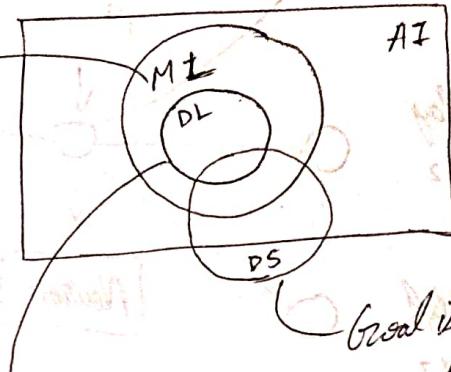


- AI :- Application which can do its own task without human intervention.
1. Netflix app (provides automatic recommendation).
 2. Self Driving cars Amazon app

AI which integrates w/ that is already existing to make the experience better.

Stats tool to analyse the data,
visualize data, predictions,
forecasting, clustering



Researchers (1958)

Main aim mimic the human brain

Goal is to
create AI
application.

⇒ Why Deep Learning is becoming popular?

① 2005 → FB, Insta, WhatsApp, LinkedIn, Twitter (We will be interacting with data)

DATA → Exponentially

2008 → {Big Data} → Efficiently

2013 → Company had huge amount of Data

{AI → popular}

Damless Products

Panasonic :- AC's, TV's, Refrigerator & Data's

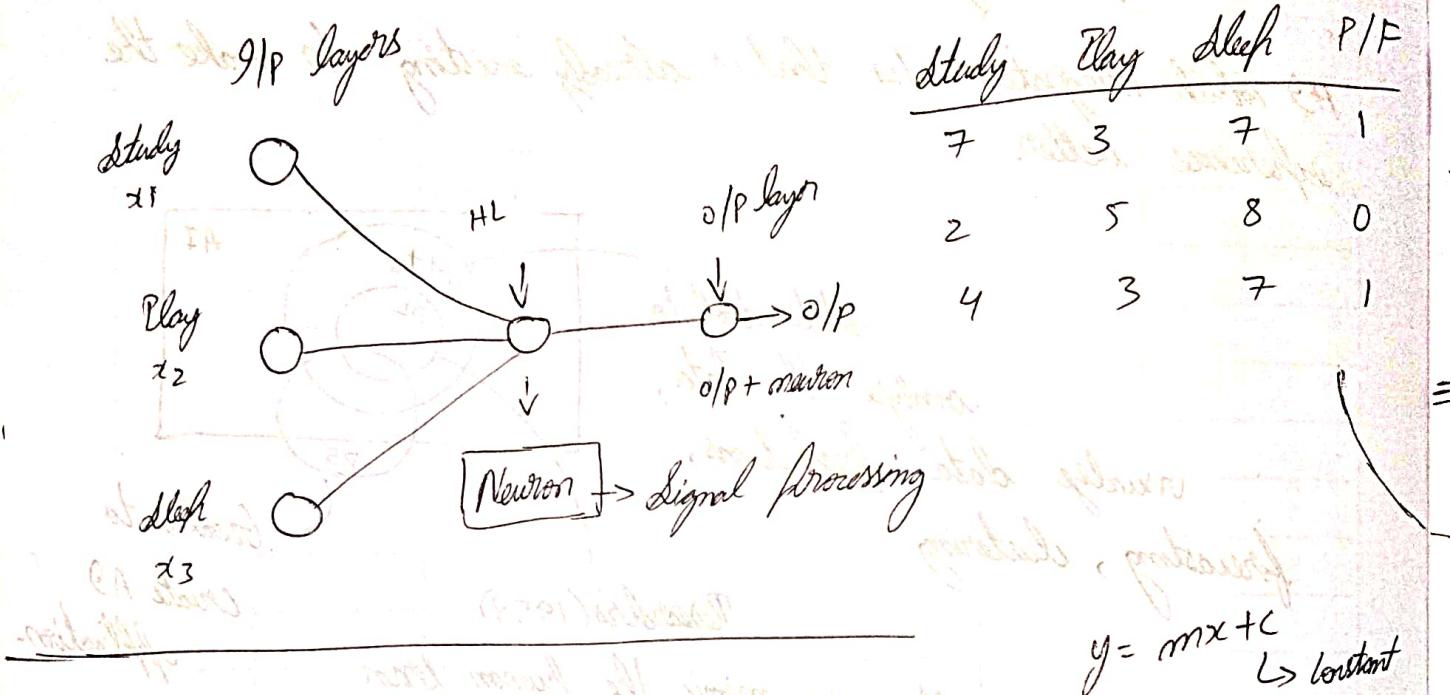
Model → Reduce the electricity Bills ↓↓

Sold in subscription basis ⇒ Company generate revenue, Better Decisions

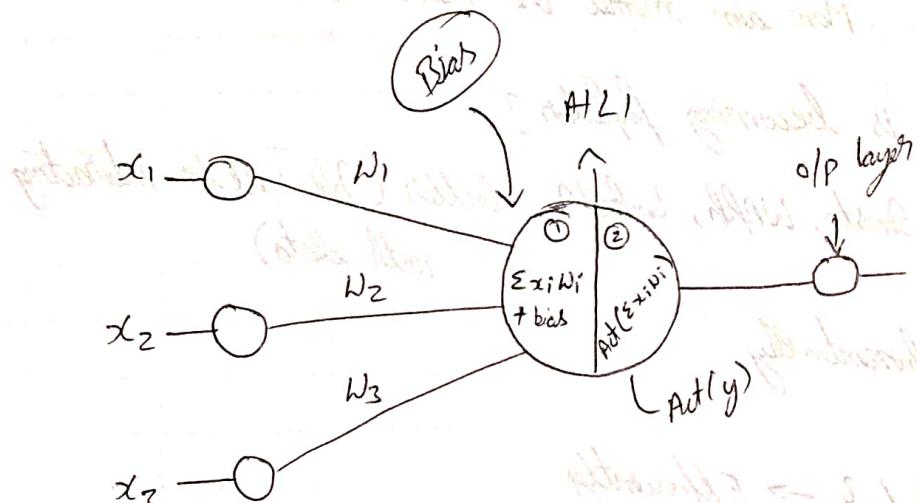
② Hardware Advancement (NVIDIA) \rightarrow GPU's \rightarrow Training the model

* Perceptron of single layered Neural Network

Binary Classification



$$y = mx + c \quad \hookrightarrow \text{constant}$$



$$\begin{aligned} \sum x_i w_i &= x_1 w_1 + x_2 w_2 + x_3 w_3 \\ &\downarrow \\ &= w^T x \end{aligned}$$

\Rightarrow Putting a hot obj in right hand, suddenly I'll move my right hand away.
Bcoz the neurons are passing from here those are getting activated.
Neurons passes the signals to the brain.

\Rightarrow During the training of neural net, we will be taking care that which weight should have value so that neuron should activate on till that level.

When calculating the o/p of node, the i/p are multiplied by weights, & a bias value is added to the result.

The Bias allows the activation function to be shifted to left or right, to better fit the data.

→ Putting hot obj of moving the bond; this basically shows that whether the neurons should activate or deactivate.

⇒ We took i/p, added a bias & activate.

→ Forward Propagation.

$$3 - \frac{7}{7} = 0 \rightarrow \hat{y}$$

$$y = 1$$

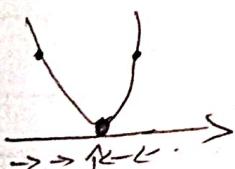
$$(y - \hat{y}) = 1$$

⇒ loss function

To Reduce the Error

Back Propagation & update weights

gradient descent \Rightarrow kind of optimizer



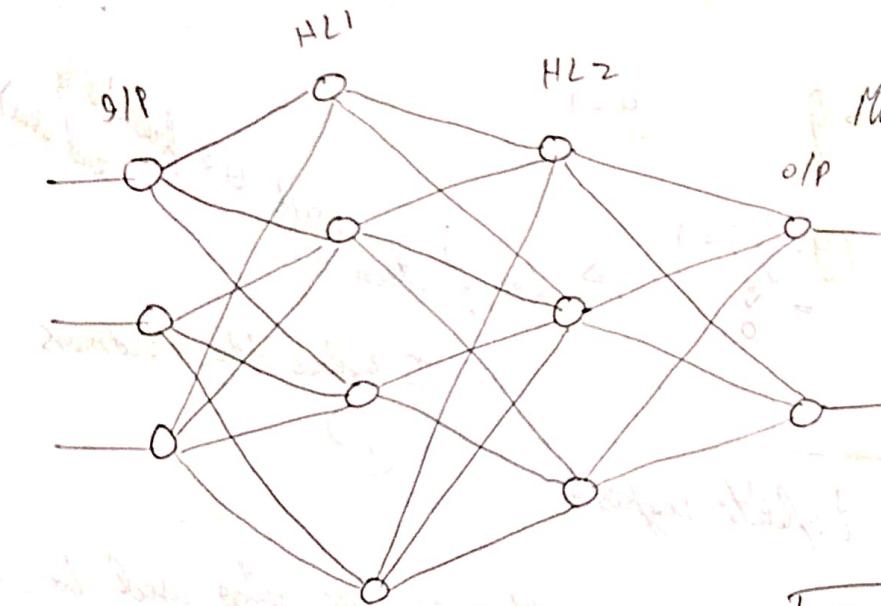
$$y - \hat{y} \Rightarrow \text{Mistake}$$

Optimizers are also used to adjust the parameters (weights & bias) of neural net to minimize the loss function.

Ex:- Gradient descent (It is changing the coefficient).

Conclusion

- ① 91P layers
 - ② Weights
 - ③ Bias
 - ④ Activation
 - ⑤ Loss function $\frac{1}{2} \sum (y - \hat{y})^2$
 - ⑥ Optimizers
 - ⑦ Update the weight
- Forward propagation Backward propagation
- Everything works on
- ① ANN
② CNN
③ RNN
④ Object Detection



Multi Layered Neural Net

$$y = w^T x + b$$

$$L \cdot (y - \hat{y})^2$$

Weight update formula

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

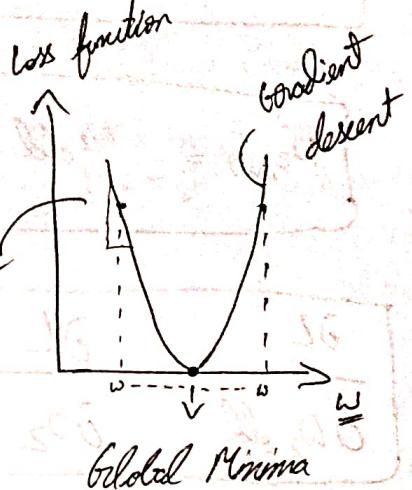
slope
-ve slope

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

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

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

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



$\eta \approx \text{small number}$

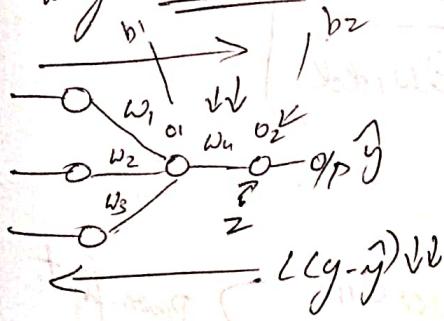
$$\eta = 0.01$$

If it may large
then the point can jump
here & there.

Chain Rule of Differentiation

To update weights during training, we need compute gradient of loss with respect to each parameter, across many layers; can be done

using Chain Rule.



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

$$\frac{\partial L}{\partial w_{\text{old}}} = \frac{\partial L}{\partial o_2} + \frac{\partial o_2}{\partial w_{\text{old}}}$$

{Chain Rule of Differentiation}

$$z = o_1 w_4 + b_4$$

→ loss is dependent on o_2 .

→ o_2 is dependent on o_1 .

→ o_1 is dependent on w_1 .

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

$$b_{2\text{new}} = b_{2\text{old}} - \frac{\partial L}{\partial b_{2\text{old}}}$$

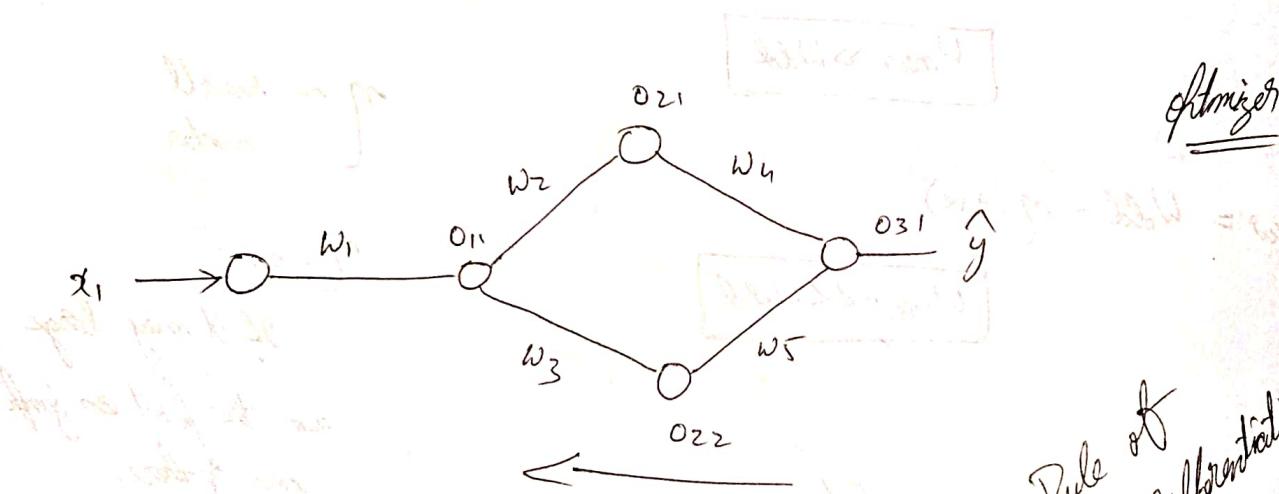
If wanted to update bias,

$$\frac{\partial L}{\partial w_{1,old}} = \frac{\partial L}{\partial o_2} + \frac{\partial o_2}{\partial o_1} + \frac{\partial o_1}{\partial w_{1,old}}$$

\Rightarrow to update w_1 , this term become chain rule of differentiation.

$$w_{1,new} = w_{1,old} - \eta \frac{\partial L}{\partial w_{1,old}} \quad \frac{\partial o_2}{\partial w_4} + \frac{\partial w_4}{\partial o_1}$$

$$\frac{\partial L}{\partial w_{2,old}} = \frac{\partial L}{\partial o_2} + \frac{\partial o_2}{\partial o_1} + \frac{\partial o_1}{\partial w_{2,old}}$$

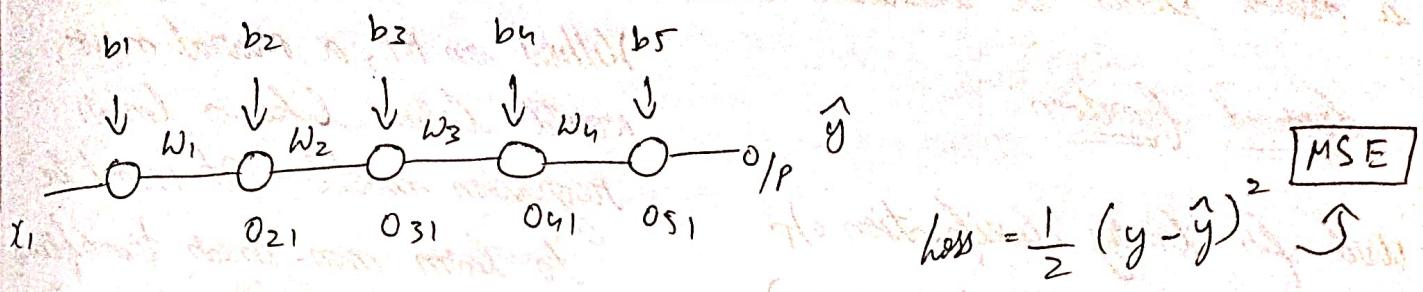


$$w_{1,new} = w_{1,old} - \eta \frac{\partial L}{\partial w_{1,old}}$$

$$\frac{\partial L}{\partial w_{1,old}} = \left[\frac{\partial L}{\partial o_{31}} + \frac{\partial o_{31}}{\partial o_{21}} + \frac{\partial o_{21}}{\partial o_{11}} + \frac{\partial o_{11}}{\partial w_{1,old}} \right] \text{ after part}$$

$$\left[\frac{\partial L}{\partial o_{31}} + \frac{\partial o_{31}}{\partial o_{22}} + \frac{\partial o_{22}}{\partial o_{11}} + \frac{\partial o_{11}}{\partial w_{1,old}} \right] \text{ down part}$$

③ Vanishing Gradient Problem



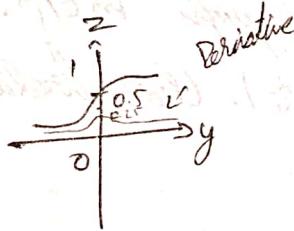
$$w_{1,\text{new}} = w_{1,\text{old}} - \eta \frac{\partial L}{\partial w_{1,\text{old}}}$$

$$\begin{aligned} \frac{\partial L}{\partial w_1} &= \frac{\partial L}{\partial o_{51}} * \frac{\partial o_{51}}{\partial o_{41}} + \frac{\partial o_{41}}{\partial o_{31}} * \frac{\partial o_{31}}{\partial o_{21}} + \frac{\partial o_{21}}{\partial w_1} \\ &= 0.25 * 0.15 + 0.10 * 0.05 + 0.02 \end{aligned}$$

\Rightarrow Assumption

Sigmoid AF

$$y = \frac{1}{1+e^{-x}} \rightarrow o_{51}$$



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

Derivative condition

$$o_{51} = \sigma \left[(o_{41} + w_4) + b \right]$$

Sigmoid AF

$$w_{\text{new}} = w_{\text{old}} - \eta (\text{small not})$$

$w_{\text{new}} \approx w_{\text{old}}$ \Rightarrow Vanishing gradient problem

sol is to use another AF.



No change in weights

Activation Function, is a mathematical function that determines whether a neuron should be "activated" or not.

1. Sigmoid function

Used for binary classification o/p layers (eg probability b/w 0 & 1).

Without an AF, a neural net would behave like a linear regression model - It won't be able to learn non-linear transformations.

→ The function o/p is not centered on 0, which will reduce the efficiency of weight.

→ Sigmoid function performs exponential operations, which is slower for computers.

Advantages

- * Smooth gradient, prevents "jumps" in o/p.
- * o/p values bind b/w 0 & 1. Clear predictions, i.e. very close to 100%.

Disadvantages

- * Prone to gradient vanishing.
- * Force o/p is not 0 centered.
- * Power operations were really time consuming.

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

Used in o/p layers (binary classification)

2. tanh function (Hyperbolic tangent function)

- + Value ranges -1 to 1.
- + The whole func is 0-centred, which is better than Sigmoid.
- + But still vanishing gradient.

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

→ For binary classification problems, the tanh func is used for hidden layers & sigmoid function for o/p layers.

3. ReLU function (Rectified Linear Unit)

$$f(x) = \max(0, x)$$

- + Positive values stay the same, -ve values become 0. Range [0, ∞]
- + But may cause dead neurons
- When i/p is +ve, there is no gradient saturation problem.
- Diseadvantages
 - + When i/p is -ve, ReLU is completely inactive.
 - + One -ve not entered ReLU will die.

Used for most hidden layers.

4. Leaky ReLU function

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

- + The i/p of -ve value, cannot 0 in this function
- + It solves dead neuron problem, by allowing small gradient for -ve i/p.

Range [-∞, ∞]

5. ELU (Exponential Linear Units) function.

* Designed to solve ReLU problems.

→ No Dead ReLU issues.

→ The mean of o/p is close to zero. 0-centered.

* Computationally expensive, bcz we are using exponential value.

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

6. Softmax function

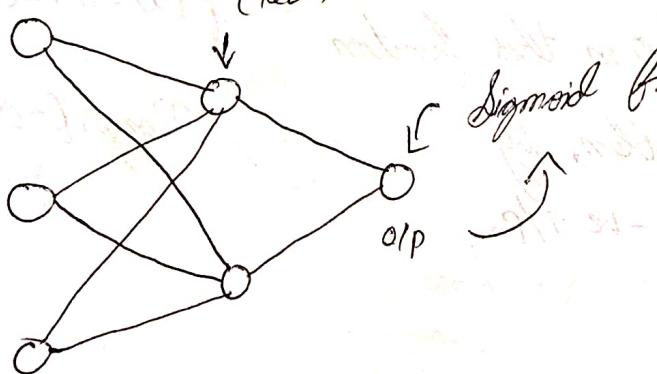
Used for o/p layers of multiclass classification

$$f(x_i) = \frac{e^{x_i}}{\sum_j e^{x_j}}$$

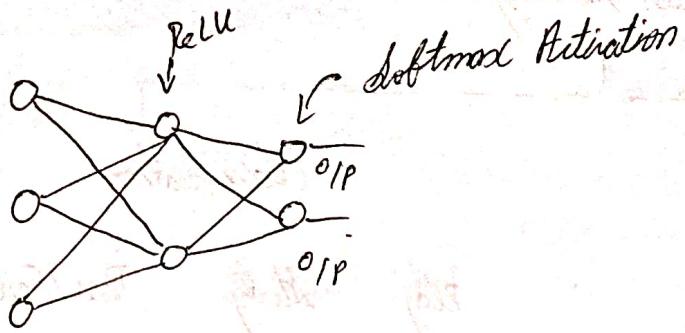
Technique which activation fn we should use;

Binary Classification

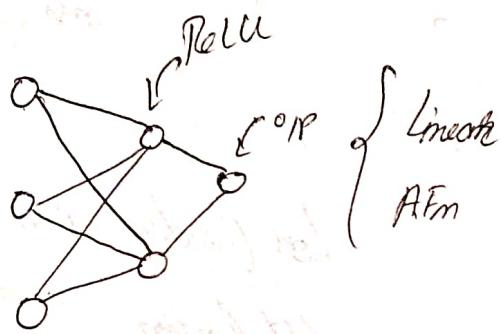
(ReLU) (PreReLU) (ELU)



Multiclass



Regression (continuous o/p)



Deep Learning (ANN)

Regression

	Exp	Degree	Salary
10	PhD	-	70k
-	-	-	-
-	-	-	-

Classification

	Play	Study	Pass/Fail
10	2	-	Fail
4	3	-	Fail
5	5	-	Maybe
2	7	-	Pass

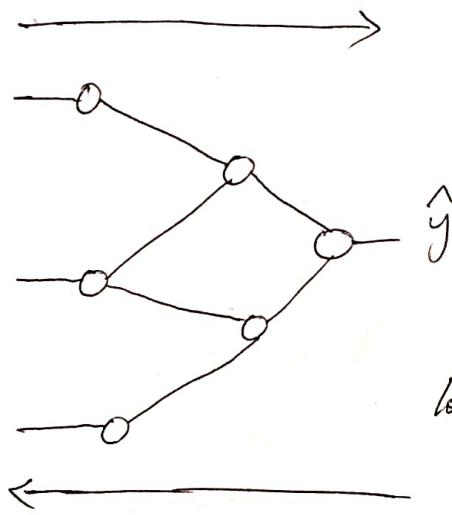
Regression

- ① MSE
- ② MAE
- ③ Huber loss

Loss function AND
Cost function

10 rows
↓

Dataset = 100 results



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

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

* In loss fm, we will be providing one data point.

* In cost fm, I specifically provide batch of data points.

① Mean Squared Error (MSE) → is a loss fn used to measure diff b/w predicted avg squared value \hat{y}_i & actual values.

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

Quadratic eqⁿ

Advantages

- + Differentiable
- + It has only 1 local or global minima.
- + It converges faster

Disadvantages

- + Not robust to outlier.

② Mean Absolute Error (MAE) → is regression loss fn that calculates the avg of absolute diff b/w predicted & actual values.

$$\text{loss fn} = \frac{1}{2^n} |y - \hat{y}|$$

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

- + Robust to outlier.

③ Huber loss → combination of MSE & MAE.

when outlier not present

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

) value set before training that controls how ML model learns.
It is not learned from data.

Classification

Gross Entropy Binary cross entropy (Binary classification)
Categorical cross entropy (Multiclass classification)

① Binary Gross entropy

$$\text{loss} = -y \cdot \log(\hat{y}) - (1-y) \cdot \log(1-\hat{y}) \Rightarrow \text{logistic Regression}$$

$$\text{loss} = \begin{cases} -\log(1-\hat{y}) & \text{if } y=0 \\ -\log(\hat{y}) & \text{if } y=1 \end{cases} \quad \left. \begin{array}{l} \text{Binary classification} \\ \downarrow \end{array} \right.$$

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

② Categorical Gross entropy

one hot encoding

	f_1	f_2	f_3	0/p	Good	Bad	Neutral
1	2	3	4	Good	[1 0 0]		
2	5	6	7	Bad	0 1 0		
3	8	9	10	Neutral	0 0 1		

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

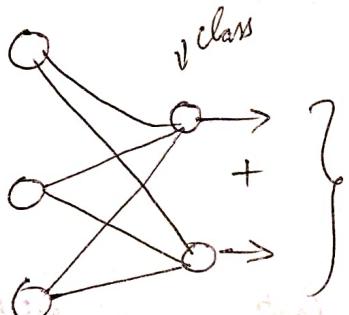
$i = \text{row}$

$j = \text{col}$

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

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

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



for multiclass classification

$$f(x_i) = \frac{e^{x_i}}{\sum_j e^{x_j}}$$

→ categorical cross entropy (loss)

Conclusion :-

ReLU, Softmax \Rightarrow Multiclass classification

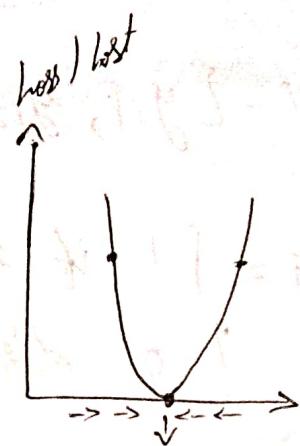
ReLU, Sigmoid \Rightarrow Binary \Rightarrow Binary cross entropy (loss)

ReLU, Linear AF \rightarrow $\frac{\text{MAE}}{\text{MSE}}$ \Rightarrow Linear Regression
Huber loss

Optimizers

- Gradient Descent.
- SGD (Stochastic GD)
- Mini Batch SGD
- SGD with Momentum
- Adagrad
- RMS Prop
- Adam Optimizers

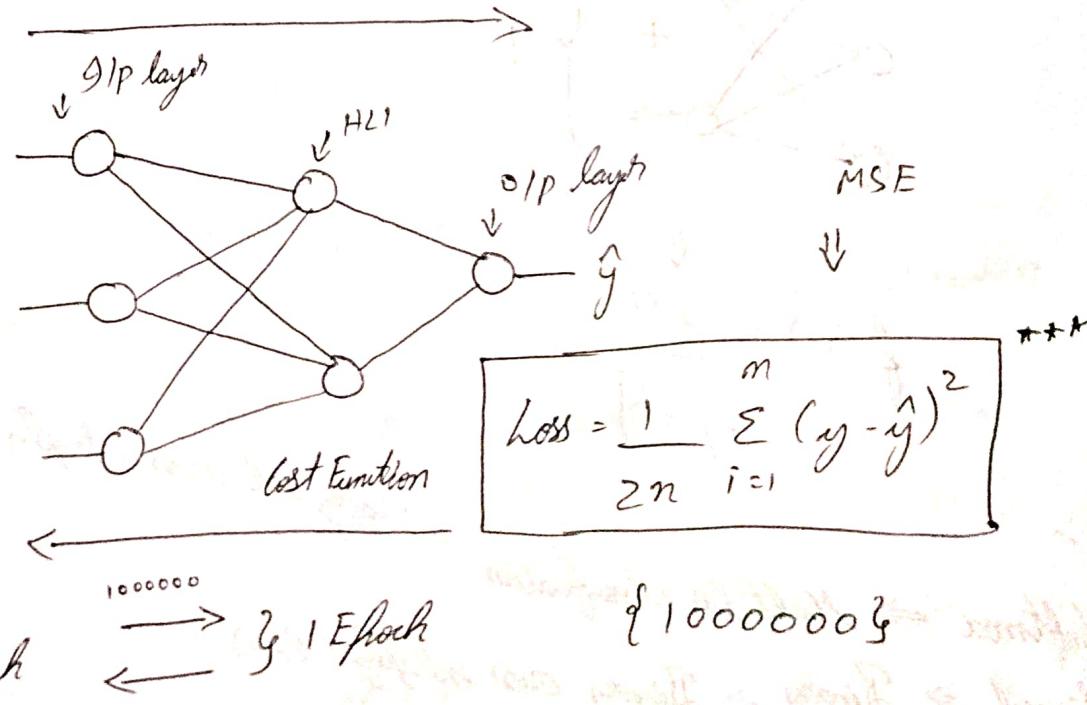
Batch, Epochs, Iterations



① Gradient Descent

Weight Updation formula

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



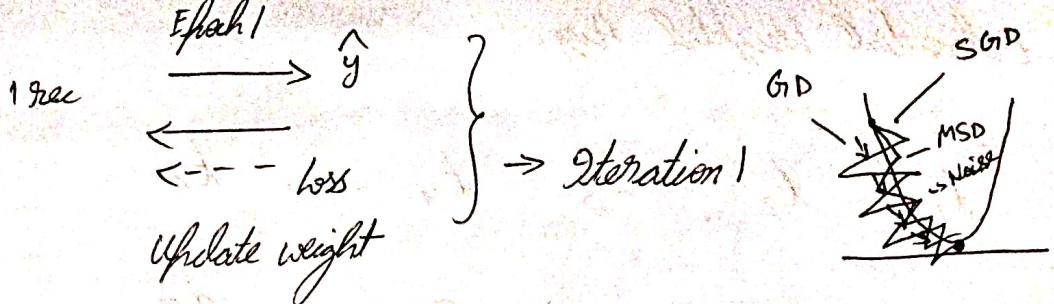
↓
→ Resource extensive & Huge RAM

② Stochastic GD

\Rightarrow RAM $\downarrow\downarrow$

Discadv

- + convergence will be very slow
- + TC will be high



③ Mini batch SGD

\Rightarrow Resource Intensive.

\Rightarrow convergence will be better.

\Rightarrow TC will be improved.

batch size = 1000

$$E_1 \xrightarrow{1000} g_1$$

$$E_2 \xrightarrow{1000} g_2$$

$$\dots \xrightarrow{1000} g_{1000}$$

④ SGD with Momentum

{ Exponential Weighted Avg's

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

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

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

Time Series

ARIMA, ARMA

Exponential Weighted Avg.

$t_1, t_2, t_3 \leftarrow t_4, \dots, t_n$

$a_1, a_2, a_3, a_4, \dots, a_m$

$\beta \Rightarrow \text{Hyperparameter}$

$$V_{t_1} = a_1$$

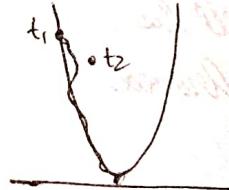
$$\beta = 0 \text{ to } 1$$

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

$$0.95$$

$$= (0.95) * V_{t_1} + (0.05) * a_2$$

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



Exponential Weighted Average

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

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

→ Noise will reduced.

→ Quick conversion.

⑤ Adagrad (Adaptive Gradient Descent)

$\eta = \text{fixed} \Rightarrow \text{adaptive} \Rightarrow \text{Learning rate} \Rightarrow \text{decreasing} \Rightarrow \text{Global minima}$

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

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

↓

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

Huge note

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

$$t=1 \quad t=2 \quad t=3$$

$$\eta = 0.01 \quad \eta = 0.05 \quad \eta = 0.002$$

⑥ Adadelta And RMSProp

$$\eta' = \frac{\eta}{\sqrt{Sd\omega + \epsilon}}$$

$$Sd\omega = 0$$

$$Sd\omega_t = \beta Sd\omega_{t-1} + (1-\beta) \left(\frac{\partial L}{\partial w_t} \right)^2$$

$$\beta = 0.95$$

$$Sd\omega_t = (0.95) Sd\omega_{t-1} + (1-\beta) \left(\frac{\partial L}{\partial w_t} \right)^2$$

⑦ Adam Optimizer (Best optimizer)

Momentum + RMS Prop (Adaptive learning rate)

$$V_{\partial w} = 0 \quad V_{\partial b} = 0 \quad S_{\partial w} = 0$$

$$S_{\partial b} = 0$$

$$w_t = w_{t-1} - \eta' V_{\partial w}$$

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

$$b_t = b_{t-1} - \eta' V_{\partial b}$$

$$V_{\partial w} = \beta + V_{\partial w_{t-1}} + (1-\beta) \left(\frac{\partial L}{\partial w_{t-1}} \right)$$

$$V_{\partial b} = \beta + V_{\partial b_{t-1}} + (1-\beta) \frac{\partial L}{\partial b_{t-1}}$$

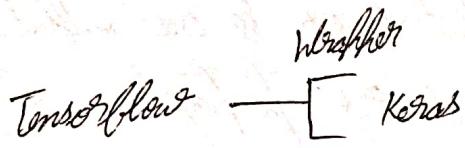
→ Smoothing

→ Learning rate becomes adaptive

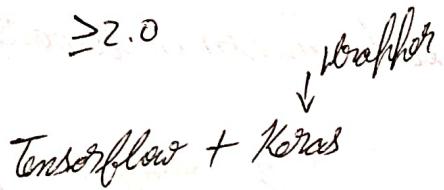
open-sourced by

Tensorflow \Rightarrow Google, Deepmind

< 2.0



≥ 2.0



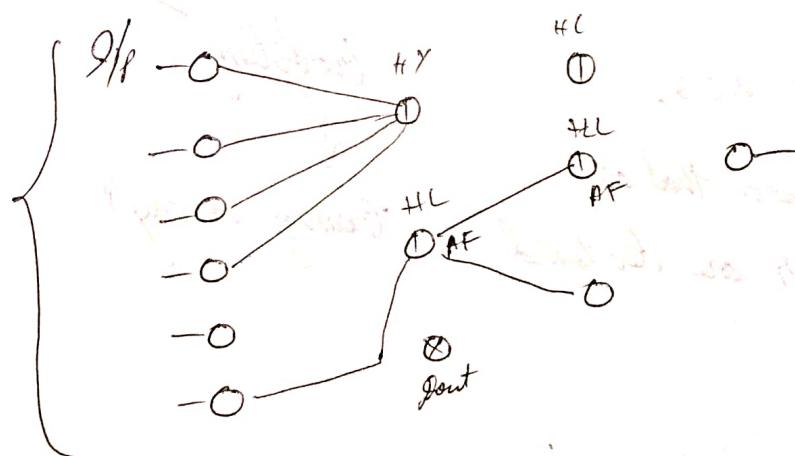
Tensorflow is used for

developing complex neural nets.

Keras is a high-level API built on top of Tensorflow that make it easier & faster to build deep learning models with simple code.

```
from tensorflow.keras.models import Sequential  
layers import Dense  
layers import ReLU, PReLU, LeakyReLU, ELU  
layers import Dropout
```

of Sequential...



Sequential is a class from Keras (part of Tensor Flow) used to build feed forward neural nets where layers are added one after another, in sequence.

→ Where each layer has exactly one input tensor & one output tensor.
We can also perform forward & backward propagation.

Dense is used to create layers where each neuron is connected to every neuron in previous layer.

→ Using Dense we will be used to create i/p layer, hidden layer, o/p.

AF, introduces non-linearity to the model, allowing it to learn complex patterns.

Droptout is a regularization technique used during training the model to prevent overfitting in NN.

If I specify dropout is 0.3,

so 1. of entire neurons that are

present in these layers are deactivated

while training.

Overfitting
↓
Test Accuracy ↑

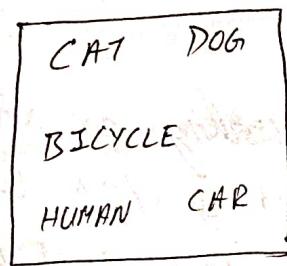
Training Accuracy ↑

Early Stopping is a regularization technique used during training to stop the model early when the performance on the validation set starts to degrade, to prevent overfitting.

Black box model is a ML model whose internal workings are not easily interpretable by humans. We can see the i/p's & o/p's but not how model make its decision. Ex:- Neural networks, SVM etc

White box model is fully transparent & interpretable, where we can understand how the model make decisions. Ex:- Linear Regression.

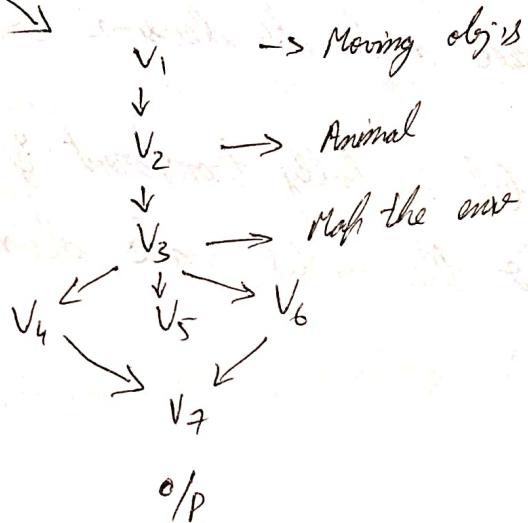
CNN vs Human Brain \rightarrow Always understand how does it work



Cerebral cortex



Visual cortex \rightarrow responsible for identifying the items.



CNN

Convolution

Images =

0-255
5x5

10	0	255		
19	17	240	241	
.	.	.	.	

1 channel \leftarrow B & W
&
RGB

0 - Black

255 - White color

6x6

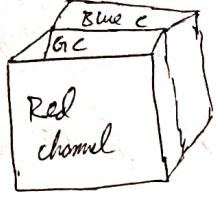
3x3

4x4

, what happening here is image size is decreasing.
there losing some kind of info. So to prevent
that will can apply padding to it.

Padding is about protecting the image by adding another layer on top of it.
we can fill 0's or apply nearest value.

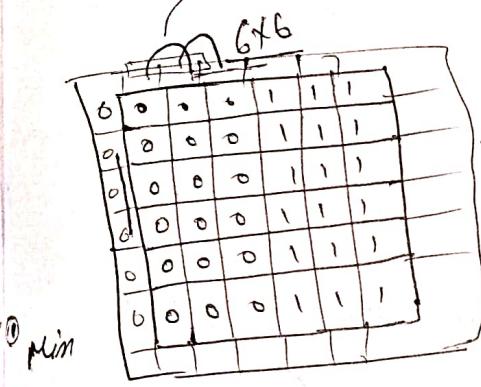
RGB \rightarrow 3 channel



0×255
each channel, when combine this we
will getting diff color.

$5 \times 5 \times 3$

① Convolution operation



min

Max

scaling ***

$$\begin{aligned} 1. & 0 + 0 + 0 + 0 + 0 = 0 \\ 2. & 0 + 0 + 1 + 0 + 0 - 1 = 0 \\ 3. & 0 + 0 + 2 + 1 - 2 - 1 = 0 \end{aligned}$$

stride = 1, to move filter 1 step to right

$$\text{After padding} = m + 2p - f + 1$$

$$\text{By padding} = m - f + 1$$

3×3
Horizontal edge filter

filter / kernel

$$\begin{bmatrix} 1 & 2 & 1 \\ 0 & 0 & 0 \\ -1 & -2 & -1 \end{bmatrix}$$

$$\Rightarrow \begin{bmatrix} 0 & 0 & 0 & 0 \\ -4 & -4 & -4 & -4 \\ -4 & -4 & -4 & -4 \\ 0 & 0 & 0 & 0 \end{bmatrix}$$

$0/p$

$$m = 6, f = 3$$

$$6 \times 6 \quad 3 \times 3$$

4×4

By folding

$$m - f + 1 = 6 + 1 - 3 = 4$$

$$\begin{bmatrix} 1 & 0 & -1 \\ 2 & 0 & -2 \\ 1 & 0 & -1 \end{bmatrix}$$

$$\begin{array}{l} m=6 \\ p=1 \\ f=3 \end{array}$$

$$\text{After folding } 8 - 3 + 1 = 6$$

$$m + 2p - f + 1 = 6 + 2(1) - 3 + 1 = 6$$

6×6

$0/p$

Vertical edge filter

$$\Rightarrow \begin{bmatrix} 0 & -4 & -4 & 0 \\ 0 & -4 & -4 & 0 \\ 0 & -4 & -4 & 0 \\ 0 & -4 & -4 & 0 \end{bmatrix}$$

filters can take out the information from the specific image.

Importance of padding is that to prevent the loss of image we make use of info padding.

ANN is the basic form of NN, inspired by how human brain works. It is made of input, hidden & output layers & each neuron is connected with a weight. Ex:- hand force prediction.

CNN (Convolutional NN) is deep learning algo that can automatically learn spatial data of image. It uses convolutional layers that can detect features like edges, shapes, textures.

CNN Architecture

1. Convolutional layer :- extracts features from image & applies filter (kernel) that slide over image & detect patterns.
2. Apply AF, to learn complex patterns. converts -ve value to 0. ReLU.

3. Pooling layer is used to downsamples the features with to reduce spatial size & computation. Keeps most imp feature.

Reduces

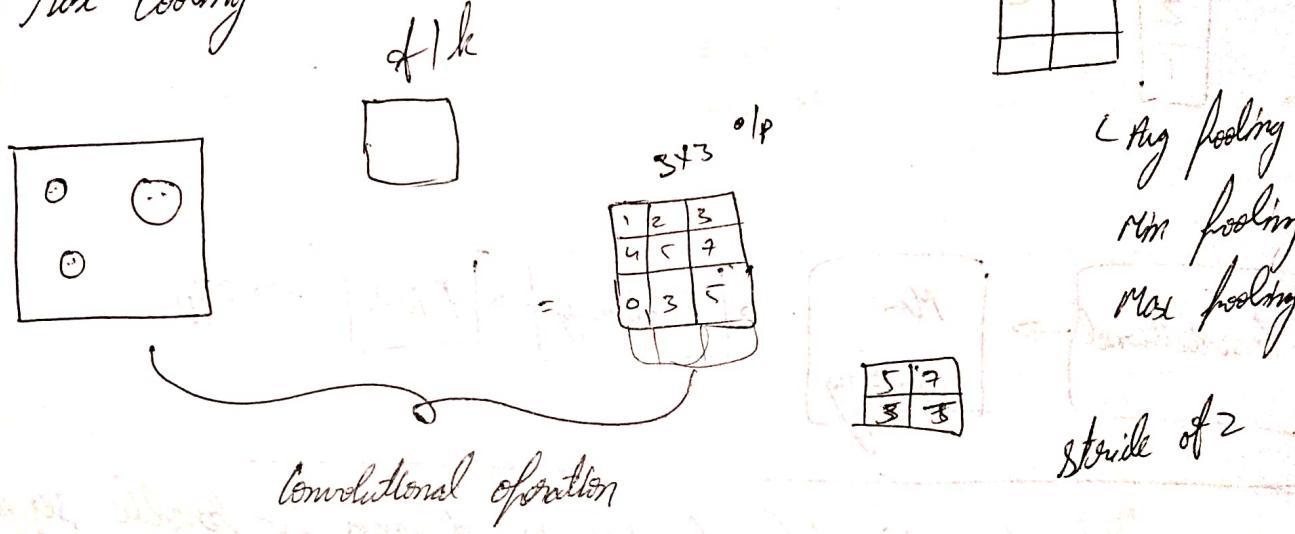
Dimensions

Computational cost
overfitting

4. Flatten layer, converts the multi-dimensional tensor (2D or 3D) into 1D vector. After flattening the vector is then passed to fully connected layer.

5. Fully connected (Dense) layer :- performs classification based on extracted features.

② Max Pooling

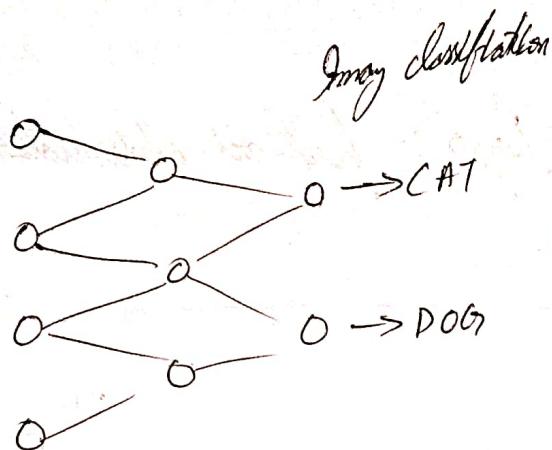


localization invariance refers to the ability to recognize features like edges, shapes, objects regardless of their position of image.
Essentially means if I have multiple obj also if 1 of the filter can determine that obj.

Flattening Layer - similar to ANN Dense layers

Max Pool

5
7
3
5
7
9
2
1

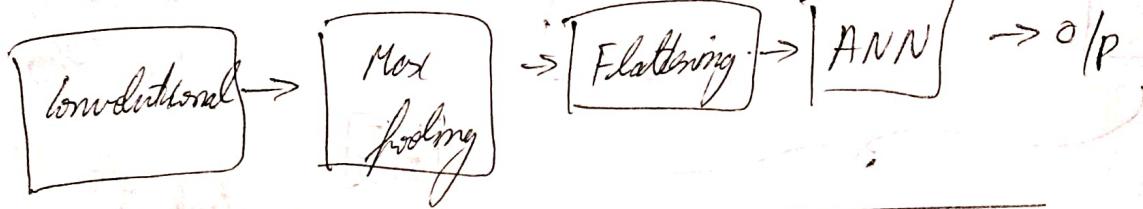


Max

5	7
3	5

7	9
2	1

6	5
3	1



RNN (Recurrent NN) type of NN is designed to handle sequential data, such as:-

- data, such as:-
- > Time series data (eg stock prices)
- > Natural language (eg text or speech)
- > Video frames or sensor data over time.

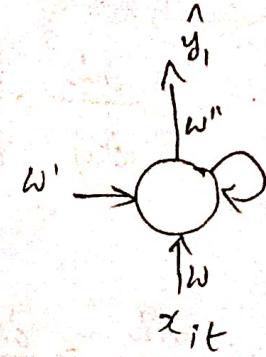
RNN

① NLP \rightarrow 91% data

\hookrightarrow BOW, TF-IDF, WORD2VEC

Vectors

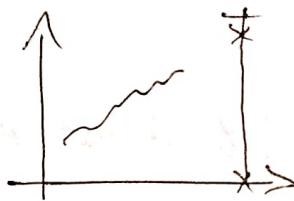
	w_1	w_2	w_3	\dots
s_1	1	0	1	0
	0.3	0.75	-	-



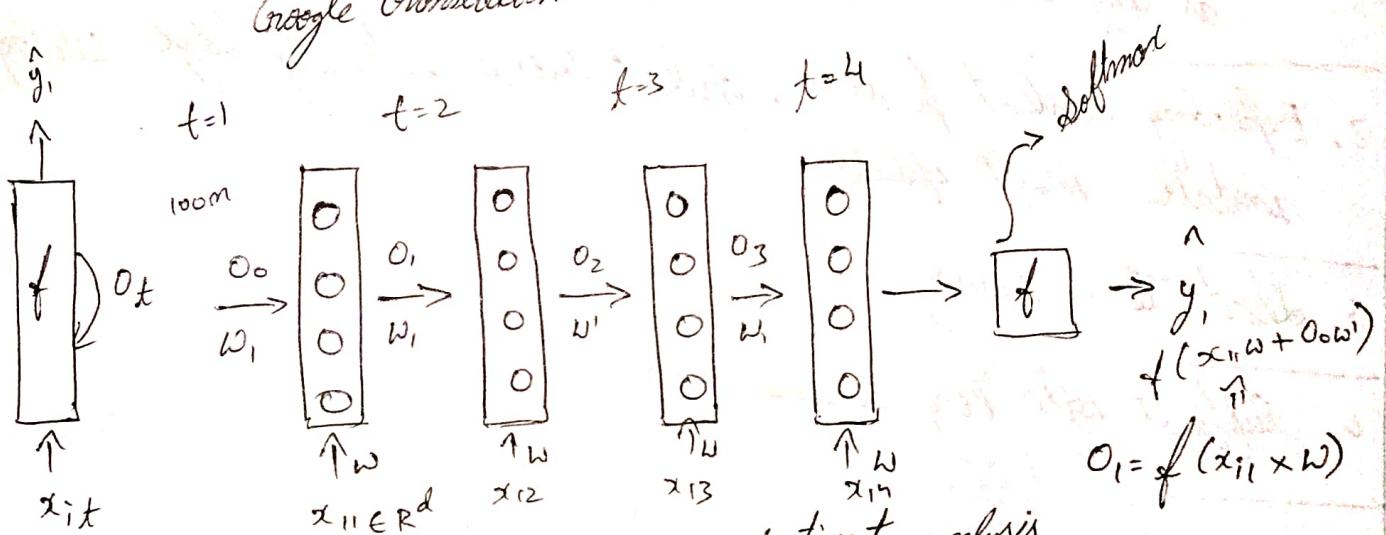
Time series

ARIMA

SARIMAX



Google Translation



NLP Usecase sentiment analysis

$$x_1 = \langle x_{11}, x_{12}, x_{13}, x_{14} \rangle$$

$$o_2 = f(x_{12}w + o_1w_1)$$

$$o_3 = f(x_{13}w + o_2w_1)$$

$$o_4 = f(x_{14}w + o_3w_1)$$

\longrightarrow

Forward propagation

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

$$\frac{\partial L}{\partial \hat{y}}, \frac{\partial L}{\partial w''} = \frac{\partial L}{\partial \hat{y}_i} \cdot \frac{\partial \hat{y}_i}{\partial w''}$$

$$w'' = w'' - \frac{\partial L}{\partial w''}$$

$$\frac{\partial L}{\partial \hat{y}}, \frac{\partial L}{\partial w} = \frac{\partial L}{\partial \hat{y}_i} \cdot \frac{\partial \hat{y}_i}{\partial o_i} \cdot \frac{\partial o_i}{\partial w}$$

The RNN has many problems like

$$w = w - \frac{\partial L}{\partial w}$$

1. Vanishing gradient problem

During backward propagation, gradients becomes very small (almost 0) as they are multiplied each step.

2. Exploding gradient problem, gradient become extremely large, causing unstable weight updates.

3. Short-term memory.

4. Difficult with long sequences

LSTM - RNN (It decides what to remember, what to forget & what to update @ each step.)

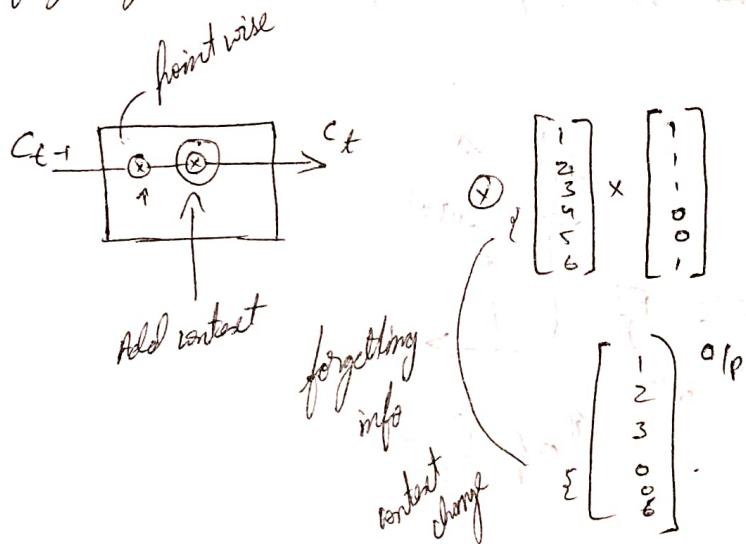
long-short term memory, it's a special kind of RNN having the capability to learn long-term dependencies.

Introduced by Hochreiter & Schmidhuber (1997).

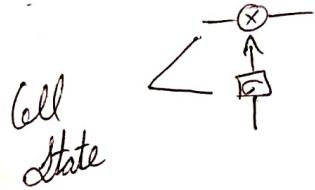
It resolves the vanishing gradient problem which occurs in RNN.

① Memory Cell, basically used for remembering & forgetting based on context of i/p.

1. Memory Cell
2. Forget Gate
3. GIP Gate
4. O/P Gate



② Forget Gate, decides what to forget from old memory.



$$f_t = \sigma(W_f \cdot [h_{t-1}, x_t] + b_f)$$

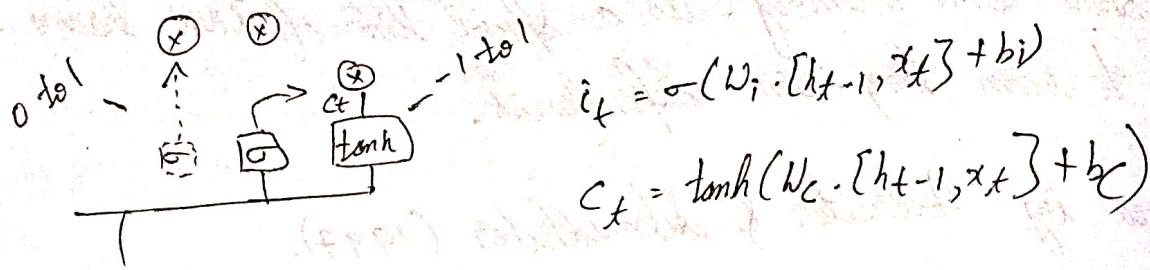
$$y = m x + c$$

$$y = w^T x + b$$

King Queen
[] []

King Queen
[] []

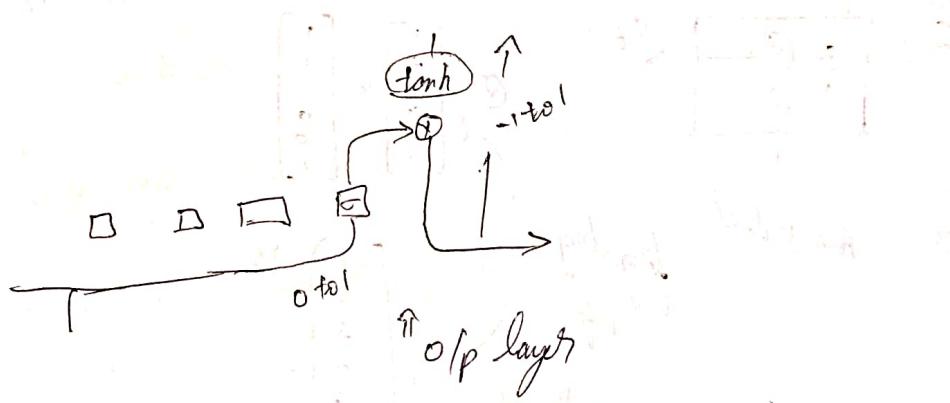
③ Input Gate, decides what new info to add.



Adding information

$$c_t = f_t + c_{t-1} \cdot i_t + c_t$$

④ Output Gate, decides what info to pass to the next step.



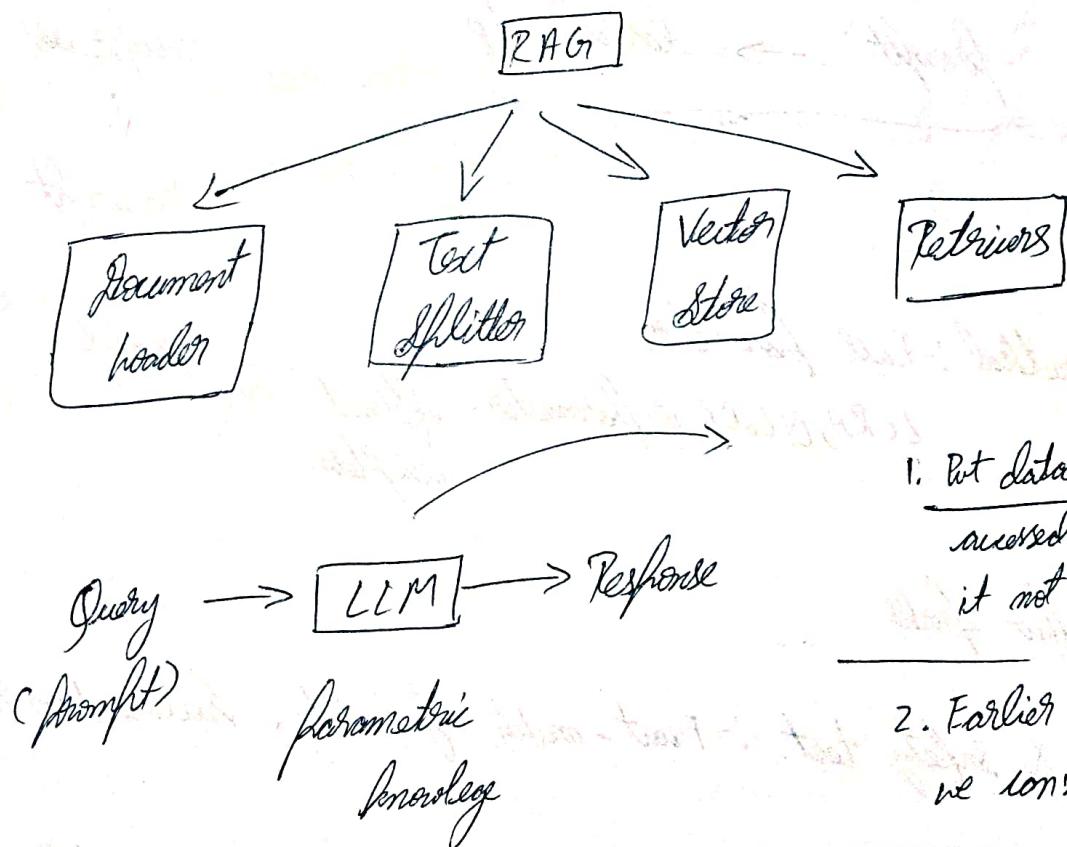
cell state is a long-term memory of LSTM. It keeps track of all 3 gates.

RAG (Retrieval-Augmented Generation) a modern technique of NLP that combines retrieval based models with generative models to provide more accurate & informative responses.

→ It retrieve relevant documents from knowledge source to generate answers based on both, retrieved content & LLM query.

Tools commonly used in RAG systems:-

- + Vector DB's :- FAISS, Pinecone, Weaviate
- + LLM's :- OpenAI's GPT, Hugging face Transformer.
- + Framework :- langchain, LlamaIndex.



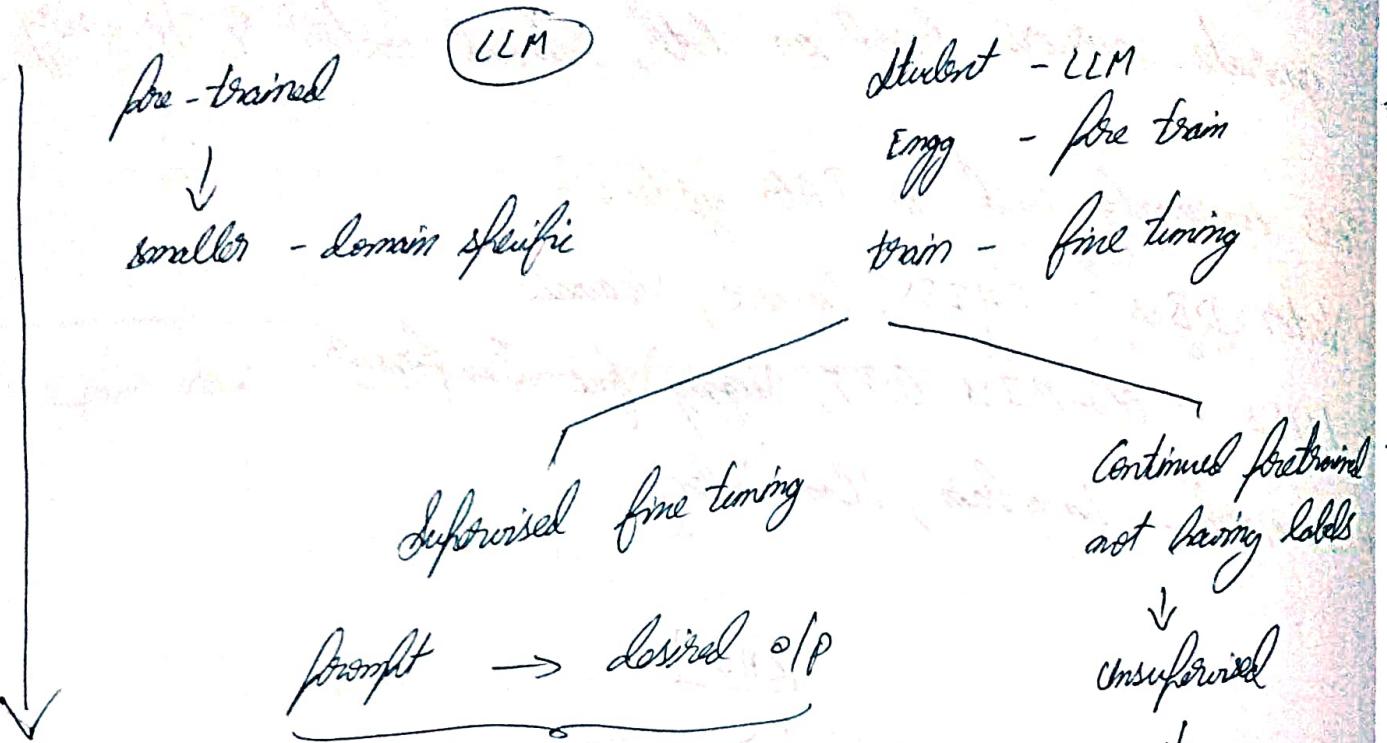
1. Put data one can't be accessed. Bcz pre-trained it not seen that data.

2. Earlier or recent data we won't get response.

3. Hallucination, it factually gives incorrect answer using parametric by confident.

We can solve that 3 LLM problem atleast using fine-tuning.

Fine-tuning is the process in ML, where a pre-trained model is further trained on specific dataset to adapt new domain.



1. Collect data

2. Choose a method :- Full parameter,

LORA, QLORA or parameter-efficient adapters.

3. Train for few epochs

4. Evaluate & safety test. :- Exact-match, factuality, hallucination rate.

Hallucination solved giving examples.

↳ I don't know

Problems in FT

1. LLM \rightarrow train (computationally expensive).
2. Strong technical expertise.
3. High data rate; updating data again & again.

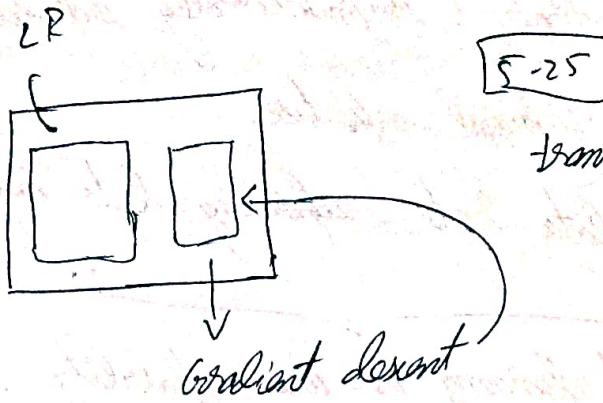
In context learning is one capability of LLM like GPT 3/4, where models learn to solve a task purely by seeing examples in the prompt - without updating its weights.

→ few short prompting.

In context-learning is a LLM's emergent property - is ability that suddenly appears in a system when it reaches a certain scale. Even though it was not explicitly programmed or expected from individual components.

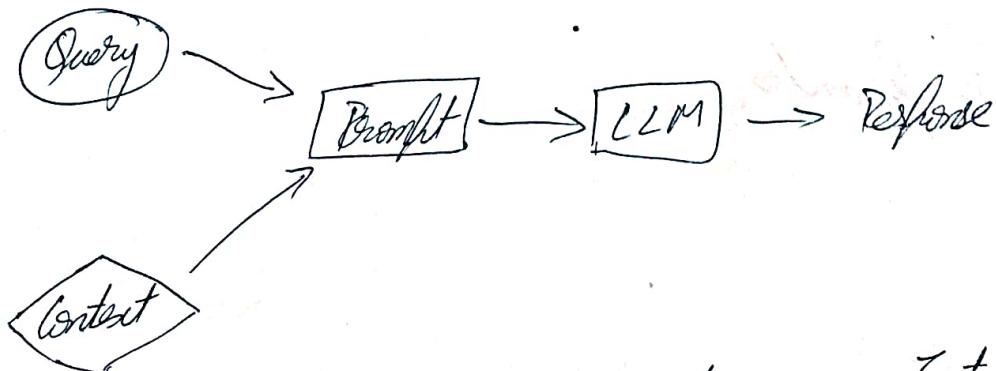
GPT-3 - 175B Parameters.

context



- Instead of just ex tasks, retrieve background info, facts, documents, product manuals etc.
- Inject that into prompt to augmented the model's knowledge.

RAG is a way to make language model smarter by giving it extra information at the time we ask the question.



- * Indexing
- * Retrieval
- * Augmentation
- * Generation

RAG → Information +
Retrieval Text
generation

Indexing is a preprocessing step, large datasets are broken down into smaller chunks, Each chunk is converted into a vector embedding using a model (like BERT). These embeddings are stored in vector database.

Retrieval at query time, the user can also embedded into a vector. The system retrieves the most similar chunks from indexed database using similarity search (like cosine similarity).

Augmentation, the retrieved documents are combined with user query. This augmented i/p is then sent to language model. The model now has contextual information from the knowledge base to answer more accurately.

Generation, this is where actual o/p is created. Generates informed responses.

* Indexing takes 4 steps:-

- Document Ingestion :- load our knowledge into memory.
- Text Chunking :- Break large doc into small, semantically meaningful chunks.
- Embedding generation :- converts each chunk into dense vector that capture its meaning.
- Storage in Vector store :- store the vectors along with original chunk text + metadata in a vector database.