

- How to treat the text data in columns ?

1. Bag of Words (Bow) :

	Reviews	Sentiments	Document
0			
1			
.			
.			
1000			

Combinations of Documents are called Corpus.

- There are simple text cleaning techniques that can be used as a preprocessing step, such as:
 - Ignoring case
 - Ignoring punctuation
 - Ignoring frequent words that don't contain much information, called stop words, like "a," "of," etc.
 - Fixing misspelled words.
 - Reducing words to their stem (e.g. "play" from "playing") using stemming algorithms.

- Second, we try to learn Vocabulary. Vocabulary contains unique words of corpus.

For Example,

It was the best of times,

it was the worst of times,

it was the age of wisdom,

it was the age of foolishness,

['best', 'times', 'worst', 'age', 'wisdom', 'foolishness']

- We have to make Documents Vector from Corpus.

best	times	worst	age	wisdom	foolishness
1	1	0	0	0	0
0	1	1	0	0	0
0	0	0	1	1	0
0	0	0	1	0	1

Disadvantages of Bag of Words :

- It just takes frequency of words, it has no meaningful extraction like other advanced methods. It does not make semantic relationships with other tokens.

TF-IDF(Term Frequency - Inverse Document Frequency) :

Sentences

$D_1 = \text{Prime Minister India}$
 $D_2 = \text{Prime Minister Nepal}$
 $D_3 = \text{Corona Virus India}$

freq
 Prime $\rightarrow 2$
 Minister $\rightarrow 2$
 India $\rightarrow 2$
 Nepal $\rightarrow 1$
 Corona $\rightarrow 1$
 Virus $\rightarrow 1$

Term Freq \Rightarrow
$$\frac{\text{No of Reptetion of Word in a Sentence}}{\text{No of Word in a Sentence}}$$

Inverse Document

$$\text{Freq} = \log \left[\frac{\text{No of sentences}}{\text{No of sent Containing tokens}} \right]$$

TF	Word	D ₁			D ₂			D ₃			IDF
		0.33	0.33	0	0.33	0	0.33	0	0	0.33	
	Prime	0.33	0.33	0	0	0	0	0	0	0.176	0.176
	minister	0.33	0.33	0	0	0	0.33	0	0	0.176	0.176
	India	0.33	0	0	0	0	0	0	0	0.477	0.477
	Nepal	0	0	0	0.33	0	0	0.33	0	0.477	0.477
	Corona	0	0	0	0	0	0	0.33	0	0.477	0.477
	virus	0	0	0	0	0	0.33	0	0	0	0

$$\frac{\text{TF-IDF}}{D_1} = \frac{\text{Prime}}{0.058} = \frac{\text{Minister}}{0.058} = \frac{\text{India}}{0.058} = \frac{\text{Corona}}{0} = \frac{\text{Virus}}{0} = \frac{\text{Nepal}}{0}$$

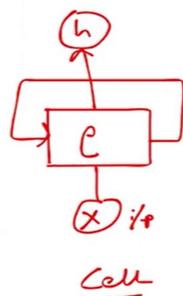
<u>TF-IDF</u>	<u>Prime</u>	<u>Minister</u>	<u>India</u>	<u>Corona</u>	<u>Virus</u>	<u>Nepal</u>
D1	0.058	0.058	0.058	0	0	0.15
D2	0.058	0.058	0	0	0	0.15
D3	0	0	0.058	0.15	0.15	0
	0.116	0.116	0.116	0.15	0.15	0.15

Those who have same values are interlinked with each other and higher values are most frequent.

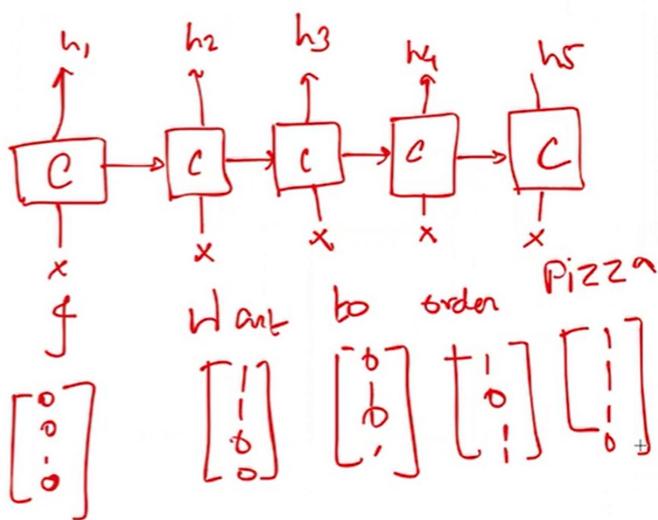
RNN(Recurrent Neural Network):

It is deep learning approach to solve text related problems. It provides semantic relationships between tokens that vectorizers not provide. It helps to achieve meaningful extraction from words.

Recurrent Neural Network



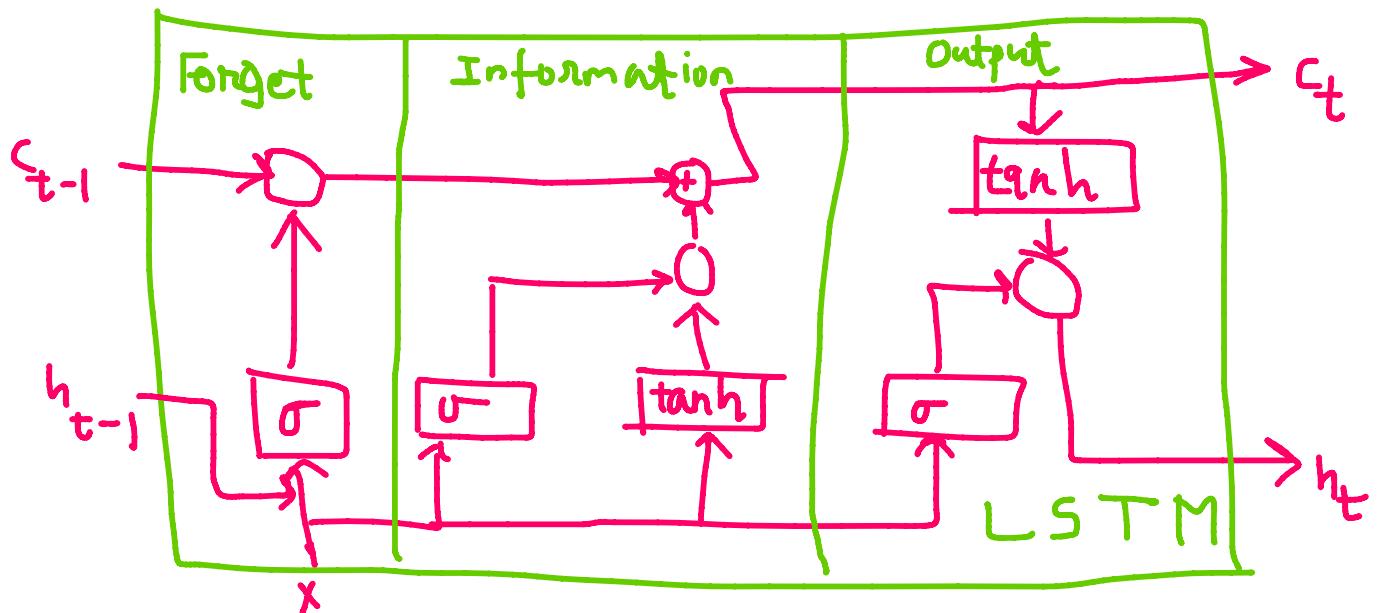
Cell



If we get very big sentence, it cannot handle such sentences. It starts getting vanishing gradients because it not remember large sentences.

LSTM(Long Short Term Memory):

It has ability to forgot unwanted tokens and stores only necessary tokens. It provides weightage to necessary tokens.



C_t :- Cell State or weights

h_{t-1} :- Hidden State or inputs from previous cell

X :- inputs

Forget layer :- used to forget unwanted tokens

Information layer :- stores important tokens

Output layer :- combines current input with previous input and provides o/p to next cell

$$\begin{aligned}
 f_t &= \sigma(W_f[h_{t-1}, x_t] + b_f) \quad \begin{cases} \rightarrow 0 \rightarrow \text{forget} \\ \downarrow 1 \rightarrow \text{remember} \end{cases} \\
 i_t &= \sigma(W_i[h_{t-1}, x_t] + b_i) \quad \begin{cases} \rightarrow \text{All the values} \\ \rightarrow i_t \times c_t \rightarrow \text{to update} \end{cases} \\
 c_t &= \tanh(W_c[h_{t-1}, x_t] + b_c) \quad \begin{cases} \rightarrow \text{new vector to add} \\ \rightarrow \text{it state} \end{cases} \\
 l_t &= f_t \times l_{t-1} + i_t \times c_t \longrightarrow \text{OL} \\
 h_t &= \sigma(W_h[h_{t-1}, x_t] + b_h) \quad | \quad h_t^{\text{Output}} = h_t \times \tanh(c_t)
 \end{aligned}$$

Disadvantages of LSTM:

LSTM is good when we need only sequence of sentence but it fails in case we want to go reverse. We can't able to relate a token with its previous token.

We cannot change vector of previous token based on current token.

For Example,

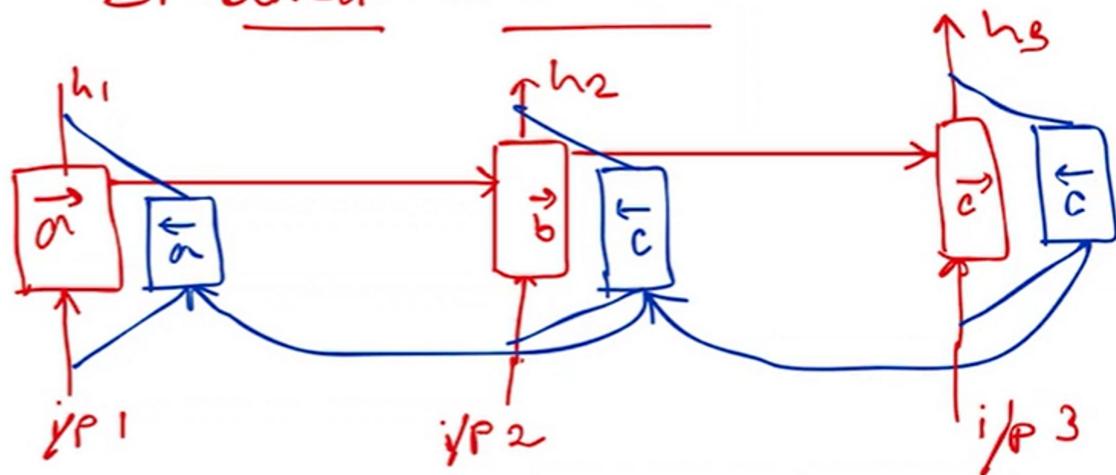
Tite Kubo's earliest influence is from [Shigeru Mizuki's manga Gegege no Kitaro](#).

He remembers trying to sketch its characters and found his own designs to be simpler than that of Mizuki's.

In above case, LSTM cannot relate 'He' with 'Tite Kubo' here. It only moves in forward direction.

Bidirectional LSTM :

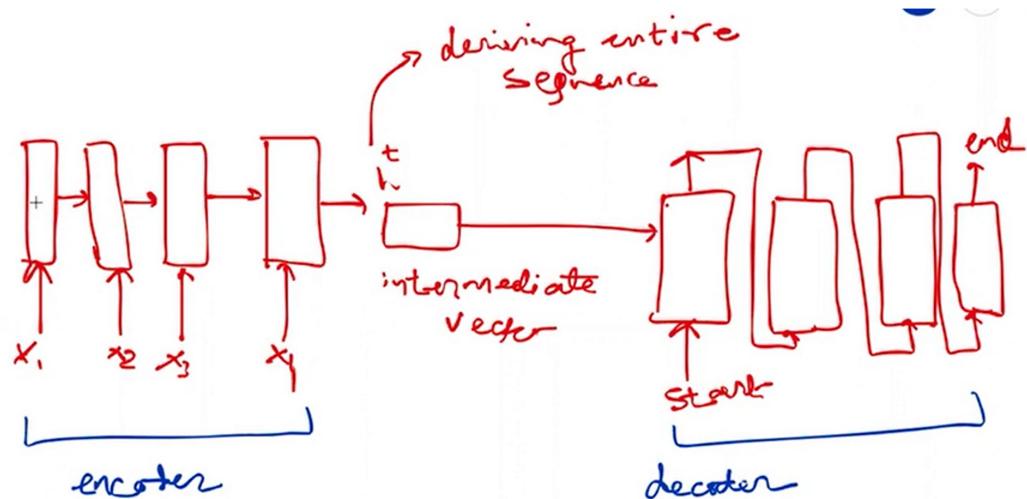
bi-directional LSTM



Disadvantages of Bidirectional LSTM:

If we have very long sentences or if we have multiple references in text. It couldn't able relate it with text. For Example, if we have multiple person name in text it not able to relate it with text.

Encoder & Decoder :

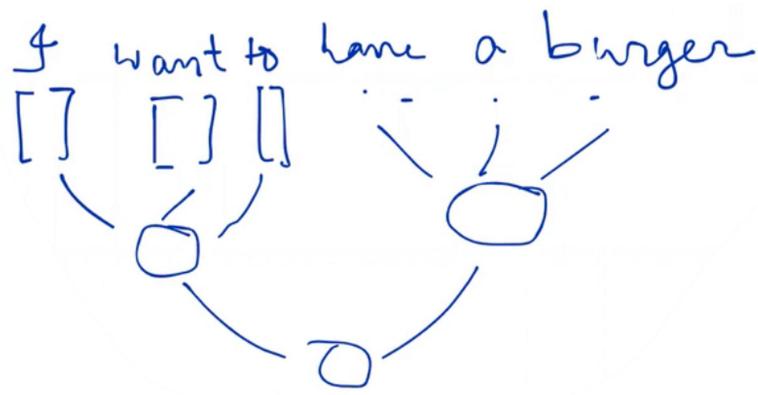


Encoder can be RNN or LSTM used to take input, it not generate any intermediate hidden state output and directly generate final output which is intermediate vector, is transferred to decoder.

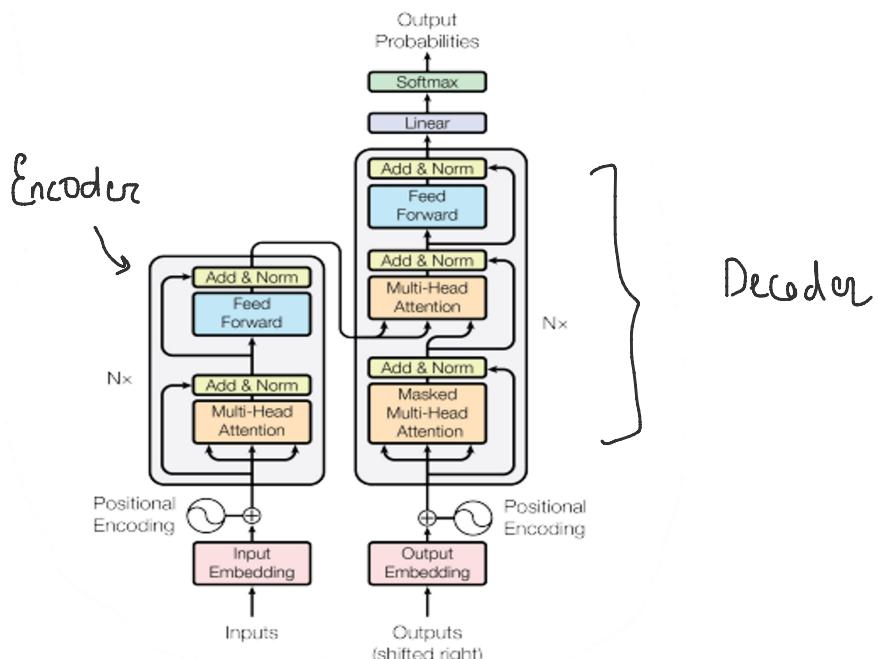
Decoder has multiple LSTM cells which are connected which helps us to generate final output.

The Application of Encoder and Decoder is Neural Machine Translation which is translation of language to other.

Attention is all you need:



Like this we have multiple sentences where we assign each token a attention vector, Multiple tokens forms another vector and multiple sentences forms another vector. In this way, attention vector gives some weightage to each token and more focus on important vectors.



Positional Encoding : includes (i) Sentence embedding : Finds no of sentences
(ii) Positional Embedding : finds semanticness between the tokens
(iii) Token Embedding : token embedding

Multi-Head Attention : It provides some attention vector to each tokens, multiple tokens having different attention. That's why it is called Multi Head Attention.

Masked-Multi Head Attention : It mask or omit the important token/ word to decoder so it can learn it. Some time it omit few tokens next time it omit next few token.

i) query

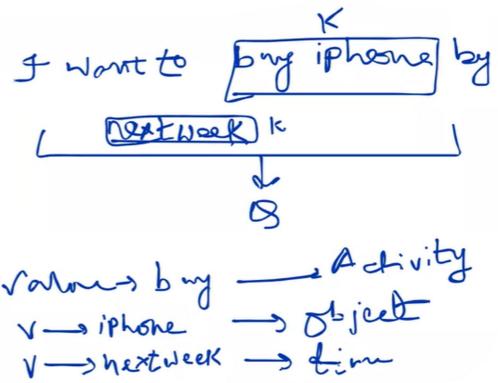
ii) key

iii) value

$$\alpha = \sum \alpha_j h_j$$

↓
Weight

↑ context
element which you want to



a - attention vector

$$e_{ij} = f(s; h_i)$$

↓
decode encoder
Seg Seg

$f(s_i) \cdot g(h_i)$

$$\alpha_{ij} = \frac{\exp(e_{ij})}{\sum \exp(e_{ij})}$$

Softmax = $\frac{QK}{\sqrt{d}} \cdot a$ Q - no of query , K - no of keys, d - no of tokens
and a - attention vector