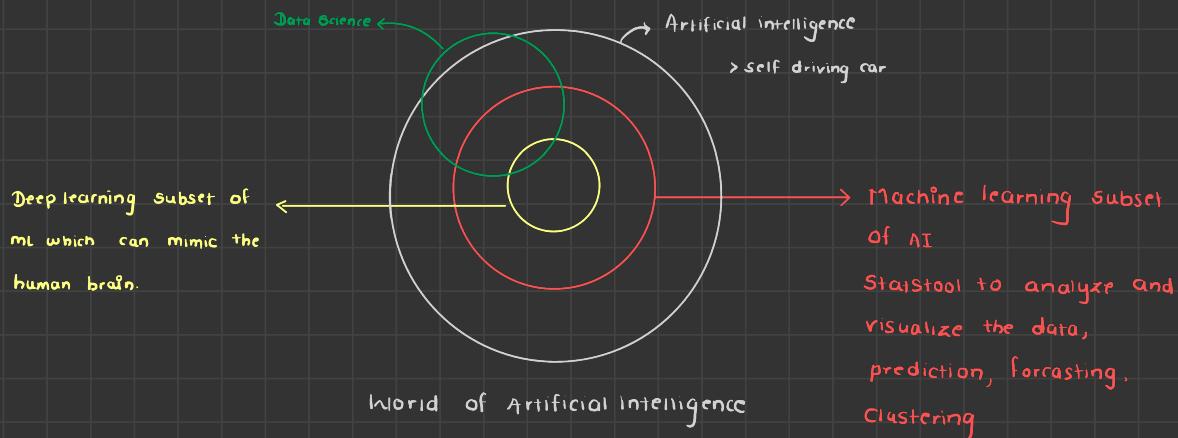
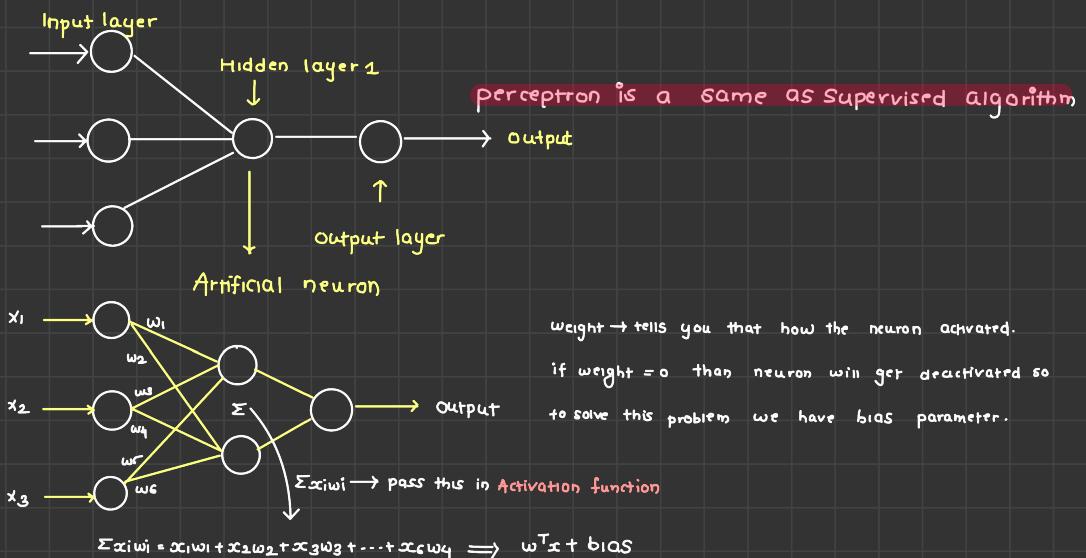


— DEEP LEARNING —

AI VS ML VS DL VS DS



Perceptron → A single layer neural network



Activation function

Basically activation decide the whether the neuron should be activated or not

lets take sigmoid activation function (Binary classification)

$$\sigma = \frac{1}{1 + e^{-x}} = \frac{1}{1 + (x \text{ with})} \Rightarrow \text{output b/w 0 to 1}$$

threshold value = 0.5

$\sigma \geq 0.5 \rightarrow \text{Activated neuron}$

$\sigma < 0 \rightarrow \text{deactivated neuron}$

Forward propagation

Passing inputs with weight through activation function to generates output.

Inputs \rightarrow Activation \rightarrow output

Backward propagation

Suppose we got the output 0 but actual it is 1

Now calculate error or difference between $y=1$ and $\hat{y}=0$

: loss function = $y - \hat{y}$

= 1 - 0

= 1 \rightarrow it should be a minimal

To reduce this we need to update the weights & biases

Reverse the process of forward propagation

We need optimization algo to update weights

GRADIENT DESCENT

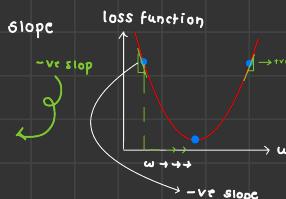
update formula for weight

$$w_{\text{new}} = w_{\text{old}} - \alpha \frac{\partial L}{\partial w_{\text{old}}} \rightarrow \text{derivation of loss with respect to } w_{\text{old}}$$

$$\frac{\partial L}{\partial w_{\text{old}}} \rightarrow \text{calculating slope}$$

$$w_{\text{new}} = w_{\text{old}} - \alpha (-\text{ve})$$

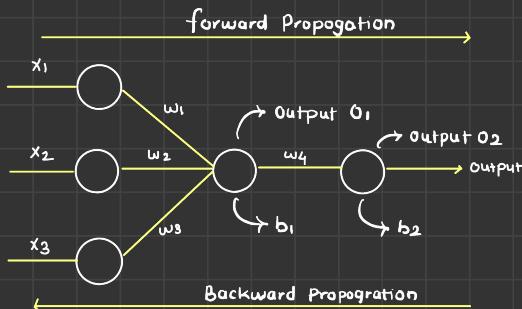
$$w_{\text{new}} = w_{\text{old}} + \alpha (\text{ve})$$



$$w_{\text{new}} = w_{\text{old}} - \alpha (-\text{ve})$$

$$w_{\text{new}} = w_{\text{old}} - \alpha (\text{ve})$$

Chain rule of differentiation:



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

lets say that we want to update w_4

$$w_{4 \text{ new}} = w_{4 \text{ old}} - \eta \frac{\partial L}{\partial w_{4 \text{ old}}} \quad \text{By chain rule}$$

$$\text{Weight} \rightarrow \frac{\partial L}{\partial w_{4 \text{ old}}} = \boxed{\frac{\partial L}{\partial o_2} * \frac{\partial o_2}{\partial w_{4 \text{ old}}}} = \frac{\partial L}{\partial w_{4 \text{ old}}}$$

want to update bias

$$b_{2 \text{ new}} = b_{2 \text{ old}} - \eta \frac{\partial L}{\partial b_{2 \text{ old}}} \quad \text{By chain rule}$$

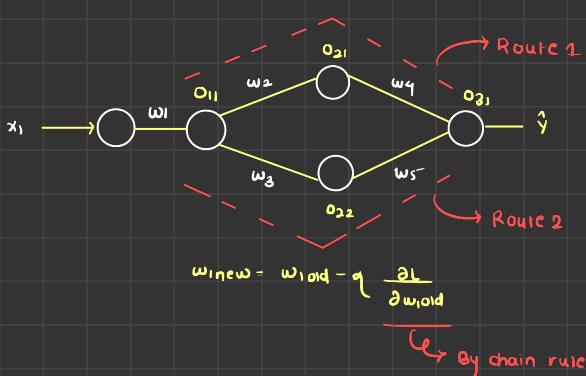
$$w_{1 \text{ new}} = w_{1 \text{ old}} - \eta \frac{\partial L}{\partial w_{1 \text{ old}}} \quad \text{By chain rule}$$

$$\frac{\partial L}{\partial w_{1 \text{ old}}} = \boxed{\frac{\partial L}{\partial o_2} * \frac{\partial o_2}{\partial o_1} * \frac{\partial o_1}{\partial w_{1 \text{ old}}}}$$

$$w_{2 \text{ new}} = w_{2 \text{ old}} - \eta \frac{\partial L}{\partial w_{2 \text{ old}}} \quad \text{By chain rule}$$

$$\frac{\partial L}{\partial w_{2 \text{ old}}} = \frac{\partial L}{\partial o_2} * \boxed{\frac{\partial o_2}{\partial o_1} * \frac{\partial o_1}{\partial w_{2 \text{ old}}}} \quad \text{By chain rule}$$

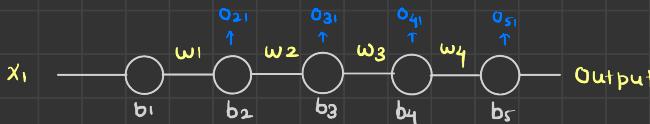
$$\frac{\partial o_2}{\partial w_4} * \frac{\partial w_4}{\partial o_1}$$



$$\frac{\partial L}{\partial w_{1\text{old}}} = \left[\frac{\partial L}{\partial o_{31}} * \frac{\partial o_{31}}{\partial o_{21}} * \frac{\partial o_{21}}{\partial o_{11}} * \frac{\partial o_{11}}{\partial w_{1\text{old}}} \right] + \left[\frac{\partial L}{\partial o_{32}} * \frac{\partial o_{32}}{\partial o_{22}} * \frac{\partial o_{22}}{\partial o_{11}} * \frac{\partial o_{11}}{\partial w_{1\text{old}}} \right]$$



VANISHING GRADIENT PROBLEM

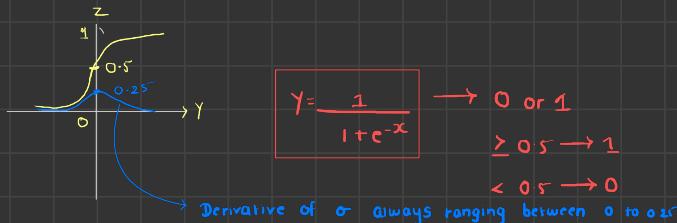


$$\text{mse loss} = \frac{1}{2} (y - \hat{y})^2$$

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

$$\frac{\partial L}{\partial w_{1\text{new}}} = \frac{\partial L}{\partial o_5} * \frac{\partial o_5}{\partial o_4} * \frac{\partial o_4}{\partial o_3} * \frac{\partial o_3}{\partial o_2} * \frac{\partial o_2}{\partial o_1} * \frac{\partial o_1}{\partial w_{1\text{old}}}$$

Sigmoid activation



lets say we got derivative of $\sigma = 0.25$

$0.25 \rightarrow$ this will get reduce in backward propagation

$$0.25 * 0.15 * 0.10 * 0.05 * 0.02 = \text{very small value}$$

$$w_{1\text{new}} = w_{1\text{old}} - \eta (\text{very small value})$$

↳ also small

Than the result $\rightarrow w_{new} = w_{old} - 0$

\hookrightarrow Not changing any more

And this is problem called **Vanishing gradient descent** that lead not updating weights any more.

To solve this problem we have different types of activation functions.

1.) Tanh

2.) Relu

3.) Leaky Relu

and many more are there.

A SIGMOID ACTIVATION FUNCTION

Advantages:

- > Convergence become more faster due to smoothness of curve
- > Output between 0 to 1.
- > Also give clear prediction.

Disadvantage:

- > prone to vanishing gradients.
- > function output is not zero-centered.
- > computationally cost \rightarrow perform exponential operation

Derivation of Sigmoid:

Let's denote the sigmoid function as the following.

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

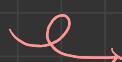
Another way to express the sigmoid function

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

You can easily derive the second equation from the first equation:

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

Since $\frac{d}{dx} = 1$, so in essence, we're just multiplying $\frac{1}{1+e^{-x}}$ by 1.



Sigmoid derivative

The derivative of the sigmoid function $\sigma(x)$ is the sigmoid function $\sigma(x)$ multiplied by $1 - \sigma(x)$.

$$\sigma'(x) = \frac{1}{1+e^{-x}} \cdot \frac{d}{dx}(e^{-x}) = e^{-x}(1 - \sigma(x))$$

Before we begin, here is a reminder of how to find the derivative of exponentials for dummies:

$$\frac{d}{dx} e^{kx} = k e^{kx}$$



Sigmoid derivative via chain rule

Chain rule: $\frac{d}{dx}[f(g(x))] = f'[g(x)] \cdot g'(x)$.

Example: Find the derivative of $f(x) = (x^2 + 1)^3$:

$$f'(x) = 3(x^2 + 1)^{3-1} \cdot 2x^2 \cdot 1 \\ = 3(x^2 + 1)^2(2x) \\ = 6x(x^2 + 1)^2$$



sigmoid derivation below uses this rule.

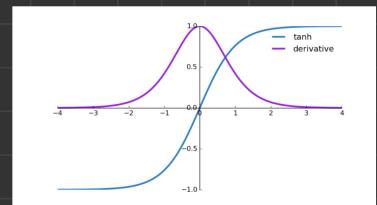
$$\begin{aligned} \frac{d}{dx} \sigma(x) &= \frac{d}{dx} \left[\frac{1}{1+e^{-x}} \right] = \frac{d}{dx} (1+e^{-x})^{-1} \\ &= -1 * (1+e^{-x})^{-2} (-e^{-x}) \\ &= \frac{-e^{-x}}{(1+e^{-x})^2} \\ &= \frac{1}{1+e^{-x}} \frac{e^{-x}}{1+e^{-x}} \\ &= \frac{1}{1+e^{-x}} \frac{e^{-x} + (1-1)}{1+e^{-x}} \\ &= \frac{1}{1+e^{-x}} \frac{(1+e^{-x}) - 1}{1+e^{-x}} \\ &= \frac{1}{1+e^{-x}} \left[\frac{(1+e^{-x})}{1+e^{-x}} - \frac{1}{1+e^{-x}} \right] \\ &= \frac{1}{1+e^{-x}} \left[1 - \frac{1}{1+e^{-x}} \right] \\ &= \sigma(x)(1 - \sigma(x)) \end{aligned}$$

B. TANH ACTIVATION FORMULA

function output lies between -1 to 1

derivative of tanh → 0 to 1

Basically when you create deep neural network it also creates problem of vanishing gradient descent.



$$\tanh(x) = \frac{e^x - e^{-x}}{e^x + e^{-x}}$$

$$\tanh(x) = \frac{e^x - e^{-x}}{e^x + e^{-x}}$$

it can be written as trigonometric function

$$\tanh(x) = \frac{\sinh(x)}{\cosh(x)}$$

first order derivative is as follows

$$\frac{d}{dx} \tanh(x) = \frac{d}{dx} \frac{\sinh(x)}{\cosh(x)}$$

Applying rule of differentiation rule of division

$$\frac{d}{dx} \frac{\sinh(x)}{\cosh(x)} = \frac{\cosh(x) \frac{d}{dx} \sinh(x) - \sinh(x) \frac{d}{dx} \cosh(x)}{\cosh^2(x)}$$

Division Rule

$$\frac{dy}{dx} = \frac{v \frac{du}{dx} - u \frac{dv}{dx}}{v^2}$$

$$\frac{d}{dx} \frac{\sinh(x)}{\cosh(x)} = \frac{\cosh^2 x - \sinh^2 x}{\cosh^2 x} = \frac{\cosh^2 x - \sinh^2 x}{\cosh^2 x} = 1 - \tanh^2(x)$$

3. ReLU function: Rectified Linear Unit

$$\text{ReLU} = \max(0, x)$$

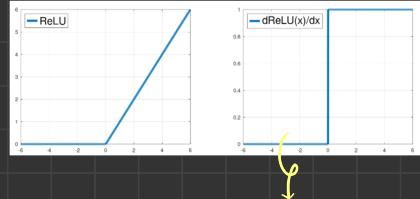
Whenever x is negative it return 0 otherwise it return 1

derivative of relu → 0 to 1

When the input is neg → ReLU will die

Calculation is faster than tanh and sigmoid function

ReLU is not 0-centric function

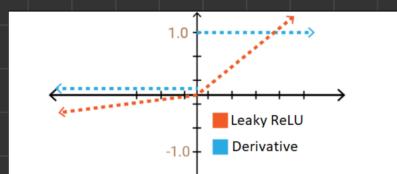


derivation of relu

4. LEAKY RELU FUNCTION:

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

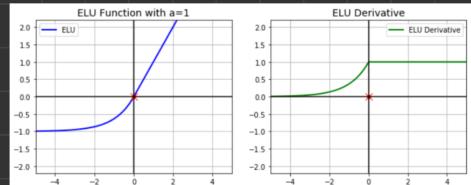
Solve dead neuron problem



5 ELU (Exponential linear units) function:

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

- Computationally expensive
- close to zero-centered
- No dead ReLU issues

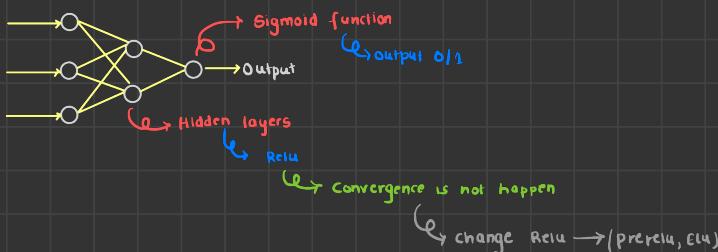


6 SOFTMAX Activation function

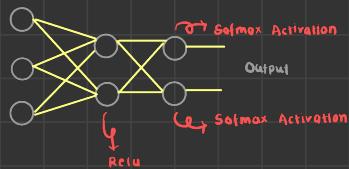
Vectorized version of sigmoid activation function whereas sigmoid take scalar value.

Technique to choose activation function

Binary Classification



Multiclass Classification

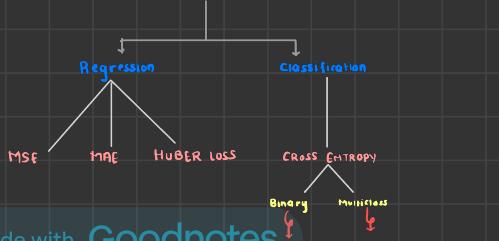


Regression



Loss Functions

Deep Learning (AI)



Cost function vs loss function

↳ loss cause by batch data

→ not whole data but few data from whole data you have.

Overall loss cause by the data you pass into a forward propagation

$$\text{Loss function} = \frac{1}{2} (y - \hat{y})^2 \quad (\rightarrow \text{MSE})$$

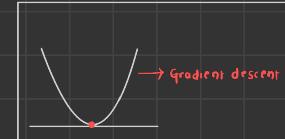
$$\text{Cost function} = \frac{1}{2} \sum_{i=1}^n (y_i - \hat{y}_i)^2$$

Mean Squared Error (MSE):

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

↳ quadratic equation

$$(a - b)^2 = a^2 - 2ab + b^2 \\ ax^2 + bx + c$$



Advantage:

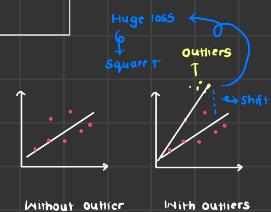
differentiable

It has only local or global minima

Convergence much faster

Disadvantage:

Not robust to outliers



Mean Absolute Error (MAE)

$$\text{loss function} = \frac{1}{2} |y - \hat{y}| \quad \text{cost function} = \frac{1}{2} \sum_{i=1}^n |y_i - \hat{y}_i|$$

Advantage:

Robust to outliers

Disadvantage:

Time consuming



HUBER LOSS

Combination of MSE and MAE

$$\text{Loss} : \begin{cases} \frac{1}{2} (y - \hat{y})^2 & \text{if } |y - \hat{y}| \leq \delta \\ \delta |y - \hat{y}| - \frac{1}{2} \delta^2 & \text{otherwise} \end{cases} \quad \begin{matrix} \text{When outliers are not present} \\ \hookrightarrow \text{hyperparameter} \end{matrix}$$

Binary Cross Entropy

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

↓

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

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

Categorical Cross Entropy

$f_1 \quad f_2 \quad f_3 \quad \text{output} \rightarrow \text{multiclass classification}$

2 3 5 Good

6 0 1 Neutral

1 4 1 Bad

↳ convert using one hot encoding

| Good | Neutral | Bad |
|------|---------|-----|
| 1 | 0 | 0 |
| 0 | 1 | 0 |
| 0 | 0 | 1 |

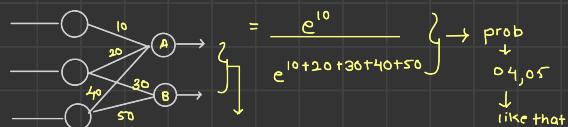
$$L(x_i, y_i) = - \sum_{j=1}^C y_{ij} * \ln(\hat{y}_{ij}) \quad i = \text{rows}, j = \text{column}$$

$$y_i = [y_{i1}, y_{i2}, y_{i3}, \dots, y_{ic}]$$

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

$\hat{y}_{ij} = \text{Softmax Activation layers} \rightarrow \text{output layer}$

$$\text{Softmax}(z) = \sigma(z) = \frac{e^{z_i}}{\sum_{j=1}^K e^{z_j}}$$



Sum of these = 1

if A's prob is high, then it is output

Conclusion.

ReLU, Softmax \rightarrow multiclass classification

ReLU, Sigmoid \rightarrow Binary classification

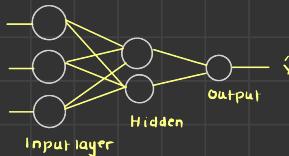
ReLU, Linear activation \rightarrow regression

↳ loss function

OPTIMIZER

GRADIENT DESCENT

weight update formula $\rightarrow w_{\text{new}} = w_{\text{old}} - \eta \frac{dL}{dw_{\text{old}}}$



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

\hookrightarrow To minimize the loss we need optimizers

Epochs

\hookrightarrow Training neural network with training data within one cycle

Disadvantage

Resource extensive (Huge RAM)

\hookrightarrow In gradient descent

STOCHASTIC GRADIENT DESCENT

Suppose I have 2000000's data

Advantage

Require less RAM

Disadvantage:

Convergence will take too much time

epoch 1

1 records $\rightarrow \hat{y}$
 \leftarrow calculate loss than update weights \hookrightarrow 1st iteration

epoch 2

2 records $\rightarrow \hat{y}$
 \leftarrow calculate loss than update weights \hookrightarrow 2nd iteration

|
|
|
|
|
n numbers of epochs

MINI BATCH SGD

Instead of passing 1000000 record once

\hookrightarrow BATCH SIZE = 10000

In every epoch I will send 10000

Advantage:

less resources intensive

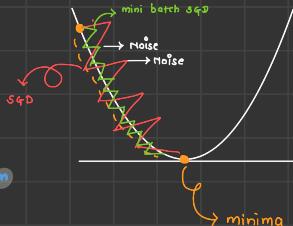
Convergence will be a faster

Time Complexity \rightarrow improved

Aim To reach global minima

How remove noise ?

↳ we use concept of momentum



SGD with Momentum

Momentum makes the smoother the zig zag curve

Exponential weighted average

↳ Time Series → ARIMA, ARMA

Update Rule

$$w_{\text{new}} = w_{\text{old}} - \eta \frac{\partial L}{\partial w_{\text{old}}} , \quad b_{\text{new}} = b_{\text{old}} - \eta \frac{\partial L}{\partial b_{\text{old}}}$$

$$w_t = w_{t-1} - \eta \frac{\partial L}{\partial w_{t-1}} \quad t \Rightarrow \text{current time}$$

$t-1 \Rightarrow \text{previous time}$

Exponential weighted average

$$\begin{matrix} t_1 & t_2 & t_3 & t_4 & \dots & t_n \\ a_1 & a_2 & a_3 & a_4 & & a_n \end{matrix}$$

$$v_{t_4} = a_1$$

Value of $t_1 = a_1$

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

β - hyperparameter $\rightarrow 0$ to 1

↳ θ goes

$$v_{t_2} = 0.5 * v_{t_1} + 0.5 * a_2$$

↳ We are giving more importance to a_2

Removing noise and smoothening the curve

Exponential weighted avg

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

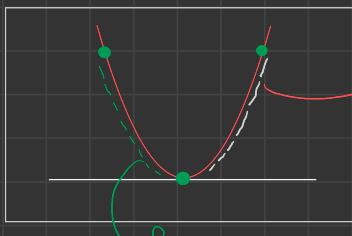
↳ value of derivation

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



- Reducing noise
- Also work for mini batch
- convergence will be faster

Ayngrad (Adaptive Gradient Descent)



we are changing the speed initially if it is too fast taking big step as it move to closer the global minima taking smaller steps.

This step control by learning rate
which constant in entire process

$$\omega_L = \omega_{t-1} - \alpha \frac{\partial L}{\partial \omega_{t-1}}$$

fixed value

Change → we will make it adaptive

$$\omega_L = \omega_{t-1} - \alpha' \frac{\partial L}{\partial \omega_t}$$

$$\alpha' = \frac{\alpha}{\sqrt{\alpha_t + \epsilon}}$$

if we divide by

Small huge

Suppose this is zero you will get infinity

huge

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

This value keep on increasing

$$\begin{array}{ccccccccc} & & & & & l=1 & t=2 & t=3 \\ \alpha & = & 0.01 & & & & \alpha = 0.005 & & \alpha = 0.002 \\ \downarrow & & & & & & & & \\ \text{o} & & & & & & & & \text{reducing} \end{array}$$

and your weight will not going to update

ADAM and RMSprop

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

Applying weighted average

$S_{dw} = \beta$

Initialize $S_{dw} = 0$

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

$\beta = 0.95$

$$S_{dw_t} = 0.95 * S_{dw_{t-1}} + 0.05 \left(\frac{\partial L}{\partial w_{t-1}} \right)^2$$

controlling now not gonna very high

ADAM Optimizer (Best optimizer)

Momentum + RMSprop (Adaptive learning)

$$\begin{array}{ll} \sigma \rightarrow \text{bias} & \rho \rightarrow \text{bias} \\ Vdw=0, Vdb=0, Sdw=0, Sdb=0 \end{array}$$

$$\begin{aligned} w_t &= w_{t-1} - \eta' Vdw \\ b_t &= b_{t-1} - \eta' Vdb \end{aligned}$$

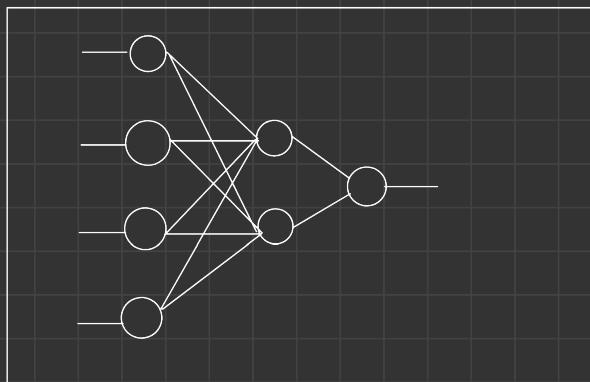
$$\eta' = \frac{\eta}{\sqrt{S_{dw} + \epsilon}}$$

$$Vdw_t = \beta * Vw_{t-1} + (1-\beta) \frac{\partial L}{\partial w_{t-1}}$$

$$Vdb_t = \beta * Vb_{t-1} + (1-\beta) \frac{\partial L}{\partial b_{t-1}}$$

- > Smoothening the curve to reduce noise make convergence faster.
- > Also making learning rate adaptive

ARTIFICIAL NEURAL NETWORK



Dense

neurons getting input from previous neuron.

Sequential

forward

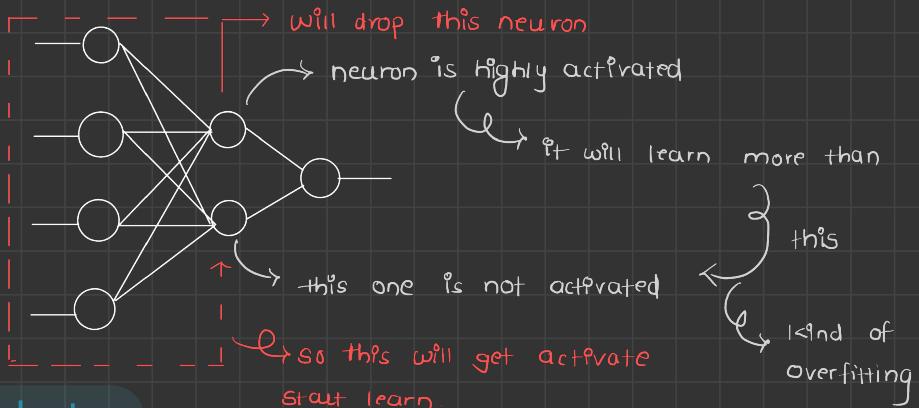
Backward

inputs ——> output layer

Activation

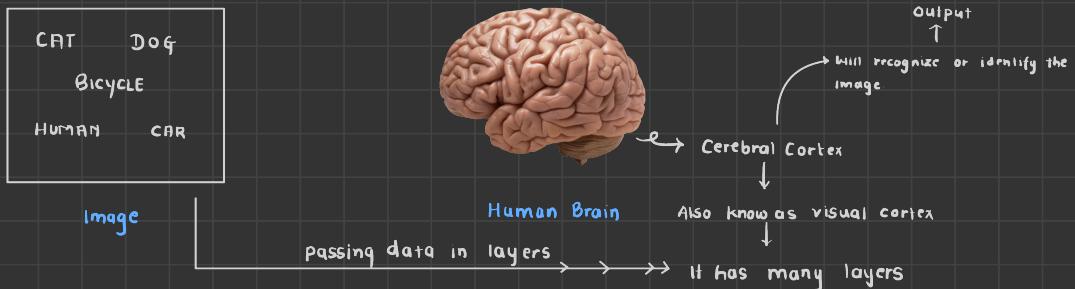
whether the neurons are activated or deactivated

Dropout

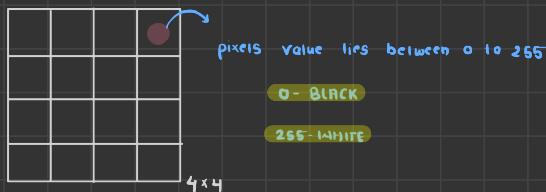
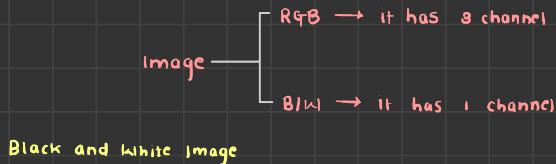


CONVOLUTION NEURAL NETWORK

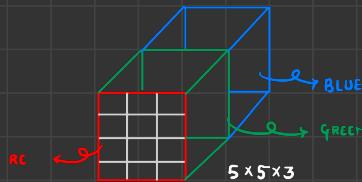
CNN vs Human brain

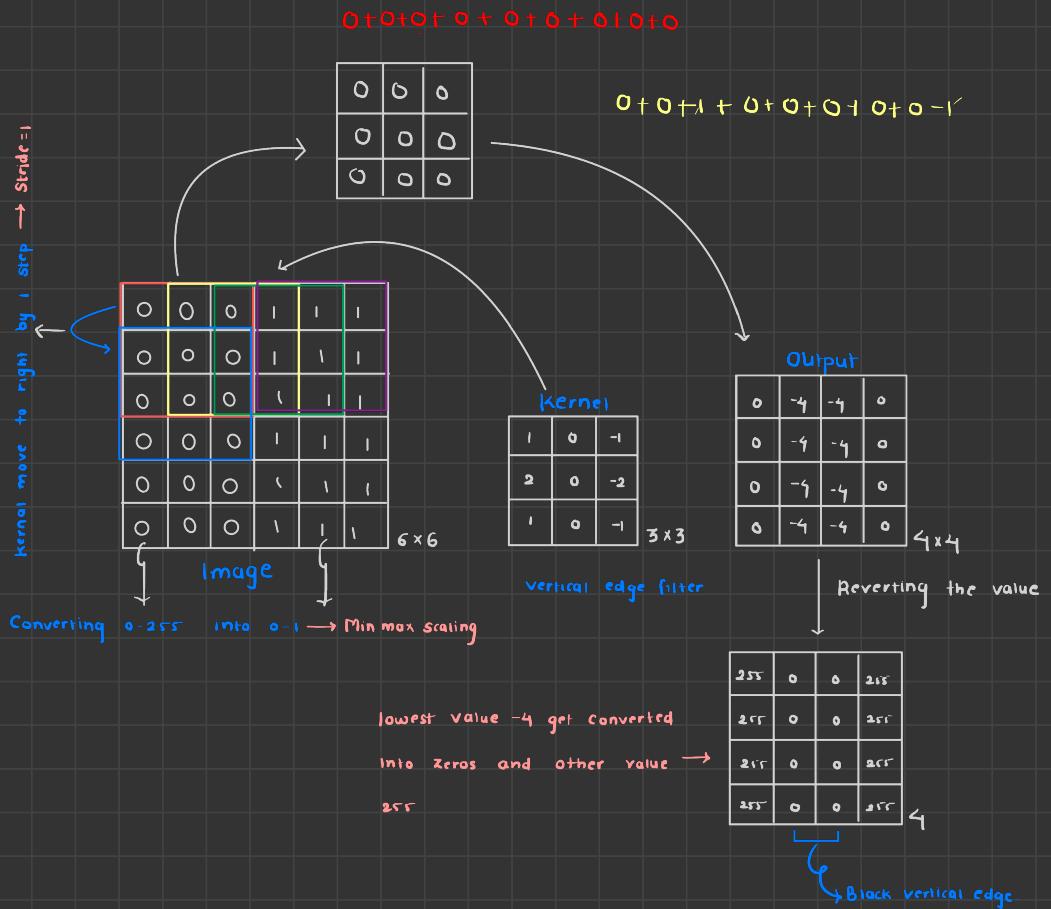


Convolution neural network

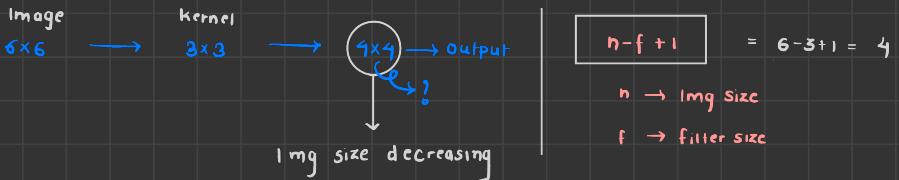


RGB Image





Basically we took image and applied filter over the image and got vertical line as output.

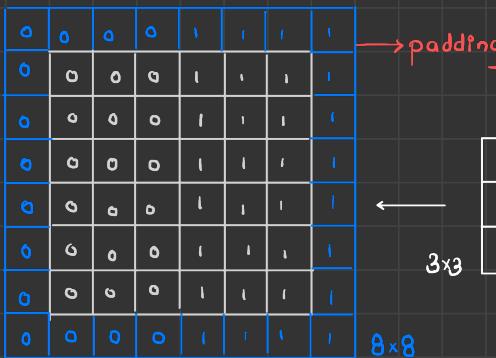


$6 \times 6 \longrightarrow 3 \times 3 \longrightarrow 4 \times 4$

We are losing information

To avoid this

PADDING



$$\begin{aligned} & n - f + 1 \\ &= 8 - 3 + 1 \\ &= 6 \end{aligned}$$

Output

Padding help to avoid the problem of losing information in convolution

Now after padding formula = $n + 2p - f + 1 \rightarrow \frac{n + 2p - f + 1}{s}$

n = input size

p = padding

f = kernel size

Stride

$$\begin{aligned} n &= 6, \quad k = 3, \quad p = 1 \\ & 6 + 2(1) - 3 + 1 \\ &= \underline{\underline{6}} \end{aligned}$$

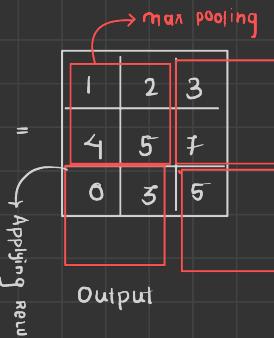
MAXPOOLING



Image



Kernel



Applying ReLU

Max pooling



| | |
|---|---|
| 5 | 7 |
| 3 | 5 |

most

prominent

Solve the problem of
location invariance

Convolution neural network

FLATTENING

like Anti Dense layer

| | |
|---|---|
| 5 | 7 |
| 3 | 5 |

