unless regularized/trained well
Interpretability	Very high (coefficients are clear)	Low (black box)
Speed	Very fast	Slower, needs more compute

Example Scenario:

Logistic Regression Fails:

You have a classification problem where the classes are **not linearly separable** (e.g., spiral or concentric circles).

👉 Logistic regression draws a straight line. It fails.

ANN Wins:

A neural network can learn **non-linear decision boundaries** and separate the classes better.

TL;DR:

Use **Logistic/Linear Regression** when:

- "Data is small"
- "Relationships are linear"
- "You want speed & interpretability"

Use **ANN** when:

- "Data is large or complex"
- "You suspect non-linear interactions"
- "You're hitting performance limits with traditional models"



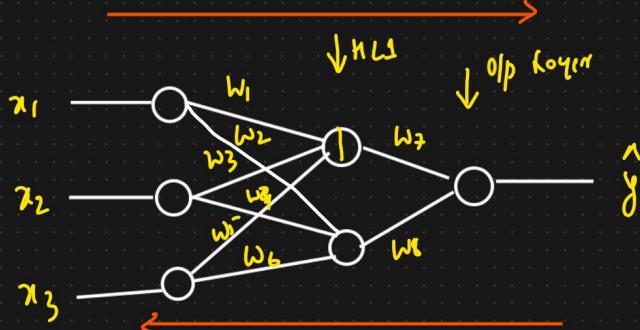
Pro tip: Always **start simple** (like Logistic Regression).

Only add complexity (ANN) when it's needed.

Optimizers → used for reducing cost function
or reach Global minima.

- ① Gradient Descent
 - ② Stochastic Gradient Descent (SGD)
 - ③ Mini batch SGD
 - ④ SGD With Momentum
 - ⑤ Adagrad and RMS Prop
 - ⑥ Adam Optimizers
- Epoch → completed iteration i.e. | Forward P.
| Backward P.
- Best of both worlds
but bit noisy

Gradient Descent Optimizer



$$\text{loss} = [\quad] \downarrow \downarrow \downarrow$$

Optimizers ↑

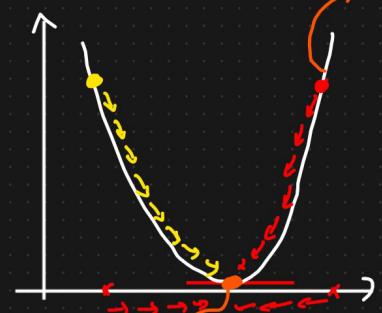
loss or cost

Gradient Descent

loss ↓ ↓

$w_{\text{new}} \approx w_{\text{old}}$

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



MSE

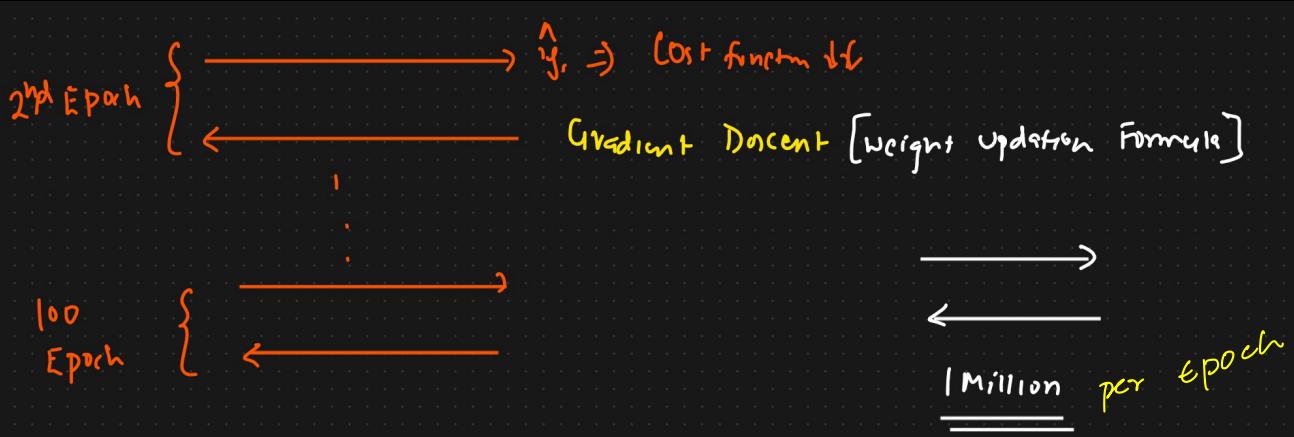
$$\text{Loss fn} = (y - \hat{y})^2 \quad \text{Cost fn} = \frac{1}{h} \sum_{i=1}^n (y_i - \hat{y}_i)^2$$

DataSet = 1000 Data Points

Epochs, Iteration
1 Epoch = 1 Iteration
1 Epoch → $\hat{y}_i \Rightarrow$ Cost function w.
Weights will get updated

Gradient Descent optimizer [Weight update]

10 Iteration
 $\frac{100}{10} = 100$
100 →
←
←
←
←
←
←
←
←
←



Advantages

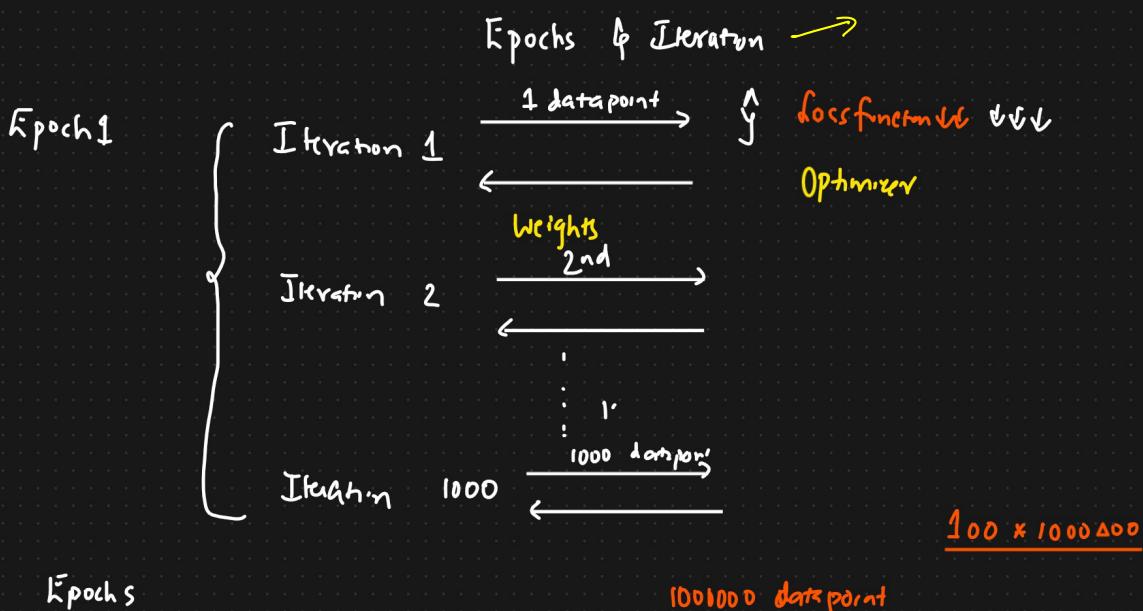
- ① Convergence will happen.
 - Accurate Gradients
 - Good for Small Datasets
- ② Huge Resource RAM, GPU

Disadvantage

Resource Intensive -

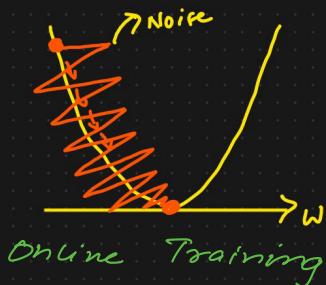
(2) Stochastic Gradient Descent (SGD)

1000 datapoint



Advantage

- ① Solve Resource Issue



Disadvantage

- ② Time Complexity \uparrow

- ③ Convergence will also take More time.
- ④ Noise gets Introduced

→ may skip Global minima.

③ Mini Batch SGD

Batch size

$$\text{No. of iterations} = \frac{100000}{1000} = 100 \text{ iterations}$$

Epoch, Iteration, Batch-size

Data points = 100000

batch.size = 1000

MSG

$$\text{Cost fn} = \sum_{i=1}^{1000} (y_i - \hat{y}_i)^2$$

Epoch 1

Iteration 1

Optimizer \Rightarrow Mini Batch SGD

change the weight

1000

↓↓

[8gb]

1

Iteration 2

[16gb]

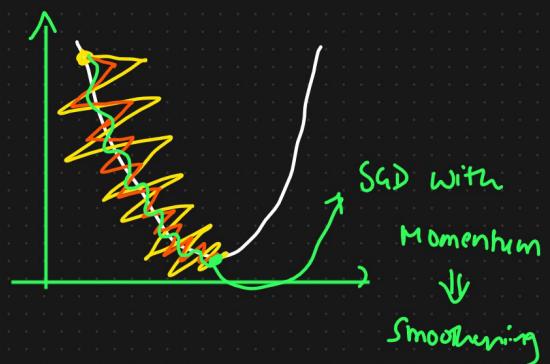
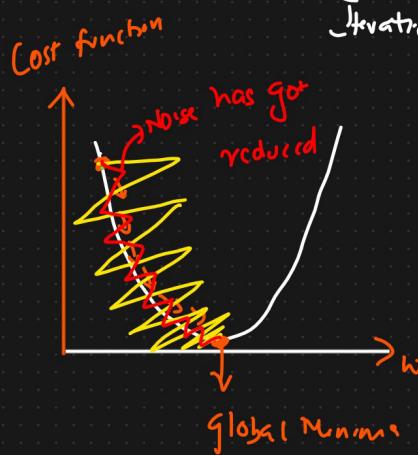
5000

Iteration 3

⋮

Iteration 100

Batch Size



Advantages

- ① Convergence speed will increase
- ② Noise will be less when compared to SGD
- ③ Efficient Resource Usage (RAM)

Disadvantage

- ① Noise still exists

in every epoch of 1000 iterations
1000 DP ob

CONCLUSION

Optimizers = Reduce cost function

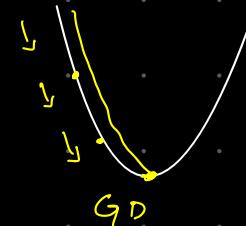
working



1. Batch Gradient Descent (GD)

You taste all 100 cakes at once, take notes, average out all your feedback, and then make one big adjustment to the recipe.

- One big decision after thorough review
- Takes time (you must finish all cakes)
- Very precise, but slow
- ✗ 1 update per round (epoch)
- ✗ Uses all 100 cakes each time



2. Stochastic Gradient Descent (SGD)

You're in a hurry — you taste one cake, make a quick adjustment. Then you taste the next cake, make another tweak.

- Fast but very chaotic
- Some tweaks may help, others might hurt
- But over time, you may stumble upon the perfect recipe
- ✗ 100 updates per round (epoch)
- ✗ Uses 1 cake at a time, very noisy

Summary (Cake Edition):

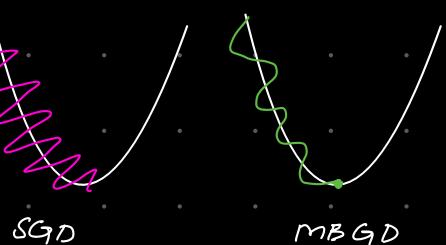
Method	How You Taste	How You Adjust	Feeling
GD	Taste all 100	One big tweak	Slow but sure
SGD	Taste 1 cake	Constant tweaking	Fast but shaky
Mini-Batch	Taste 10 at a time	Small steady tweaks	Best of both worlds



3. Mini-Batch Gradient Descent (batch size = 10)

You group the 100 cakes into 10 plates (batches) of 10 cakes each. You taste a plate, make a small adjustment, and repeat.

- Balanced speed
- Not too noisy, not too slow
- Gets better with practice
- ✗ 10 updates per round (epoch)
- ✗ Each update uses 10 cakes



EXPLANATION ;

Gradient Descent Variants Explained

Type	How it Works	Speed	Stability	Used When	Pros	Cons
Batch Gradient Descent	Uses the entire dataset to compute gradient before updating weights.	✗ Slow (per update)	✓ Very stable	Small datasets	+ Accurate gradients + Converges smoothly	- Very slow on large data - High memory usage
Stochastic GD (SGD)	Uses 1 data point at a time to update weights.	✓ Fast (per update)	✗ Very noisy	Large datasets, online learning	+ Fast updates + Good for huge/streaming data + Can escape local minima	- Noisy convergence - May overshoot or bounce around
Mini-Batch SGD	Uses a small subset of data (e.g., 32, 64) for each update.	✗ Balanced	✓ Smooth + efficient	Most common in deep learning	+ Fast and stable + Efficient on GPUs + Generalization-friendly	- Needs tuning of batch size - Still a bit noisy

we need Momentum biased SGD to smooth noise (exponential weighted) average

(4) SGD With Momentum



Weight Updation formula

$$w_{\text{new}} = w_{\text{old}} - \eta \left[\frac{\partial h}{\partial w_{\text{old}}} \right]$$

$$b_{\text{new}} = b_{\text{old}} - \eta \left[\frac{\partial h}{\partial b_{\text{old}}} \right]$$

$$w_t = w_{t-1} - \eta \left[\frac{\partial h}{\partial w_{t-1}} \right]$$

Exponential Weight Average {Smoothing} \Rightarrow ARIMA, SARIMAX

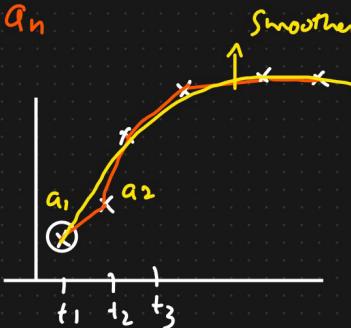
Time = $t_1 \ t_2 \ t_3 \ t_4 \ \dots \ t_n$ Time Series

Values = $a_1 \ a_2 \ a_3 \ a_4 \ \dots \ a_n$

$$V_{t_1} = a_1$$

$$\beta = 0.95$$

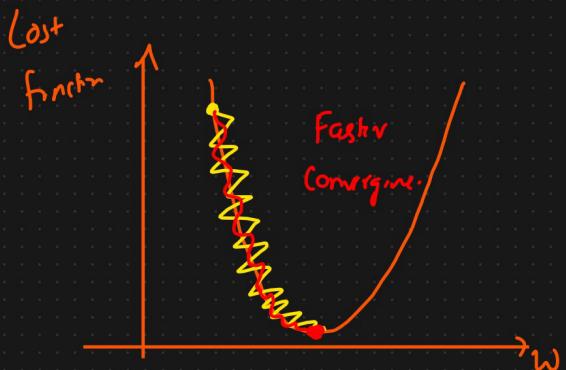
$$V_{t_2} = \boxed{\beta} * V_{t_1} + (1-\beta) * a_2$$



$$\underline{\beta = 0.95} = \underline{0.95 * a_1 + (0.05) a_2}$$

$$V_{t_3} = \beta * V_{t_2} + (1-\beta) * a_3$$

$$= 0.95 \left[0.95 * a_1 + (0.05) a_2 \right] + (0.05) * a_3$$



Advantage

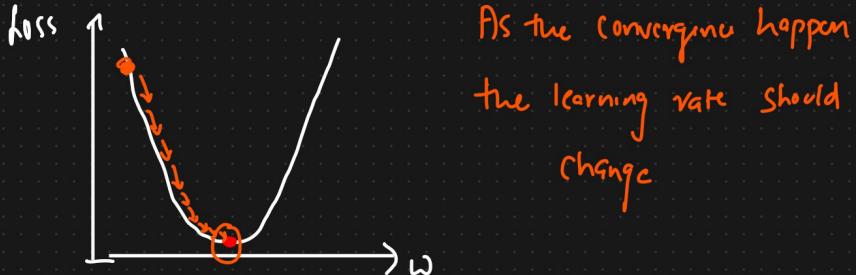
- ① Reduces the noise
- ② Quick Convergence

⑤ Adagrad = Adaptive Gradient Descent

$\eta = \text{fixed} \Rightarrow \text{Dynamic learning}$

$$w_t = w_{t-1} - \eta \left[\frac{\partial h}{\partial w_{t-1}} \right]$$

Learning Rate = 0.001



$0.00001 \approx 0$

$$w_t = w_{t-1} - \eta' \left[\frac{\partial h}{\partial w_{t-1}} \right]$$

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

$d_t \uparrow \uparrow \quad \epsilon \text{ small value}$

$t=1 \quad t=2 \quad t=3 \quad w_t \approx w_{t-1}$

$$d_t = \sum_{i=1}^t \left(\frac{\partial h}{\partial w_i} \right)^2$$

$$\eta = 0.01 \quad \eta = 0.005 \quad \eta = 0.003$$

Disadvantage

- ① $\eta' \rightarrow$ Possibility to become a very small value ≈ 0

⑥ Adadelta And RMSprop

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

$\stackrel{=}{t}$

Exponential Weight Average

$$\beta = 0.95 \quad S_{dw_t} = 0$$

$$S_{dw_t} = \beta * S_{dw_{t-1}} + (1-\beta) \left(\frac{\partial h}{\partial w_{t-1}} \right)^2$$

[Dynamic LR + Smoothening (EWA)]

$$w_t = w_{t-1} - \eta' \frac{\partial h}{\partial w_{t-1}}$$

⑦ Adam Optimizer

SGD with Momentum + RMSprop [Dynamic LR + Smoothening]

$$w_t = w_{t-1} - \eta' V_{dw} \quad \text{Weight Update}$$

$$b_t = b_{t-1} - \eta' V_{db} \quad \text{Bias Update}$$

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

EWA

$$S_{dw_t} = 0$$

$$S_{dw_t} = \beta * S_{dw_{t-1}} + (1-\beta) \left(\frac{\partial h}{\partial w_{t-1}} \right)^2$$

$$V_{dw_t} = \beta * V_{dw_{t-1}} + (1-\beta) \frac{\partial h}{\partial w_{t-1}}$$

$$V_{db_t} = \beta * V_{db_{t-1}} + (1-\beta) \frac{\partial h}{\partial b_{t-1}}$$

\Rightarrow Momentum
 \hookrightarrow Smoothening

Final Conclusion

✓ Gradient Descent Variants & Optimizers — Full Comparison

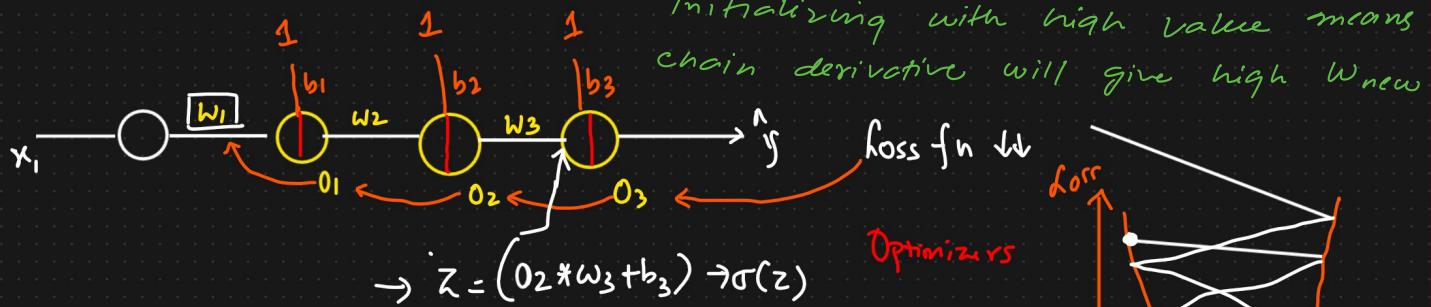
Type	How it Works	Speed	Stability	Used When	Pros	Cons
Batch Gradient Descent	Uses the entire dataset to compute gradient before each weight update.	✗ Slow (per update)	✓ Very stable	Small datasets	+ Accurate gradients + Smooth convergence	– Very slow on large data – High memory usage
Stochastic GD (SGD)	Uses 1 data point at a time to compute and apply updates.	✓ Fast (per update)	✗ Very noisy	Large datasets, online learning	+ Very fast updates + Good for real-time/streaming + Can escape local minima	– Noisy updates – May bounce or overshoot
Mini-Batch SGD	Uses a small batch (e.g., 32–128 samples) per update.	⚖ Balanced	✓ Smooth + efficient	Deep learning, GPU training	+ Combines speed & stability + Efficient on hardware + Improves generalization	– Needs batch size tuning – Some noise still persists
SGD with Momentum	Like Mini-Batch SGD, but adds momentum : a fraction of past updates to current updates — helps smooth curves and accelerate convergence.	✓ Faster than SGD	✓ More stable	Deep learning, non-convex optimization	+ Speeds up convergence + Reduces oscillation + Helps escape flat regions (plateaus)	– Needs momentum tuning – Risk of overshooting
AdaGrad	Adapts learning rate per parameter based on past squared gradients — good for sparse data.	⚡ Fast initially	⚖ Decent	Sparse features, NLP tasks	+ No LR tuning needed + Great for sparse data like text (NLP)	– Learning rate decays too fast – Stops learning early
Adam Optimizer	Combines momentum + adaptive learning rate per parameter using moving averages of gradients.	✓ Fast & efficient	✓ Very stable	Deep learning (default choice)	+ Fast convergence + Works well out-of-the-box + Handles sparse & noisy gradients	– May overfit – Sensitive to hyperparameters – May not generalize best

🧠 TL;DR:

Situation	Best Choice
Small, simple dataset	Batch GD
Real-time or massive streaming data	<i>Online Train</i> SGD
General deep learning tasks	Mini-Batch SGD / Adam
Complex loss surfaces, faster training	<i>Less Noise</i> SGD with Momentum
NLP or sparse data	AdaGrad / Adam

Imp. Adam optimizer = Best optimizer (moving avg of gradients).
 Combine Adaptive Learning Rate + momentum

Exploding Gradient Problem \Rightarrow Weight Initialization Technique



$$w_{1,new} = w_{1,old} - \eta \boxed{\frac{\partial h}{\partial w_{1,old}}}$$

$\left\{ \begin{array}{l} w_{new} \ggg w_{1,old} \\ w_{new} \lll w_{1,old} \end{array} \right.$
 Chain Rule of Derivatives

$$\frac{\partial h}{\partial w_{1,old}} = \frac{\partial h}{\partial o_3} * \boxed{\frac{\partial o_3}{\partial o_2}} * \frac{\partial o_2}{\partial o_1} * \frac{\partial o_1}{\partial w_{1,old}}$$

big * big * big * big \rightarrow big value

$$\frac{\partial o_3}{\partial o_2} = \boxed{\frac{\partial \sigma(\omega)}{\partial \omega}} * \frac{\partial \omega}{\partial o_2}$$

$$= [0 - 0.25] * \frac{\partial (o_2 * w_3 + b_3)}{\partial o_2} = [0 - 0.25] * w_3 \Rightarrow \underline{\underline{500 - 1000}}$$

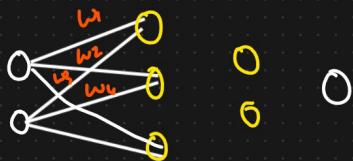
Exploding gradients happen when gradients grow too large during backpropagation, often due to deep networks or bad initialization. They can break the training process. Fix it with good initialization, gradient clipping, and careful architecture design.

Gradient explode =

- Bad initialization
- Deep networks

Weight Initialising Techniques

- ① Uniform Distribution ✓
- ② Xavier/Glorot Initialization ✓
- ③ Kaiming (He) Initialization ✓

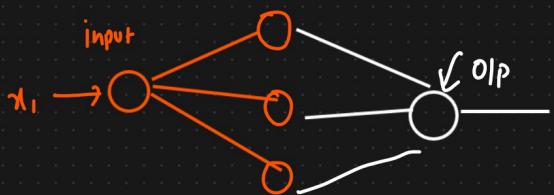


- ① Weights should be small ✓
- ② Weights should not be same ✓
- ③ Weights should have good variance

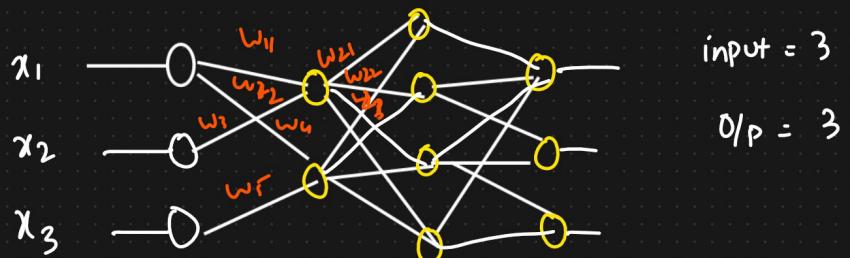
Key Points

Solutions:

Problem	Solution
Gradients getting too large	Gradient Clipping
Poor initialization	Use He/Xavier Initialization
Very deep network	Use architectures like ResNet (skip connections)
High learning rate	Lower learning rate



input = 1
Output = 1



input = 3

O/P = 3

① Uniform Distribution

$$W_{ij} \sim \text{Uniform Distribution} \left[\frac{-1}{\sqrt{\text{input}}}, \frac{1}{\sqrt{\text{input}}} \right]$$

$$\left[\frac{-1}{\sqrt{3}}, \frac{1}{\sqrt{3}} \right]$$

② Xavier/Glorot Initialization

Recurrent \rightarrow Xavier Glorot

① Xavier Normal Init

$$W_{ij} \sim N(0, \sigma)$$

$$\sigma = \sqrt{\frac{2}{(\text{input} + \text{output})}}$$

② Xavier Uniform

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

$$\left[\frac{-\sqrt{6}}{\sqrt{\text{input} + \text{output}}}, \frac{\sqrt{6}}{\sqrt{\text{input} + \text{output}}} \right]$$

③ Kaiming He Initialization

① He Normal

$$W_{ij} \sim N(0, \sqrt{\sigma^2})$$

$$\sigma = \sqrt{\frac{2}{\text{input}}}$$

② He uniform

$$W_{ij} \sim \text{Uniform Distribution} \left[-\sqrt{\frac{6}{\text{input}}}, \sqrt{\frac{6}{\text{input}}} \right]$$

Conclusion:

Q: why do we need weight initialization Techniques?

when training a neural network, how we initialize the weights matters a lot - it can make or break our model's ability to learn.

✓ Key Reasons:

Reason	Explanation
Breaks Symmetry	If all weights are the same (e.g., zeros), neurons learn the same features — no diversity.
Avoids Vanishing/Exploding Gradients	Poor initialization can cause gradients to become too small or too large during backprop.
Speeds Up Convergence	Good initialization helps the model learn faster and more reliably.

Imp

✓ Essential Weight Initialization Techniques

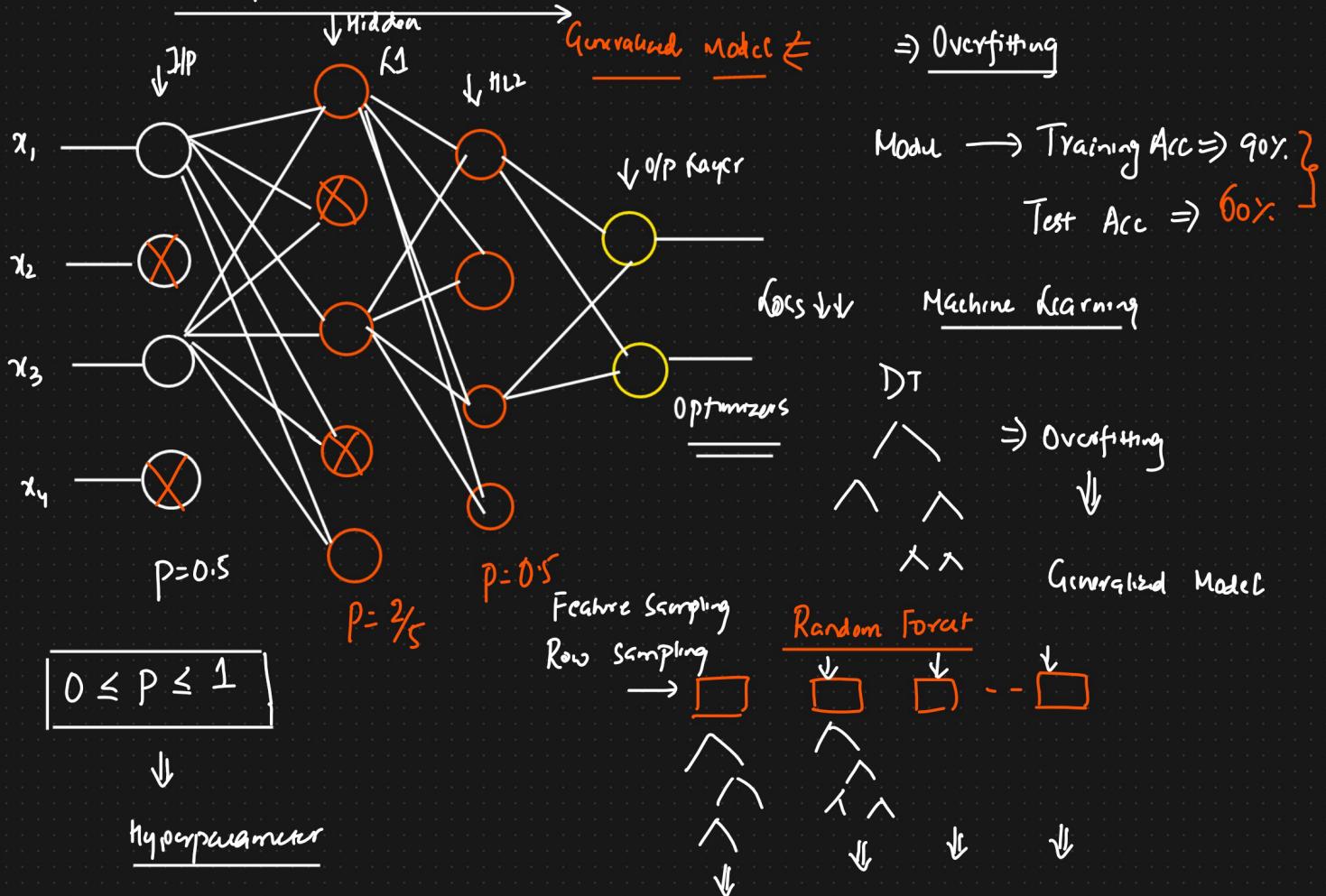
Technique	Use When	Activation Function
Xavier (Glorot)	For balanced input/output layers	tanh, sigmoid
He Initialization	For deep networks with ReLU-based activations	ReLU, Leaky ReLU
Random Uniform	Simple or shallow networks	Any (but not recommended for deep networks)



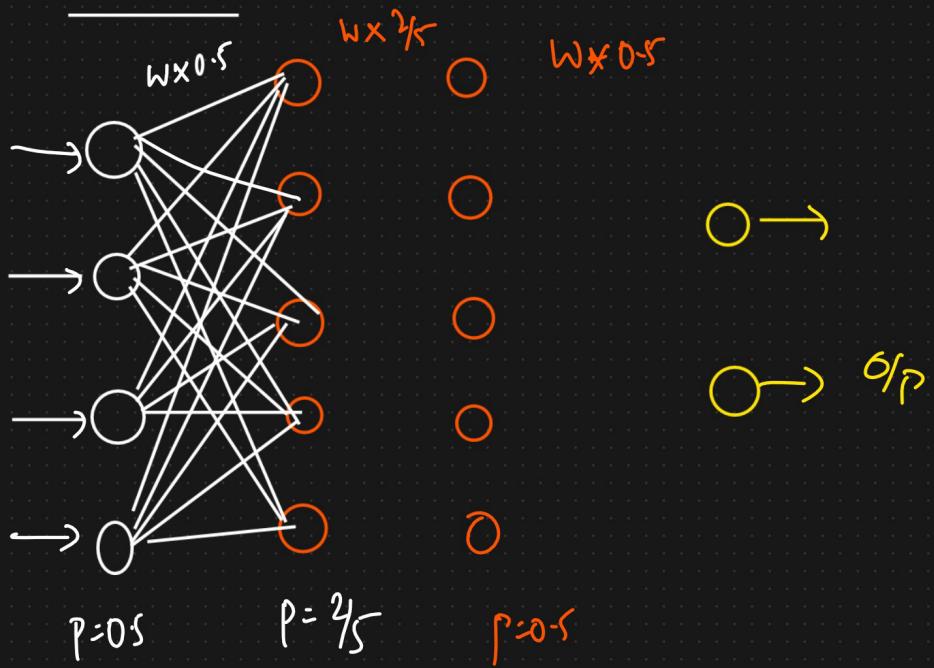
Quick Recap:

- 🎯 Xavier → good for smooth activations
- ⚡ He → best for ReLU and its variants
- ⚠ Avoid zeros or constant initializations

Drop Out layer



Test data



🧠 What is a Dropout Layer (in Deep Learning)?

A Dropout layer is a regularization technique used in neural networks to prevent overfitting.

✍ How it works:

During training, the Dropout layer:

- Randomly turns off (drops) a percentage of neurons in a layer.
- This means those neurons don't contribute to forward or backward passes for that batch.

For example, if you set:

```
python
```

Copy Edit

```
Dropout(rate=0.5)
```

➡ 50% of the neurons in that layer are randomly "shut off" during each training step.

Very important

📌 Why use Dropout?

- It forces the network to not rely too heavily on specific neurons.
- Encourages the network to learn redundant representations, making it more robust.
- Great for reducing overfitting, especially in large networks.

💻 In Code (Keras Example):

```
python
```

Copy Edit

```
from tensorflow.keras.layers import Dense, Dropout

model.add(Dense(128, activation='relu'))
model.add(Dropout(0.5)) # Drop 50% of neurons
model.add(Dense(10, activation='softmax'))
```

ANN PRACTICAL

FOR CLASSIFICATION

Tensorflow → library for Deep learning by Google

Keras → API on top of Tensorflow

Important; we can use either Tensorflow / Pytorch.

STEP I; Setting up Environment

```
Bash → Conda Create -p venv python == 3.13  
conda activate venv  
pip install -r requirements.txt  
install ipykernel
```

PREPROCESSING

1. Remove unnecessary columns
2. Use Label Encoder → Binary Categories
3. Use OHE for → Multi Categories

Imp- If we use Label Encoder on multi class, it will assign a=1, b=2, c=0 ; which will mean for bc b>a & c , so OHE is preferred.

4. Save Encoders in pickle file.
5. Split the Data into train-test
6. Scale the Data (Standard Scaler)
→ Save this also as pickle

VV Imp In Tensorflow, we call ANN or sequential model.

Useful parameters ; SKPS

1. sequential
2. Dense → 64 means 64 neuron
3. Activation fn → Sigmoid, softmax → O/P
ReLU, Tanh → hidden
4. Optimizers →

5. Loss \rightarrow GD, SGD

6. Evaluation Metrics

7. Training → Logs → folders
↓
Tensorboard

• ANN

Libraries → tensorflow.keras.models import sequential

`tensorflow.keras.layers import Dense` → Nodes

```
tensorflow.keras.callbacks import EarlyStopping  
TensorBoard
```

TensorBoard

↳ logs

Build ANN Model

```
# Build ANN Model
model = Sequential()
model.add(Dense(64, activation='relu', input_shape=(X_train_scaled.shape[1],))) ## 1st Input layer connected to 64 neurons - HL1
model.add(Dense(32, activation='relu')) # no shape needed here as it is a hidden layer - Hidden layer with 32 neurons - HL2
model.add(Dense(1, activation='sigmoid')) # Output layer with 1 neuron for binary classification
```

→ Setup Optimizer, Loss function, metrics

```
# compile the model with optimizer, loss function and metrics
from tensorflow.keras.optimizers import Adam
model.compile(optimizer=Adam(learning_rate=0.01), loss='binary_crossentropy', metrics=['accuracy'])
```

→ set up Tensorboard & Early stopping

- **TensorBoard**: For real-time monitoring of training (loss, accuracy, etc.)
 - **EarlyStopping**: To automatically stop training when performance stops improving (prevents overfitting and saves time) \longrightarrow *stops when model reaches certain threshold*

```
log_dir = "logs/fit/" + datetime.datetime.now().strftime("%Y%m%d-%H%M%S")
```

1

→ Train the Model file

```
history = model.fit(X_train_scaled, y_train,
                     validation_data=(X_test_scaled, y_test),
                     epochs=100,
                     callbacks=[tensorboard_callback, early_stopping])

✓ 16.5s

Epoch 1/100
250/250 [=====] - 2s 7ms/step - loss: 0.5535 - accuracy: 0.7806 - val_loss: 0.5473 - val_accuracy: 0.8030
Epoch 2/100
250/250 [=====] - 1s 6ms/step - loss: 1.2733 - accuracy: 0.7465 - val_loss: 1.9880 - val_accuracy: 0.6815
Epoch 3/100
250/250 [=====] - 1s 6ms/step - loss: 2.8616 - accuracy: 0.7368 - val_loss: 3.8312 - val_accuracy: 0.7565
Epoch 4/100
250/250 [=====] - 1s 6ms/step - loss: 3.3271 - accuracy: 0.7303 - val_loss: 17.8678 - val_accuracy: 0.6560
Epoch 5/100
250/250 [=====] - 1s 6ms/step - loss: 6.7133 - accuracy: 0.7289 - val_loss: 3.3616 - val_accuracy: 0.8045
Epoch 6/100
250/250 [=====] - 1s 6ms/step - loss: 8.6875 - accuracy: 0.7306 - val_loss: 9.0012 - val_accuracy: 0.7670
Epoch 7/100
250/250 [=====] - 1s 6ms/step - loss: 13.4958 - accuracy: 0.7394 - val_loss: 7.7190 - val_accuracy: 0.7380
Epoch 8/100
250/250 [=====] - 1s 6ms/step - loss: 24.6641 - accuracy: 0.7144 - val_loss: 17.9428 - val_accuracy: 0.7920
Epoch 9/100
250/250 [=====] - 1s 6ms/step - loss: 18.2670 - accuracy: 0.7521 - val_loss: 6.6817 - val_accuracy: 0.8015
Epoch 10/100
250/250 [=====] - 1s 6ms/step - loss: 19.5793 - accuracy: 0.7377 - val_loss: 43.4841 - val_accuracy: 0.6100
Epoch 11/100
250/250 [=====] - 1s 6ms/step - loss: 25.3717 - accuracy: 0.7340 - val_loss: 17.2696 - val_accuracy: 0.7635
```

→ Save model as .h5

```
# Save model file as .h5 is compatible with Keras
model.save('customer_churn_model.h5')
```

→ Monitor logs in Tensorboard.

```
# Load TensorBoard Extension
%load_ext tensorboard

# Launch TensorBoard
%tensorboard --logdir logs/fit
```



Finally

Get Predictions

- Import PKL files & .h5 file
- Get prediction → Ready for production-
- Create app.py which runs all the trained models

FOR REGRESSION

```
### Load the trained model, scaler pickle, onehot
model=load_model('customer_churn_model.h5')

## load the encoder and scaler
with open('ohe_geo.pkl','rb') as file:
    label_encoder_geo=pickle.load(file)

with open('label_encoder_gender.pkl', 'rb') as file:
    label_encoder_gender = pickle.load(file)

with open('scaler.pkl', 'rb') as file:
    scaler = pickle.load(file)
```

```
In [65]: ## Predict churn
prediction=model.predict(input_scaled)
prediction

WARNING:tensorflow:5 out of the last 5 calls to <function Model.make_predict_function.<locals>.predict_function at 0x33227f560> triggered tf.function retracing. Tracing is expensive and the excessive number of tracings could be due to (1) creating @tf.function repeatedly in a loop, (2) passing tensors with different shapes, (3) passing Python objects instead of tensors. For (1), please define your @tf.function outside of the loop. For (2), @tf.function has reduce_retracing=True option that can avoid unnecessary retracing. For (3), please refer to https://www.tensorflow.org/guide/function#controlling_retracing and https://www.tensorflow.org/api_docs/python/tf/function for more details.

WARNING:tensorflow:5 out of the last 5 calls to <function Model.make_predict_function.<locals>.predict_function at 0x33227f560> triggered tf.function retracing. Tracing is expensive and the excessive number of tracings could be due to (1) creating @tf.function repeatedly in a loop, (2) passing tensors with different shapes, (3) passing Python objects instead of tensors. For (1), please define your @tf.function outside of the loop. For (2), @tf.function has reduce_retracing=True option that can avoid unnecessary retracing. For (3), please refer to https://www.tensorflow.org/guide/function#controlling_retracing and https://www.tensorflow.org/api_docs/python/tf/function for more details.

1/1 [=====] - 0s 37ms/step
Out[65]: array[[0.01600381]], dtype=float32

In [66]: prediction_proba = prediction[0][0]

In [67]: prediction_proba

Out[67]: 0.016003806

In [68]: if prediction_proba > 0.5:
    print('The customer is likely to churn.')
else:
    print('The customer is not likely to churn.')

The customer is not likely to churn.
```

CNN

Convolutional Neural Network is the type of deep learning model that is especially powerful for working with images, videos etc.

💡 Why "Convolutional"?



Because it uses a mathematical operation called **convolution** to extract **features** (like edges, textures, shapes) from input data—especially images.

Input Image Patch:

```
[1, 2, 0]  
[4, 5, 6]  
[7, 8, 9]
```

And a 3x3 filter (also called kernel):

```
csharp  
Filter:      Kernel  
[1, 0, -1]  
[1, 0, -1]  
[1, 0, -1]
```

Purpose of Convolution =

Detect edges, corners, texture

💡 Convolution means:

Multiply each value from the image and the filter element-wise, then sum them up.

markdown

Result:

$$\begin{aligned} &= 1*1 + 2*0 + 0*(-1) \\ &+ 4*1 + 5*0 + 6*(-1) \\ &+ 7*1 + 8*0 + 9*(-1) \\ \\ &= 1 + 0 + 0 + 4 + 0 - 6 + 7 + 0 - 9 = -3 \end{aligned}$$

That -3 becomes a single pixel in the output feature map.

You move the filter over the whole image to generate more values.

🔍 Key Components of a CNN:

1. Input Layer – Raw image data (e.g., 224x224 pixels).
2. Convolutional Layer – Applies filters to detect patterns (like edges or corners). → (Kernel)
3. Activation Function – Usually ReLU, adds non-linearity.
4. Pooling Layer – Reduces dimensionality (e.g., max pooling).
5. Fully Connected Layer – Final layer(s) to make predictions. → Flattening
6. Output Layer – Gives final result (e.g., classification label like "cat" or "dog").

💡 What CNNs Are Used For:

- Image classification (e.g., recognizing objects in photos)
- Face recognition
- Self-driving cars (lane detection, obstacle recognition)

🧠 Simple Analogy:

Think of CNNs like a brain that first **detects edges**, then **shapes**, then **objects** — like how humans recognize patterns step by step.

Main concept

is to mimim how human brain analyses images

Cerebral cortex → Visual (layers) cortex

CNN TERMINOLOGY

Q: What is an image?

Image is a grid of pixel values.

→ For B & W image → each pixel = 0 - 255

0 = Black, 255 = White

Single channel (one value for each pixel)

e.g., 5×5 image = 5×5 vector

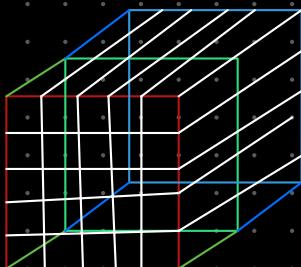
→ For a colored image (RGB)

Each pixel has 3 values (3 channels)

e.g., 5×5 image = $\underline{5 \times 5 \times 3}$ tensor

Imp.

R G B



$5 \times 5 \times 3 \rightarrow 3$ channel

◆ For Grayscale Images:

- Pixel value = 0 → black
- Pixel value = 255 → white
- Values in between represent shades of gray
 - 127 ≈ medium gray

◆ For RGB Images:

Each pixel has three channels:

- Red (R), Green (G), Blue (B)
- Each ranges from 0 (no intensity) to 255 (full intensity)

Example:

- (0, 0, 0) → Black (no light at all)
- (255, 255, 255) → White (full intensity of all colors)
- (255, 0, 0) → Red
- (0, 255, 0) → Green
- (0, 0, 255) → Blue

Q: What is Tensor?

A tensor is a container that can hold data in dimensions (any dimensions)

T2
34

Examples:

- 0D Tensor → Scalar → 7
- 1D Tensor → Vector → [3, 5, 7]
- 2D Tensor → Matrix → [[1, 2], [3, 4]]
- 3D Tensor → Stack of matrices → Like a color image [Height x Width x Channels]

Steps

→ Step I ; Input image pixel info

Step II ; Normalise → divide by 255

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

0	0	0	1	1	1
0	0	0	1	1	1
0	0	0	1	1	1
0	0	0	1	1	1
0	0	0	1	1	1
0	0	0	1	1	1

Step III ;

Convolution Operation in CNN to detect edges, corners and texture.

stride = 1 means
Right by 1 skip

I/P

$$\begin{array}{c}
 \text{stride=1} \\
 \begin{array}{|c|c|c|c|c|c|c|} \hline
 0 & 0 & 0 & 1 & 1 & 1 & 1 \\ \hline
 0 & 0 & 0 & 1 & 1 & 1 & 1 \\ \hline
 0 & 0 & 0 & 1 & 1 & 1 & 1 \\ \hline
 0 & 0 & 0 & 1 & 1 & 1 & 1 \\ \hline
 0 & 0 & 0 & 1 & 1 & 1 & 1 \\ \hline
 0 & 0 & 0 & 1 & 1 & 1 & 1 \\ \hline
 \end{array} \\
 6 \times 6 \times 1
 \end{array}
 \quad
 \begin{array}{c}
 n=6 \quad f=3 \\
 * \quad \begin{array}{|c|c|c|} \hline
 +1 & 0 & -1 \\ \hline
 +2 & 0 & -2 \\ \hline
 +1 & 0 & -1 \\ \hline
 \end{array} \\
 3 \times 3 \quad \text{filter} \\
 \text{filters / Kernel} \\
 \text{vertical edge filters}
 \end{array}
 \quad
 \begin{array}{c}
 0/p=4 \\
 \begin{array}{|c|c|c|c|} \hline
 0 & -4 & -4 & 0 \\ \hline
 0 & -4 & -4 & 0 \\ \hline
 0 & -4 & -4 & 0 \\ \hline
 0 & -4 & -4 & 0 \\ \hline
 \end{array} \\
 4 \times 4 \quad O/P
 \end{array}$$

Imp

Q, how come we got O/p as 4x4?

so we have formula ;

$$\hookrightarrow n - f + 1 \quad \begin{matrix} n \\ \downarrow \\ \text{Input} \end{matrix} \quad \begin{matrix} f \\ \downarrow \\ \text{filter} \end{matrix} \Rightarrow 6 - 3 + 1 = 4$$

Filters = are used to get info out of images

Q, what is padding ; → As seen above $6 \times 6 \rightarrow 4 \times 4$

padding is used to prevent the loss of image information (pixels) during convolution operations.

✓ In short: Padding adds extra pixels (usually zeros) around the image to preserve size and information during convolution.

$$\begin{array}{c}
 \begin{array}{|c|c|c|c|c|c|c|} \hline
 0 & 0 & 0 & 0 & 1 & 1 & 1 & 1 \\ \hline
 0 & 0 & 0 & 0 & 1 & 1 & 1 & 1 \\ \hline
 0 & 0 & 0 & 0 & 1 & 1 & 1 & 1 \\ \hline
 0 & 0 & 0 & 0 & 1 & 1 & 1 & 1 \\ \hline
 0 & 0 & 0 & 0 & 1 & 1 & 1 & 1 \\ \hline
 0 & 0 & 0 & 0 & 1 & 1 & 1 & 1 \\ \hline
 0 & 0 & 0 & 0 & 1 & 1 & 1 & 1 \\ \hline
 0 & & & & & & & \\ \hline
 \end{array} \\
 n=6 \\
 \downarrow \\
 p=1
 \end{array}
 \quad
 \begin{array}{c}
 f=3 \\
 * \\
 \begin{array}{|c|c|c|} \hline
 +1 & 0 & -1 \\ \hline
 +2 & 0 & -2 \\ \hline
 +1 & 0 & -1 \\ \hline
 \end{array} \\
 3 \times 3 \quad \text{filter} \\
 \text{filters / Kernel} \\
 \text{vertical edge filters}
 \end{array}
 \quad
 \begin{array}{c}
 6-3+2p+1=6 \\
 \hline
 \end{array}$$

$$\begin{array}{c}
 \text{new formula} \\
 = n + 2p - f + 1 \\
 = 6 + 2(1) - 3 + 1 \\
 = 8 - 3 + 1 = 6
 \end{array}$$

Zero Padding

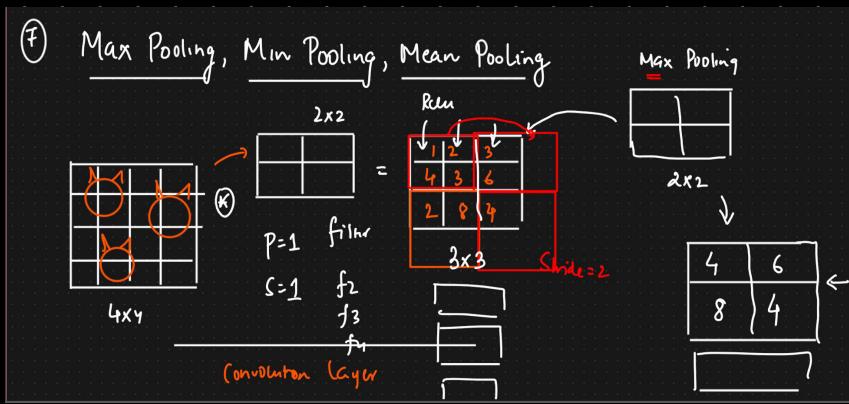
Pads with zeros around the image. Most common.

Reflect Padding

Pads with reflected values of the border.

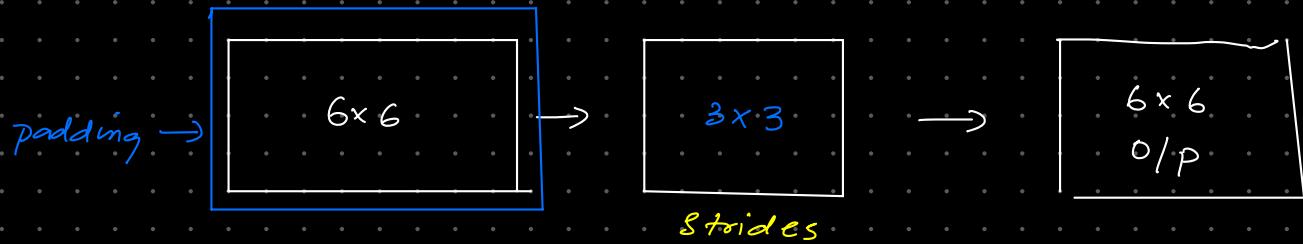
POOLING = is a technique in CNN to reduce dimensionality

Types =



Conclusion

1. Input Image
2. Normalization $(0-1)$ where $1 = 255$
3. Convolution method (apply filter / kernel)



4. Activation function (ReLU)
 5. Pooling \rightarrow to reduce dimensionality
- Max, Min or Mean pooling

6. Flattening \rightarrow elongating pooling output

$$\begin{array}{|c|c|} \hline 2 & 5 \\ \hline 1 & 6 \\ \hline \end{array} = \begin{array}{c} 2 \\ 5 \\ 1 \\ 6 \end{array}$$

7. Finally \rightarrow ANN classification (Dense layers).

CONCLUSION

✓ Step-by-Step: Your CNN Understanding

✓ 1. Input Image

Let's say a grayscale image: 6×6 matrix of pixel values (0–255)

✓ Correct.

✓ 2. Normalization

Normalize pixel values: $255 \rightarrow 1, 0 \rightarrow 0$

✓ Correct. Sometimes values are scaled between $0\text{--}1$ or $-1\text{--}1$.

✓ 3. Apply Convolution (Filter/Kernel 3x3)

To detect edges, textures, corners, etc.

✓ Correct.

Input Image (6×6) → Normalize →

Convolution (3x3 filter) →

ReLU →

Pooling (e.g., Max Pooling) →

Flatten →

Fully Connected Layers (ANN) →

Output with Activation (sigmoid/softmax)

✗ 4. You said: Output will be 6×6

Not always — usually the output size shrinks (unless you use padding).

Formula (without padding):

arduino

Copy Edit

$$\begin{aligned} \text{Output size} &= (\text{Input size} - \text{Kernel size}) + 1 \\ &= (6 - 3) + 1 = 4 \\ \rightarrow \text{So, output} &= 4 \times 4 \end{aligned}$$

✓ But if padding = same, then output stays 6×6 .

V.v. imp

? CONFUSION: Why ReLU after Convolution?

Let's clarify:

✓ After each convolution, we apply ReLU (or another activation function).

🔍 Why?

Because convolution may produce negative values, and ReLU keeps only the positive ones:

lua

Copy Edit

```
ReLU(x) = max(0, x)
```

This introduces non-linearity into the model. Without this, your entire CNN would just be one big linear operation, and it wouldn't be able to learn complex patterns.

So yes — we apply ReLU right after convolution, before pooling.

✓ 5. Pooling Layer (e.g., Max Pooling)

To reduce dimensions, keep important features

✓ Correct. Typically reduces size (e.g., $4 \times 4 \rightarrow 2 \times 2$)

✓ 6. Flatten the Output

Convert the 2D/3D tensor to 1D vector for input into fully connected (ANN) layer

✓ Correct.

✓ 7. Apply Fully Connected (Dense) Layers

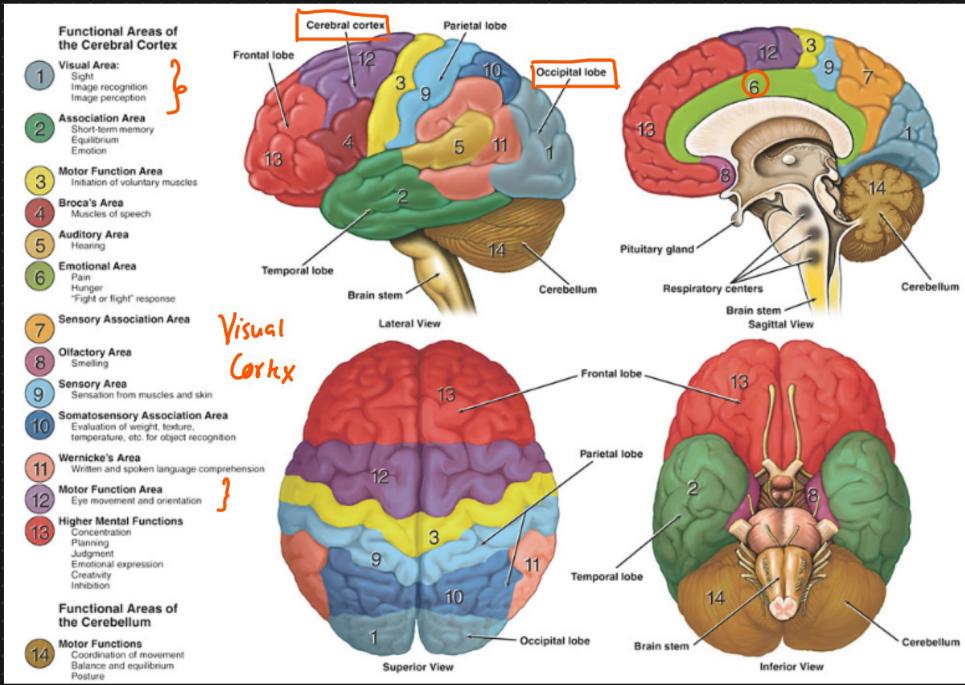
Feed into ANN, apply activation functions (like ReLU, softmax, etc.)

✓ Correct.

✓ Final Layer Activation:

- For binary classification: sigmoid
- For multi-class classification: softmax

Convolutional Neural N/w



<https://www.dana.org>

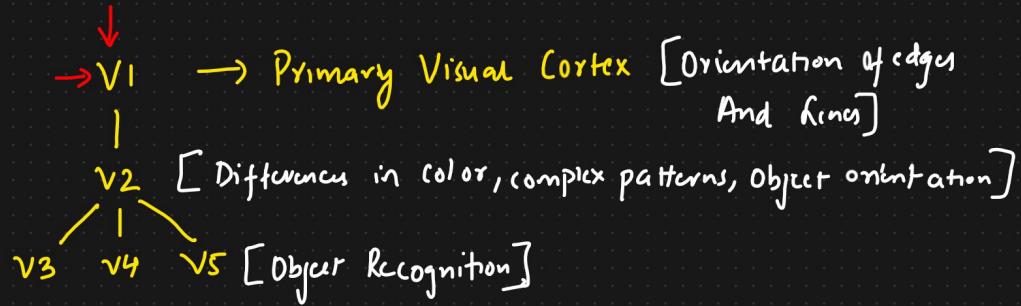


Dataset : I/p features O/p

② CNN : I/p ⇒ Images Eg: Image classification,
 Object Detection, Segmentation

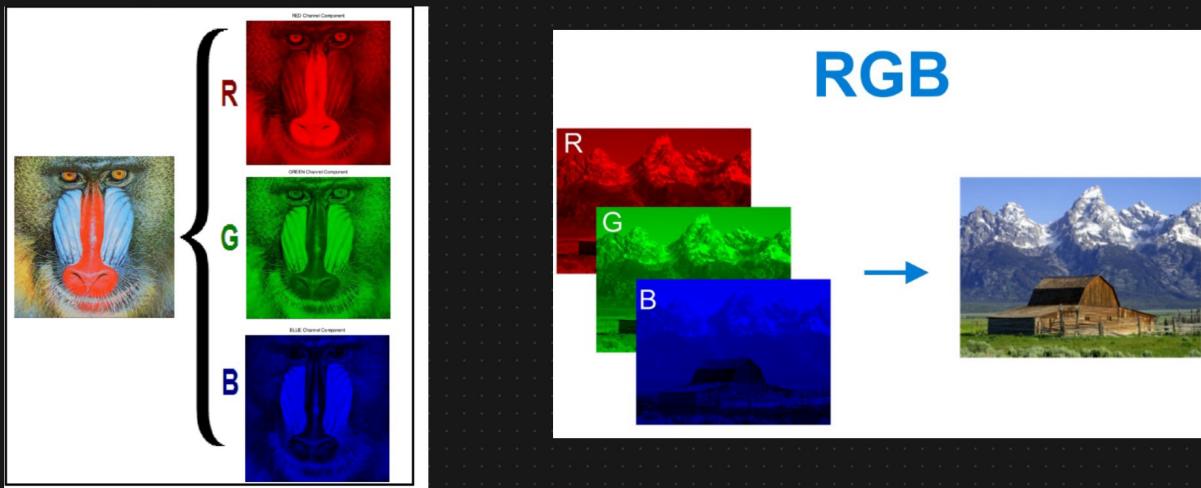
② Cerebral Cortex And Visual Cortex

Visual Cortex (V1-V5) [Region of the brain that receives, integrates and processes visual information relayed from the retinas].

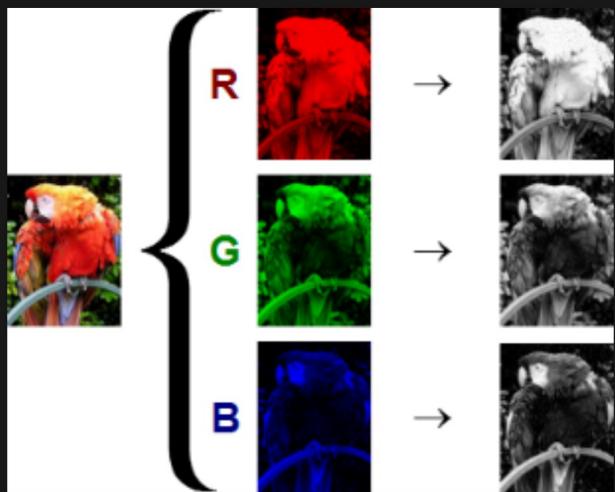


Visualize the Image

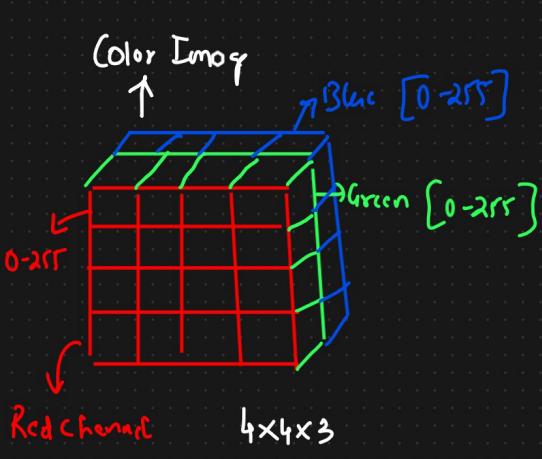
③ RGB Images And Gray Scale Images



<https://www.researchgate.net/>



$0-255 \leftarrow$ \Rightarrow Gray Scale Image.



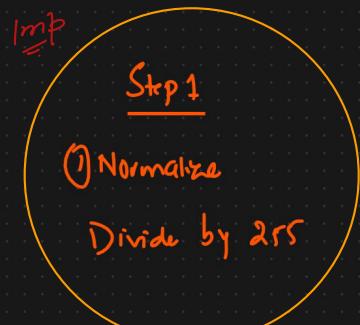
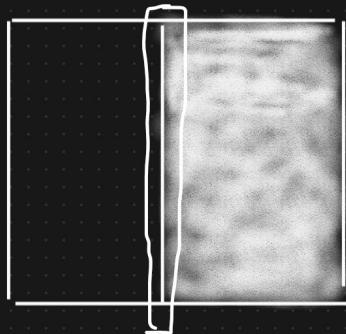
<https://commons.wikimedia.org/>

④ Convolution Operation In CNN

$\rightarrow (0,1)$

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

\Rightarrow



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

$6 \times 6 \times 1$

Imp Convolution operation

$s\text{tride}=1$

0	0	0	1	1	1
0	0	0	1	1	1
0	0	0	1	1	1
0	0	0	1	1	1
0	0	0	1	1	1
0	0	0	1	1	1

$n=6$

$f=3$

$o/p=4$

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

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

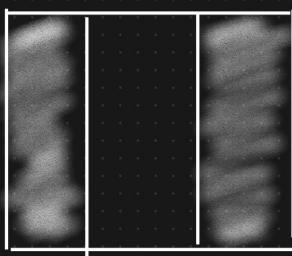
3×3

filters

Vertical edge filters

4×4

$$\begin{aligned} h - f + 1 &= \\ &= 6 - 3 + 1 = 4 \end{aligned}$$



arr	0	0	2rr
2rr	0	0	2rr
2rr	0	0	2rr
2rr	0	0	2rr

⑤ Padding In CNN

0	0	0	0	1	1	1	1
0	0	0	0	1	1	1	1
0	0	0	0	1	1	1	1
0	0	0	0	1	1	1	1
0	0	0	0	1	1	1	1
0	0	0	0	1	1	1	1
0	0	0	0	1	1	1	1
0							

6×6

$n=6$

8×8

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

$n-f+2p+1$

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

6×6

\equiv

$$6 - 3 + 2p + 1 = 6$$

$$3 + 2p + 1 = 6$$

$$2p = 6 - 4$$

$$p = \frac{2}{2} = 1$$

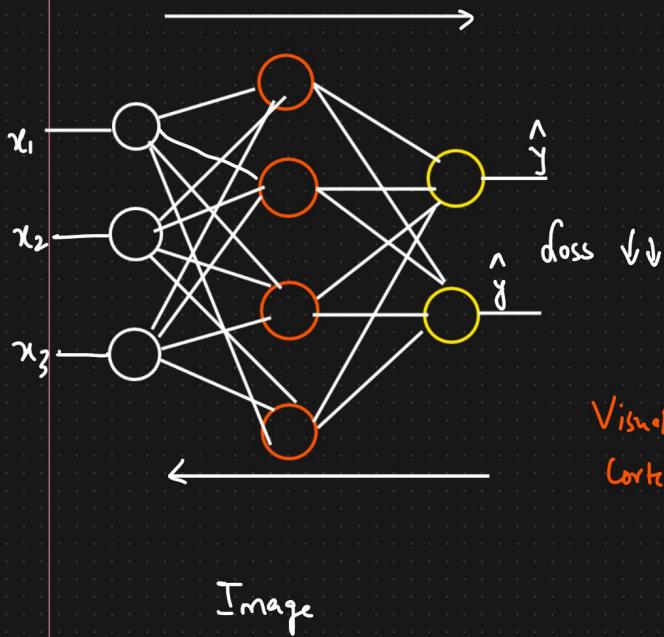
7×7

3×3

7×7

How much padding you need to apply?

⑥ Operation of CNN Vs ANN



$$z = w^T x_i + b$$

$$\text{ReLU}(z)$$

Vision
Cortex



→ ReLU operation $\max(0, x)$

$$\begin{array}{c} \downarrow \quad \downarrow \quad \downarrow \quad \downarrow \\ \begin{array}{|c|c|c|c|} \hline 0 & -4 & -4 & 0 \\ \hline 0 & -4 & -4 & 0 \\ \hline 0 & -4 & -4 & 0 \\ \hline 0 & -4 & -4 & 0 \\ \hline \end{array} \end{array}$$

$$\begin{array}{c} \begin{array}{|c|c|c|c|} \hline - & - & - & - \\ \hline - & - & - & - \\ \hline - & - & - & - \\ \hline - & - & - & - \\ \hline \end{array} \end{array}$$

*

$$\begin{array}{|c|c|c|} \hline +1 & 0 & -1 \\ \hline +2 & 0 & -2 \\ \hline +1 & 0 & -1 \\ \hline \end{array}$$

f_1

$$\begin{array}{|c|c|} \hline & \\ \hline & \\ \hline \end{array}$$

f_2

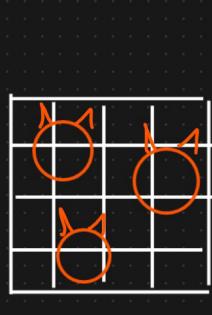
f_3

\vdots

f_n

Convolution Layer

⑦ Max Pooling, Min Pooling, Mean Pooling



2×2

\Rightarrow

$$\begin{array}{|c|c|} \hline & \\ \hline & \\ \hline \end{array}$$

\Downarrow

$P=1$ filter
 $S=1$

Convolution Layer

$$\begin{array}{c} \text{ReLU} \\ \hline \begin{array}{|c|c|c|} \hline 1 & 2 & 3 \\ \hline 4 & 3 & 6 \\ \hline 2 & 8 & 4 \\ \hline \end{array} \end{array}$$

3×3 stride=2

$$\begin{array}{|c|c|} \hline & \\ \hline & \\ \hline \end{array}$$

2×2

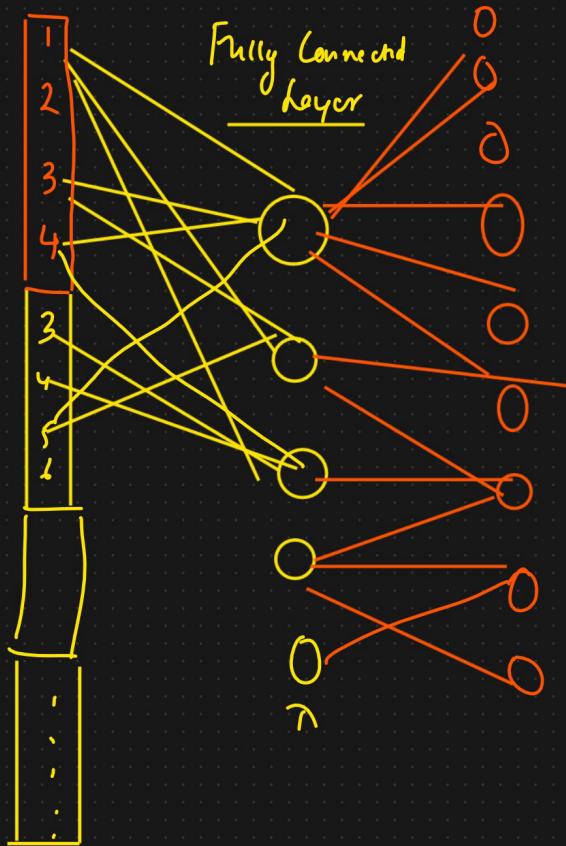
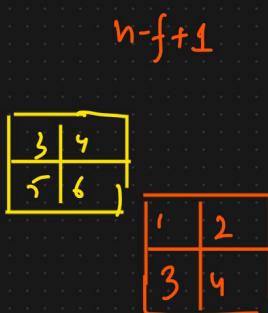
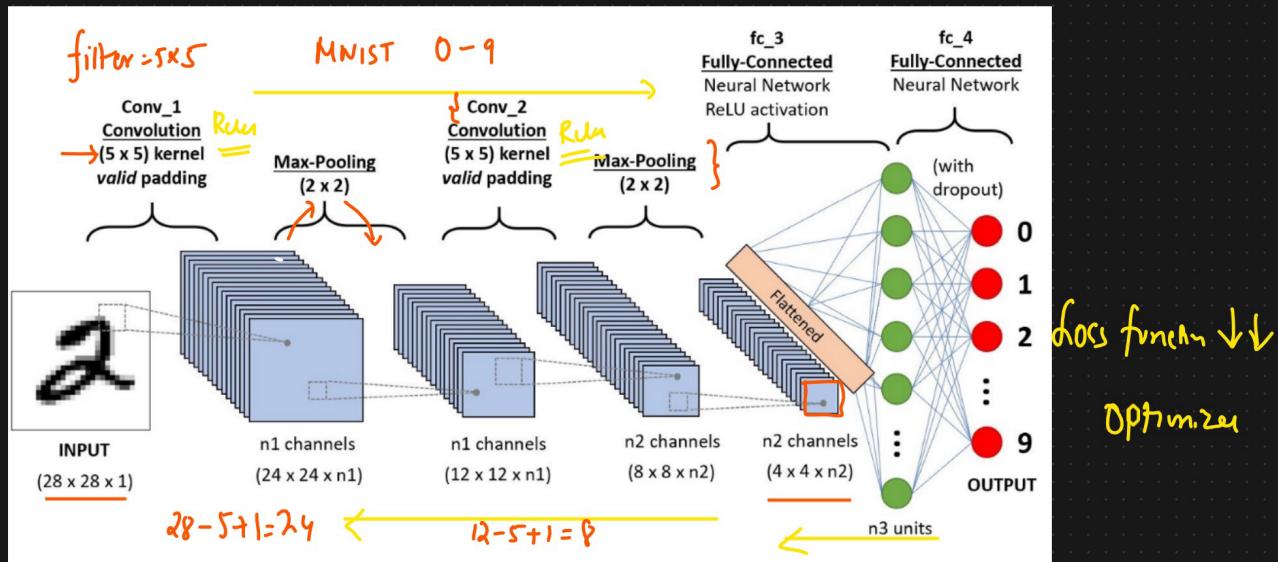
\Downarrow

$$\begin{array}{c} \begin{array}{|c|c|} \hline 4 & 6 \\ \hline 8 & 4 \\ \hline \end{array} \end{array}$$

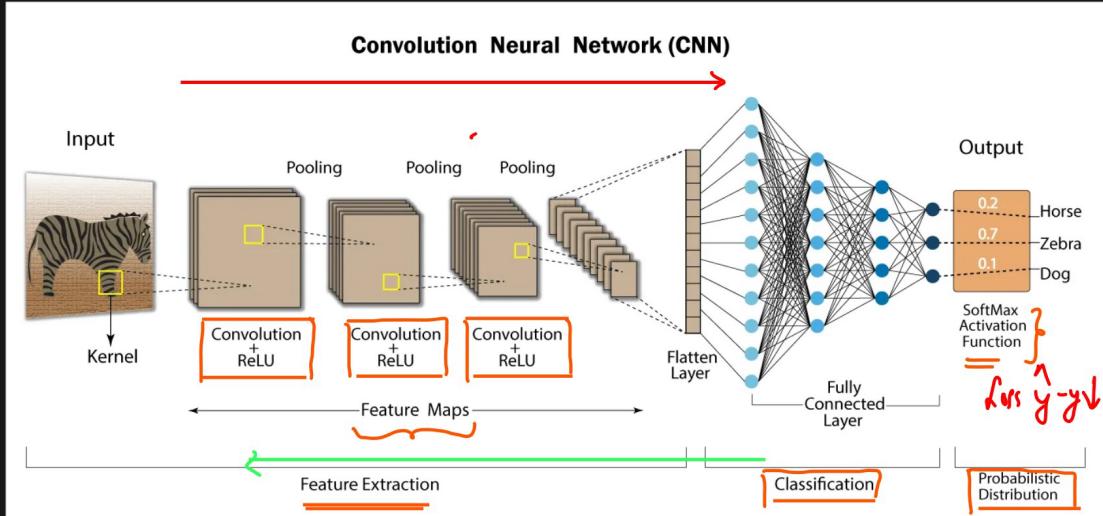
Max Pooling

⑦ Location Invariant

⑧ Fully Connected Layer In CNN [Flattened Layer]



④ CNN Complete Example



<https://developersbreach.com/convolution-neural-network-deep-learning/>

0	0	0		
0	0	0		
0	0	0		
0	0	0		
0	0	0		

*

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

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