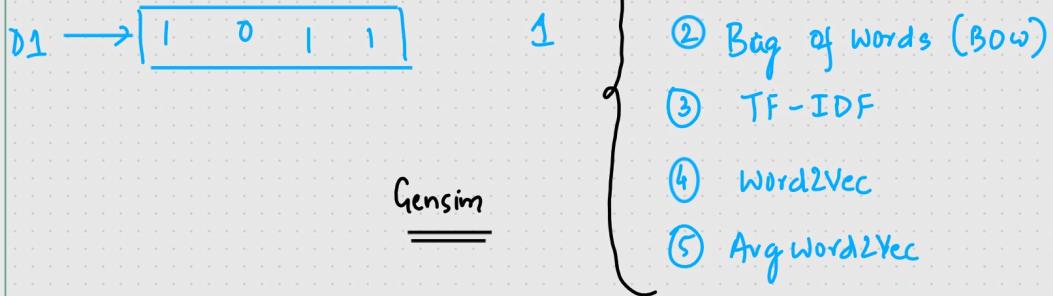
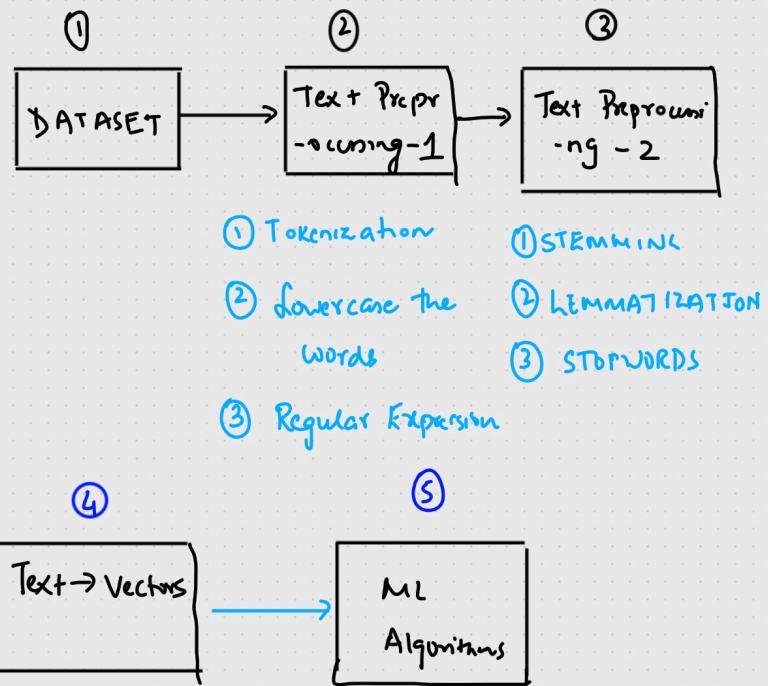


Text Preprocessing → What we have learnt?

Sentiment Analysis

	<u>Text</u>	<u>O/P</u>
D1	The food is good	1
D2	The food is bad	0
D3	Pizza is Amazing	1
D4	Burger is bad	0



① One Hot Encoding

Vocabulary size: 7

Vocabulary {unique words}

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

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

{50K vocabulary size}

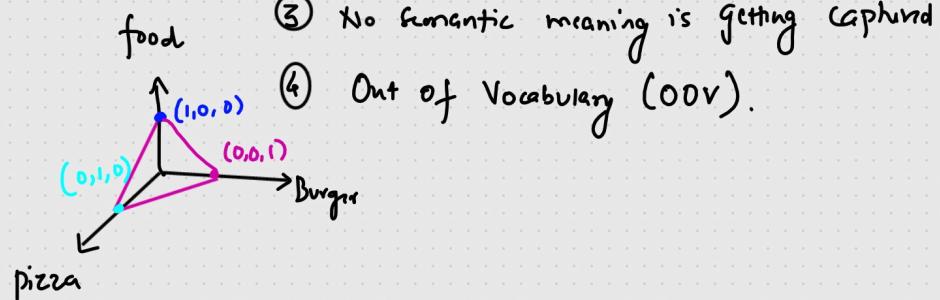
Disadvantages

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 IP

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	Lower all the words case	S1 → good boy
She is a good girl	1	⇒ Stopwords	S2 → good girl good
Boy and girl are good	1		S3 → Boy girl goods [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}

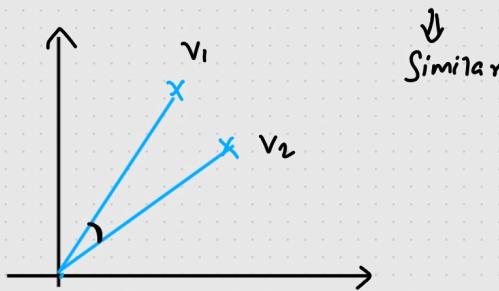
Advantages

- 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 $\rightarrow [1 \ 1 \ 1 \ 1 \ 1] \rightarrow v_2$



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

[food not good food good food not not good]
 S1 1 0 1 1 0 0 0 }
 S2 1 1 1 0 1 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]

$S_1 \rightarrow \text{good boy}$

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

$S_2 \rightarrow \text{good girl}$

$S_3 \rightarrow \text{boy girl good}$

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

Term Frequency * IDF

	S_1	S_2	S_3	Words	IDF
good	$\frac{1}{2}$	$\frac{1}{2}$	$\frac{1}{3}$	good	$\log_e(3/3) = 0$
boy	$\frac{1}{2}$	0	$\frac{1}{3}$	boy	$\log_e(3/2)$
girl	0	$\frac{1}{2}$	$\frac{1}{3}$	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

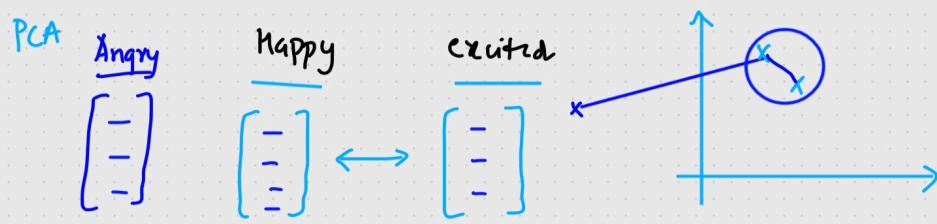
- ① Intuitive
- ② Fixed Size \rightarrow Vocab size
- ③ Word Importance is getting captured

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]

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
Royal	0.01	0.02	0.95	0.96	-0.02	0.02
Age	0.03	0.02	0.75	0.68	0.95	0.96
Food	-	-	-	-	0.91	0.92

300 dimension
num

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

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

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

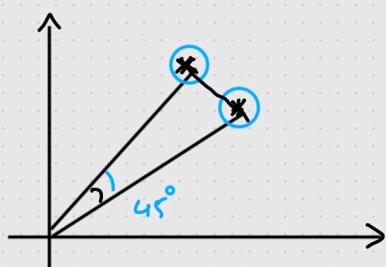
$$\text{Man} [0.95 \quad 0.98]$$

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

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

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

Cosine Similarity



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

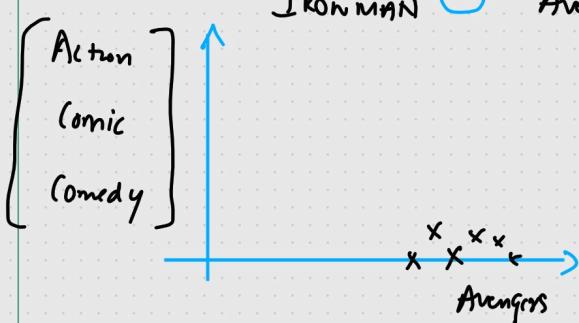
$$\text{Cosine-Sim} = \cos 45^\circ = \frac{1}{\sqrt{2}} = 0.7071$$

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

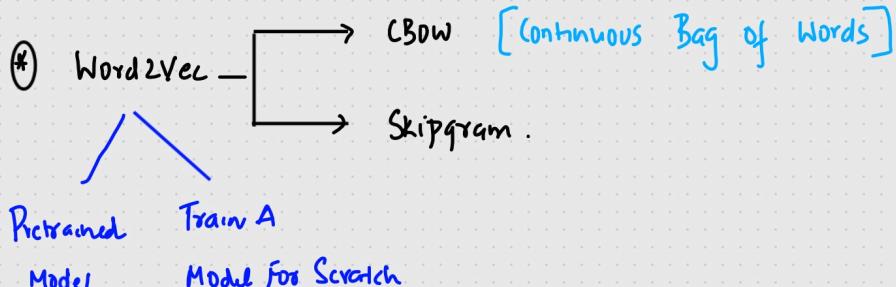
$$\boxed{1-1=0}$$

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

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

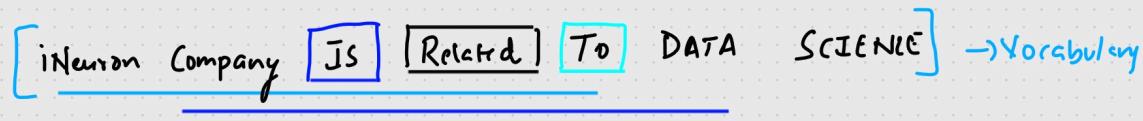


ANN, loss, Optimizers



① CBOW [Continuous Bag of Words]

CORPUS / DATASET



I/p

O/p

→ [iNeuron, Company, Related, To]

IS

iNeuron [1 0 0 0 0 0 0]

→ Company, IS, To, DATA

Related

Company [0 1 0 0 0 0 0]

→ IS, Related, DATA, SCIENCE

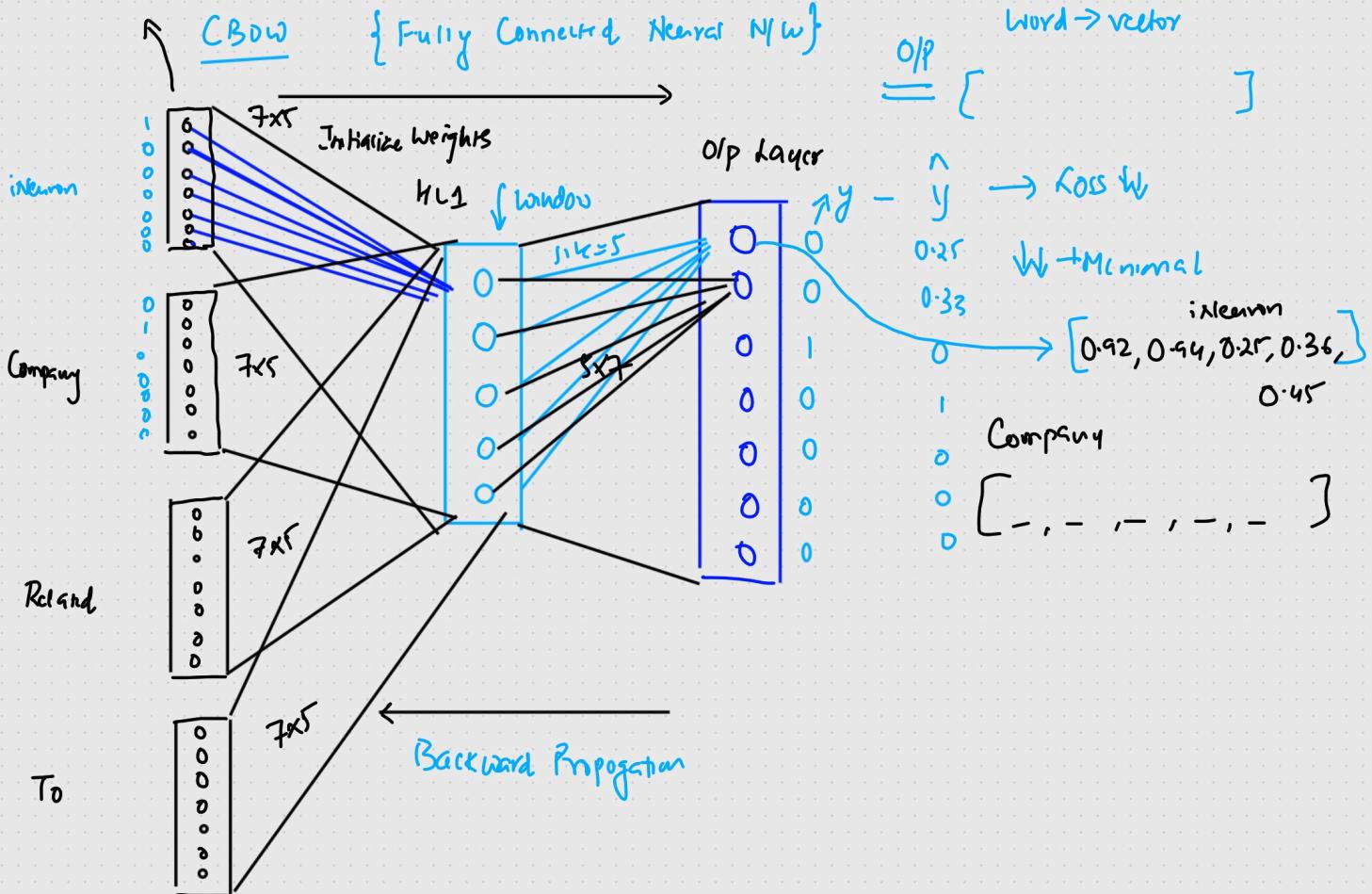
To

Related [0 0 0 1 0 0 0]

To [0 0 0 0 0 1 0 0]

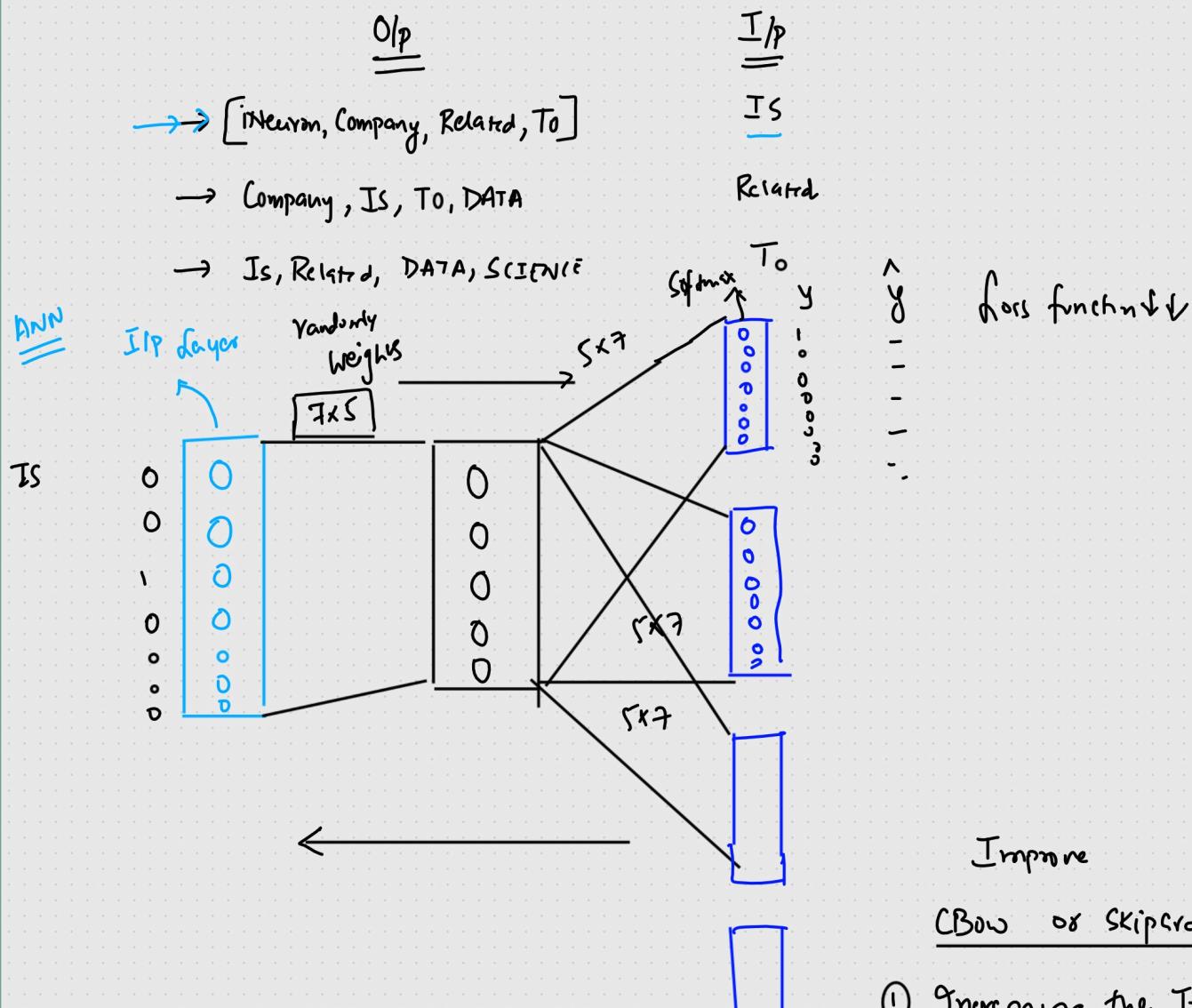
I/p Layer

ANN



② Skipgram - Word2Vec

Window Size=5



When Should we apply CBOW or SkipGram.

- ② Increase the window size

-vector dimension is also increasing.

{ Small Dataset → CBow
Huge Dataset → SkipGram }

Google Word2vec

Gunsim

3 billion words → Google News

feature representation of 300 dimension vectors]

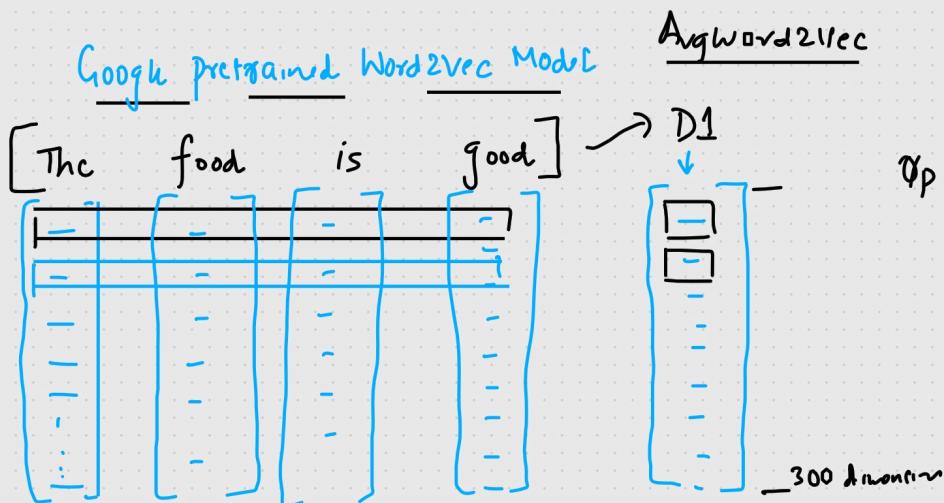
Cricket → [- . - - - - - - - -].

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	$\begin{bmatrix} \quad \quad \quad \quad \quad \quad \quad \quad \end{bmatrix}$	

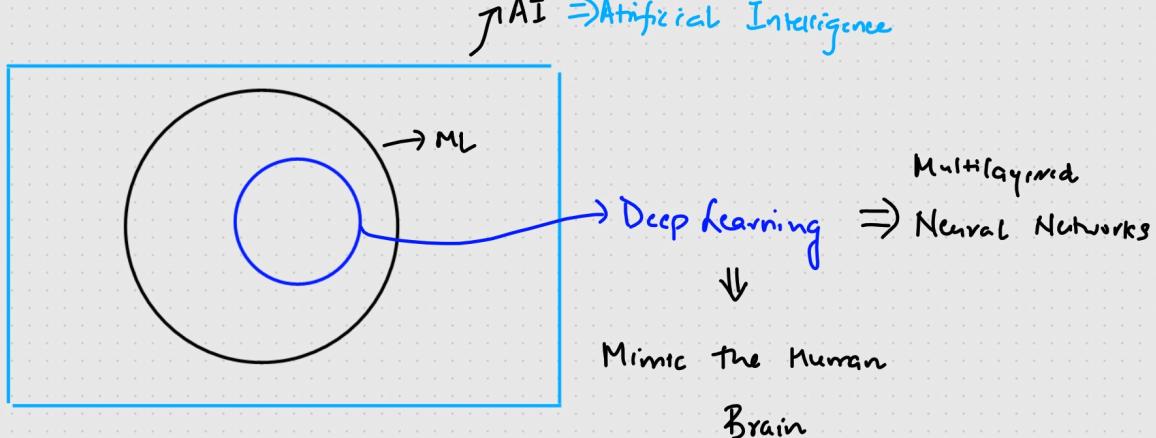


(*) Gensim, Glove

\hookrightarrow pretrained Google Word2Vec

↳ Train A Word2Vec From
Scratch

Deep learning



Deep learning

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

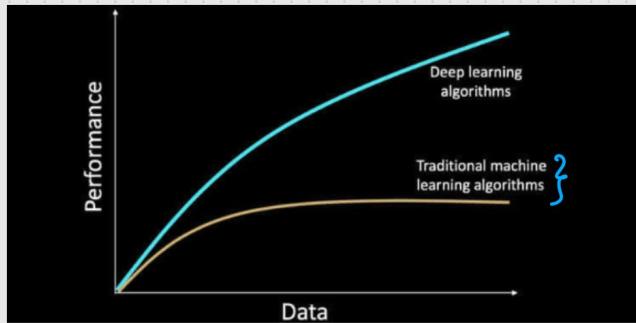
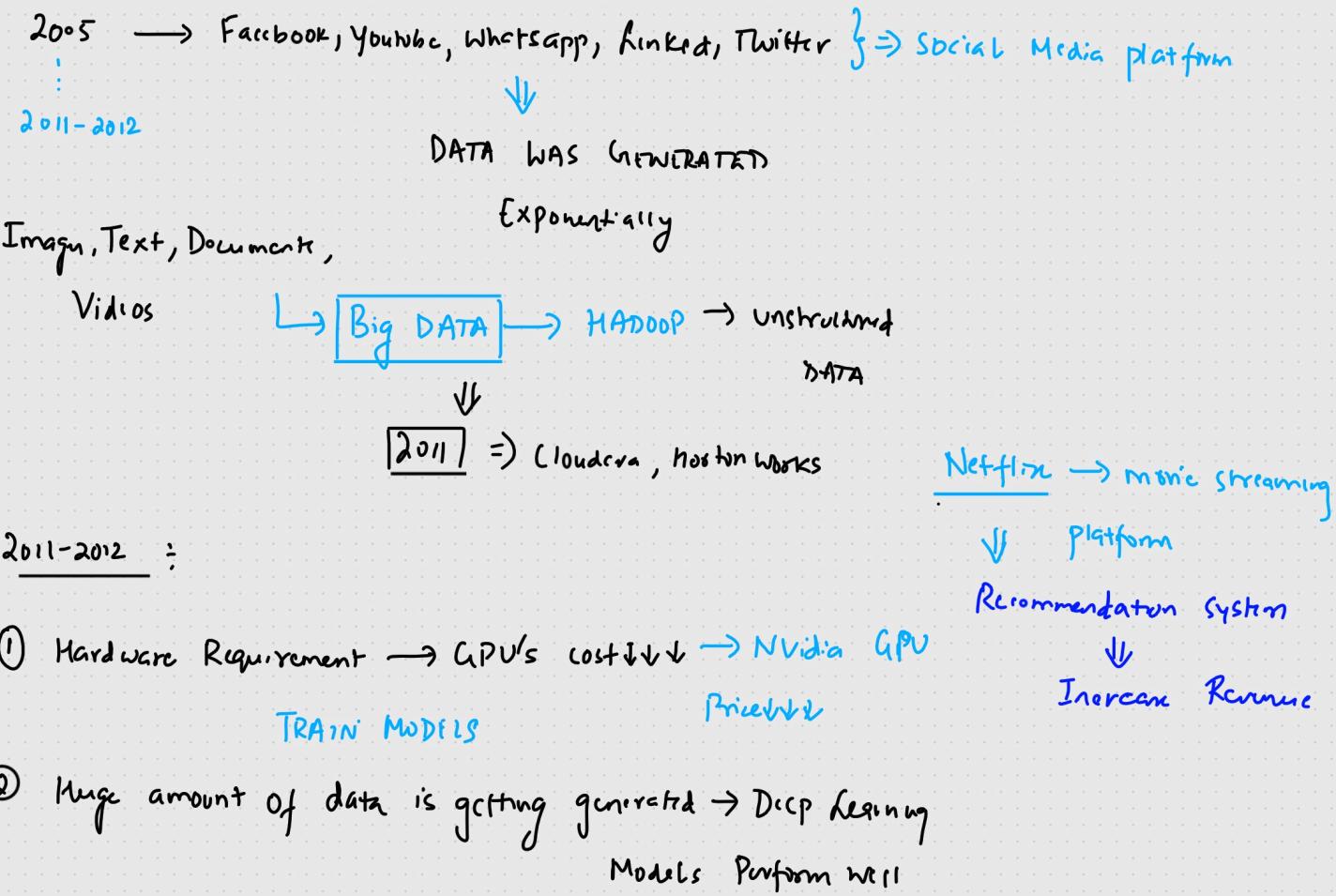
FRAMEWORK

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 ✓

[Single Layered NN]

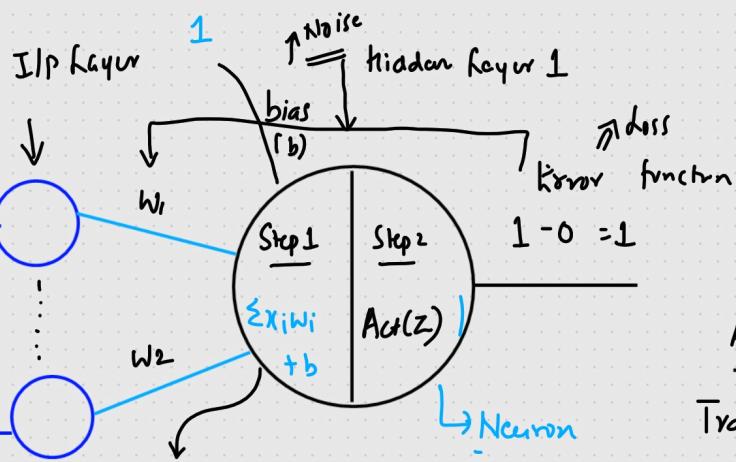
② Hidden layer ✓

③ Weights ✓

Binary classifier

④ Activation function ✓

DATASET



x ₁ ID	x ₂ No. of study hours	O/P Pers/Fac
→ 95	3	0
→ 110	4	1
→ 100	5	1

Activation function

Transform the output

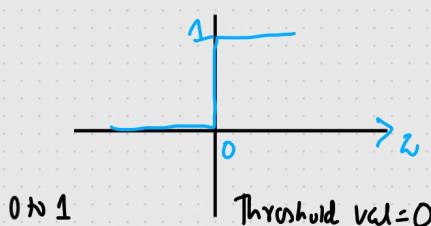
between 0 to 1

-1 to 1

$$z = w_1x_1 + w_2x_2 + b$$

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

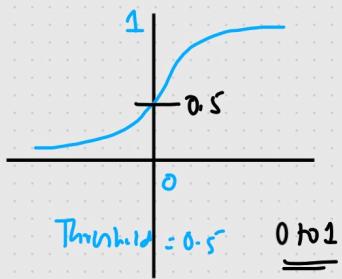
Step Function



0 to 1

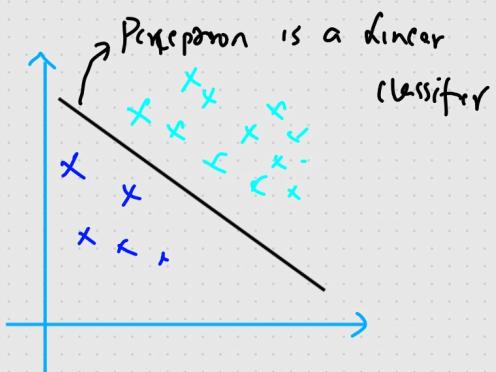
$$\begin{cases} 0 & z \leq 0 \\ 1 & z > 0 \end{cases}$$

Sigmoid Function



0 to 1

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



Step 1

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

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

$$y = mx + c$$

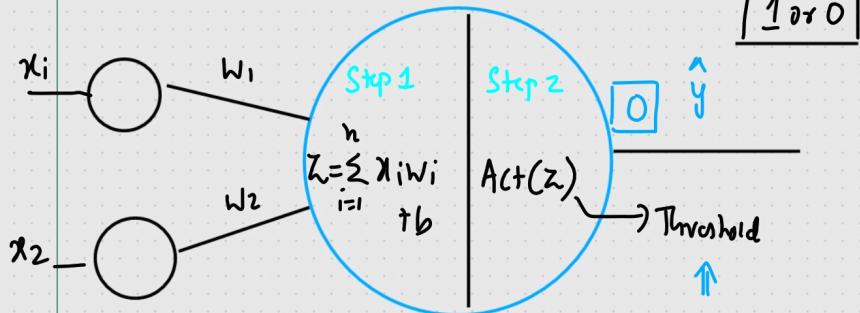
$$y = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \beta_3 x_3 + \dots + \beta_n x_n$$

Linear
problem
Statement

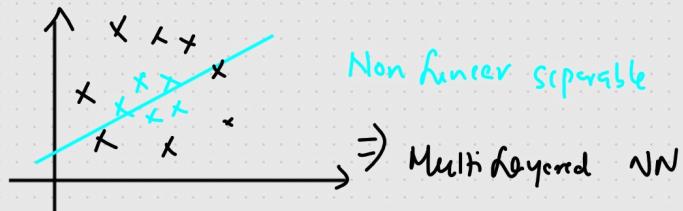
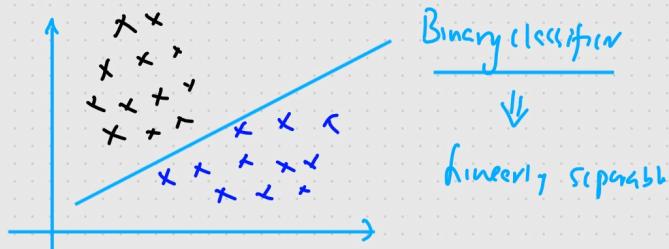
Perceptron Models

Single Layer Perceptron Model

Feed Forward Neural NW



Linear Separable problem



[ANN]

Multi Layered Perceptron Model

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

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

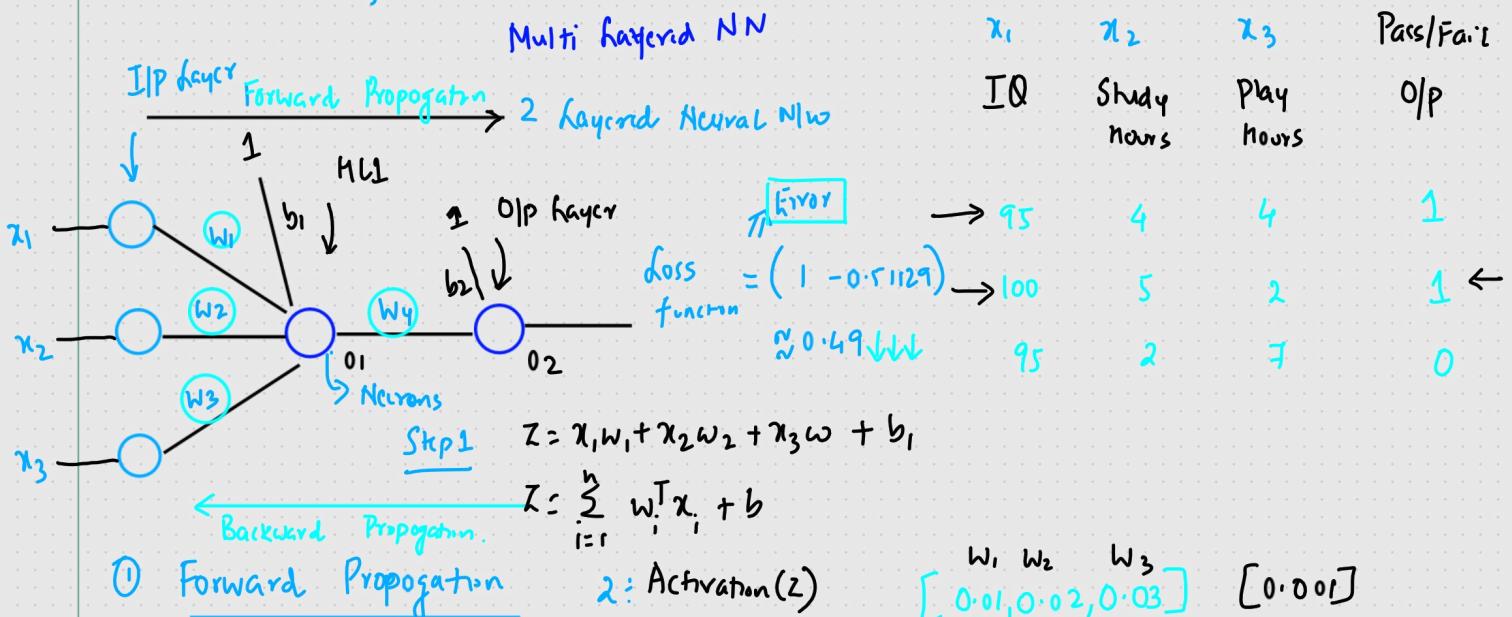
① Forward Propagation

② Backward Propagation → Geoffrey Hinton →

③ Loss function

④ Optimizers ✓

⑤ Activation function.



Hidden Layer 1

$$\text{Step 1: } Z = 95 \times 0.01 + 4 \times 0.02 + 4 \times 0.03 + 1 \times 0.01$$

$$= \underline{\underline{1.151}}$$

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

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

$$O_1 = 0.759$$

Hidden Layer 2

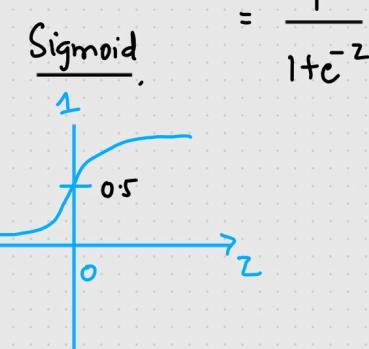
$$w_4 = 0.02$$

$$b_2 = 0.03$$

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

$$= 0.759 \times 0.02 + 0.03$$

$$= 0.04518$$

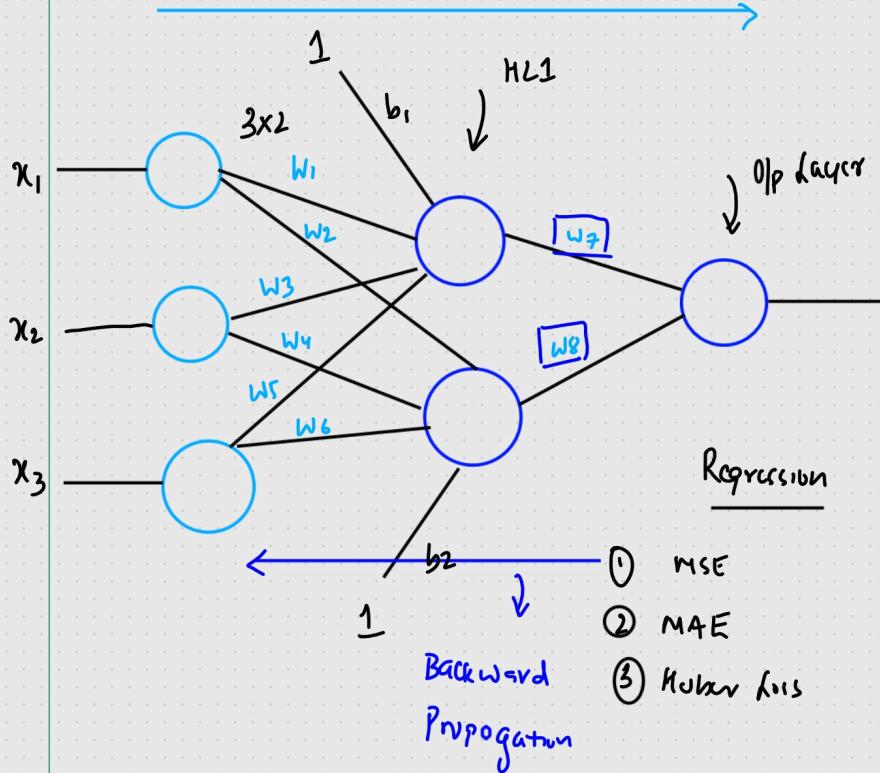


Step 2 : Activation(z)

$$\frac{1}{1+e^{-(0.04518)}} = 0.51129$$

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

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

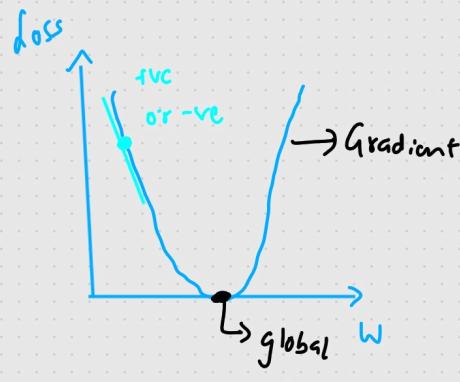
Weight update Formula

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

slope

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

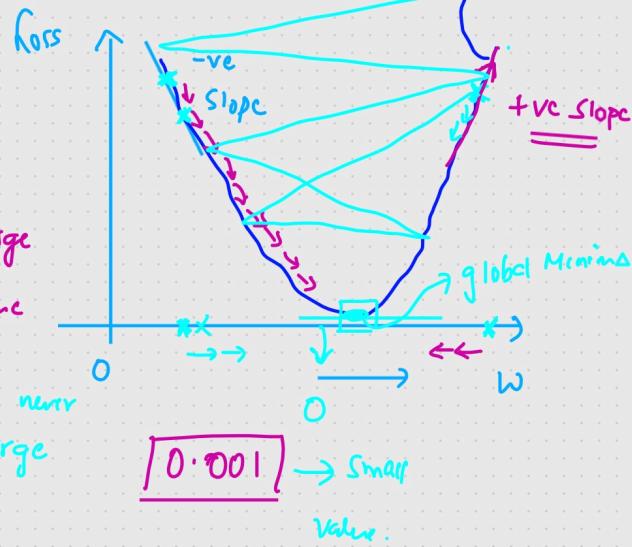
learning Rate



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

Gradient Descent
Optimizers

Optimizers : To reduce the loss value



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

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

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

↳ Learning Rate

η = large value

It may never converge

$$\boxed{0.001}$$

→ Small value.

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

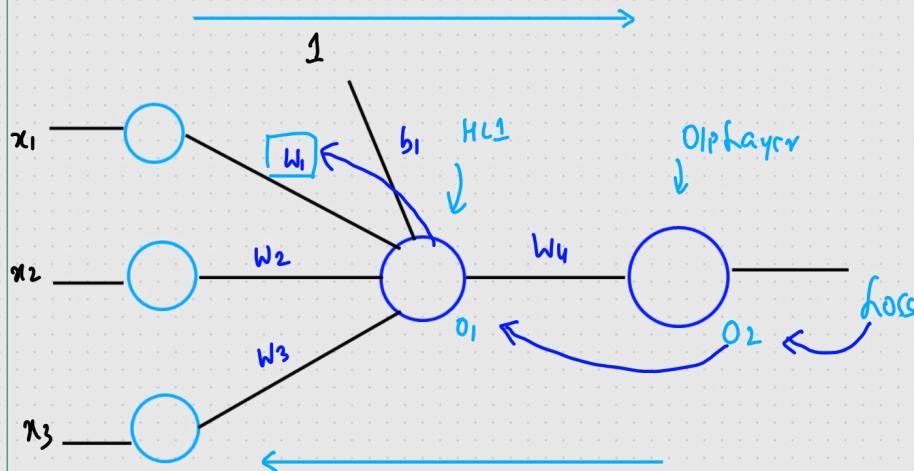
$$= W_{\text{old}} - \eta \quad (+\text{ve})$$

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

When W reaches Global Minima

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

⑦ Chain Rule of Derivative



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

$$W_{4,\text{new}} = W_{4,\text{old}} - \eta$$

$$\boxed{\frac{\partial h}{\partial W_{4,\text{old}}}}$$

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

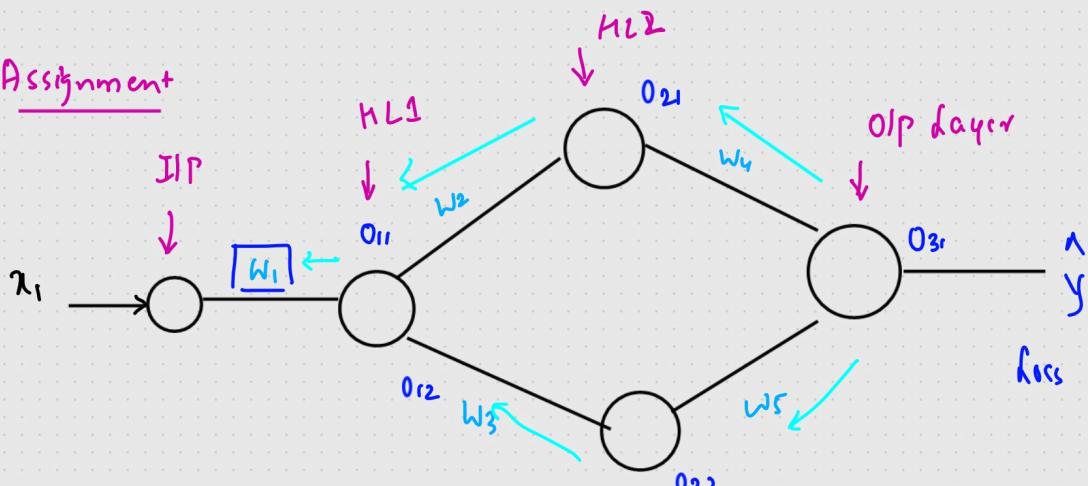
$$w_{1\text{new}} = w_{1\text{old}} - \eta \left[\frac{\partial h}{\partial w_{\text{old}}} \right]$$

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

$w_{2\text{new}}$

$w_{3\text{new}}$

Assignment



$$w_{1\text{new}} = w_{1\text{old}} - \eta \left[\frac{\partial h}{\partial w_{\text{old}}} \right]$$



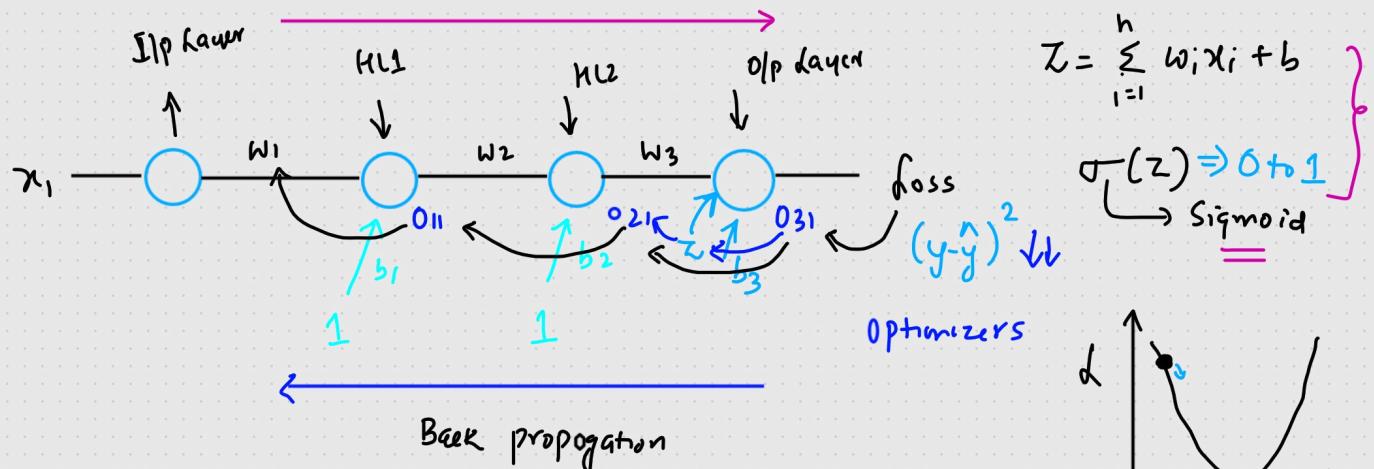
$$\frac{\partial h}{\partial w_{\text{old}}} = \left[\frac{\partial h}{\partial o_{31}} * \frac{\partial o_{31}}{\partial o_{21}} * \frac{\partial o_{21}}{\partial o_{11}} * \frac{\partial o_{11}}{\partial w_{\text{old}}} \right]$$

+

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

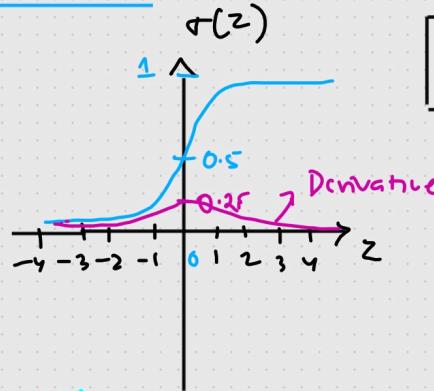
⑧ Vanishing Gradient Problem And Activation functions

$$\phi(z)$$



f) Sigmoid Activation function

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



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

$$\text{Derivative } (\sigma(z)) \\ 0 \leq \sigma'(z) \leq 0.25$$

Vanishing Gradient Problem.

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

$$\frac{\partial h}{\partial w_{1\text{ old}}} = \frac{\partial h}{\partial w_{1\text{ old}}} * \left[\frac{\partial o_{31}}{\partial o_{21}} * \frac{\partial o_{21}}{\partial o_{11}} + \frac{\partial o_{11}}{\partial w_{11}} \right] \Rightarrow \text{smaller value}$$

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

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

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

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

↓

Derivative of
Sigmoid

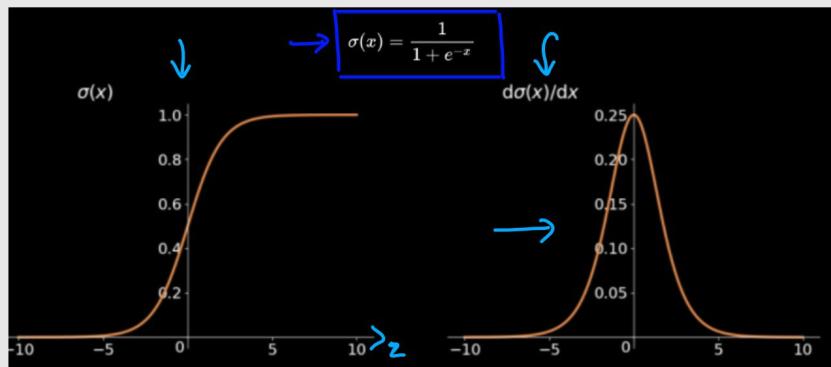
$$\frac{\partial o_{31}}{\partial o_{21}} = 0 \leq \sigma(z) \leq 0.25 \quad * \quad w_3_{\text{old}}$$

- ④ To fix this problem Researchers started exploring other Activation function
- ④ Tanh ④ ReLU ④ Prelu ④ Swiss

Activation Functions

① Sigmoid Activation function [0 to 1]

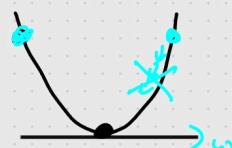
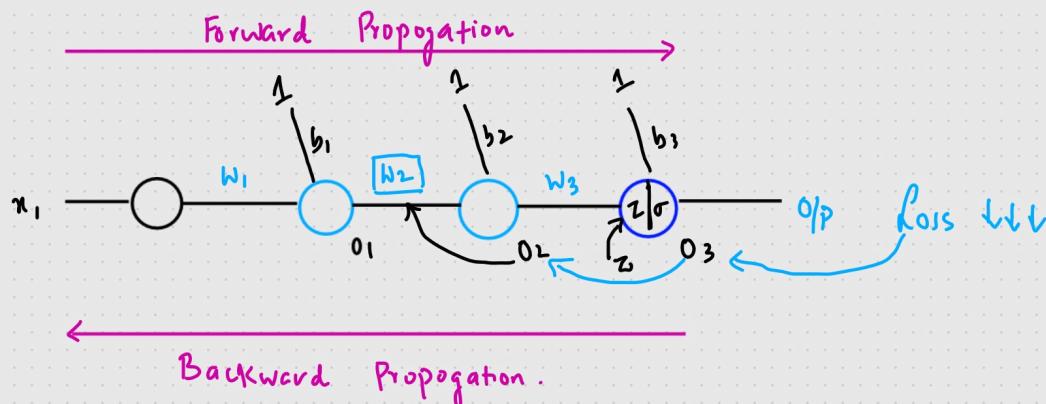
$$z = \sum_{i=1}^n w_i x_i + b$$



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

$$\phi(z) \Rightarrow$$

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



$$w_{2\text{new}} = w_{2\text{old}} - \eta \frac{\partial h}{\partial w_{2\text{old}}} \quad \begin{matrix} 0.001 \\ \uparrow \end{matrix} \quad \begin{matrix} 0.0001 \\ \uparrow \end{matrix} \quad \Rightarrow w_{2\text{new}} \approx w_{2\text{old}}$$

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

$0.20 \downarrow * 0.01 \times \quad \text{det } z = (o_2 * w_3) + b_3$

$$\frac{\partial o_3}{\partial o_2} = \frac{\partial (\sigma(z))}{\partial z} * \frac{\partial z}{\partial o_2} \quad [0 \text{ to } 1]$$

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

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

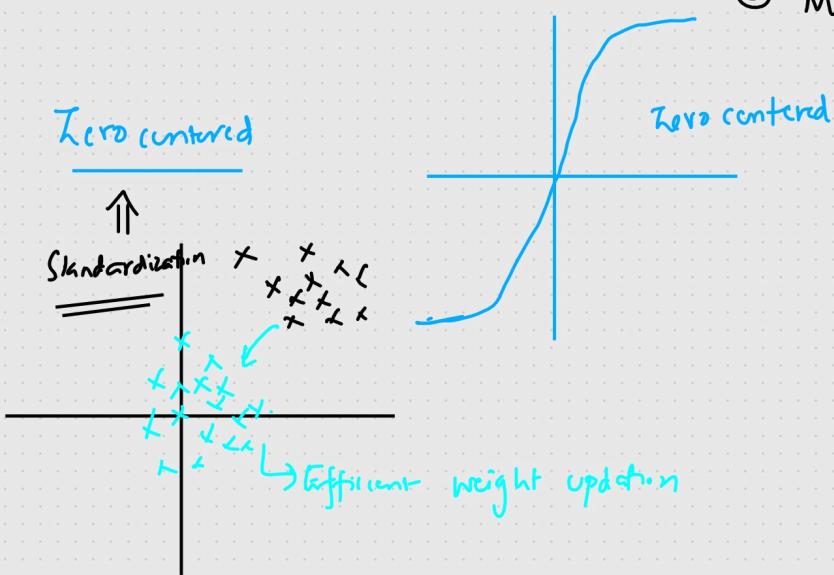
Advantages

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

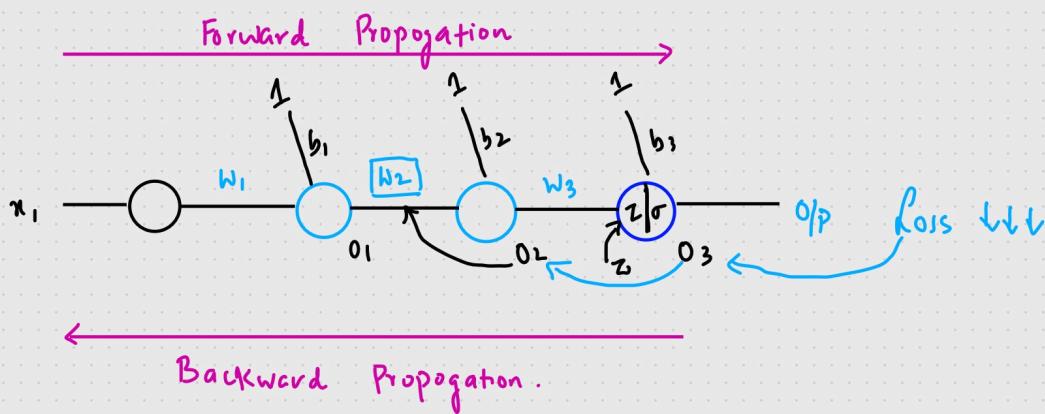
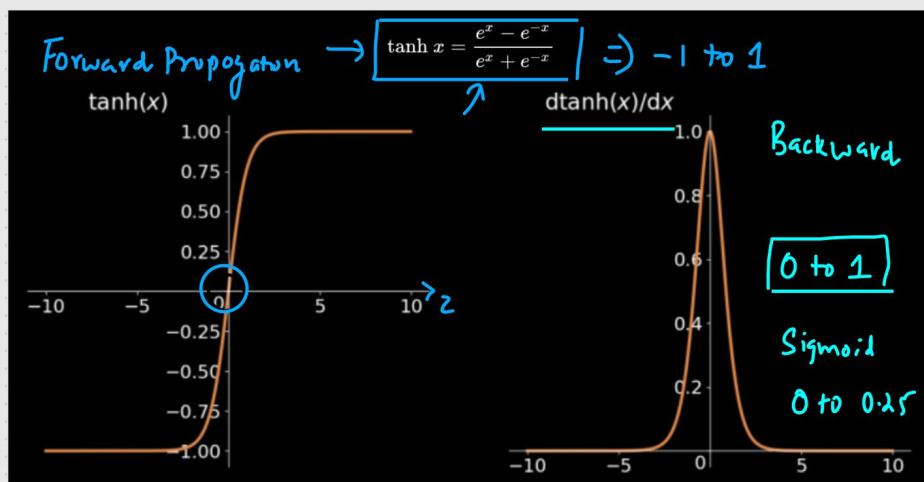
Disadvantages

- 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



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

↓↓↓

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

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

Advantages

① Zero Centric \Rightarrow Weight Updation is Efficient

Disadvantages

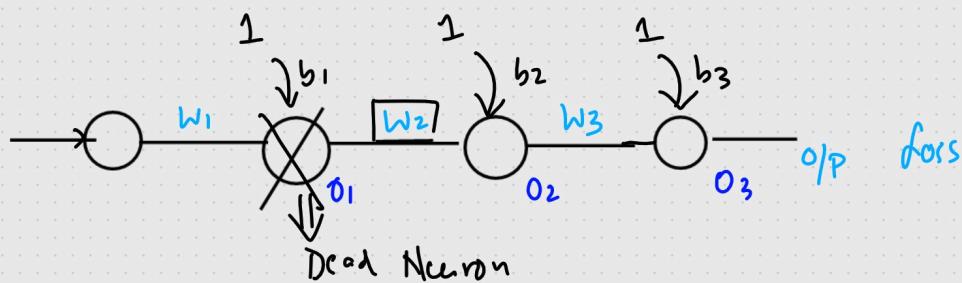
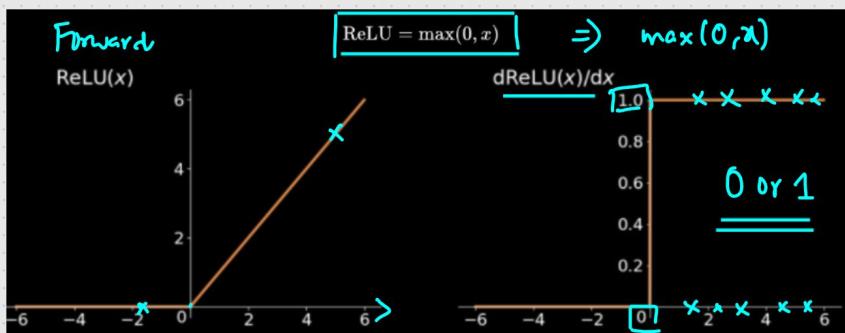
- ① Prone to Vanishing Gradient Problem
- ② Time complexity

[Rectified Linear Unit].

③ ReLU Activation Function

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

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



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

↓↓↓

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

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

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

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

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

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

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

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

Advantages

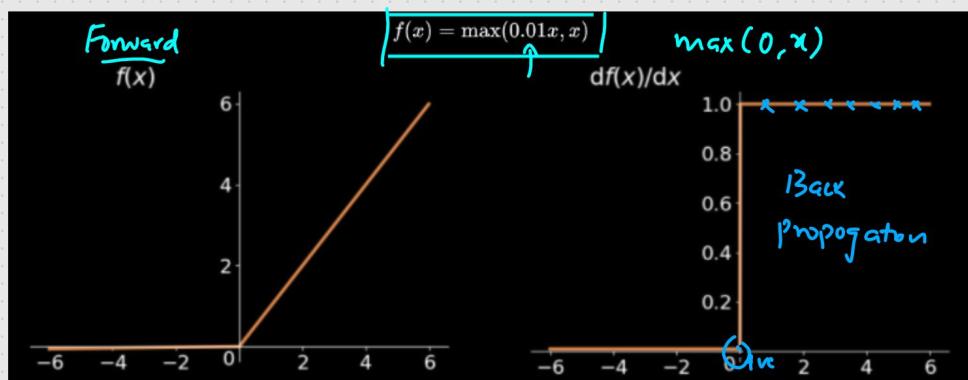
- ① Solving Vanishing Gradient Problem
- ② $\text{Max}(0, x) \rightarrow$ Calculation is Superfast. The ReLU function has a linear relationship.
- ③ It is much faster than Sigmoid or Tanh.

Disadvantages

- ① Dead Neuron
- ② ReLU function O/P $(0, x) \Rightarrow 0$ or +ve numbers
↓
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

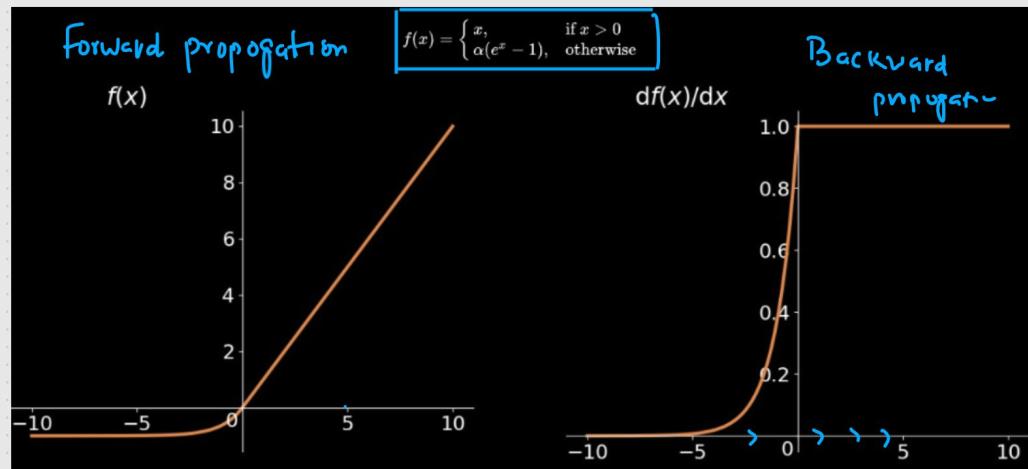
Advantages

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

Disadvantage

- ① It is not zero centric

⑤ ELU (Exponential Linear Units)



Advantages

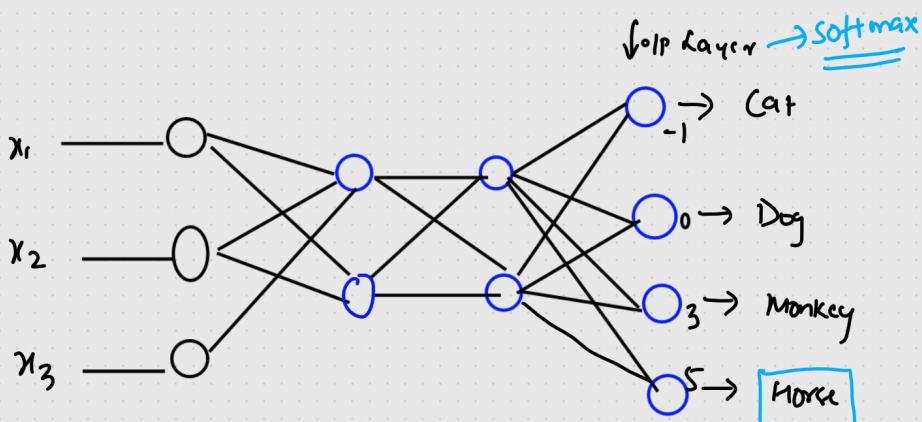
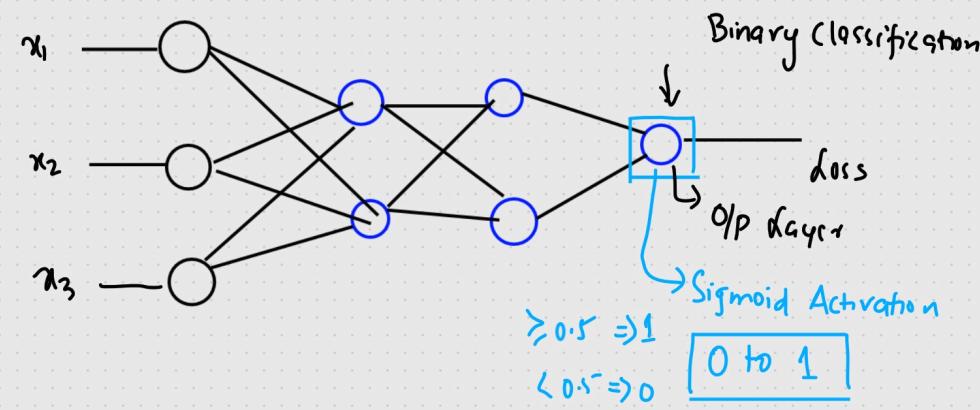
[It is used to solve
ReLU problems]

Disadvantage

- ① No Dead ReLU Issues
- ② Zero centered

i) Slightly more computationally intensive.

⑥ Softmax Activation function [Multiclass classification problem]



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

$$y_i = \mathbf{O} \cdot \mathbf{W} + b$$

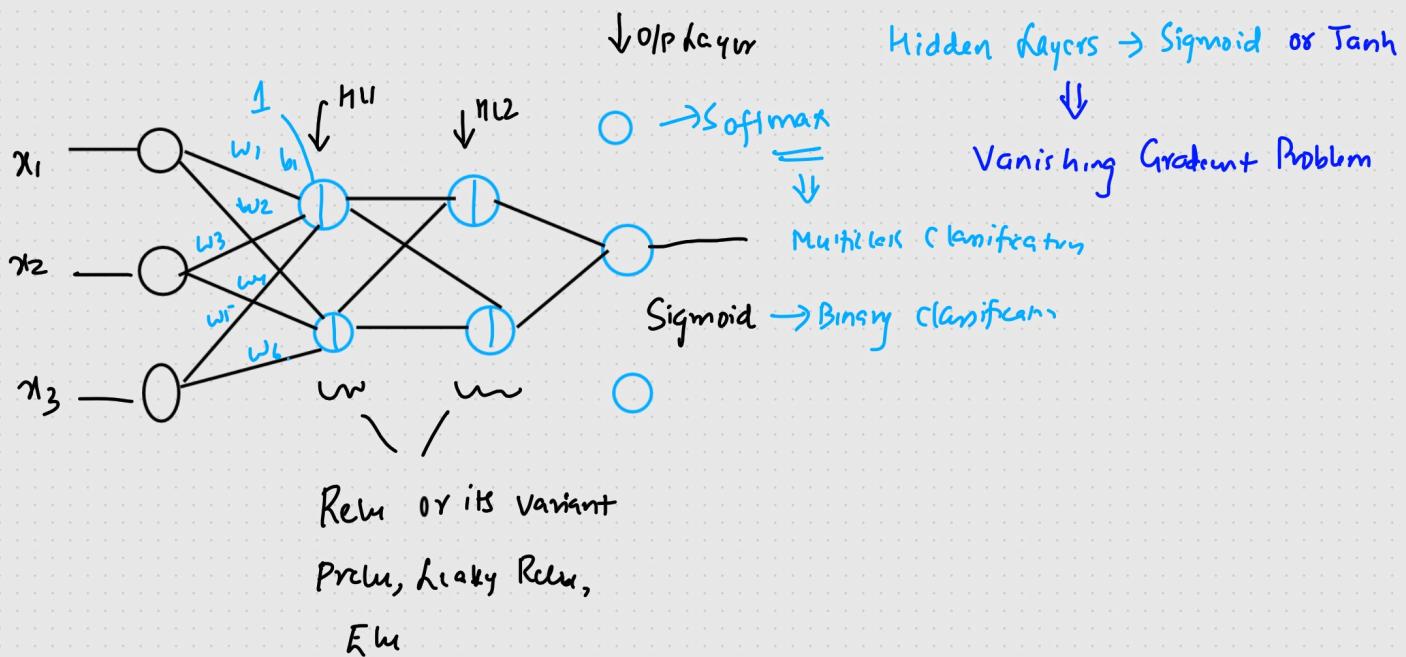
$$\text{Softmax} \Rightarrow \text{Cat} = \frac{e^{-1}}{e^{-1+0+3+5}} = 0.00033 \quad \Pr(\text{Horse}) = \frac{0.1353}{0.00033 + 0.0024 + 0.0183 + 0.1353}$$

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

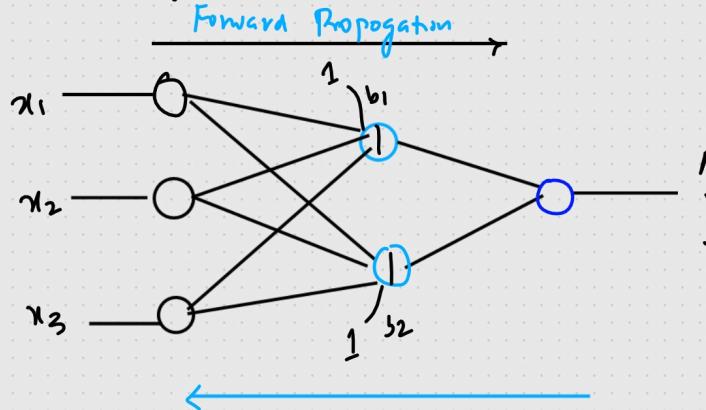
$$\text{Monkey} = \frac{e^3}{e^{-1+0+3+5}} = 0.0183$$

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

⑦ Which Activation Function To Use When?



Loss function And Cost Function



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

↑ Error

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

Optimizers

Gradient Descent

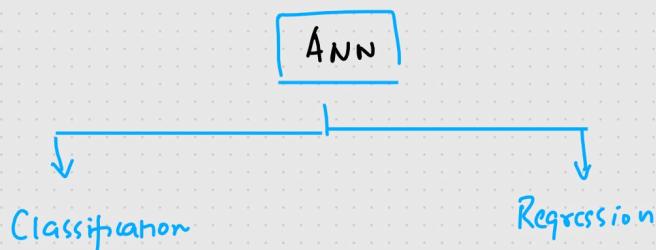
x ₁	x ₂	x ₃	O/P
-	-	-	0
-	-	-	1
-	-	-	0
-	-	-	1

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

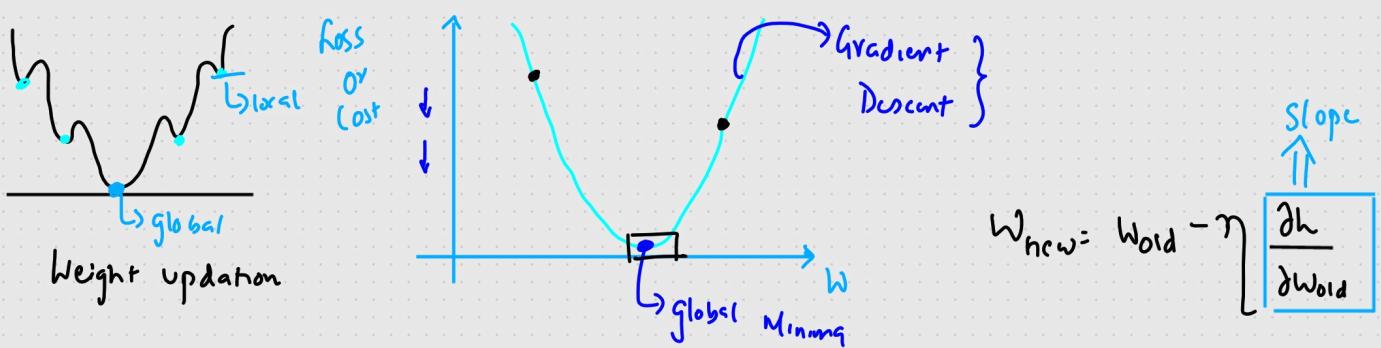
① Mean Squared Error (MSE)

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

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

↳ Quadratic Equations



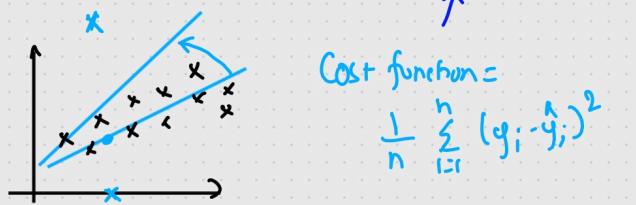


Advantages

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

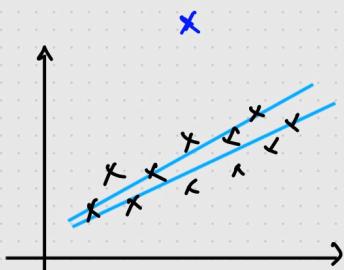
$$\text{Loss fn} = |y - \hat{y}|$$

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



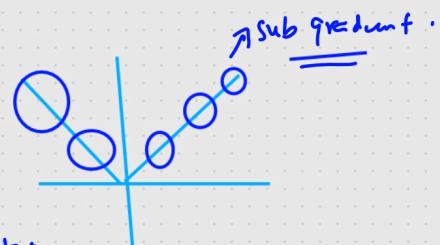
Advantages

- ① Robust to outliers



Disadvantages

- ① Convergence usually takes time in MAE



③ Huber loss

$$\text{① MSE}$$

$$\text{② MAE}$$

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

↑ No outlier ↑ Hyperparameter

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

④ RMSE (Root Mean Squared Error)

=

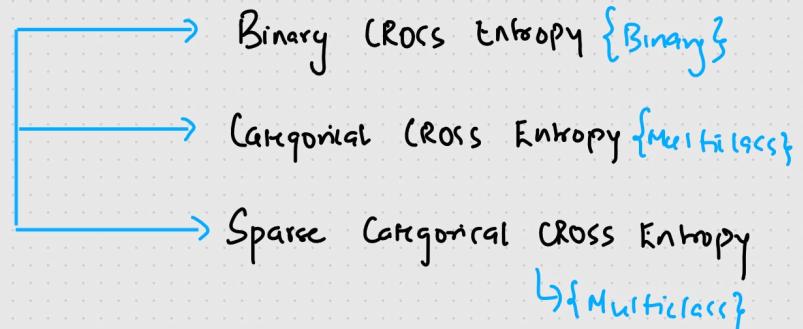
Cost function = $\sqrt{\frac{\sum_{i=1}^N (\hat{y}_i - y_i)^2}{N}}$

Advantages

Disadvantages

② Loss or Cost function For Classification Problems

Classification \rightarrow Cross Entropy



① Binary Cross Entropy

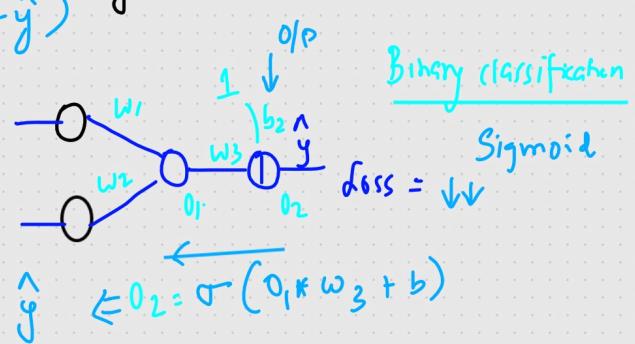
\downarrow log loss

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

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

y = Actual Value

\hat{y} = Predicted Value



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

\Rightarrow Sigmoid Activation function

② Categorical Cross Entropy (Multiclass classification)

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

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

Actual value $\leftarrow y_{ij} = [y_{i1} \ y_{i2} \ y_{i3} \ \dots \ y_{ic}]$

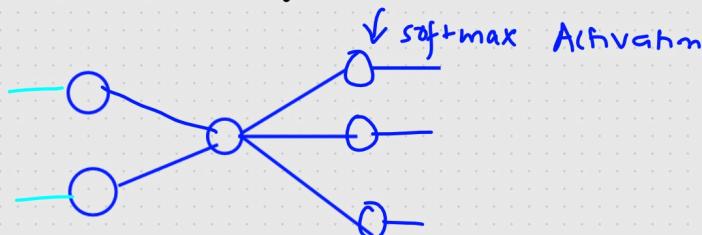
$$[y_{21}, y_{22}, y_{23}, \dots, y_{2c}]$$

:

:

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

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



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

$$\text{Op } \hat{y}_{ij} = \text{Probabilities} \quad [0.1, 0.2, 0.3, 0.2, 0.2] \leftarrow 1$$

\downarrow

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

(Ross Entropy) $\overbrace{\hspace{10em}}$.

[This also gives the probability of other categories]

③ Sparse Categorical Cross Entropy

$$\left[\begin{matrix} 0 & 1^{\text{st}} & 2^{\text{nd}} \\ 0.2, 0.3, & 0.5 \end{matrix} \right] \rightarrow (\text{Gregorius})$$

1

$$\boxed{2^{\text{nd}} \text{ Index}} \Rightarrow \underline{0/p}$$

$\left[\begin{matrix} 0 & 1 & 2^{\text{nd}} & 3^{\text{rd}} & 4^{\text{th}} \\ 0.2, 0.3, 0.1, 0.2, 0.2 \end{matrix} \right] \Rightarrow \begin{matrix} 1^{\text{th}} \\ \Downarrow \\ \text{Index} \end{matrix}$

(category \Rightarrow o/p)

Disadvantages

① Loosening info about the probability of other category.

f) Right Combination

Activation applied

Hidden Layers

O/P Layer

Problem Statement

Loss function

1

Rehn 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

fineray

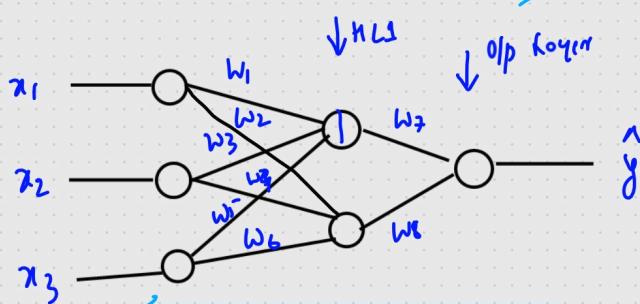
Regression

MSE, MAE, RMSE ders,

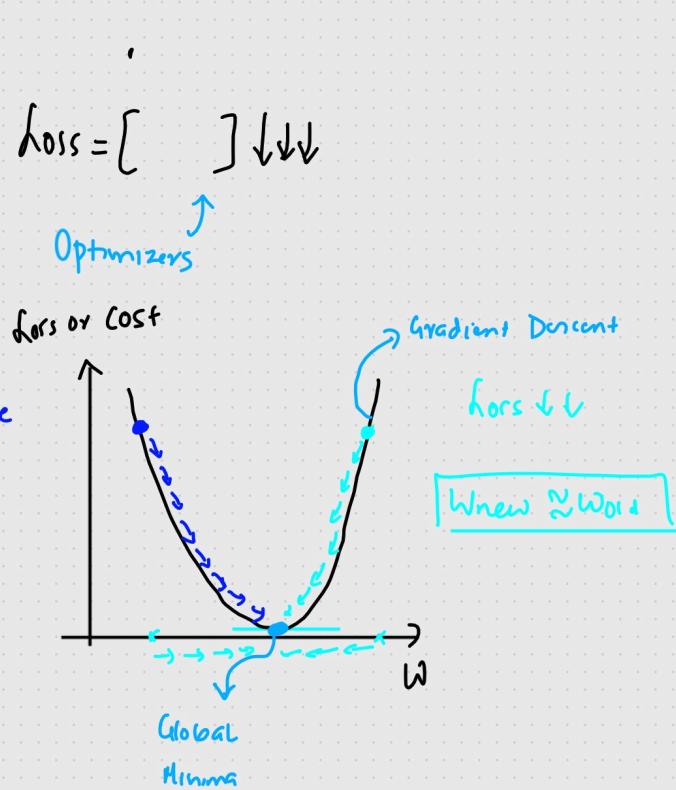
Optimizers

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

Gradient Descent Optimizer



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



MSE

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

DataSet = 1000 Data Points

Epochs, Iteration

1 Epoch = 1 Iteration

1000 datapoint → $\hat{y}_i \Rightarrow$ Cost function wh

Weights will get updated

1 Epoch = 1 Iteration

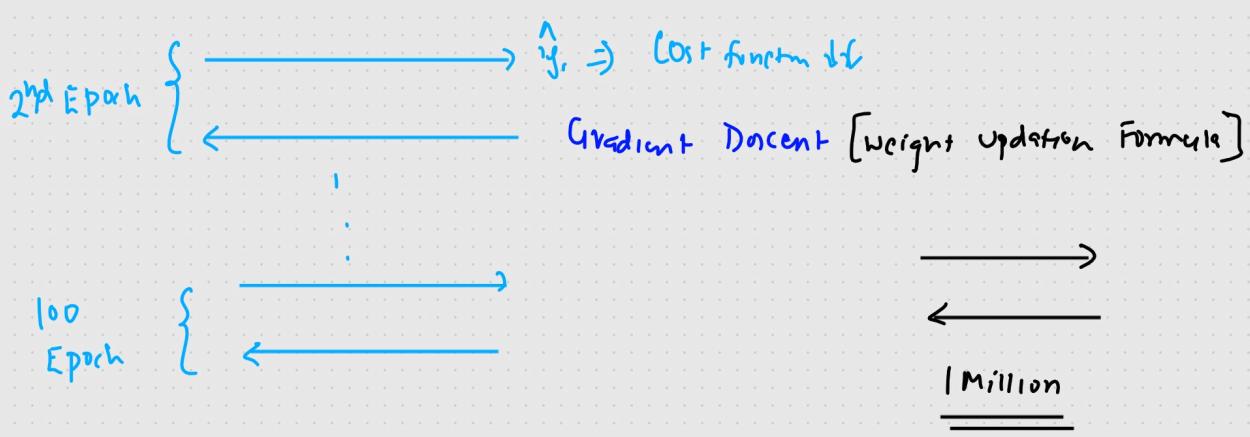
Gradient Descent optimizers [Weight updation]

10 Iteration

$100/10 = 100$

100 →

← 100



Advantages

① Convergence will happen.

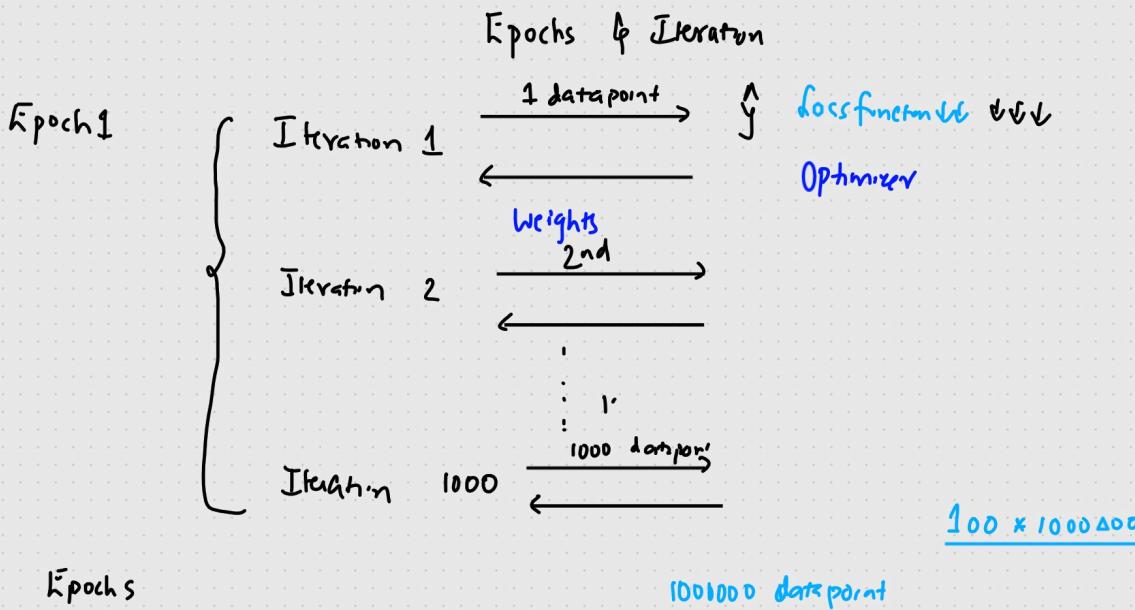
Disadvantage

① Huge Resource RAM, GPU
 \Downarrow

Resource Intensive -

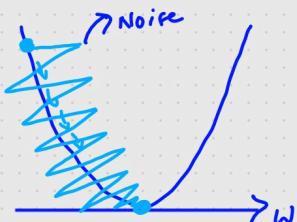
② Stochastic Gradient Descent (SGD)

1000 datapoint



Advantage

① Solve Resource Issue



Disadvantage

① Time Complexity $\uparrow \uparrow$

- ④ Convergence will also take more time.
- ⑤ Noise gets Introduced

$\hat{y} \Rightarrow \text{Cost}$

③ Mini Batch SGD

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

Epoch, Iteration, Batch-size

Data points = 100000

batch.size = 1000

Epoch 1

Iteration 1

MSGD

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

Optimizer \Rightarrow Mini Batch SGD

change the weight

8gb 1

Iteration 2

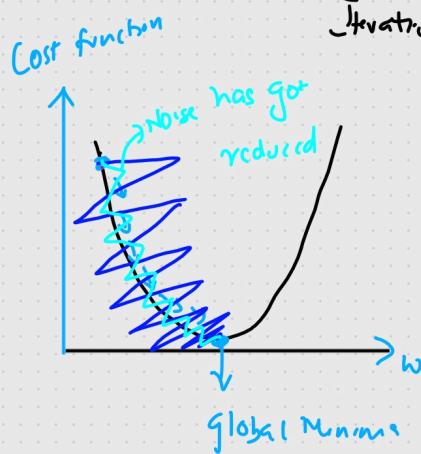
\downarrow

16gb 500

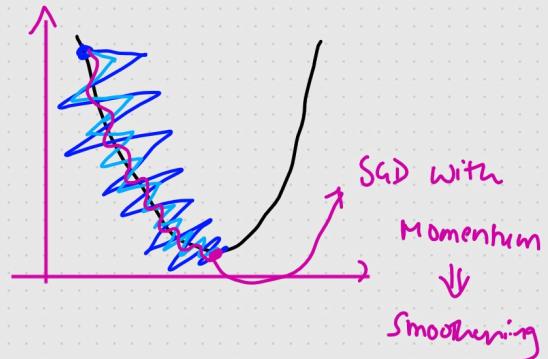
Iteration 3

!

Iteration 100



Batch Size



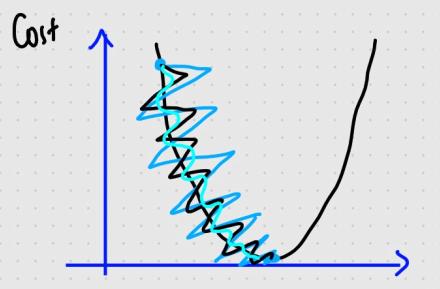
Advantages

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

Disadvantage

- ① Noise still exists

④ SGD With Momentum



Weight Update formula

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

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

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

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

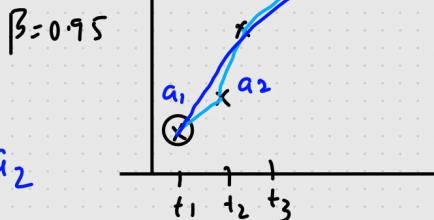
Time = $t_1, t_2, t_3, t_4, \dots, t_n$ Time Series

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

$$V_{t_1} = a_1$$

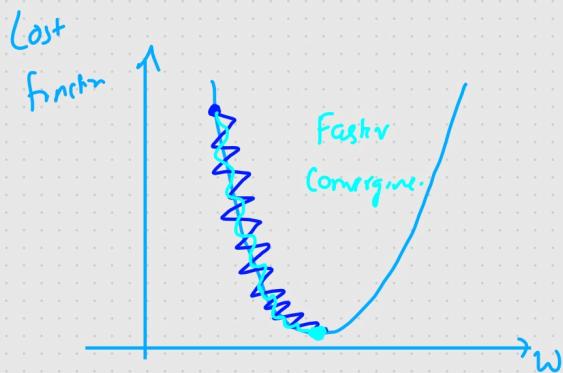
$$V_{t_2} = \beta * V_{t_1} + (1-\beta) * a_2$$

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



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

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



Advantage

- ① Reduces the noise
- ② Quick Convergence

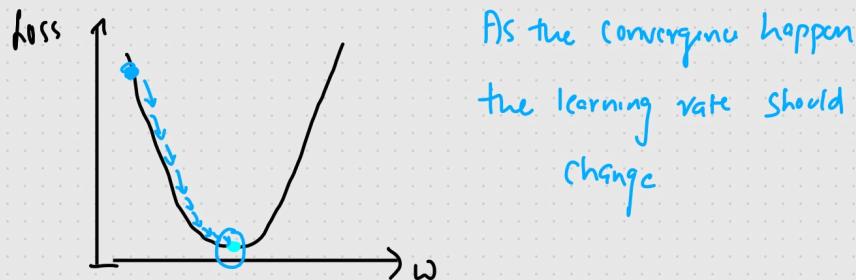
⑤ Adagrad = Adaptive Gradient Descent

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

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

Learning Rate = 0.001

Rate



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

$0.00001 \approx 0$

$\eta' = \frac{\eta}{\sqrt{d_t + \epsilon}}$ $\eta' \downarrow \text{and } d_t \uparrow$

$d_t \approx 1$ $\epsilon \approx \text{small value}$

\downarrow

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

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

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

Disadvantage

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

⑥ Adadelta And RMSprop

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

Exponential Weight Average
 $\beta = 0.95$ $S_{dw_t} = 0$

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

[Dynamic LR + Smoothening (EWA)]

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

⑦ Adam Optimizer

SGD with Momentum + RMSprop [Dynamic LR + Smoothening]

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

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

weight
Bias
Updation

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

$S_{dw_t} = 0$ EWA

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

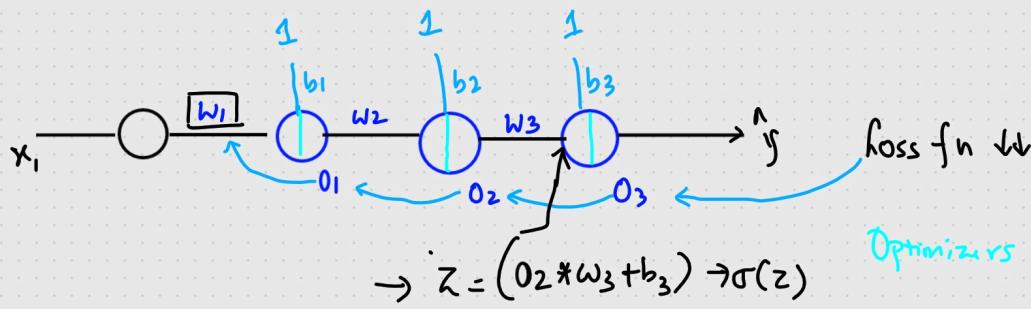
[EWA]

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

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

\Rightarrow Momentum
 \hookrightarrow Smoothening

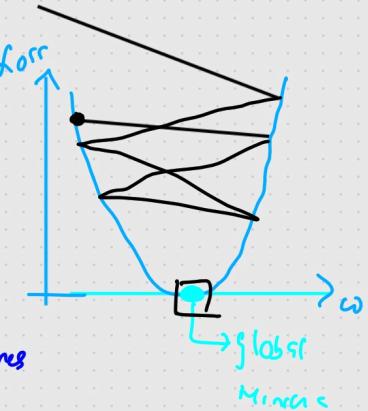
Exploding Gradient Problem \Rightarrow Weight Initialization Technique



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

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

\Rightarrow Chain Rule of Derivatives



$$\frac{\partial h}{\partial w_{\text{old}}} = \frac{\partial h}{\partial o_3} * \frac{\partial o_3}{\partial o_2} * \frac{\partial o_2}{\partial o_1} * \frac{\partial o_1}{\partial w_{\text{old}}}$$

big * big * big * big \Rightarrow big value

Weight initialization
↓

Very high value

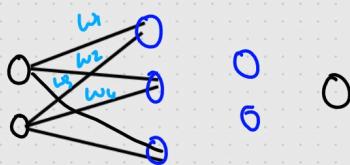
$$\frac{\partial o_3}{\partial o_2} = \frac{\partial \sigma(z)}{\partial z} * \frac{\partial z}{\partial o_2}$$

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

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

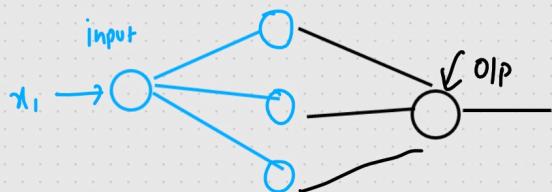
Weight Initializing Techniques

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

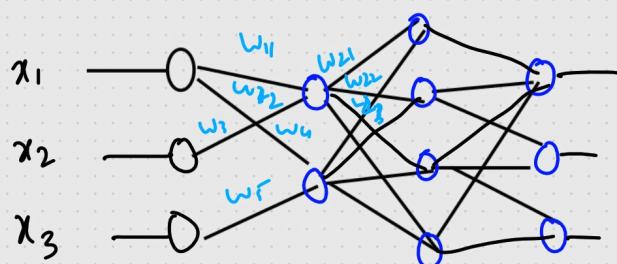


Key Points

- 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

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

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

② Xavier/Glorot Initialization

Researcher → 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

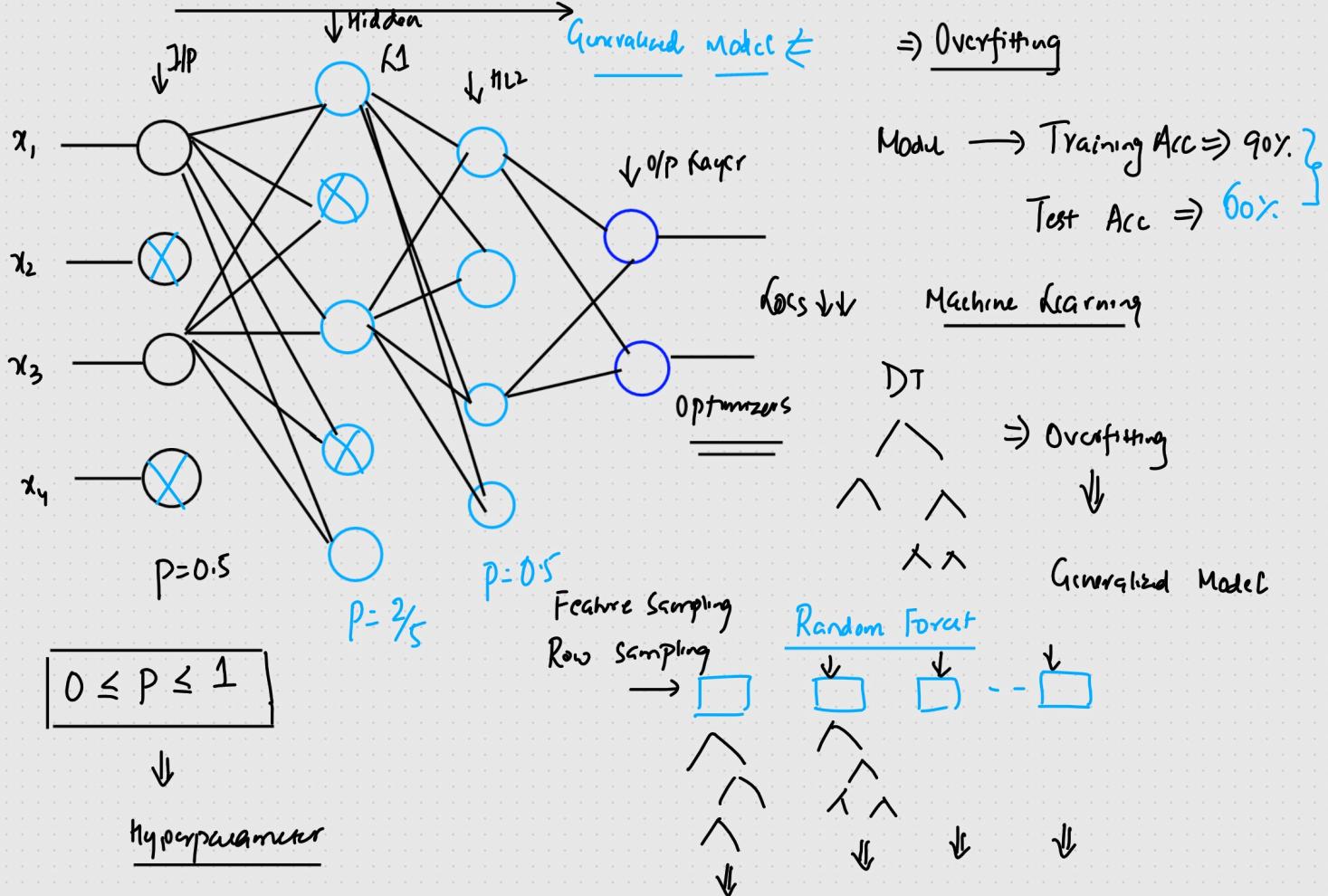
$$w_{ij} \sim N(0, \sigma)$$

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

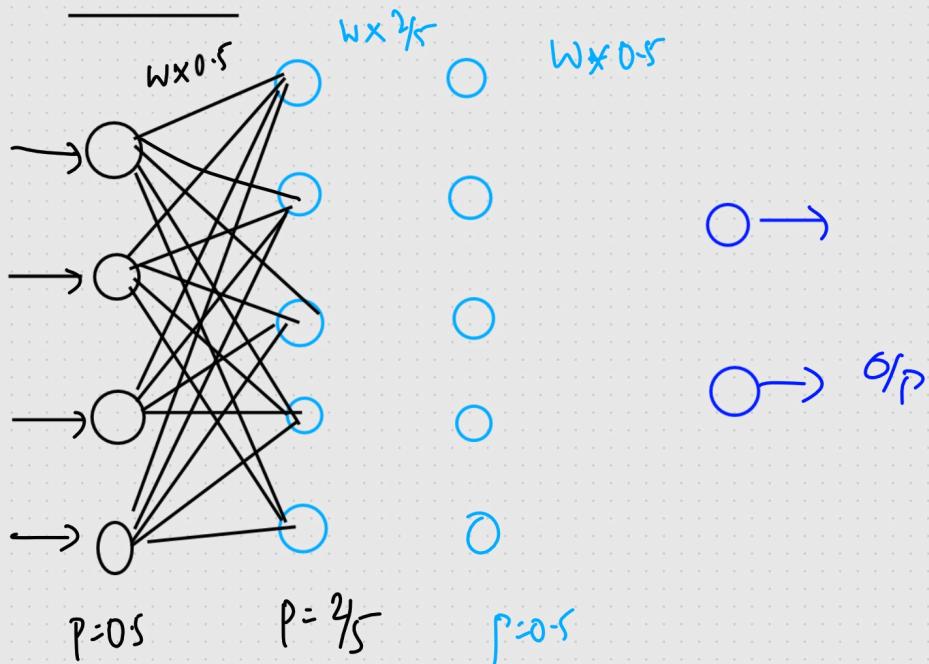
② He uniform

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

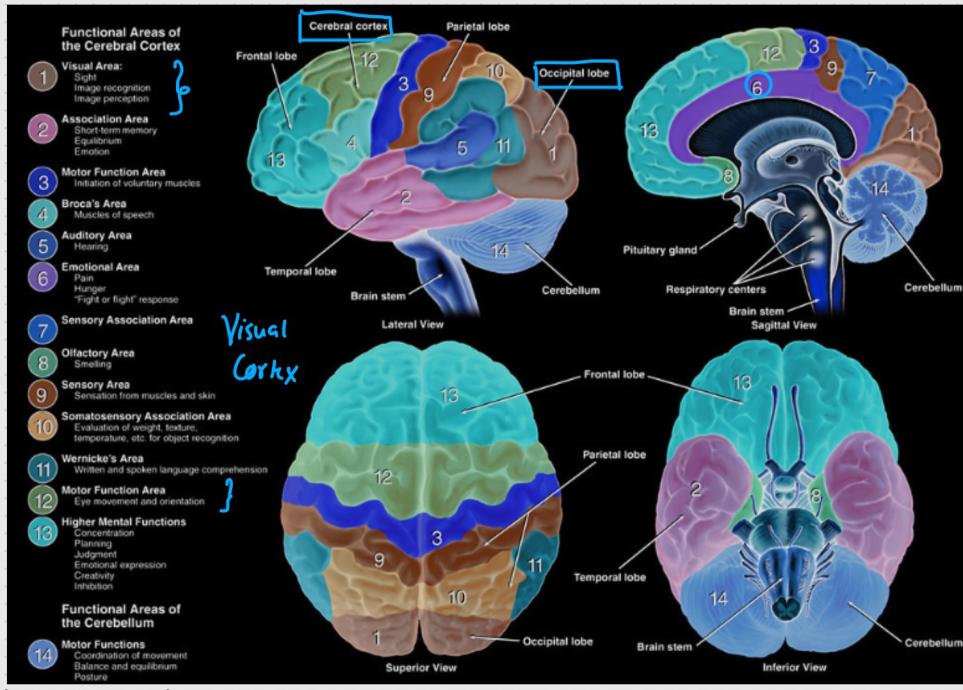
Drop Out layer



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

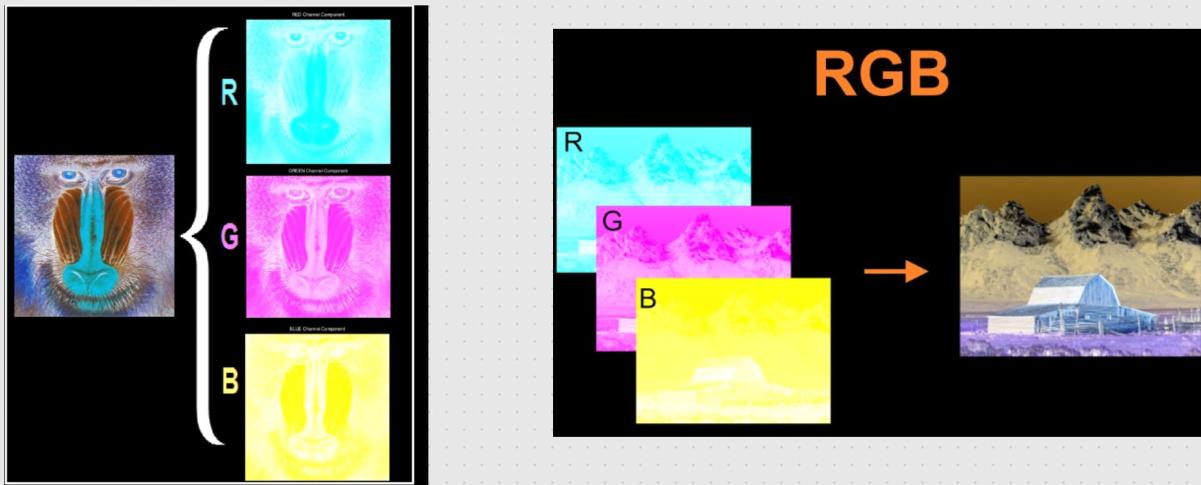
Visual Cortex (VI-V5) [Region of the brain that receives, integrates and processes visual information relayed from the retinas].

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

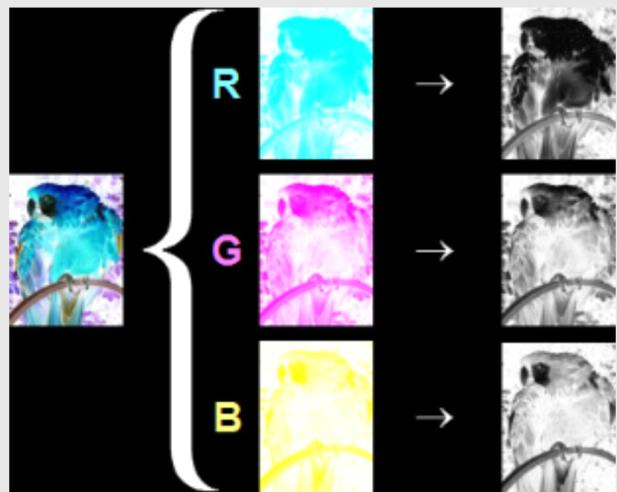
V2
|
V3 V4 V5 [Differences in color, complex patterns, object orientation]
[Object Recognition]

Visualize the Image

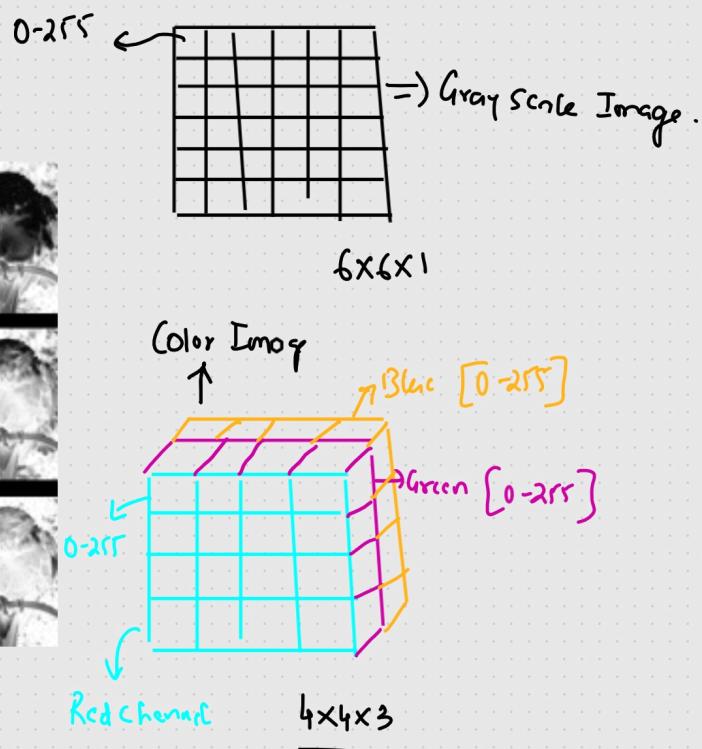
③ RGB Images And Gray Scale Images



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



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



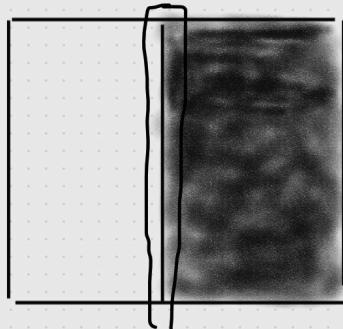
④ Convolution Operation In CNN

$\rightarrow (0,1)$

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

$6 \times 6 \times 1$

\Rightarrow



Step 1

① Normalize

Divide by 255

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

Convolution operation

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

$h=6$

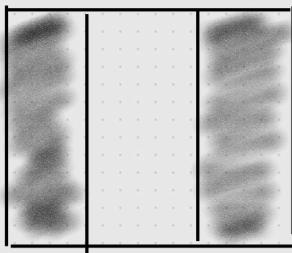
$f=3$

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

$O/p = 4$

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

$6 \times 6 \times 1$



*

Vertical edge filters



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

$$h - f + 2p + 1 = 6$$

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

$$2p = 6 - 4$$

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

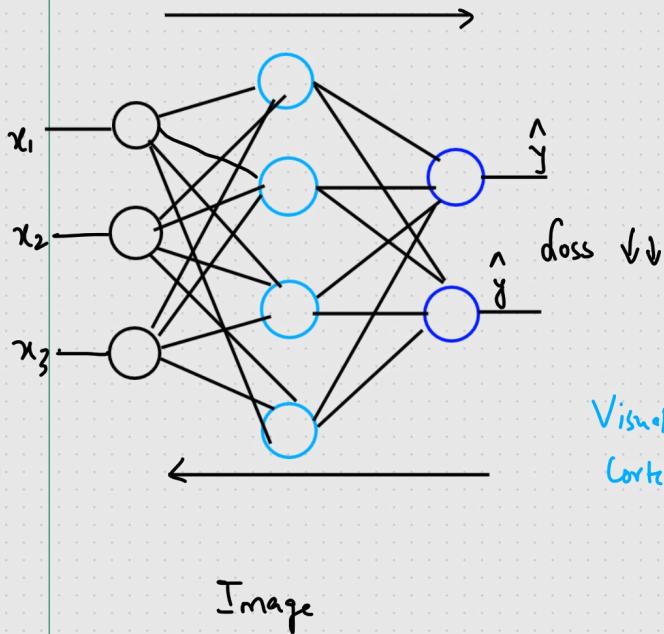
7×7

3×3

7×7

How much padding you need to apply?

⑥ Operation of CNN Vs ANN

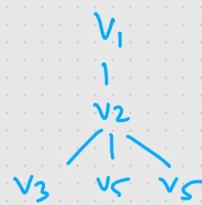


$$z = w^T x_i + b$$

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

loss ↓↓

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



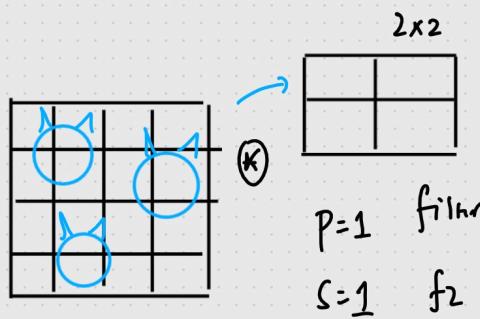
f_2
 f_3
⋮
 f_n

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

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

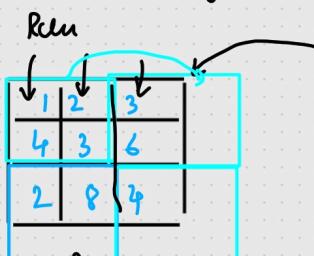
Convolution Layer

⑦ Max Pooling, Min Pooling, Mean Pooling



$P=1$
 $S=1$
 f_2
 f_3
 f_1

Convolution Layer



Max Pooling

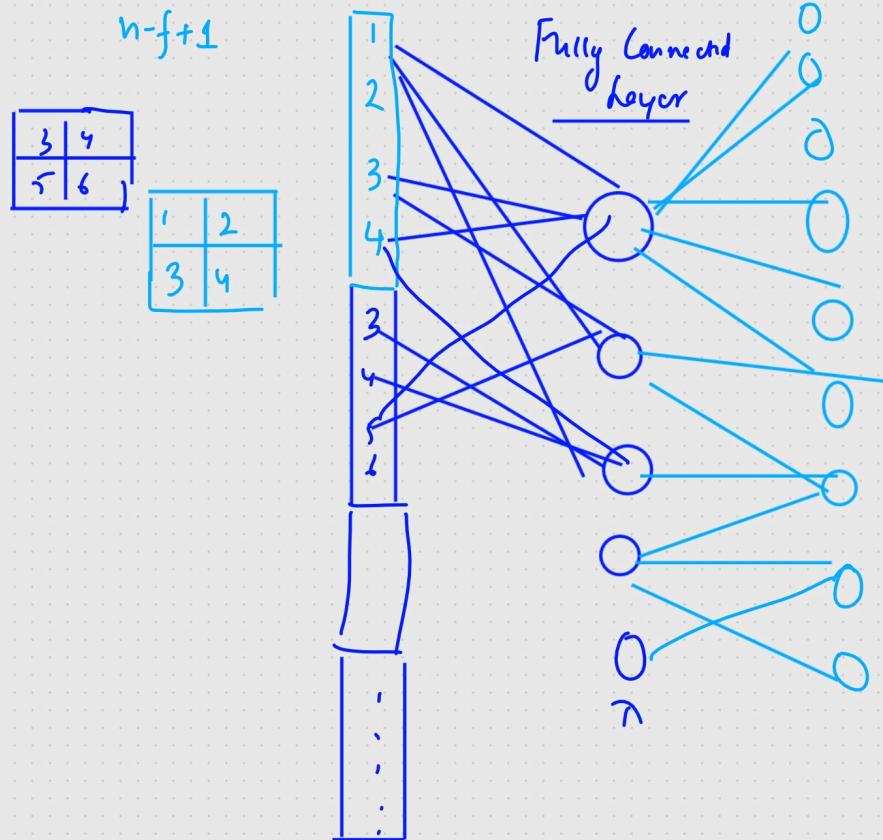
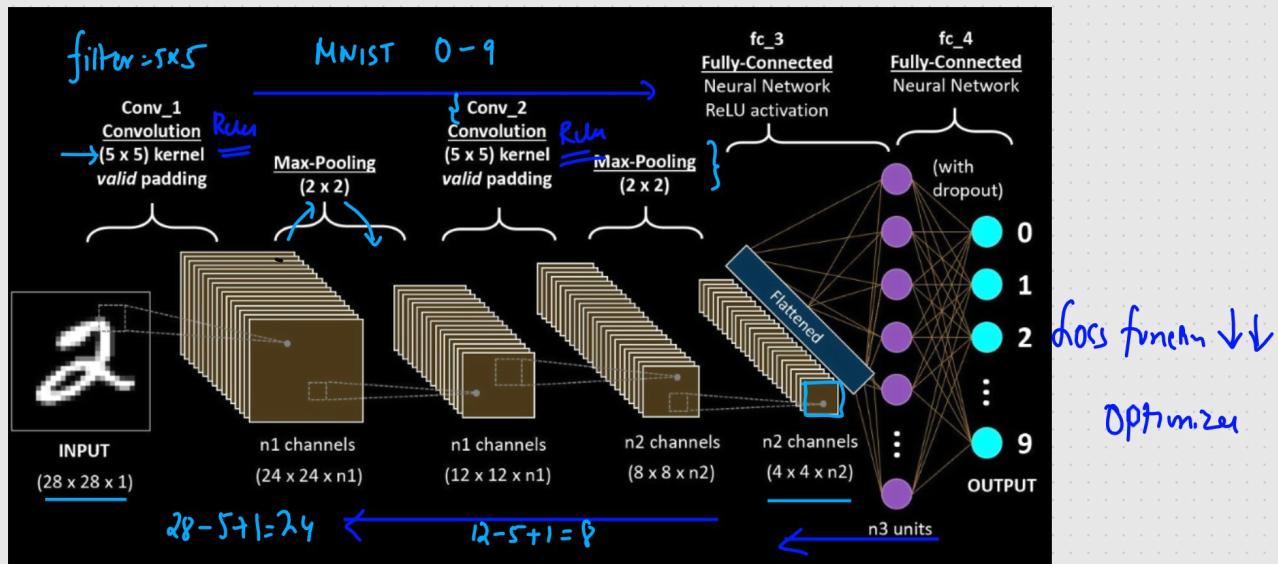
1
2

2×2

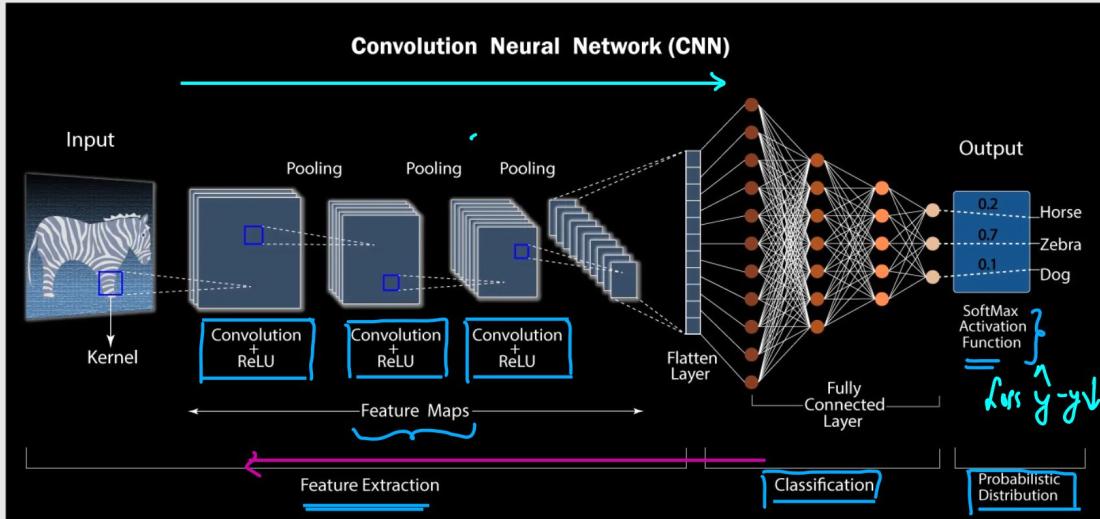
4	6
8	4

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		

*

+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}

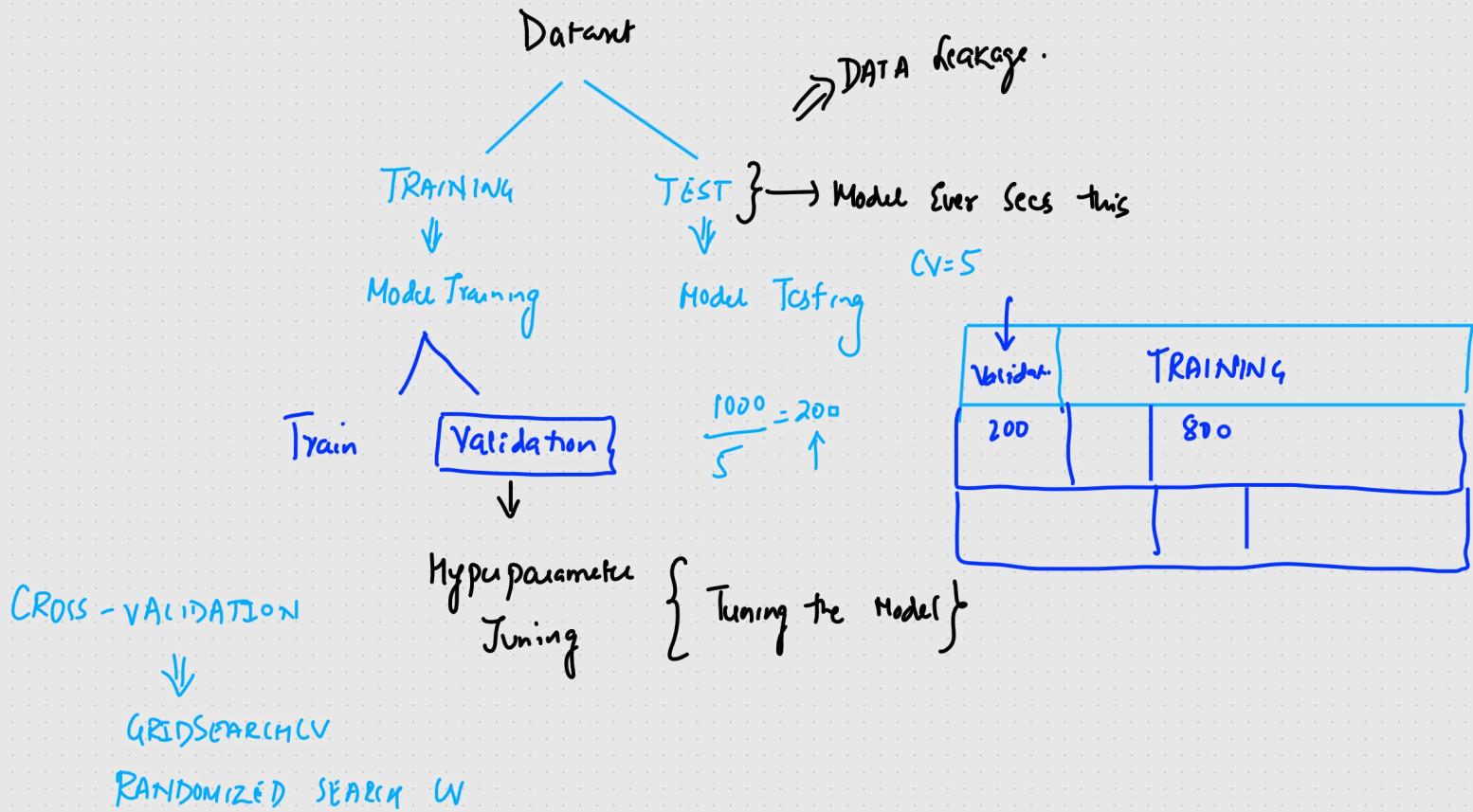
① BOW, TFIDF, Word2vec, Avg Word2Vec {Python Practical Implementation}



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



Hindi

English

Suggestion \uparrow Grammaticality
Completion of sentences

Deep learning

① RNN ② LSTM-RNN ③ Transformers ④ BERT

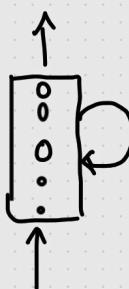
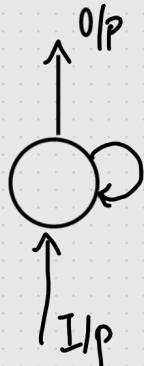


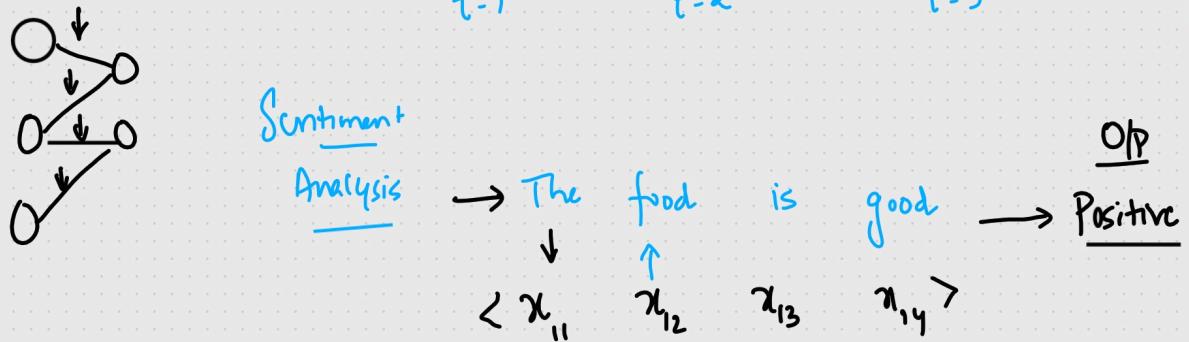
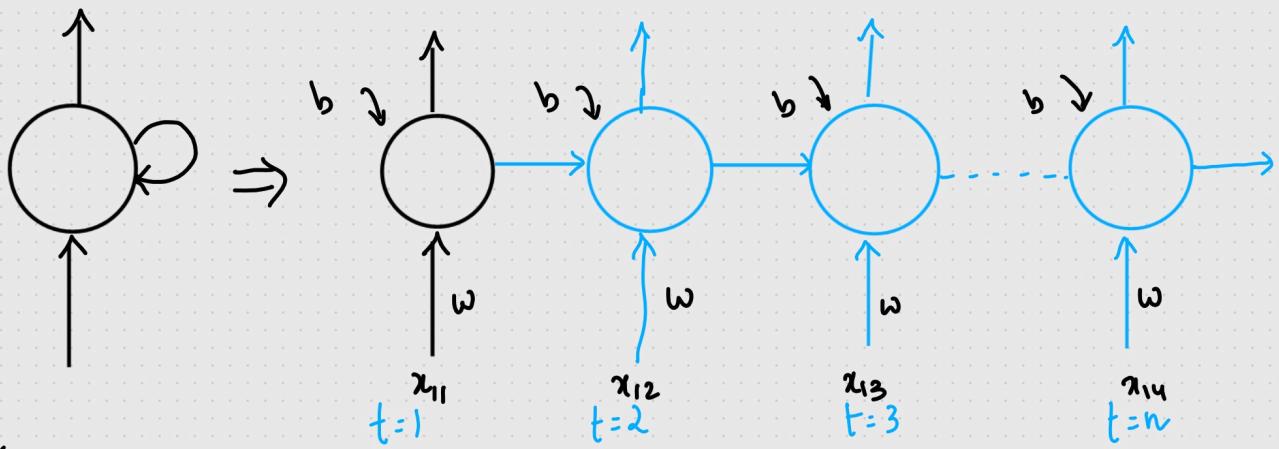
Word2vec

\downarrow Deep learning

Words \rightarrow Vec

① Recurrent Neural N/w

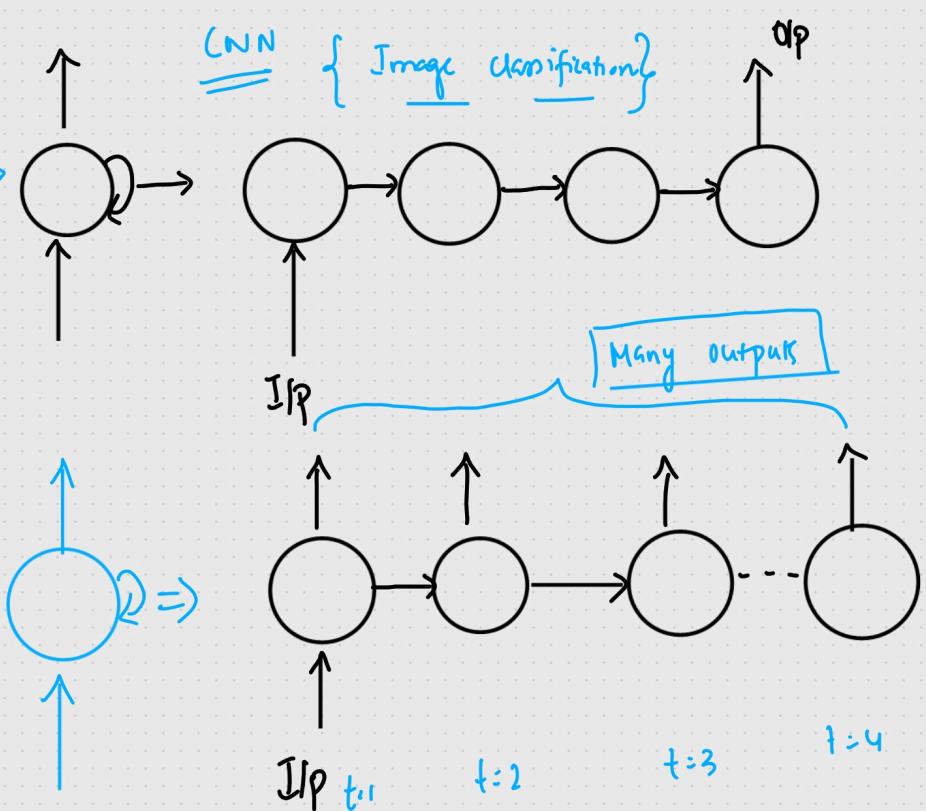




$x_{11} \rightarrow \text{Word2vec} \rightarrow \text{Vectors} \rightarrow d = \underline{\underline{300}}$

Types of RNN

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

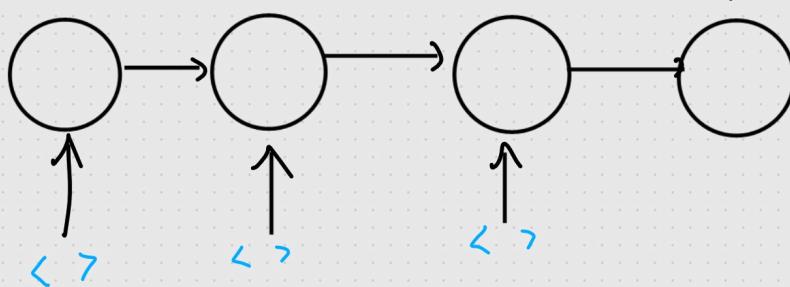


{ Eq: Music generation, Text generation }

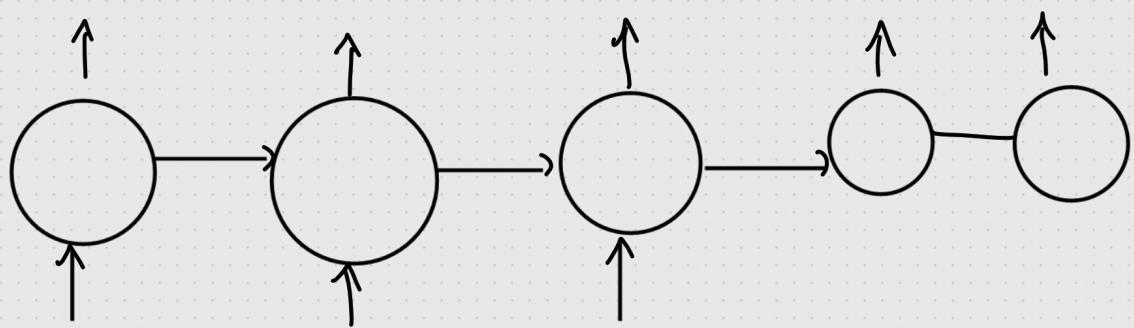
Google Search Suggestion, Movie Recommendation

③ Many to One

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



④ Many to Many



① Language Translation

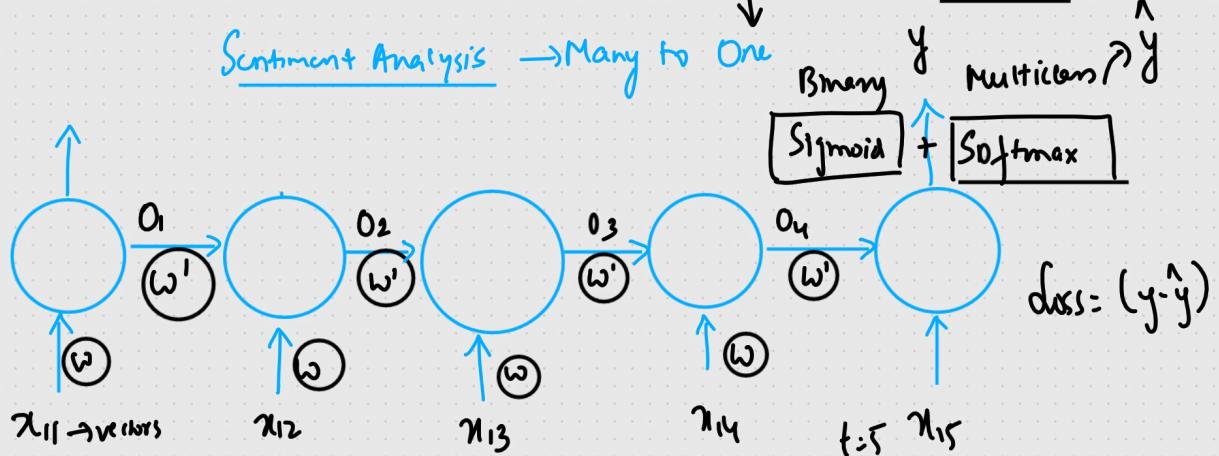
② Question answers

③ Chatbots.

{ ① Forward Propagation
② Backward Propagation }

① Forward Propagation In RNN

Time Series Data



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

The food is very good Positive.

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

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

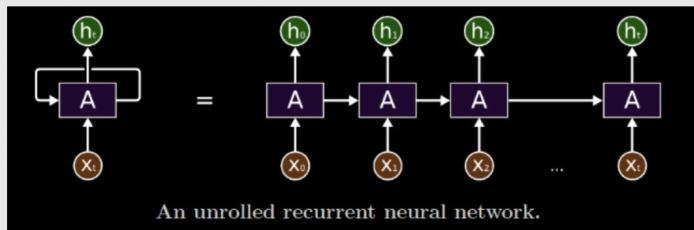
$$O_3 = f(x_{13} * \omega + O_2 * \omega')$$

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
 → Long Term Memory ✓ ↗
 → Short Term Memory ✓ ↗
- ④ LSTM Architecture
- ⑤ Working of LSTM RNN

Problems With RNN → Long Term Dependency

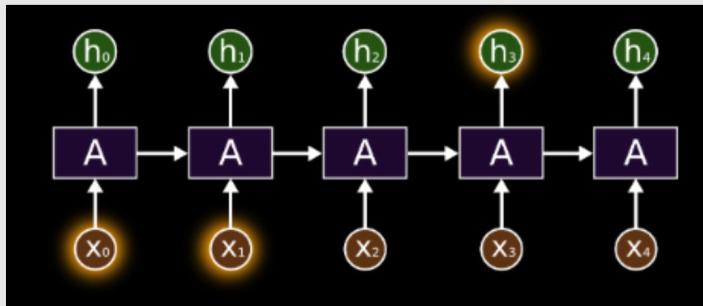


Vanishing Gradient Problem

Task:

Next Word In a Sentence

The color of the sky is blue
— further context

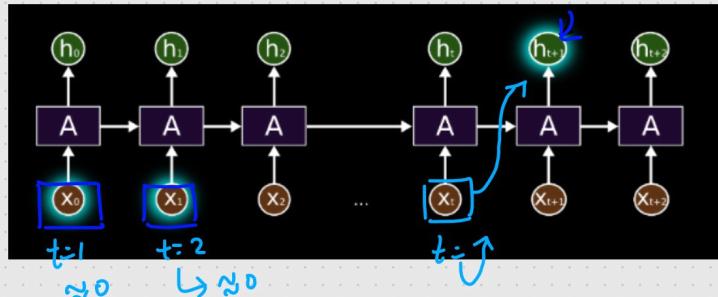


Gap is 1cm

Huge gap O/p ← Context

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

Name of language
↓
further context



huge gap → Long Term Dependency

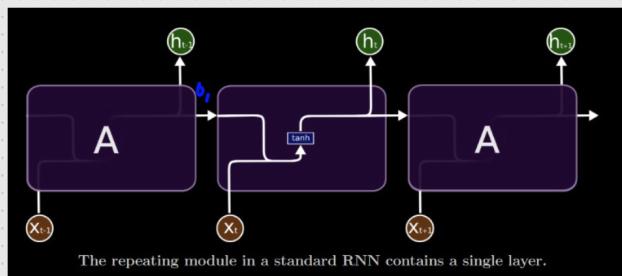
0-1

RNN → Long Term Dependency → Vanishing Gradient Problem

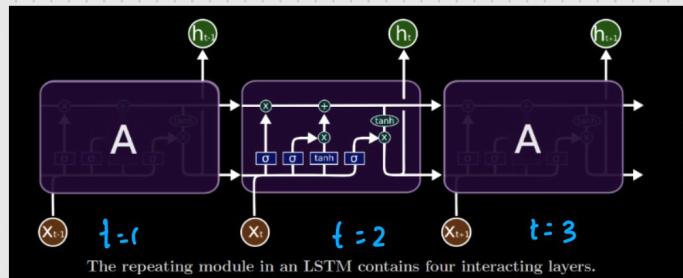
Chain Rule → ≈ 0 .

Basic Representation of RNN And LSTM RNN

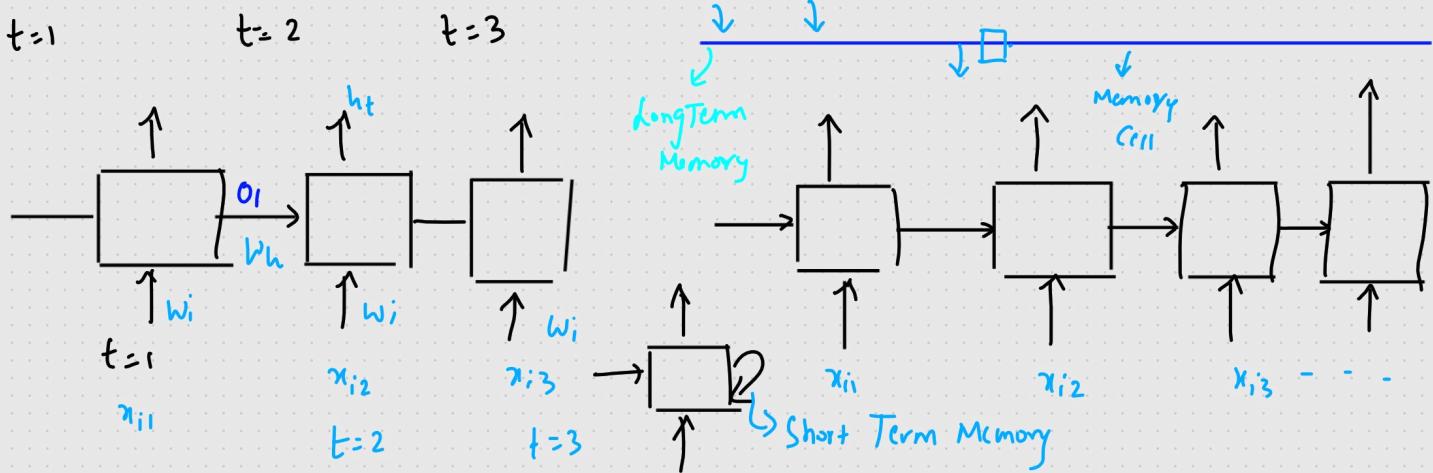
LSTM RNN



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



The repeating module in an LSTM contains four interacting layers.

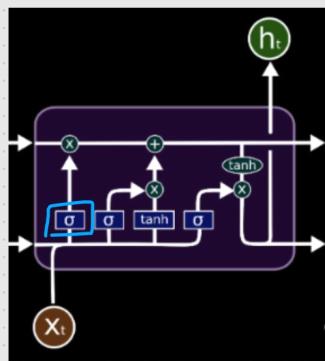


LSTM RNN → Long Term Memory
→ Short Term Memory

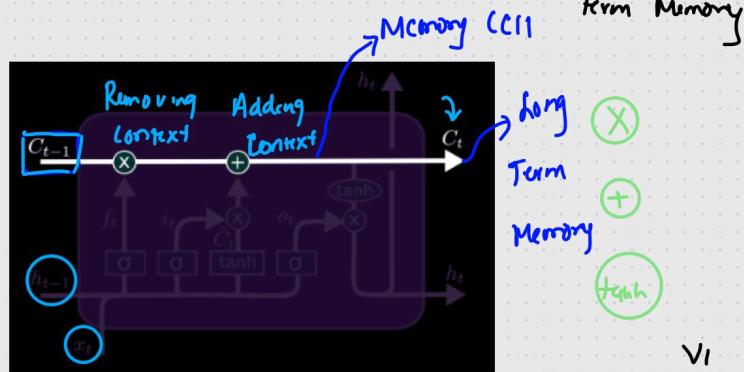
Convayana But : fuggages



LSTM Architecture



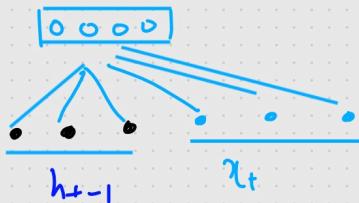
Basic Architecture



Combining 2 Vectors

$$h_{t-1} = \begin{bmatrix} 1 & 2 & 3 \end{bmatrix}$$

$$x_t = \begin{bmatrix} 2 & 3 & 4 \end{bmatrix}$$



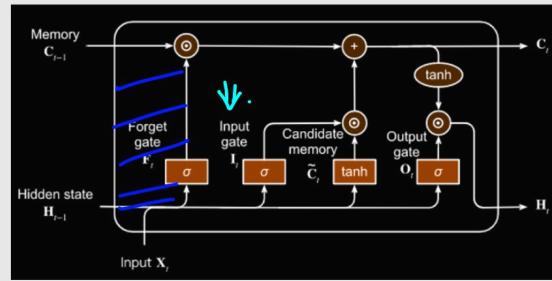
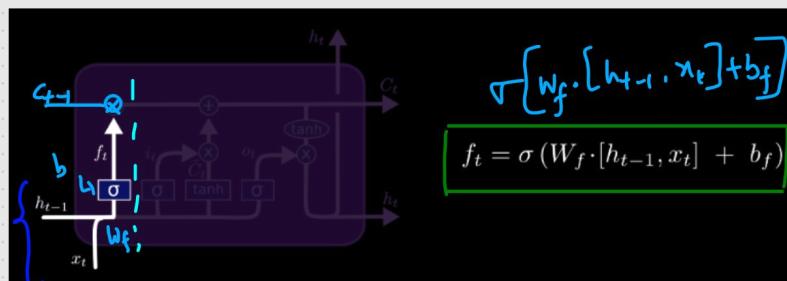
v_1

$$\begin{bmatrix} 1 & 2 & 3 \\ 4 & 5 & 6 \end{bmatrix} \xrightarrow{\text{tanh}} \begin{bmatrix} 4 & 10 & 18 \\ 5 & 7 & 9 \end{bmatrix}$$

$\otimes \Rightarrow [4 \ 10 \ 18]$

$\oplus \Rightarrow [5 \ 7 \ 9]$

$\tanh \Rightarrow [\tanh(1) \ \tanh(2) \ \tanh(3)]$



Forget Gate

Text Next Word

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

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



$$\begin{bmatrix} 0 & 2 & 4 & 1 \end{bmatrix} \quad \begin{bmatrix} 4 & 5 & 12 \end{bmatrix} \quad \dots$$

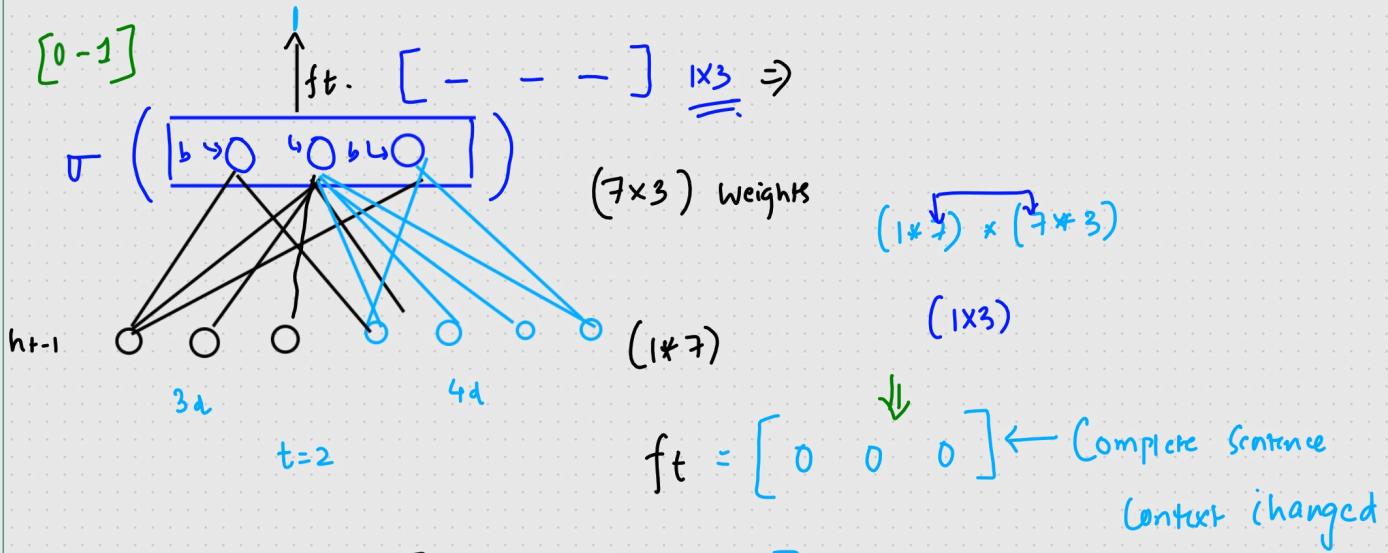
x_t

$h_{t-1} = \begin{bmatrix} 1 & 2 & 4 \end{bmatrix}$

$3d$

$c_{t-1} = \begin{bmatrix} 4 & 2 & 1 \end{bmatrix} \Leftarrow$

$3d$



$$\textcircled{1} \quad c_{t-1} = [6 \quad 8 \quad 9] \otimes [0 \quad 0 \quad 0]$$

$$= [0 \quad 0 \quad 0] \leftarrow \text{Removing all the previous context}$$

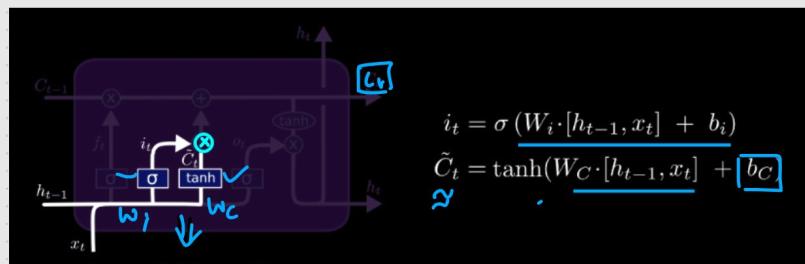
$$ft = [1 \quad 1 \quad 1]$$

$$\textcircled{2} \quad c_{t-1} = [6 \quad 8 \quad 9] \otimes [1 \quad 1 \quad 1] = [6 \quad 8 \quad 9]$$

$$\textcircled{3} \quad c_{t-1} = \begin{matrix} 6 \\ 8 \\ 9 \end{matrix} \otimes \begin{matrix} 0.5 \\ 1 \\ 0.5 \end{matrix} = \begin{matrix} 3 \\ 8 \\ 4.5 \end{matrix}$$

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

$$I_t = \begin{bmatrix} 2 & 4 & 1 \end{bmatrix} \xrightarrow{\oplus} \begin{bmatrix} 0.8 & 0 \end{bmatrix} \xrightarrow{\text{tanh}} c_t = \begin{bmatrix} 0 & 2 & 0 \end{bmatrix} \Rightarrow \underline{\text{Input Gate}}$$

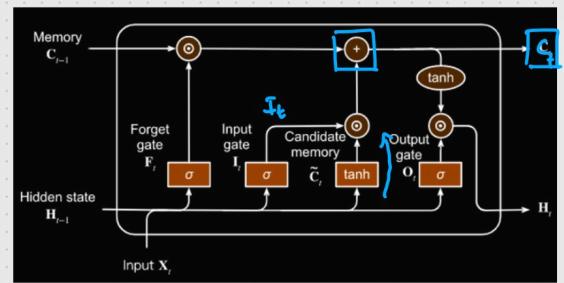
$b \rightarrow \begin{bmatrix} 0 & 0 & 0 \end{bmatrix}$

$x_t \rightarrow \begin{bmatrix} 0 & 0 & 0 \end{bmatrix}$

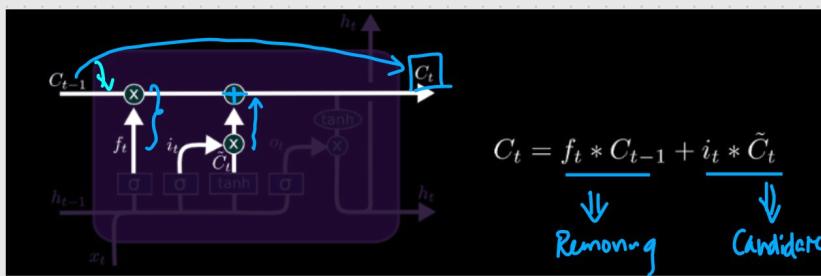
$h_{t-1} \rightarrow \begin{bmatrix} 0 & 0 & 0 \end{bmatrix}$

$W_i \rightarrow \begin{bmatrix} 1 & 1 & 1 & 1 \end{bmatrix}$

$w_c \rightarrow \begin{bmatrix} 1 & 1 & 1 & 1 & 1 & 1 & 1 \end{bmatrix}$

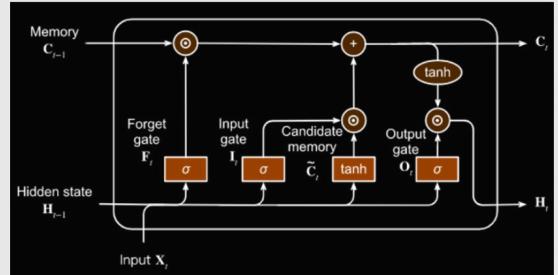


Context = If any information needed to be added to the memory
 $c_{t-1} \rightarrow$ The information will be added



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

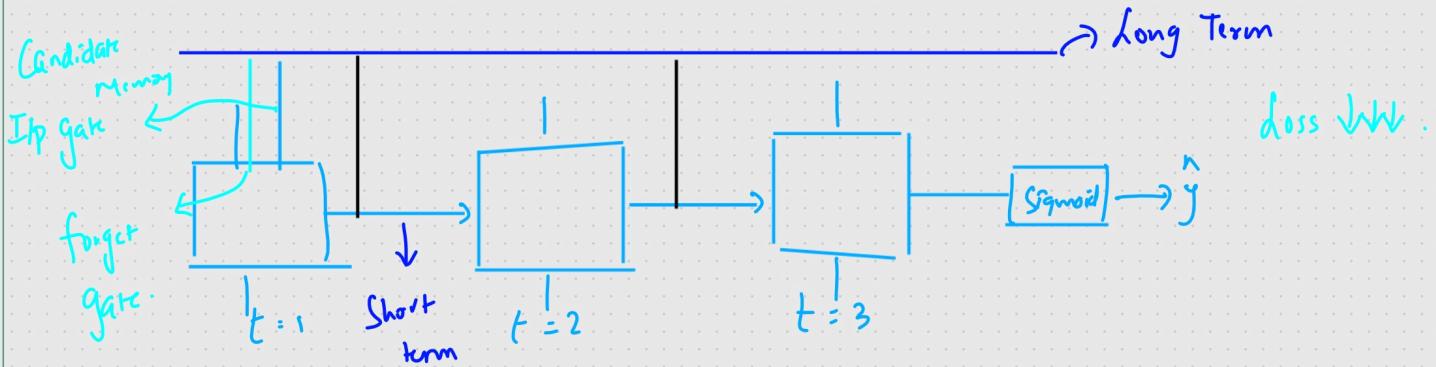
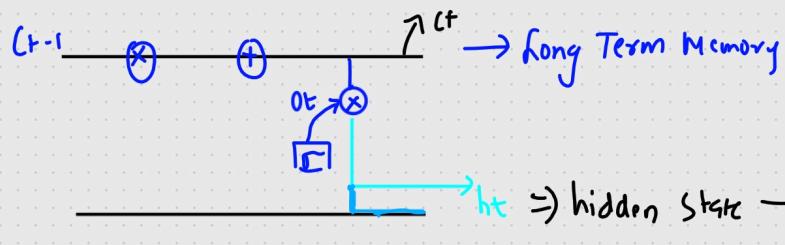
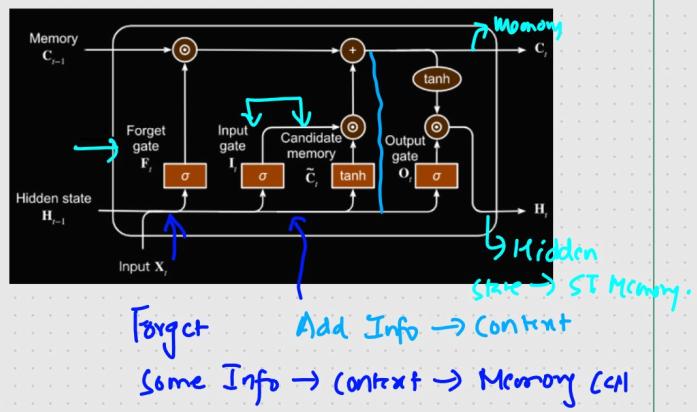
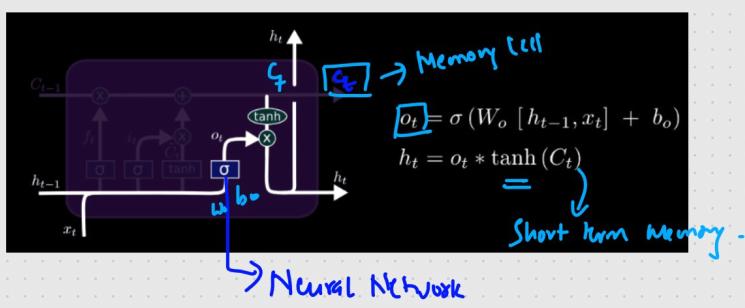
OR
Forgetting
Some info



Forget Gate

I/p Gate \otimes Candidate memory
+
 $c_{t-1} \Rightarrow c_t$

Output gate LSTM RNN



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

GRU RNN \Rightarrow LSTM Variant

Training Data With LSTM RNN

{Training} Text Paragraph

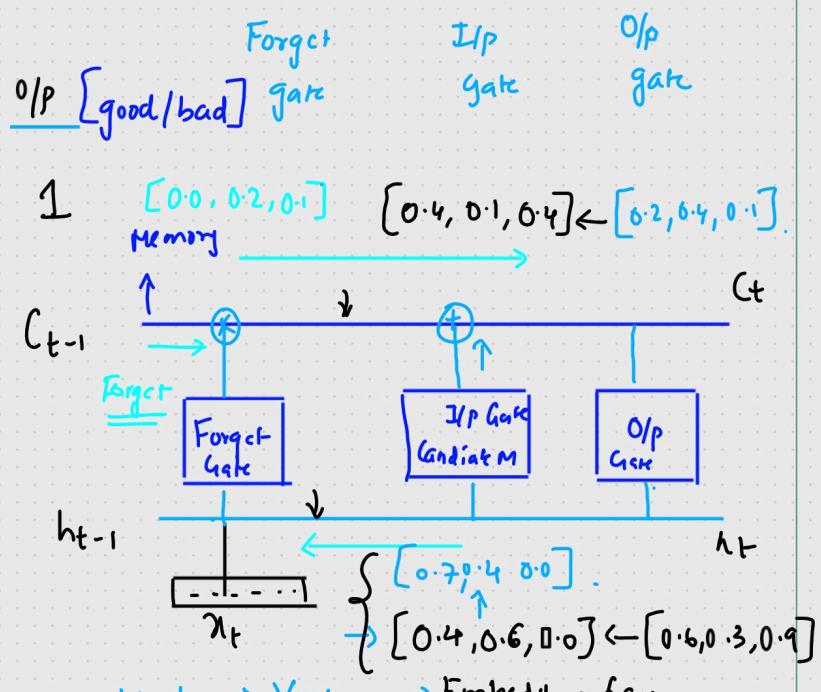
I Went to Restaurant and order burger

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.



Word2Vec [3 dimension-vector]

→ $[\text{Good}] [\text{Bad}] [\text{Healthy}] \leftarrow \text{Black Box}$

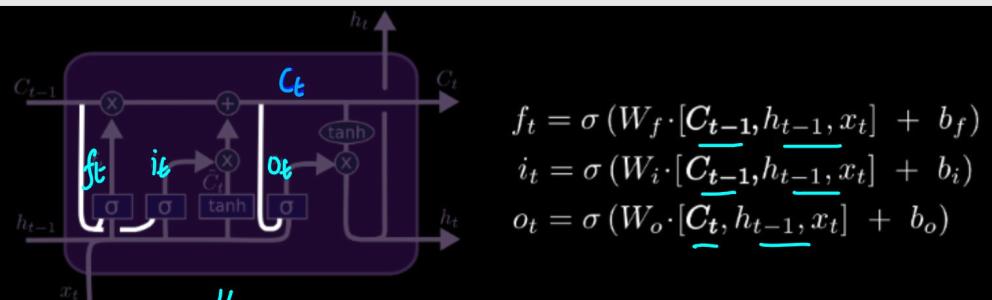
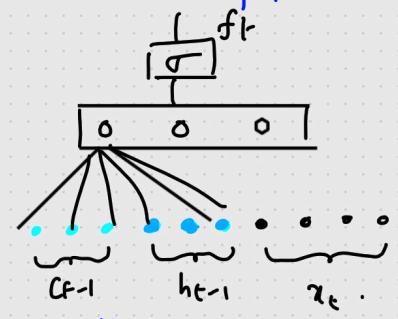
Tasty $[0.9 \quad 0.0 \quad 0.1] \leftarrow 3 \text{ d.}$

Variants of LSTM RNN

LSTM Variants Introduced By Gers & Schmidhuber [2000]

LSTM RNN [1970-80]

↳ Research paper



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

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

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

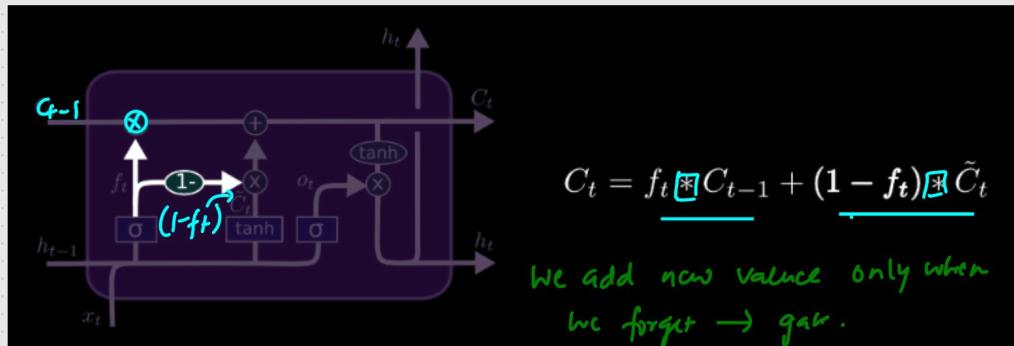
Connections → From Memory cell to

forget gate Peephole
i/p gate \Rightarrow 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.

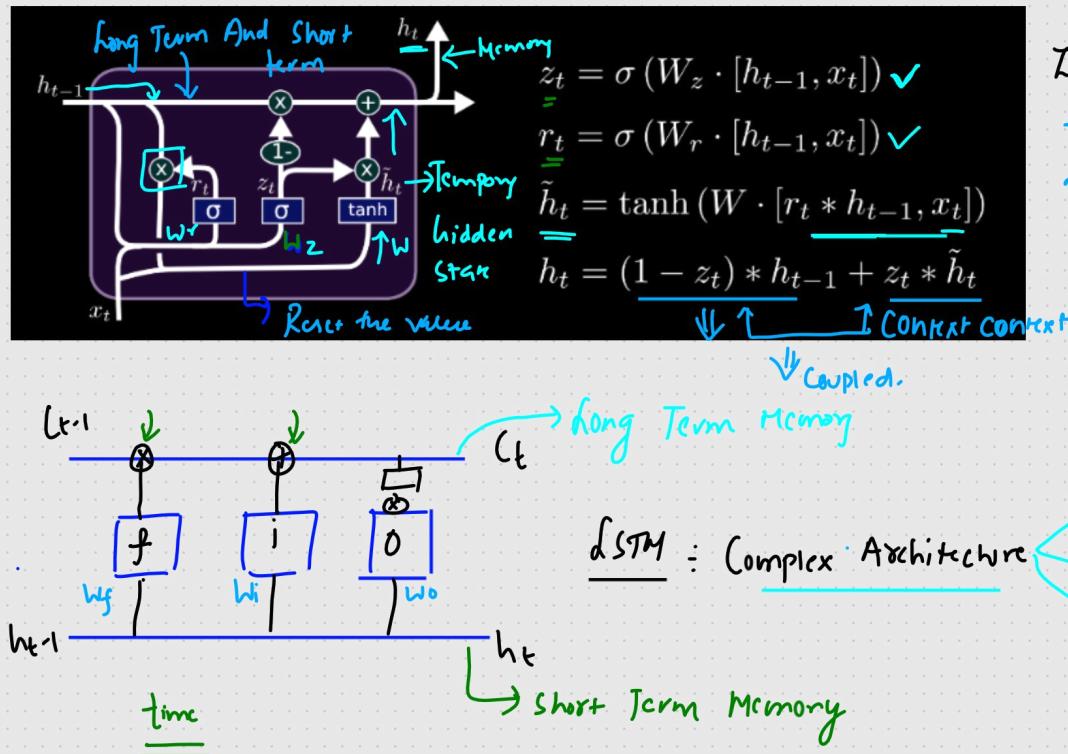
We only i/p new values to the state when we forget something older.

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]]

1990 \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.

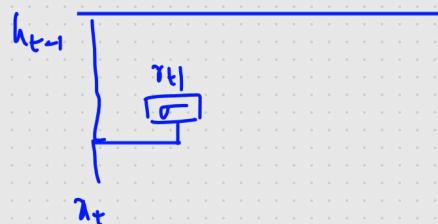
Trainable parameters
[w_f, w_i, w_o]

forget
i/p + candidate memory
O/p

Training Time $\uparrow \uparrow$

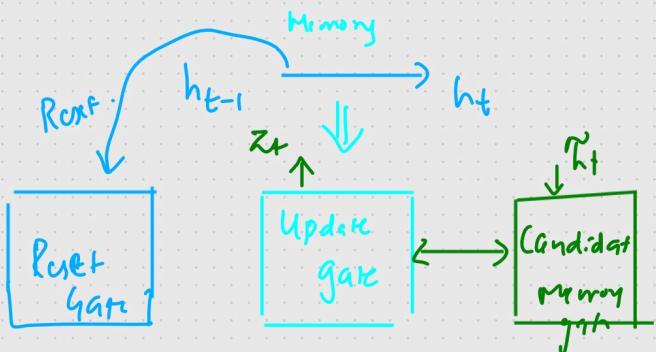
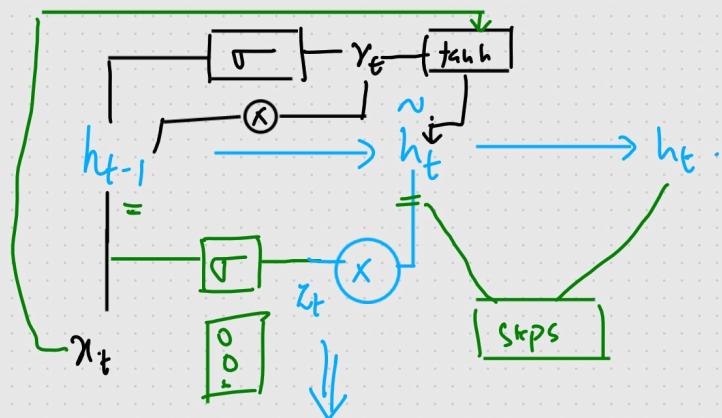
z_t ,

Reset gate = $r_t =$



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

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



What Context Info needs to be Added



Candidate hidden state · [Current Context]

↓
Imp → Add Info

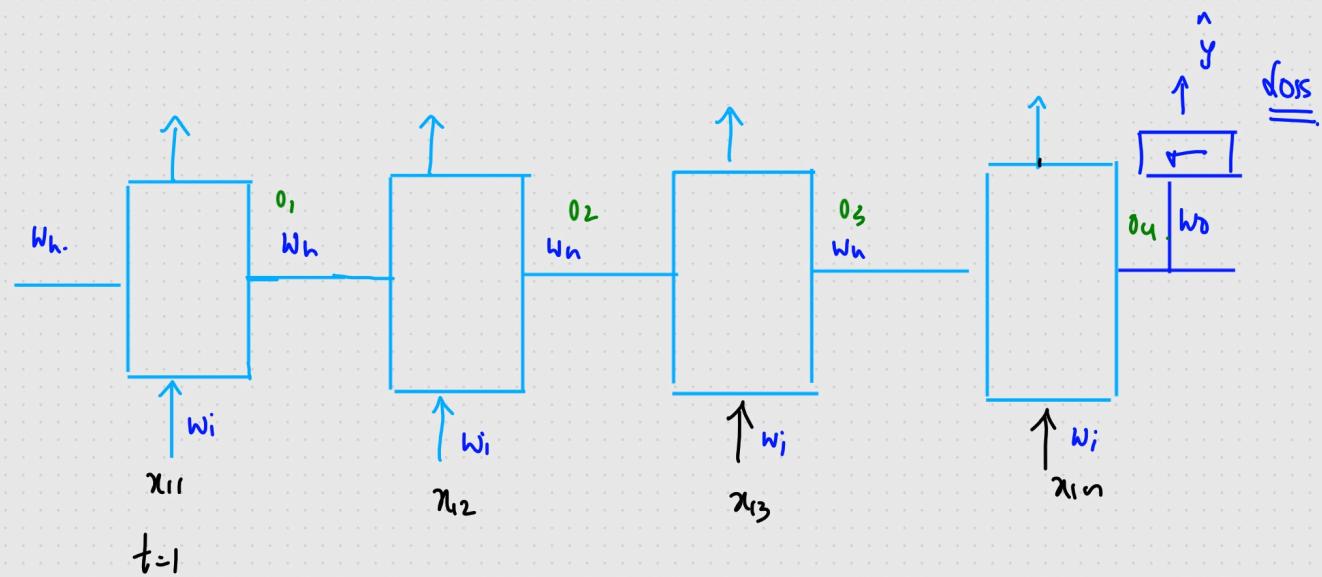
$h_t' \Rightarrow$ New Info

Bidirectional RNN

① Simple RNN → practical Implement → Embedding layers

② LSTM, GRU Variants RNN → Practical Implement

③ Bidirectional RNN / LSTM RNN



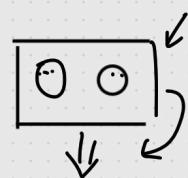
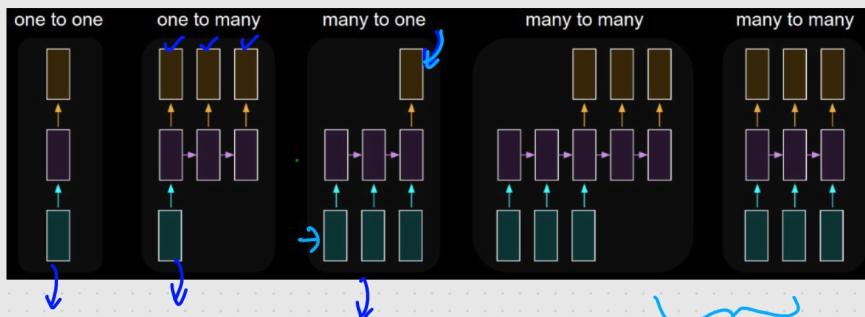
Types of RNN

① One to Many RNN

② Many to One RNN

③ Many to Many RNN

④ One to One RNN



Eg:

Eg: Image
Captioning

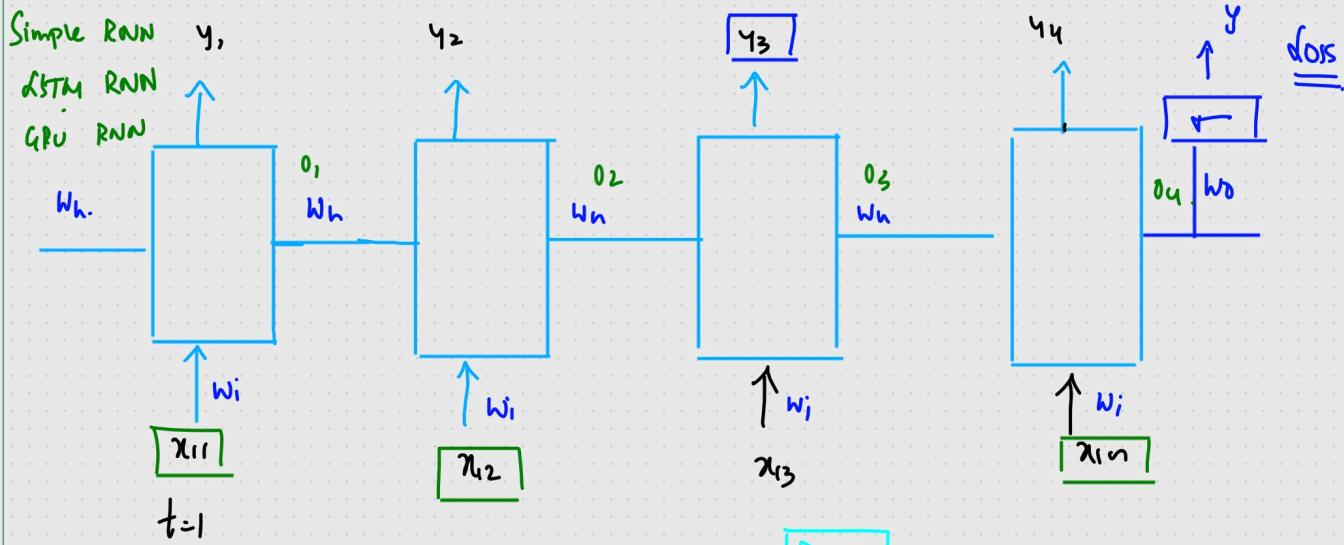
There is a dog and cat

One i/p
One o/p

One i/p
o/p Many

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

Language Translation
Many I/P
Many O/P.



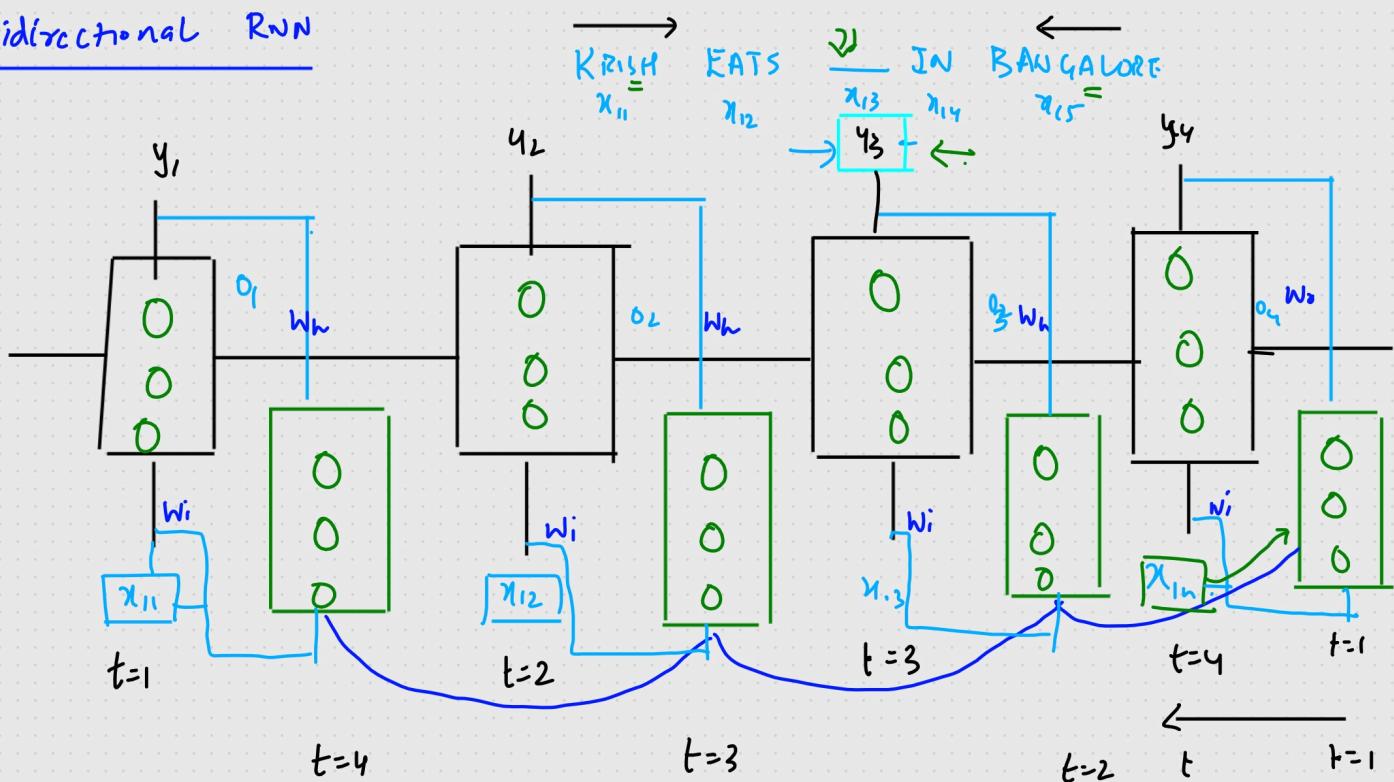
Text Example = KRISM ← EATS → IN → BANGALORE

KRISH ← EATS → IN → PARIS

PIZZA

DOSA

Bidirectional RNN



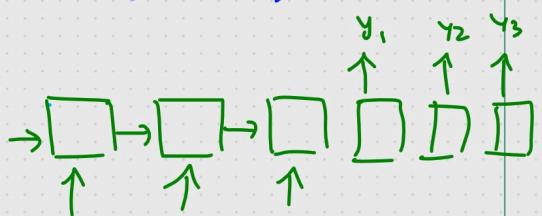
① Forward propagation → Equation ←

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 Others

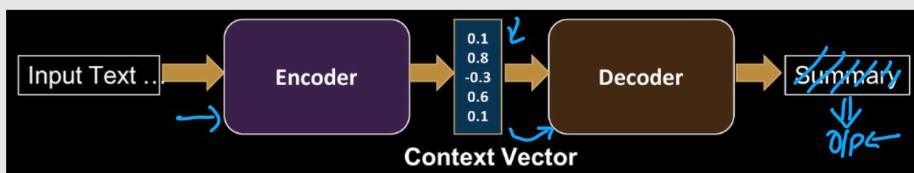
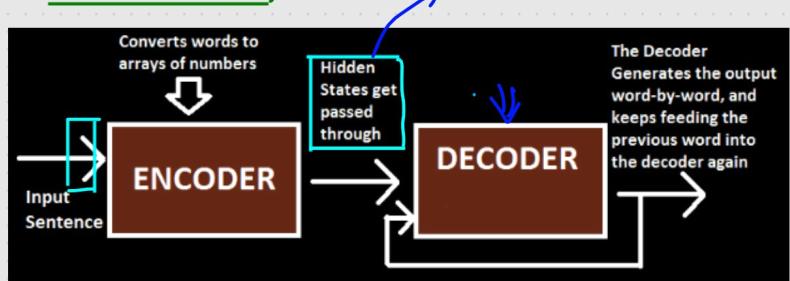
English → French

Eg: Filtered 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

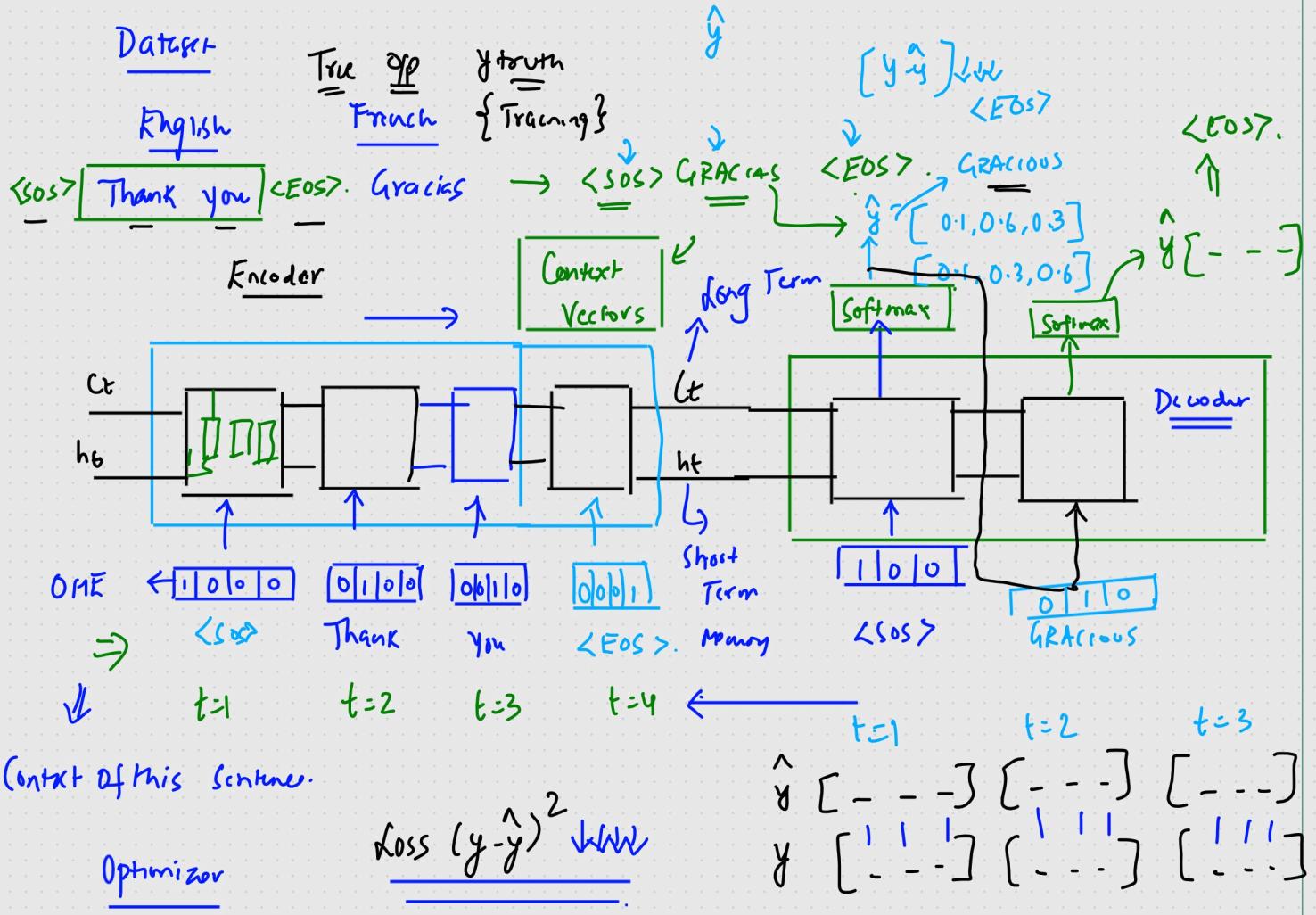
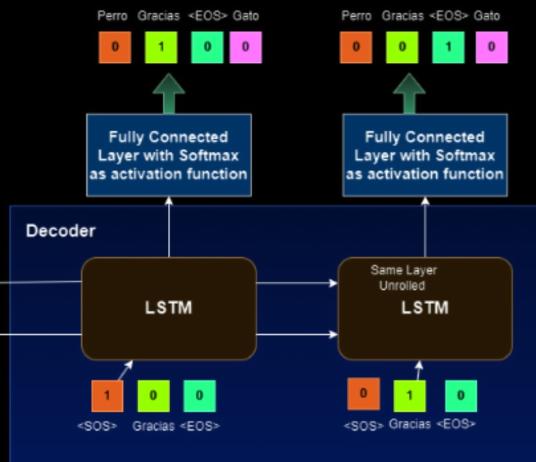
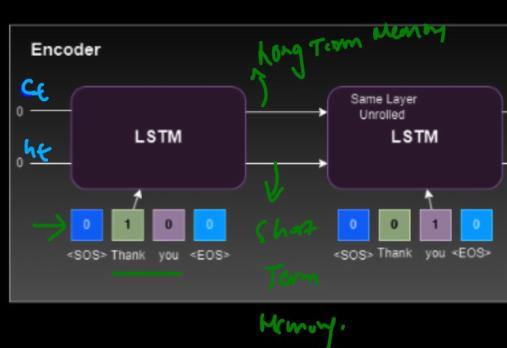
RNN → Vanishing Gradient Problem

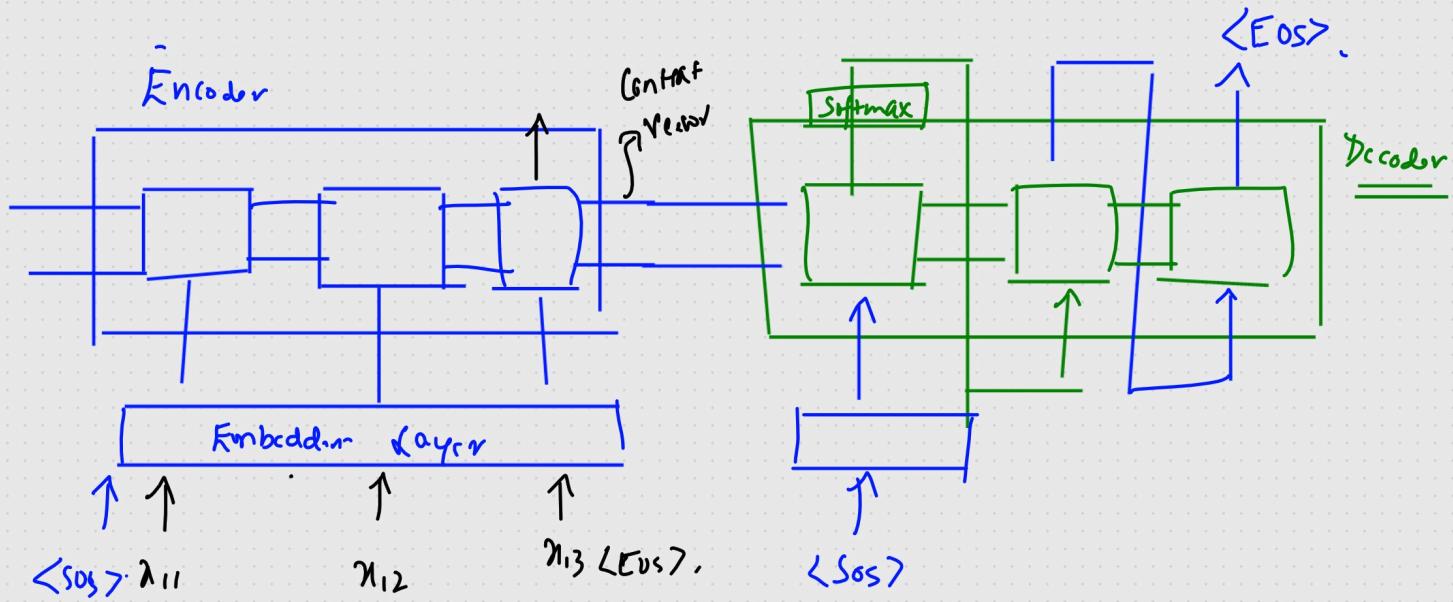
① forget gate

② I/p gate & candidate i/p

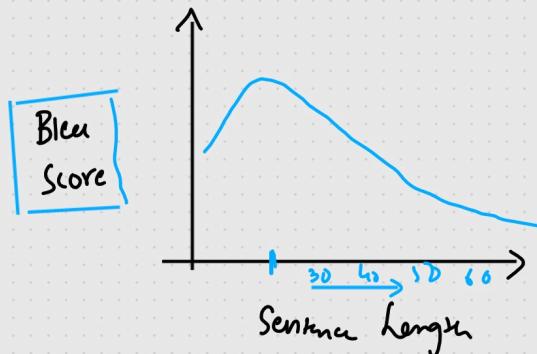
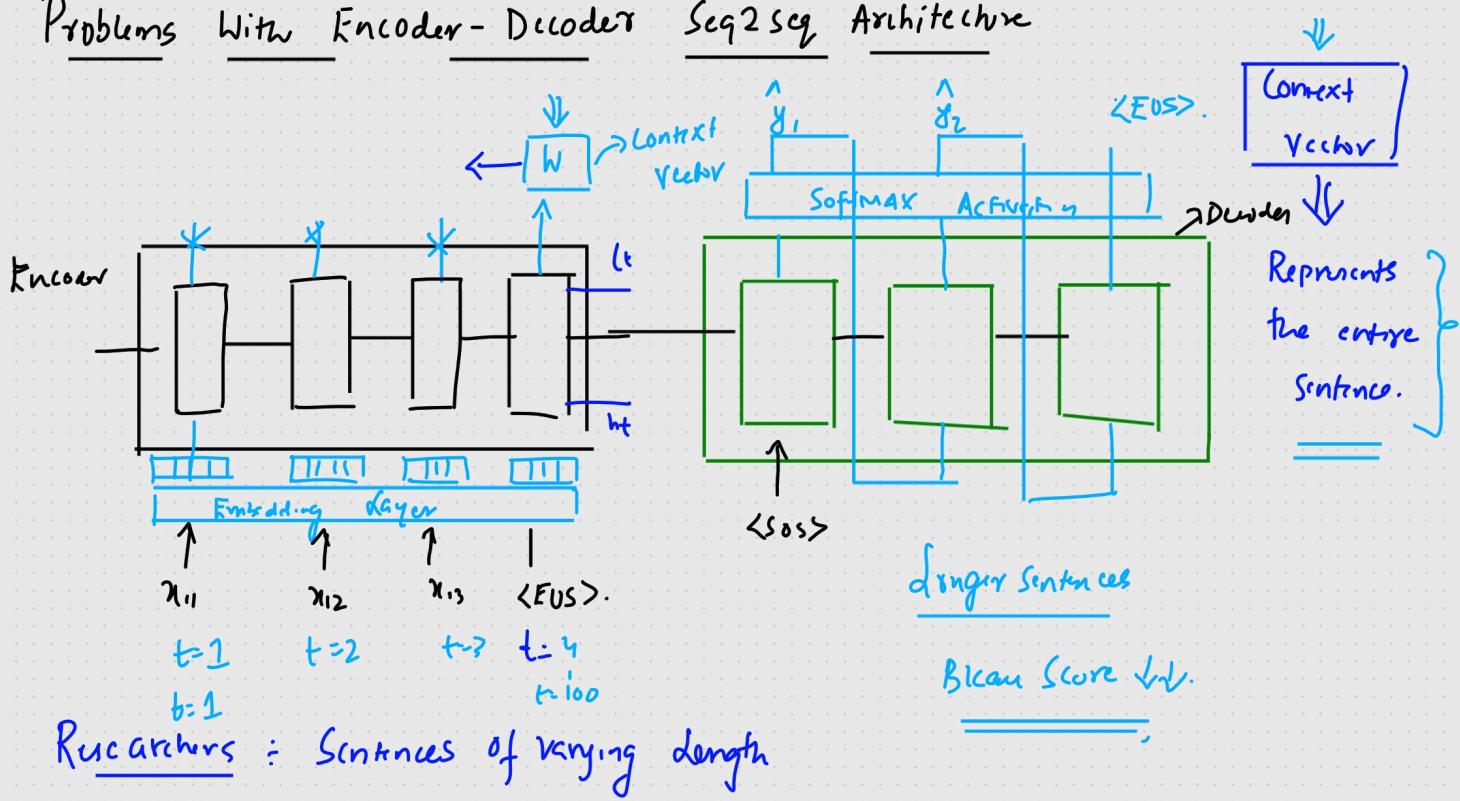
③ O/P

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





Problems With Encoder-Decoder Seq2Seq Architecture



\Rightarrow Seq to Seq Data.

* Attention Mechanism \rightarrow Seq2Seq Network

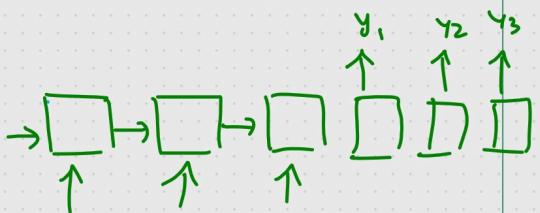
longer paragraph $\rightarrow \{ \text{Context Vector} \}$
+
 $\{ \text{Context} \}$

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 Others

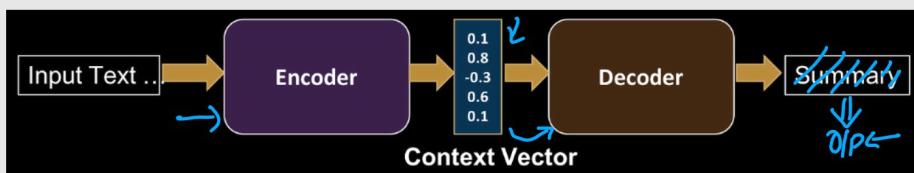
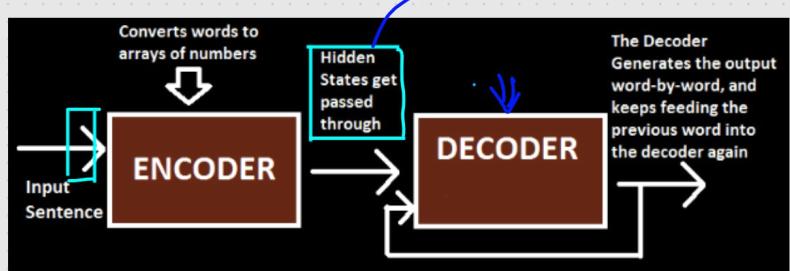
English → French

Eg:- Filtered 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

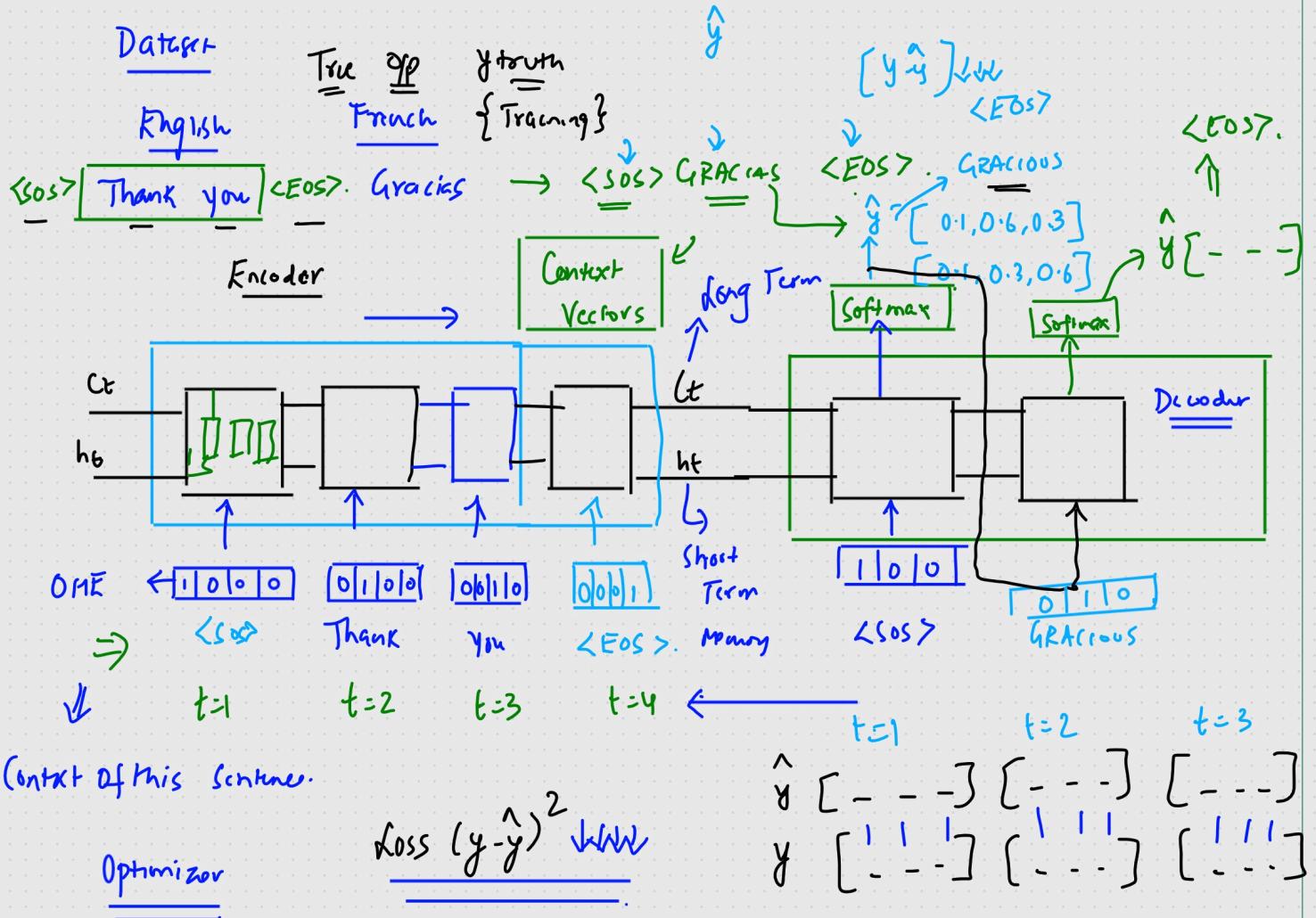
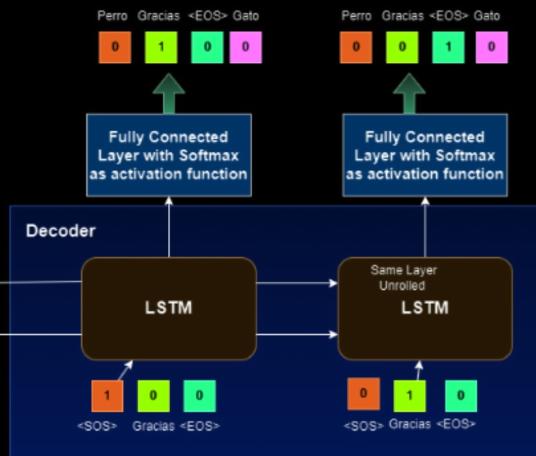
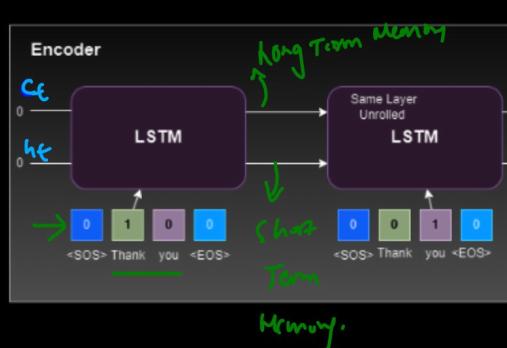
RNN → Vanishing Gradient Problem

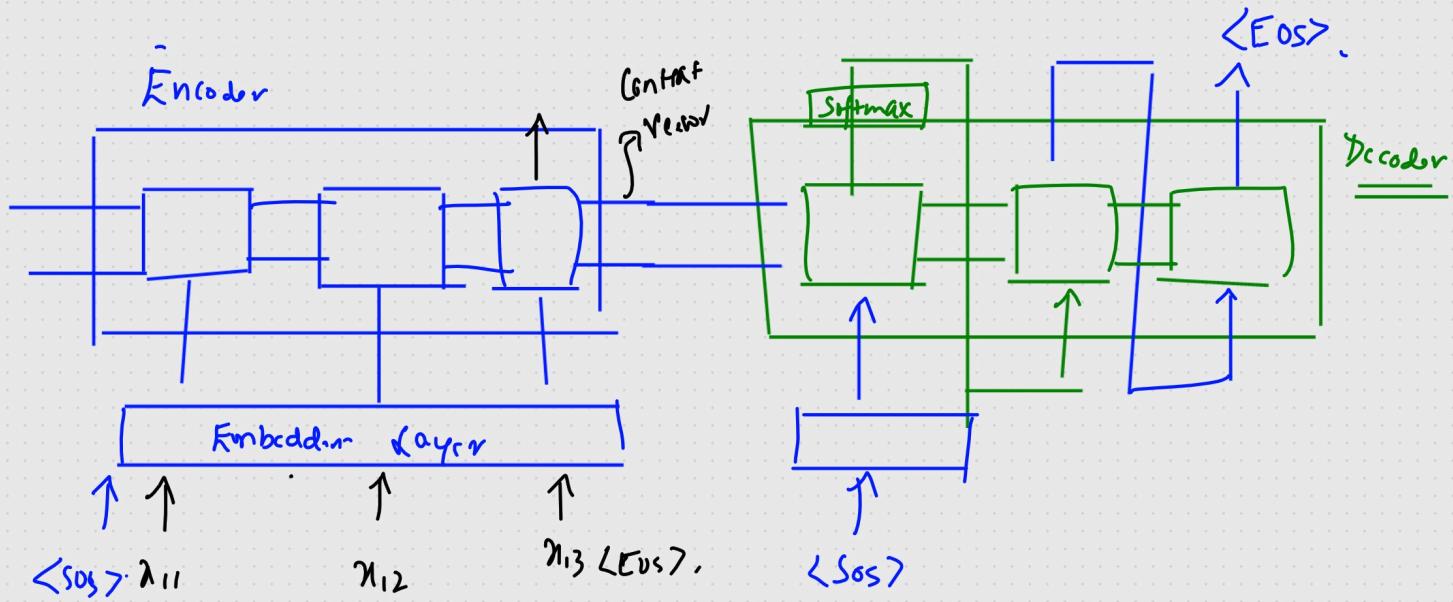
① forget gate

② I/p gate & candidate i/p

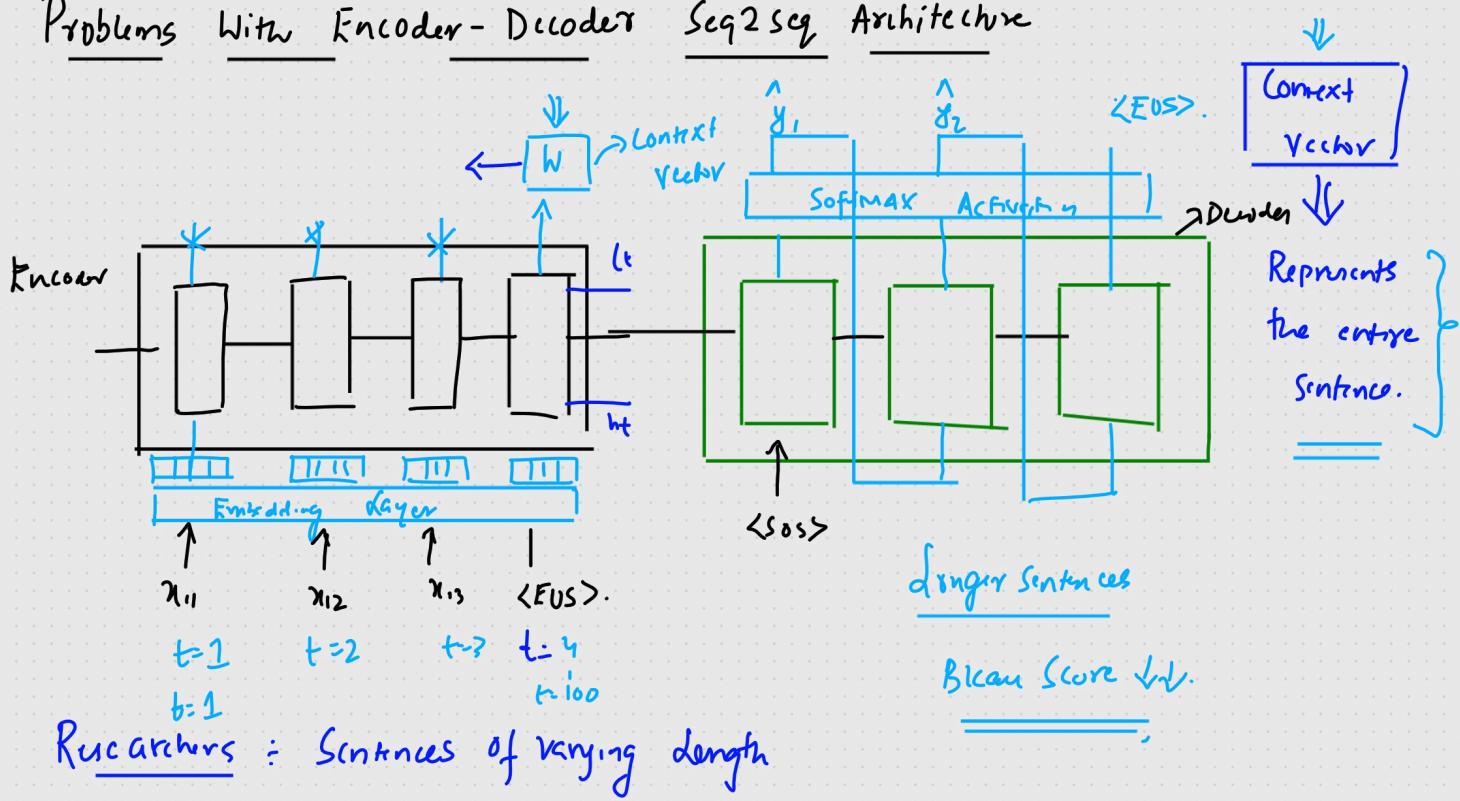
③ O/P

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



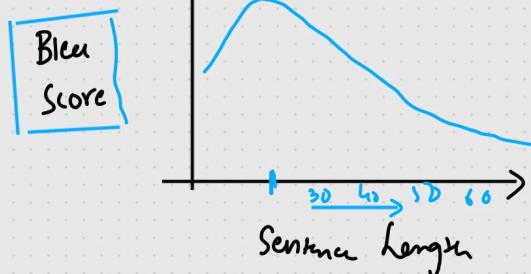


Problems With Encoder-Decoder Seq2Seq Architecture



Rucahars : Sentences of varying length

\Rightarrow Seq to Seq Data.



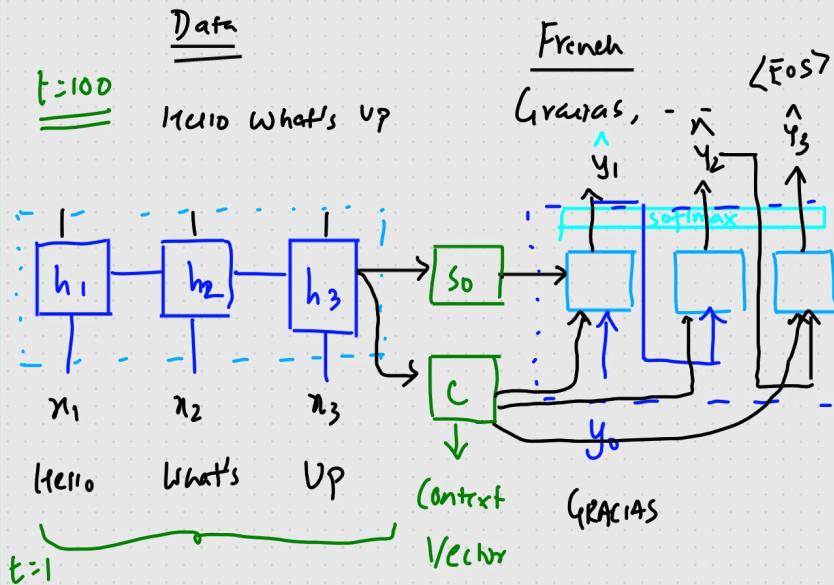
(*) Attention Mechanism → Seq2Seq Network

longer paragraph → {Context Vector}.

+

{Context}

Attention Mechanism | Seq2Seq Networks



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:

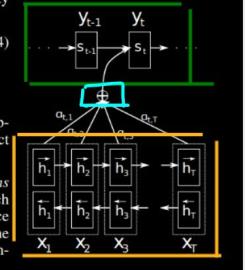
$$p(y_i|y_1, \dots, y_{i-1}, \mathbf{x}) = g(y_{i-1}, s_i, c_i), \quad (4)$$

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_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 :

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

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

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

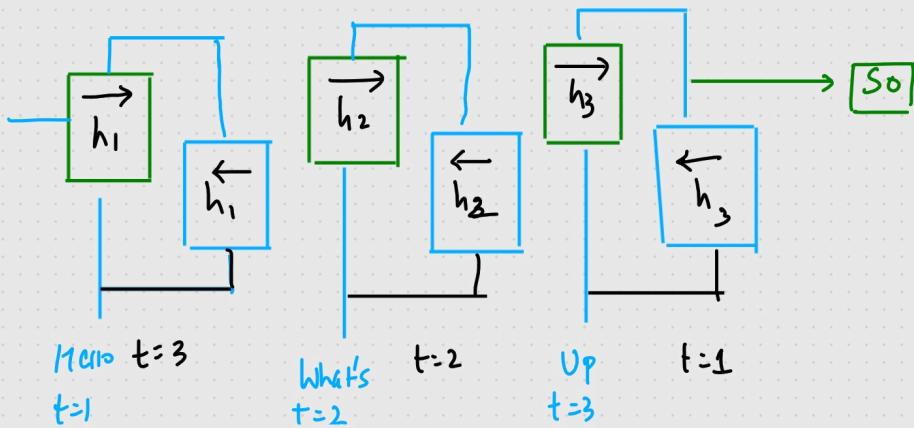
where

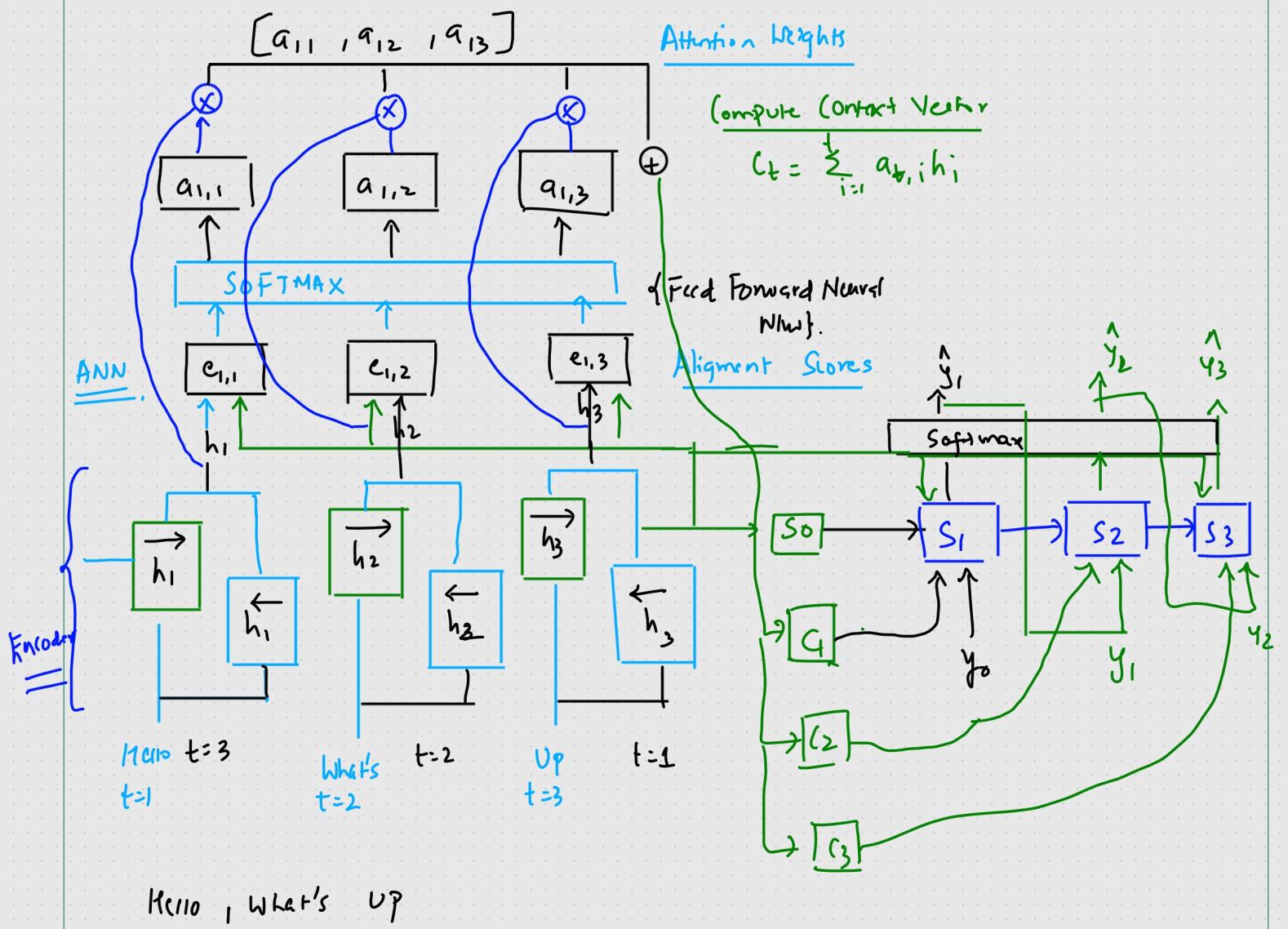
$$e_{ij} = a(s_{i-1}, h_j)$$

Encoder Decoder Architecture

Attention Mechanism

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





Introduction To Transformers

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

① Why Transformers?

② Architecture of Transformers?

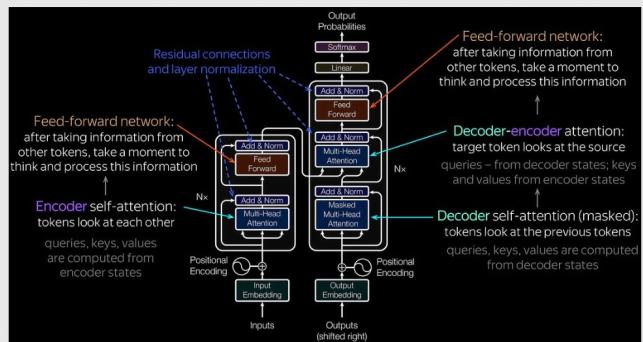
③ SELF ATTENTION $\rightarrow Q, K, V$

④ Positional Encoding

⑤ Multi Head ATTENTION

⑥ Combining the Working of Transformers

Architecture.



Generative AI \rightarrow LM, Multimodel

BERT, GPT \leftarrow

OpenAI \rightarrow ChatGPT

GPT-4o

① What And Why \rightarrow Transformers

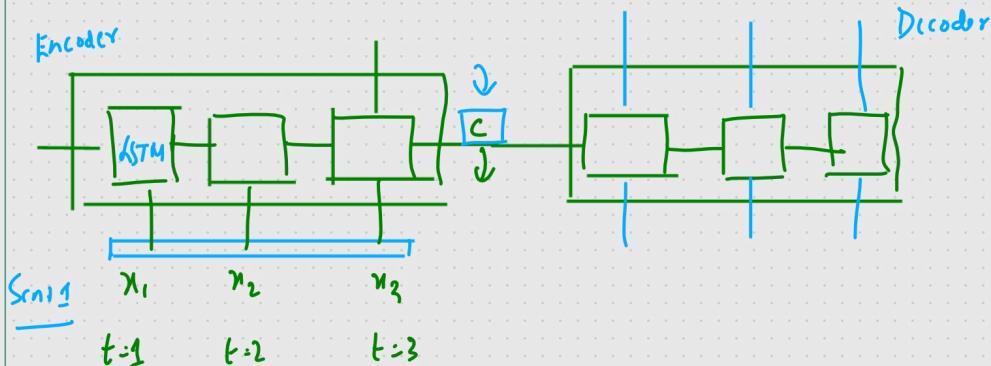
Transformers in natural language processing (NLP) are a type of deep learning model that use self-attention mechanisms to analyze and process natural language data. They are encoder-decoder models that can be used for many applications, including machine translation. \Rightarrow Seq2Seq Task

Eg: Language Translation \longrightarrow Google Translation

English \rightarrow French

i/p \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_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 :

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

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 \uparrow DATASET \rightarrow Amazing SOTA \leftarrow NLP

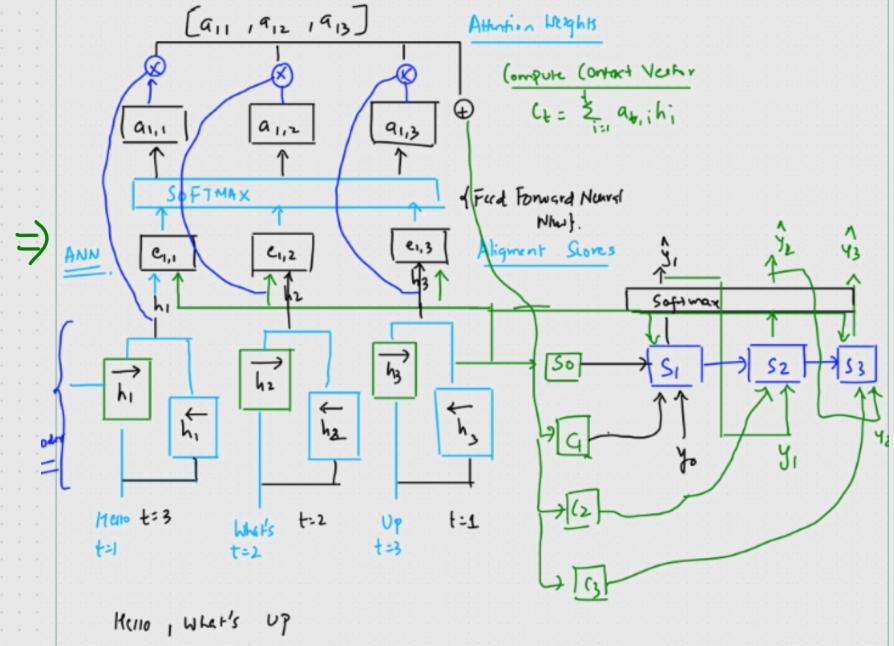
Transfer Learning \rightarrow MultiModal Task \rightarrow NLP + Image \leftarrow

Transformers \div AI Space \rightarrow SOTA Model \rightarrow



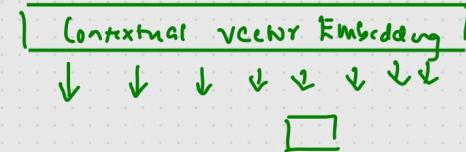
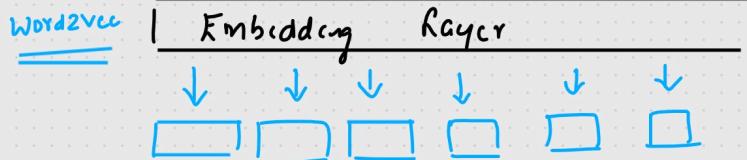
BERT GPT \rightarrow Transfer learning \rightarrow SOTA Models \rightarrow DALLE \leftarrow Generative AI

Train huge Data



② Contextual Embedding → Self Attention

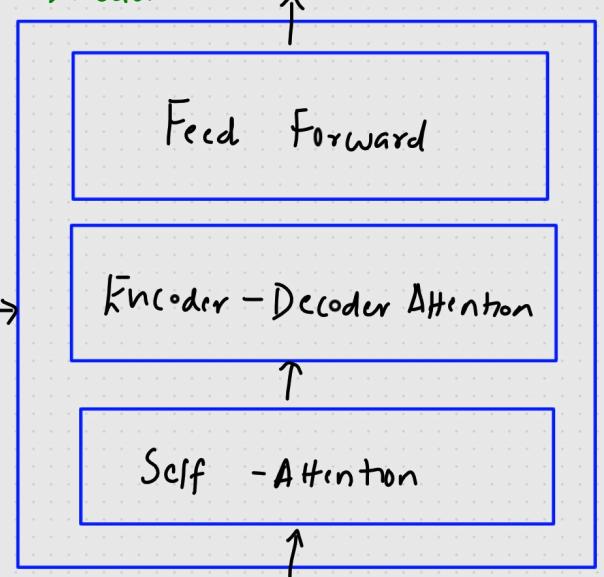
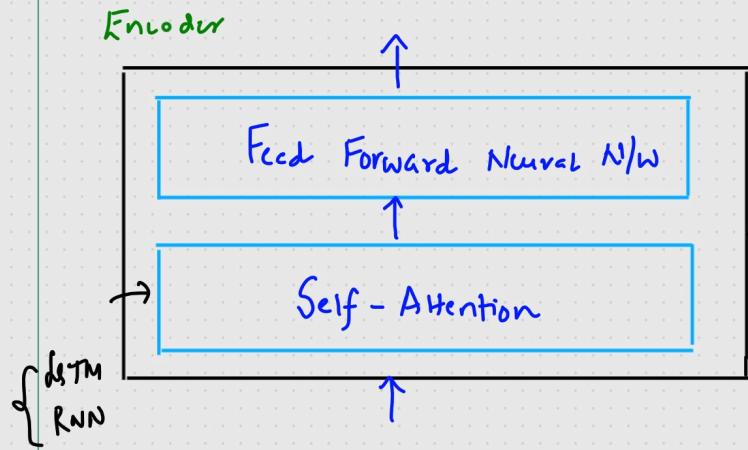
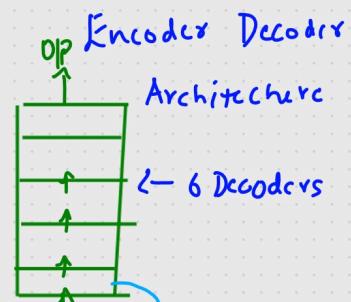
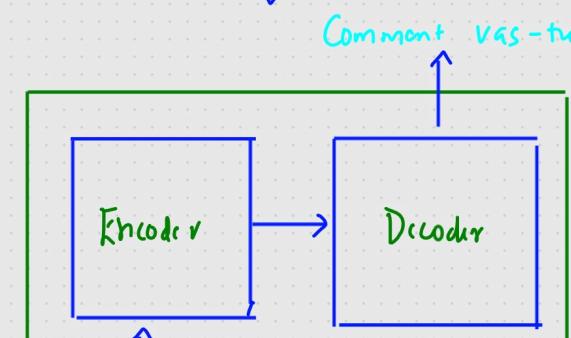
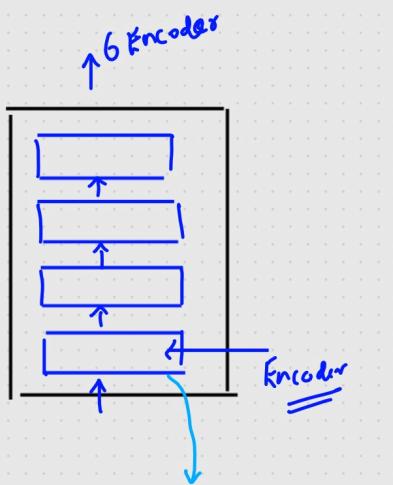
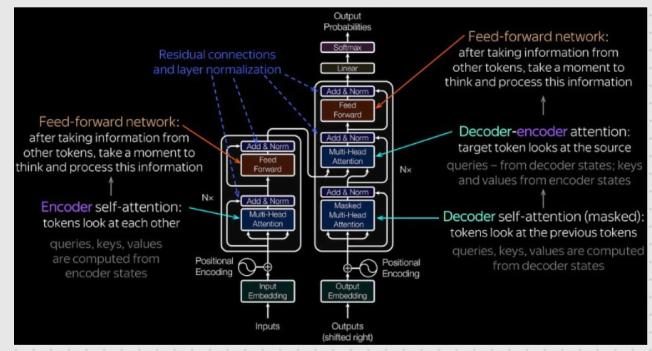
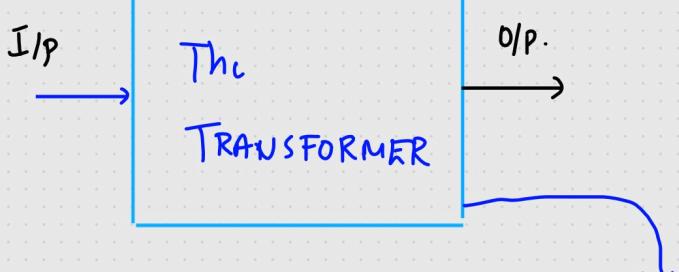
Eg: My name is KRISH And I play CRICKET



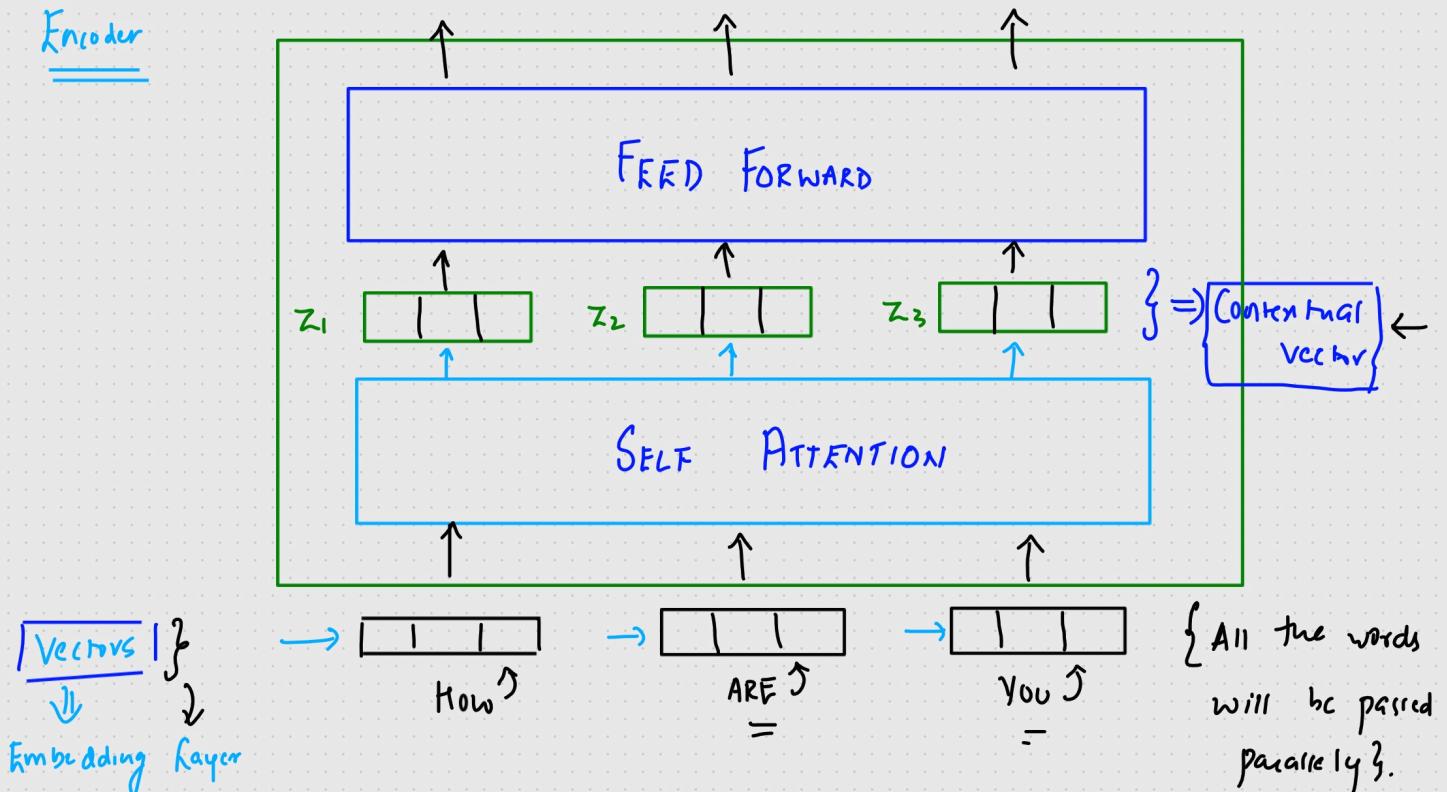
Contextual Vector

② Basic Transformer Architecture

{Seq2Seq Task} → Language Translation {Eng → French}

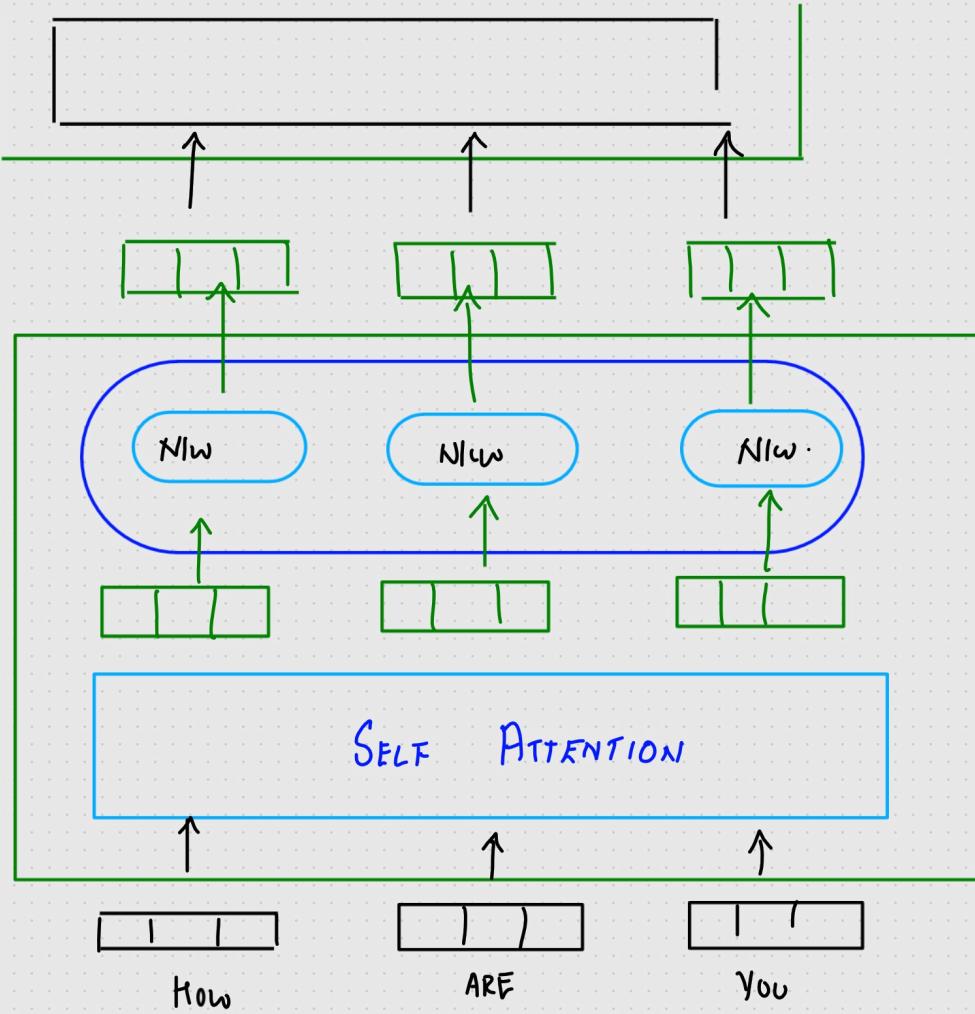


Encoder



Encoder²

Encoder¹

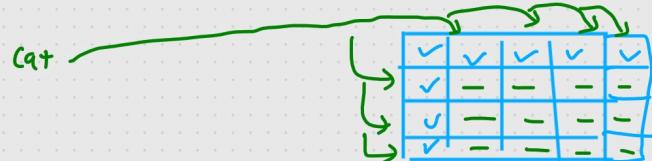
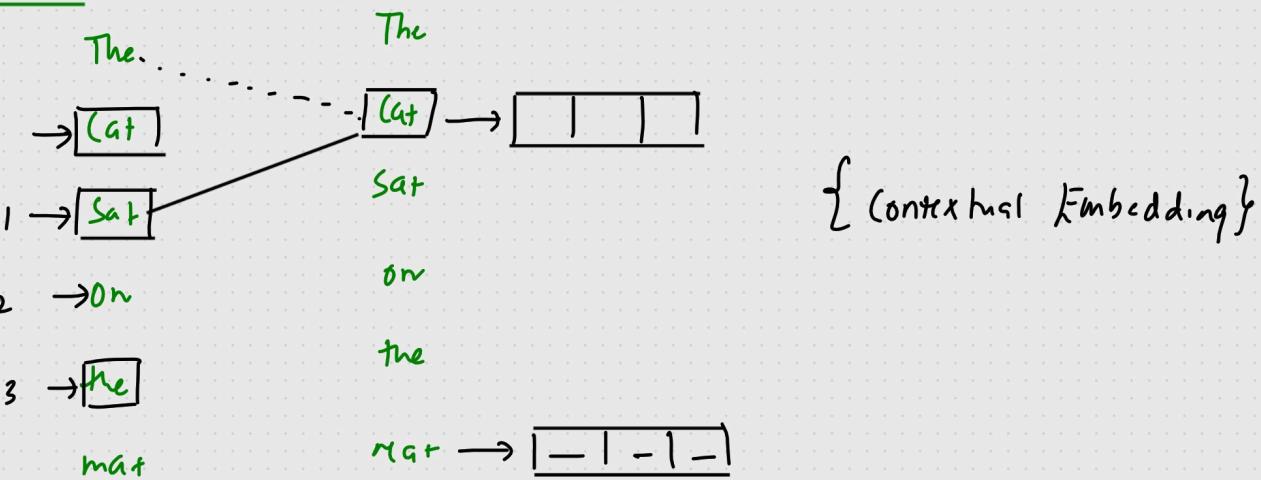


Self Attention At a Higher Level

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

Word Embedding
↓ ↓ ↓ ↓ ↓ ↓ ↓ ↓ ↓ ↓ ↓ ↓

SELF ATTENTION



Self Attention In Detail

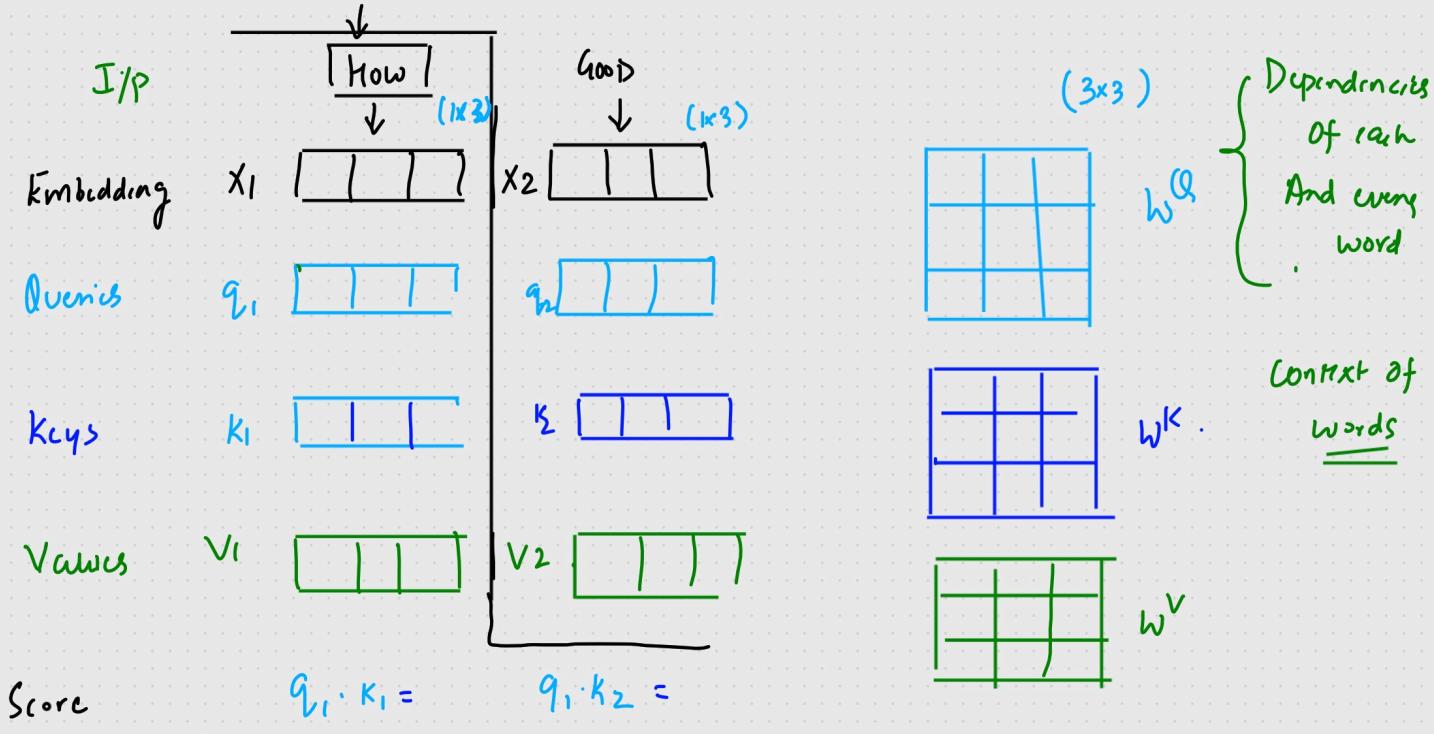
- To create 3 vectors from each of the encoder i/p. Query vector, Key vector, Value vector. \Rightarrow Contextual Embedding

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

Query \rightarrow {Key} $\xrightarrow{\text{Tags}}$ Tags
 $\xrightarrow{\text{Description}}$ Description $\xrightarrow{\text{Title}}$ Title \Rightarrow Value \Rightarrow O/P Video

- 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.



④ Focus to place more ^{How much} on other part of the I/P Sentence as we encode a word at certain position

Divide \sqrt{dk}

More stable gradients

Soft+Max

Dimension=3 Dimension $\frac{1}{\sqrt{dk}}$

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

$\begin{bmatrix} 0-1 \end{bmatrix} \quad \begin{bmatrix} 0-1 \end{bmatrix} \cdot \begin{bmatrix} 0-1 \end{bmatrix}$

$0.88 + 0.12 = 1$

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

x

Value Vectors $V_1 \begin{bmatrix} | & | \end{bmatrix}$

$V_2 \cdot \begin{bmatrix} 1 & 0.88 & 0.12 \end{bmatrix}$

$0.88 \begin{bmatrix} | & | & | \end{bmatrix}$

$\begin{bmatrix} 0.88 & 0.88 & 0.88 \end{bmatrix} + \begin{bmatrix} 0.12 & 0.12 & 0.12 \end{bmatrix}$

$\xrightarrow{T_{11}} \begin{bmatrix} | & | \end{bmatrix}$

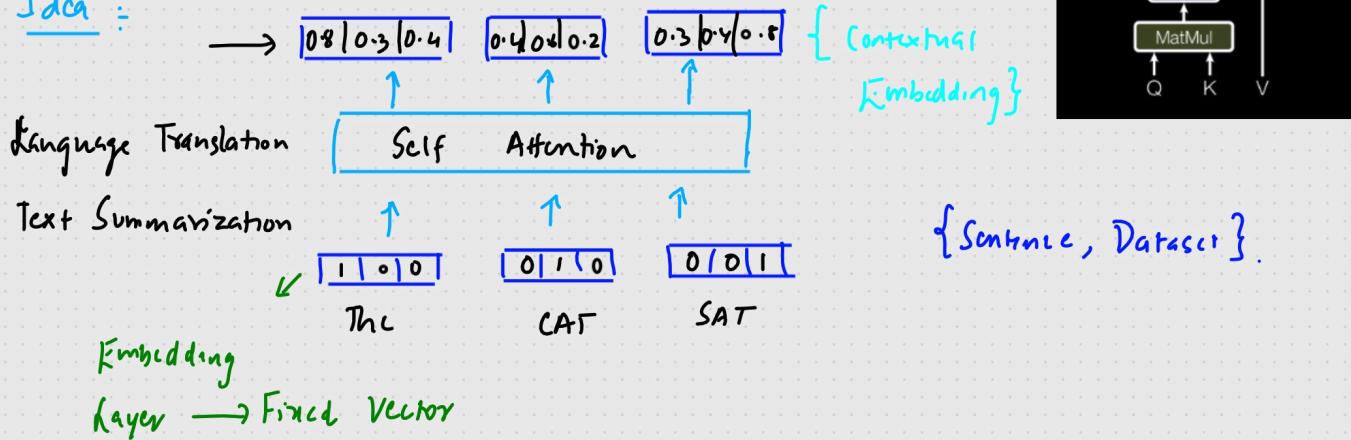
$\xrightarrow{T_{22}} \begin{bmatrix} | & | & | \end{bmatrix}$

Contextual Vector

Self Attention At Higher And Detailed Level

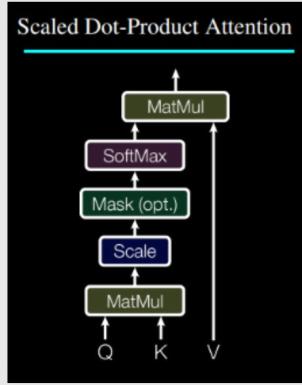
Self-attention, also known as scaled dot-product attention, is a crucial mechanism in the transformer architecture that allows the model to weigh the importance of different tokens in the input sequence relative to each other

Idea :



1) Inputs: Queries, Keys, And Values

Model → 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:

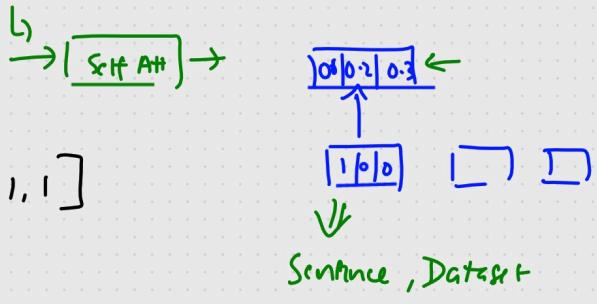
Information Aggregation: Values contain the data that will be weighted by the attention scores. The weighted sum of values forms the output of the self-attention mechanism, which is then passed on to the next layers in the network.

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.

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

$$\text{Embedding Size} = 4$$

$$Q, K, V \Rightarrow 4$$



① Token Embedding

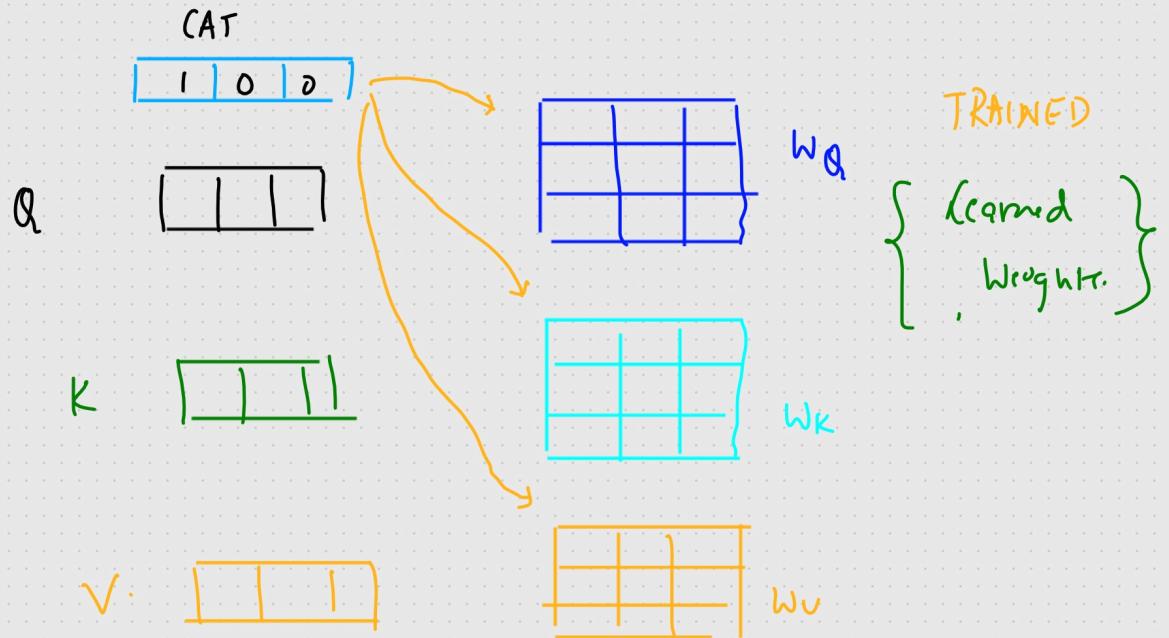
$$E_{\text{The}} = [1 \ 0 \ 1 \ 0]$$

$$E_{(\text{AT})} = [0 \ 1 \ 0 \ 1]$$

$$E_{\text{SAT}} = [1, 1, 1, 1]$$

② Linear Transformation

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



Let's consider

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

$$Q_{\text{The}} = [1 \ 0 \ 1 \ 0] \cdot \begin{bmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \end{bmatrix} = [1 \ 0 \ 1 \ 0]$$

$$K_{\text{The}} = [1 \ 0 \ 1 \ 0]$$

$$V_{\text{The}} = [1 \ 0 \ 1 \ 0]$$

$$\textcircled{1} \quad Q_{\text{The}} = K_{\text{The}} = V_{\text{The}} = [1 \ 0 \ 1 \ 0]$$

$$\textcircled{2} \quad Q_{\text{CAT}} = K_{\text{CAT}} = V_{\text{CAT}} = [0 \ 1 \ 0 \ 1]$$

$$\textcircled{3} \quad Q_{\text{SAT}} = K_{\text{SAT}} = V_{\text{SAT}} = [1, 1, 1, 1]$$

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

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

$$\text{Score}(Q_{\text{SAT}}, K_{\text{CAT}}) = 2$$

$$\text{Score}(Q_{\text{SAT}}, K_{\text{SAT}}) = 4$$

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

$$[[10, 1]] = [0.99, 0.01]$$

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

$$\approx [0.88, 0.12].$$

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

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

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

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

will be done for
all other tokens.

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

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

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

$$⑤ \text{Apply Softmax}$$

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

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

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

$$⑥ \text{Weight Sum of Values}$$

We multiply the attention weights by corresponding value vectors

For the Token The =

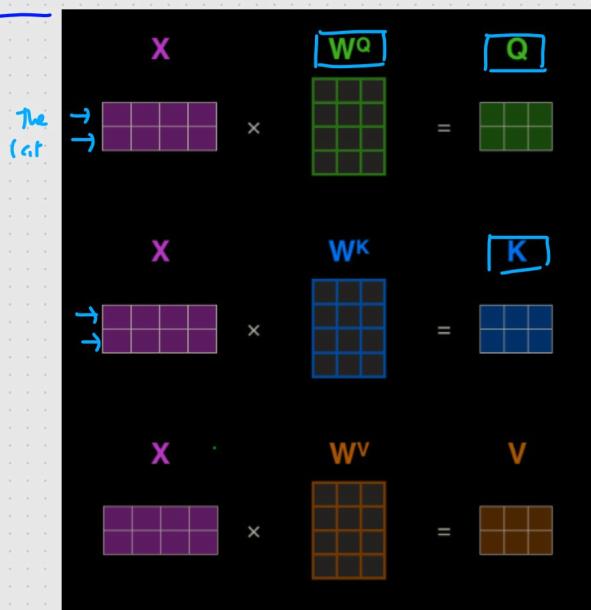
$$\begin{aligned}\text{Output}_{(\text{the})} &= 0.4223 * V_{\text{The}} + 0.1554 * V_{\text{CAT}} + 0.4223 * V_{\text{Sat}} \\ &= 0.4223 [1 \ 0 \ 1 0] + 0.1554 [0 \ 1 \ 0 1] + 0.4223 [1 \ 1 \ 1 \ 1] \\ &= [0.4223, 0, 0.4223, 0] + [0, 0.1554, 0, 0.1554] + [0.4223, 0.4223, \\ &\quad 0.4223, 0.4223] \\ &= [1.2669, 0.9999, 1.2669, 0.9999].\end{aligned}$$

↙ Contextual
vector v

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).

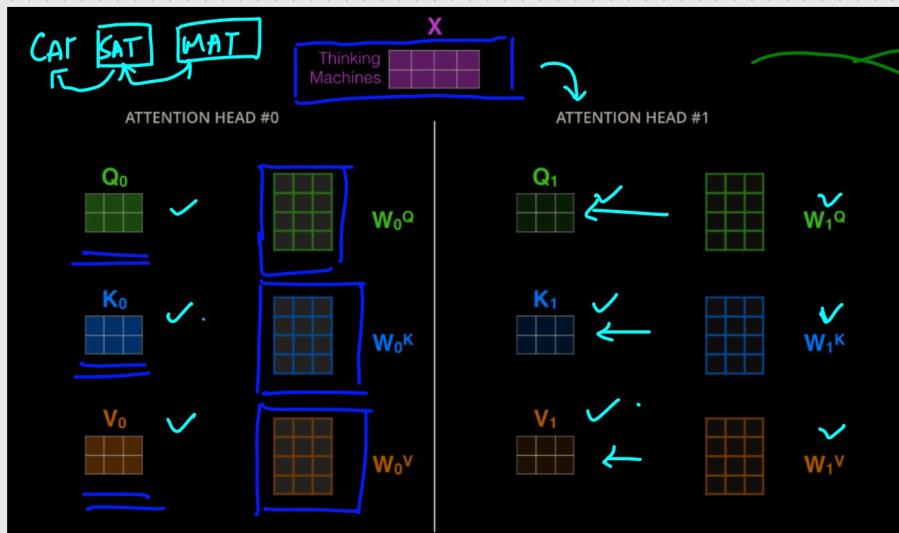
④ Multi Head Attention



$$\text{softmax} \left(\frac{\mathbf{Q} \times \mathbf{K}^T}{\sqrt{d_k}} \right) \cdot \mathbf{V}$$

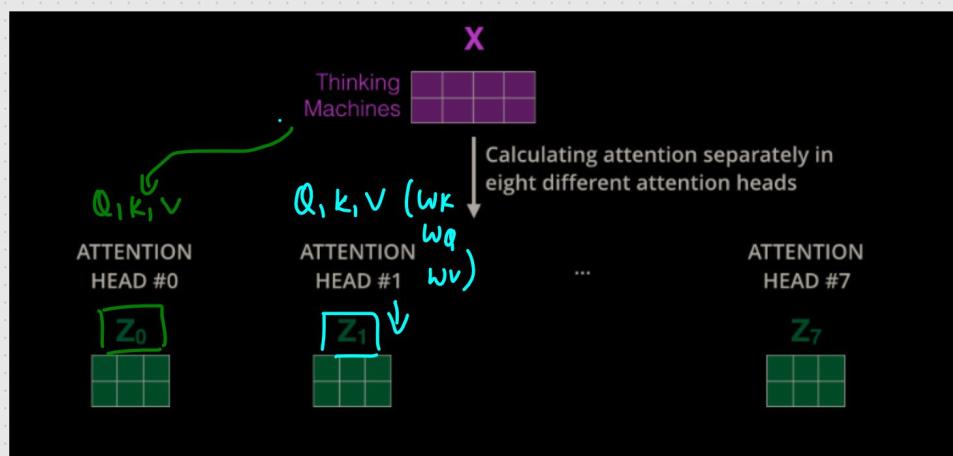
↳ Scaling
↳ Vectors
Attention Head.

→ Self Attention with Multi Heads



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

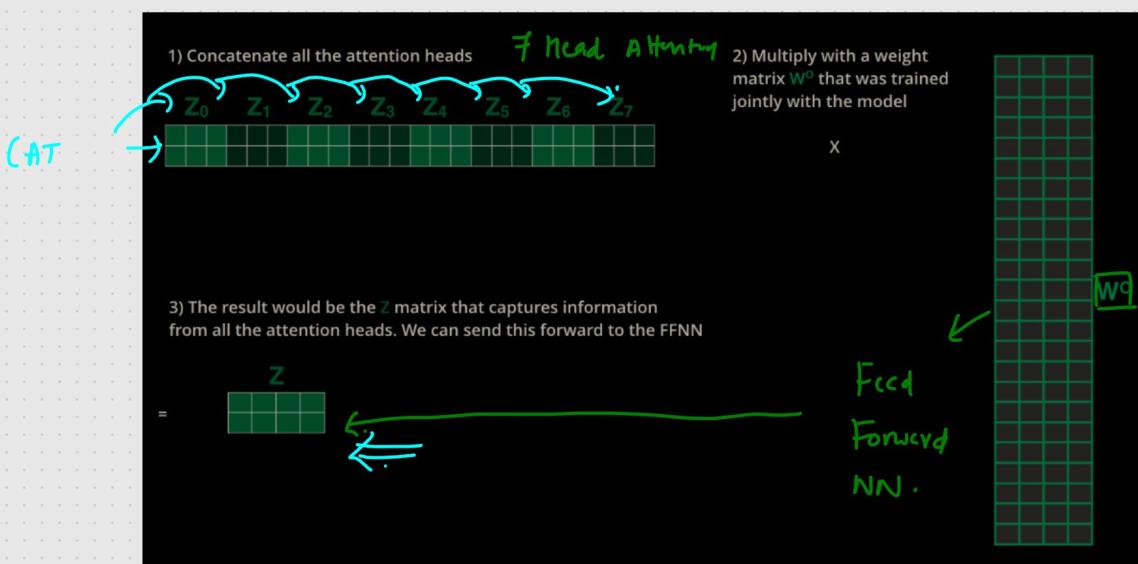
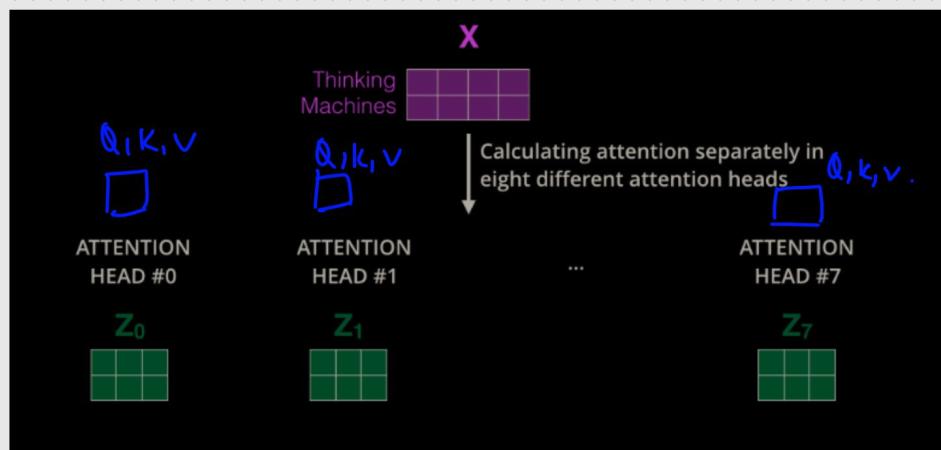
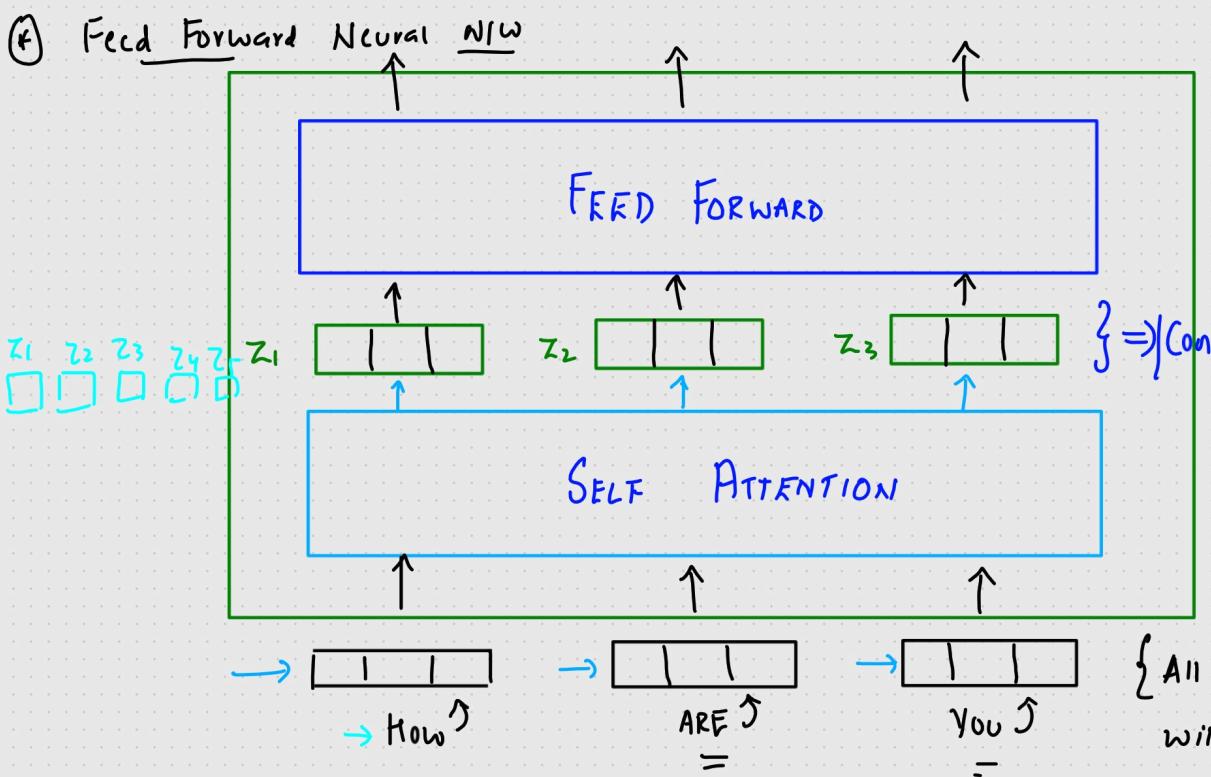
Score, Softmax $\times V$
 \downarrow \downarrow \downarrow \downarrow
 Z_0 $Vectors$ Z_1
Multi Head Attention

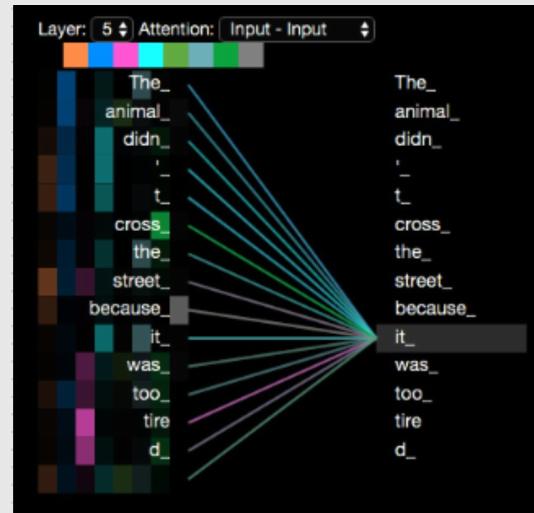
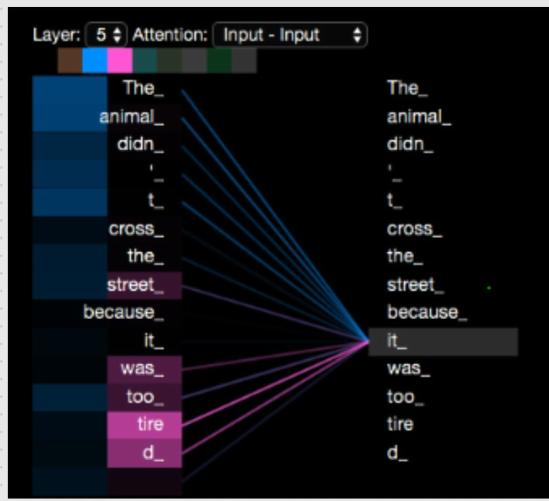
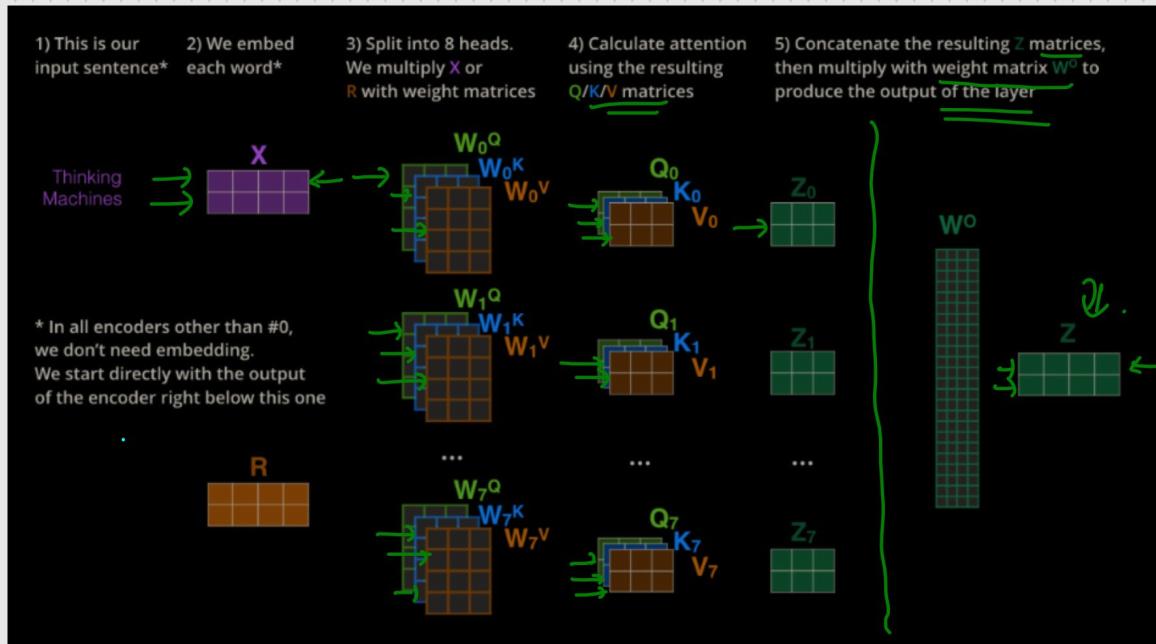


$[Z]$

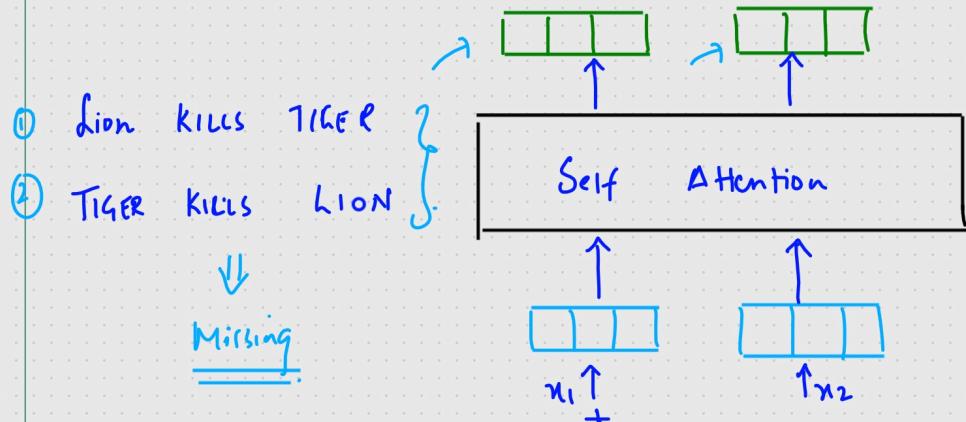
$[Z]$

$[Z]$





⑥ Positional Encoding - Representing Order of Sequence



Advantages

- 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

Novel

↳ 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

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

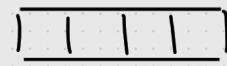
d_{model} is the dimensionality of the embeddings.

Eg: The Cat Sat

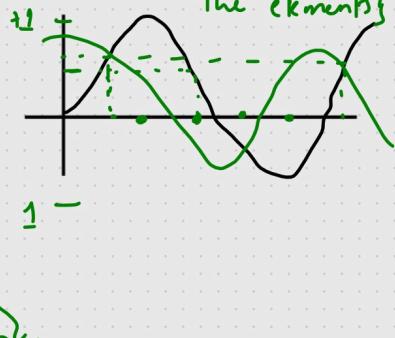
$$\text{The} \rightarrow [0.1 \ 0.2 \ 0.3 \ 0.4]$$

$$\text{CAT} \rightarrow [0.5 \ 0.1 \ 0.7 \ 0.8]$$

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



{ Miss the order of the elements }



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

For our example $d_{model} = 4$

For position $pos = 0$

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

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

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

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

$$P.E = [0, 1, 0, 1]$$

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

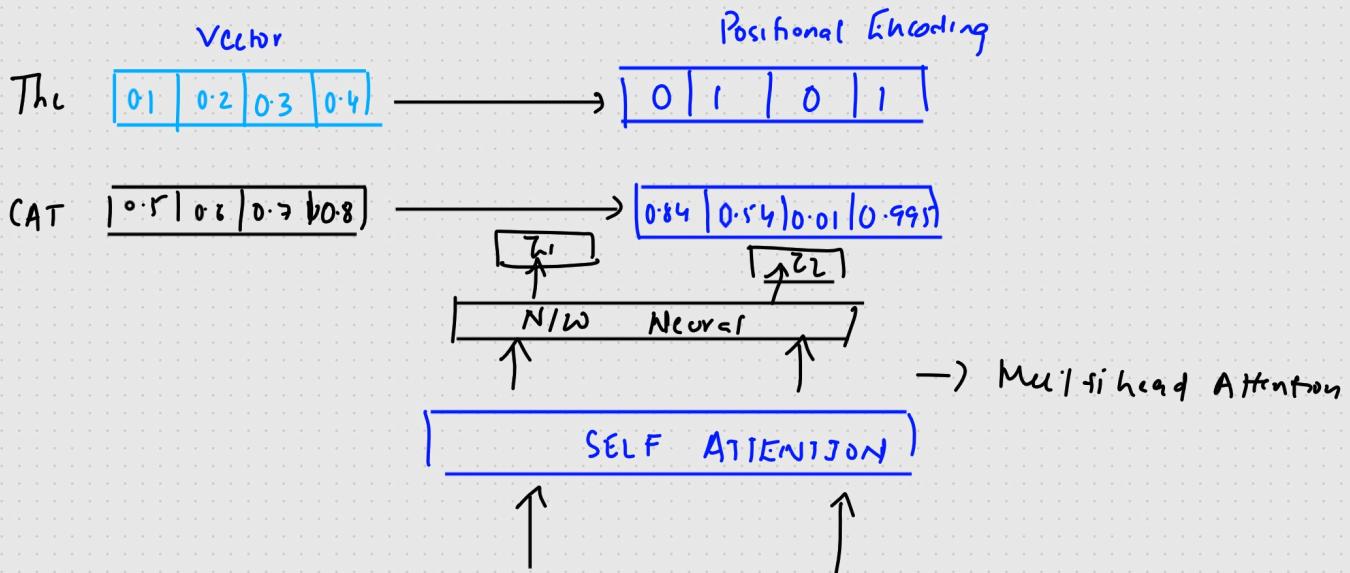
For $pos = 1$

$$P.E(1,0) = \sin\left(\frac{1}{10000^0/4}\right) = \sin(1) = 0.8415 \quad P.E(pos, 2i+1) = \cos\left(\frac{pos}{10000^{2i+1}/d_{model}}\right)$$

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

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

$$P.E(1,3) = \cos\left(\frac{1}{10000^3/4}\right) \approx 0.99995$$



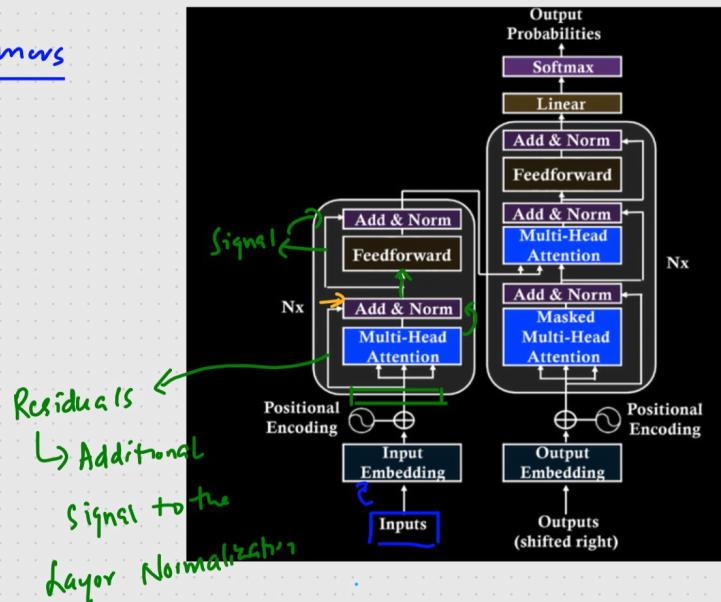
Positional Encoding

$$\begin{array}{c}
 \text{Positional Encoding} \\
 \left\{ \begin{array}{|c|c|c|c|c|} \hline 0.84 & 0.34 & 0.01 & 0.99 \\ \hline \end{array} \right\} + \begin{array}{|c|c|c|c|c|} \hline 0.5 & 0.6 & 0.7 & 0.8 \\ \hline \end{array} \\
 \uparrow \quad \uparrow \\
 \text{CAT} \quad \text{Thc.}
 \end{array} \quad \leftarrow \text{Positional Encoding}$$

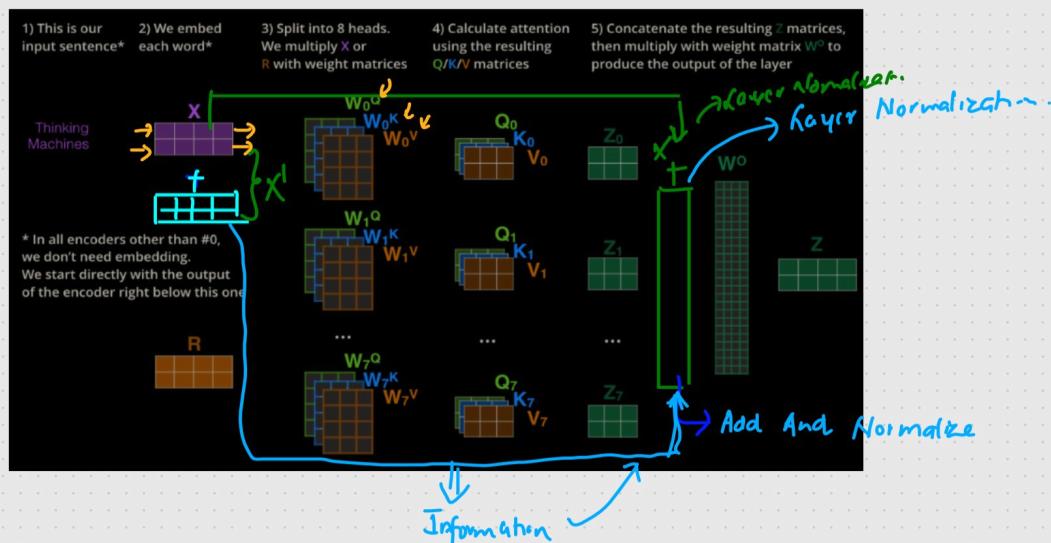
④ Layer Normalization In Transformers

Transformers

Residuals

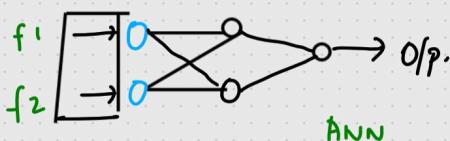


- ① Self Attention layer
 - ② Multi head Attention
 - ③ Positional Encoding
 - ④ Layer Normalization
- ADD AND Normalize



Normalization

Batch Normalization
Layer Normalization



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

Normalization : Standard Scaling

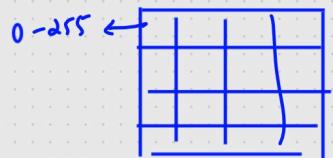
$$\mu, \sigma \leftarrow | 2000 |$$

5.5

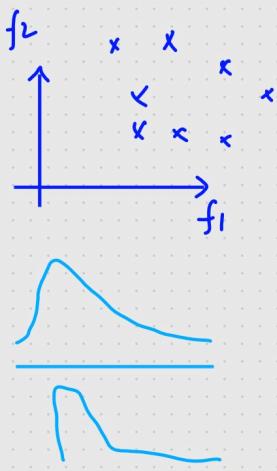
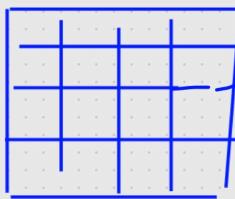
80

$$z_{score} = \frac{x_i - \mu}{\sigma} \Rightarrow \boxed{\mu=0, \sigma=1} \quad f_1 \rightarrow f'_1 \Leftarrow \mu=0, \sigma=1$$

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



\Rightarrow Min Max Scaler \Rightarrow



Advantages

① Improved Training Stability

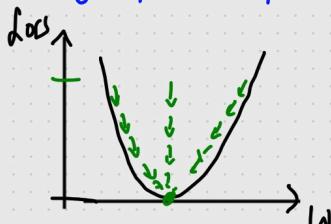


Vanishing and exploding Gradient problem

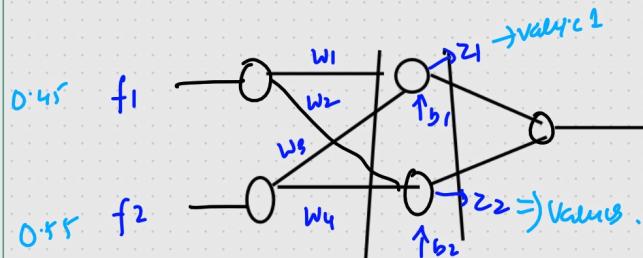
$$\begin{cases} \mu=0 \\ \sigma=1 \end{cases}$$



② Faster Convergence



\Rightarrow Back propagation \Rightarrow Stable update.



$$z_1 = \sigma[(0.45 \cdot w_1 + 0.55 \cdot w_3) + b_1] = \text{Value 1} \quad \text{Does this distribution} \\ z_2 = \sigma[(0.45 \cdot w_2 + 0.55 \cdot w_4) + b_2] = \text{Value 2} \quad \boxed{\mu=0, \sigma=1} \quad \begin{matrix} \mu_1, \sigma_1 \\ \mu_2, \sigma_2 \end{matrix}$$

{Batch Normalization}

Normalizes (standardizes)

$$\begin{matrix} z_1 & z_2 \\ \mu_1 & \mu_2 \\ \sigma_1 & \sigma_2 \end{matrix}$$

Batch Normalization VS Layer Normalization

f1	f2	z_1	z_2
-	-	-	-
-	-	-	-
-	-	-	-

$$z_{score} = \frac{x_i - \mu_1}{\sigma_1}$$

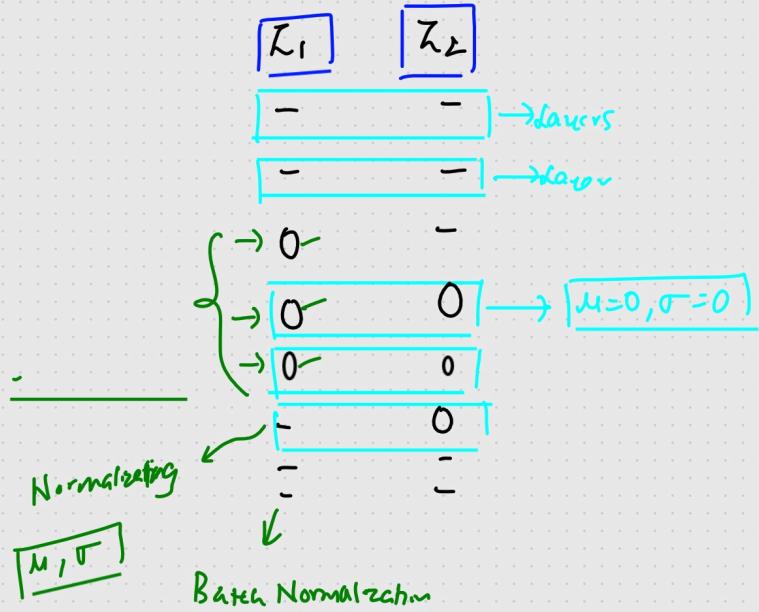
$$z_{score} = \frac{x_i - \mu_2}{\sigma_2}$$

$$z_{score} = \frac{x_i - \mu_3}{\sigma_3}$$

Layer

Normalization

$\gamma, \beta \rightarrow$ learnable parameters



Normalization



γ, β

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

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

Normalized.

Scale And Shift parameters

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

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

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

\Rightarrow Scale And Shift param

i) Compute the mean

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

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

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

ii) Compute the variance (σ^2)

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

iii) Normalize the i/p

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

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

$$\hat{x}_1 = \frac{2.0 - 5.0}{2.236} \approx -1.34$$

Normalized vector

$$\hat{x}_2 = \frac{4.0 - 5.0}{2.236} \approx -0.45 \quad \hat{x} = [-1.34, -0.45, 0.45, 1.34]$$

$$\hat{x}_3 = \frac{6.0 - 5.0}{2.236} \approx 0.45$$

$$\hat{x}_4 = \frac{8.0 - 5.0}{2.236} \approx 1.34.$$

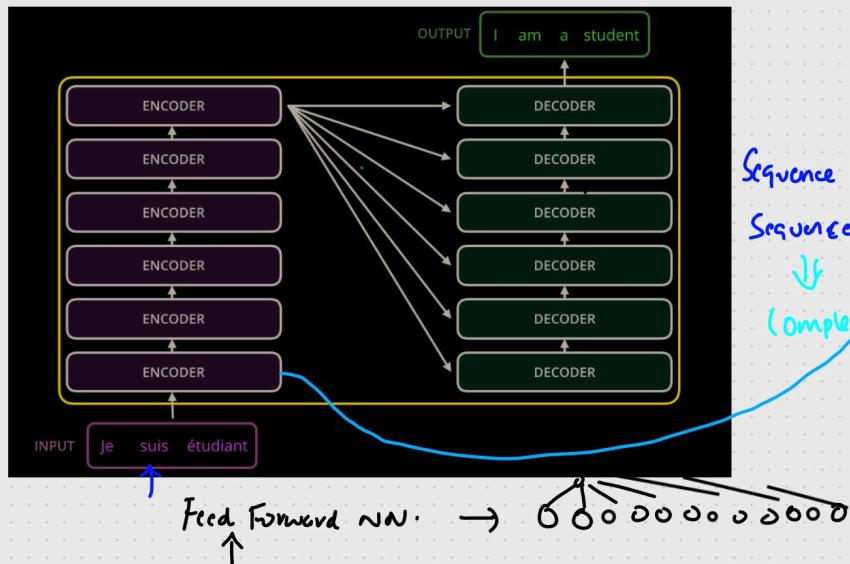
4) Scale And Shift

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

$$r = [1.0, 1.0, 1.0, 1.0] \quad \beta = [0.0, 0.0, 0.0, 0.0].$$

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

① Encoder Architecture [Research Paper]



Residuals
 Text Embeddings + Positional Encoding
 J/I/P Sequence
 $\Rightarrow 512 \rightarrow \{\text{Research paper}\}$.
 Every word = 512.
 $\sqrt{512} = 8\sqrt{2}$.

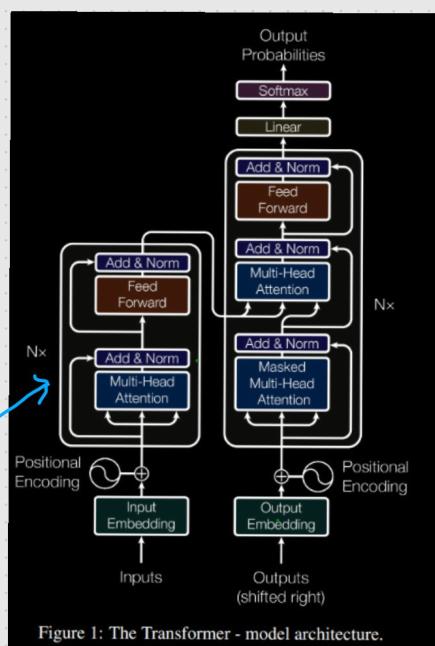


Figure 1: The Transformer - model architecture.

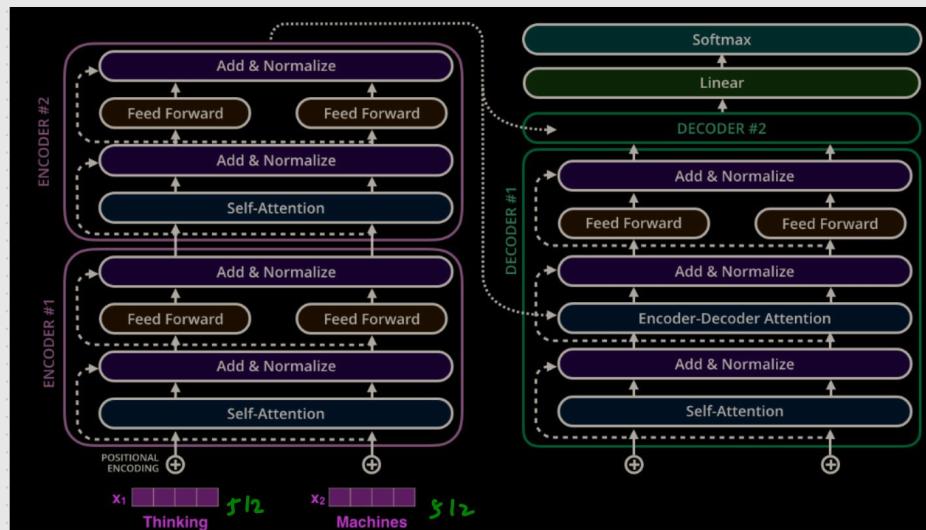
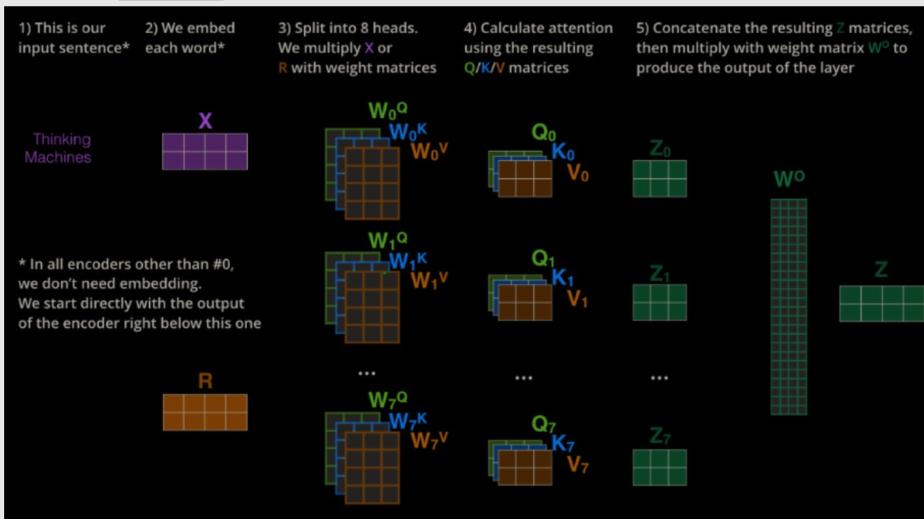
$$d_h = 64$$

$$K = 64$$

$$V = 64$$

$$\sqrt{64} = 8\sqrt{2}$$

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

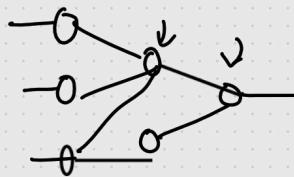
Residuals: Residual connection create a short paths for gradients to flow directly through the n/w. Gradient remains sufficiently large.

2) Improve Gradient flow

Convergence will be faster.

3) Enables Training of Deeper Networks.

① 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 these representation further

and allows the model to learn

Richer Representation.

AoNN ⇒

- ③ FFN → Depth ⇒ Adds Depth to the Model.

Depth ↑↑ ⇒ More learnings → DATA

④ 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 (Encode Decoder Attention)

③ Feed Forward Neural Network.

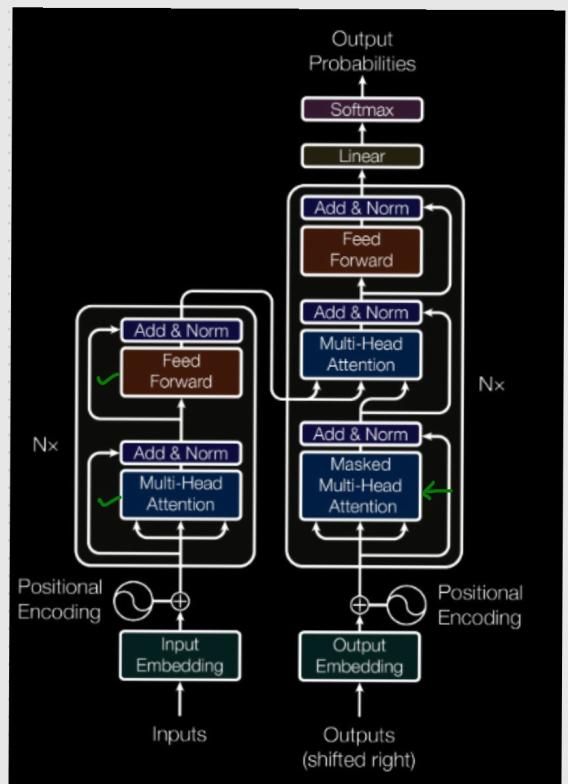
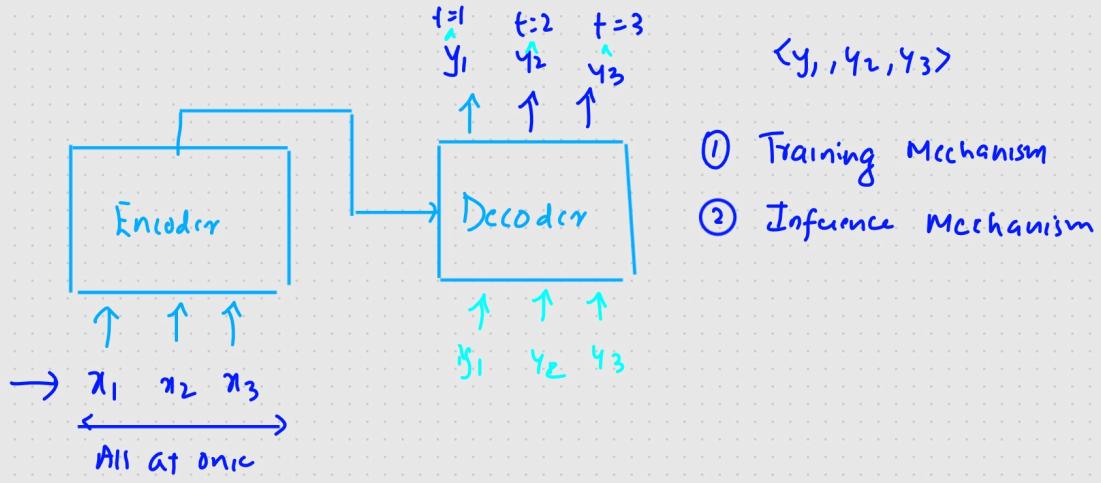
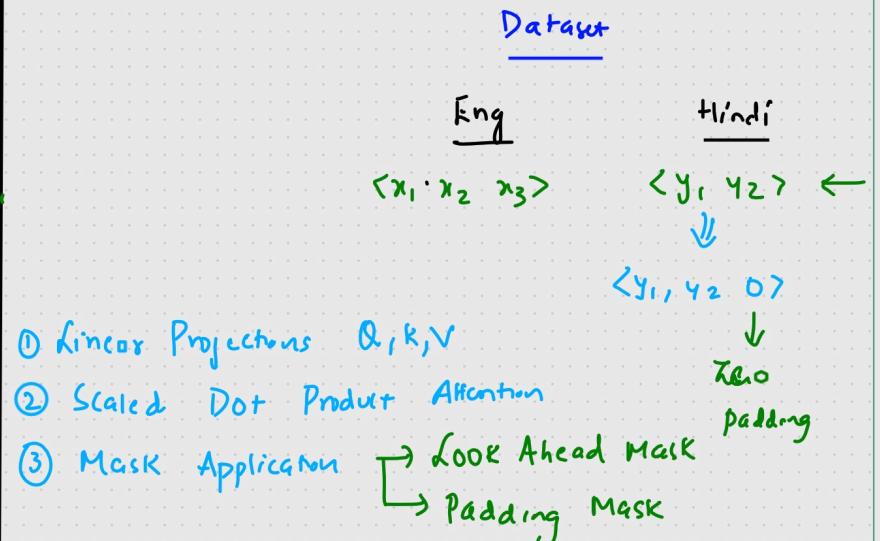
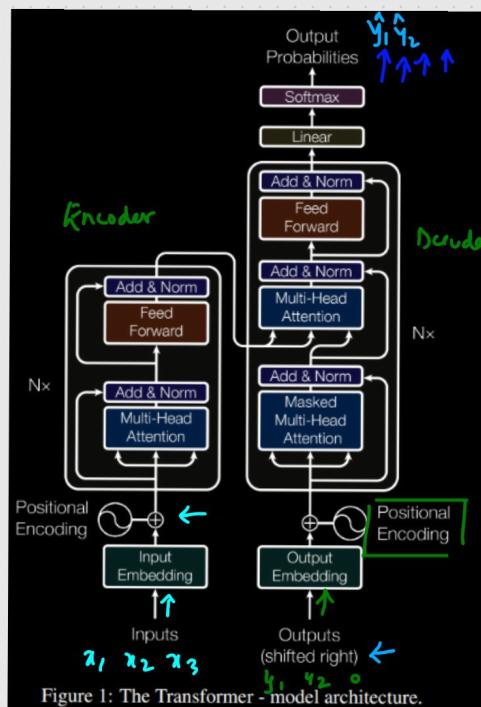


Figure 1: The Transformer - model architecture.



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



Masked

Multi Head
Attention.

$$[4 \ 5 \ \frac{I/P}{6} \ 7] \quad [O/P \ 1 \ 2 \ 3]$$

$$\begin{bmatrix} 1 & 2 & 3 & 0 \\ \downarrow \end{bmatrix}$$

4 dimension vector

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

$Q = \text{Output Embedding} + W_Q = \text{Output Embedding}$

$K = " + W_K = "$

$V = " + W_V = "$

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

③ Scaled Dot Product Attention Calculation

$$\text{Scores} = Q * K^T / \sqrt{dk}$$

$$= Q * K^T / 2.$$

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

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

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: $\text{inp} \leftarrow \text{Sequence 1 } [1, 2, 3] \rightarrow y_1, y_2, 0.0000$

$\text{Olp} \leftarrow \text{Sequence 2 } [4, 5, \boxed{0}]$ 0 is the padding token

$\uparrow \rightarrow$ Influence the Attention Mechanism

$\rightarrow 100.$

\Downarrow

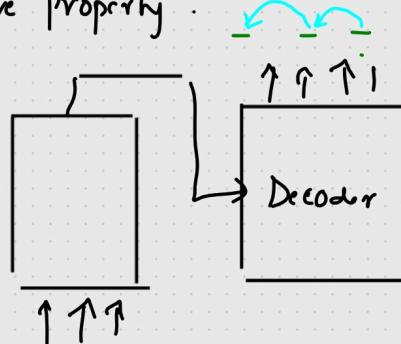
lead to Incorrect or biased predictions.

A padding mask \Rightarrow The tokens are ignored.

MASKING \rightarrow Padding Mask $\left\{ \begin{array}{l} \text{Padding Mask } [1, 1, 1] \\ \text{Look Ahead Mask } [1, 1, 0] \end{array} \right.$

② Look Ahead Mask → Maintain Auto Regressive Property

- ① To ensure that each position in the decoder output sequence can only attend to the previous position, but no future position



How Anyo

- ② Sequence → Language Modelling, Translation

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

Attention ← to token 1, 2 Mechanism 1, 2

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

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

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

④ 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

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

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

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

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.

* ①

Softmax

$$\text{Softmax Score} = \text{softmax}(\text{masked Scores})$$

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

* ② Weight Sum of Value

$$\text{Attention O/p} = \text{Softmax Scores} * V.$$

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

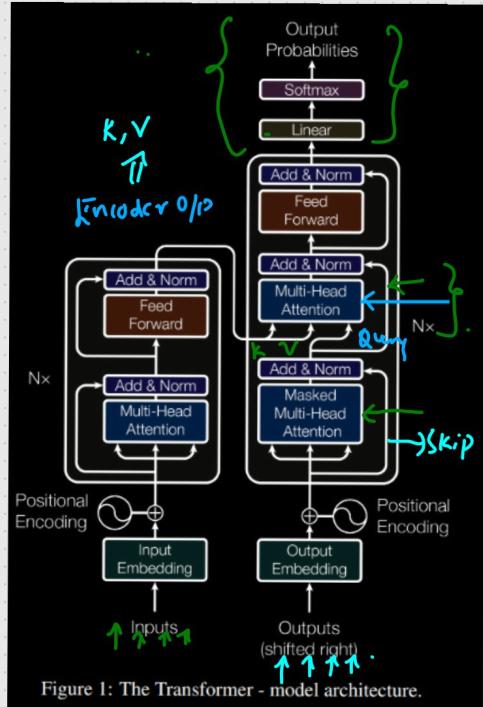


Figure 1: The Transformer - model architecture.

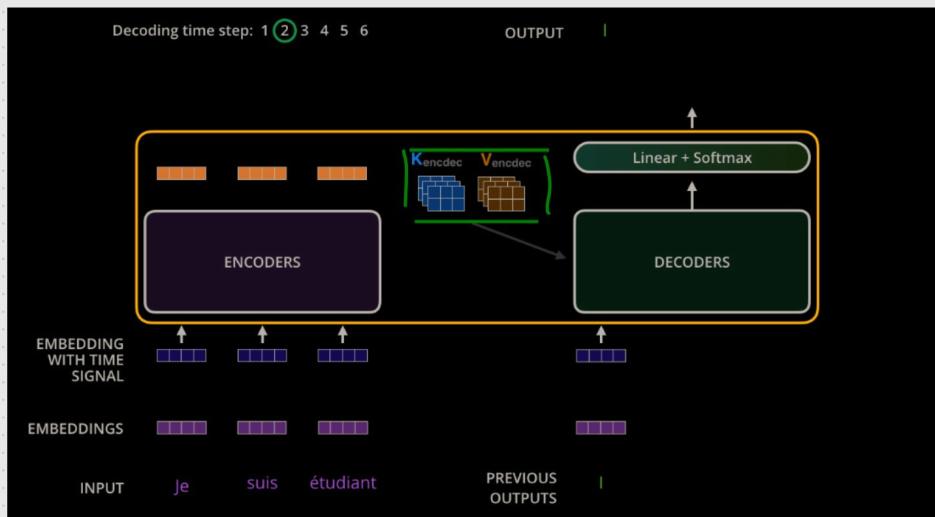
- ① Encoder O/p \rightarrow Set of Attention vector K & V
- ② Masked Multihead \rightarrow Attention Vector Q {Query Vector}

These are to be used by each decoder in its "Encoder-decoder" Attention layer



Helps the Decoder to focus on appropriate places in the i/p Sequence

Datactf
I/P O/P
 $\langle x_1, x_2, x_3 \rangle$ $\langle y_1, y_2, y_3 \rangle$.



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

{ BLOG } ⇒ Transformers

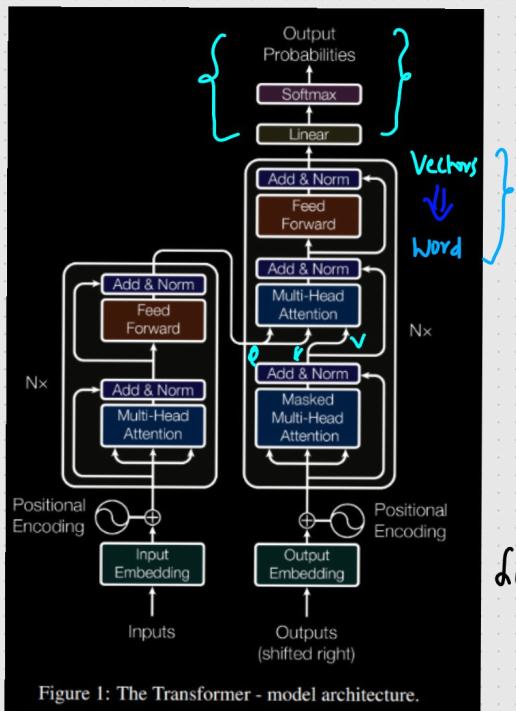
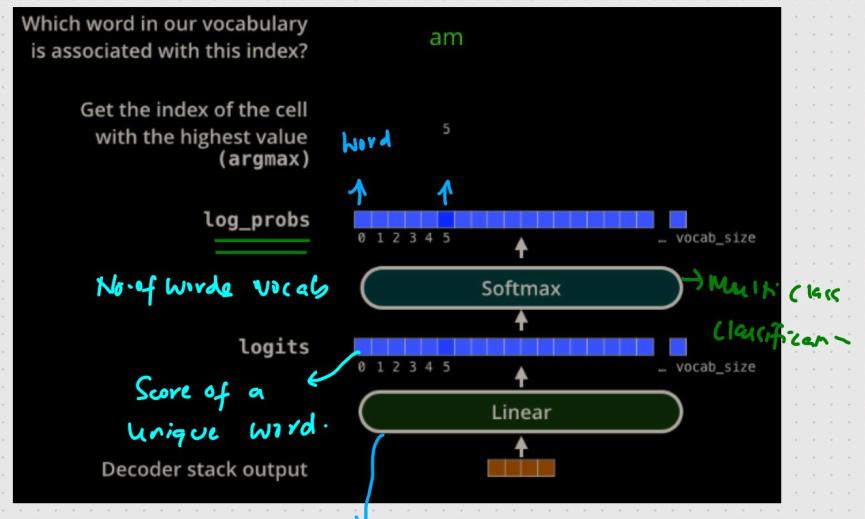


Figure 1: The Transformer - model architecture.

Which word in our vocabulary is associated with this index?

Get the index of the cell with the highest value (argmax)



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

② 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.

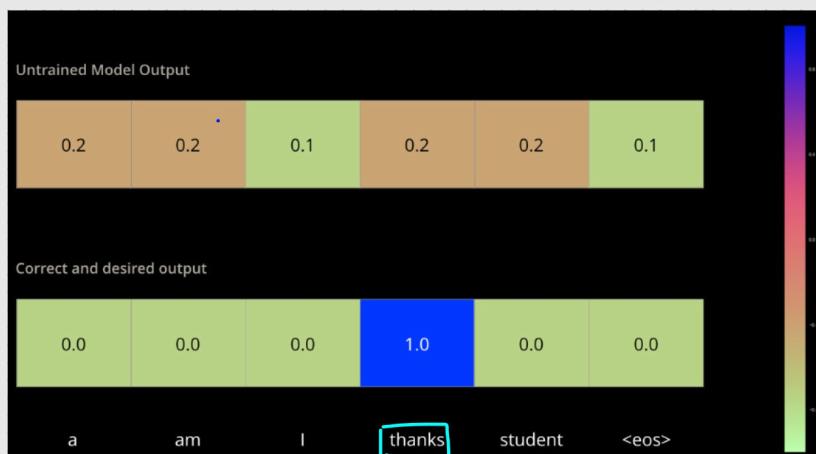
Recap of Training

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

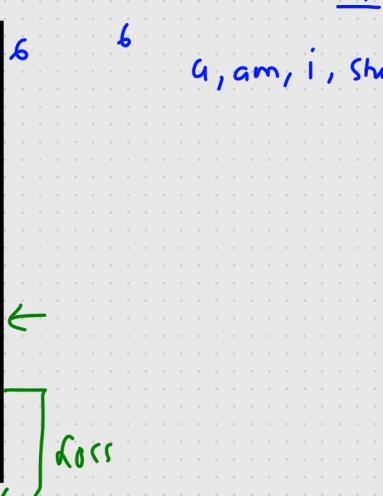
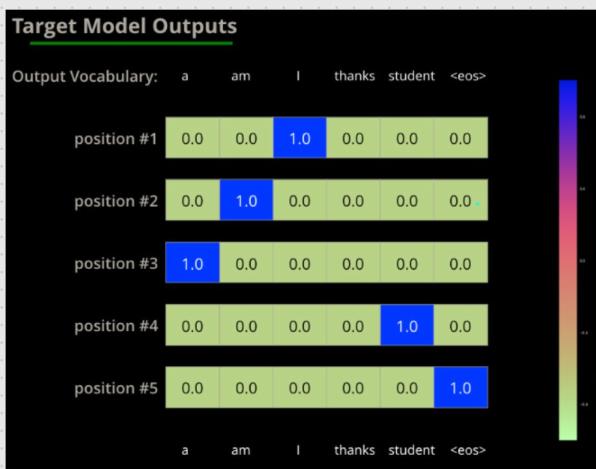
Merci → Thanks

I am a student <eos>.



⇒ Loss function ↓.

Back Propagation



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

