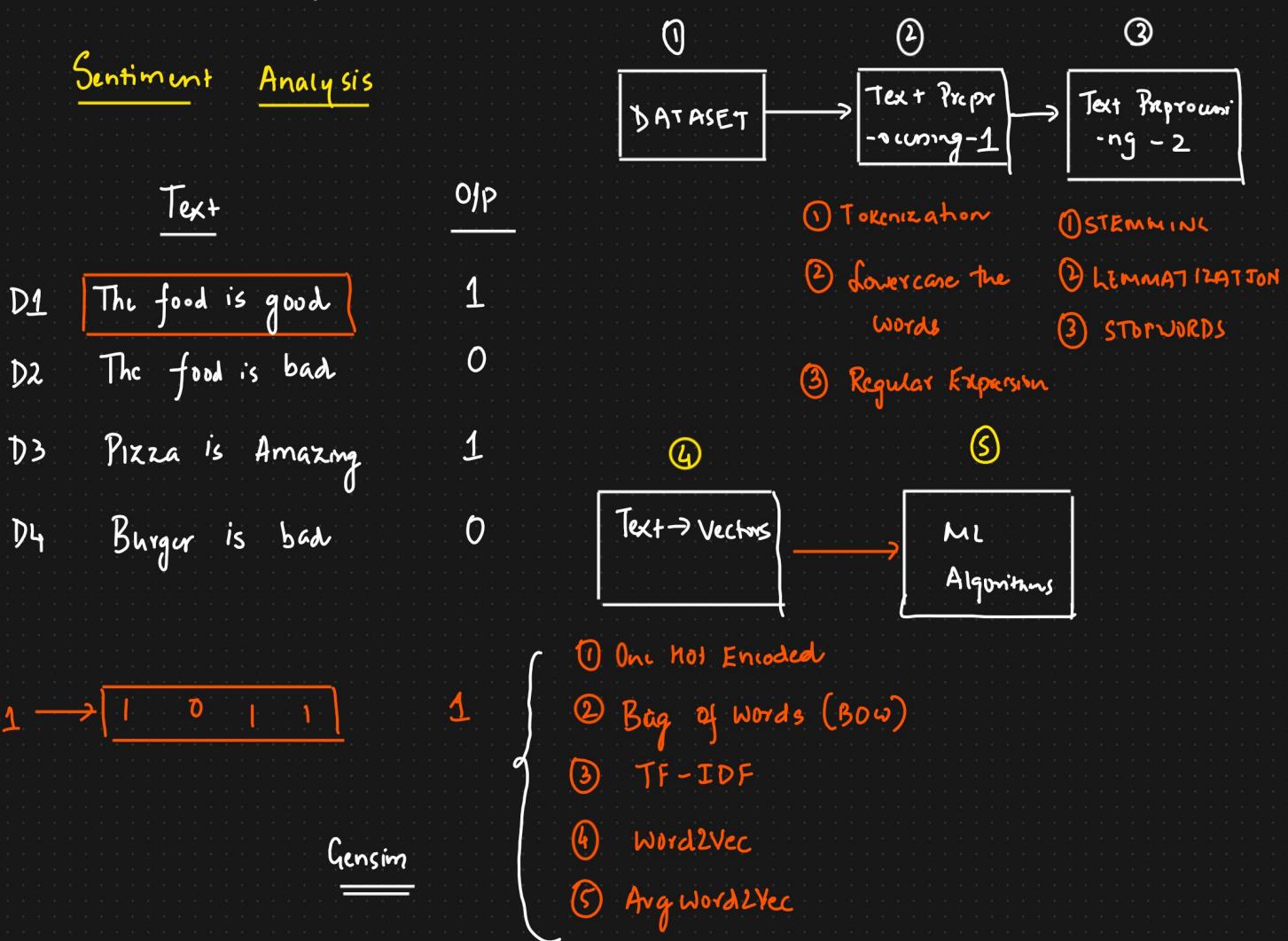


Text Preprocessing → What we have learnt?



① One Hot Encoding

Vocabulary size: 7

Vocabulary {unique words}

	<u>Text</u>	<u>O/P</u>	The food is good bad Pizza Amazing
D1	The food is good	1	1 0 0 0 0 0 0
D2	The food is bad	0	0 1 0 0 0 0 0
D3	Pizza is Amazing	1	0 0 1 0 0 0 0
Test	Burger is bad		
D1	[1 0 0 0 0 0 0],		D2 [1 0 0 0 0 0 0], D3 [0 0 0 0 1 0],
	[0 1 0 0 0 0 0],		[0 0 1 0 0 0 0], [0 0 1 0 0 0 0],
4x7	[0 0 1 0 0 0 0].	4x7.	3x7 [0 0 0 0 0 1]

$$\begin{bmatrix} 0 & 0 & 0 & 1 & 0 & 0 \end{bmatrix} \quad \begin{bmatrix} 0 & 0 & 0 & 0 & 1 & 0 \end{bmatrix}$$

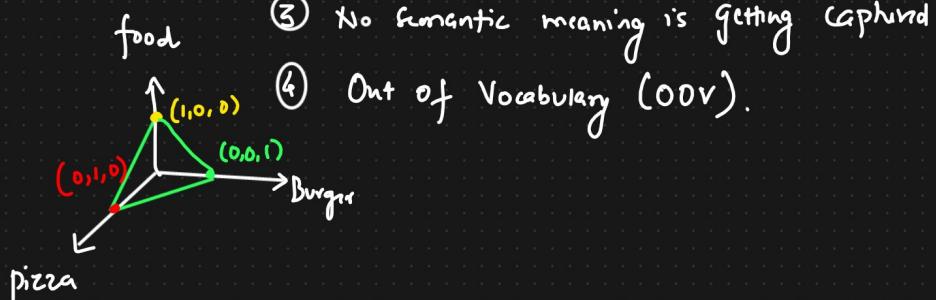
$\left\{ \text{50K vocabulary size} \right\}$ —
Disadvantages

Advantages

- ① Easy to implement with python

[sklearn OneHotEncoder, pd.get_dummies()]

food	pizza	burger
1	0	0
0	1	0
0	0	1



- ① Sparse matrix → Overfitting

- ② ML Algorithm → Fixed Size IIP

- ③ No semantic meaning is getting captured

- ④ Out of Vocabulary (OOV).

② Bag of Words

Dataset

Text	O/P		
He is a good boy	1	Counts all the words case	S1 → good boy
She is a good girl	1	⇒	S2 → good girl good
Boy and girl are good	1	Stopwords	S3 → Boy girl good [School] [Test]

Vocabulary	frequency		[good boy girl]	O/P
good	3	↓	S1 [1 1 0]	1
boy	2	↓	S2 [1 0 1]	1
girl	2	↓	S3 [1 1 1]	1

Binary Bow and BOW

{ 1 and 0 }

{ Count will get updated
based on frequency }

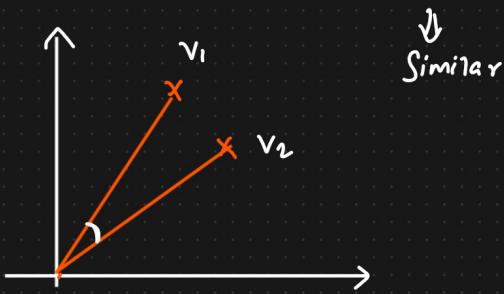
Advantages

- ① Simple and Intuitive
- ② Fixed Sized I/p \rightarrow ML Algorithms

Disadvantages

- ① Sparse matrix or array \rightarrow Overfitting
- ② Ordering of the word is getting changed
- ③ Out of Vocabulary (OOV).
- ④ Semantic meaning is still not captured.

{ The food is good $\rightarrow [1 \ 1 \ 1 \ 0 \ 1] \rightarrow v_1$
 The food is not good $\rightarrow [1 \ 1 \ 1 \ 1 \ 1] \rightarrow v_2$



N-grams Eg: bigrams, trigrams

S1 \rightarrow The food is good
 S2 \rightarrow The food is not good

Bigram

food not good

1 0 1

1 1 1

$[\text{food} \ \text{not} \ \text{good} \ \text{food good} \ \text{food not} \ \text{not good}]$

S1	1	0	1	1	0	0	}
S2	1	1	1	0	1	1	

Sklearn \rightarrow n-grams = (1,1) \rightarrow unigrams

= (1,2) \rightarrow unigram, bigram

= (1,3) \rightarrow unigram, bigram,
trigram

= (2,3) \rightarrow Bigram, trigram.

④ TF-IDF [Term Frequency - Inverse Document Frequency]

S₁ → good boy

$$\text{Term Freq(TF)} = \frac{\text{No. of rep. of words in sentence}}{\text{No. of words in sentence}}$$

S₂ → good girl

S₃ → boy girl good

$$\text{IDF} = \log_e \left(\frac{\text{No. of sentences}}{\text{No. of sentences containing the word}} \right)$$

Term Frequency * IDF

	S ₁	S ₂	S ₃	Words	IDF
good	1/2	1/2	1/3	good	$\log_e(3/3) = 0$
boy	1/2	0	1/3	boy	$\log_e(3/2)$
girl	0	1/2	1/3	girl	$\log_e(3/2)$

Final TF-IDF

	good	boy	girl	<u>Op</u>	<u>Bow</u>
Sent 1	0	$\frac{1}{2} \times \log_e(3/2)$	0	1	1 0
Sent 2	0	0	$\frac{1}{2} \log_e(3/2)$	1	0 1
Sent 3	0	$\frac{1}{3} \log_e(3/2)$	$\frac{1}{3} \log_e(3/2)$	1	1 1

Advantages

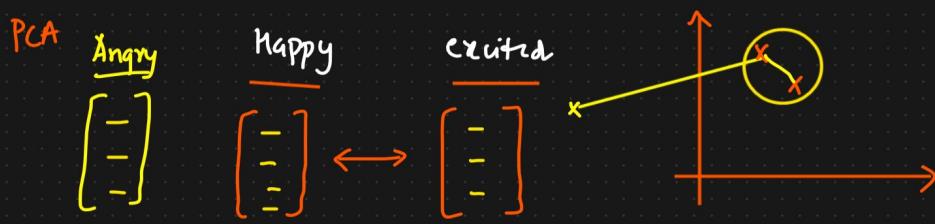
- ① Intuitive
- ② Fixed Size → Vocab size
- ③ Word Importance is getting captured

Disadvantages

- ① Sparsity still exists
- ② OOV

Word Embeddings [Wikipedia]

In natural language processing (NLP), word embedding is a term used for the representation of words for text analysis, typically in the form of a real-valued vector that encodes the meaning of the word such that the words that are closer in the vector space are expected to be similar in meaning.



Words → vectors

Sentence → vectors

Word Embeddings

[ANN]

Count or frequency

- {
- ① One
- ② Bow
- ③ TF-IDF

Deep Learning Trained
Model

Word2Vec

Cbow Skipgram

[Continuous BAG OF WORDS]

Word2Vec → Feature Representation

Google

Word2vec is a technique for natural language processing published in 2013. The word2vec algorithm uses a neural network model to learn word associations from a large corpus of text. Once trained, such a model can detect synonymous words or suggest additional words for a partial sentence. As the name implies, word2vec represents each distinct word with a particular list of numbers called a vector.

Vocabulary → Unique Words → Corpus.

	Boy	Girl	KING	QUEEN	Apple	Mango
--	-----	------	------	-------	-------	-------

Gender	-1	1	-0.92	0.93	0.01	0.05
Feature Representation	0.01	0.02	0.95	0.96	-0.02	0.02
Royal	0.03	0.02	0.75	0.68	0.95	0.96
Age	-	-	-	-	0.91	0.92
Food	-	-	-	-	-	-

300 dimension !

$$\begin{bmatrix} - \\ - \\ - \\ - \\ - \end{bmatrix}$$

$$[\text{KING} - \text{BOY} + \text{QUEEN} = \text{GIRL}]$$

$$\text{KING} [0.95, 0.96]$$

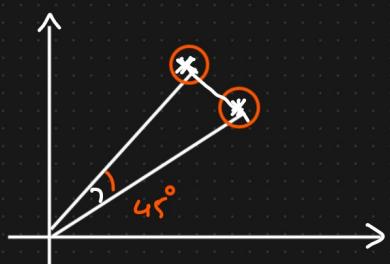
$$\text{Man} [0.95, 0.98]$$

$$\text{QUEEN} [-0.96, 0.95]$$

$$\text{Women} [-0.94, -0.96]$$

$$\text{KING} - \text{MAN} + \text{QUEEN} = \text{WOMEN}$$

Cosine Similarity



$$\text{Distance} = 1 - \text{Cosine Similarity}$$

$$\text{Cosine-Sim} = \cos 45^\circ = \frac{1}{\sqrt{2}} = 0.7071$$

$$\text{Distance} = 1 - 0.7071 \\ \hookrightarrow = 0.29$$

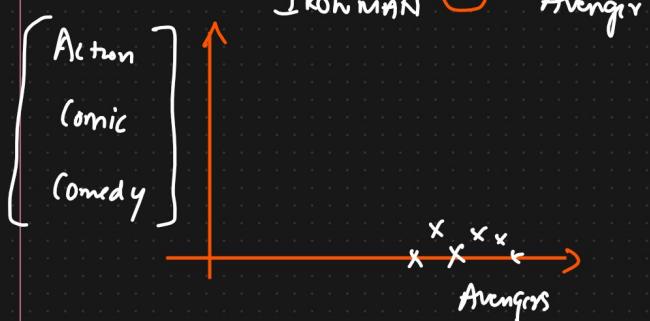
$$\boxed{\sqrt{1-1}=0}$$

$$\text{Distance} = 1 - 0 \\ = 1$$

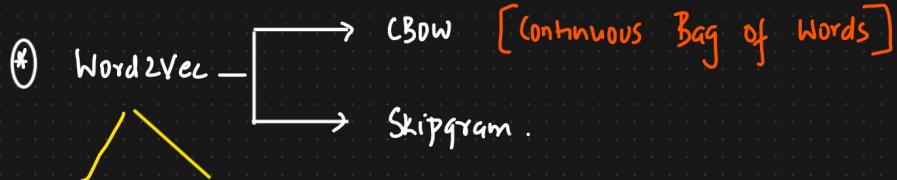


$$\text{Distance} = 1 - \cos 0^\circ$$

$$= 1 - 1 \\ = 0\%$$



ANN, Loss, Optimizers



Pretained Train A

Model Model For Scratch

① CBOW [Continuous Bag of Words]

CORPUS / DATASET



300 dimension

300 ← Window Size = 5 ⇒ Feature Representation [- - - -] ONE

I/P

O/P

IS

iNeuron

[1 0 0 0 0 0 0]

Company [0 1 0 0 0 0 0]

Related [0 0 0 1 0 0 0]

To [0 0 0 0 1 0 0]

Related

To

→ [iNeuron, Company, Related, To]

→ Company, IS, To, DATA

→ IS, Related, DATA, SCIENCE

I/P Layer

ANN

{ Fully Connected Neural Net }

iNeuron

Company

Related

To

CBOW

Initialize Weights

Window

h_{l1}

→

O/P Layer

O/P

Word → vector

[]

y -

0.25

0.33

0

0

0

0

0

0

0

0

0

0

0

0

0

0

0

0

0

0

0

0

0

y - y

→ Loss W

↓ W → Minimal

iNeuron

0.45

Company

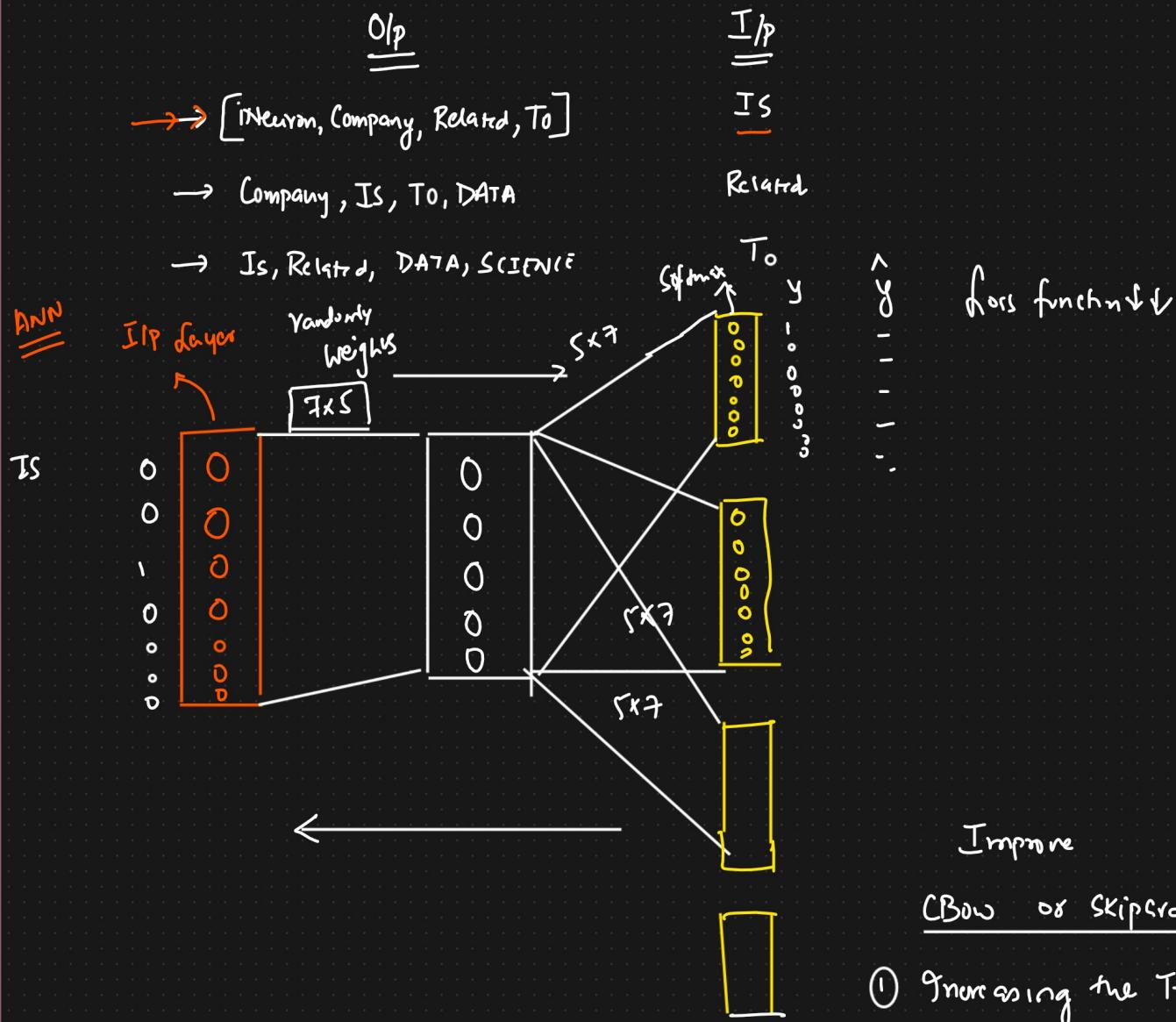
0.45

[0.92, 0.94, 0.25, 0.36,]

Backward Propagation

② Skipgram - Word2Vec

Window Size = 5



When Should we apply CBOW or SkipGram.

- { Small Dataset \rightarrow CBOW }
Huge Dataset \rightarrow SkipGram }

Google Word2vec

Gunsim

3 billion words → Google News

feature representation of 300 dimension vectors]

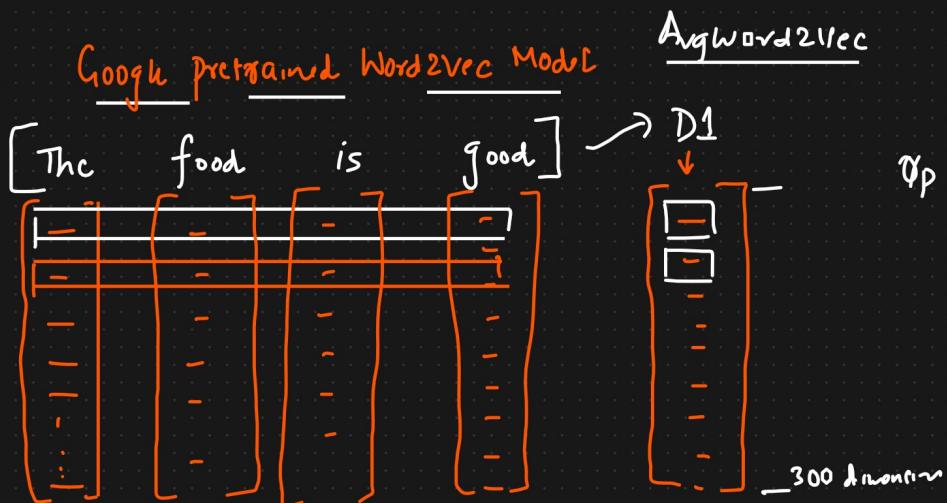
Cricket → [- - - - - - - - - -].

Advantages of Word2Vec

- f) Sparse Matrix \longrightarrow Dense Matrix
- f) Semantic Info is getting captured $\begin{bmatrix} \text{Moust}, \text{good} \end{bmatrix}$
- f) Vocabulary Size \longrightarrow Fixed set of dimension vectors
Google Word2Vec $[300 \text{ dimension}]$
- f) OOV is also solved

Avg Word2Vec

<u>Text</u>	<u>O/P</u>	<u>Avg word2vec</u>	<u>O/P</u>
D1 The food is good	1	$\begin{bmatrix} - & - & - & - & - & - & - \end{bmatrix}$	$\frac{1}{1}$
D2 The food is bad	0	$\begin{bmatrix} \quad \quad \quad \quad \quad \quad \quad \quad \end{bmatrix}$	
D3 Pizza is Amazing	1	$\begin{bmatrix} \quad \quad \quad \quad \quad \quad \quad \quad \end{bmatrix}$	

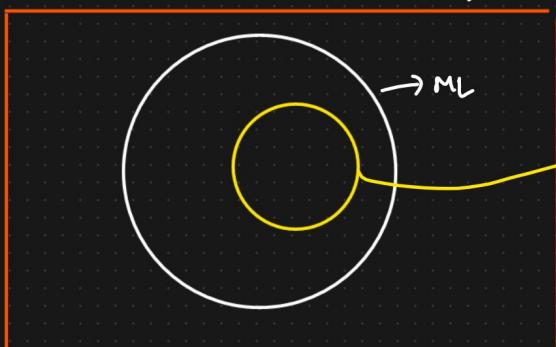


Gensim, Glove

\hookrightarrow pretrained Google Word2Vec

Deep learning

AI \Rightarrow Artificial Intelligence



Deep Learning \Rightarrow Neural Networks
Multi-layered
↓
Mimic the Human Brain

Deep learning

- ① ANN \rightarrow Artificial Neural N/w $\begin{cases} \rightarrow \text{Classification} \\ \rightarrow \text{Regression} \end{cases}$
- ② CNN \rightarrow Convolutional Neural N/w \rightarrow I/P : Image, Video \rightarrow RCNN, MASKED RCNN, frames
Computer Vision
Object Detection
↑
- ③ RNN \rightarrow Recurrent Neural N/w \rightarrow NLP \rightarrow NLP, Time Series
I/P: Text, Time Series
Detection, Yolo V3, V6, V7

FRAMEWORK

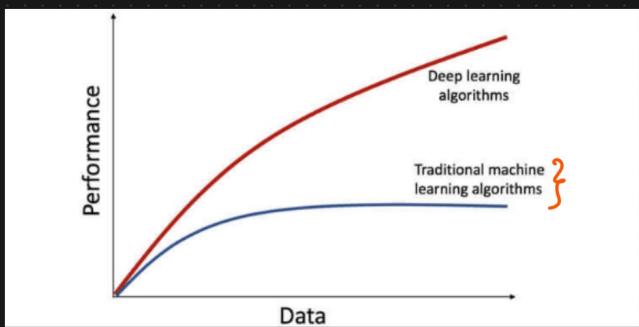
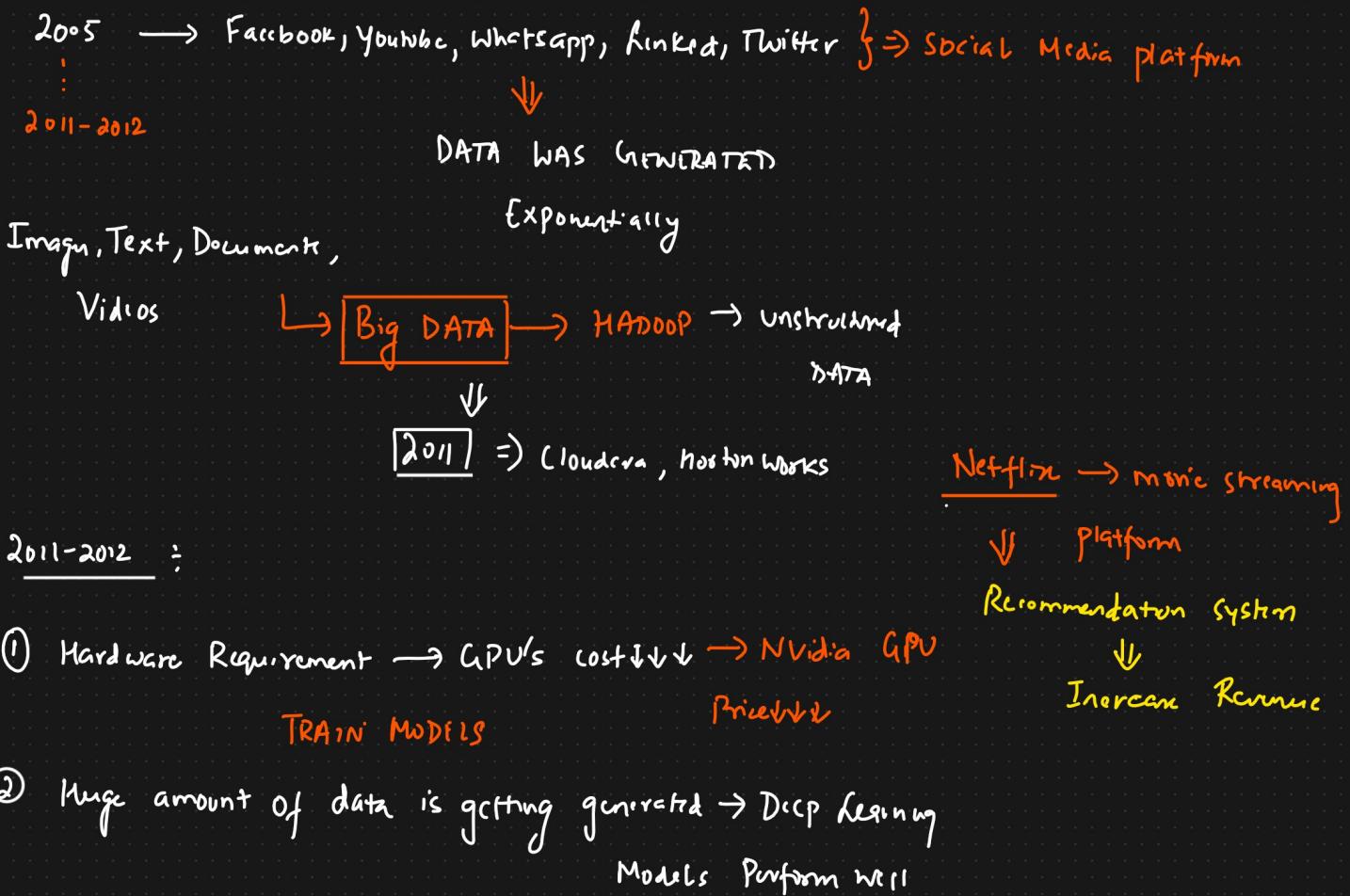
TENSORFLOW

End-to-End Project

{ Word Embedding, LSTM RNN, GRU RNN,

Bidirectional LSTM RNN, Encoder Decoder,
Transformers, BERT }

② Why Deep Learning Is Becoming popular?



③ Deep learning is been used in Many Domains

- ① Medical
 - ② E-commerce
 - ③ Retail
 - ④ Marketing . . .
- ④ Frameworks Open source

Community size ↑↑

Tensorflow



Google

Pytorch



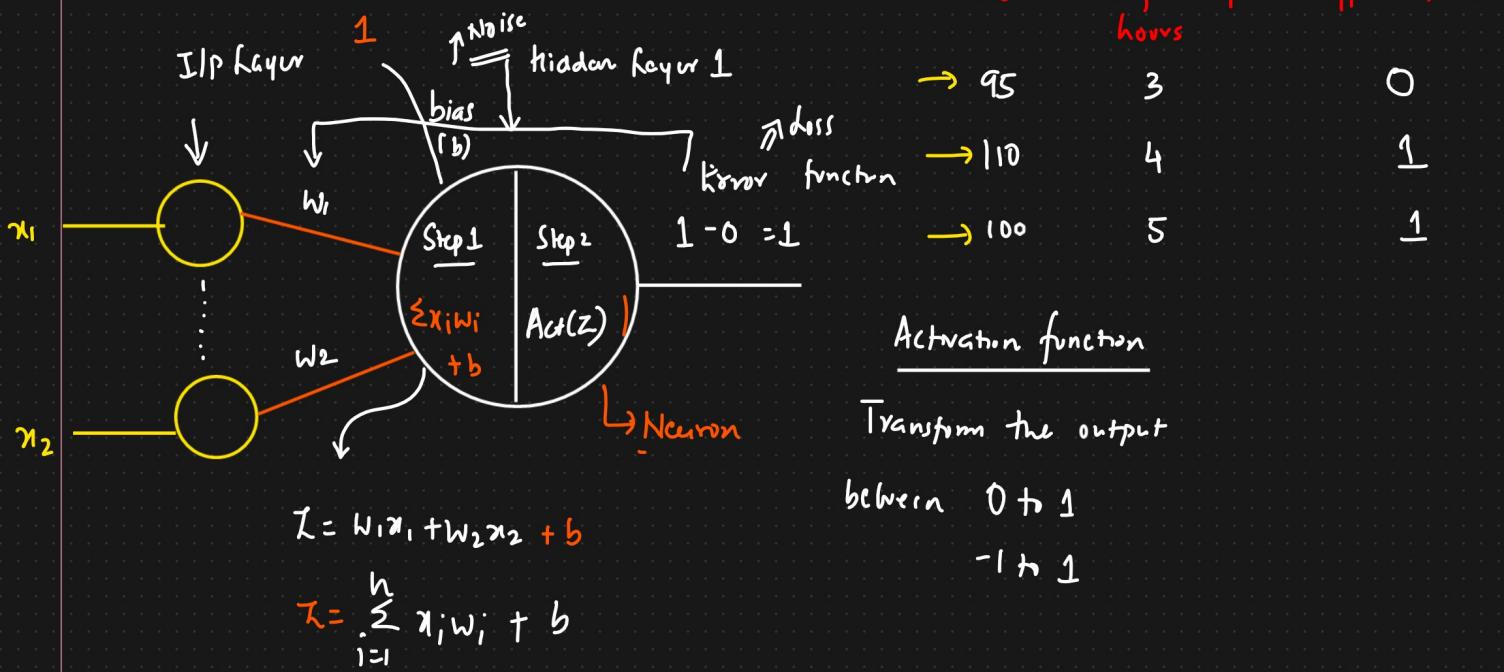
Facebook



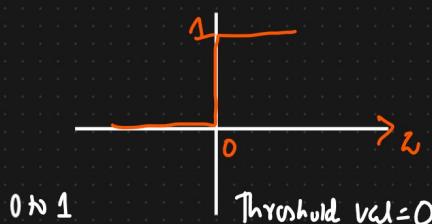
More Research

③ Perceptron [Artificial Neuron or Neural Network Unit]

- ① Input layer ✓
 - ② Hidden layer ✓
 - ③ Weights ✓
 - ④ Activation function ✓
- [Single Layered NN]
- Binary classifier
- DATA SET

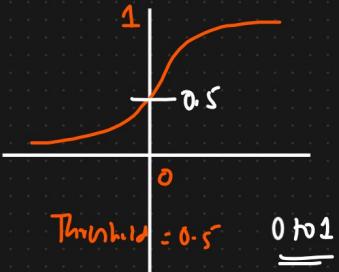


Step Function

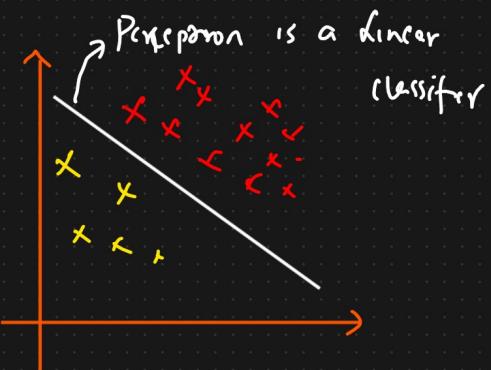


$$0 \leq 1 = \begin{cases} 0 & Z \leq 0 \\ 1 & Z > 0 \end{cases}$$

Sigmoid Function



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



Step 1

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

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

$$y = mx + c$$

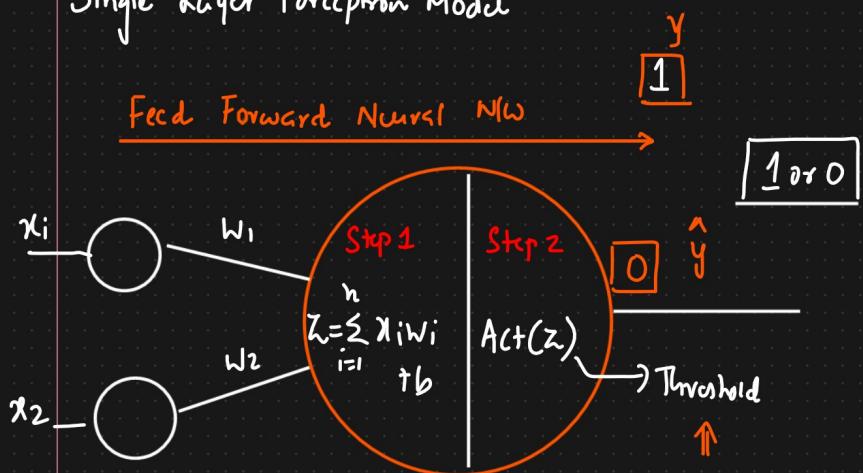
$$y = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \beta_3 x_3 + \dots + \beta_n x_n$$

↓
linear
problem
statement

Perceptron Models

Single Layer Perceptron Model

Feed Forward Neural NW



Multi Layered Perceptron Model

① Forward Propagation

② Backward " "

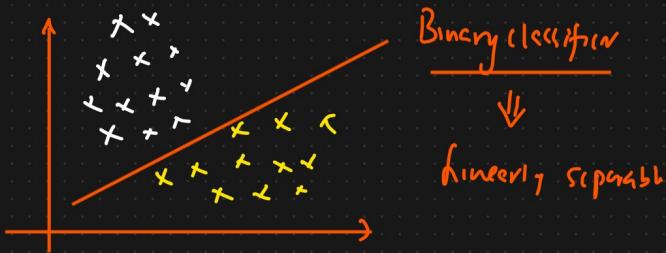
③ Loss functions

④ Activation functions

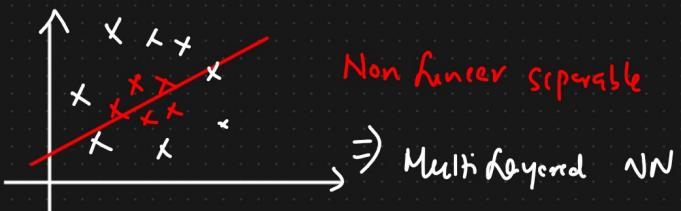
⑤ Optimizers

[ANN]
↓

Linear Separable problem



Non linearly separable



⇒ Multi layered NN

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



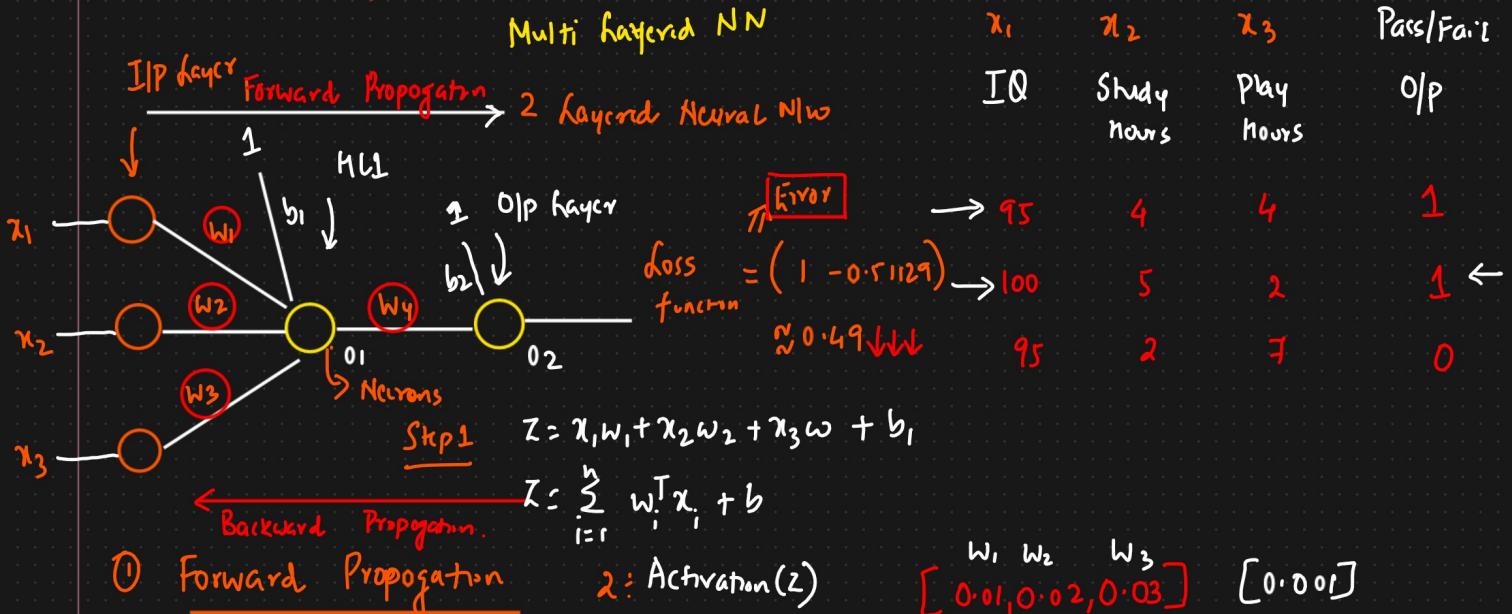
① Forward Propagation

② Backward Propagation → Geoffrey Hinton →

③ Loss function

④ Optimizers ✓

⑤ Activation function.



Hidden Layer 1

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

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

$$\text{Step 2: Activation (Z)}$$

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

$$O_1 = 0.759$$

Hidden Layer 2

$$w_4 = 0.02$$

$$b_2 = 0.03$$

$$\text{Step 1: } Z = O_1 \times w_4 + b_2$$

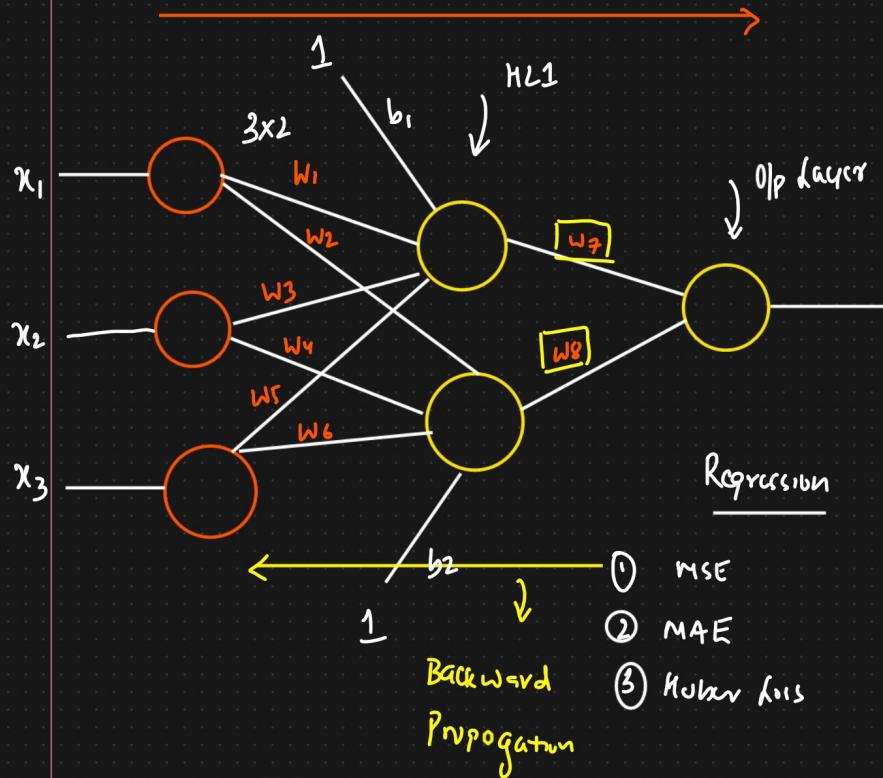
$$= 0.759 \times 0.02 + 0.03$$

$$= 0.04518$$

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



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

Regression

- ① MSE
- ② MAE
- ③ Huber Loss

Classification

- ① Binary Cross Entropy
- ② Categorical Cross Entropy

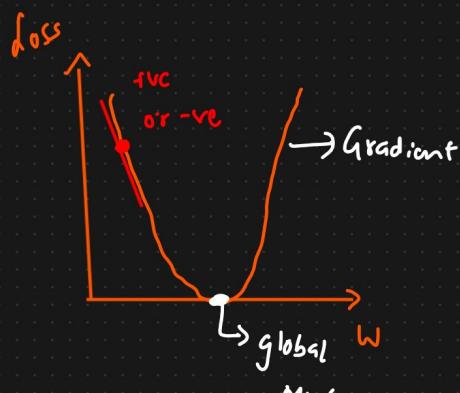
Weight update Formula

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

slope

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

Rate



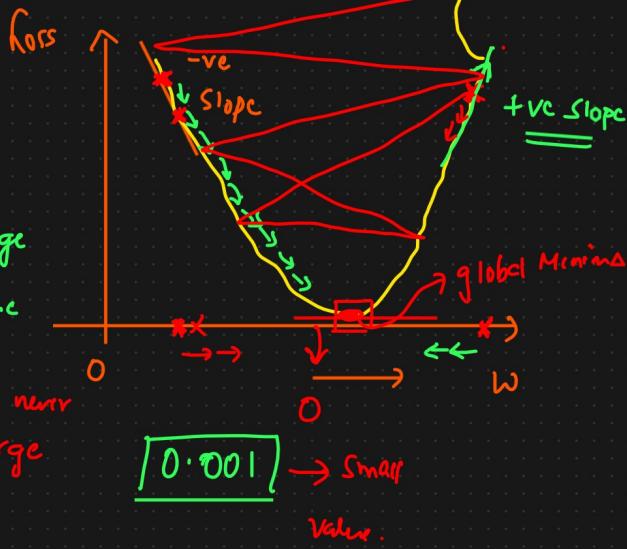
$$w_{\text{new}} = w_{\text{old}} - \eta \left[\frac{\partial h}{\partial w_{\text{old}}} \right] \Rightarrow \text{Weight Updation Formula.}$$

Gradient Descent

↗ Optimizers

I/P featur.s			Pass/Fail
x_1	x_2	x_3	
IQ	Study hours	Play hours	O/P
> 95	4	4	1
→ 100	5	2	1 ←
95	2	7	0

Optimizers : To reduce the loss value



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

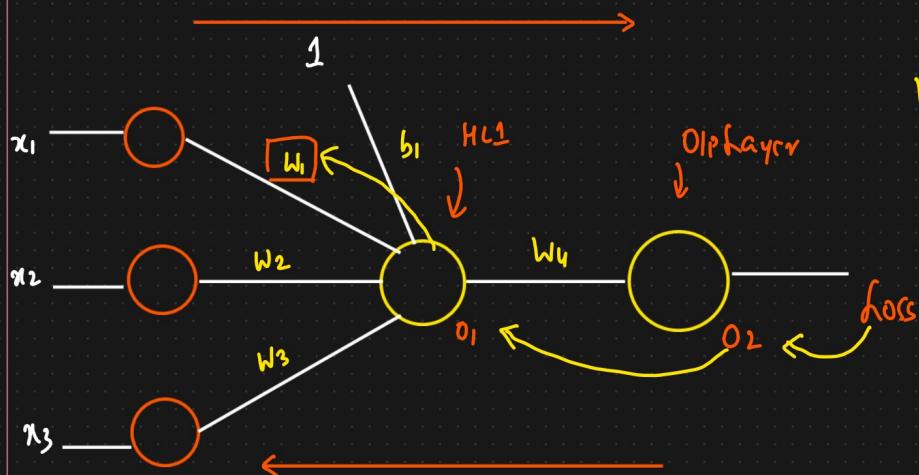
$$= W_{\text{old}} + \eta (+\text{ve})$$

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

When w reaches Global Minima

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

⑦ Chain Rule of Derivative



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

$$W_{4,\text{new}} = W_{4,\text{old}} - \eta \left[\frac{\partial h}{\partial W_{4,\text{old}}} \right]$$

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

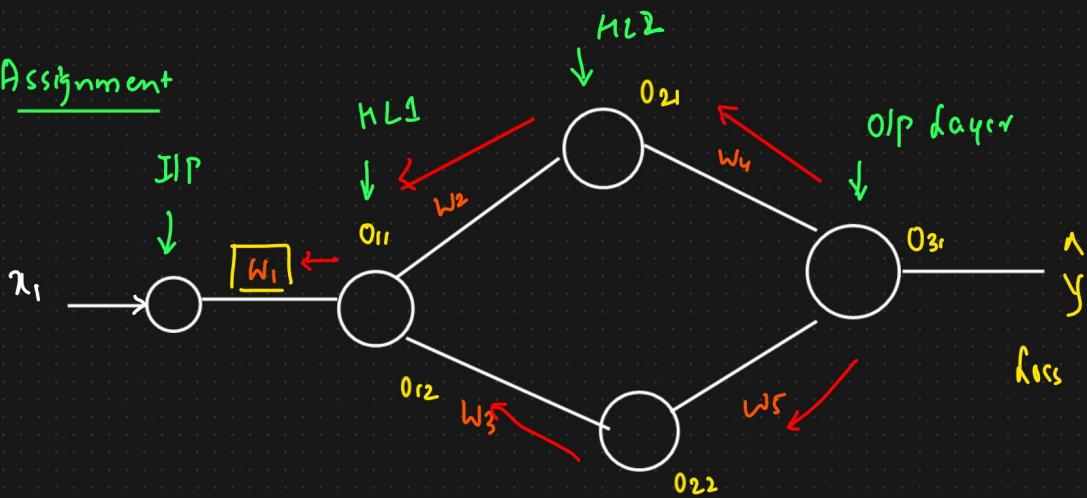
$$w_{i,new} = w_{i,old} - \eta \left[\frac{\partial h}{\partial w_{i,old}} \right]$$

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

$w_{2,new}$

$w_{3,new}$

Assignment

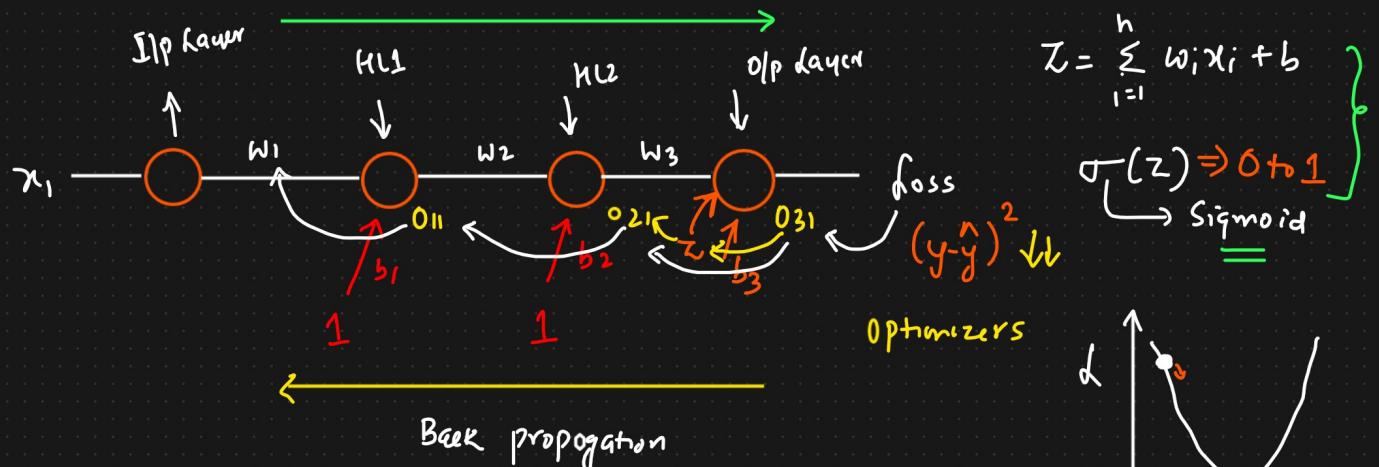


$$w_{i,new} = w_{i,old} - \eta \left[\frac{\partial h}{\partial w_{i,old}} \right]$$

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

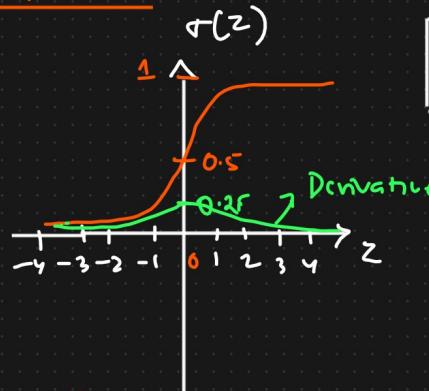
$$+ \left[\frac{\partial h}{\partial o_{31}} * \frac{\partial o_{31}}{\partial o_{22}} * \frac{\partial o_{22}}{\partial o_{12}} * \frac{\partial o_{12}}{\partial w_{i,old}} \right]$$

⑧ Vanishing Gradient Problem And Activation functions



⑨ Sigmoid Activation function

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



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

Derivative ($\sigma(z)$)

Vanishing
Gradient
Problem.

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

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

$$\frac{\partial h}{\partial w_{1,\text{old}}} = \frac{\partial h}{\partial w} * \boxed{\frac{\partial w}{\partial o_{31}}} * \boxed{\frac{\partial o_{31}}{\partial o_{21}}} * \boxed{\frac{\partial o_{21}}{\partial o_{11}}} * \boxed{\frac{\partial o_{11}}{\partial w_{11}}} \Rightarrow \text{smaller value}$$

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

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

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

$$0 \leq \sigma(z) \leq 0.25 \quad * \quad \frac{\partial((w_3 * o_{21}) + b_3)}{\partial(o_{21})}$$

↓

Derivative of
Sigmoid

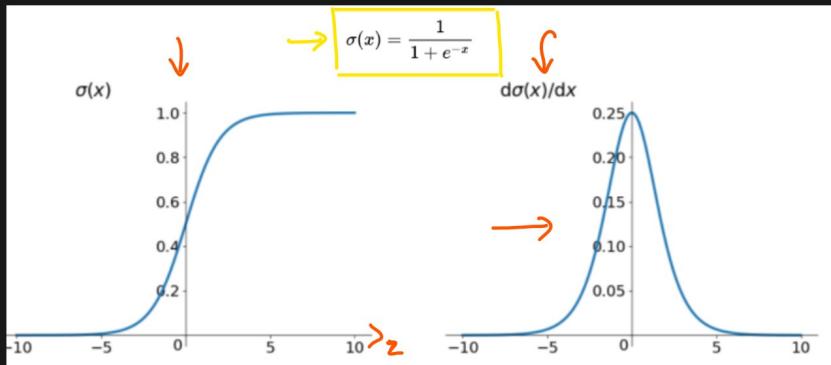
$$\left| \frac{\partial o_{31}}{\partial o_{21}} = 0 \leq \sigma(z) \leq 0.25 * w_3 \right|_{\text{old}}$$

④ To fix this problem Researchers started exploring other Activation function

① Tanh ② ReLU ③ PReLU ④ Swiss

Activation Functions

① Sigmoid Activation function $[0 \text{ to } 1]$ $z = \sum_{i=1}^n w_i^T x + b$

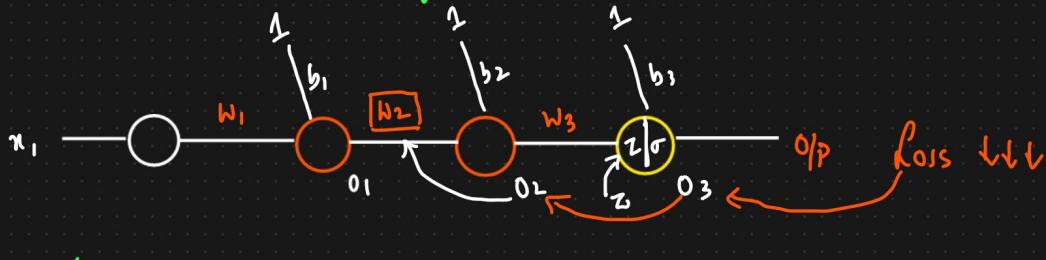


$$\sigma(z) \Rightarrow 0 \text{ to } 1$$

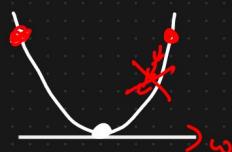
$$\phi(z) \Rightarrow$$

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

Forward Propagation



Backward Propagation



$$w_{2\text{new}} = w_{2\text{old}} - \eta \left[\frac{\partial h}{\partial w_{2\text{old}}} \right] \quad \boxed{0.001} \quad \boxed{0.0001} \quad \Rightarrow w_{2\text{new}} \approx w_{2\text{old}}$$

$$\frac{\partial h}{\partial w_{2\text{old}}} = \frac{\partial h}{\partial o_3} * \boxed{\frac{\partial o_3}{\partial o_2}} * \frac{\partial o_2}{\partial w_2}$$

$\downarrow \downarrow \downarrow$

$$0.20 * 0.01 * \frac{\partial h}{\partial o_3} \quad \det z = (o_2 * w_3) + b_3$$

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

$[0 \text{ to } 1]$

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

$$= [0 - 0.25] * w_3 \Rightarrow \text{Small value} \Rightarrow$$

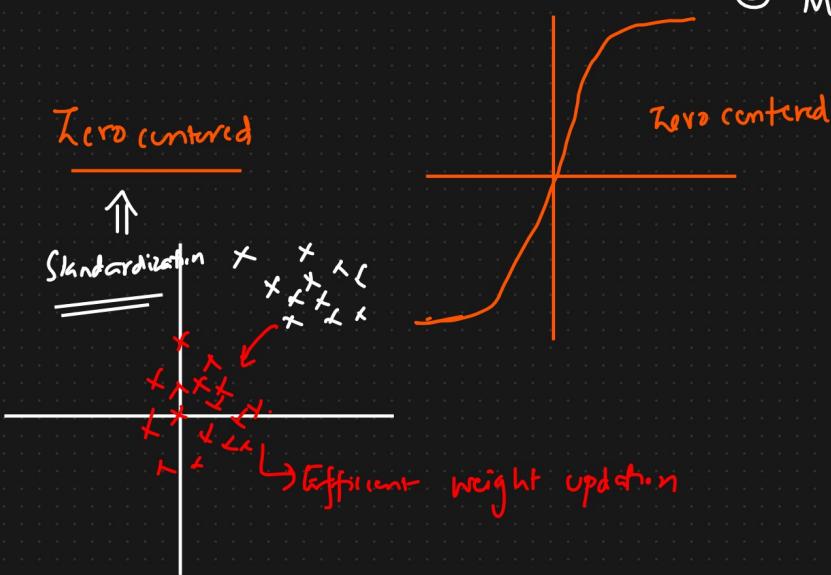
Advantages

- ① Binary Classification Suitable.
- ② clear prediction i.e. very close 1 or 0

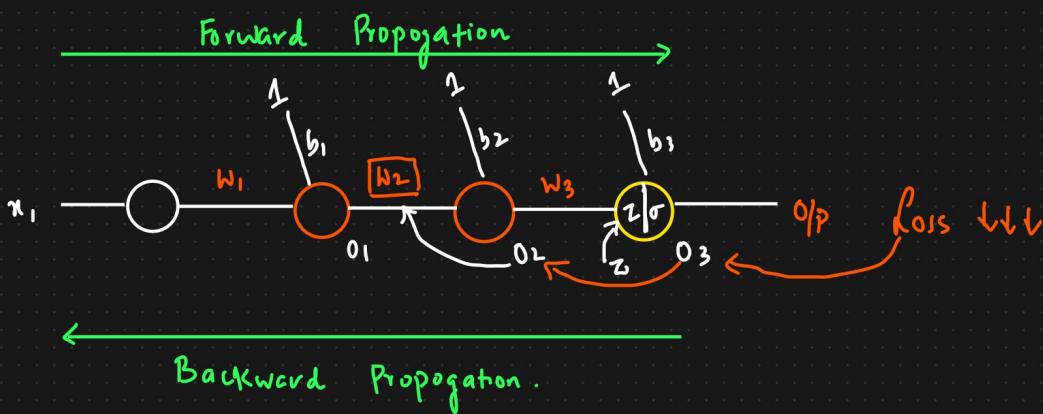
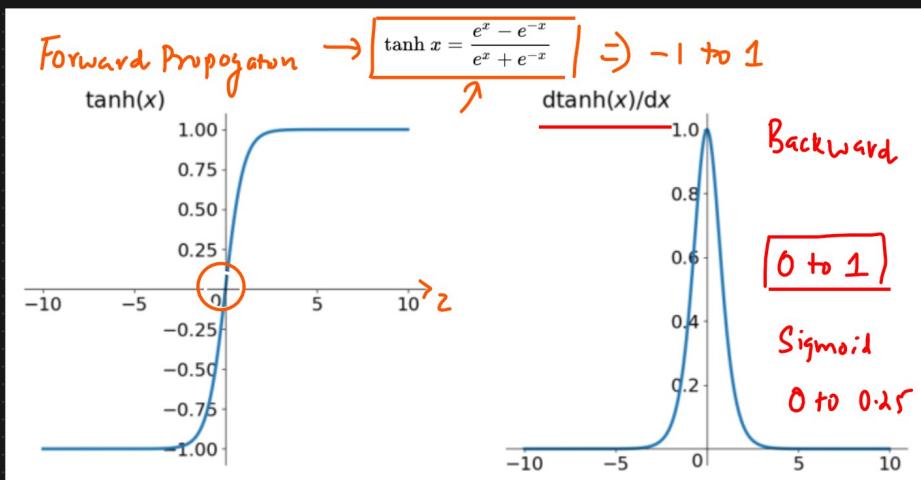
Disadvantages

- 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



$$\frac{\partial h}{\partial w_{2012}} = \frac{\partial h}{\partial o_3} \neq \boxed{\frac{\partial o_3}{\partial o_2}} * \frac{\partial o_2}{\partial w_2}$$

↓↓↓

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

$$\begin{aligned}\frac{\partial o_3}{\partial o_2} &= \boxed{\frac{\partial (\tanh(z))}{\partial z}} * \frac{\partial z}{\partial o_2} && [0 \text{ to } 1] \\ &= [0 - 1] * \frac{\partial [(o_2 * w_3) + b_3]}{\partial o_2} \\ &= [0 - 1] * w_3 \Rightarrow \text{Small value} \Rightarrow\end{aligned}$$

Advantages

① Zero Centric \Rightarrow weight updation is efficient

Disadvantages

① Prone to Vanishing Gradient Problem

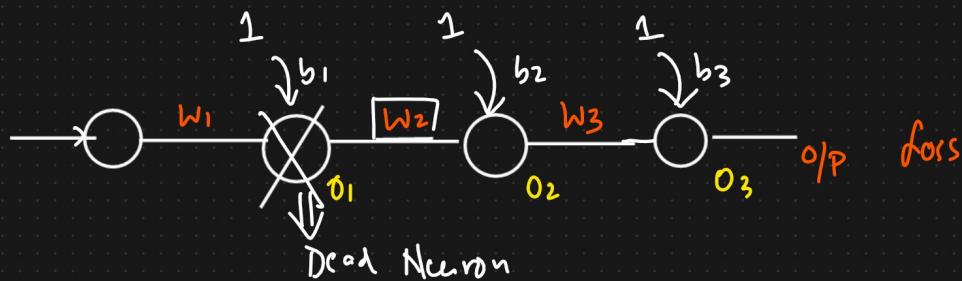
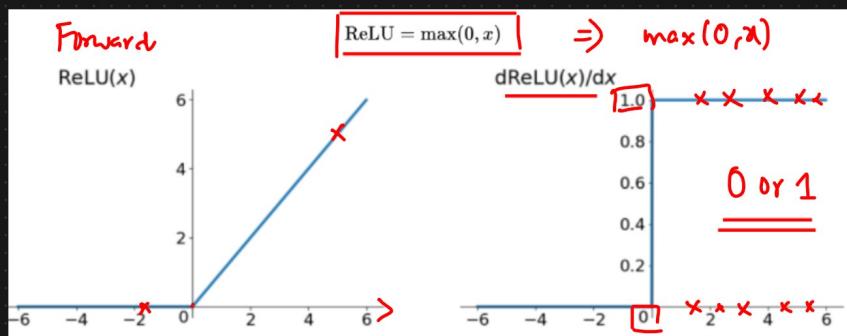
② Train Complexity

[Rectified Linear Unit].

③ ReLU Activation Function

Tanh $\Rightarrow 0 \text{ to } 1$

Sigmoid $\Rightarrow 0 \text{ to } 0.25$



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

$\downarrow \downarrow \downarrow \downarrow$

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

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

$[0 \rightarrow 1]$

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

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

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

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

Advantages

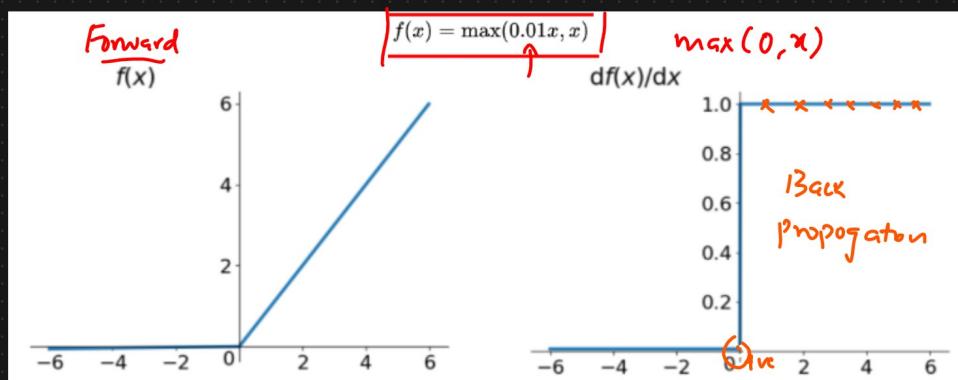
- ① Solving Vanishing Gradient Problem
- ② $\text{Max}(0, x) \rightarrow$ Calculation is Superfast. The ReLU function has a linear relationship.
- ③ It is much faster than Sigmoid or Tanh.

Disadvantages

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

④ Leaky ReLU And Parametric ReLU

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



ReLU → Dead Neuron → Dead ReLU Problem

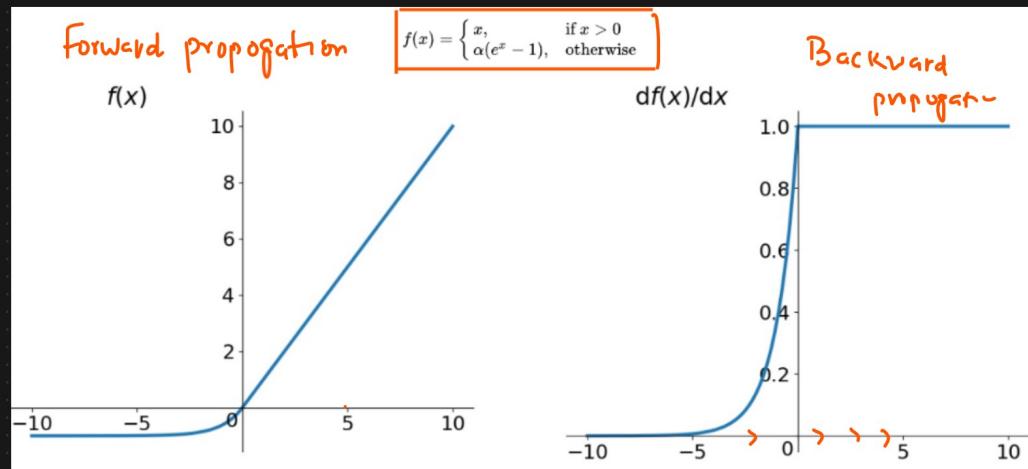
Advantages

- 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)



Advantages

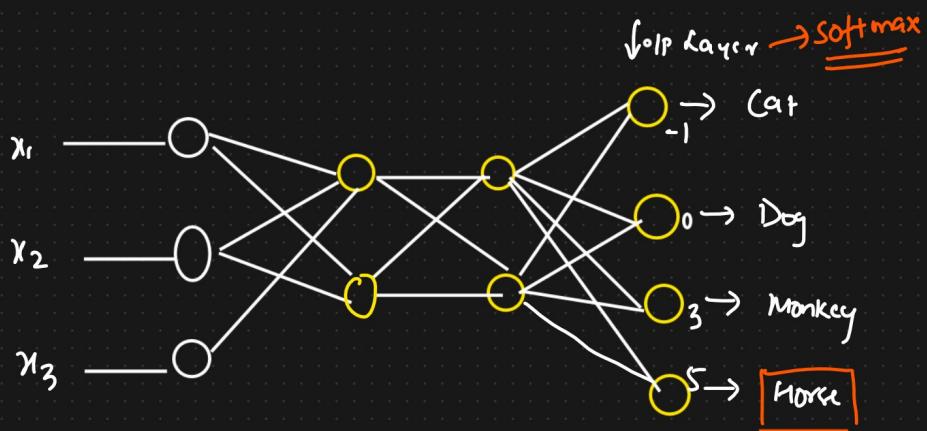
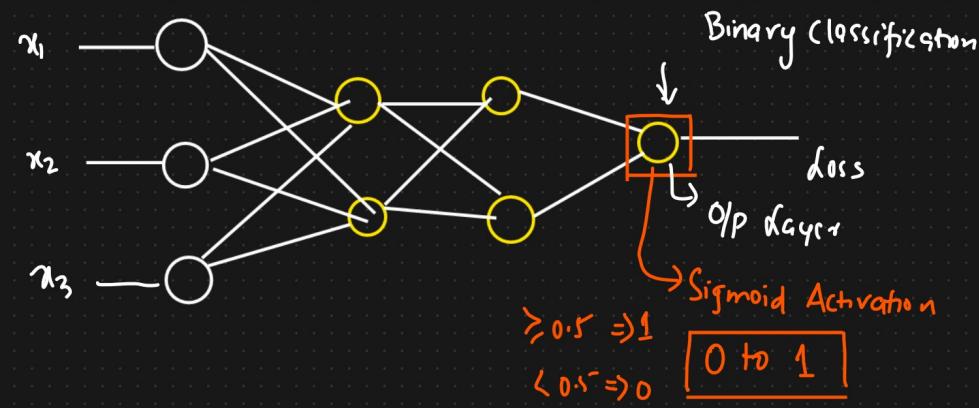
[It is used to solve
ReLU problems]

Disadvantage

- ① No Dead ReLU Issues
- ② Zero centered

i) Slightly more computationally intensive.

⑥ Softmax Activation function [Multiclass classification problem]



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

$$y_i = \theta \cdot w + b$$

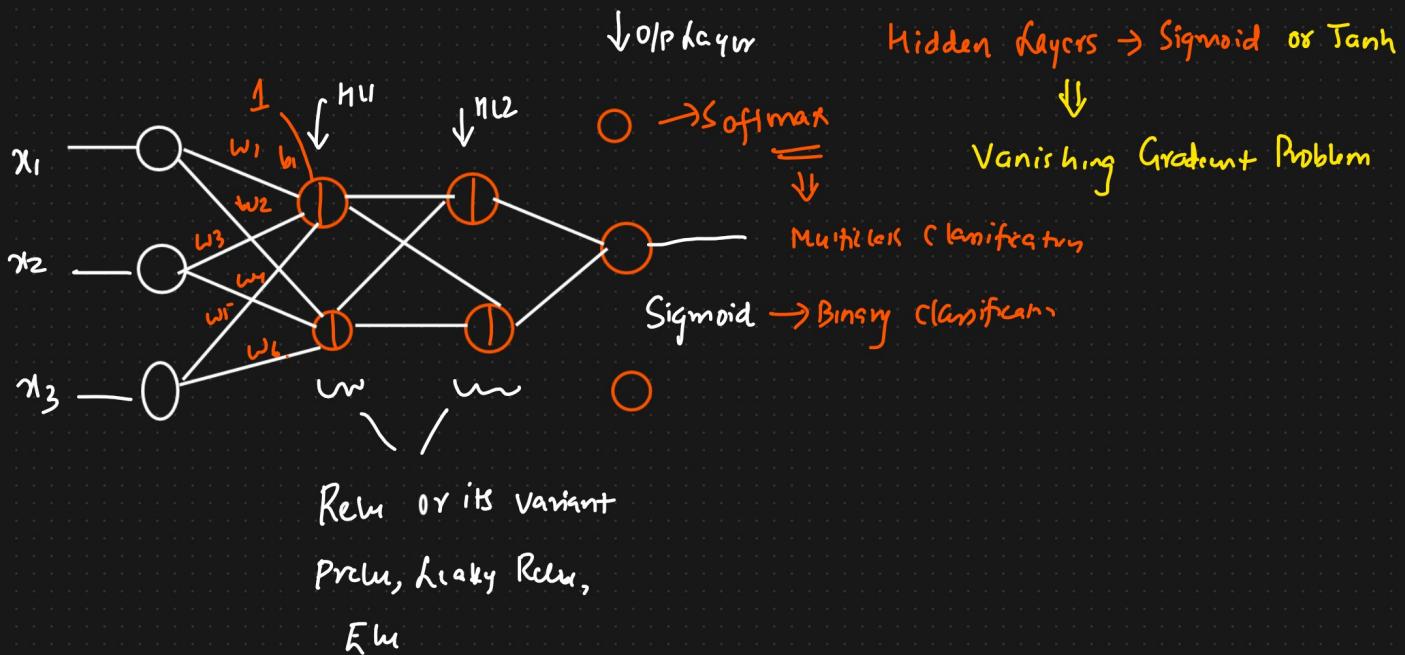
$$\text{Softmax} \Rightarrow \text{Cat} = \frac{e^{-1}}{e^{-1+0+3+5}} = 0.00033 \quad \Pr(\text{Horse}) = \frac{0.1353}{0.00033 + 0.0024 + 0.0183 + 0.1353}$$

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

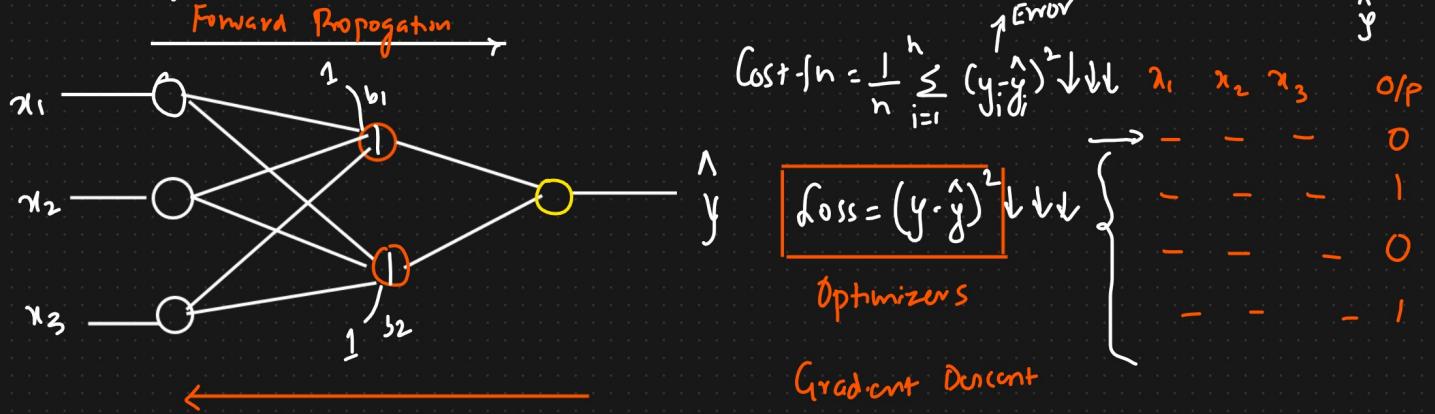
$$\text{Monkey} = \frac{e^3}{e^{-1+0+3+5}} = 0.0183$$

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

⑦ Which Activation Function To Use When?



Loss function And Cost Function

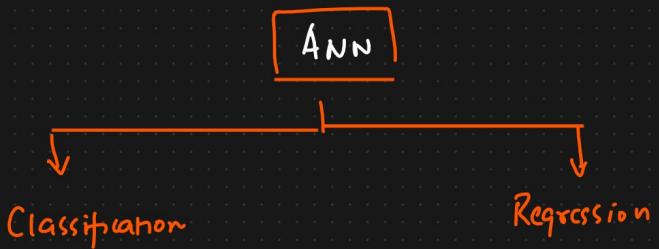


Loss function

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

Cost function

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



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

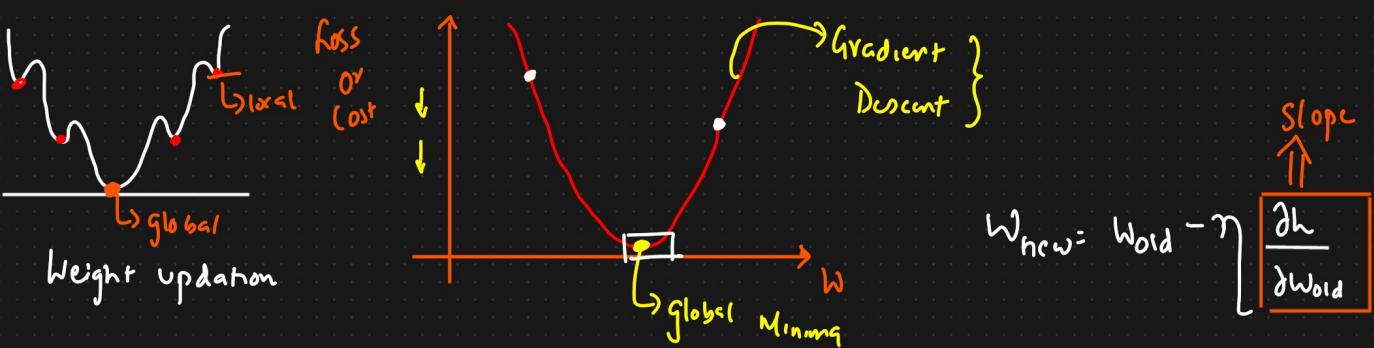
① Mean Squared Error (MSE)

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

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

↳ Quadratic Equations



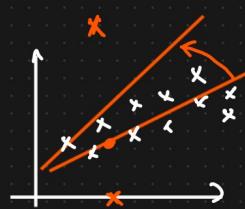


Advantages

- ① MSE is Differentiable
- ② It has 1 local or global Minima
- ③ It converges faster
- ④ Mean Absolute Error (MAE)

Disadvantages

- ① Not Robust to outliers penalizing the error



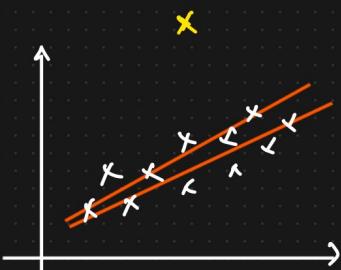
Cost function = $\frac{1}{n} \sum_{i=1}^n (y_i - \hat{y}_i)^2$

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

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

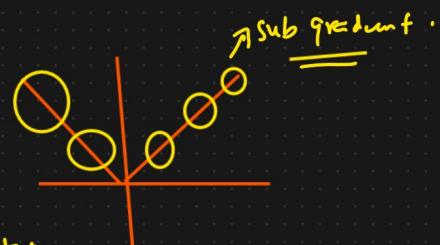
Advantages

- ① Robust to outliers



Disadvantages

- ① Convergence usually takes time in MAE



③ Huber loss

① MSE

② MAE

Cost fn = $\begin{cases} \text{MSE} & \text{if } |y - \hat{y}| \leq \delta \\ \text{MAE} & \text{if } |y - \hat{y}| > \delta \end{cases}$ No outlier \uparrow Hypo-parameter \uparrow

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

④ RMSE (Root Mean Squared Error)

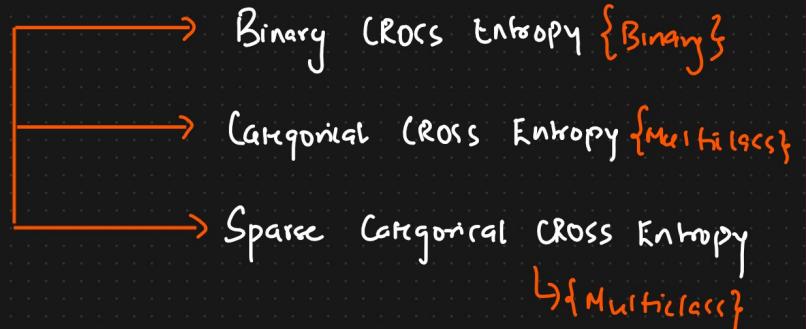
$$\text{Cost function} = \sqrt{\frac{\sum_{i=1}^N (y_i - \hat{y}_i)^2}{N}}$$

Advantages

Disadvantages

② Loss or Cost function For Classification Problems

Classification \rightarrow Cross Entropy



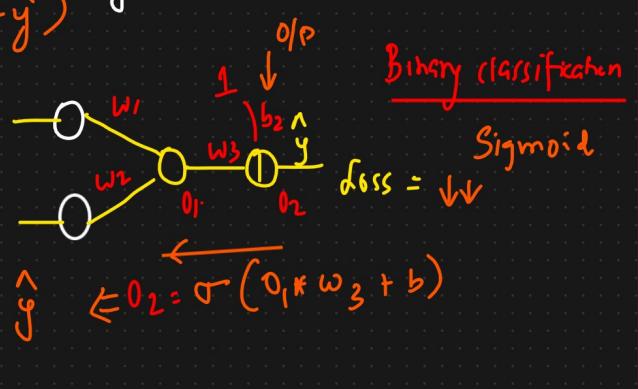
① Binary Cross Entropy

\downarrow Log Loss

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

y = Actual Value

\hat{y} = Predicted Value



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

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

② Categorical Cross Entropy (Multiclass classification)

ONE → One Hot Encoding

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

$i = 1 \rightarrow n$

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

Actual value $\leftarrow y_{ij} = [y_{11} \ y_{12} \ y_{13} \ \dots \ y_{1c}]$

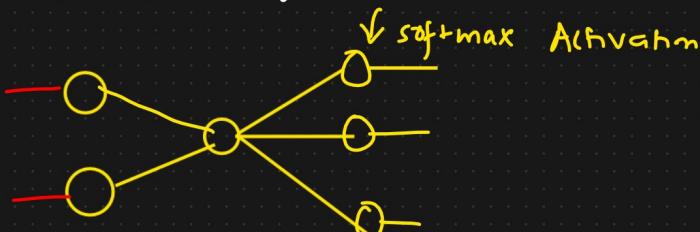
$$[y_{21}, y_{22}, y_{23}, \dots, y_{2c}]$$

⋮

⋮

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

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



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

O/p \hat{y}_{ij} = Probabilities

$$[0.1, 0.2, 0.3, 0.2, 0.2] \leftarrow 1$$

Categorical $\Rightarrow [0.2, 0.3, 0.5]$

(Cross Entropy)

 .

\downarrow
[This also gives the probability of other categories]

③ Sparse Categorical Cross Entropy

$[0, 1^{\text{st}}, 2^{\text{nd}}]$ \rightarrow Categories
 $[0.2, 0.3, 0.5]$

Disadvantage
① losing info about the probability of other category.

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

$[0, 1, 2^{\text{nd}}, 3^{\text{rd}}, 4^{\text{th}}]$
 $[0.2, 0.3, 0.1, 0.2, 0.2]$ $\Rightarrow 1^{\text{st}} \text{ Index}$
Category $\Rightarrow \underline{\text{O/p}}$

④ Right Combination

Activation applied

\downarrow

Hidden layers

O/p layer

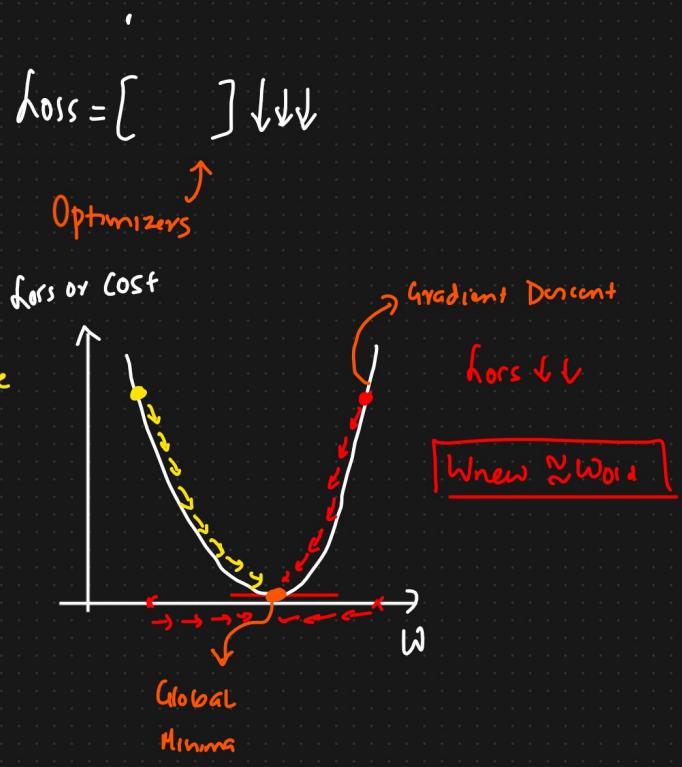
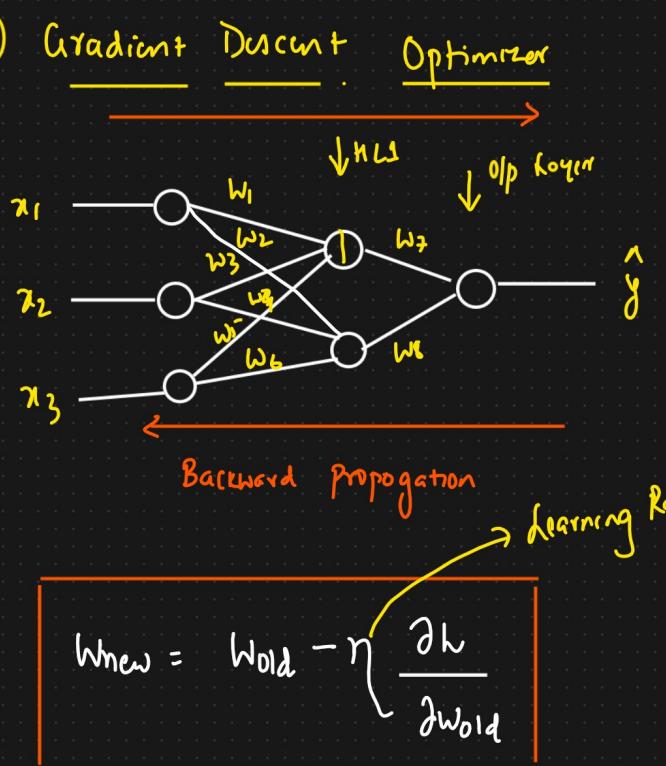
Problem Statement

Loss function

①	ReLU or its variants	Sigmoid	Binary Classification	Binary Cross Entropy
②	ReLU or its variants	Softmax	Multi Class	Categorical or Sparse CE
③	ReLU or its variants	Linear	Regression	MSE, MAE, L1/L2 loss, RMSE

Optimizers

- ① Gradient Descent
 - ② Stochastic Gradient Descent (SGD)
 - ③ Mini batch SGD
 - ④ SGD With Momentum
 - ⑤ Adagrad and RMSprop
 - ⑥ Adam Optimizers



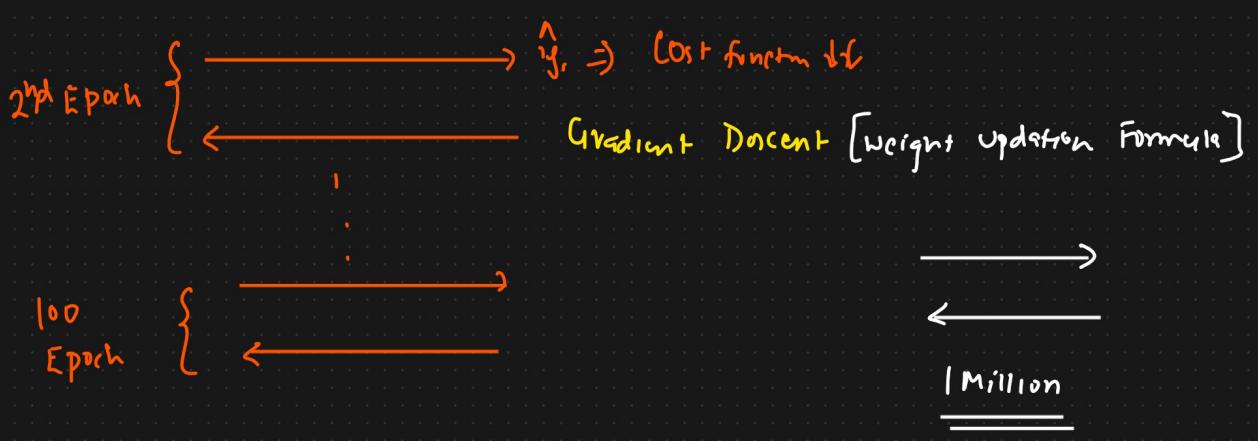
$$\text{Loss fn} = (y - \hat{y})^2 \quad \text{Cost fn} = \frac{1}{n} \sum_{i=1}^n (y_i - \hat{y}_i)^2 \quad \text{Data size} = 1000 \text{ Data points}$$

Epochs, Iteration

1 Epoch = 1 Iteration

10 Iteration

Gradient Descent optimizers [Weight update]



Advantages

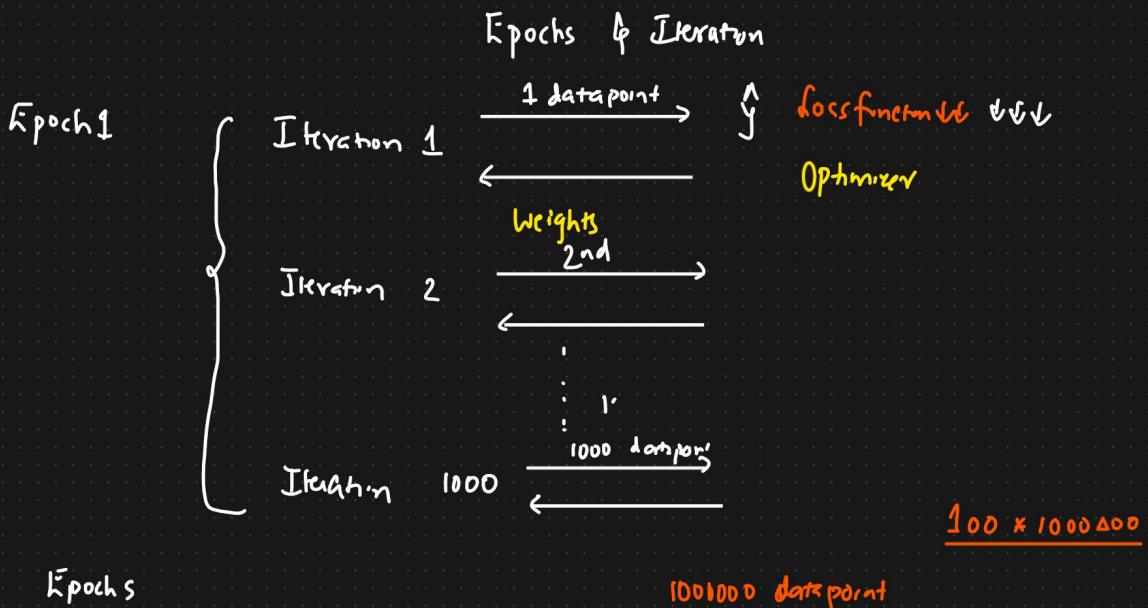
① Convergence will happen.

① Huge Resource RAM, GPU
 \Downarrow

Resource Intensive -

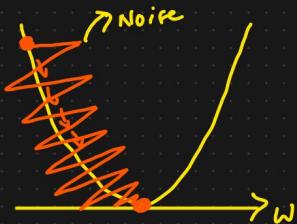
② Stochastic Gradient Descent (SGD)

1000 datapoint



Advantage

① Solve Resource Issue



Disadvantage

① Time Complexity \uparrow

- ④ Convergence will also take more time.
- ④ Noise gets introduced

$\rightarrow \hat{y}$ Cost

③ Mini Batch SGD

Epoch, Iteration, Batch-size

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

Data points = 100000

batch.size = 1000

MSGD

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

Epoch 1

Iteration 1

Optimizer \Rightarrow Mini Batch SGD

change the weights

1000

\downarrow

[8gb] 1

Iteration 2

[16gb] 5000

Iteration 3

Iteration 100

1000

1000

1000

1000

1000

1000

1000

1000

1000

1000

1000

1000

1000

1000

1000

1000

1000

1000

1000

1000

1000

1000

1000

1000

1000

1000

1000

1000

1000

1000

1000

1000

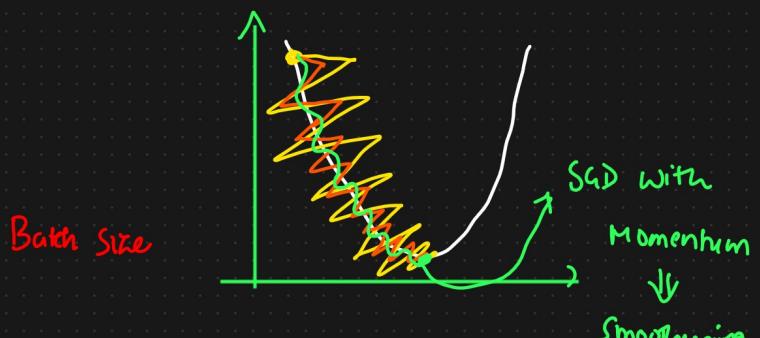
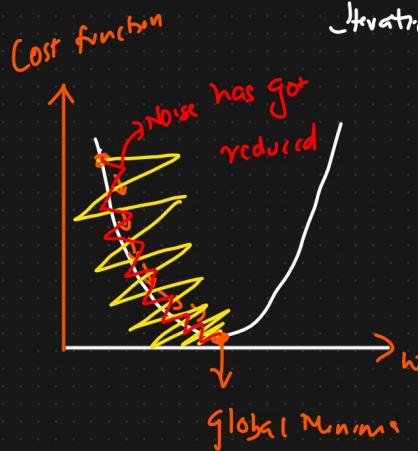
1000

1000

1000

1000

1000



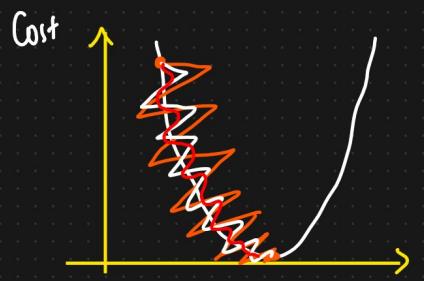
Advantages

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

Disadvantage

- ① Noise still exists

④ SGD With Momentum



Weight Update formula

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

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

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

Exponential Weight Average {Smoothing} } \Rightarrow ARIMA, SARIMAX

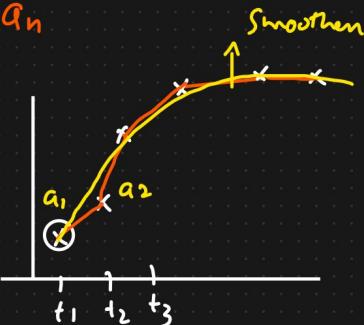
Time = $t_1, t_2, t_3, t_4, \dots, t_n$ Time Series

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

$$V_{t_1} = a_1$$

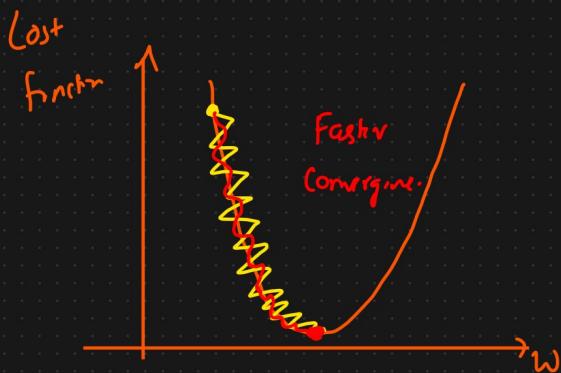
$$V_{t_2} = \boxed{\beta} * V_{t_1} + (1-\beta) * a_2$$

$$\begin{aligned} \beta = 0.95 \\ V_{t_2} = 0.95 * a_1 + (0.05) a_2 \end{aligned}$$



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

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

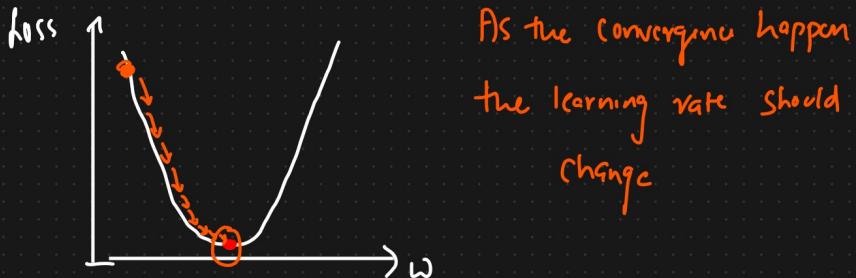


Advantage

- ① Reduces the noise
- ② Quick Convergence

⑤ Adagrad = Adaptive Gradient Descent $\eta = \underline{\text{fixed}} \Rightarrow \underline{\text{Dynamic learning}}$

$$w_t = w_{t-1} - \eta \left[\frac{\partial h}{\partial w_{t-1}} \right] \quad \begin{array}{l} \text{Learning} \\ \text{Rate} = 0.001 \end{array}$$



$$w_t = w_{t-1} - \eta' \left[\frac{\partial h}{\partial w_{t-1}} \right] \quad \begin{array}{l} 0.00001 \approx 0 \\ d_t \uparrow \quad \eta' \downarrow \quad d_t \uparrow \end{array}$$

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

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

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

Disadvantage

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

⑥ Adadelta And RMSprop

Exponential Weight Average

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

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

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

Dynamic LR + Smoothing [EWA]

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

⑦ Adam Optimizer

SGD with Momentum + RMSprop [Dynamic LR + Smoothing]

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

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

Weight Updation

Bias Updation

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

EWA

$S_{dw_t} = 0$

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

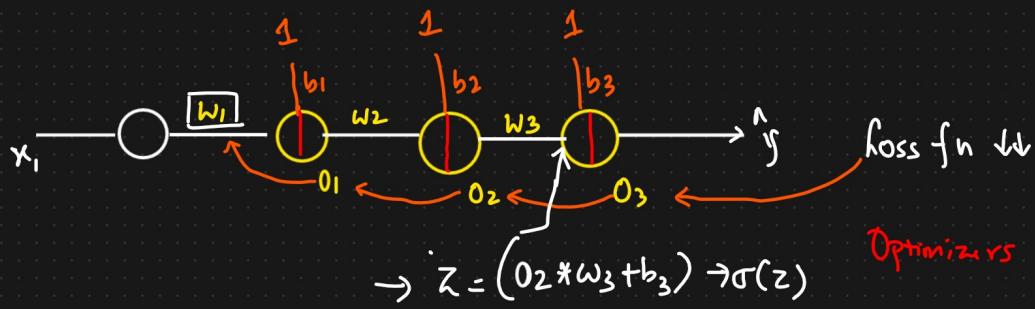
[EWA]

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

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

\Rightarrow Momentum
↳ Smoothing

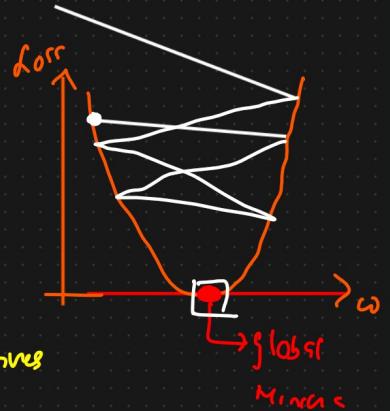
Exploding Gradient Problem \Rightarrow Weight Initialization Technique



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

$\left\{ \begin{array}{l} w_{\text{new}} \ggg w_{\text{old}} \\ w_{\text{new}} \lll w_{\text{old}} \end{array} \right\}$

Chain Rule of Derivatives



$$\frac{\partial h}{\partial w_{\text{old}}} = \frac{\partial h}{\partial o_3} * \left[\frac{\partial o_3}{\partial o_2} \right] * \frac{\partial o_2}{\partial o_1} * \frac{\partial o_1}{\partial w_{\text{old}}}$$

big * big * big * big \Rightarrow big value

Weight initialization



Very high value

$$\frac{\partial o_3}{\partial o_2} = \left[\frac{\partial \sigma(z)}{\partial z} \right] * \frac{\partial z}{\partial o_2}$$

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

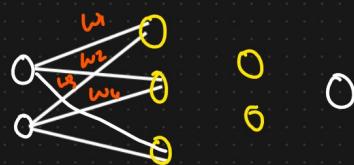
$$= [0 - 0.25] * \frac{w_3}{w_3} = \underline{\underline{5.00 - 1000}}$$

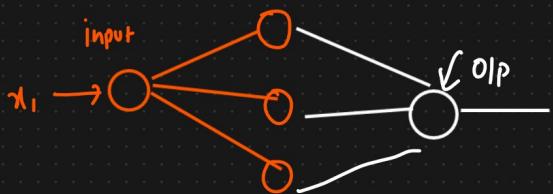
Weight Initializing Techniques

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

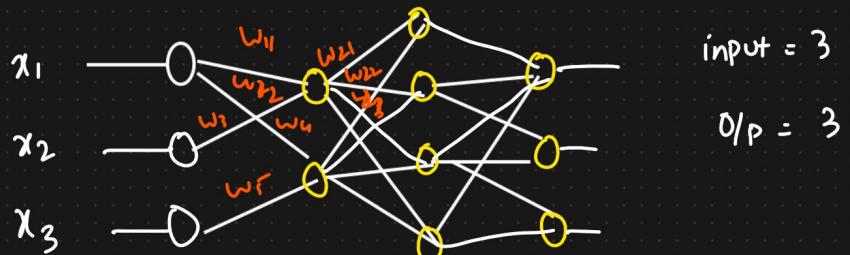
Key Points

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





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 → 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

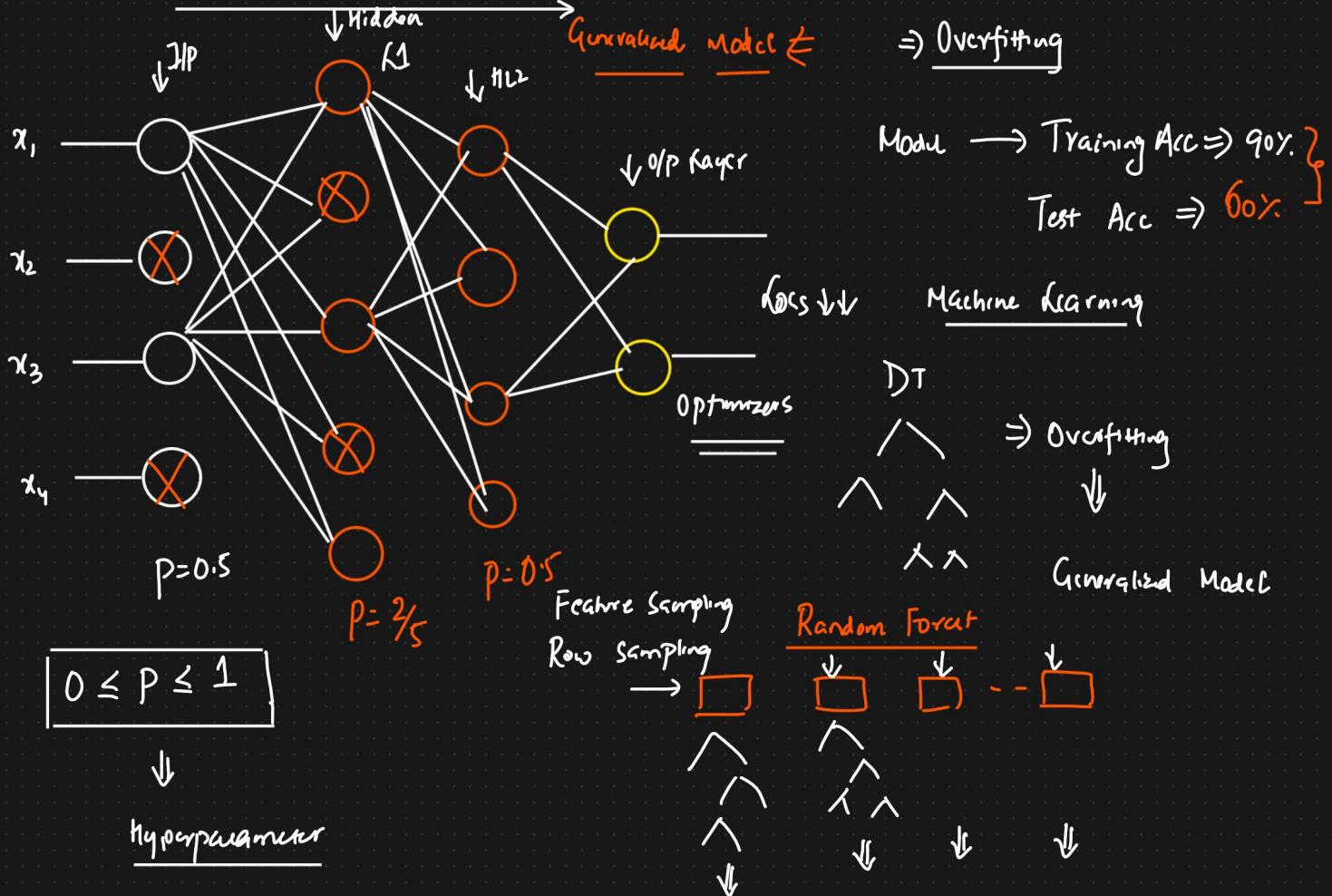
$$w_{ij} \sim N(0, \sigma)$$

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

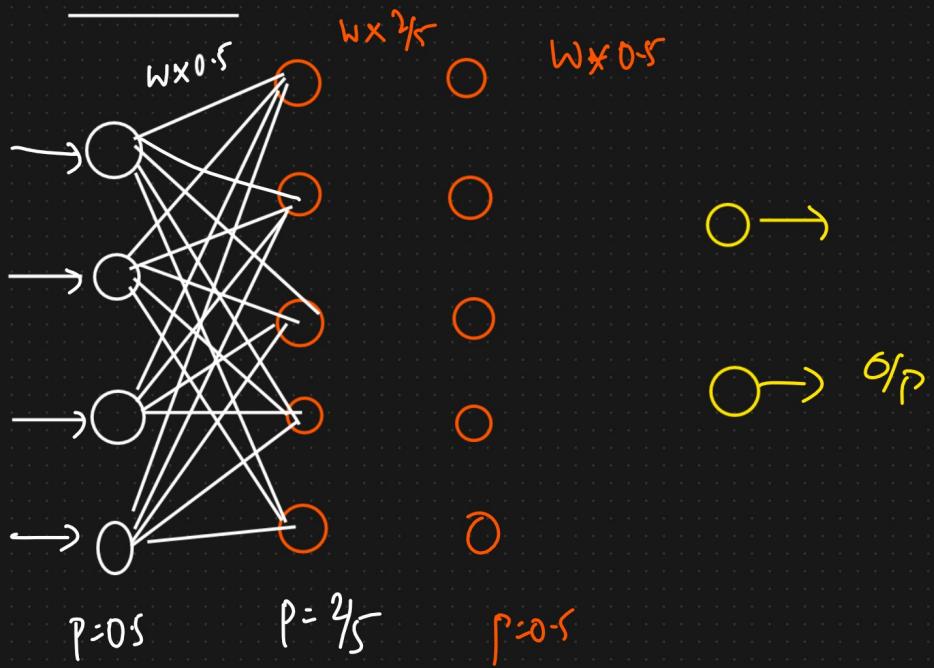
② He uniform

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

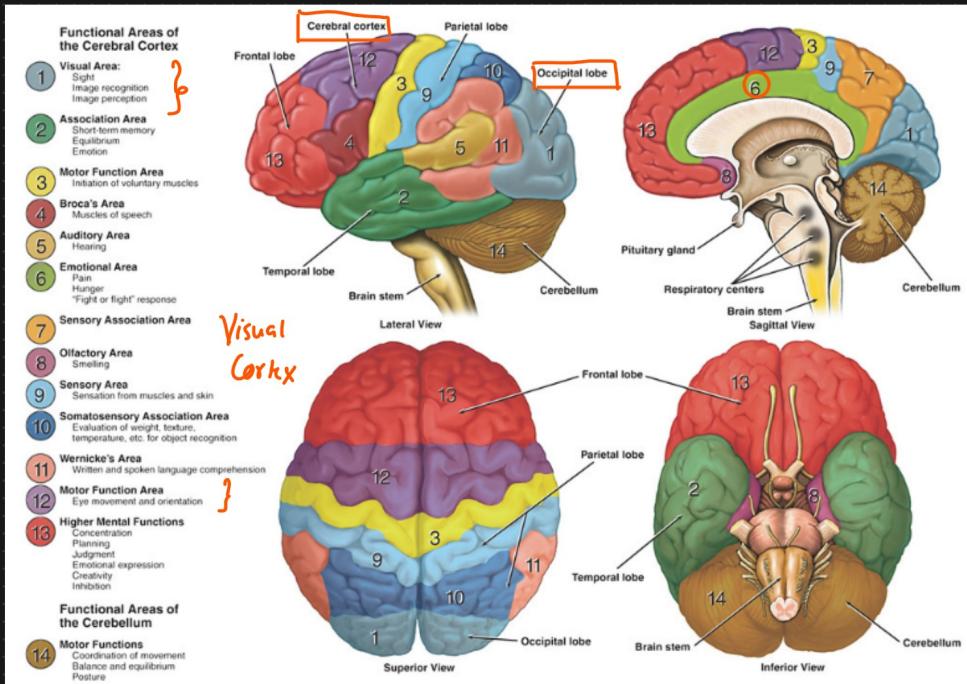
Drop Out layer



Test data



Convolutional Neural N/w



<https://www.dana.org>

① ANN → Supervised Learning →
 Classification
 Regression

Dataset : I/p features O/P

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

② Cerebral Cortex And Visual Cortex

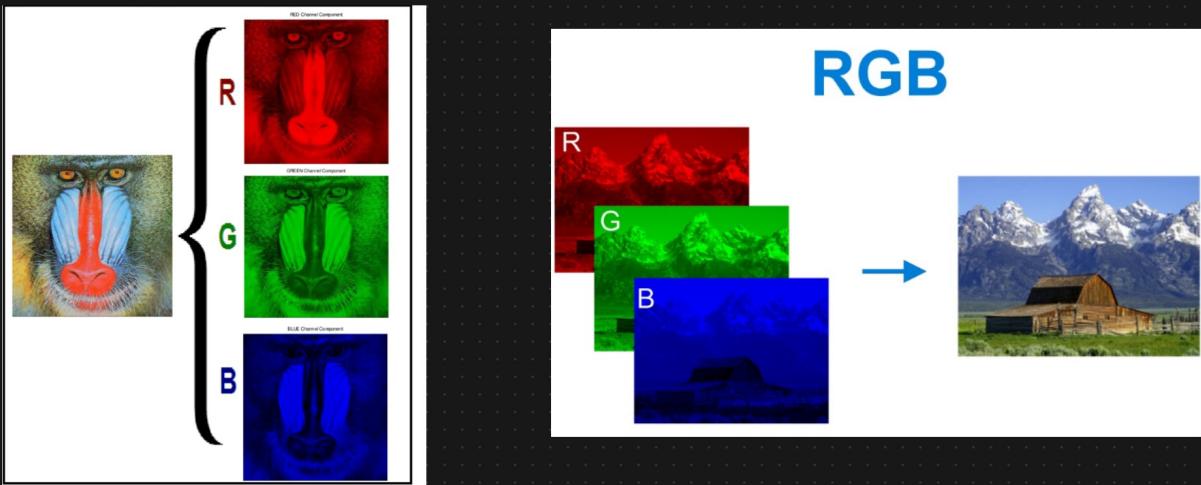
Visual Cortex (VI-V5) [Region of the brain that receives, integrates and processes visual information relayed from the retinas].

↓
→ VI → Primary Visual Cortex [Orientation of edges And lines]

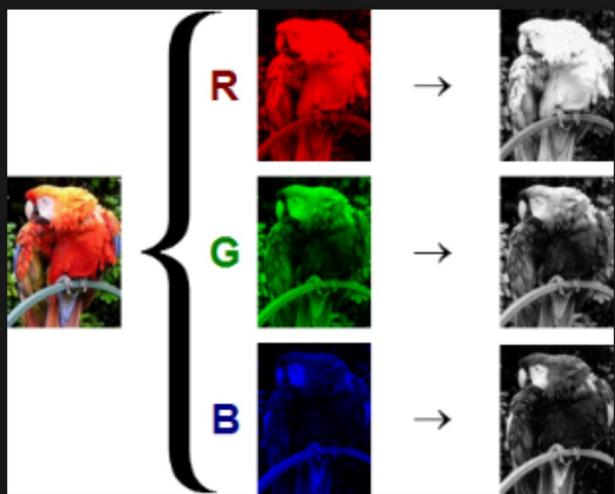
V1
V2 [Differences in color, complex patterns, object orientation]
V3 V4 V5 [Object Recognition]

Visualize the Image

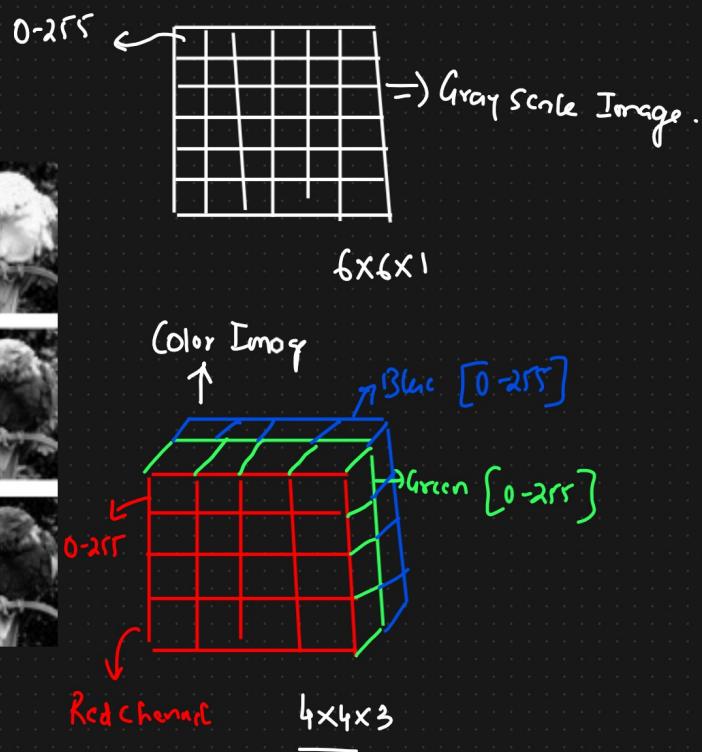
③ RGB Images And Gray Scale Images



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



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



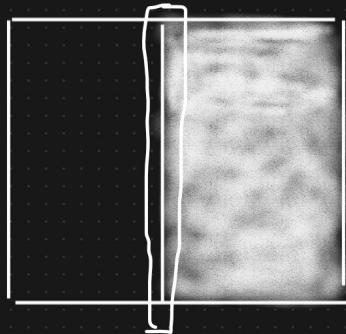
④ Convolution Operation In CNN

→ (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

→

6x6x1



Convolution operation

Step 1

① Normalize

Divide by 255

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

Stride=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

$n=6$

*

$f=3$

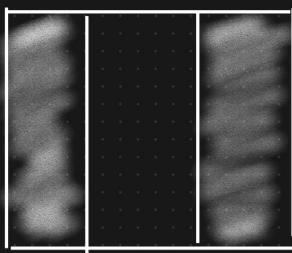
+1	0	-1
+2	0	-2
+1	0	-1

=

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

4×4

$6 \times 6 \times 1$



filters
vertical edge filters

\leftarrow

arr	0	0	arr
arr	0	0	arr
arr	0	0	arr
arr	0	0	arr

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

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 Zero padding

② Neighbour padding

$f=3$

+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	

6×6
=

$n-f+2p+1$

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

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

$$2p = 6 - 4$$

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

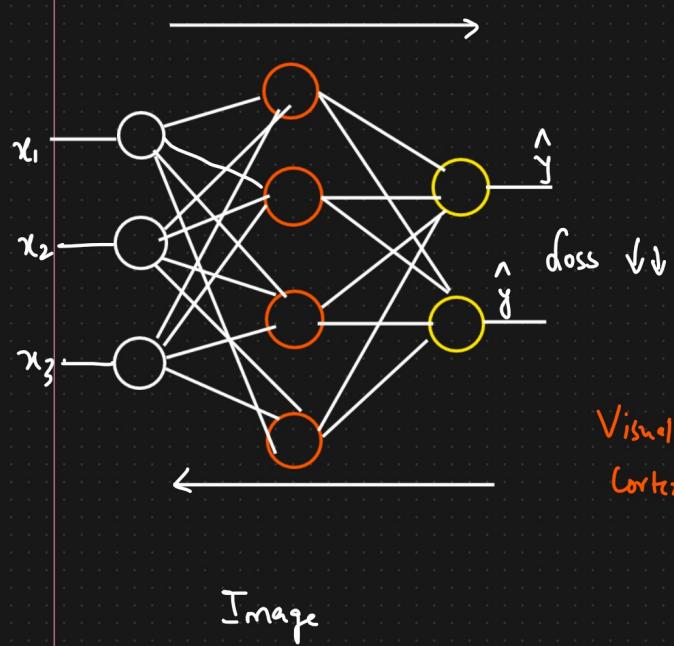
7×7

3×3

7×7

How much padding you need to apply?

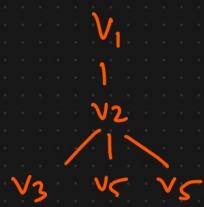
⑥ Operation of CNN vs ANN



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

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

Visual
Cortex



→ ReLU operation $\max(0, x)$

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

*

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

f_1

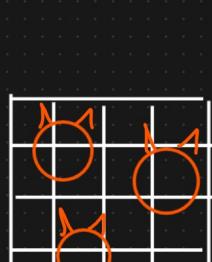
0	-4	-4	0
0	-4	-4	0
0	-4	-4	0
0	-4	-4	0

-	-	-	-
-	-	-	-
-	-	-	-
-	-	-	-

f_2
 f_3
⋮
 f_n

Convolution Layer

⑦ Max Pooling, Min Pooling, Mean Pooling



$P=1$ filter
 $S=1$

Convolution Layer

1	2	3
4	3	6
2	8	4

3×3 $\text{Slide}=2$

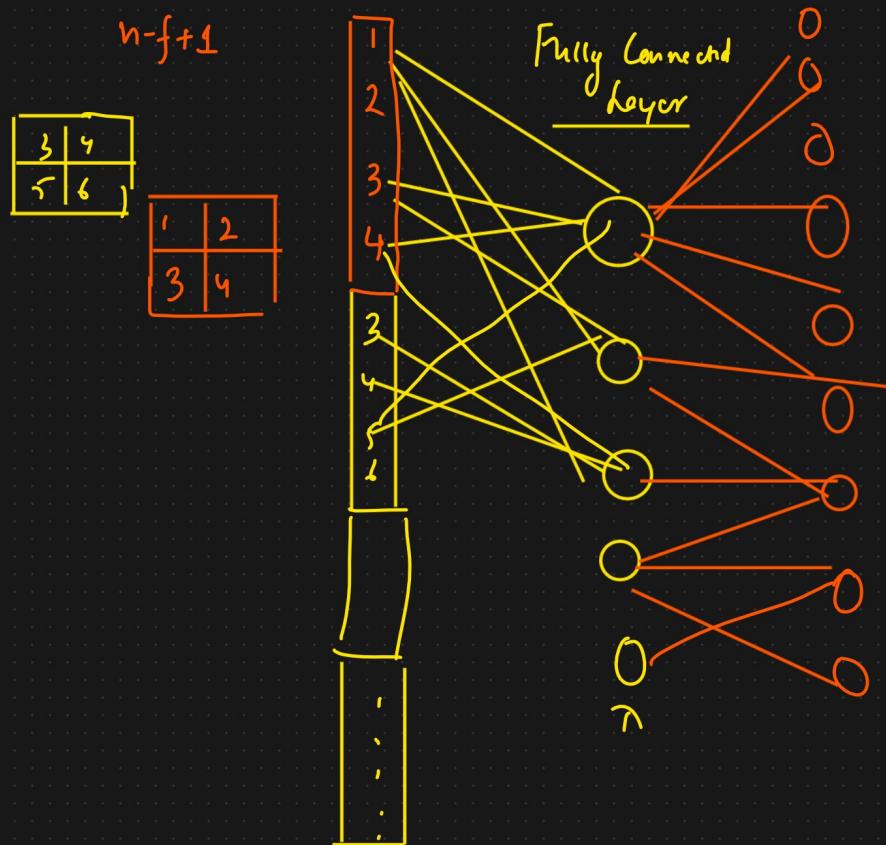
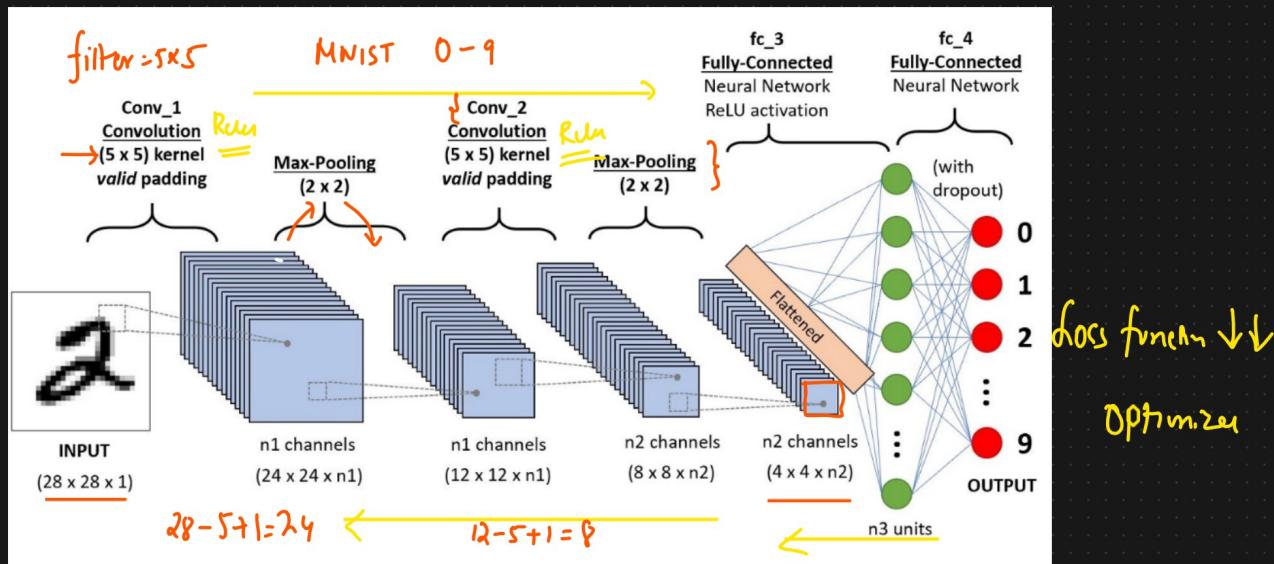
1	

4	6
8	4

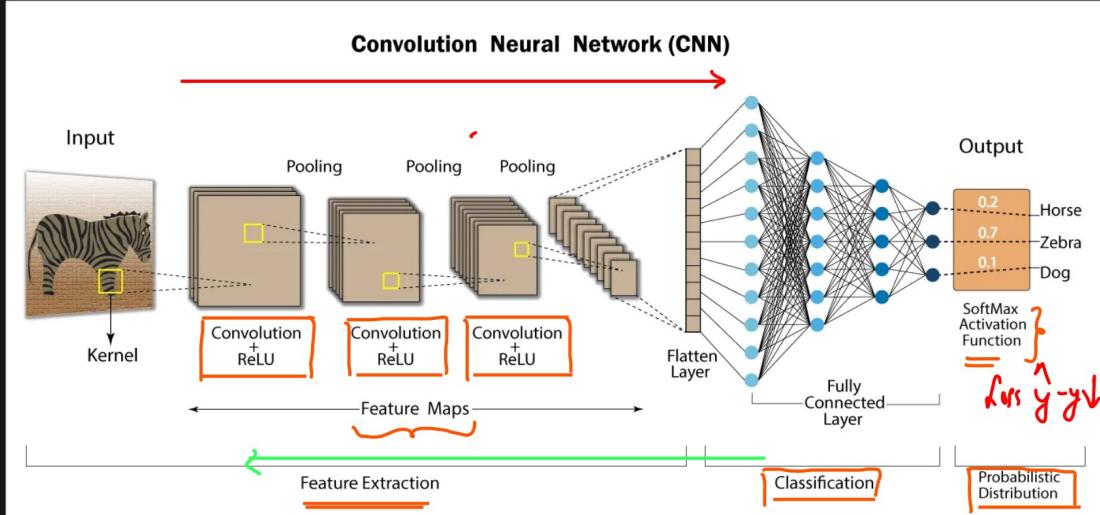
Max Pooling

⑦ Location Invariant

⑧ Fully Connected layer In CNN [Flattened Layer]



⑨ 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	
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

Day 6 - NLP { Recurrent Neural N/w }

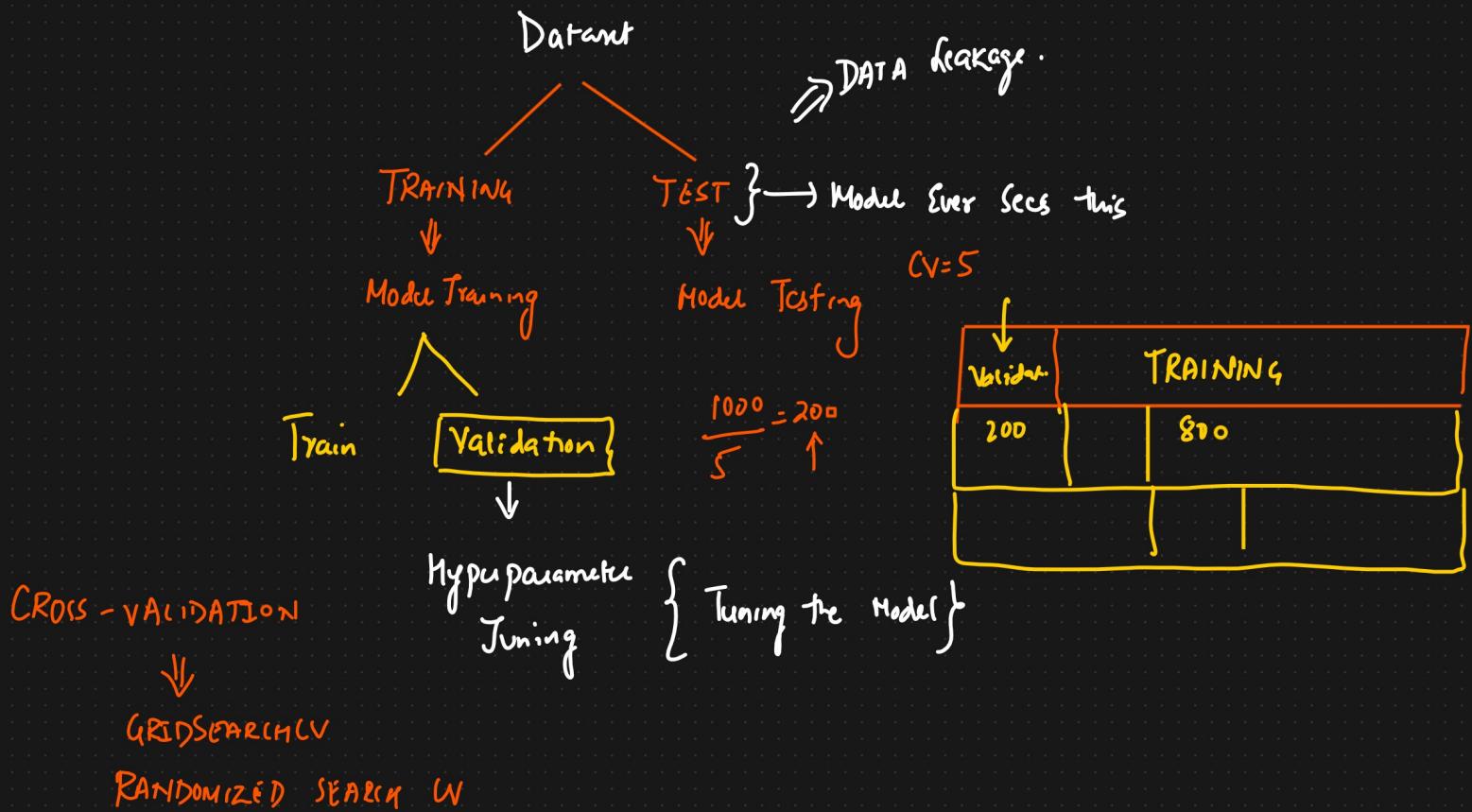
① BOW, TFIDF, Word2vec, Avg Word2Vec { Python Practical Implementation }



② Deep learning → ① RNN ② LSTM RNN ③ GRU RNN ④ Bidirectional LSTM RNN
⑤ Encoders - Decoders ⑥ TRANSFORMERS ⑦ BERT

Interview Question

① Train vs Test Vs Validation →



② Why RandomForest instead of Decision Tree { Answer }.

DT → { Low Bias High Variance }

{ Low Bias Low Variance } \rightarrow Random Forest }.

① Recurrent Neural N/w \Rightarrow Text \rightarrow Vectors

Machine Learning

Word

Embedding

\leftarrow Word2Vec, AvgWord2Vec



Chatbot

: [Question]

and

[Answer]



[Sequence of words]

Language Translation

\rightarrow [Hindi] \longrightarrow [English]

Grammaticality

Text generation

\rightarrow A Sentence $\xrightarrow{\text{Suggestion}}$ Completion of

Sentences

GMAIL

AutoSuggestion

Word2Vec

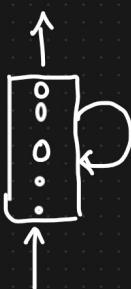
\Downarrow Deep learning

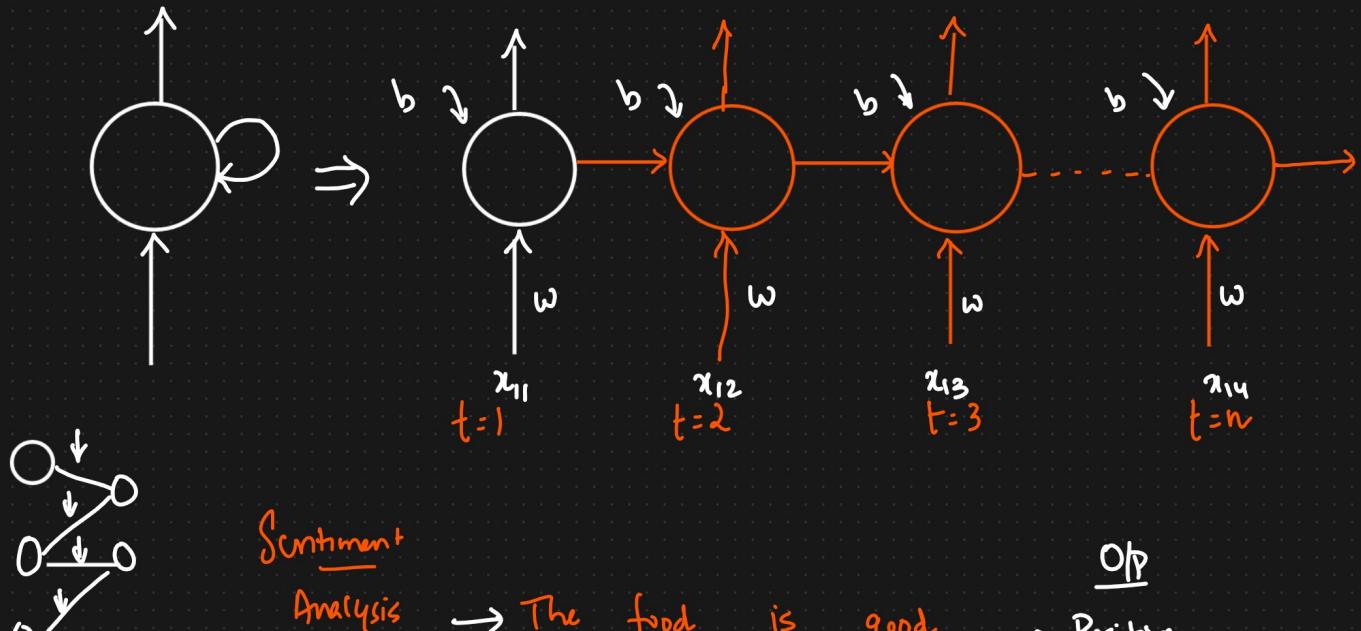
① RNN ② LSTM RNN ③ Transformers ④ BERT



[Words \rightarrow Vec]

① Recurrent Neural N/w





Sentiment

Analysis

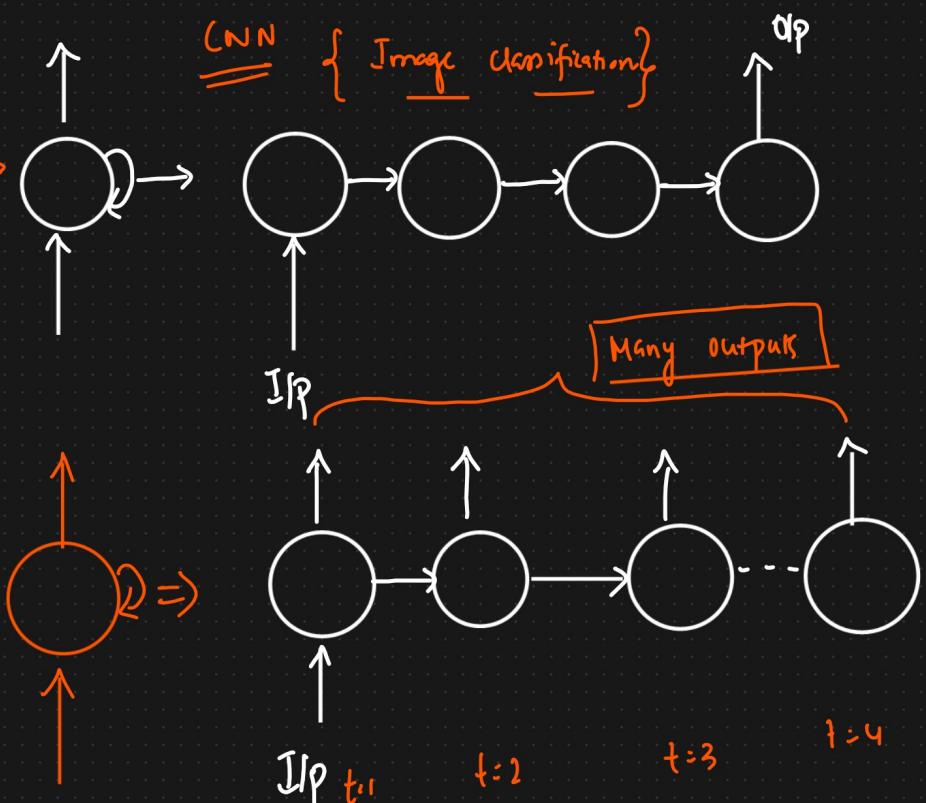
The food is good \rightarrow Positive

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

$x_{11} \rightarrow \text{Word2Vec} \rightarrow \text{Vectors} \rightarrow d = 300$

Types of RNN

- ① One to One RNN
- ② One to Many RNN
- ③ Many to One RNN
- ④ Many to Many RNN

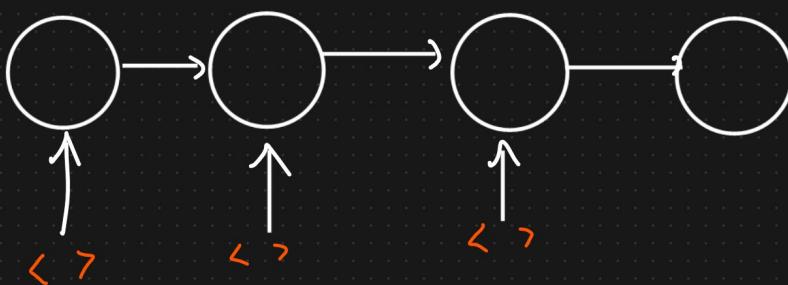


$\left\{ \begin{array}{l} \text{Eq: Music Generation, Text Generation} \\ \text{of} \end{array} \right.$

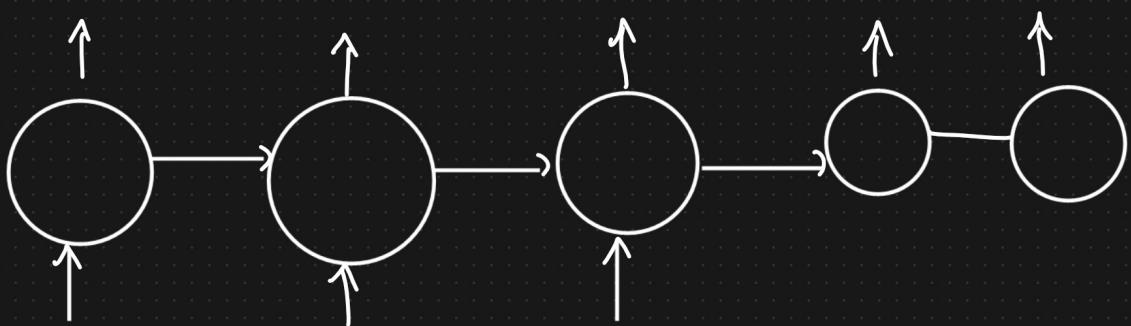
Google Search Suggestion, Movie Recommendation

③ Many to One

Sentiment Analysis ↗ O/P { Predict Next Day }
Sales



④ Many to Many



① Language Translation

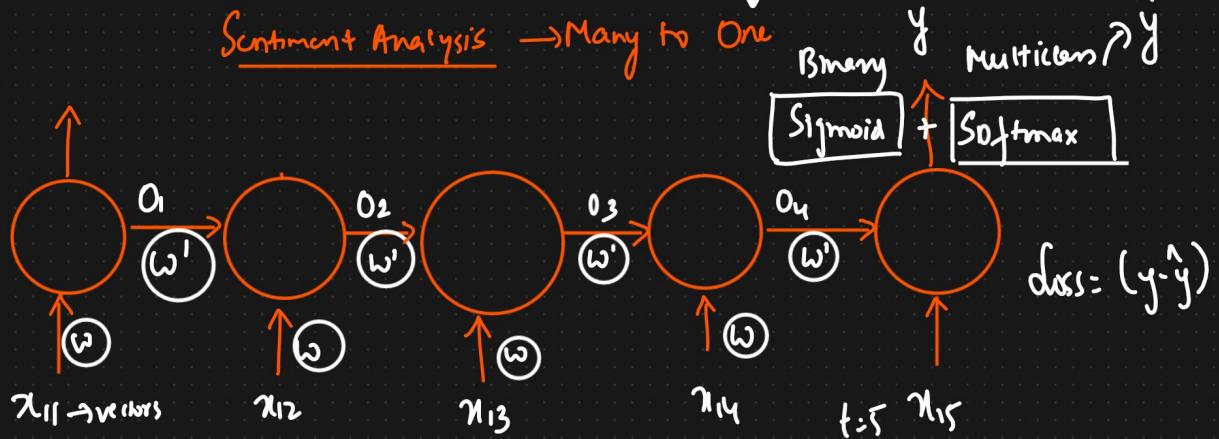
② Question answers

③ Chatbots.

{ ① Forward Propagation
② Backward Propagation }

① Forward Propagation In RNN

Time Series Data



$t=1$ $t=2$ $t=3$ $t=4$ $O_1 = f(x_{11} * \omega)$

The food is very good Positive.

x_{11} x_{12} x_{13} x_{14} x_{15}

O/P

$O_2 = f(x_{12} * \omega + O_1 * \omega_i)$

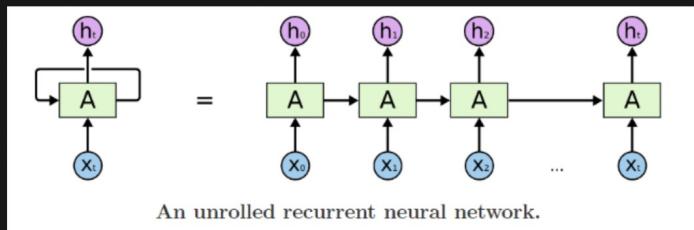
$O_3 = f(x_{13} * \omega + O_2 * \omega')$

LSTM RNN [Long Short Term Memory RNN]

RNN → Long Term Dependencies → Vanishing Gradient Problem

- ① RNN → Problem? ✓
- ② Why LSTM RNN? ✓ Basic Representation.
- ③ How LSTM RNN works
 
- ④ LSTM Architecture
- ⑤ Working of LSTM RNN

Problems With RNN → Long Term Dependency

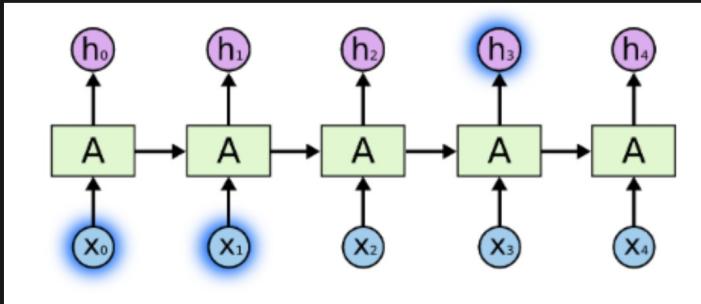


Vanishing Gradient Problem

Task:

Next Word In a Sentence

The color of the Sky is blue
— further context



Gap is 1cm

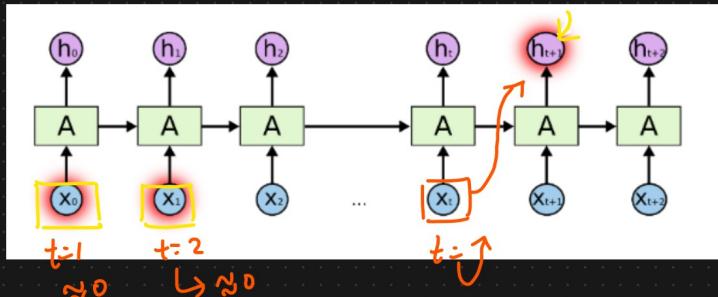
Huge gap O/P ← Context

I grew up in India ... I speak
fluent language ← Context

Name of language



further context



10-0.2r)

0-1)

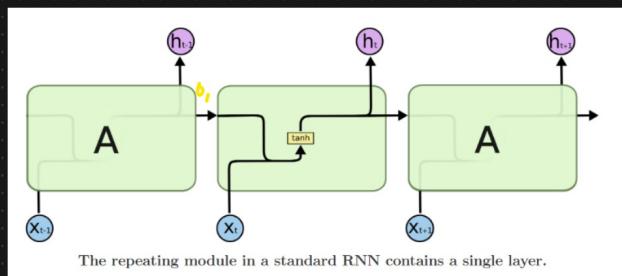
Huge Gap → Long Term Dependency

RNN → Long Term Dependency → Vanishing Gradient Problem

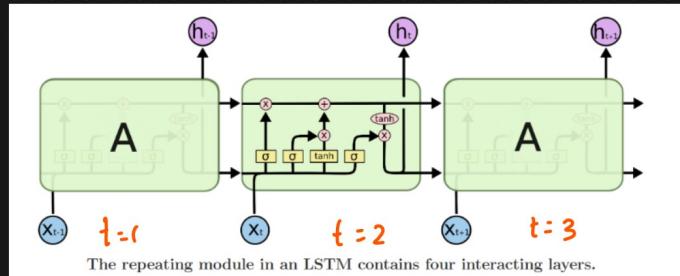
Chain Rule → ≈ 0 .

Basic Representation of RNN And LSTM RNN

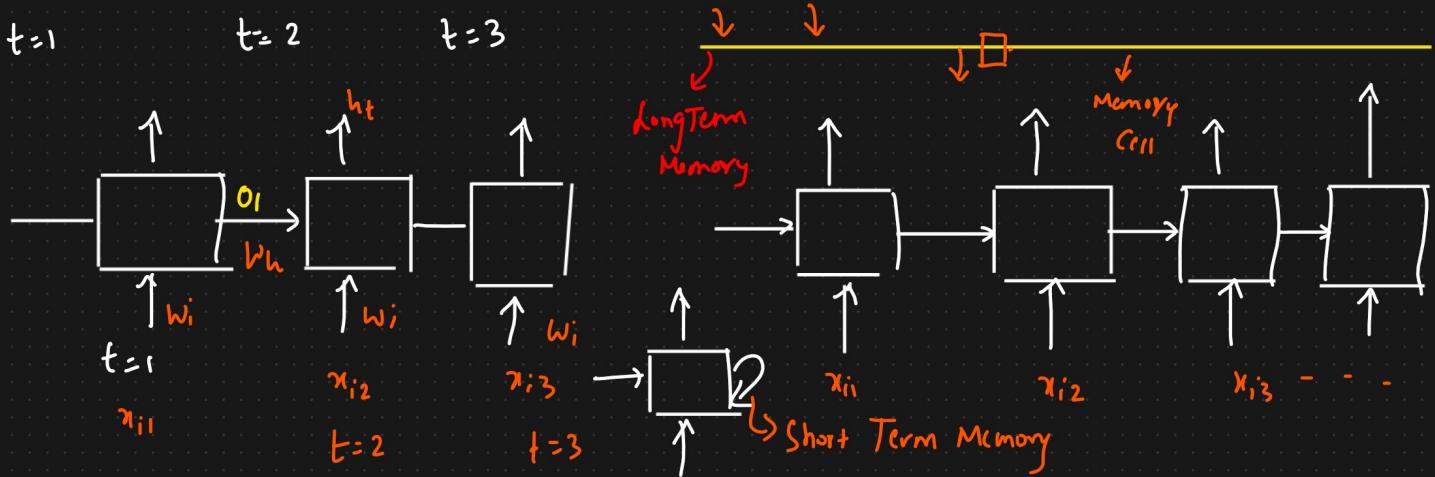
LSTM RNN



The repeating module in a standard RNN contains a single layer.



The repeating module in an LSTM contains four interacting layers.



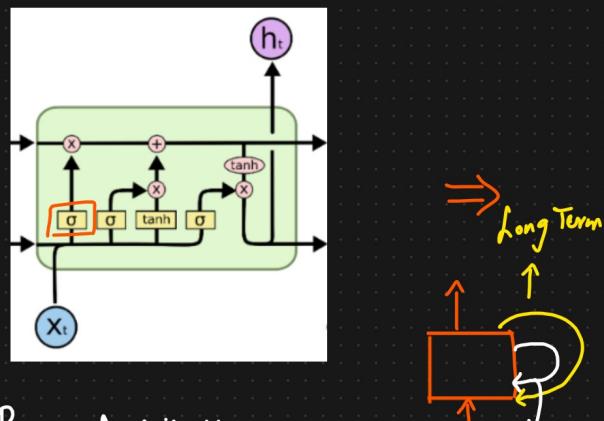
LSTM RNN → Long Term Memory
 → Short Term Memory

Convoyana But : fuggages

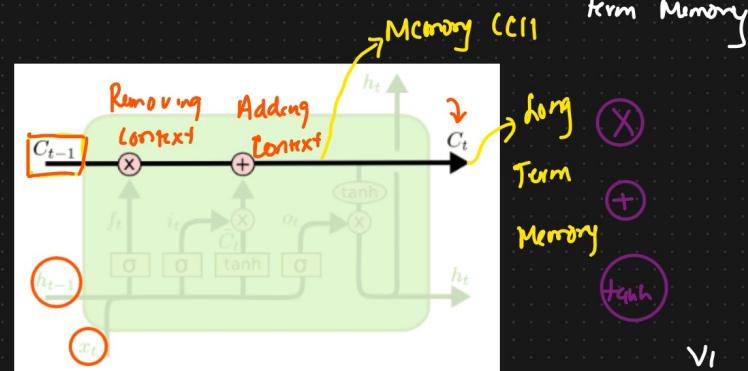


Human

LSTM Architecture



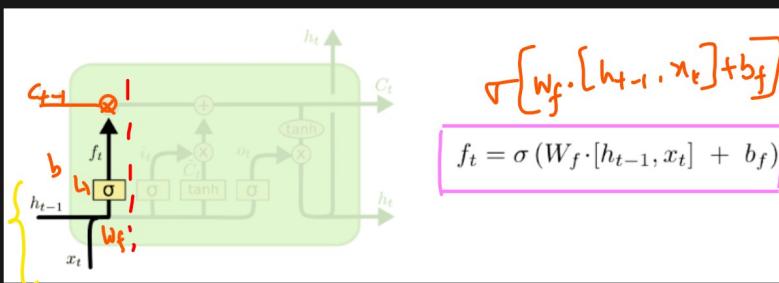
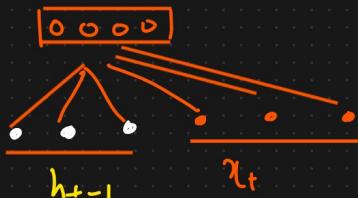
Basic Architecture



Combining 2 vectors

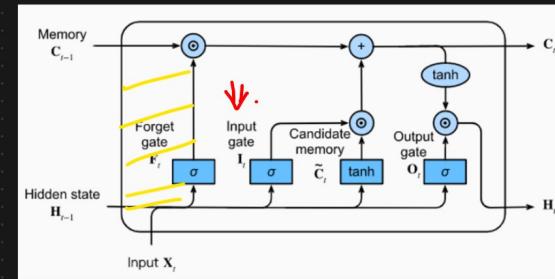
$$h_{t-1} = [1 \ 2 \ 3]$$

$$x_t = [2 \ 3 \ 4]$$



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

$$\sqrt{w_f \cdot [h_{t-1}, x_t] + b_f}$$



Forget Gate

Text Next Word
 $x_1 \ x_2 \ x_3 \ x_4 \quad y_5$



h_{t-1} = Hidden state of previous time stamp
 x_t = Word passed as i/p in the current time stamp

$$x_t = [0 \ 2 \ 4 \ 1]$$

$$h_{t-1} = [1 \ 2 \ 4]$$

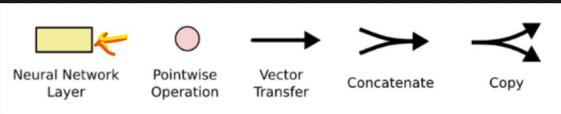
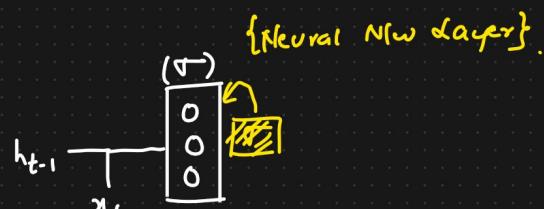
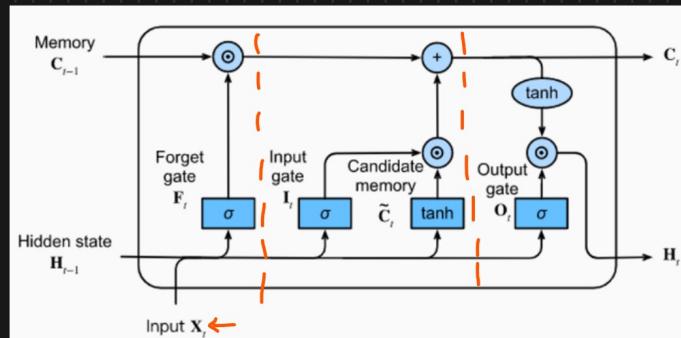
$$c_{t-1} = [3d]$$

$$y_5 = [4 \ 5 \ 1 \ 2]$$

$$x_{t+1} = [5 \ 7 \ 9]$$

$$h_t = [tanh(1) \ tanh(2) \ tanh(3)]$$

LSTM RNN



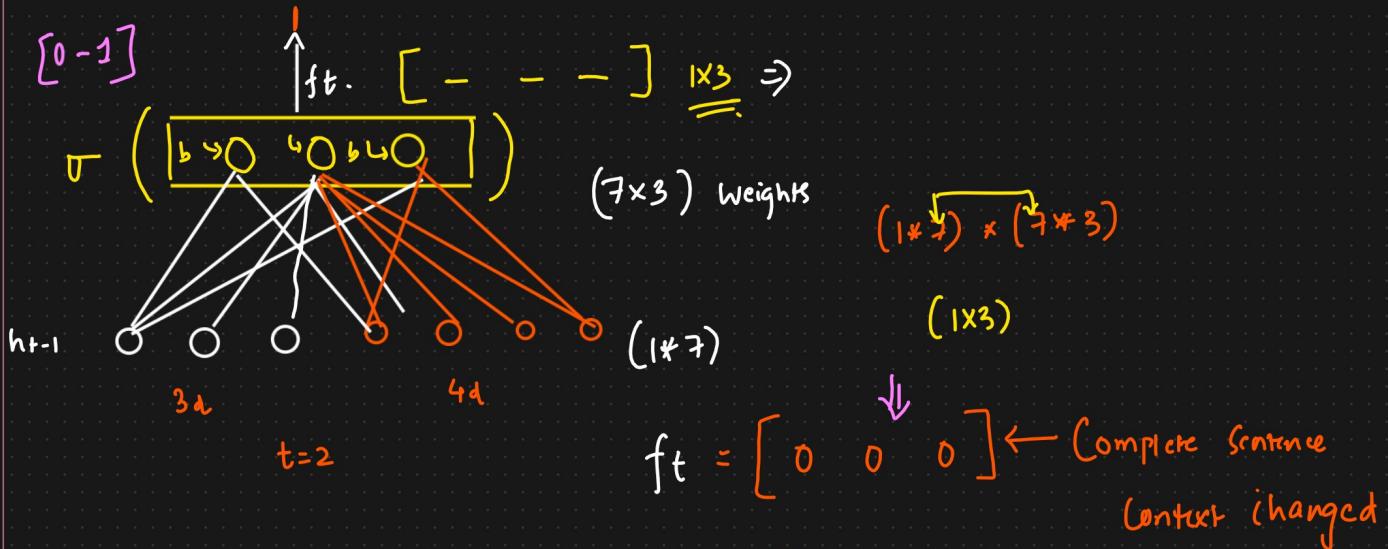
$$v_1 = [1 \ 2 \ 3] \leftarrow \text{tanh}$$

$$v_2 = [4 \ 5 \ 6] \leftarrow \text{tanh}$$

$$\times = [4 \ 10 \ 18]$$

$$+ = [5 \ 7 \ 9]$$

$$\text{tanh} = [\tanh(1) \ \tanh(2) \ \tanh(3)]$$



$$\textcircled{1} C_{t-1} = \begin{bmatrix} 6 & 8 & 9 \end{bmatrix} \otimes \begin{bmatrix} 0 & 0 & 0 \end{bmatrix}$$

$= \begin{bmatrix} 0 & 0 & 0 \end{bmatrix} \leftarrow \text{Removing all the previous context}$

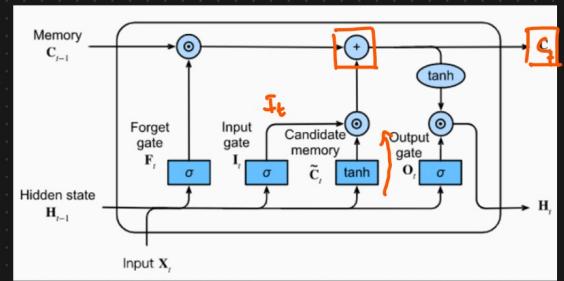
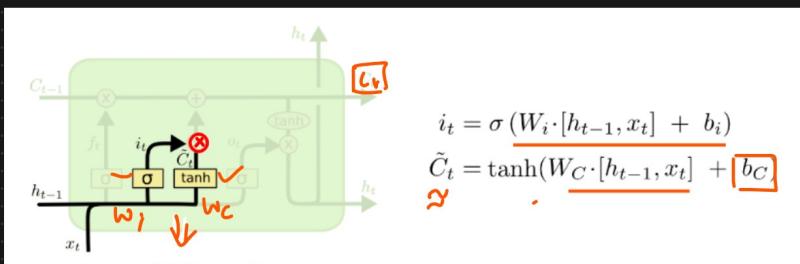
$$ft = [1 \ 1 \ 1]$$

$$\textcircled{2} C_{t-1} = \begin{bmatrix} 6 & 8 & 9 \end{bmatrix} \otimes \begin{bmatrix} 1 & 1 & 1 \end{bmatrix} = \begin{bmatrix} 6 & 8 & 9 \end{bmatrix}$$

$$\textcircled{3} C_{t-1} = \begin{bmatrix} 6 & 8 & 9 \\ \hline = & = & = \end{bmatrix} \otimes \begin{bmatrix} 0.5 & 1 & 0.5 \end{bmatrix} = \begin{bmatrix} 3 & 8 & 4.5 \end{bmatrix}$$

Conclusion : Based on the context \rightarrow Forget gate will let go some information or will not let go some info {Forgetting}.

② Input Gate And Candidate Memory



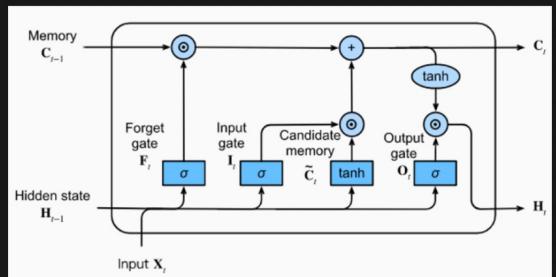
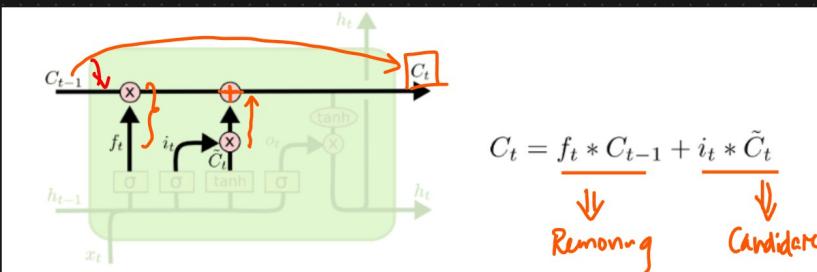
Adding Info

$$I_t = [0 \ 1 \ 0] \xrightarrow{\oplus} [0 \ 2 \ 0] \Rightarrow \text{Input Gate}$$

$$b \rightarrow [0 \ 0 \ 0] \xrightarrow{W_I} [0 \ 8 \ 0]$$

Diagram illustrating the computation of the input gate I_t from bias b and hidden state h_{t-1} using weight matrix W_I .

Context = If any information needed to be added it the memory
 $c_{t-1} \rightarrow$ The information will be added



I stay in India - - - - -
 and I speak English Hindi

or
 Forgetting
 Some info

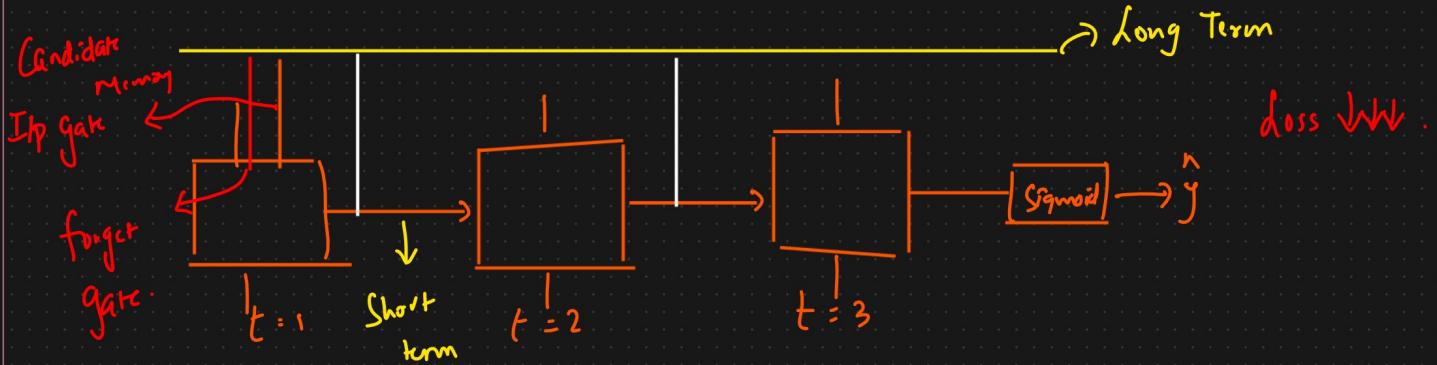
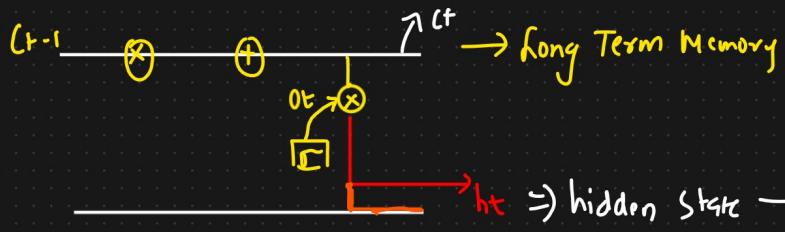
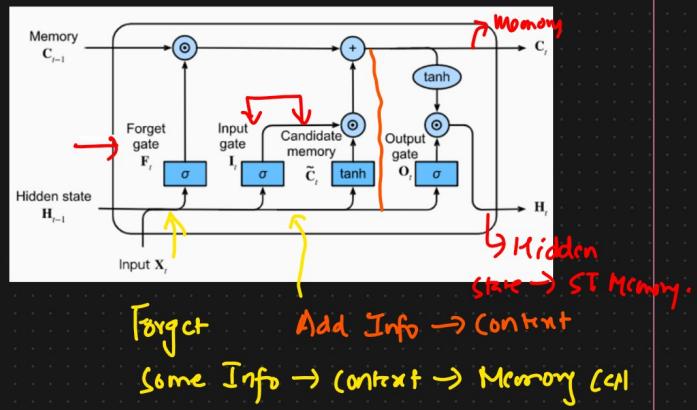
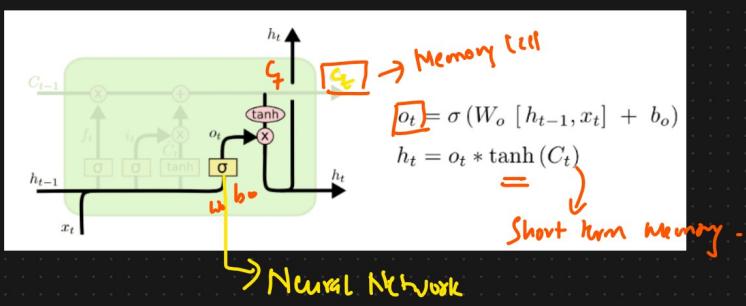
Memory -

Forget Gate

I/P Gate \otimes Candidate memory

+
 $C_{t-1} \Rightarrow C_t$

Output gate LSTM RNN



$[W_i, W_c, W_o]$ \rightarrow Updating \leftarrow Back Propagation

GRU RNN \Rightarrow LSTM Variant

Training Data With LSTM RNN

{Training} Text Paragraph

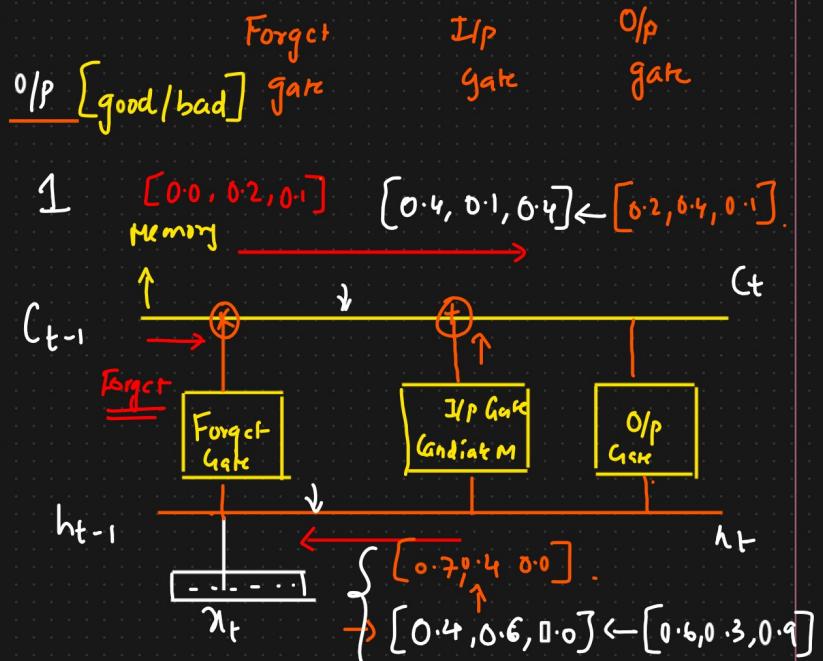
I Went to Restaurant and order burger

The burger looked tasty and crispy

→ But burger is not good for health

→ It has lot of fats, cholesterol

→ But this burger was made with Whey protein and only vegetables were used, so it was good



Word → Vectors → Embedding Layer

Word2Vec [3 dimension - vector]

→ [[Good] [Bad] [Healthy]] ← Black Box

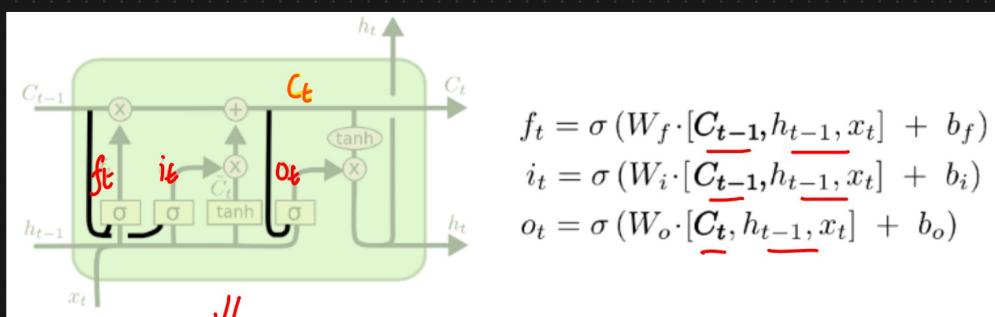
Tasty [0.9 0.0 0.1] ← 3 d.

Variants of LSTM RNN

LSTM Variants Introduced By Gers & Schmidhuber [2000]

LSTM RNN [1970-80]

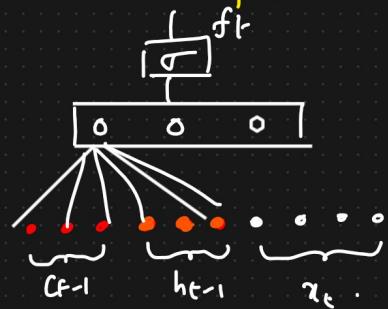
↳ Research paper



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

$$i_t = \sigma(W_i \cdot [C_{t-1}, h_{t-1}, x_t] + b_i)$$

$$o_t = \sigma(W_o \cdot [C_t, h_{t-1}, x_t] + b_o)$$



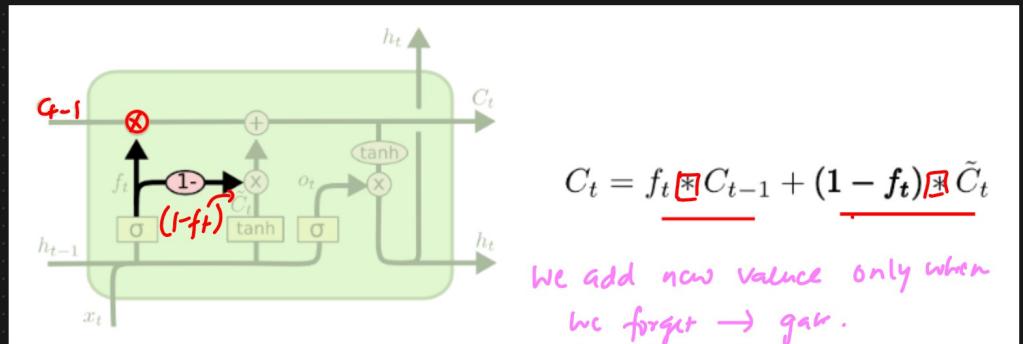
Connections → From memory cell to forget gate

i/p gate ⇒ Peephole connections

O/p gate

Peephole Connections: We let the gate layers look at the cell state

Another variation \rightarrow Coupling Forget And I/p Gates



Goal: We only forget

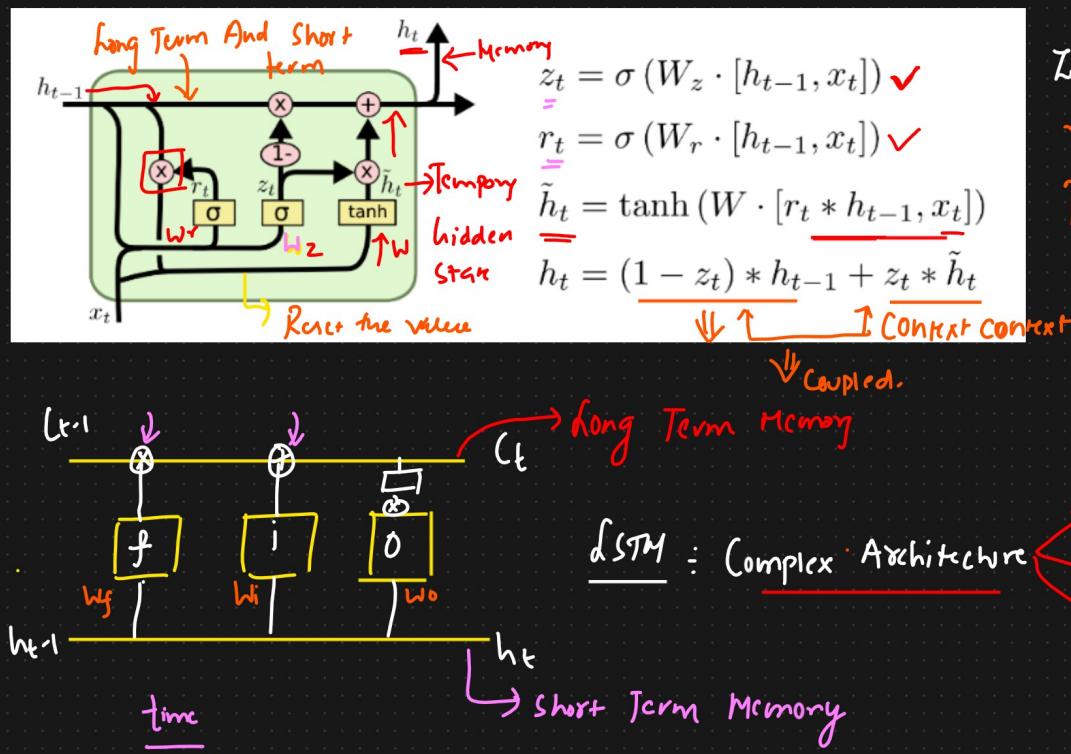
when we're going to i/p something in its place.

We only i/p new values to the state when we forget something older.

Instead of separately deciding what to forget and what we should add new Info, we make this decision together.

GRU \rightarrow Gated Recurrent Unit [Cho, et al [2014]]

1980 \rightarrow LSTM
2000 - variants
2014 \rightarrow GRU



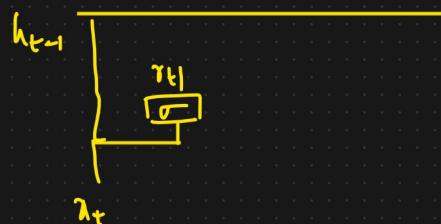
$z_t \Rightarrow$ Update Gate \leftarrow
 $r_t \Rightarrow$ Reset Gate \leftarrow
 $\tilde{h}_t \Rightarrow$ Temporary hidden state.

Trainable parameters
[w_f, w_i, w_o]

forget
i/p + candidate memory
O/p

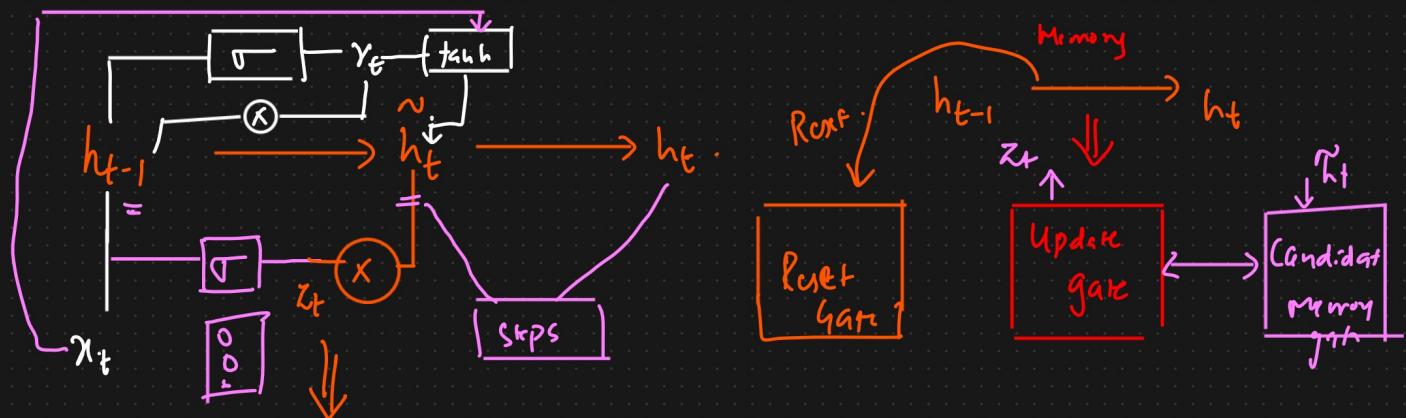
Training Time $\uparrow \uparrow$

$$\text{Reset gate} = r_t =$$



\rightarrow Resetting some info from \$h_{t-1} \Rightarrow\$ Memory \rightarrow LTM + STM

$$\begin{aligned}
 h_{t-1} &= [0.6 \quad 0.5 \quad 0.3 \quad 0.9] \\
 r_t &\leftarrow [0.2 \quad 0.4 \quad 0.8 \quad 0.2] \\
 \downarrow & \quad \downarrow \quad \downarrow \quad \downarrow \\
 x_t &\rightarrow [0.12 \quad 0.20 \quad 0.24 \quad 0.18] \leftarrow \text{Rescaling} \rightarrow \text{Context}
 \end{aligned}$$



What Context Info needs to be Added



Candidate hidden State · [Current Context]

↓
Imp → Add Info

$\tilde{h}_t \Rightarrow$ New Info

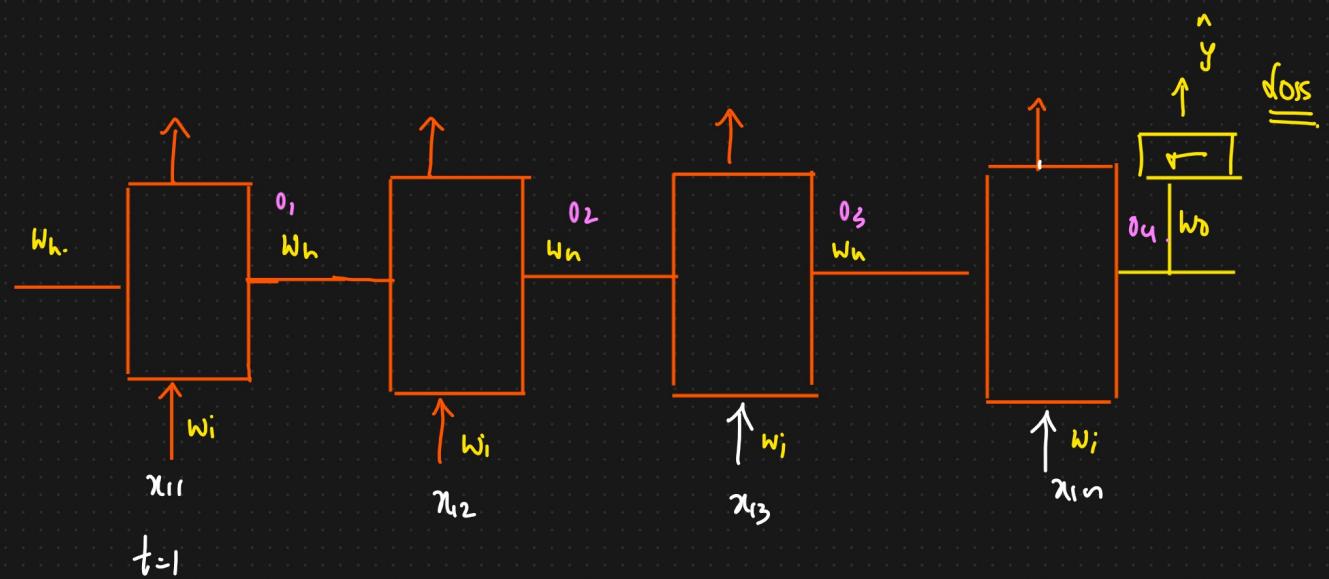
—

Bidirectional RNN

① Simple RNN → practical Implement → Embedding layers

② LSTM, GRU Variants RNN → Practical Implement

③ Bidirectional RNN / LSTM RNN



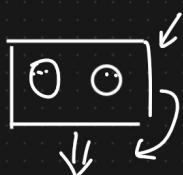
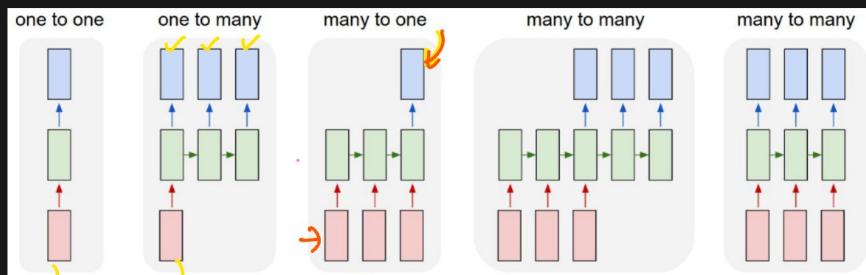
Types of RNN

① One to Many RNN

② Many to One RNN

③ Many to Many RNN

④ One to One RNN



Eg:

There is a dog and cat

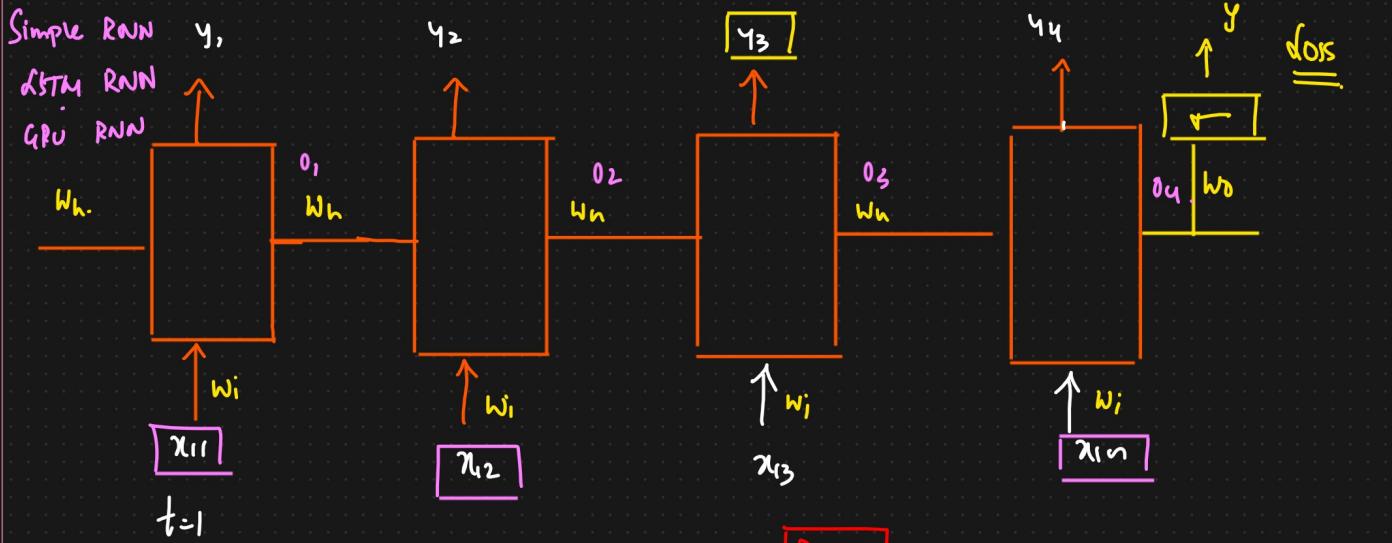
One i/p
One o/p

Eg: Image
Captioning

Image search
I/P Many
O/P One.
A cat eating
food
o catting

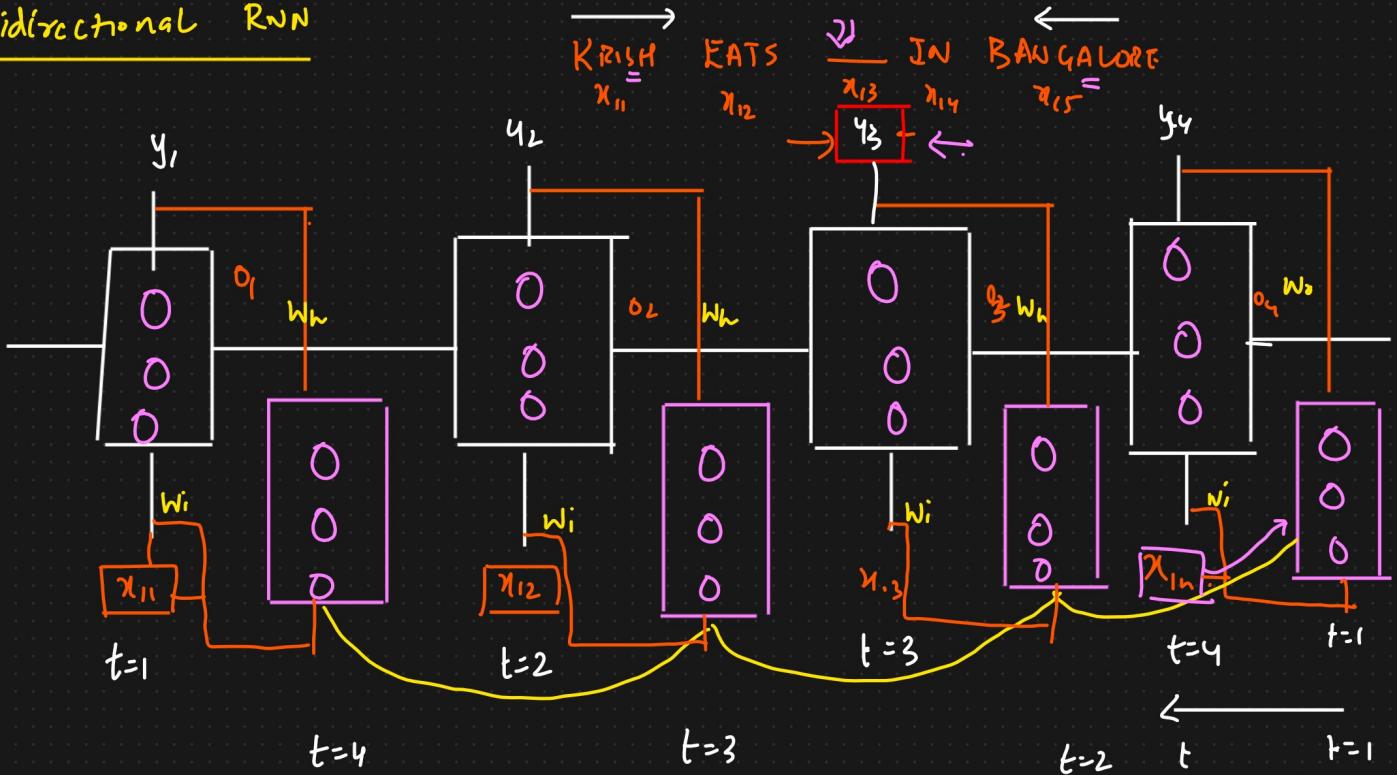
Many I/P
Many O/P.

Language Translation



Text Example = KRISH EATS IN BANGALORE
 KRISH EATS $\xrightarrow{\text{DOSA}}$ IN $\xrightarrow{\text{DOSA}}$ BANGALORE
 PIZZA $\xrightarrow{\text{DOSA}}$ PARIS

Bidirectional RNN



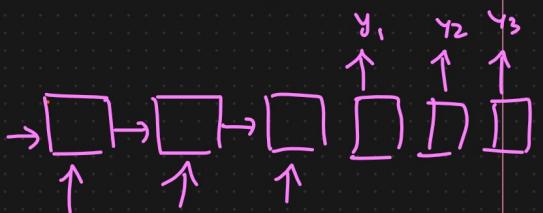
① Forward propagation \rightarrow Equation \leftarrow

Encoder And Decoder

- ① Simple RNN → Vanishing Gradient Problem
- ② LSTM RNN → ↗ Long Short Term Memory.
- ③ GRU RNN → ↗
- ④ Bidirectional RNN ←

Type of RNN

- ① Many to Many RNN



Encoder And Decoder }

Eg:- One language To Other

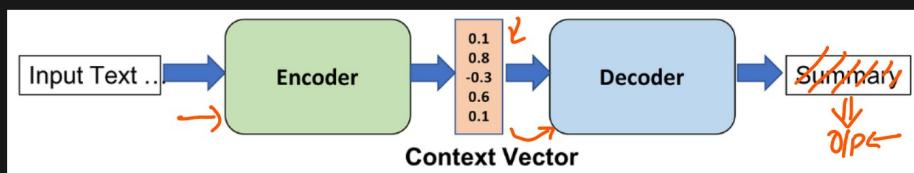
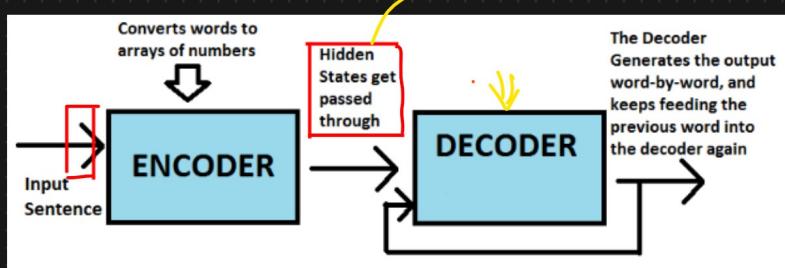
English → French

Eg:- French Chat → Hi, how are you?

Sequences I/p

O/p Sequence Of Words

Simple Working



- ① Encoder ⇒ I/p ⇒ Context Vector ⇒ Vectors
- ② Decoder ⇒ ← ⇒ O/p

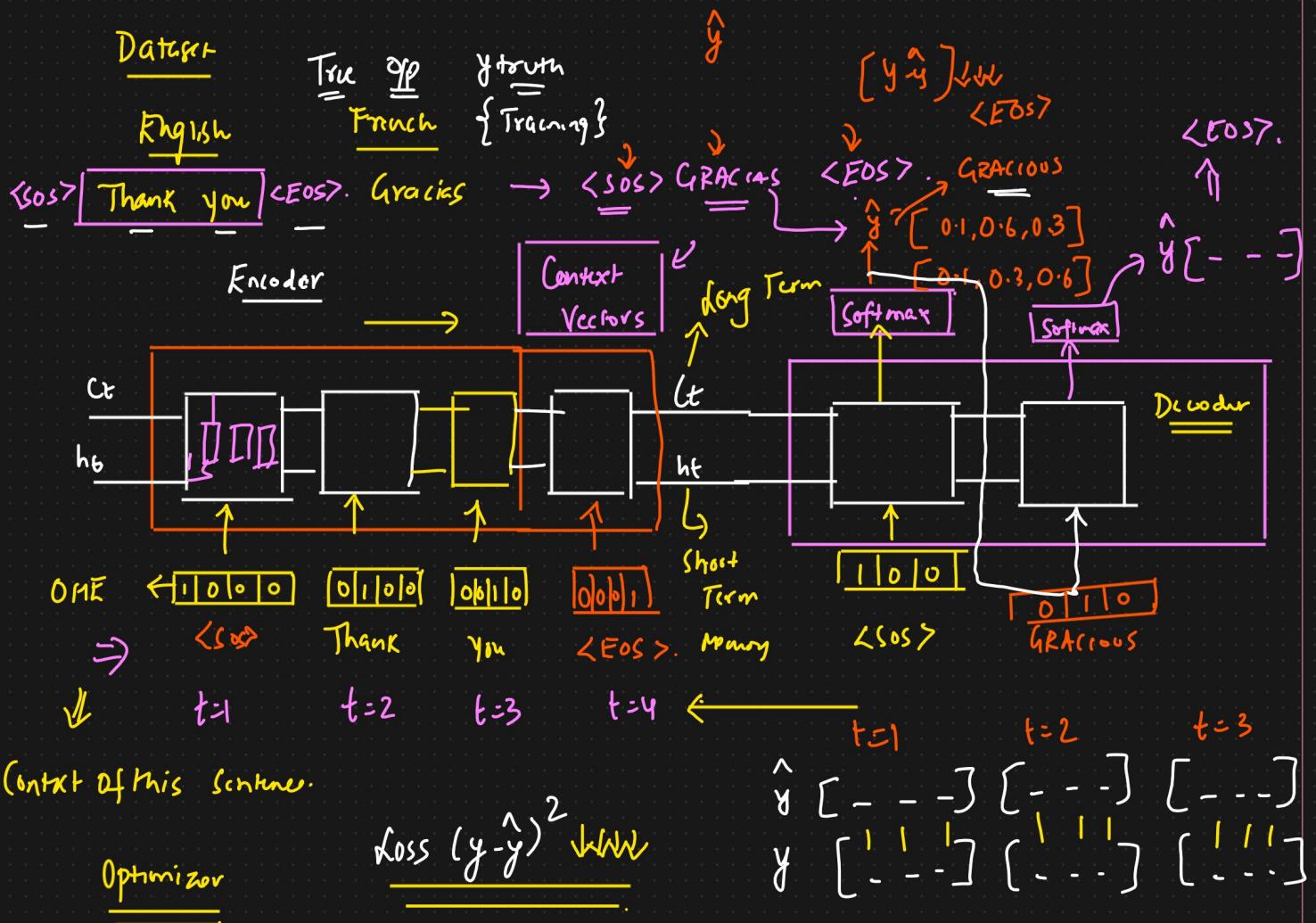
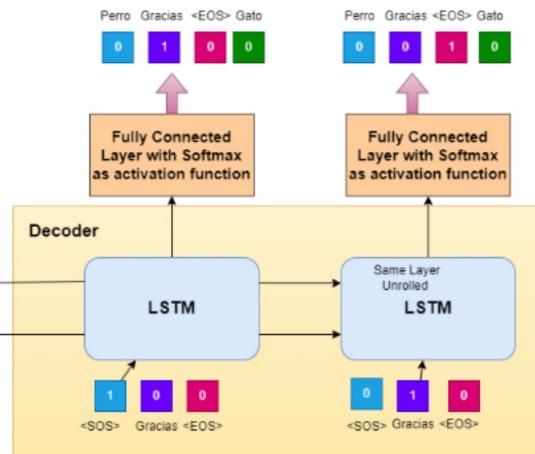
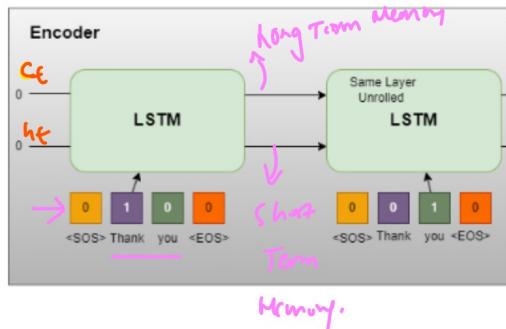
usecase

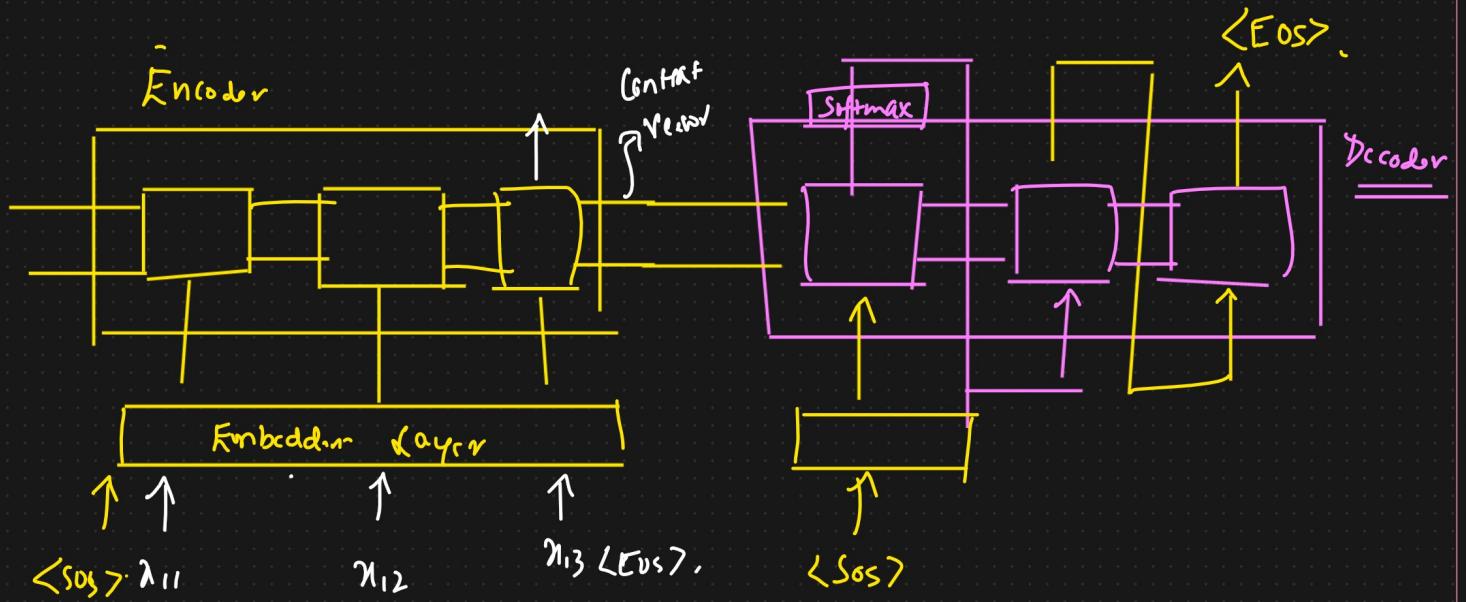
- ① Language Translation
- ② Text Generation
- ③ Text Suggestion

RNN → Vanishing Gradient Problem

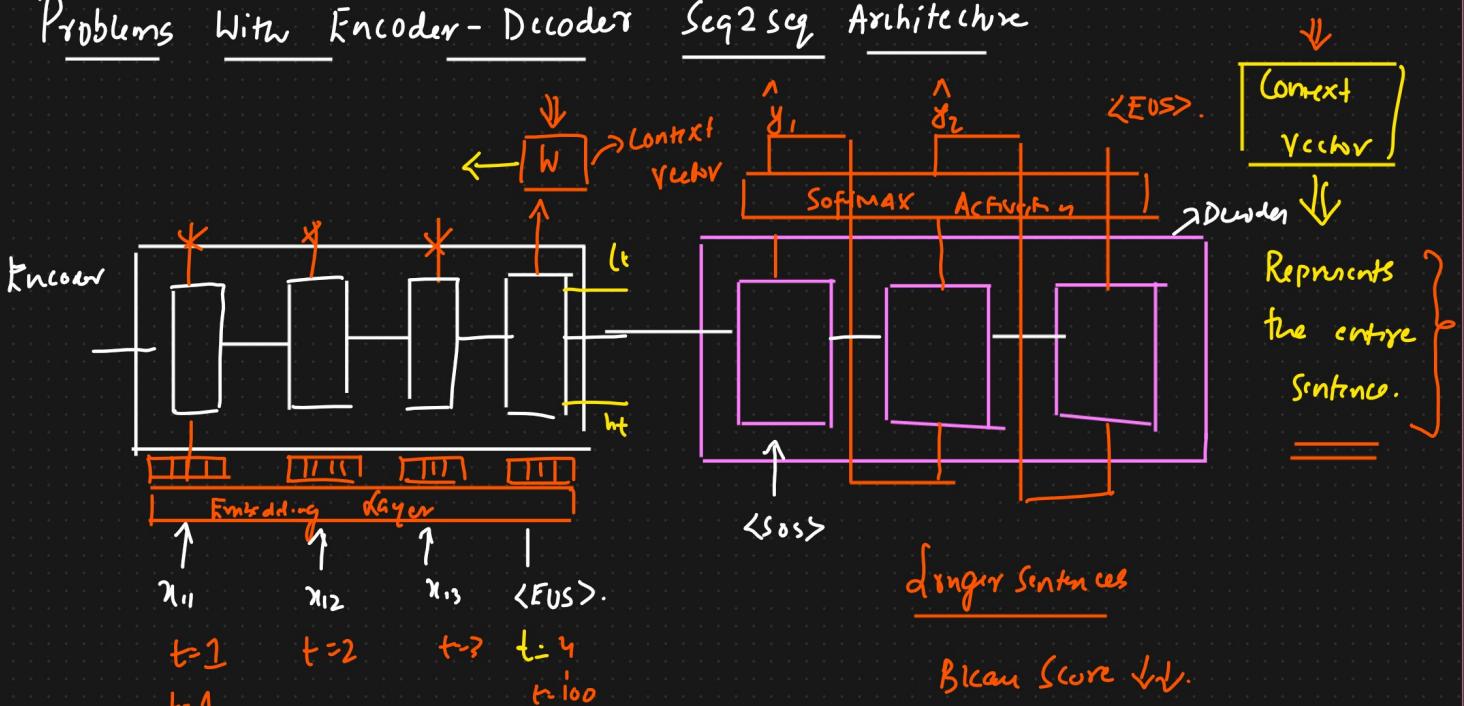
- ① forget gate
- ② I/P gate & candidate i/p
- ③ O/P

Sequence-to-Sequence (seq2seq) Encoder-Decoder Neural Network

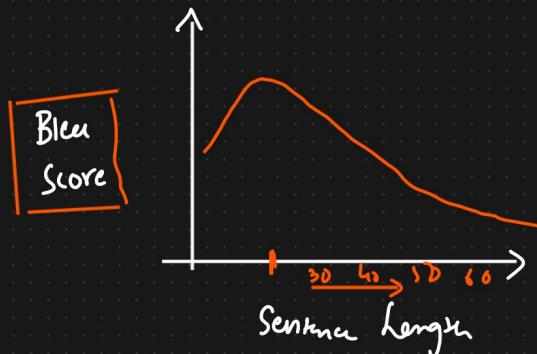




Problems With Encoder-Decoder Seq2Seq Architecture



Ruarchers ÷ Sentences of varying length



\Rightarrow Seq to Seq Data

* Attention Mechanism \rightarrow Seq2Seq Network

longer paragraph $\rightarrow \{ \text{Context Vector} \}$.

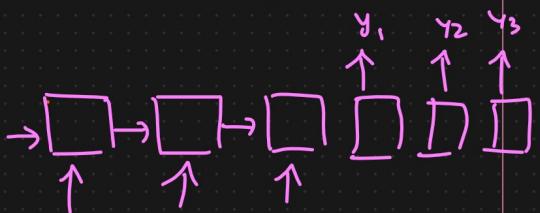
$$+ \\ \{ \text{Context} \}$$

Encoder And Decoder

- ① Simple RNN → Vanishing Gradient Problem
- ② LSTM RNN → ↗ Long Short Term Memory.
- ③ GRU RNN → ↗
- ④ Bidirectional RNN ←

Type of RNN

- ① Many to Many RNN



Encoder And Decoder }

Eg:- One language To Other

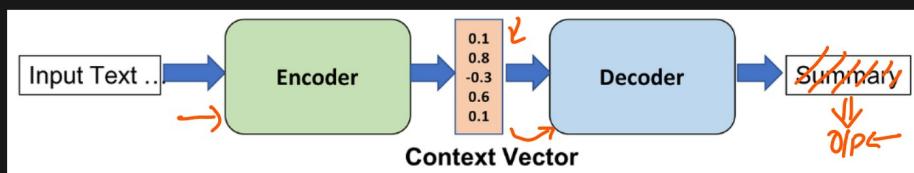
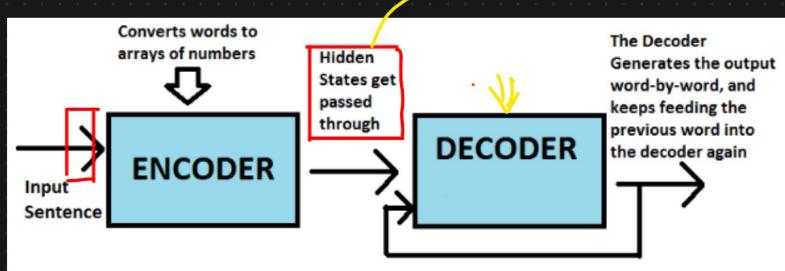
English → French

Eg:- French Chat → Hi, how are you?

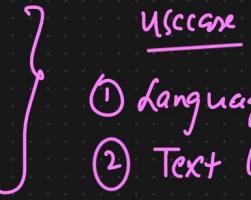
Sequences I/p

O/p Sequence Of Words

Simple Working



- ① Encoder ⇒ I/p ⇒ Context Vector ⇒ Vectors
- ② Decoder ⇒ ⇒ O/p

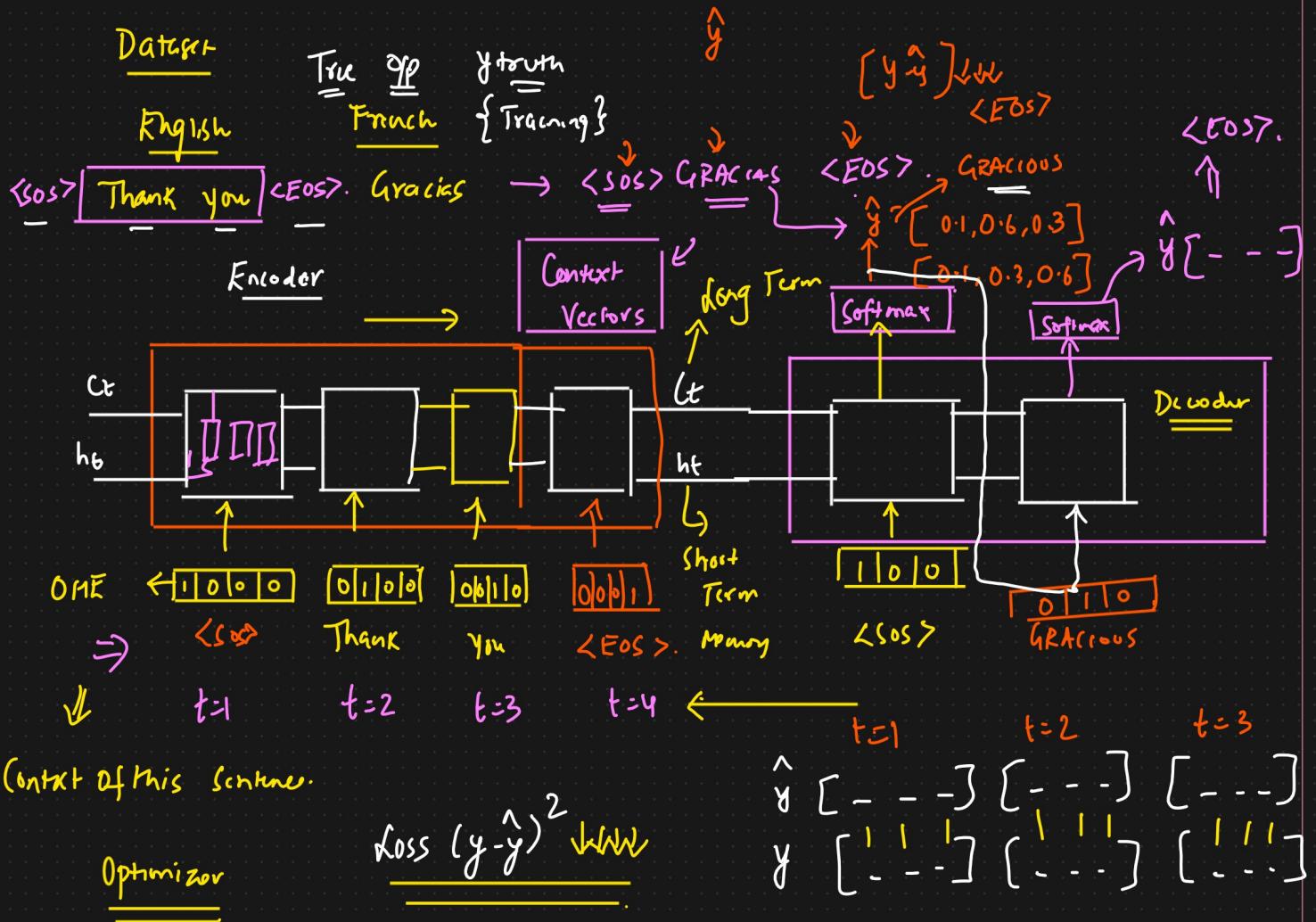
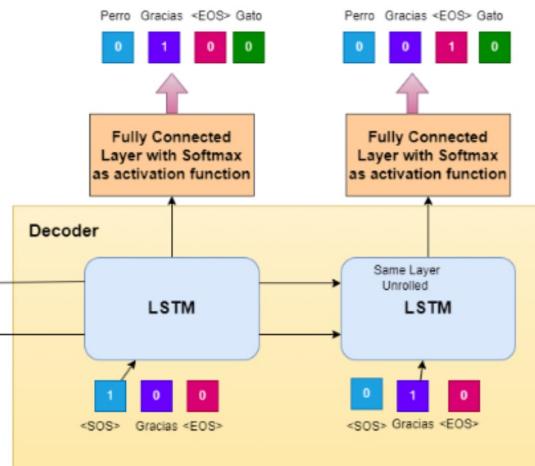
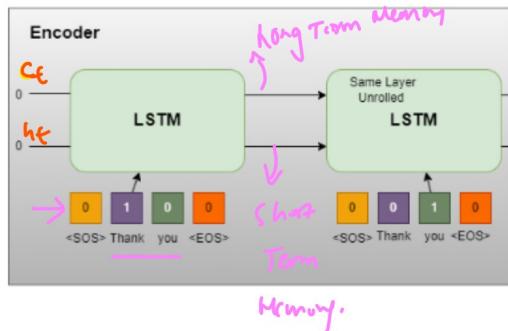


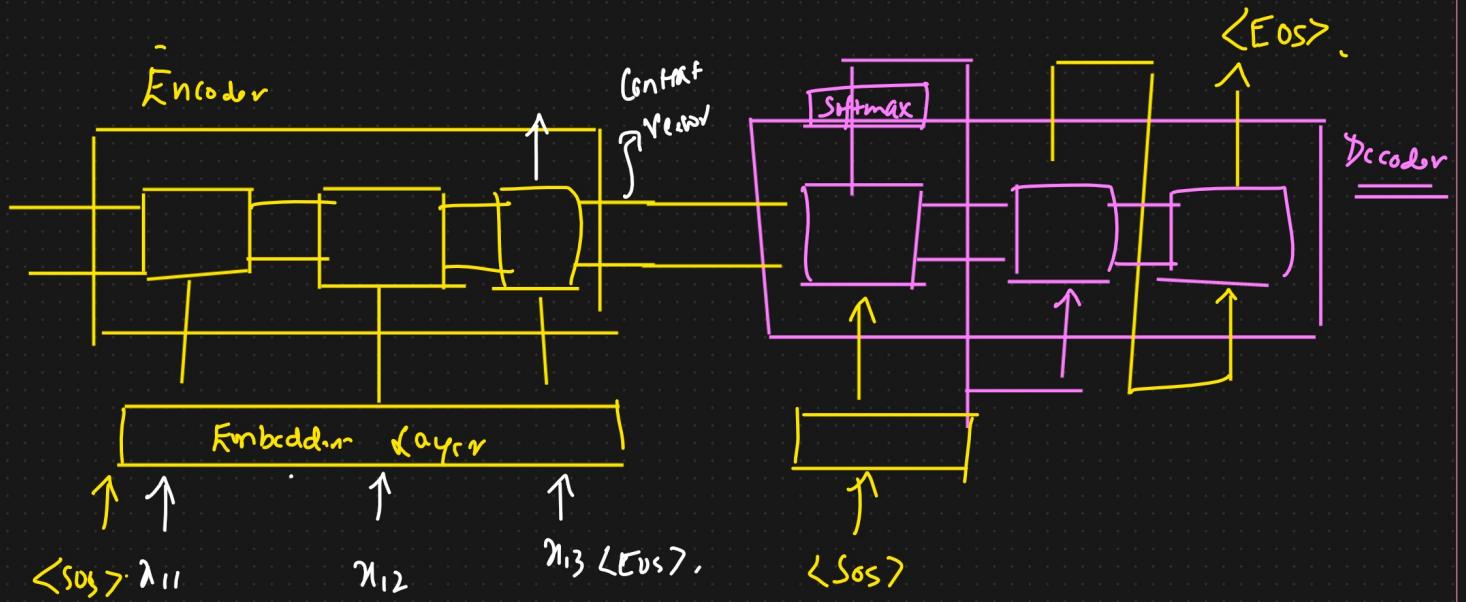
- ① Language Translation
- ② Text Generation
- ③ Text Suggestion

RNN → Vanishing Gradient Problem

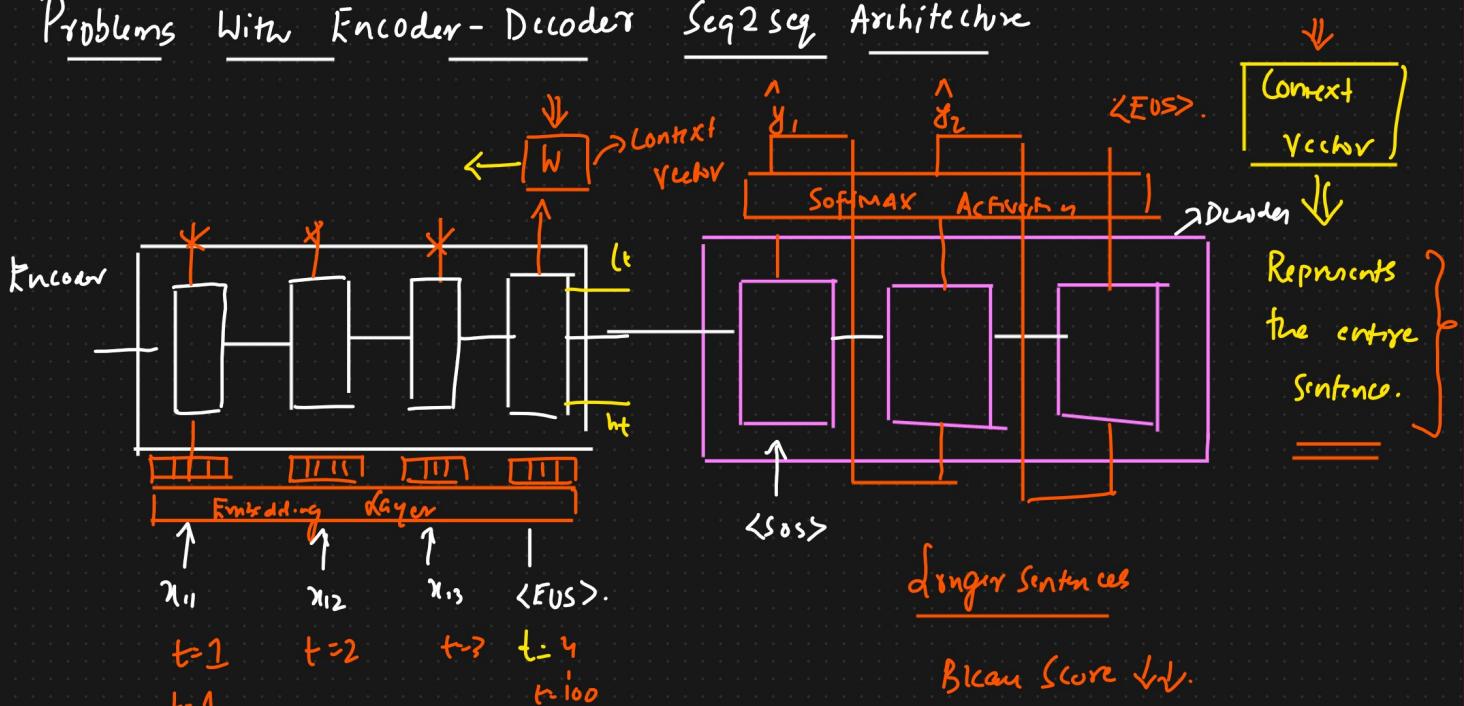
- ① forget gate
- ② I/P gate & candidate i/p
- ③ O/P

Sequence-to-Sequence (seq2seq) Encoder-Decoder Neural Network

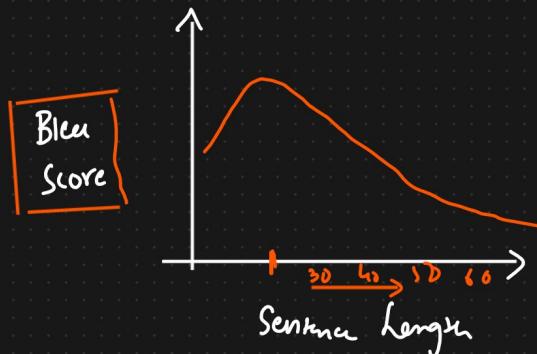




Problems With Encoder-Decoder Seq2Seq Architecture



RuArchers : Sentences of varying length



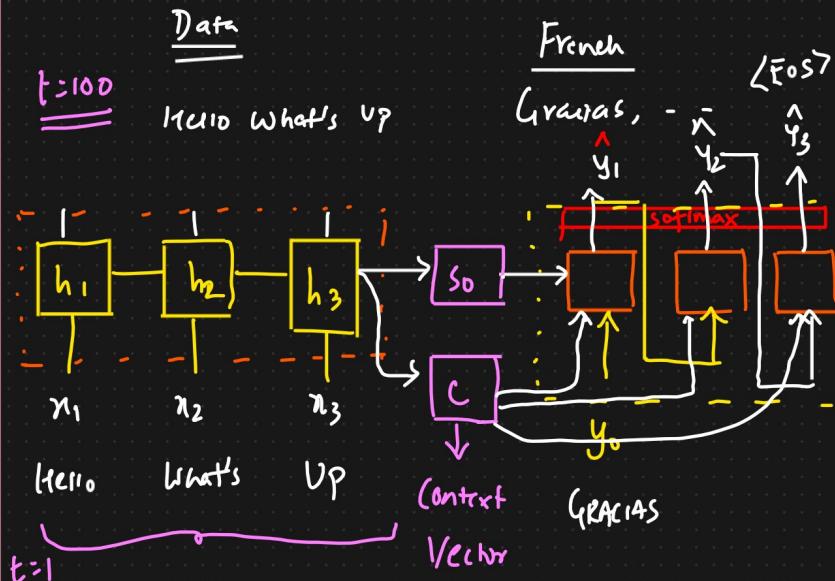
\Rightarrow Seq to Seq Data

* Attention Mechanism → Seq2Seq Network

longer paragraph → {Context Vector}.

+
 {Context}

Attention Mechanism | Seq2Seq Networks



Encoder Decoder Architecture

Attention Mechanism

<https://erdem.pl/2021/05/introduction-to-attention-mechanism>

3 LEARNING TO ALIGN AND TRANSLATE

In this section, we propose a novel architecture for neural machine translation. The new architecture consists of a bidirectional RNN as an encoder (Sec. 3.2) and a decoder that emulates searching through a source sentence during decoding a translation (Sec. 3.1).

3.1 DECODER: GENERAL DESCRIPTION

In a new model architecture, we define each conditional probability in Eq. (2) as:

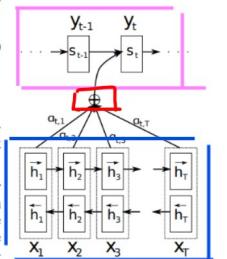
$$p(y_i|y_1, \dots, y_{i-1}, \mathbf{x}) = g(y_{i-1}, s_i, c_i), \quad (4)$$

where s_i is an RNN hidden state for time i , computed by

$$s_i = f(s_{i-1}, y_{i-1}, c_i).$$

It should be noted that unlike the existing encoder-decoder approach (see Eq. (2)), here the probability is conditioned on a distinct context vector c_i for each target word y_i .

The context vector c_i depends on a sequence of annotations (h_1, \dots, h_{T_x}) to which an encoder maps the input sentence. Each annotation h_i contains information about the whole input sequence with a strong focus on the parts surrounding the i -th word of the input sequence. We explain in detail how the annotations are computed in the next section.



The context vector c_i is, then, computed as a weighted sum of these annotations h_j :

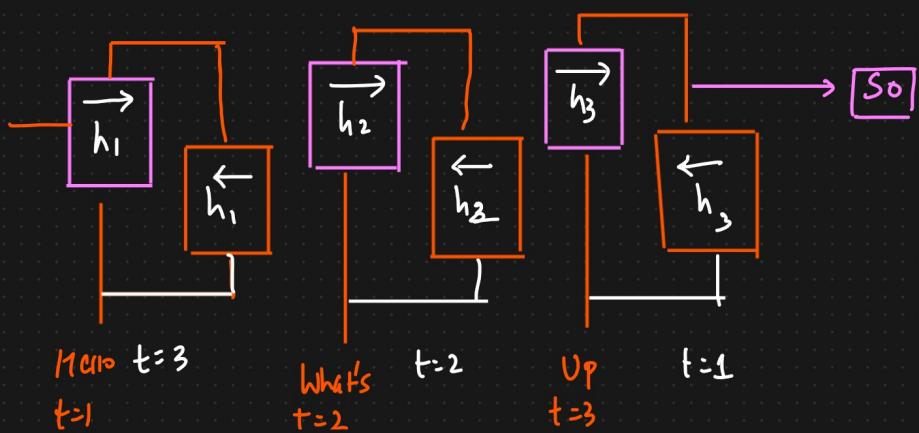
$$\{ c_i = \sum_{j=1}^{T_x} \alpha_{ij} h_j. \} \quad (5)$$

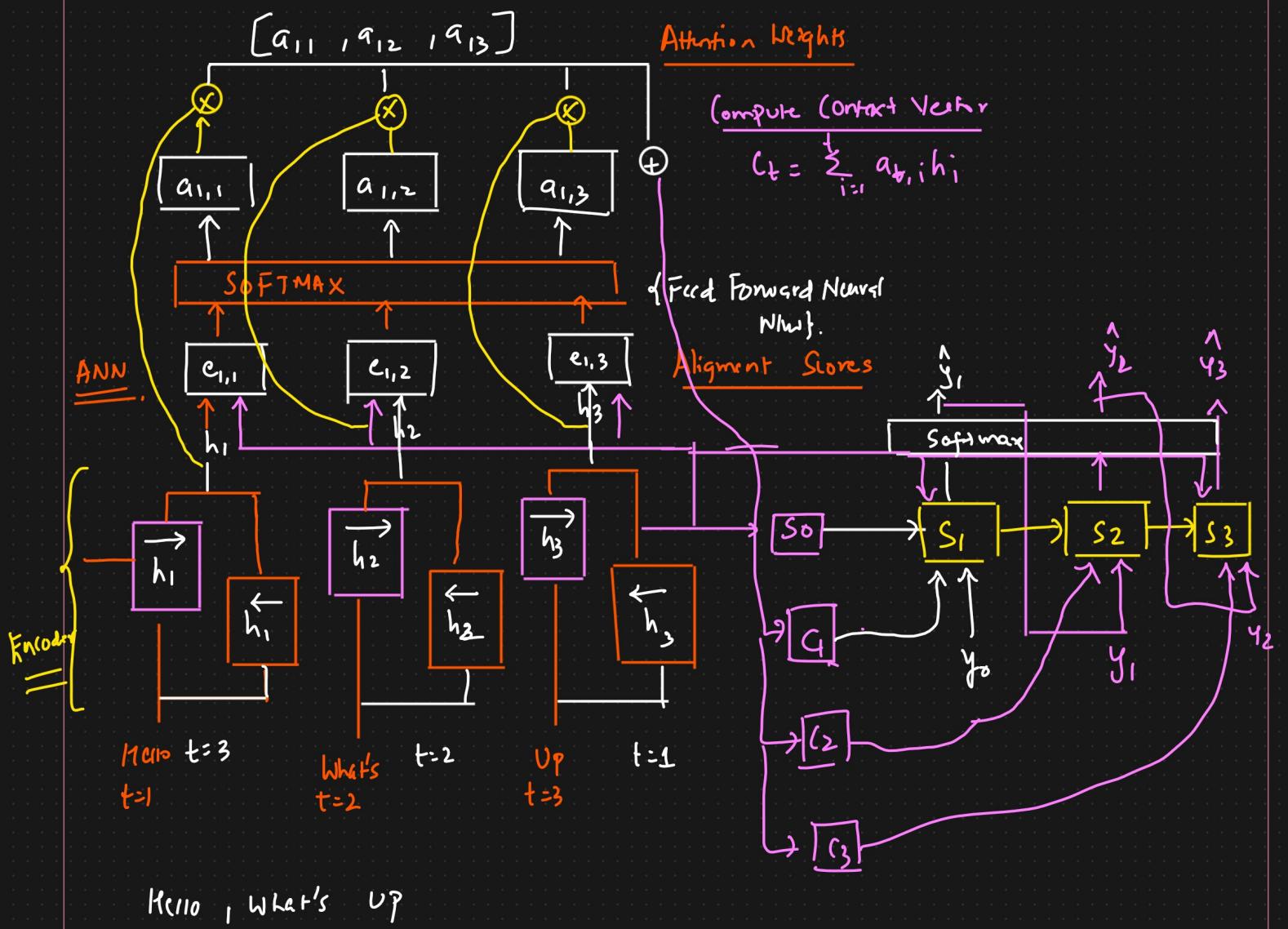
The weight α_{ij} of each annotation h_j is computed by

$$\alpha_{ij} = \left\{ \frac{\exp(e_{ij})}{\sum_{k=1}^{T_x} \exp(e_{ik})} \right\}$$

where

$$e_{ij} = a(s_{i-1}, h_j)$$





Introduction To Transformers

- 1) RNN / LSTM / GRU RNN
- 2) Encoder Decoder Architecture
- 3) ATTENTION MECHANISM
- 4) TRANSFORMERS

① Why Transformers?

② Architecture of Transformers?

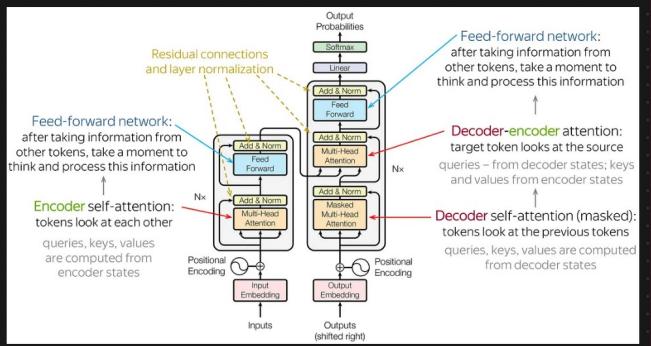
③ SELF ATTENTION $\rightarrow Q, K, V$

④ Positional Encoding

⑤ Multi Head ATTENTION

⑥ Combining the Working of Transformers

Architecture



Generative AI \rightarrow LM, Multimodel

BERT, GPT \leftarrow

Open AI \rightarrow ChatGPT

GPT-4o

① What And Why \rightarrow Transformers

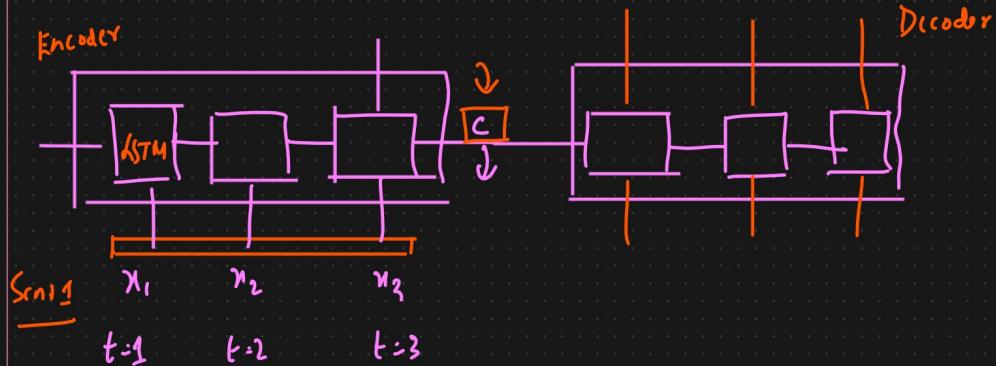
Transformers in natural language processing (NLP) are a type of deep learning model that use self-attention mechanisms to analyze and process natural language data. They are encoder-decoder models that can be used for many applications, including machine translation. \Rightarrow Seq2Seq Task

Eg: Language Translation \rightarrow Google Translation

English \rightarrow French

if \Rightarrow Many \rightarrow O/p: Many. $\{$ Length of the sentence $\}$.

Encoder - Decoder



Sentence length $\uparrow \uparrow$

Beam Score $\downarrow \downarrow$

Length Sentence $\uparrow \uparrow$

3.1 DECODER: GENERAL DESCRIPTION

In a new model architecture, we define each conditional probability in Eq. (2) as:

$$p(y_i | y_1, \dots, y_{i-1}, s_i, c_i), \quad (4)$$

where s_i is an RNN hidden state for time i , computed by

$$s_i = f(s_{i-1}, y_{i-1}, c_i).$$

It should be noted that unlike the existing encoder-decoder approach (see Eq. (2)), here the probability is conditioned on a distinct context vector c_i for each target word y_i .

The context vector c_i depends on a sequence of annotations (h_1, \dots, h_{T_s}) to which an encoder maps the input sentence. Each annotation h_i contains information about the whole input sequence with a strong focus on the parts surrounding the i -th word of the input sequence. We explain in detail how the annotations are computed in the next section.

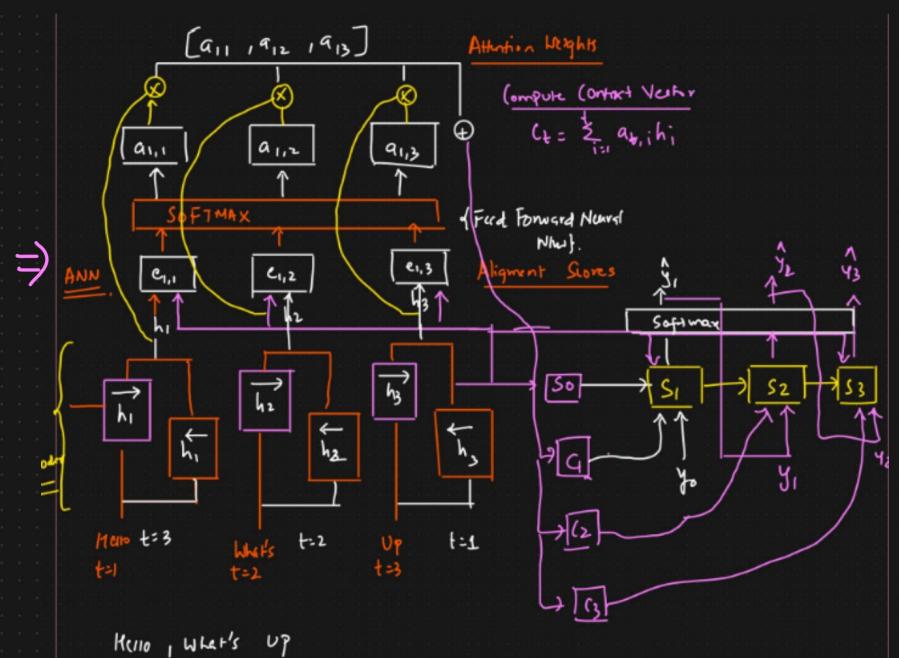
The context vector c_i is, then, computed as a weighted sum of these annotations h_i :

$$c_i = \sum_{j=1}^{T_s} \alpha_{ij} h_j. \quad (5)$$

Figure 1: The graphical illustration of the proposed model trying to generate the i -th target word y_i given a source sentence (x_1, x_2, \dots, x_T) .

Additional Context \rightarrow Decoder

long Sentence Accuracy ↑↑



Attention Mechanism

① Parallelly we cannot send all the words in a sentence \rightarrow Scalable

DATASET \rightarrow Huge \rightarrow Scalable With Respect to Training.

TRANSFORMERS \neq LSTM RNN

Self Attention Module \leftarrow All the words will be parallelly sent to encoder.



Positional Encoding

Transformer ↑ DATASET \rightarrow Amazing SOTA \leftarrow NLP

Transfer Learning \rightarrow MultiModal Task \rightarrow NLP + Image \leftarrow

Transformers \div AI Space \rightarrow SOTA Model \rightarrow



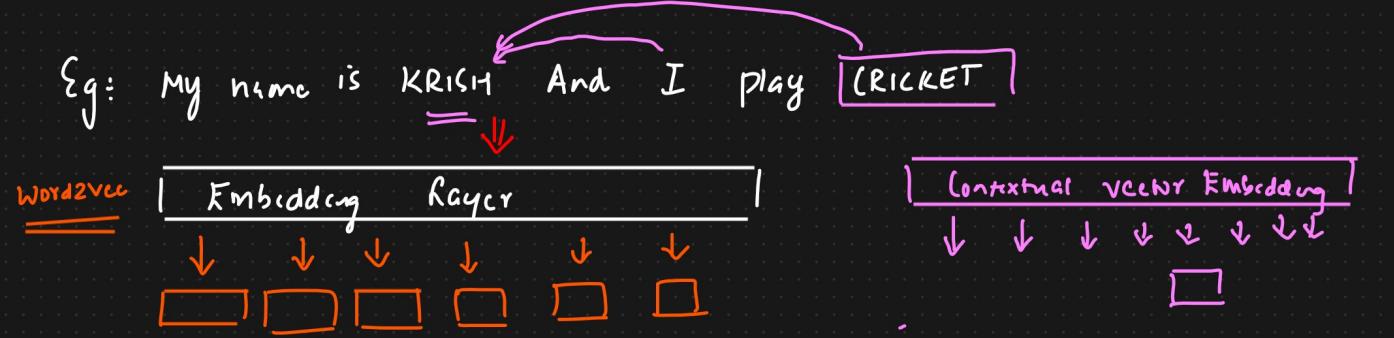
BERT GPT \rightarrow Transfer learning \rightarrow SOTA Models \rightarrow DALLE \leftarrow Generating AI

Train huge Data

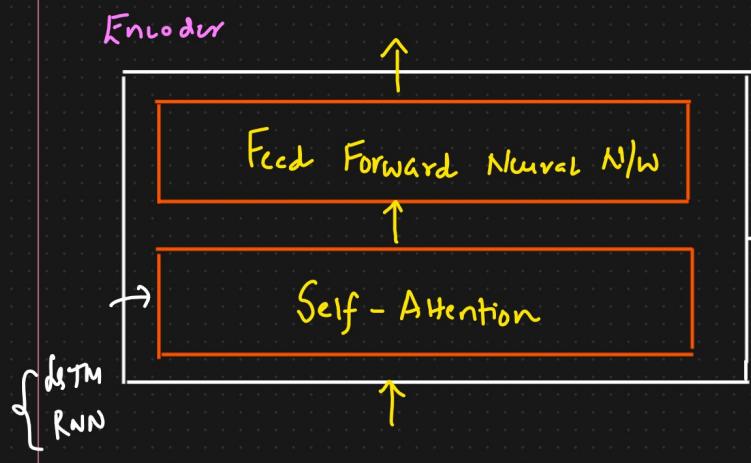
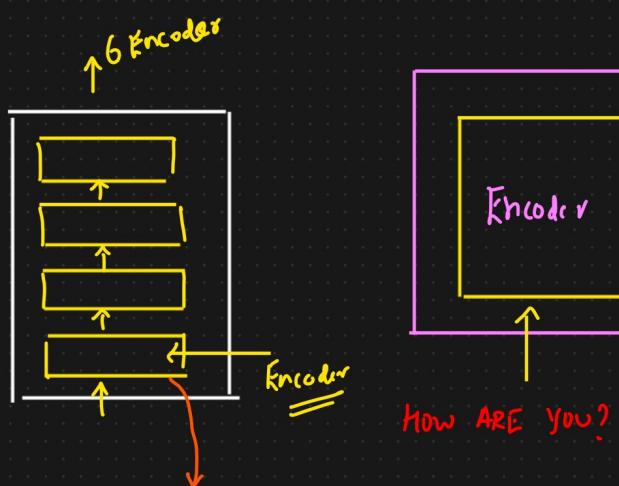
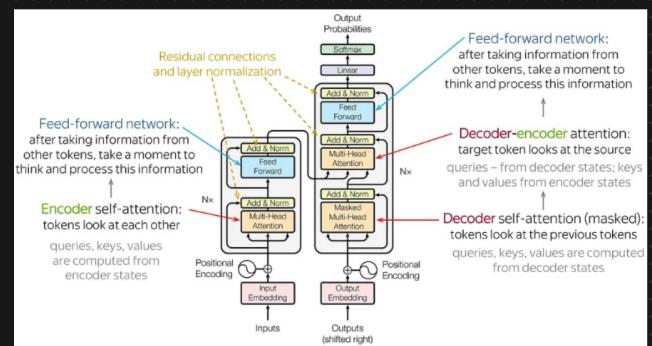
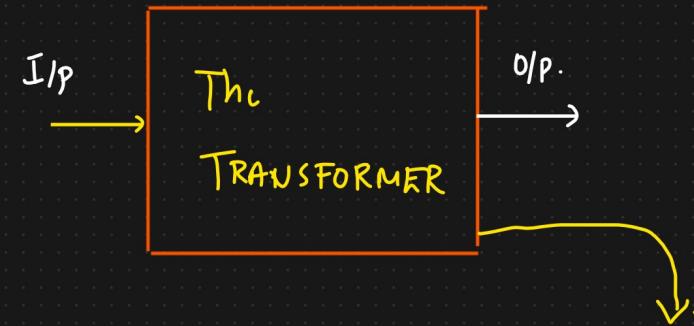
AI's

Generating AI

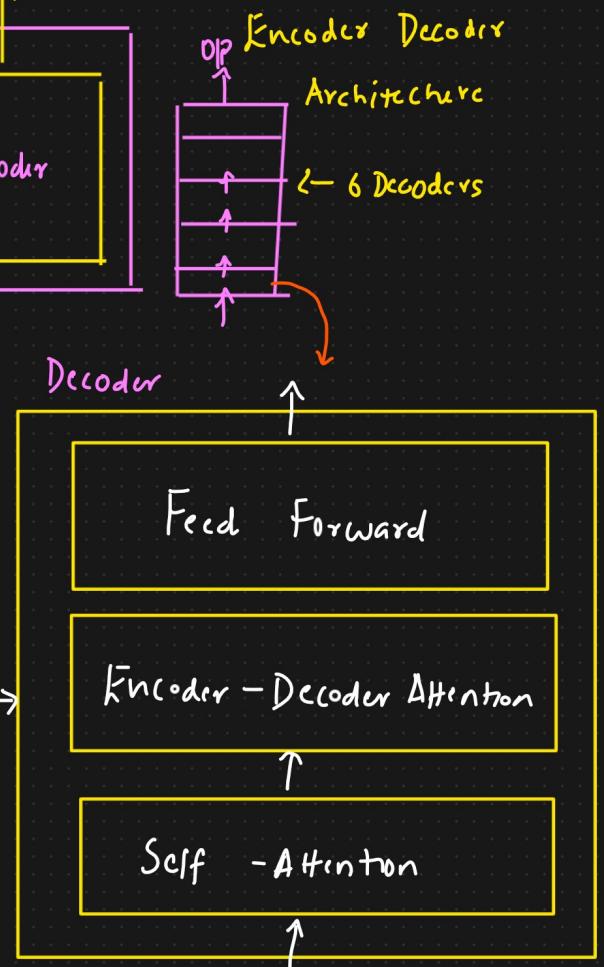
② Contextual Embedding → Self Attention



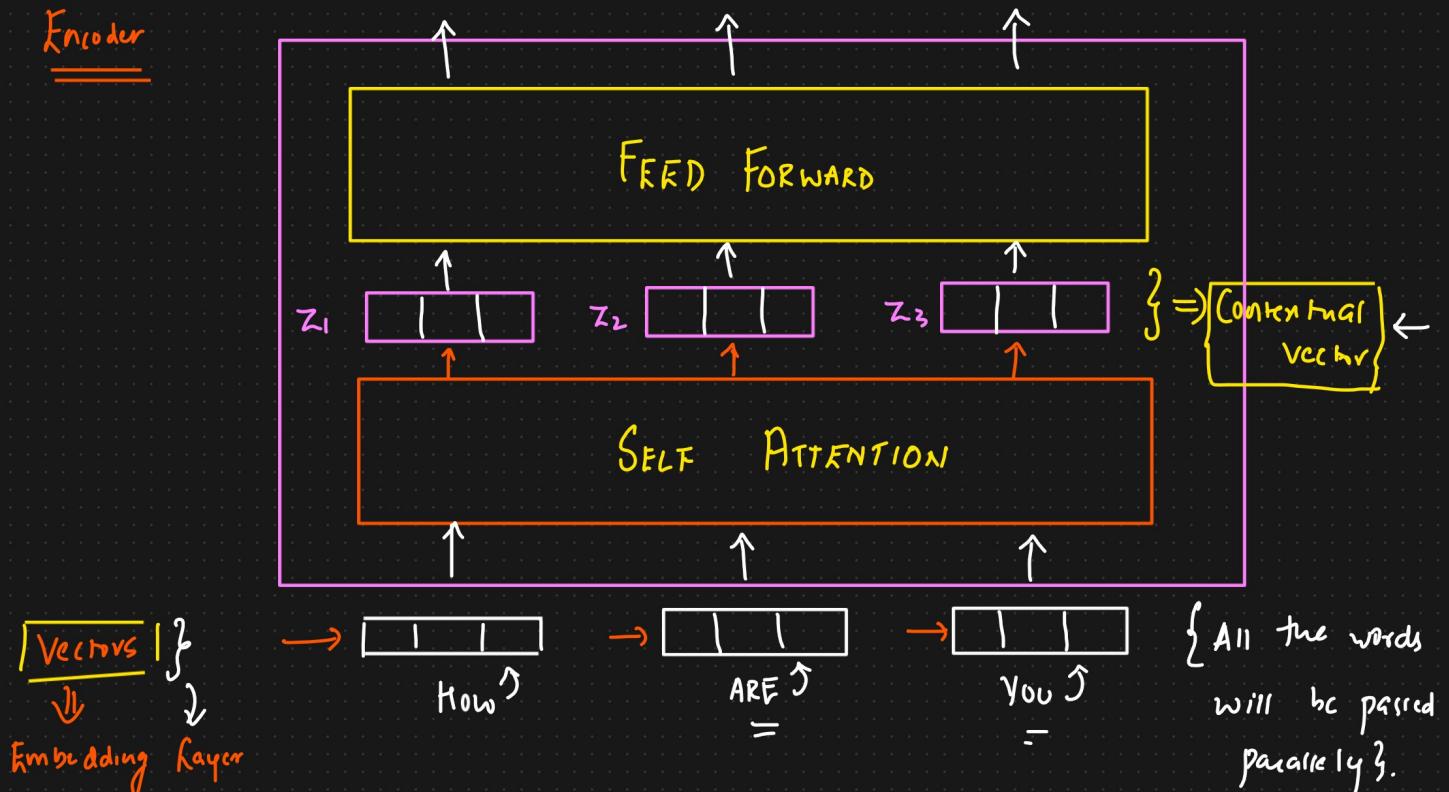
② Basic Transformer Architecture {Seq2Seq Task} → Language Translation {Eng → French}



Comment vas-tu?

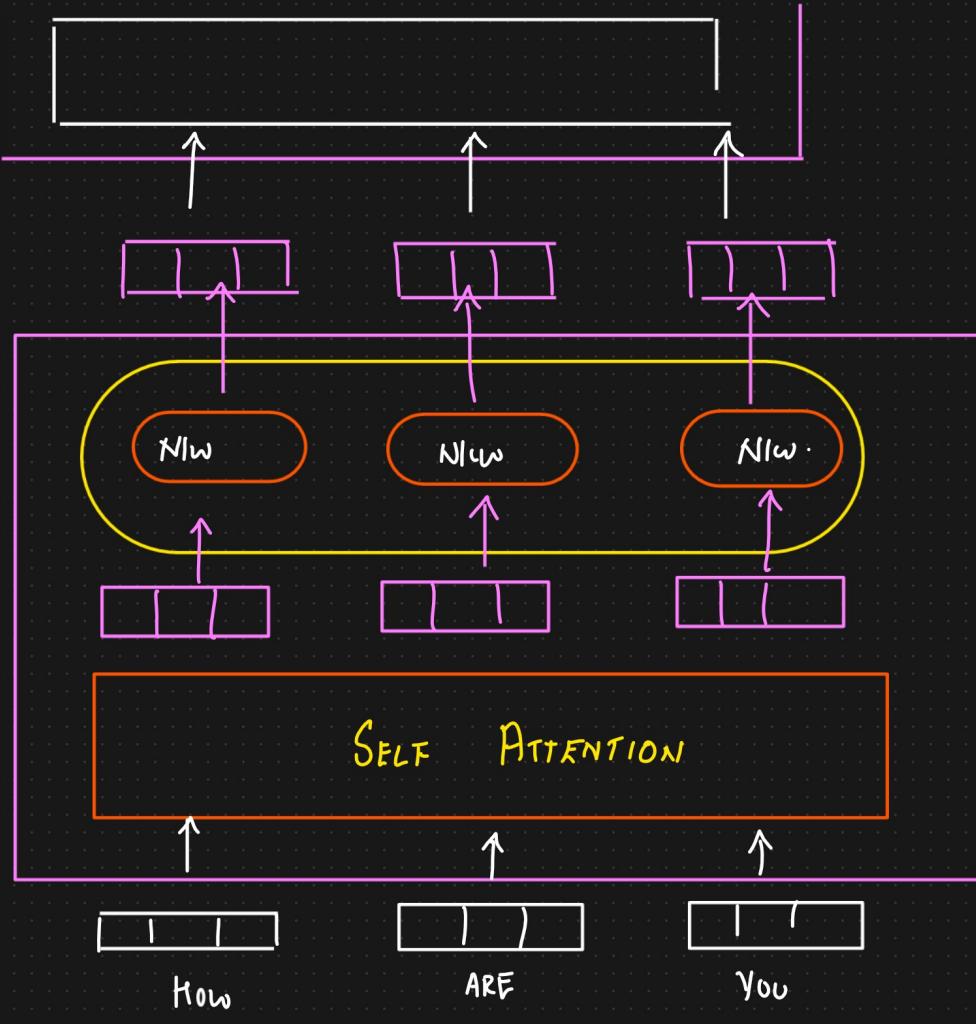


Encoder



Encoder¹

Encoder



Self Attention At a Higher Level

Eg.: The cat sat on the mat, the cat lay on the rug.

Word Embedding
↓ ↓ ↓ ↓ ↓ ↓ ↓ ↓ ↓ ↓ ↓ ↓ .

SELF ATTENTION

[] [] [] [] [] []
↓ .
[]

The. The

→ [Cat] --- [Cat] → [] [] []
1 → [Sat] Sat

= Rank 2 → On
3 → [he]
mat

{ Contextual Embedding }

mat → [- | - | -]



Self Attention In Detail

① To Create 3 vectors from each of the encoder i/p. Query vector,

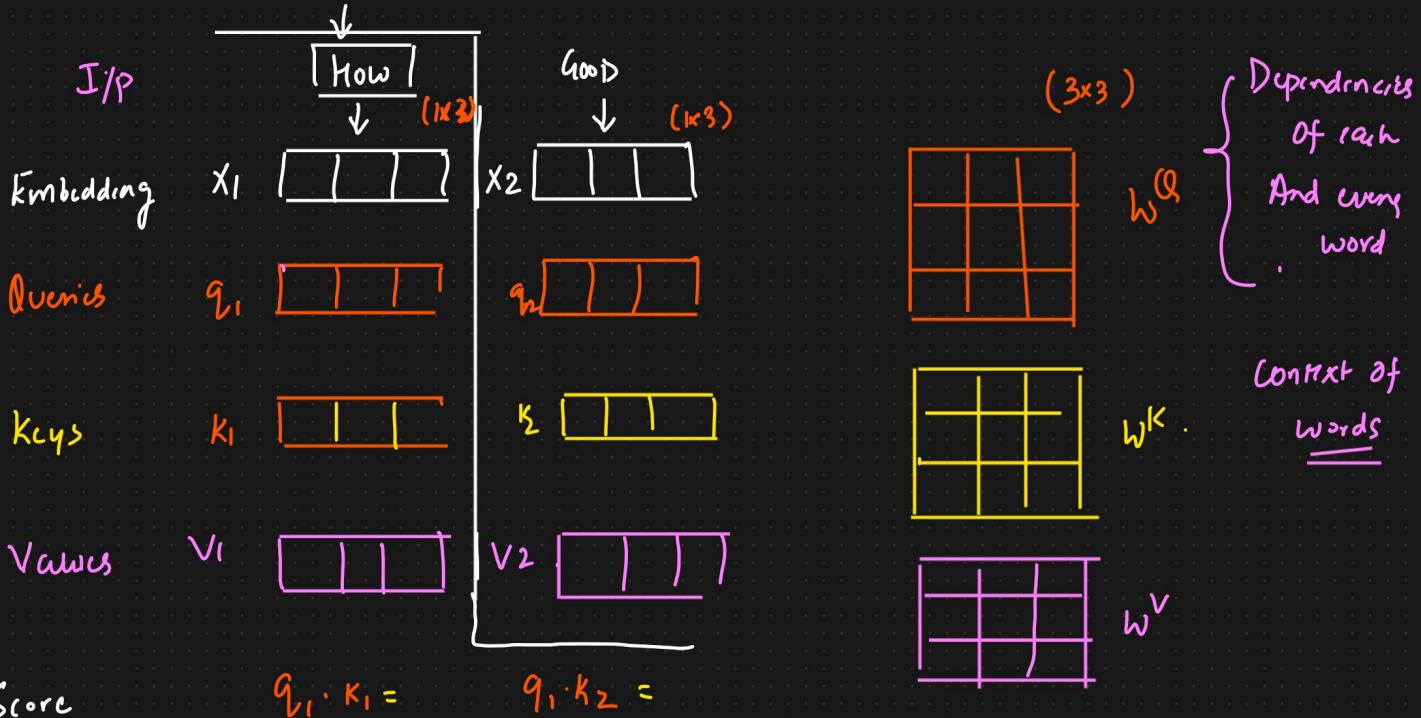
Key vector, value vector. ⇒ Contextual Embedding

Eg.: YT \Rightarrow Search keywords $\Rightarrow \{ \text{Query} \}$.
Topic

Query $\rightarrow \{ \text{Key} \} \rightarrow \begin{array}{l} \text{Tags} \\ \text{Description} \\ \text>Title \end{array} \Rightarrow \text{Value} \Rightarrow \text{o/p Video}$

② Second step in calculating self attention is to calculate the Score.

The Score determines how much focus to place on the other part of the sentence.



$$\text{Score} \quad q_1 \cdot k_1 = \quad q_1 \cdot k_2 =$$

④ How much focus to place on other part of the I/P Sentence as we encode a word at certain position

Divide \sqrt{dk}

More stable gradients

SoftMax

Dimension = 3 Dimension $\uparrow = T/2$

\downarrow

$\frac{\text{Value}}{\sqrt{dk}}$ \approx $\frac{\text{Variance}}{\sqrt{dk}} \Rightarrow$

$$\frac{[0-1]}{\sqrt{dk}} + \frac{[0-1]}{\sqrt{dk}} = 1$$

Softmax \Rightarrow The Softmax score determines how much each word will be expressed at this position

X X X

Value Vectors v_1 [| |] v_2 [| | |]

0.88 [| | |]

$\frac{0.88}{0.88+0.88+0.88} + \frac{0.12}{0.88+0.88+0.88}$

$\downarrow \quad \downarrow$

Contextual Vector $\xrightarrow{T_1}$ [| | |] $\xrightarrow{T_2}$ [| | |]

Self Attention At Higher And Detailed Level

Self-attention, also known as scaled dot-product attention, is a crucial mechanism in the transformer architecture that allows the model to weigh the importance of different tokens in the input sequence relative to each other.

Idea :



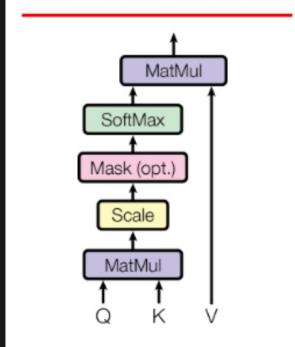
Language Translation Self Attention

Text Summarization Self Attention $\{ \text{Sentence, Dataset} \}$

Embedding \leftarrow The CAT SAT

layer \rightarrow Fixed Vector

Scaled Dot-Product Attention



i) Inputs: Queries, Keys, And Values

Model \rightarrow Queries, Keys And Values

1. Query Vectors (Q):

Role: Query vectors represent the token for which we are calculating the attention. They help determine the importance of other tokens in the context of the current token.

Importance:

Focus Determination: Queries help the model decide which parts of the sequence to focus on for each specific token. By calculating the dot product between a query vector and all key vectors, the model assesses how much attention to give to each token relative to the current token.

Contextual Understanding: Queries contribute to understanding the relationship between the current token and the rest of the sequence, which is essential for capturing dependencies and context.

2. Key Vectors (K):

Role: Key vectors represent all the tokens in the sequence and are used to compare with the query vectors to calculate attention scores.

Importance:

Relevance Measurement: Keys are compared with queries to measure the relevance or compatibility of each token with the current token. This comparison helps in determining how much attention each token should receive.

Information Retrieval: Keys play a critical role in retrieving the most relevant information from the sequence by providing a basis for the attention mechanism to compute similarity scores.

3. Value Vectors (V):

Role: Value vectors hold the actual information that will be aggregated to form the output of the attention mechanism.

Importance:

Information Aggregation: Values contain the data that will be weighted by the attention scores. The weighted sum of values forms the output of the self-attention mechanism, which is then passed on to the next layers in the network.

Context Preservation: By weighting the values according to the attention scores, the model preserves and aggregates relevant context from the entire sequence, which is crucial for tasks like translation, summarization, and more.

$$\text{Input Sequence} = \left[\text{"The"}, \text{"(A)"}, \text{"SAT"} \right]$$

Embedding size = 4

$$E_{\text{Thc}} = [1 \ 0 \ 1 \ 0]$$

$$E_{CAT} = \begin{bmatrix} 0 & 1 \\ 0 & 1 \end{bmatrix}$$

$$E_{Sqr} = [1, 1, 1, 1]$$

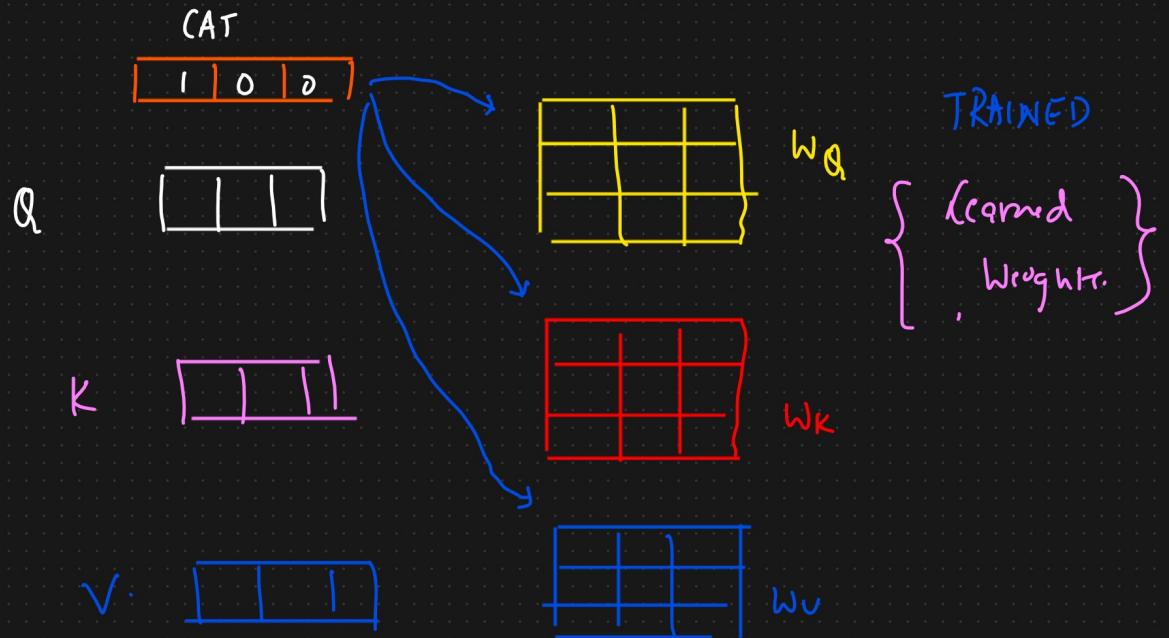
$$Q_{1K}, V \Rightarrow 4$$



Sentiment, Dataset

② Linear Transformation

We create Q, K, V by multiplying the embeddings by learned weights matrices W_Q , W_K and W_V .



lets consider

$$W_Q = W_K = W_V = I \quad \begin{bmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \end{bmatrix} \leftarrow$$

$$Q_{\text{Trained}} = \begin{bmatrix} 1 & 0 & 0 \end{bmatrix} \cdot \begin{bmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \end{bmatrix} = \begin{bmatrix} 1 & 0 & 0 \end{bmatrix}$$

$$K_{\text{Trained}} = \begin{bmatrix} 1 & 0 & 0 \end{bmatrix}$$

$$V_{\text{Trained}} = \begin{bmatrix} 1 & 0 & 0 \end{bmatrix}$$

$$\textcircled{1} \quad Q_{\text{Trained}} = K_{\text{Trained}} = V_{\text{Trained}} = \begin{bmatrix} 1 & 0 & 0 \end{bmatrix}$$

$$\textcircled{2} \quad Q_{\text{CAT}} = K_{\text{CAT}} = V_{\text{CAT}} = \begin{bmatrix} 0 & 1 & 0 \end{bmatrix}$$

$$\textcircled{3} \quad Q_{\text{SAT}} = K_{\text{SAT}} = V_{\text{SAT}} = \begin{bmatrix} 1 & 1 & 1 & 1 \end{bmatrix}$$

③ Compute Attention Scores

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

The

$$\text{Score}(Q_{\text{The}}, K_{\text{The}}) = \begin{bmatrix} 1 & 0 & 1 & 0 \end{bmatrix} \cdot \begin{bmatrix} 1 & 0 & 1 & 0 \end{bmatrix}^T = 2$$

$$\text{Score}(Q_{\text{The}}, K_{\text{CAT}}) = \begin{bmatrix} 1 & 0 & 1 & 0 \end{bmatrix} \cdot \begin{bmatrix} 0 & 1 & 0 & 1 \end{bmatrix}^T = 0$$

$$\text{Score}(Q_{\text{The}}, K_{\text{SAT}}) = \begin{bmatrix} 1 & 0 & 1 & 0 \end{bmatrix} \cdot \begin{bmatrix} 1 & 1 & 1 & 1 \end{bmatrix}^T = 2$$

For the token CAT

$$\text{Score}(Q_{\text{CAT}}, K_{\text{The}}) = \begin{bmatrix} 0 & 1 & 0 & 1 \end{bmatrix} \cdot \begin{bmatrix} 1 & 0 & 1 & 0 \end{bmatrix}^T = 0$$

$$\text{Score}(Q_{\text{CAT}}, K_{\text{CAT}}) = \begin{bmatrix} 0 & 1 & 0 & 1 \end{bmatrix} \cdot \begin{bmatrix} 0 & 1 & 0 & 1 \end{bmatrix}^T = 2$$

$$\text{Score}(Q_{\text{CAT}}, K_{\text{SAT}}) = \begin{bmatrix} 0 & 1 & 0 & 1 \end{bmatrix} \cdot \begin{bmatrix} 1 & 1 & 1 & 1 \end{bmatrix}^T = 2.$$

For the Token SAT

$$\left\{ \begin{array}{l} \text{Score}(Q_{\text{SAT}}, K_{\text{The}}) = \begin{bmatrix} 1 & 1 & 1 & 1 \end{bmatrix} \cdot \begin{bmatrix} 1 & 0 & 1 & 0 \end{bmatrix}^T = 2 \\ \text{Score}(Q_{\text{SAT}}, K_{\text{CAT}}) = 2 \\ \text{Score}(Q_{\text{SAT}}, K_{\text{SAT}}) = 4 \end{array} \right.$$

④ Scaling

We take up the scores and scale down by dividing the scores by the

$$\sqrt{d_K} \Rightarrow d_K = 4 \quad \sqrt{d_K} = 2.$$

Scaling in the attention mechanism is crucial to prevent the dot product from growing too large. \Rightarrow Ensure stable gradients during Training.

d_K is large \rightarrow

① Gradient Exploding

② Softmax Saturation $\{\curvearrowright\}$ \rightarrow Vanishing Gradient Problem.

$$Q = \begin{bmatrix} 2 & 3 & 4 & 1 \end{bmatrix} \quad K_1 = \begin{bmatrix} 1 & 0 & 1 & 0 \end{bmatrix} \quad K_2 = \begin{bmatrix} 0 & 1 & 0 & 1 \end{bmatrix}$$

Without Scaling

$$Q \cdot K_1^T = 2 \times 1 + 3 \times 0 + 4 \times 1 + 1 \times 0 = 2 + 4 = 6.$$

$$Q \cdot K_2^T = 2 \times 0 + 3 \times 1 + 4 \times 0 + 1 \times 1 = 0 + 3 + 0 + 1 = 4$$

* Score $[6, 4] \Rightarrow$ Scaling Not Applied

$$\text{Softmax}([6, 4]) = \left[\frac{e^6}{e^6 + e^4}, \frac{e^4}{e^6 + e^4} \right] = \left[\frac{e^6}{e^6(1 + e^{-2})}, \frac{e^4}{e^4(e^2 + 1)} \right]$$

① Property of Softmax w_Q, w_K, w_V

$$([10, 1]) = \left[\frac{0.99}{\cancel{10}}, \frac{\cancel{0.01}}{\cancel{1}} \right] = \left[\frac{1}{(1 + e^{-2})}, \frac{1}{(e^2 + 1)} \right]$$

Dot product = Large values $\approx [0.88, 0.12]$.

Most of the attention weight is assigned to the first key vector,
very little to the second vector,

With Scaling

① Compute Scaled Dot Product

$$[6, 4] \Rightarrow \text{Scale} \Rightarrow \left[\frac{6}{\sqrt{2}}, \frac{4}{\sqrt{2}} \right] = \left[\frac{3}{\sqrt{2}}, \frac{2}{\sqrt{2}} \right]. \quad \frac{\sqrt{d_K} = \sqrt{4} = \sqrt{2}}{\sqrt{d_K} \uparrow \text{same dimension} \uparrow \text{Variance} \uparrow 2}$$

$$\text{Softmax}\left(\left[\frac{3}{\sqrt{2}}, \frac{2}{\sqrt{2}} \right]\right) = \left[\frac{e^3}{e^3 + e^2}, \frac{e^2}{e^3 + e^2} \right] = \left[\frac{e^3}{e^3(1 + e^{-1})}, \frac{e^2}{e^2(e^1 + 1)} \right] = [0.73, 0.27] \Rightarrow \text{Attention weights}$$

(4) Here, the attention weights are more balanced compared to the unscaled case

Summary of Importance

Stabilizing Training: Scaling prevents extremely large dot products, which helps in stabilizing the gradients during backpropagation, making the training process more stable and efficient.

Preventing Saturation: By scaling the dot products, the softmax function produces more balanced attention weights, preventing the model from focusing too heavily on a single token and ignoring others.

Improved Learning: Balanced attention weights enable the model to learn better representations by considering multiple relevant tokens in the sequence, leading to better performance on tasks that require context understanding.

Scaling ensures that the dot products are kept within a range that allows the softmax function to operate effectively, providing a more balanced distribution of attention weights and improving the overall learning process of the model.

$$(4) \text{Scaling} = \sqrt{d_K} = \sqrt{4} \Rightarrow 2$$

Similarly Scaling

$$\text{Scaled-Score } (Q_{\text{The}}, K_{\text{The}}) = 2/2 = 1$$

will be done for
all other tokens.

$$\text{Scaled-Score } (Q_{\text{The}}, K_{\text{CAT}}) = 0/2 = 0$$

$$\text{Scaled-Score } (Q_{\text{The}}, K_{\text{SAT}}) = 2/2 = 1$$

(5) Apply Softmax

$$\text{ATTENTION WEIGHTS}_{\text{"The"}} = \text{Softmax}([1, 0, 1]) = [0.4223, 0.1554, 0.4223]$$

$$\text{ATTENTION WEIGHTS}_{\text{"CAT"}} = \text{Softmax}([0, 2, 2]) = [0.1554, 0.4223, 0.4223]$$

$$\text{ATTENTION WEIGHTS}_{\text{"SAT"}} = \text{Softmax}([2, 2, 4]) = [0.2119, 0.2119, 0.5762]$$

(6) Weight Sum of Values

We multiply the attention weights by corresponding value vectors

For the Token The =

$$\text{Output}_{(\text{The})} = 0.4223 * V_{\text{The}} + 0.1554 * V_{\text{CAT}} + 0.4223 * V_{\text{Sal.}}$$

$$= 0.4223 [1 \ 0 \ 1 0] + 0.1554 [0 \ 1 0 1] + 0.4223 [1 1 1]$$

$$= [0.4223, 0, 0.4223, 0] + [0, 0.1554, 0, 0.1554] + [0.4223, 0.4223, \\ 0.4223, 0.4223]$$

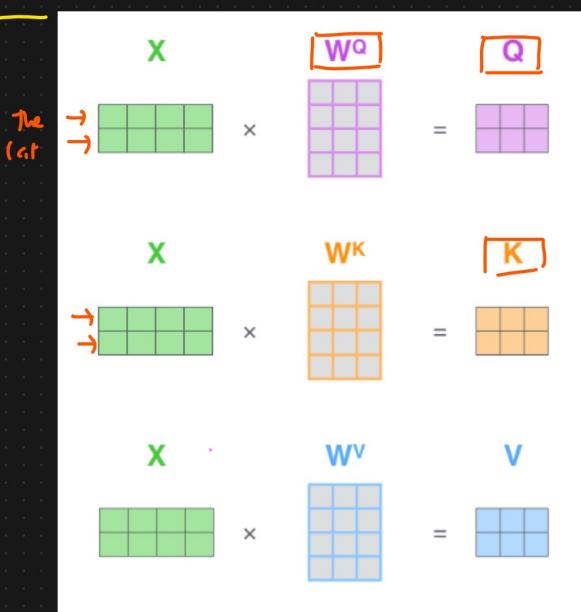
$$= [1.2669, 0.9999, 1.2669, 0.9999].$$

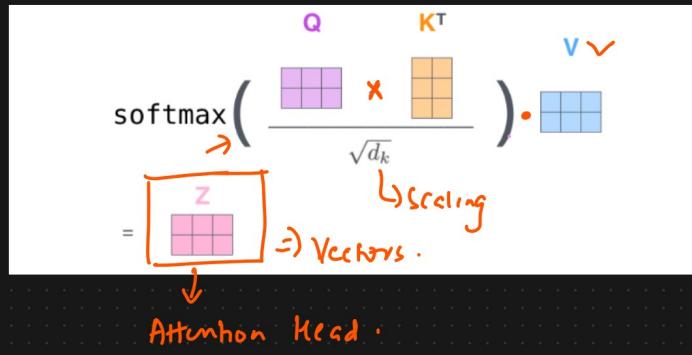
↓ Contextual
vector

The 1 1 0 1 1 0 \Rightarrow Self Attention $\Rightarrow [1.2669, 0.9999, 1.2669, 0.9999].$

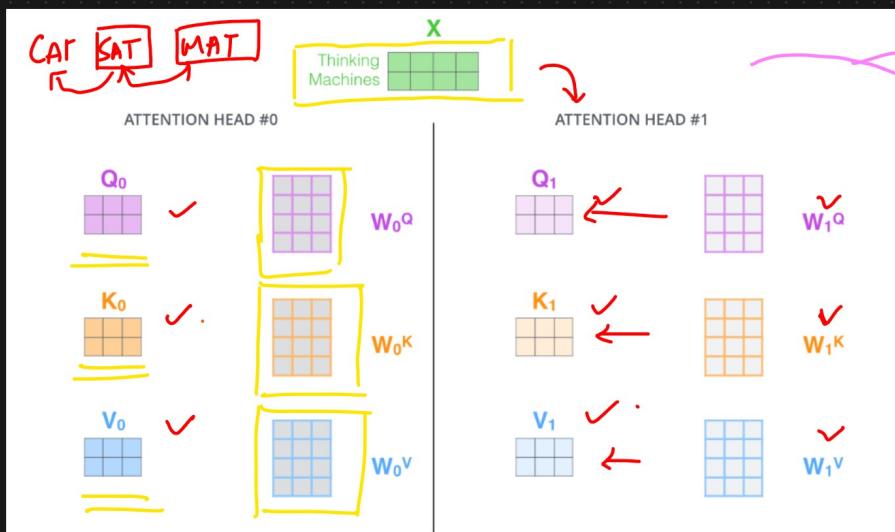
- ① $\hookrightarrow Q, K, V [w^Q, w^K, w^V]$
② \hookrightarrow Attention Score
③ \hookrightarrow Scaled
④ \hookrightarrow Softmax
⑤ \hookrightarrow Weighted Sum of Value (Softmax \times V)

④ Multi Head Attention





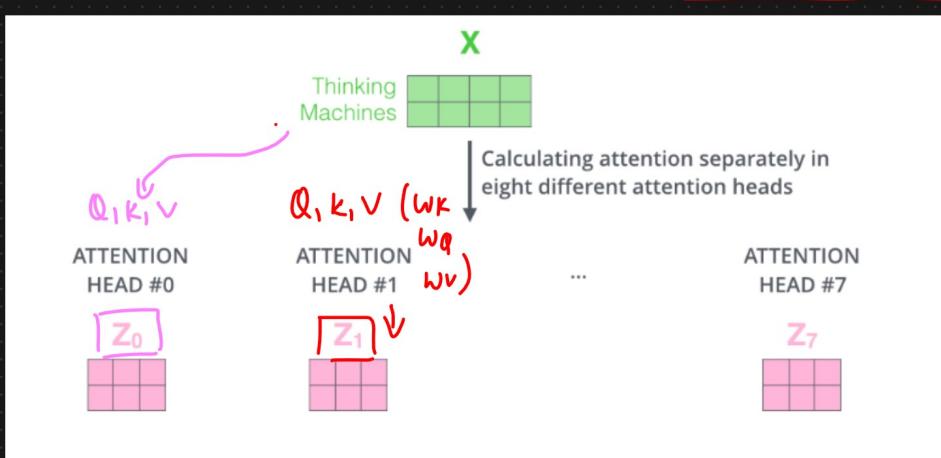
→ Self Attention with Multi Heads



It expands model's ability to focus on different position of tokens.

Score, Softmax $\times V$
 \downarrow \downarrow
 $\boxed{\dots}$ $\boxed{\dots}$
 Z_0 Z_1
 Vectors

Multi Head Attention

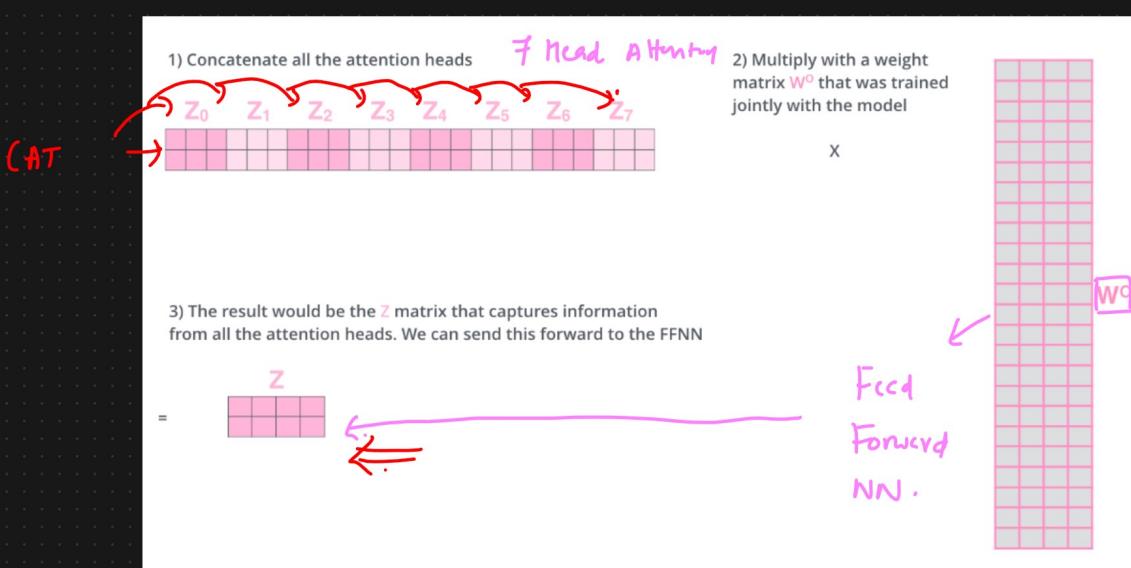
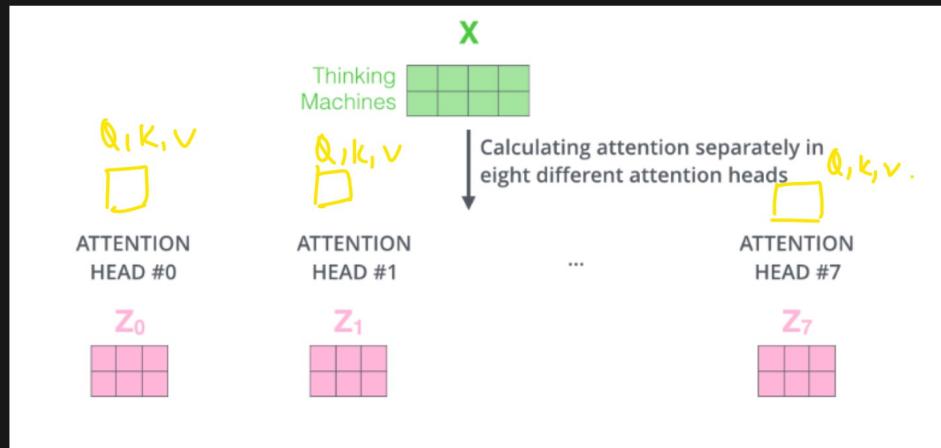
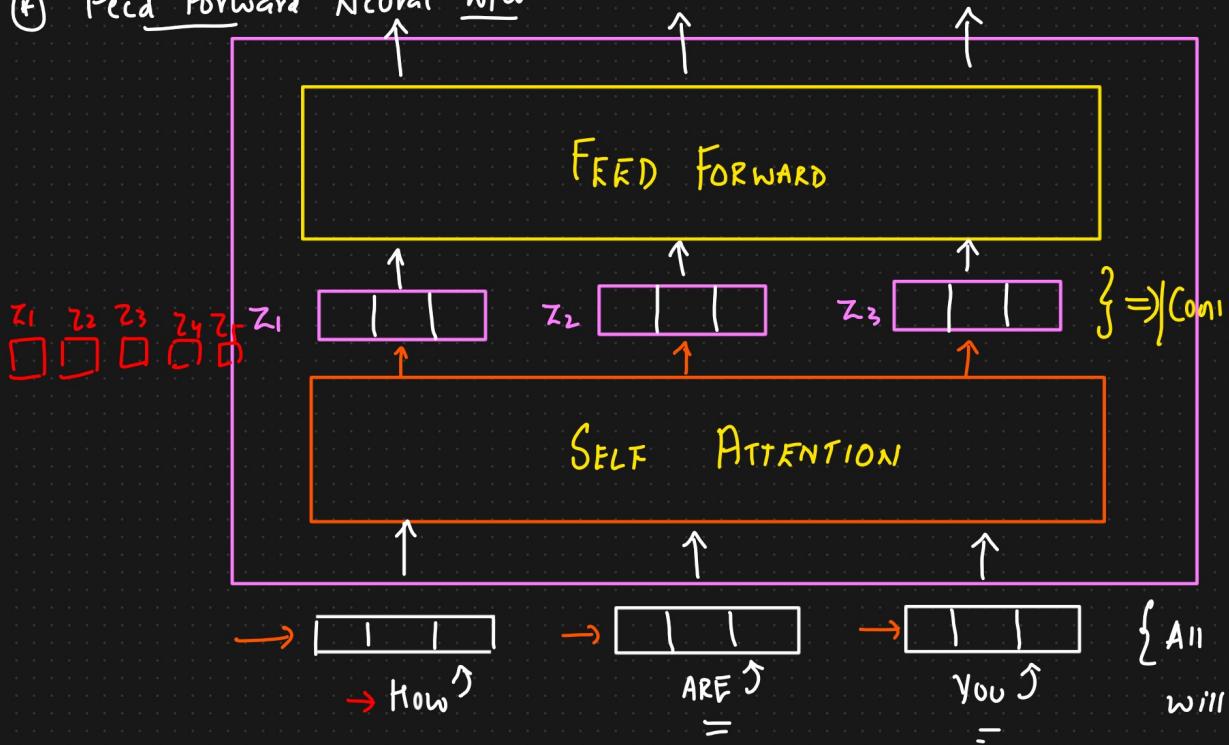


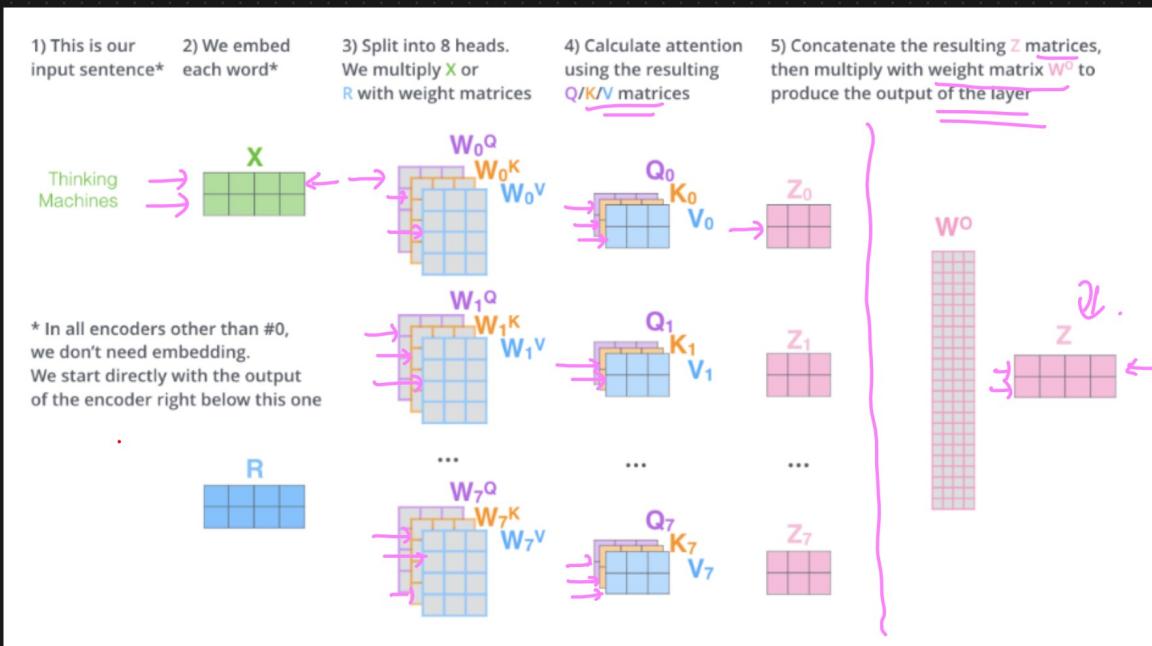
\boxed{z}

\boxed{z}

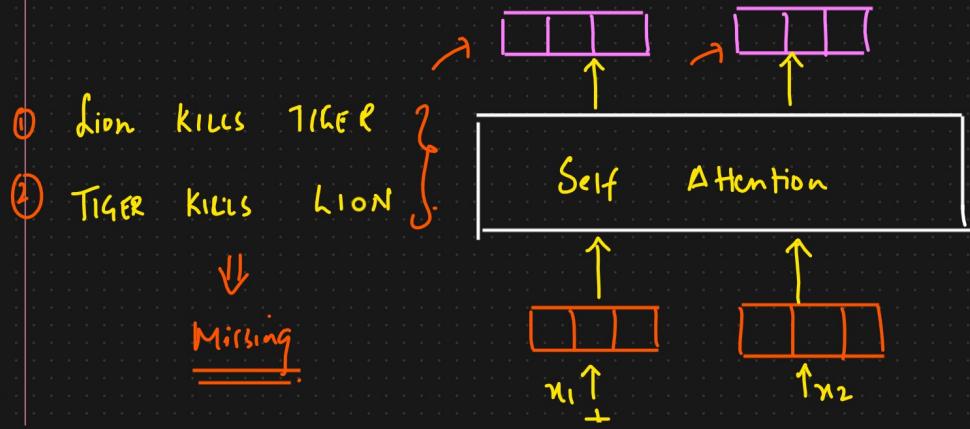
\boxed{z}

⑥ Feed Forward Neural Network





⑥ Positional Encoding - Representing Order of Sequence



Advantage

- ① Word Tokens it can process parallelly



DRAW BACK

Lack the sequential structure of the words {order}

Attention IS All

You NEED



Positional
Encoded
Vector

Book, Journal

Newspaper

↳ 1 lakh words



Backpropagation



Types of Position Encoding

1) Sinusoidal Position Encoding →

2) Learned Positional Encoding ⇒ Positional Encoding Are learned during Training ↲

① Sinusoidal Positional Encoding : It uses Sinc and Cosine functions

of different frequencies to create positional encodings

Formula :



$$P.E(pos, 2i) = \sin\left(\frac{pos}{10000^{2i/d_{model}}}\right)$$

Where pos is the position
 i is the dimension

$$P.E(pos, 2i+1) = \cos\left(\frac{pos}{10000^{2i/d_{model}}}\right).$$

d_{model} is the dimensionality of the embeddings.

Eg: The Cat Sat

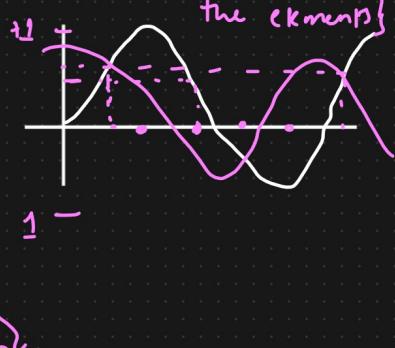
$$\text{The} \rightarrow [0.1 \ 0.2 \ 0.3 \ 0.4]$$

$$\text{CAT} \rightarrow [0.5 \ 0.6 \ 0.7 \ 0.8]$$

$$\text{SAT} \rightarrow [0.9 \ 1.0 \ 1.1 \ 1.2]$$



{ Miss the order of the elements }



$$PE(pos, 2i) = \sin\left(\frac{pos}{10000^{2i/d_{model}}}\right) \quad PE(pos, 2i+1) = \cos\left(\frac{pos}{10000^{2i/d_{model}}}\right)$$

For our example $d_{model} = 4$

For position $pos = 0$

$$PE(0,0) = \sin\left(\frac{0}{10000^0/4}\right) = \sin(0) = 0$$

$$PE(0,1) = \cos\left(\frac{0}{10000^1/4}\right) = \cos(0) = 1$$

$$PE(0,2) = \sin\left(\frac{0}{10000^2/4}\right) = \sin(0) = 0$$

$$PE(0,3) = \cos\left(\frac{0}{10000^3/4}\right) = 1$$

$$PE = [0, 1, 0, 1]$$

$$PE(pos, 2i) = \sin\left(\frac{pos}{10000^{2i/d_{model}}}\right)$$

For $pos = 1$

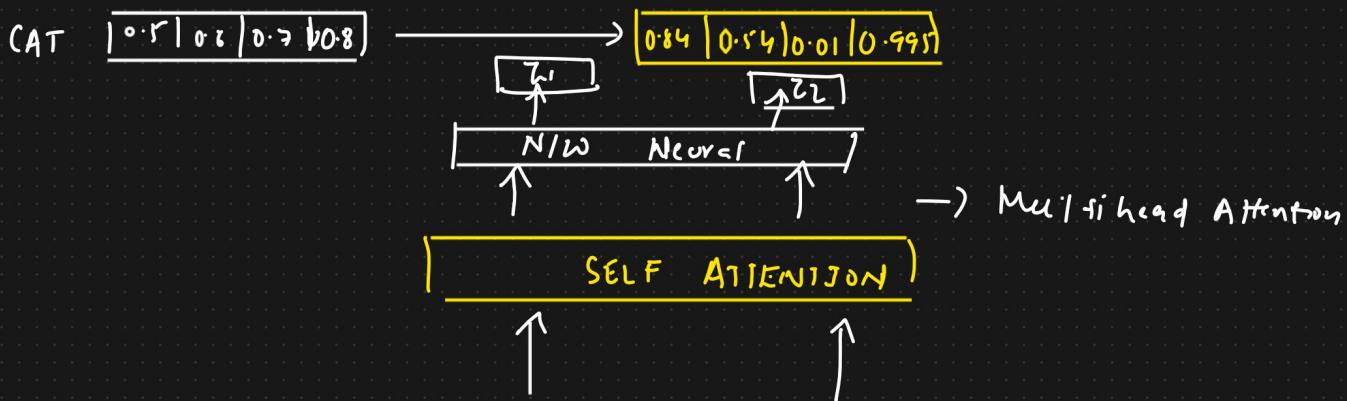
$$PE(1,0) = \sin\left(\frac{1}{10000^0/4}\right) = \sin(1) = 0.8415$$

$$PE(pos, 2i+1) = \cos\left(\frac{pos}{10000^{2i+1/d_{model}}}\right)$$

$$PE(1,1) = \cos\left(\frac{1}{10000^1/4}\right) \approx 0.5403$$

$$PE(1,2) = \sin\left(\frac{1}{10000^2/4}\right) \approx 0.01$$

$$PE(1,3) = \cos \approx 0.99995$$



Positional Encoding

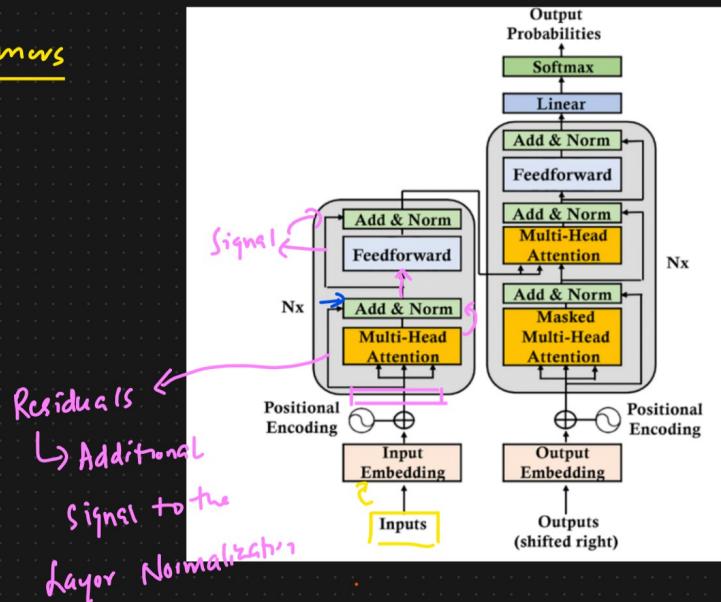
$$\begin{matrix} 0.84 & 0.79 & 0.01 & 0.99 \end{matrix} + \begin{matrix} 0.5 & 0.6 & 0.7 & 0.8 \end{matrix} = \begin{matrix} 1.0 & 1.0 & 0.1 & 0.4 \end{matrix}$$

CAT → Thc .

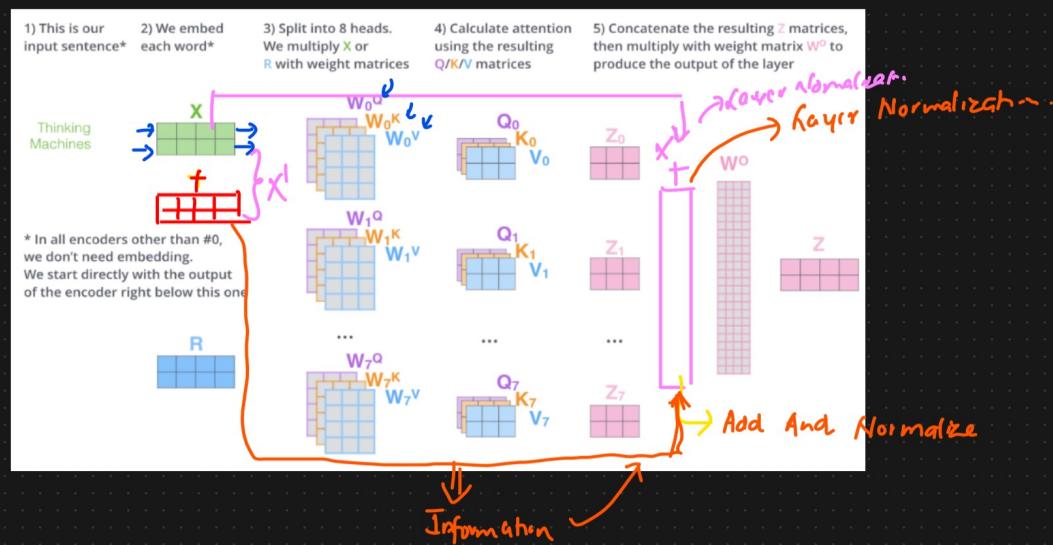
④ Layer Normalization In Transformers

Transformers

Residuals



- ① Self Attention layer
 - ② Multi head Attention
 - ③ Positional Encoding
 - ④ Layer Normalization
- ADD AND Normalize



Normalization

→ Batch Normalization
→ Layer Normalization



	f_1	f_2	Price
House Size	1200	2	45
No. of Rooms	1500	3	70

Normalization : Standard Scaling

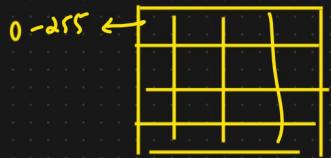
$$\text{z-score} = \frac{x_i - \mu}{\sigma}$$

$$\{\mu=0, \sigma=1\}$$

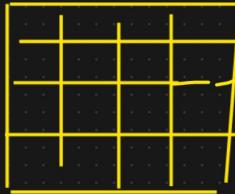
$$f_1 \rightarrow f_1' \Leftarrow \mu=0, \sigma=1$$



Deep learning : I/P Image \Rightarrow Min Max Scaler



\Rightarrow Min Max Scaler \Rightarrow



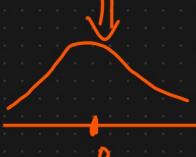
Advantages

① Improved Training Stability

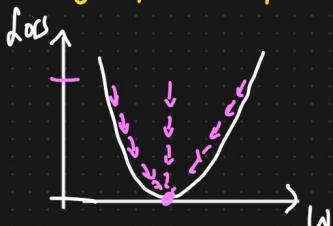


Vanishing and exploding Gradient problem

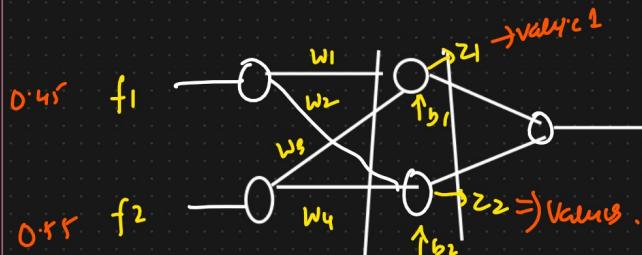
$$\left\{ \begin{array}{l} \mu=0 \\ \sigma=1 \end{array} \right.$$



② Faster Convergence



\Rightarrow Back propagation \Rightarrow Stable update.



f_1

Mouse size
0.45
0.60
-

f_2

Rooms
0.55
0.20
-

Batch Normalization

Normalizes the distribution

Price	Z_1	Z_2
45	-	-
-	-	-
-	-	-
-	-	-

$$Z_1 = \sigma \left[(0.45 \cdot w_1 + 0.55 \cdot w_3) + b_1 \right] = \text{Value 1} \quad \text{Does this distribution}$$

$$Z_2 = \sigma \left(\dots + b_2 \right) = \text{Value 2}. \quad \left. \begin{array}{l} \mu=0, \sigma=1 \\ \mu_1, \sigma_1 \\ \mu_2, \sigma_2 \end{array} \right\}$$

Batch Normalization VS Layer Normalization

f_1

f_2

Z_1 Z_2

μ_1, σ_1
μ_2, σ_2
μ_3, σ_3

$$Z_{\text{score}} = \frac{x_i - \mu_1}{\sigma_1}$$

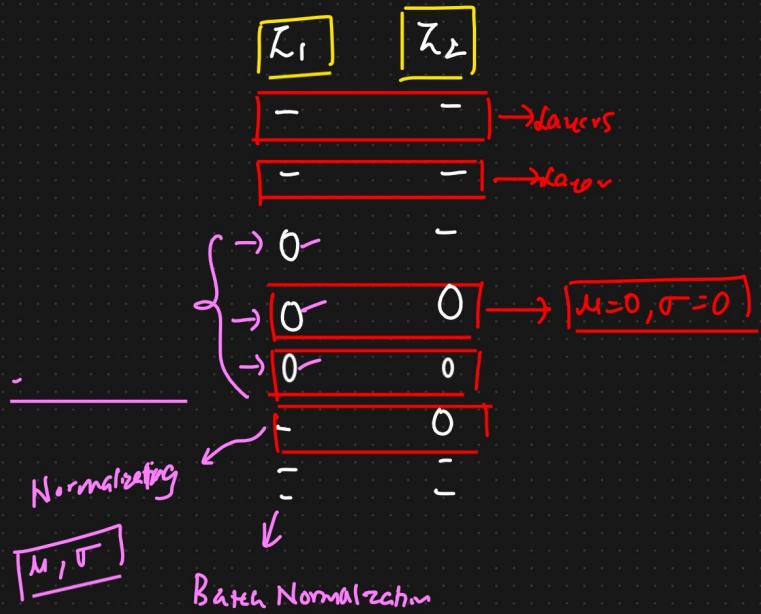
$$Z_{\text{score}} = \frac{x_i - \mu_2}{\sigma_2}$$

$$Z_{\text{score}} = \frac{x_i - \mu_3}{\sigma_3}$$

Layer

Normalization

$\gamma, \beta \rightarrow \text{learnable parameters}$



Normalization

\downarrow

γ, β

$$z_1 = \Gamma [w_1^T x + b_1]$$

$$y = \underline{\underline{\gamma}} \left[\frac{z_1 - \mu_1}{\sigma_1} \right] + \underline{\underline{\beta}}$$

Scale And Shift parameters

Normalized.

$$1) \text{ "CAT"} = [2.0, 4.0, 6.0, 8.0] \leftarrow \{ \text{Vectors} \}$$

$$2) \text{ Parameters } = \gamma = [1.0, 1.0, 1.0, 1.0] \rightarrow \text{Learned Scale}$$

$$\beta = [0.0, 0.0, 0.0, 0.0] \rightarrow \text{Shift}$$

\Rightarrow Scale And Shift param

i) Compute the mean

$$z_{score} = \frac{x_i - \mu}{\sigma}$$

$$\mu = \frac{1}{4} (2.0 + 4.0 + 6.0 + 8.0)$$

$$= \frac{20.0}{4} = 5.0$$

ii) Compute the variance (σ^2)

$$\sigma^2 = \frac{1}{4} \left[(2.0 - 5.0)^2 + (4.0 - 5.0)^2 + (6.0 - 5.0)^2 + (8.0 - 5.0)^2 \right] = 5.0$$

iii) Normalize the i/p

$$\hat{x}_i = \frac{x_i - \mu}{\sqrt{\sigma^2 + \epsilon}} \quad \epsilon = 1e^{-5} \Rightarrow \text{Avoid division by 0}$$

$$\sqrt{\sigma^2 + \epsilon} = \sqrt{5.0 + 1e^{-5}} \approx \sqrt{5.00001} = 2.236$$

$$\hat{x}_1 = \frac{2.0 - 5.0}{2.236} \approx -1.34 \quad \underline{\text{Normalized vector}}$$

$$\hat{x}_2 = \frac{4.0 - 5.0}{2.36} \approx -0.45 \quad \hat{x} = [-1.34, -0.45, 0.45, 1.34]$$

$$\gamma_3 = \frac{6.0 - 5.0}{2.236} \approx 0.45$$

$$\chi_4 = \frac{80 - 5.0}{2.236} \approx 1.34.$$

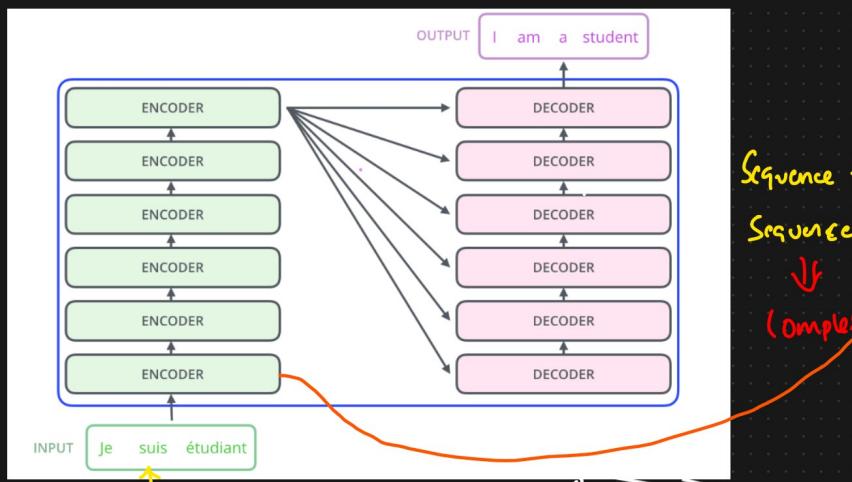
4) Scale And Shift

$$Y = \begin{bmatrix} 1.0, 1.0, 1.0, 1.0 \end{bmatrix} \quad J^3 = \begin{bmatrix} 0.0, 0.0, 0.0, 0.0 \end{bmatrix}.$$

$$y_i = \gamma_i \hat{x}_i + \beta_i$$

$$y = [-1.34, -0.45, 0.45, 1.34].$$

④ Encoder Architecture [Research Paper]



The diagram illustrates the Transformer architecture's processing flow:

- Input Sequence:** J/P Sequence is processed by **Text Embeddings + Positional Encoding**.
- Multi-Head Attention:** The embedded sequence is processed by Multi-Head Attention, resulting in outputs Z_1, Z_2, Z_3, Z_4, Z_5 . A residual connection (indicated by a skip connection line) adds the original input embeddings to the Multi-Head Attention output.
- (Add And Norm):** The output of Multi-Head Attention is passed through a layer labeled **(Add And Norm)**, which performs residual connection and layer normalization.
- Feed Forward NN:** The normalized output is processed by a Feed Forward NN, resulting in the final output sequence $\text{Every word} = \underline{512}$.

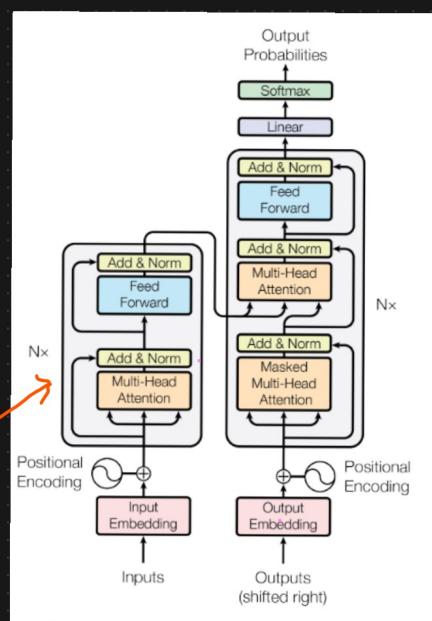
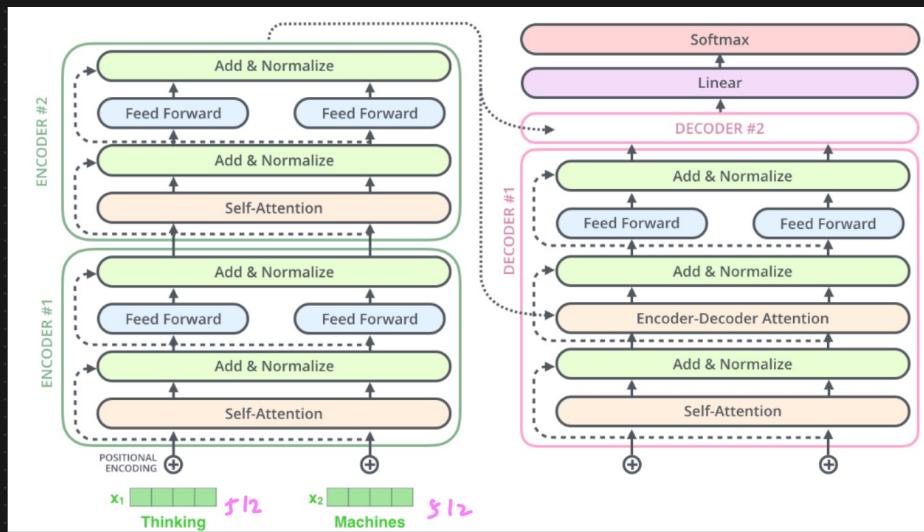
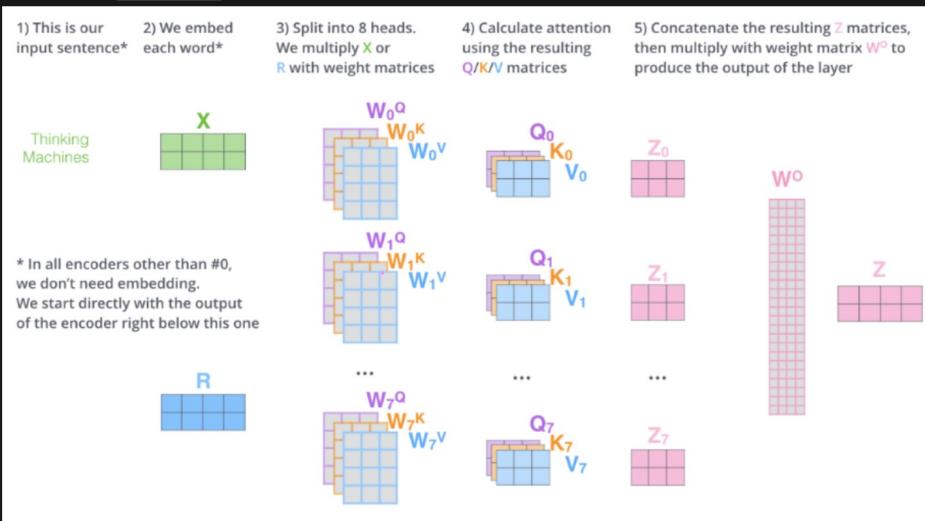


Figure 1: The Transformer - model architecture.

Self Attention to Feed Forward Neural N/W



Embedding Vector: 512

$Q, K, V = 64$

Head Attention: 8

① Residual connection : Skip connection NN.

i) Addressing the Vanishing Gradient Problem

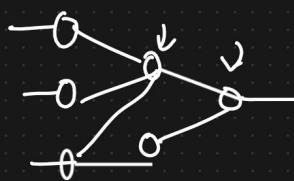
Residuals : Residual connection create a short paths for gradients to flow directly through the n/w. Gradient remains sufficiently large.

2) Improve Gradient flow

Convergence will be faster.

3) Enables Training of Deeper Networks.

④ Feed Forward NN



dinner function problem

{Non linear functions}.

① Adding Non linearity

② Processing Each position Independently.

Self Attention → (capture relationships)

FFN → Each token representation independently.



Transforming those representation further

and allows the model to learn

Richer Representation.

A^{NN} =>

③ FFN → Depth => Adds Depth to the Model.

Depth ↑↑ => More learnings → DATA

⑤ Decoders In Transformers

3 main Components

The transformer decoder is responsible for generating the output sequence one token at a time, using the encoder's output and the previously generated tokens.

① Masked Multi Head Self Attention ✓.

② Multi Head Attention (Encoder Decoder Attention)

③ Feed Forward Neural Network.

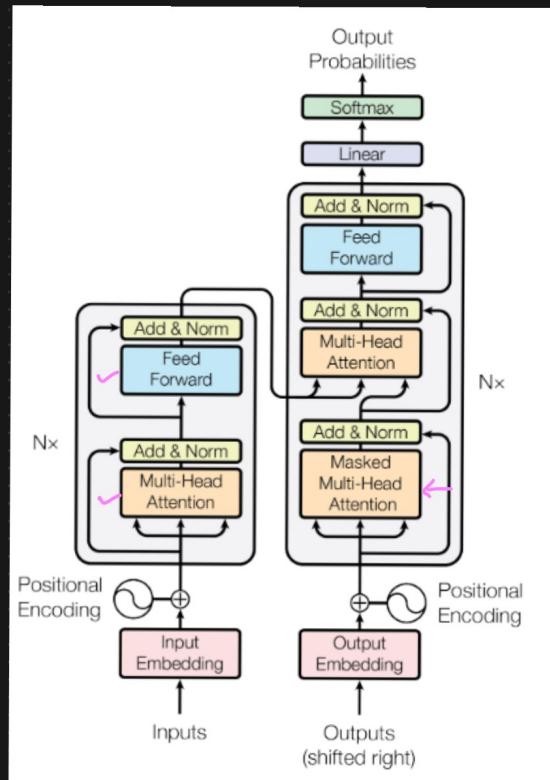
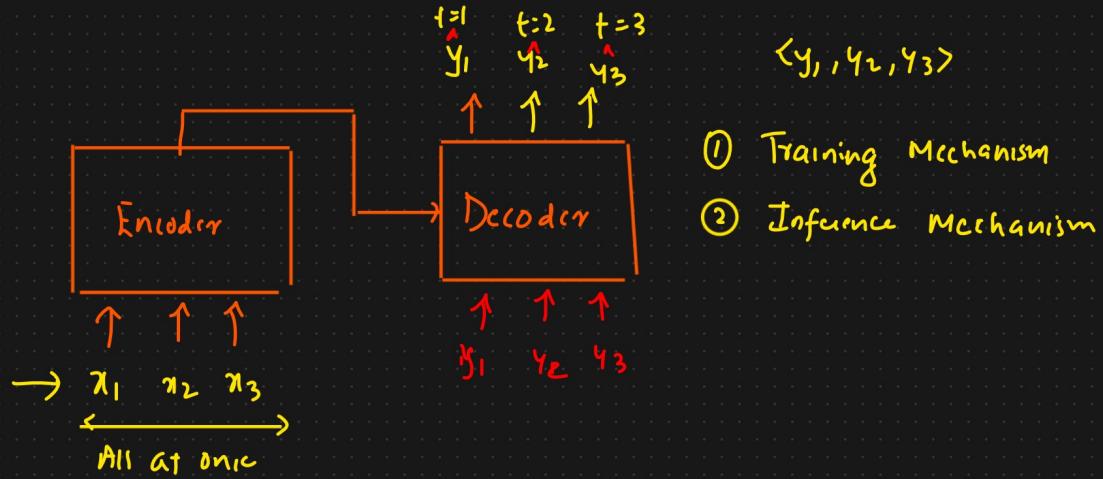


Figure 1: The Transformer - model architecture.



Masked Multi Head Self Attention

- ① I/P Embedding And Positional Embedding ✓ → Two padding → Sequence length equal
- ② Linear Projection for Q,K,V
- ③ Scaled Dot Product Attention
- ④ Mask Application } => Try to understand the imp.
- ⑤ Multi Head Attention
- ⑥ Concatenation And Final linear projection
- ⑦ Residual connection And Layer Normalization

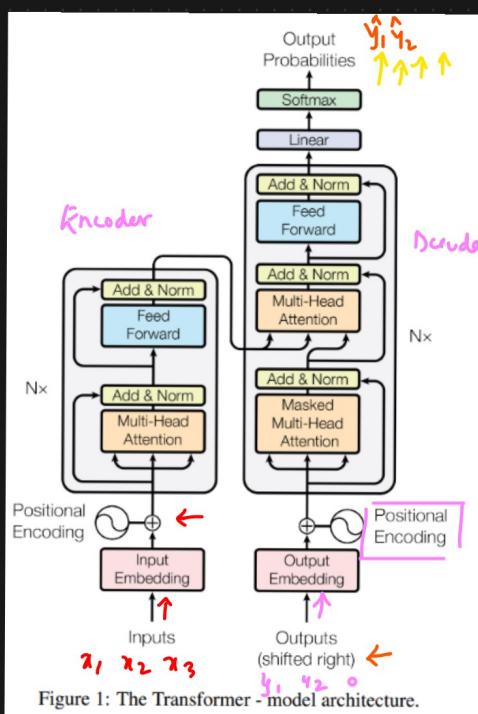
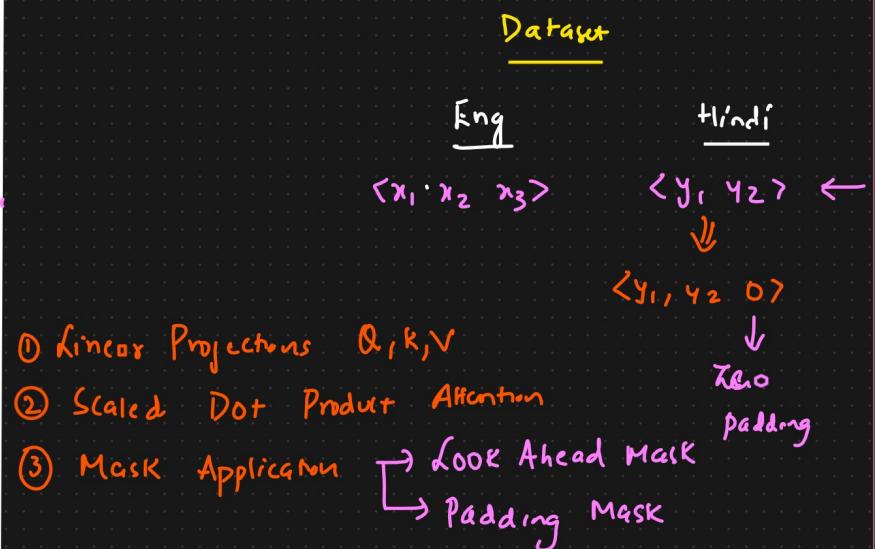


Figure 1: The Transformer -model architecture.



Masked
Multi Head
Attention.

$$\begin{array}{c} \text{I/P} \\ [4 \ 5 \ 6 \ 7] \\ \downarrow \\ [1 \ 2 \ 3 \ 0] \end{array} \quad \begin{array}{c} \text{O/P} \\ [1 \ 2 \ 3] \end{array}$$

i) Input Embedding and Positional Encoding

Output Embedding [STEP 1]

$$\begin{bmatrix} [0.1, 0.2, 0.3, 0.4], \\ [0.5, 0.6, 0.7, 0.8], \\ [0.9, 1.0, 1.1, 1.2], \\ [0.0, 0.0, 0.0, 0.0] \end{bmatrix}$$

Step 2 : Linear Projection for Q, K, and V.

$$WQ = WK = WV = I$$

Create query (Q), key (K) and value (V) vectors

Q = Output Embedding $\times W_Q$ = Output Embedding

K = " $\times W_K$ = "

V = " $\times W_V$ = "

$$\begin{aligned} Q = K = V &= \begin{bmatrix} [0.1, 0.2, 0.3, 0.4], \\ [0.5, 0.6, 0.7, 0.8], \\ [0.9, 1.0, 1.1, 1.2], \\ [0.0, 0.0, 0.0, 0.0] \end{bmatrix}, \quad \begin{bmatrix} 0.1 \\ 0.2 \\ 0.3 \\ 0.4 \end{bmatrix} \times \begin{bmatrix} 0.5 \\ 0.6 \\ 0.7 \\ 0.8 \end{bmatrix}^T \cdot \begin{bmatrix} 0.9 \\ 1.0 \\ 1.1 \\ 1.2 \end{bmatrix} = \begin{bmatrix} 0 \\ 0 \\ 0 \\ 0 \end{bmatrix} \end{aligned}$$

③ Scaled Dot Product Attention Calculation

$$\text{Scores} = Q \times K^T / \sqrt{d_K}$$

$$= Q \times K^T / 2$$

$$\begin{bmatrix} 0.1 \times 0.1 + 0.2 \times 0.2 + 0.3 \times 0.3 + 0.4 \times 0.4, \\ 0.1 \times 0.5 + 0.2 \times 0.6 + 0.3 \times 0.7 + 0.4 \times 0.8, \\ 0.1 \times 0.9 + 0.2 \times 1.0 + 0.3 \times 1.1 + 0.4 \times 1.2 \\ 0.1 \times 0 + 0.2 \times 0 + 0.3 \times 0 + 0.4 \times 0 \end{bmatrix}$$

Score : $\left[\begin{bmatrix} 0.3, 0.7, 1.1, 0.0 \\ 0.7, 1.9, 3.1, 0.0 \\ 1.1, 3.1, 5.1, 0.0 \\ 0.0, 0.0, 0.0, 0.0 \end{bmatrix} \right]$

④ Masked Application

It helps manage the structure of the sequences being processed and ensures the model's behavior is correct during training and inference.

Reasons

① Handling Variable length Sequences with Padding MASK

Purpose

- ① To handle sequences of different length in batch
- ② To ensure that padding tokens, which are added to make sequences of uniform length, do not affect the model prediction

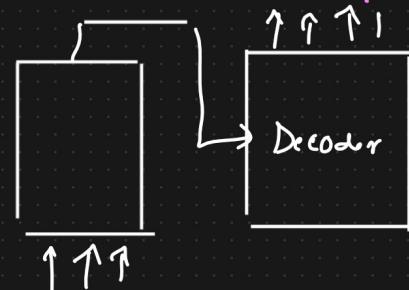
Eg: $\text{inp} \leftarrow \text{Sequence 1} [1, 2, 3] \rightarrow \text{Op} [y_1, y_2, 0, 0, 0]$

$\text{inp} \leftarrow \text{Sequence 2} [4, 5, \boxed{0}]$ 0 is the padding token
 → Influence the Attention Mechanism
 $\rightarrow 100.$ ↓
 lead to Incorrect or biased predictions.

A padding mask \Rightarrow The tokens are ignored.

Masking \rightarrow Padding Mask } } Padding Mask $\left[\begin{array}{ccc} 1 & 1 & 1 \\ 1 & 1 & 0 \end{array} \right]$
 \rightarrow Look Ahead Mask .

② Look Ahead Mask → Maintain Auto Regressive Property



How Anyo

② Sequence → Language Modelling, Translation

Eg: $[4, 5, 0] \rightarrow [1, 1, 0]$ → 1D MASK.

Attention ← token 1 attends to token 1, 2
Mechanism 1, 2

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

Convert 1D to 2D MASK

For each token in the sequence,
the mask should indicate
which tokens it can attend
to.

* Look Ahead Mask → Decoder Output

$$\left[\begin{bmatrix} 1 & 0 & 0 \\ 1 & 1 & 0 \\ 1 & 1 & 1 \end{bmatrix} \right]$$

A Combine Padding And Looking Ahead Mask

Element wise multiplication of 2 mask

Combine Mask = $\begin{bmatrix} [1, 0, 0] \\ [1, 1, 0] \\ [0, 0, 0] \end{bmatrix}$

④ Mask

Score : $\begin{bmatrix} [0.3, 0.7, 1.1, 0.0] \\ [0.7, 1.9, 3.1, 0.0] \\ [1.1, 3.1, 5.1, 0.0] \\ [0.0, 0.0, 0.0, 0.0] \end{bmatrix}$

Look Ahead Mask

$\begin{bmatrix} [1, 0, 0, 0] \\ [1, 1, 0, 0] \\ [1, 1, 1, 0] \\ [1, 1, 1, 1] \end{bmatrix}$

Padding Masking [extended to 2D Format]

$$\begin{bmatrix} [1, 1, 1, 0] \\ [1, 1, 1, 0] \\ [1, 1, 1, 0] \\ [0, 0, 0, 0] \end{bmatrix}$$

Combined Mask = Look Ahead Mask + Padding Mask

$$\begin{bmatrix} [1*1, 0*1, 0*1, 0*0] \\ [1*1, 1*1, 0*1, 0*0] \\ [1*1, 1*1, 1*1, 0*0] \\ [1*1, 1*1, 1*1, 1*0] \end{bmatrix} = \begin{bmatrix} [1, 0, 0, 0] \\ [1, 1, 0, 0] \\ [1, 1, 1, 0] \\ [1, 1, 1, 0] \end{bmatrix}, \quad \begin{bmatrix} [1, -\infty, -\infty, -\infty] \\ [1, 1, -\infty, -\infty] \\ [1, 1, 1, -\infty] \\ [1, 1, 1, -\infty] \end{bmatrix}$$

Masked Score

$$\begin{bmatrix} [0.3, -\infty, -\infty, -\infty] \\ [0.7, 1.9, -\infty, -\infty] \\ [1.1, 3.1, 5.1, -\infty] \\ [0.0, 0.0, 0.0, -\infty] \end{bmatrix}$$

Zero out the influence when the Softmax is applied.

Attention weight.



Softmax

Softmax Score = Softmax (Masked Scores)

$$= \left[[1.0, 0.0, 0.0, 0.0], \right. \\ \left[0.3, 0.7, 0.0, 0.0 \right], \\ \left[0.1, 0.3, 0.6, 0.0 \right], \\ \left. [1.0, 0.0, 0.0, 0.0] \right]$$

Weight Sum of Value

Attention O/p = Softmax Scores * V.

Masking

Masking in the transformer architecture is essential for several reasons. It helps manage the structure of the sequences being processed and ensures the model behaves correctly during training and inference. Here are the key reasons for using masking:

1. Handling Variable-Length Sequences with Padding Mask

Purpose

To handle sequences of different lengths in a batch.

To ensure that padding tokens, which are added to make sequences of uniform length, do not affect the model's predictions.

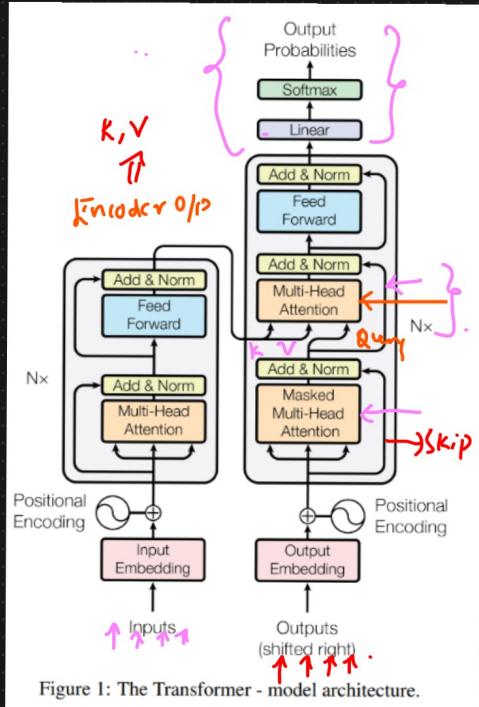
2. Maintaining Autoregressive Property with Look-Ahead Mask

Purpose

To ensure that each position in the decoder output sequence can only attend to previous positions and itself, but not future positions.

This is crucial for sequence generation tasks like language modeling and translation, where the model should not have access to future tokens when predicting the current token.

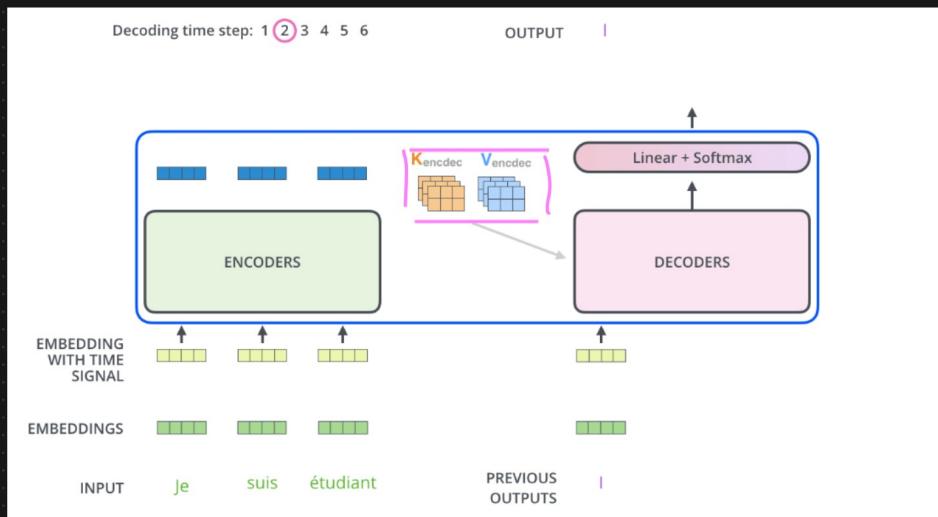
④ Encoder Decoder Multi Head Attention



- ① Encoder Out \rightarrow Set of Attention Vector K & V
 - ② Masked Multihed \rightarrow Attention Vector Q {Query Vector}

These are to be used by each decoder in its
"Encoder-decoder" attention layer
 \Downarrow
mention Helps the Decoder to focus on
 appropriate places in the i/p Sequence

$$\underline{Dg + \epsilon g f}$$



④ The Final Linear And Softmax Layer { Vectors → o/p Word }

{ BLOG } ⇒ Transformers

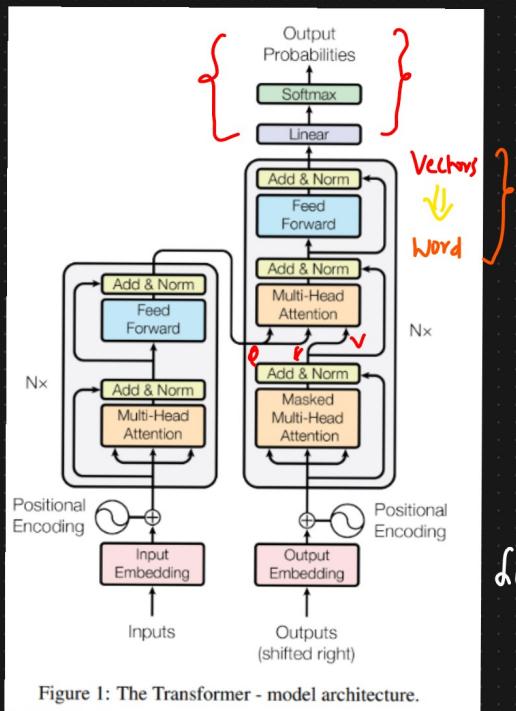
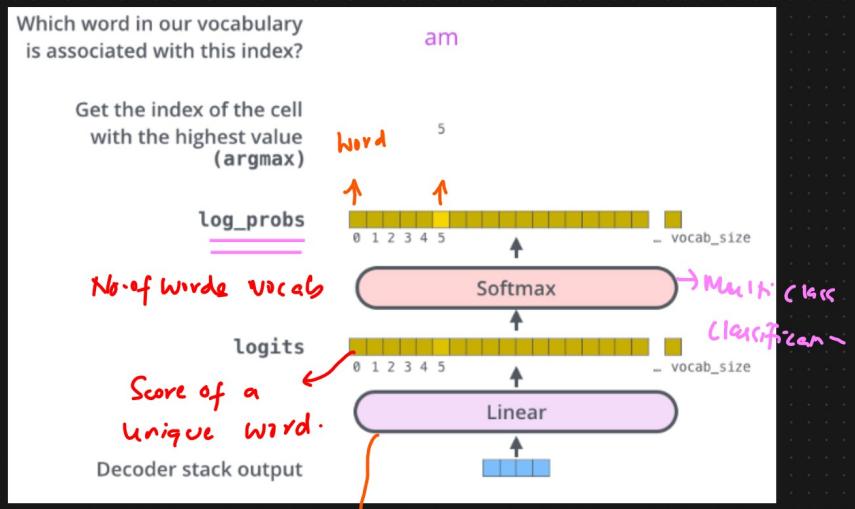


Figure 1: The Transformer - model architecture.



linear ⇒ The linear layer is a simple fully connected neural net that projects the vector produced by the stack of Decoder ⇒ logits vector ⇒

Model = 10,000 ⇒ Vocabulary ⇒ logits vector = 10000 cells wide

⑤ Softmax layer turns those scores into probabilities (all add upto 1.0).

The cell with the highest probability is chosen, and the word associated with it is produced as the o/p. = time stamp.

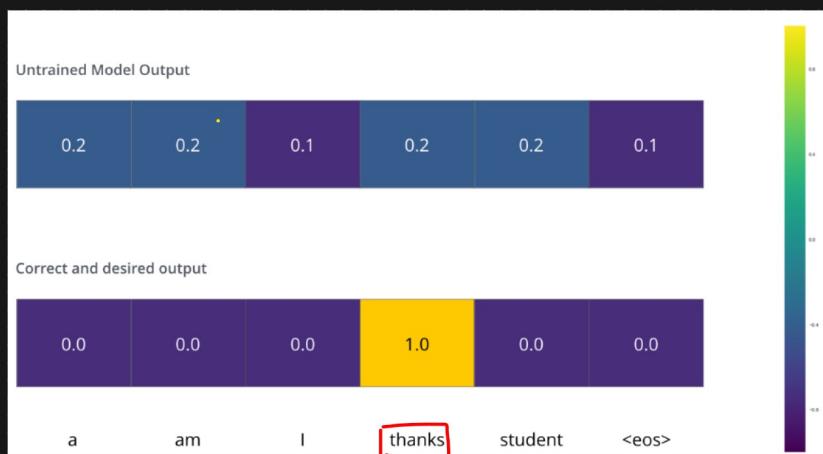
Recap of Training

Output Vocabulary						
WORD	a	am	I	thanks	student	<eos>
INDEX	0	1	2	3	4	5

Output Vocabulary						
WORD	a	am	I	thanks	student	<eos>
INDEX	0	1	2	3	4	5
One-hot encoding of the word "am"						
	0.0	1.0	0.0	0.0	0.0	0.0

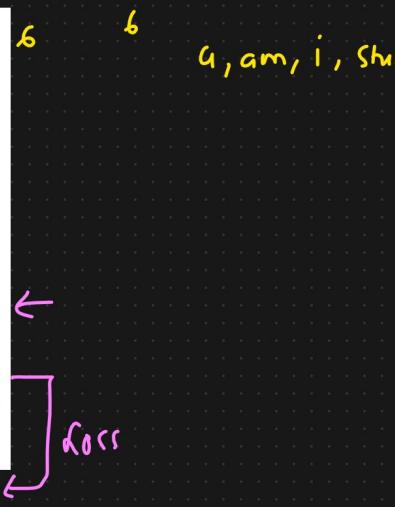
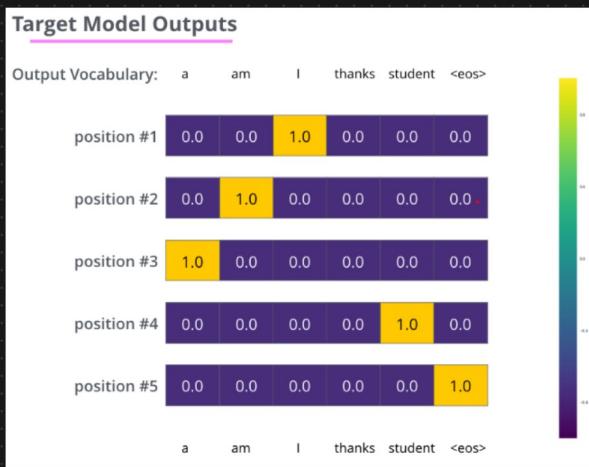
Merci → Thanks

I am a student <eos>.



⇒ Loss function ↓.

Back Propagation



D/P

I am a student
↓ ↓ ↓ ↓
OME OME OME OME

