

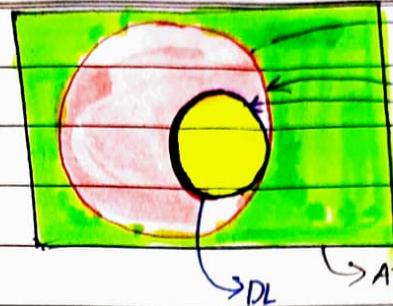
Date: _____

Topic: NATURAL LANGUAGE PROCESSING

Krish Naik

Agenda:

- Roadmap of Natural Language Processing
- Why NLP?
- Lots of examples
- Tokenization, stemming, Lemmatization
- Bag of Words.



ML

NLP

DL

Pre-requisites

- Python
- Stats
- ML Algo
- ANN, Optimizers, Loss func

NLP can be used both in ML and DL because we deals with text dataset in NLP.

Our aim w.r.t NLP should be that we need to create a model in such a way that the model should be able to understand text and based on this it should give us some output.

Make the machines to do tasks efficiently

[BERT]

[TRANSFORMER]

Roadmap

Bi-directional, LSTM, Encoder

Decoders, Attention models

Advance Text Preprocessing - Word embedding | Word2Vec

RNN, LSTM, RNN, GRU, RNN → DL

ML Use-cases (Solve)

Tent Preprocessing 3 → Word2Vec, AvgWordVec

Tent Preprocessing 2 → BOW, TF-IDF, Unigrams, Bigrams

Tent Processing + BOW, TF-IDF, Word2Vec

→ Tokenization, Lemmatization, Stopwords, Stemming

→ Cleaning the data in this step.

(I) Tokenization:

Consider one ML Use-case.

Gmail → Spam Classifier Independent features: Email body, Email subject

Dataset

Email body	Email Subject	Op (Spam/Ham)
You won 1000\$	Billionaire	Spam
Hey, How are you?	Hello	Ham
Credit card worth	Winner	Spam

Here you cannot simply give these sorts of text dataset to machine for processing, you would want to perform some tasks.

Step-I: Tokenization

Tokenization: Process of converting sentence into words.

Step-II: Stemming

for example - You won 1000 \$

Stopwords

Lemmatization

Here [You] [won] [1000] [\$] {Converted the entire sentence into words separately}

Deepak ydu

Teacher's Signature _____



Date: _____

Topic: _____

(II) Stopwords

Example - 2: Hey buddy I want to go to your house.

Here in above sentence, "to" is repeated and removing words like these may not change the idea of sentence.

So in-order to remove these words. We use technique called "Stopwords".

"~~Keywords~~" "We can also have our own Keywords".

not **w** **Important word**
changes idea of any sentences.

(III) Stemming: Process of reducing words to their base word stem.

historical → base word stem may not have any particular meaning.
history history ↓
Extracting root word (or) Baseform.

finley
final
finalized

Disadvantage of Stemming

Advantages

→ Stemming is really fast.

going
goes
go
gone

→ Here we got meaningful base word.

How In-order to overcome disadvantages of stemming we use concept called "Lemmatisation".

(IV) Lemmatisation:

How we get meaningful word:

history → history
historical → finally
final → final
finalized → final

Lemmatisation has entire dictionary of words it will check w.r.t base word and then compare.

Advantages

→ Provides meaningful Word

Disadvantages

→ Slow performance.

Use-cases of Stemming, Lemmatisation

(i) Spam classification

(ii) Review classification

(iii) Text Summarization

(iv) Language Translation

(v) Chatbot

#1 Text-Preprocessing involves following steps:

- ① Tokenization
- ② Stopwords
- ③ Stemming &
- ④ Lemmatization

Here we have cleaned the entire words in the form of various methods.

#2 Convert Words to Vectors

- ① Bag of Words
- ② TF-IDF (Term Frequency, Inverse Document Frequency)
- ③ Word2Vec



Date: _____

Topic: _____

Bag of Words:

D1: He is a good boy

D2: She is a good girl

D3: Boys and girls are good.

enclosed red words are generic words and should be removed.

These words doesn't play big role they can be removed.

stopwords

D1: good boy

D2: good girl

D3: Boys girls good boy girl

Vocabulary

good

boy

girl

Vocabulary

Frequency

f₁f₂f₃Doc 1 :

3 → good

boy girl

stop :

2 Doc 1 : 1 1 0

stop

2 Doc 2 : 1 0 1

classification ML algorithm

Doc 3 : 1 1 1

Euclidean distance

cosine similarities

BOW → Binary BOW

Advantages

1) Simple and Intuitive

Disadvantages

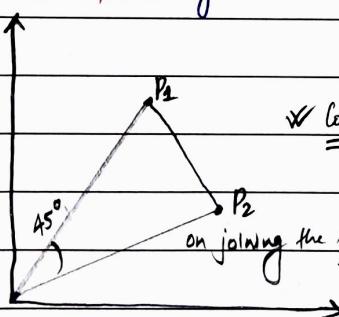
1) Sparsity

2) OOV (out of Vocabulary)

3) Ordering of word has completely changed.

4) Semantic meaning not able to capture

Cosine Similarity

Here if we need to find distance between P₁ and P₂ we can do it by two ways.

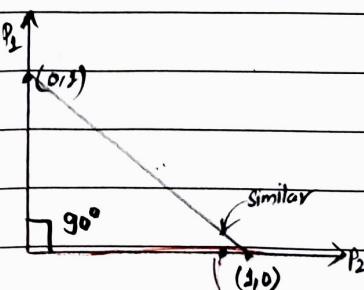
✓ (I) Euclidean distance ✓ (II) Cosine Similarity

on joining the point to base lets suppose angle is 45°.

In-order to find similarity between P₁ and P₂, we use,We have, $\cos 45^\circ = 0.707$

$$\frac{P_1 \cdot P_2}{\|P_1\| \|P_2\|} = 0.707 \Rightarrow \text{Cosine Similarity}$$

$$= 0.707 [\because P_1 \text{ & } P_2 \text{ are only } 0.707 \text{ similar}]$$

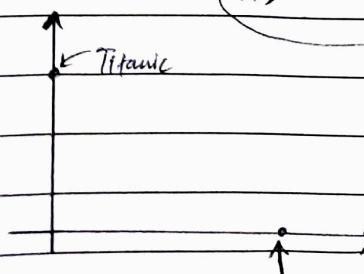


$$\cos 90^\circ = 0$$

$$1 - 0 = 1 [\text{Cosine Similarity}]$$

$$\cos 0^\circ = 1$$

$$1 - 1 = 0 \checkmark$$



In movie recommendation system we can relate that one movie

Avenger which may be similar to Ironman must be similar with the cosine as well.

Deepak ydu

Avenger Ironman

Teacher's Signature _____

So, disadvantage of BOW is that it's weak

for capturing Semantics meaning. Hence n-grams comes into picture.



Date: Day-2

Topic: _____

Agenda:

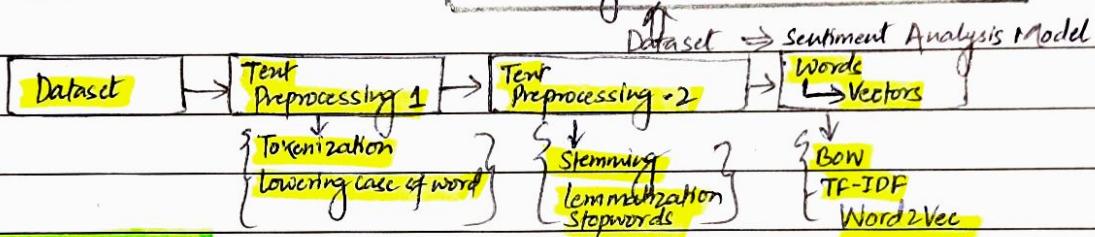
- (1) Text preprocessing → Words → Vectors
 (a) One-hot encoding
 (b) Bag of words, → n-grams
 (c) TF-IDF
 (d) Word2Vec
- (2) Term frequency Inverse Document frequency

Basic Terminology used in NLP:

- (1) Corpus → Paragraph → [D1, D2, D3, D4]
- (2) Documents → Sentence → 1st document D1
2nd " D2
etc
- (3) Vocabulary → 10K Unique words (vocabulary)
- (4) Words → Word

Sentiment Analysis:	
Tent	O/P
D1: The food is good	1
D2: The food is bad	0
D3: Pizza is amazing	1
D4: Burger is bad	0

Dictionary Book
10K Unique words

One-Hot Encoding:

Suppose a Corpus as below:

Document 1-D1 → [A man eat food]
 D2 → [Cat eat food]
 D3 → [People watch KRISH YT.]

Vocabulary:

9

"A man eat food cat people watch KRISH YT." ↗ These are unique vocabularies

D1 → [1 0 0 0 0 0 0 0 0]
 D2 → [0 1 0 0 0 0 0 0 0]
 D3 → [0 0 1 0 0 0 0 0 0]
 D4 → [0 0 0 1 0 0 0 0 0]
 D5 → [0 0 0 0 1 0 0 0 0]

In place of word replaced by 1 remaining all zeros.
 "ONE-HOT Encoding"

Advantages

↳ Simple to implement

↳ Sparse Matrix

↳ OOV (Out of Vocabulary)

↳ Not fixed size

↳ Semantic meaning b/w word is not captured.

D1 → [1 0 0 0] D2 → [1 0 0] D3 → [1 0 0] D4 → [1 0 0]
 [0 1 0 0] [0 1 0] [0 0 1 0] [0 0 1]
 [0 0 1 0] [0 0 1] [0 0 0 1] [0 0 0 1]

Teacher's Signature _____

Deepakydu



Date: _____

Topic: _____

Ngrams :- \Rightarrow (Bigrams, Trigrams, ...)

Comparing with Doc sentences from BOW
e.g. In the doc there is not word good girl so "0" but there is good boy (so "1").

	f_1	f_2	f_3	f_4	f_5
Sent-1	good	boy	girl	good boy	good girl
Sent-2	1	1	0	1	1
Sent-3	1	1	1	0	0

Example:

DOG CHASE CAT
 ↓ ↑ ↓
 1 1 1

= total \Rightarrow 2 bigrams only
 DOG CHASE CHASE CAT
 ✓ ✓

Ex-2 : KRISH IS NOT FEELING WELL (1,3) this means from unit 1 to unit 3 any group can be used

f_1 f_2 f_3 f_4 f_5

I am not feeling well
 ↑ ↑ ↑ ↑ ↑
 1 1 1 1 1

I am not = 1
 am not feeling = 1 } 3 trigrams
 not feeling well = 1

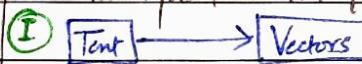
Day-3 NLP:

Agenda

- (1) Bow (Bag of words)
- (2) TF-IDF
- (3) Practical Implementation
- (4) Quiz

Bag of Words (BOW)

Technique to convert text into vectors.



Ex: Sent-1 : He is a good boy

Sent-2 : She is a good girl

Sent-3 : Boy and girl are good. is. she, is, a, and, are.

Step-1st : Apply stopwords,

Stopwords will help you remove all the unnecessary words such as he, is, a, and, are.

\rightarrow Lowercasing.

Step-2nd \leftarrow After lowercasing we obtain result shown alongside.

Step-3rd - Find out frequency of vocabulary w.r.t unique words.

II

Sent 1 good boy
 Sent 2 good girl
 Sent 3 boy girl good

III

Frequency (Vocabulary)

unique words

good 3

boy 2

girl 2

IV Sentences

(good) (boy) (girl)

f_1 f_2 f_3

Sent 1: 1 1 0

Sent 2: 1 0 1

Sent 3: 1 1 1

Step-4th : Divide sentences into features and pvt for present word in

features and pvt for present word in sentence 1 and not present in 0.

which is Sparsity Matrix.

Disadvantages

- \rightarrow If there are many words your Sparsity Matrix will be very large and its difficult in computation.
- \rightarrow OOV (Out of Vocabulary).
- \rightarrow Semantic Meaning is missing.

lets consider two sentences:

I: The food is good

II: The food is not good

food good is not The
 ↓ ↓ ↓ ↓
 1 1 1 1 1

More sentences are not similar

but if we find the cosine similarity for the sparse matrix we will come to know they are similar but in fact they are not. This is disadvantage.

Teacher's Signature

This can be fixed with TF-IDF. ✓

Deepak yd



Date: _____

Topic: _____

Term Frequency - Inverse Document Frequency (TF-IDF)

lets take both the sentences used in BOW.
Sentences after cleaning the text.

- Sent 1: good boy
- Sent 2: good girl
- Sent 3: boy girl good.

Sent 2Sent 1: ~~cat~~ eats foodSent 2: ~~bat~~ eats foodSent 3: ~~Krish~~ eats food

The issue with BOW is that it had given same weightage to all the words irrespective of meanings. For boy and girl same weightage, good also same weightage. But its not fair. Boy is repeated in sentence 1 and sent. 3 and girl in sent 2 and sent 3. So weights should be managed accordingly.

Which ever words are rarely present we should allocate more weightage.

For instance in ex-2, cat, Bat and Krish is unique (row), words. So they must be assigned with more weightage.

$$\text{Term Frequency} = \frac{\text{No. of rep. of words in sentence}}{\text{No. of words in sentence}}$$

Rare words will be captured by "Term frequency"

Common Words will be captured by "Inverse Document Frequency"

→ Sentences (TF is calculated per sentence)

On combining both we would be able to identify most rare word

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

→ Words (IDF is calculated per word)

Term Frequency :

	Sent 1	Sent 2	Sent 3
good	1/2	1/2	1/3
boy	1/2	0	1/3
girl	0	1/2	1/3

*
(multiply)

Inverse Document Frequency

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

Since good is present in all the sentences it is saying that it is playing not crucial role

boy is present rarely in 2 sentences hence boy has some

Value similarly for girls.

f_1 f_2 f_3
good boy girl

Sent 1 $1/2 * 0 = 0$ $1/2 * \log_e(3/2) = 0$

Sent 2 0 0 $1/3 * \log_e(3/2) =$

Sent 3 0 $1/3 * \log_e(3/2) = 1/3 \log_e(3/2)$

Some kind of Semantic Information

is being displayed using TF-IDF method.

Advantages

→ Intuitive

→ Word importance is getting captured.

* max_features=3 → This is going to take top 3 most occurring features.

Disadvantages

→ Sparsity

→ OOV (Out of vocabulary)

Teacher's Signature

Deepak ydu

In NLP Word embedding is a term used for the representation of words for text analysis, typically in the form of a vector-valued vectors that encode the meaning of word such that words that are closer in the vector space are expected to be similar in meaning.

Wikipedia

Date: Day 4

Topic: _____

Agenda

(1) Word Embeddings

(2) Word2Vec

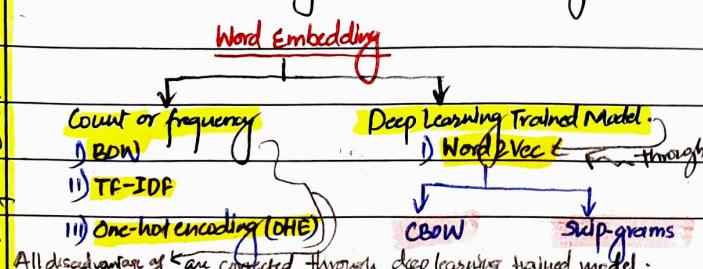
CBOW (continuous Bag of word)
Skipgram

In case of deep learning we basically use embedding layer. This is basically used for converting text to vectors.

and Bow

Word Embeddings :-

: It is a technique which converts word in vectors. Word embedding can be done on two ways:



The major disadvantages with TF-IDF was that semantic meaning is not captured.

Semantic meaning \rightarrow word depicting
Similar meaning. Honest
Good

Word2Vec :- Feature Representation

EX1: $f_1 f_2 f_3 f_4 f_5 \dots f_{60} \dots$
Boy Girl King Queen Apple Mango

In Word2Vec for every word we try to create vector. but always remember this vector will be of limited dimensions. Also Sparsity would be reduced.

Another disadvantage was that TF-IDF Sparse matrix, because of which huge dimensions were created. Many more

Vectors were created which is computationally expensive in terms of processing, model training etc also may lead to overfitting.

These all the disadvantages will be solved by Word2Vec.

This Word2Vec also make sure semantic meaning is maintained.

Suppose similar words honest and good is converted into vectors with the help of Word2Vec. We will get vectors which are very near to each other. If you try to subtract those vectors, you will get minimal difference.

Here we are assigning some vector to our feature representation for eg, Boy -1 and girl by +1 similar in case of King & Queen. But apple and mango are somehow related with one another in case of gender so it getting some random values.

Let's suppose King, Queen, Man, Women is created by following vectors.

So, from Ex1:

	f_1	f_2	f_3	f_4	$f_5 \dots$	$f_{60} \dots$
Boy	-1	1	-0.93	$+0.93$	0	0.93
Girl	0.03	0.02	0.95	0.96	-0.02	0.01
KING	0.03	0.02	0.7	0.6	0.95	0.96
QUEEN	$=$	$=$	$=$	$=$	$=$	$=$
APPLE	$=$	$=$	$=$	$=$	$=$	$=$
MANGO	$=$	$=$	$=$	$=$	$=$	$=$
Gender	-1	1	-0.93	$+0.93$	0	0.93
Royal	0.03	0.02	0.95	0.96	-0.02	0.01
Age	0.03	0.02	0.7	0.6	0.95	0.96
Food	$=$	$=$	$=$	$=$	$=$	$=$

Suppose these are our feature representation

feature representation

300 dimensions

All words are represented by 300 dimensions.

King $[0.96 \ 0.95]$ Man $[0.95 \ 0.98]$

Queen $[-0.96 \ 0.95]$ Woman $[-0.94 \ -0.95]$

King - Man + Woman = Queen.

So what ever vector you get from result that would be matching obtained feature value.

This is done by "Cosine Similarity".

X Euclidean distance
X Manhattan distance

$\Rightarrow \theta = 45^\circ$ or
 $\Rightarrow \theta = 135^\circ \Rightarrow$ "Cosine Similarity".

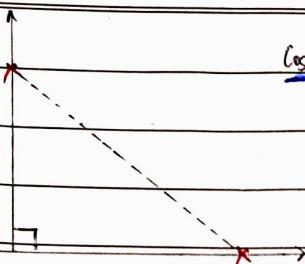
Distance = $1 - \text{Cosine Similarity}$
The more value towards 0, more similar points are. $= 1 - \cos \theta = 1 - \cos 45^\circ = 0.29$

These Values would be generated by model any pre-trained models etc.

Teacher's Signature

Date: _____

Topic: _____



$$\text{Cosine Similarity} = 1 - \text{cosine} = 1 - \cos\theta = 1 - \cos 90^\circ$$

$$= 1 - 0$$

$= 1$ (This means the points are nowhere similar to one another.)

Value towards 0 \rightarrow More Similar

Value towards 1 \rightarrow Not Similar

Suppose one's Movie Recommendation
Movie name - Vector / Vocabulary

Attribute

Action
Comedy
↓
features

Interest
Age group

Word2Vec (CBOW): (Continuous Bag of Words)

Let's suppose we have corpus.

CORPUS: Dataset Training

KRISH CHANNEL IS RELATED TO DATA SCIENCE

Since this word would bring some similarity in sentence. So it's content word

Window size = 5 \Rightarrow This window size is helpful for creating training data (i/p and o/p data)

for Word2Vec, we must have Window Size.

This Window size would be helpful for creating "Training data"

W.r.t training data, we have 2 features.

(I) Independent feature

KRISH, CHANNEL, RELATED, TO

CHANNEL, IS, TO, DATA

IS, RELATED, DATA, SCIENCE

I/P

IS

RELATED

TO

Center word

The main aim of Word2Vec is that,

will take the text and based on

the I/P as shown in example it will try to generate Vectors.

Simple BOW Representation or OHE

Word KRISH 1 0 0 0 0 0 0

Word CHANNEL 0 1 0 0 0 0 0

Word IS 0 0 1 0 0 0 0

Word RELATED 0 0 0 1 0 0 0

Word TO 0 0 0 0 1 0 0

Word DATA 0 0 0 0 0 1 0

Word SCIENCE 0 0 0 0 0 0 1

#Process In ANN (I) Assign weights

(II) Initialize the weights

(III) one hidden layer

(IV) Fw propagation

After fw propagation loss ($y - \hat{y}$) is calculated and again bw propagation would happen by updating weights

This fw and bw propagation keeps happening until we get less value very less.

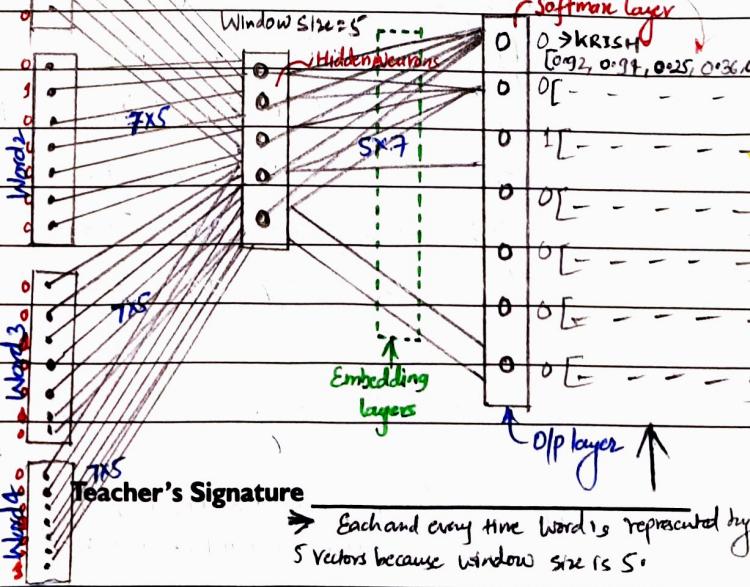
Deepak yd

This data independent features along with the I/P features and the data representation will be given to "fully Connected Layer"

fully connected layer

ANN

\rightarrow Window size is 5 mean I/P is 5 for every word



Teacher's Signature

\rightarrow Each and every time word is represented by 5 vectors because window size is 5.



Date: _____

Topic: _____

(II) Skipgram :- ? reverse of CBOW?

so Window size is define by hyper parameter tuning.

IP

OIP

IS

KRISH, CHANNEL, RELATED, TO

RELATED

CHANNEL, IS, TO, DATA

TO

IS, RELATED, DATA SCIENCE

Neural network will be reverse of CBOW

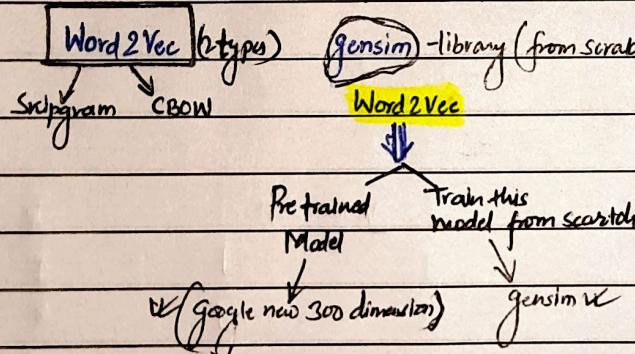
Day-5:

Agenda

① Practical implementation (Spam Classification Dataset)

→ Text preprocessing : Tokenization, Stopwords, stemming, lemmatization, NLTK

→ Text → Vectors : BOW, TF-IDF, Word2Vec, AvgWord2Vec



for example in case of Word2Vec when we are converting

word to Vector. Being single word in case of 300 dimension forms with 300 vectors.

If you have any new dataset and if it captures 75% of entire word by Pre-trained model then definitely go for Pre-trained model.

But if you see inside this pre-trained model if any such term are used wherein those terms are not present in pre-trained model at that time you can train your model from scratch.

King Queen

This problem can be actually fixed with new concept called as "AvgWord2Vec".

lets suppose we have one sentence.

"Welcome to Zambia" → We need to convert this into 300 dimensions but acc to Word2Vec individual conversion happens like Welcome gets converted to 300 dimension likewise "to" & "Zambia".

Our I/P is in this format

I/P

O/P

but our main aim is to convert them in totality.

I want to eat pizza
Convert into 300 dimensions

Spam/Ham Hence "AvgWord2Vec" is used for this.

This should be our aim.

But on using pre-trained model from scratch by using Word2Vec for each word we are getting that many no. of dimensions.

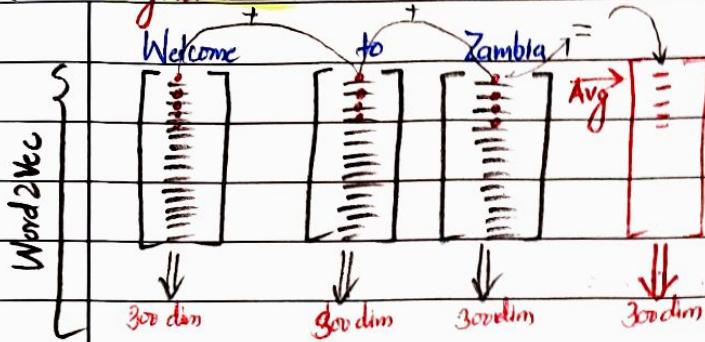
Teacher's Signature

Deepak ydu

Date: _____

Topic: _____

AvgWord2Vec



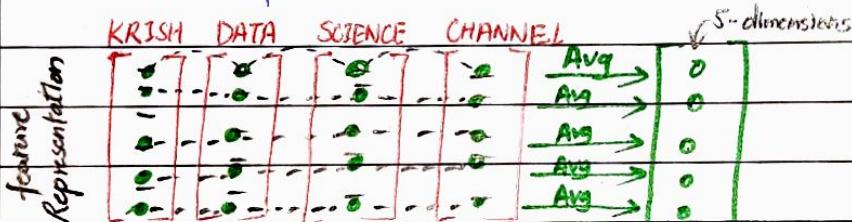
Here in case of Word2Vec. All the individual Word2Vec elements of dimensions are taken average and new Vector is created with single 300 dim

Once this is done our average 300 Dp feature and Spam/ham is o/p features



Ex 2: KRISH DATA SCIENCE CHANNEL

Word2Vec (feature representation) \rightarrow 5 dimensions

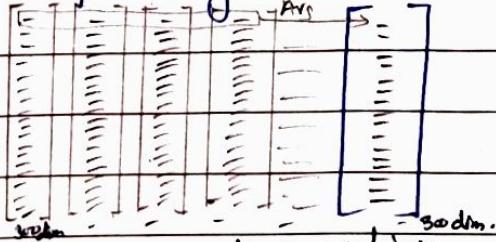


Example 2

Text	O/P	Vectors (Avg)	O/P
D1 : The food is good	1	[- - - -]	1
D2 : The food is bad	0	[- - - -]	0
D3 : Pizza is amazing	1	[- - - -]	1

Google pre-trained Word2Vec Model

[The food is good] \rightarrow D1



Here every single word is converted into vectors and taken average of it.

Our main aim should be w.r.t avgword2vec is that I can take this as an input (avgword2vec) and our provided O/P which is given above and we can basically train the model.

\rightarrow Since we are combining all the words the Semantic meaning would also be maintained. This is another advantage.

Teacher's Signature _____

Deepak ydu



Date: Day 6

Topic: NLP with DL

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

Recurrent Neural Network [RNN]

Word Embedding $\leftarrow \leftarrow$ Text \rightarrow Vectors
Word Embedding \leftarrow Word2Vec, AvgWord2Vec

Let's take an example as "chatbot"

In chatbot we are asking specific question and getting specific answers to it.

Question and Answer

Chatbot should always understand proper text of the question and then answer it.

The sequence of words are very important in case of Chatbot.

Let us take an example of "Google Translate"

In case of google translate also sequence of

word is very important on top of that grammatical corrections are equally important.

Language translation is very important thing.

Language translation \rightarrow Hindi \rightarrow English

Text Generation \rightarrow A Sentence \rightarrow Suggestions

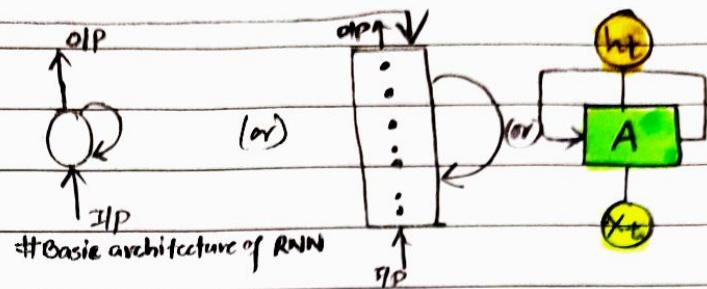
\downarrow Gmail, LinkedIn \rightarrow Auto-Suggestion

Grammatically Correct and Sentence Completion, are most imp. factors.

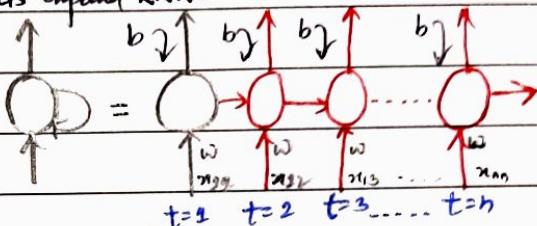
All these applications can be only trained well with the help of neural networks present in Deep Learning.

Word2Vec is a technique in deep learning which we use to train a model or we can use pre-trained model

which actually converts word into Vectors. This is important because "Semantic meaning is captured".



Let's expand RNN.



So, what ever o/p we are getting we are giving same as an input to same Neural Network in case of RNN.

Sentiment Analysis Example!

Sentences \rightarrow O/P
The food is good Positive

Let's understand how do we send this information to RNN.

So the sentence is basically indicated with $x_{1,1}, x_{1,2}, x_{1,3}, x_{1,4}$

$x_{1,4}$ notation - for example if "The" word need to be sent then its notation is assigned to the RNN as

Shown above. Some process happens inside RNN after that which is "weights are initialized".

but before that word is converted to vector and dimensions are initialized

So, when we are providing $x_{1,1}$ we will convert first $x_{1,1}$ by using Word2Vec.

$x_{1,1} \rightarrow$ Vector \rightarrow Then passed to RNN

$x_{1,1} \rightarrow$ Word2Vec \rightarrow Vectors of dimension (d)

(300)

This dimension is initialized to RNN.

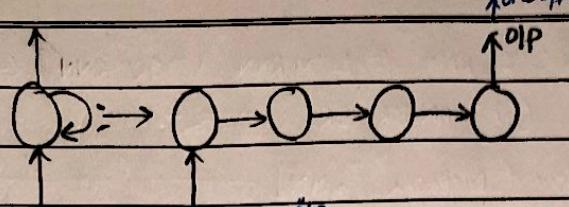


Date: _____

Topic: _____

Types of RNN :

(I) One to One RNN



Use Case (Example)
Image Classification

(II) One to Many RNN

(III) Many to One RNN

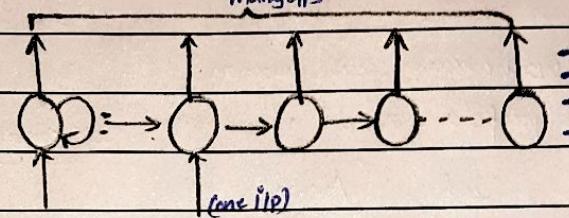
Here we are passing i/p but not retrieving it anywhere except one
So its "One to One" RNN

(IV) Many to Many RNN

There are 2 - important features
of deep learning.

(I) Forward Propagation

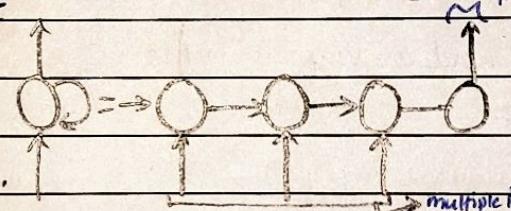
many o/p's



Examples
- Music Generation
- Text Generation
- Google Search Suggestions
- Movie Recommendations

(II) Backward Propagation

Any complex problem that are solved
by google or meta specifically w.r.t
NLP or text data where sequence is
involved, they have used RNN.

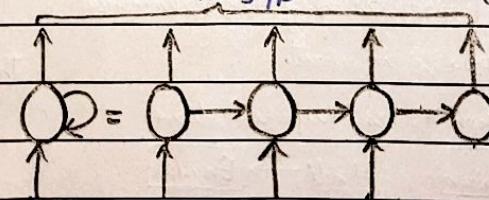


Examples:
- Sentiment Analysis
- Predict next day sale.

RNN is not only used for text data.

It is also used for "Time series Data".

Here we pass multiple i/p but retrieves only one o/p "Many to One" RNN

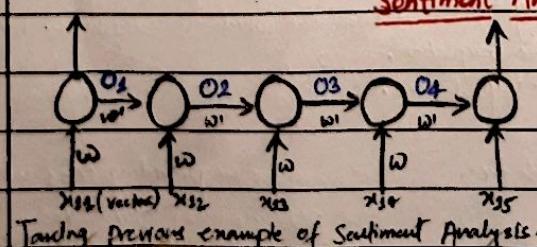


Examples
- Language translation
- Question Answering
- Chatbots etc..

Here we pass multiple i/p and retrieves multiple o/p 's, "Many to Many" RNN.

Forward Propagation (Many to One RNN)

Sentiment Analysis



In case of forward propagation we apply some specific operation. We multiply i/p s and weights

$$O_1 = f(x_{31} * w) + b$$

$$O_2 = f(x_{32} * w + O_1 * w) + b$$

$$O_3 = f(x_{33} * w + O_2 * w)$$

$$O_4 = f(x_{34} * w + O_3 * w)$$

In case of "Many to One" we apply "Softmax"

"Softmax" as its multi-class classification.

*: for binary classification we apply "Sigmoid".

Passing x_{31}, \dots, x_{35} to RNN in the form of vectors

We do that using Word2Vec Technique

Teacher's Signature

Deepak ydu



Date: Day - 7

Topic: _____

Backward Propagation:

Once we get \hat{y} after applying activation function.

We would try to calculate loss function with

\hat{y} and our original value y .

$$\text{Loss fn} = (y - \hat{y}) \text{ (Suppose)}$$

W.r.t we update all the weights (w) shown in previous fig.

$$\frac{\partial L}{\partial w} \rightarrow w' \text{ (first updation)} \quad \begin{matrix} \text{derivative of loss} \\ \text{derivative of weights} \end{matrix}$$

Weight update formula :

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

η = learning rate

$$\frac{\partial L}{\partial w} = \frac{\text{derivative of loss}}{\text{derivative of } w'}$$

$\frac{\partial L}{\partial w'}$ can be found using simple chain rule.

$$\frac{\partial L}{\partial w'} = \frac{\partial L}{\partial g} \times \frac{\partial g}{\partial w'}$$

$$w'_{\text{new}} = w'_{\text{old}} - \eta \frac{\partial L}{\partial w'}$$

Again $\frac{\partial L}{\partial w'}$ can be found by chain rule.

$$\frac{\partial L}{\partial w} = \frac{\partial L}{\partial g} \cdot \frac{\partial g}{\partial o_t} \cdot \frac{\partial o_t}{\partial w}$$

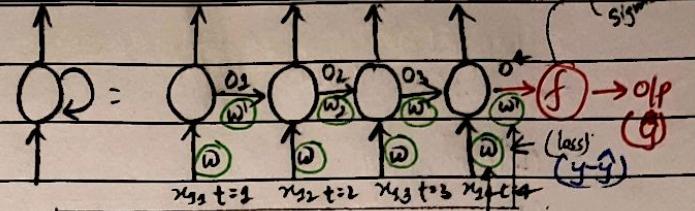
Chain rule

Similarly for $t+3$ w' , chain rule would be.

$$\frac{\partial L}{\partial w} = \frac{\partial L}{\partial g} \cdot \frac{\partial g}{\partial o_t} \cdot \frac{\partial o_t}{\partial w}$$

As we move back in backward propagation there are chances that this value will keep on getting smaller and smaller and if it is getting smaller & smaller weight updation is negligible.

In order to fix this issue we need to use LSTM RNN.



Let say there's one sentence,

"My name is Krish and I want to eat Pizza".

When we use RNN, RNN will try to make sure that each and every word will read context

in vectors and ensure every word is important.

Here "eat" is dependent on "Pizza". In this case

"eat" is just before "Pizza". So it will be able

to capture content by "I" is related to "Krish"

"I" & "My" are also related. It signifies that O_3

can be sometime dependent on O_1 and O_2 may

be dependent on O_4 . This sometimes because of

length it will not be able to capture content

properly. That is why we specifically use LSTM RNN.

⇒ RNN only remembers short term memory.

Problems with RNN :-

Suppose predict the next word.

- On Sunday I want to eat Pizza

- On Monday I want to eat _____

Q - Which RNN would be used here?

⇒ Many to One

RNN is good for short sentences, not for long.

LSTM RNN

→ LSTM RNN remembers long term memory.

→ Good for longer sentences.

Teacher's Signature

Deepak yd



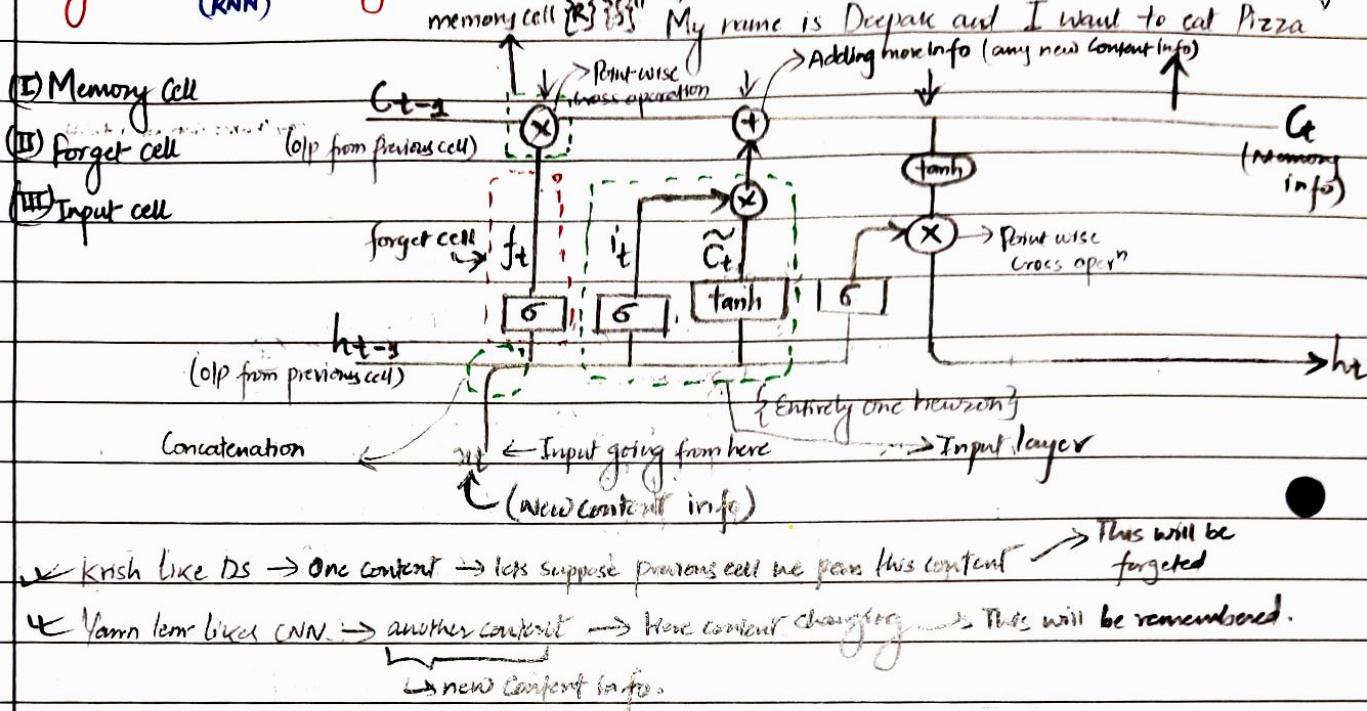
Date: Day 8

Topic: _____

* Colah blog → for detail on LSTM

Long Short Term Memory (LSTM) (RNN)

Suppose an example below:-



- ✓ Krish like DS → One content → lets suppose previous cell we pass this content → This will be forgotten
- ✓ Train like CNN → another content → Here content changing → This will be remembered.

(Q) Why LSTM RNN instead of RNN?

- Vanishing Gradient or Dead Neuron
- Content info. Deep RNN (Huge) \Rightarrow Gap.

\tanh → activation function.

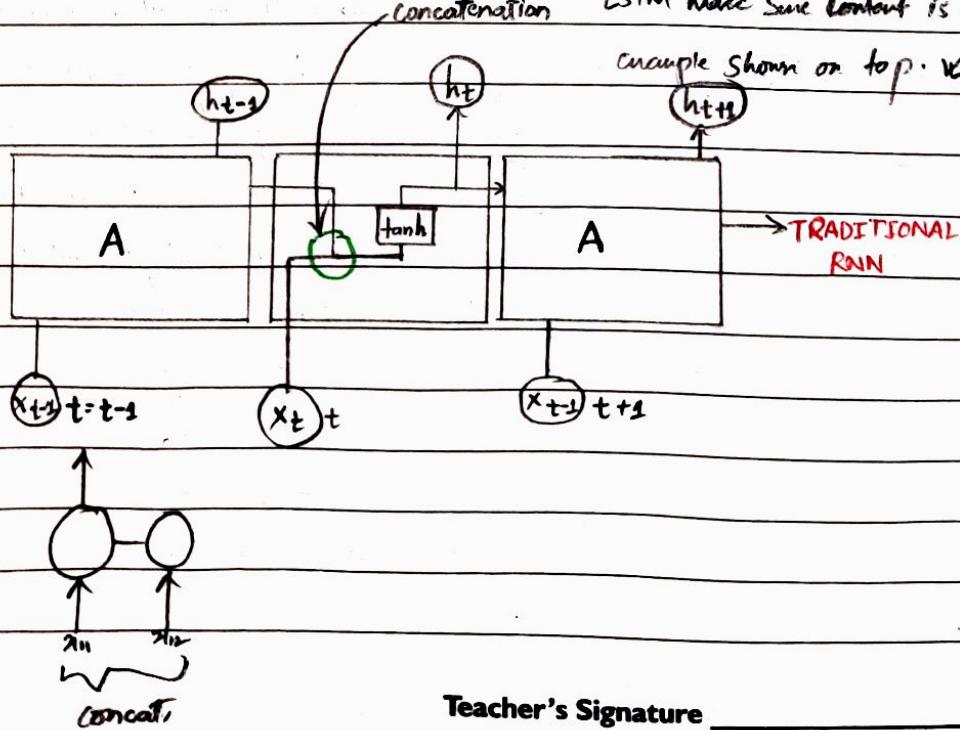
from previous state whatever o/p we get we are combining it with input we provide.

lets say if you have very deep neural network and some of the output is really dependent on first word then this content cannot be easily captured with the help of RNN, so we need LSTM RNN.

LSTM is important way to remember any important content.

LSTM make same content is capture.

Example shown on top. ✓





Date: _____

Topic: _____

(I) Why LSTM RNN?

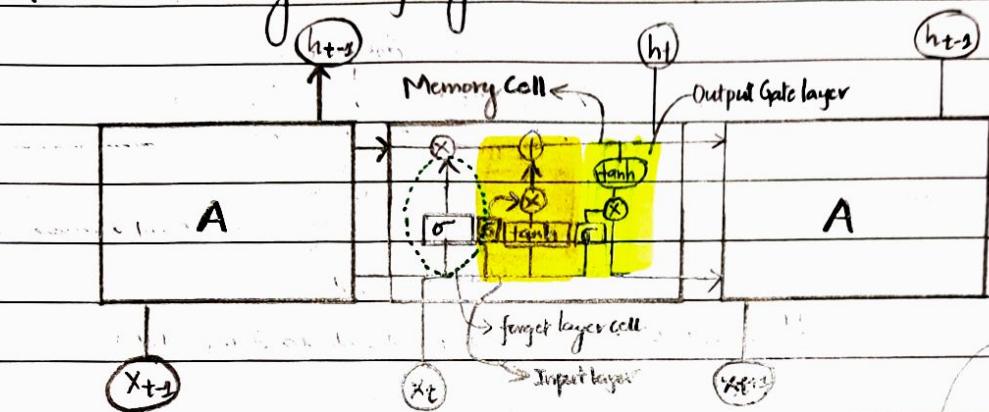
When we have long sentences and suppose there is dependency of the opf probably with the first word or second word, and we want to predict 1st word. Content will be switching in case of longer sentences.

$t=3 \quad t=2$
lets consider one long sentence,

$t=2 \quad t=2$
My name is KRISH and my friend name is

Here on reading LSTM model will find dependency Krish friend name can be related to Krish.

\rightarrow Content switching and dependency is happening



[6] \rightarrow neural network with sigmoid activation function.

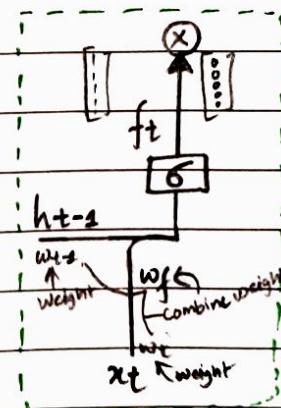
[tanh] \rightarrow neural network with tanh activation function.

\rightarrow Concatenate

\rightarrow Copy

\rightarrow Vector transform

O \rightarrow Point wise operation



forget gate layer

Memory Cell will act as conveyor belt

Similarly in memory cell you can add info or you can remove info.

Airport luggage chain

You can add remove luggage.

Forget layer cell: In this cell if previous info

and passed info are similar than most of

values will be nearer to 1 if this two will

not be similar, then most of values will be

nearer to zero. If it is nearer to zero that

Means we are making memory cell to forget the previous info.

$$\rightarrow f_t = G(w_f \cdot [h_{t-2}, x_t] + b_f)$$

Sigmoid transforms values between 0 to 1

Consider 2 sentences

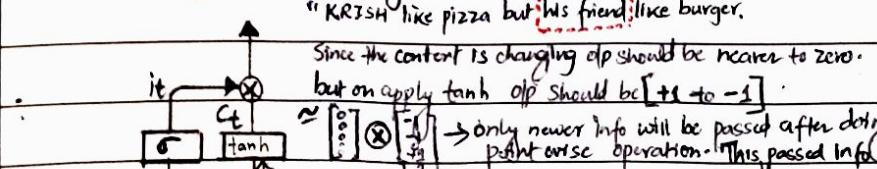
(1) Krish like pizza but he doesn't like burger

(2) Krish like pizza but his friend like burger.

Here in Sent 1 Content switching is not happening in Sent 2 Content switch is happening.

Incase of sent 3 info will not be lost because content switching is not happening.

But in sent 2 the content switching is happening so it will forget the previous content. "Krish like Pizza" will be forgotten.



Input layer cell: let us consider same example-2

"KRISH like pizza but his friend like burger."

Since the content is changing opf should be nearer to zero. but on apply tanh opf should be $[+1 \text{ to } -1]$

\rightarrow only newer info will be passed after doing point wise operation. This passed info will merged with above cell.

Important
Point wise operation \rightarrow dot product (\otimes)

\rightarrow addition.

$$it = \sigma(w_i \cdot [h_{t-2}, x_t] + b_i)$$

$$c_t = \tanh(w_c \cdot [h_{t-2}, x_t] + b_c)$$

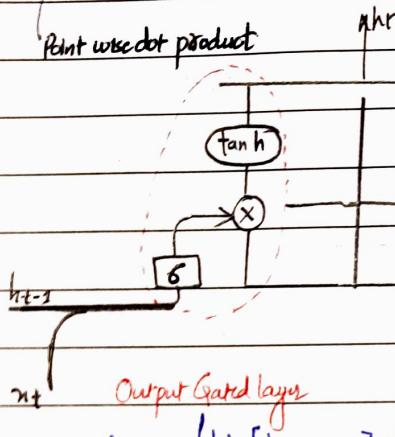
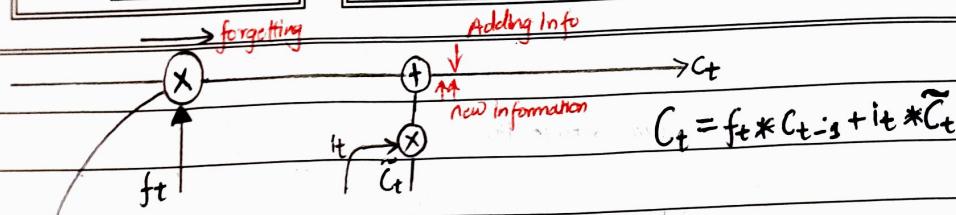
Teacher's Signature _____

Deepak ydu



Date: _____

Topic: _____



Output Gated layer: from the memory cell apply tanh operation so we get recent content info then it will get point wise dot product.

In short here we are trying to combine both the info previous forgotten info and the recent info and that is actually passed as an output to next cell. This will keep going with different - different timestamp.

In short all useless data is removed, only important data is considered and forwarded to the next cells.



Date: Day - 11

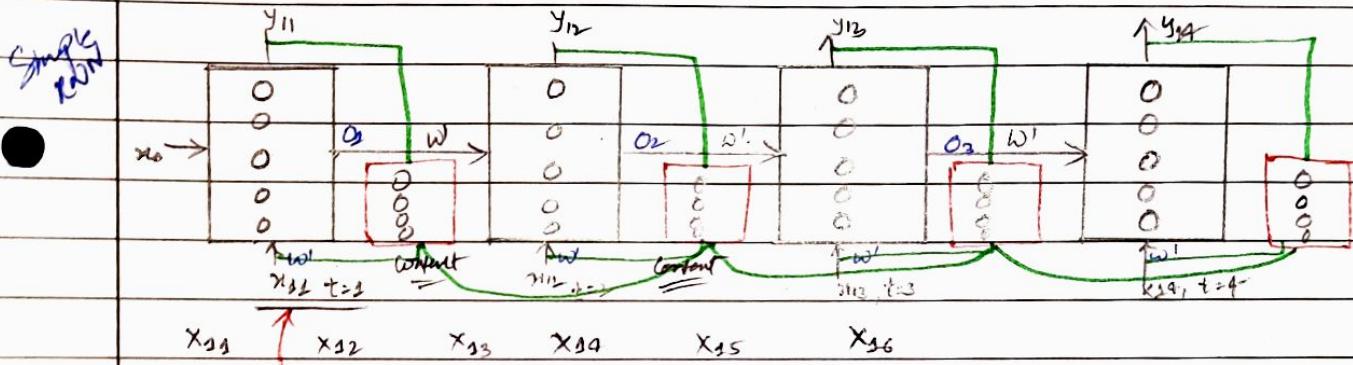
Topic: _____

Agenda

LSTM

- (I) Bidirectional LSTM RNN
- (II) Encoder and decoder
- (III) Transformer, Bert, GPT
- (IV) Practical Application.

Predicting the next word:

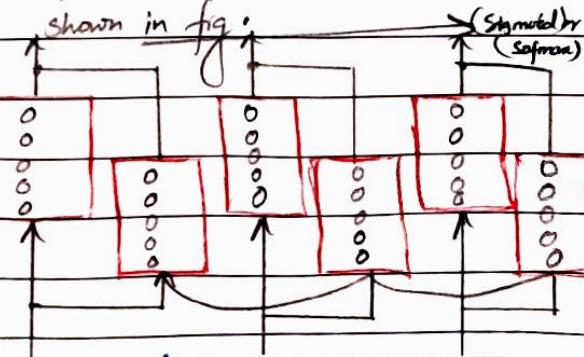
Bi-directional LSTM RNN:Example: KRISH LIKES TO EAT _____ IN BANGALORE

Embedding layers. (Many to Many RNN)

Here in case of simple RNN if we need to find y_{t+2} it will have content of all the previous layers like x_{t+1} , x_{t+2} , etc. But y_{t+3} will not have content of forward words. But in our example model should predict the next word based on one forward word, for eg, if its bangalore, as forward word model might predict food as Dosa or if its Bangalore as forward word it might predict the next word as Kabab.

Since our model has to predict the next word with respect to the content of the forward word. Simple RNN fails here and Bi-directional RNN comes into the picture.

Here bidirectional RNN is designed adding extra node next to every node of simple RNN model and adding content of previous up to the forwards as shown in fig:



Now on constructing bidirectional neural network

Now y_{t+2} will have content information of forward words and will have also content information of previous words. Now when training this model in correct way its possible to predict next of next words.

fig. Bi-directional LSTM:

Have content of both
→ forward Words
→ backward words

Teacher's Signature

Deepak ydru



Margalam

Date: Day 9

Topic: _____

Agenda

(i) Word Embedding layer → Tensorflow & Python } Suppose we have a sentence like below.

(ii) LSTM Sentiment Analysis.

Word2Vec

Word Embedding layer → (Words → Vector)
 ↓
 Training

[KRISH LIKES PIZZA]

(Vocabulary size = 500)

[SHYAM LIKES BURGERS]

Vocabulary size means
no. of unique words.

0	0
0	0
0	0
0	0
0	0

0	0
0	0
0	0
0	0
0	0

Steps
 (1) Sentences

(2) OHE - One Hot Encoding

(3) Padding → Post Padding, Pre-Padding

(4) OHE → Vectors

0	0
0	0
0	0
0	0
0	0

0	0
0	0
0	0
0	0
0	0

We do not use any pre-trained model, instead we're using dataset then we're going to train it using embedding layer and then find out the relationship. → Embedding layer and convert to vector. # Convert each and every value with some values
 frame size = 10 [0 0 0 0 0 | 120 150 160]

⇒ All the sequences length must be same in order to train into neural network.

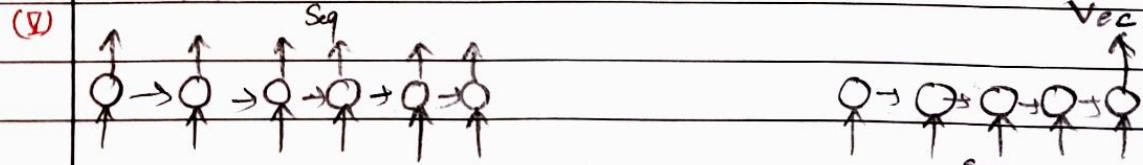
Let's look back to steps once more:

(I) Sentences

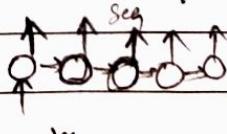
(II) Perform OHE (One Hot Encoding)

(III) Padding - (Pre and Post)

(IV) OHE - Vectors

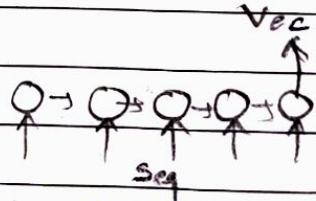


Sq Sq
 Seq - Seq (many to Many)



Vec

Seq-Vec



Seq - Vec



Vec

Seq

Teacher's Signature _____

Deepak ydu



Must Know
 ↗ RNN
 ↗ LSTM
 ↗ GRU

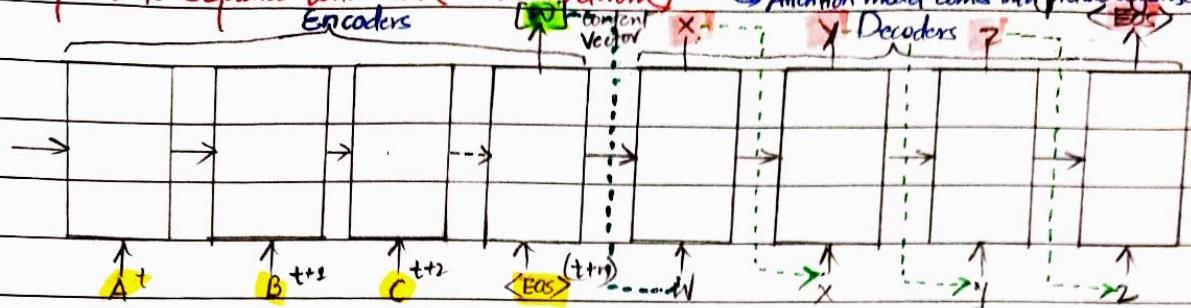
Date: NLP (27/08/18)

Topic: Encoders and Decoders:

$$x = \langle x_1, x_2, x_3 \rangle \quad y = \langle y_1, y_2, y_3 \rangle \text{ and sequence.}$$

Disadvantage: Not good accuracy with long sentences.
 ↗ Attention model comes into picture which uses bidirectional LSTM.

Sequence to Sequence with NN (Neural Network)



In sequence to sequence the model, takes in sequence of information and output sequence of information, that's one of many examples.

for ex: language translation.

for eg, English to French

If you given i/p as an image and if that particular image converted to text which is also called as "Image Captioning".

for ex: google image search.

If you search "Dog Image" in google search you can see images of dogs which is great example of "Image Captioning".

Also automated reply of "Unacademy" is an example.

Let's suppose ABC as English characters and

XYZ as French characters.

Applications:

→ grammarly

→ what app

→ Gmail.

→ Sentiment Analysis

Here encoder and decoder will be trained by pass A, B, C at different timestamps like t, t+1, t+2 etc finally end of string.

will be passed, we will be getting some content vector where all these 3 words A, B, C will be

represented in the form of content vector. After the representation it will be passed to decoder as an

input.

In case of image captioning in case of LSTM we use

Teacher's Signature

advance CNN (VGG16, RESNET) or any transfer learning etc.

at flatten layer pick up all vectors and pass them to decoders.

We can build encoder-decoder model with RNN, LSTM and GRU, mostly LSTM is considered best to build encoder-decoder model. Here the above diagram is built with LSTM model. Encoders doesn't contain o/p but decoders should contain both encoding i/p as well as o/p.

In language translation we need to pass text in the form of vectors. We may use various technique for that purpose like embedding layer, OHE, Word2Vec etc.

All the operations remain same as LSTM but only thing to remember is that there's no o/p for encoders.

After that first decoder will be giving o/p as x which will be passed to next decoder as an i/p. and then second decoder will give o/p as y which would be passed to third decoder as an i/p and so on.. and at last we get end of string.

LSTM actually works on probability. It is computed as below.

$$\Pr(\hat{y} | \langle x_1, x_2, x_3, \dots, x_n \rangle)$$

\hat{y} is combination of $y_1, y_2, y_3, y_4, \dots, y_n$

Deepak ydru

loss will be calculated b/w y and \hat{y} for this we can use Adam optimizer & loss by back propagation.

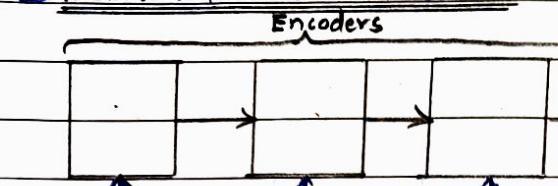


Date: _____

Topic: _____

Attention Model!

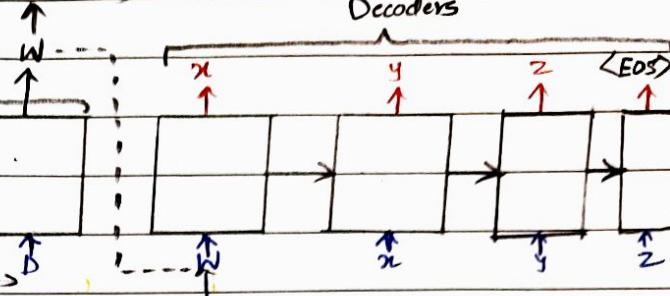
Problems with Encoder and Decoder:



The most important example in case of encoder-decoder is "Machine-Language-Translation".
Sequence to Sequence.

As we know encoder takes I/P and finally we take this particular vector "w" which we pass it to decoder for the translation.

(Content Vector)

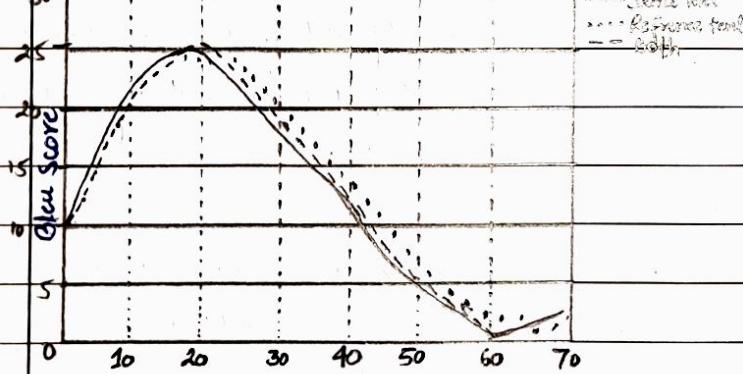


The "w" is basically called as "Content Vector". This vector is generated at the end of the string <EOS>.

Suppose if we take sentence as an I/P at the end of the string, for that particular string processing would be happening in the LSTM and finally we get 'w' vector.

This 'w' vector is then passed to decoder layer and based on that we're getting O/P w.r.t. the other language till the end of statement.

When sentence length is increasing up to specific length that time the accuracy of translation is decreasing.



"Loss in Model skill with increased sentence length"

"Sentence length ↑, accuracy ↓"

The major problem is we are not taking O/P w.r.t every I/P sentences. If we would take O/P of every I/P sentences it would capture the content properly w.r.t every individual sentences, which will obviously have content to nearest word.

This problem could be fixed using "bi-directional LSTM".

As we know that 'w' (content vector) is created based on the I/P words that we're given.

We're take just one vector at the end of the all input sentences. that we're not take

O/P w.r.t every I/P sentences. This single content vector cannot justify the content of every I/P sentences we've passed.

It is basically not being able to represent the content of whole sentences.

Which is the reason when we passed to decoder

It's not being able to translate sentences properly.

Deepak ydu

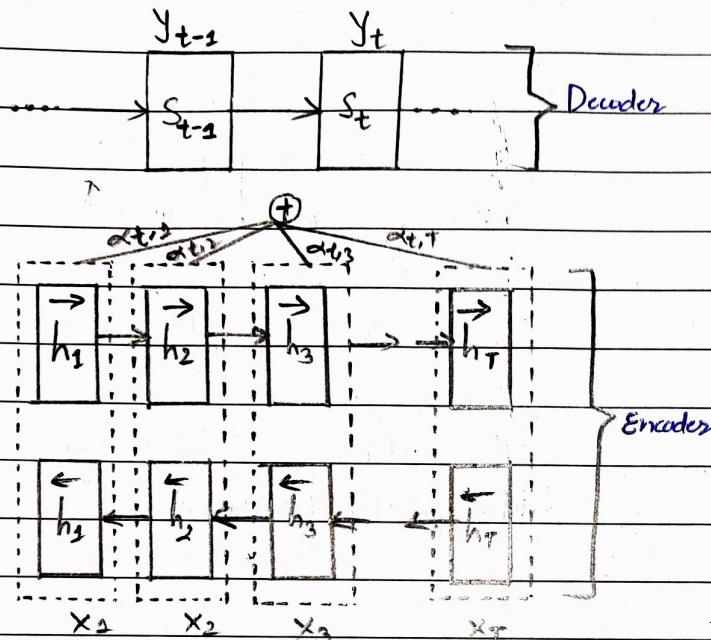
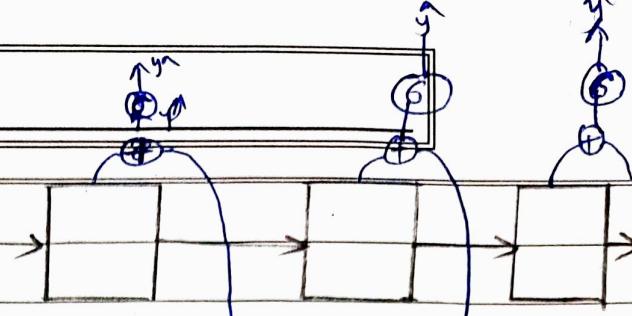
Teacher's Signature



Date: _____

Topic: _____

In case of normal Sequence-to-Sequence model
both encoder and decoder are uni-directional.

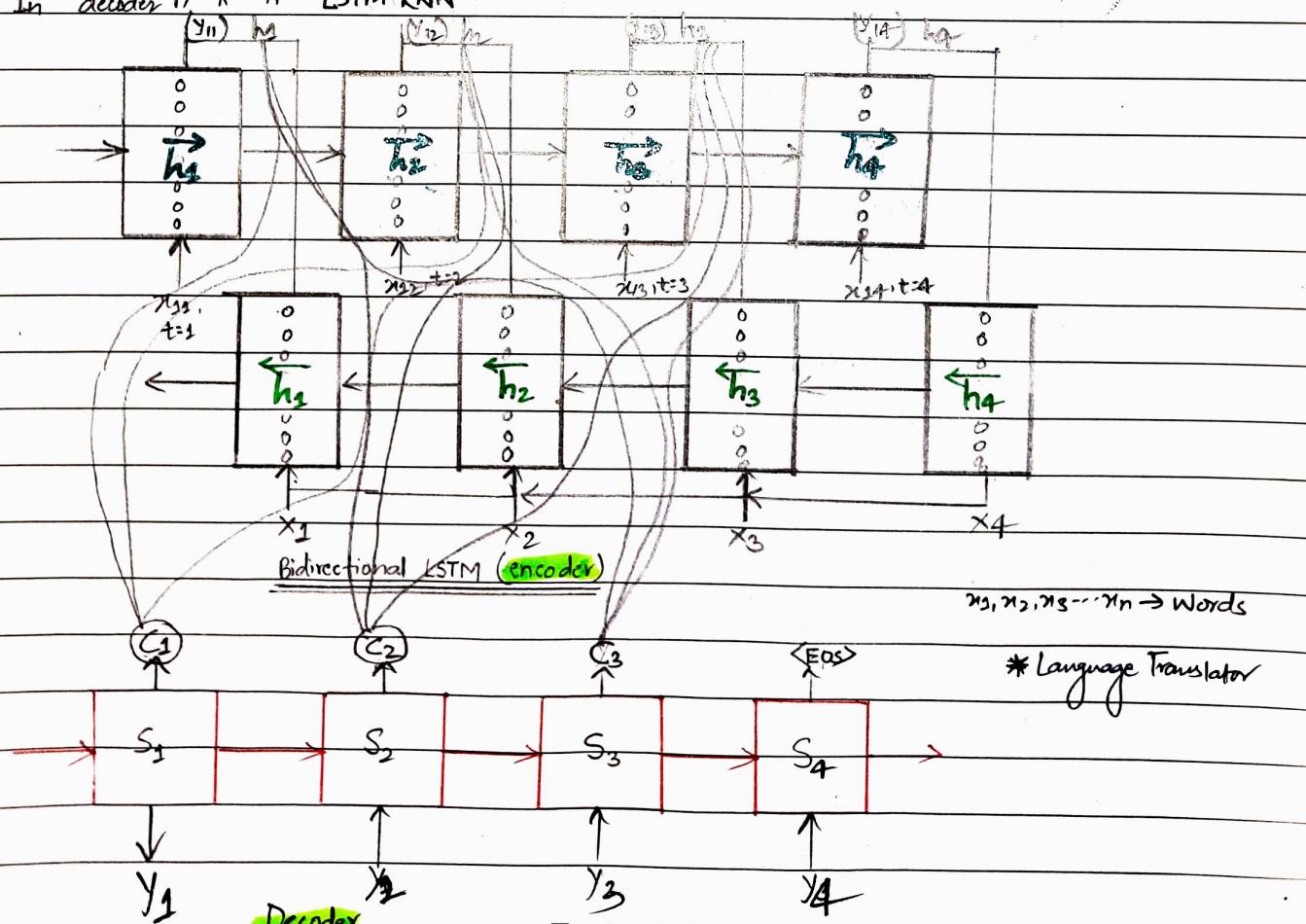


Using bi-directional LSTM we can not only predict previous content but also the future content of words.

In Encoder we will use "bi-directional LSTM RNN"

$h_1, h_2, \dots, h_n \rightarrow$ Hidden layer

In decoder we will use "LSTM RNN".



Deepak yd

Teacher's Signature _____



Date: _____

Topic: _____

→ LSTM RNN works on probability

→ When $T=3$, (window size is 3)

→ We get 3 α values. $\alpha_1, \alpha_2, \alpha_3$
for first word X_1

→ α values will be initialized by formula,

$$\sum_{i=1}^3 \alpha_i = 1$$

$\alpha_1, \alpha_2, \alpha_3 \Rightarrow$ Sum is always 1

$$\alpha_1 + \alpha_2 + \alpha_3 = 1$$

The content vector is generated by below formula.

$$C_i = \sum_{j=1}^{T_x} \alpha_{ij} h_j$$

(Multiplying o/p of bi-directional LSTM with α value)

This is the property of feed forward neural network.

The main reason of applying feed forward neural n/w
is for back propagation, because we need to update
value of $\alpha_1, \alpha_2, \alpha_3$.

Unless and until we don't update weights we can't get
good accuracy.

Feed forward neural network is same like ANN
(Artificial Neural Network)

Let's summarize,

Here consider below block as encoder

and upper block as decoder. In

decoder we are giving some
content vector ' C_i ' and in encoder

we're getting o/p as ' h '

That ' h ' based on T_x will be

generating some value $\alpha_1, \alpha_2,$
 α_3 and so on. The based on h

We will be assign $\alpha_1, \alpha_2, \alpha_3$
as weights.

We will derive these weights through
feed-forward-neural network.

$$C_i = \sum_{j=1}^{T_x} \alpha_{ij} h_j$$

Content Vector.

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

As we said earlier sum of

$\alpha_1, \alpha_2, \alpha_3$ is 1. We can use
this softmax formula to compute

This neural n/w will take 2 input

one is h and another as previous
state of decoder.

Such value.

After taking i/p it will derive $\alpha_1, \alpha_2,$

α_3 based on T_x value. Once it derives

then multiply h with $\alpha_1, \alpha_2, \alpha_3$
to get content vector ' C_i '

Where.

$$e_{ij} = \alpha(S_{i-1}, h_j)$$

α = function.

Teacher's Signature _____

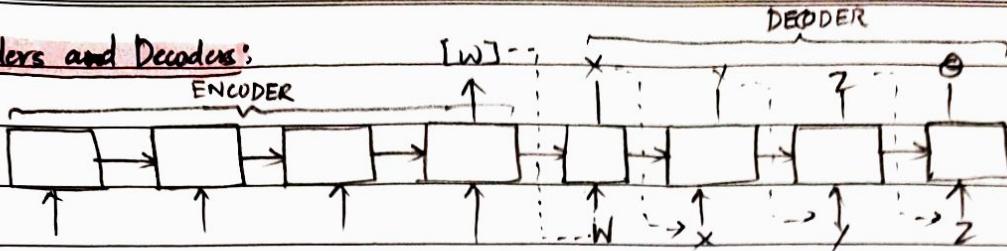
Deepak ydru



Date: _____

Topic: _____

Encoders and Decoders:

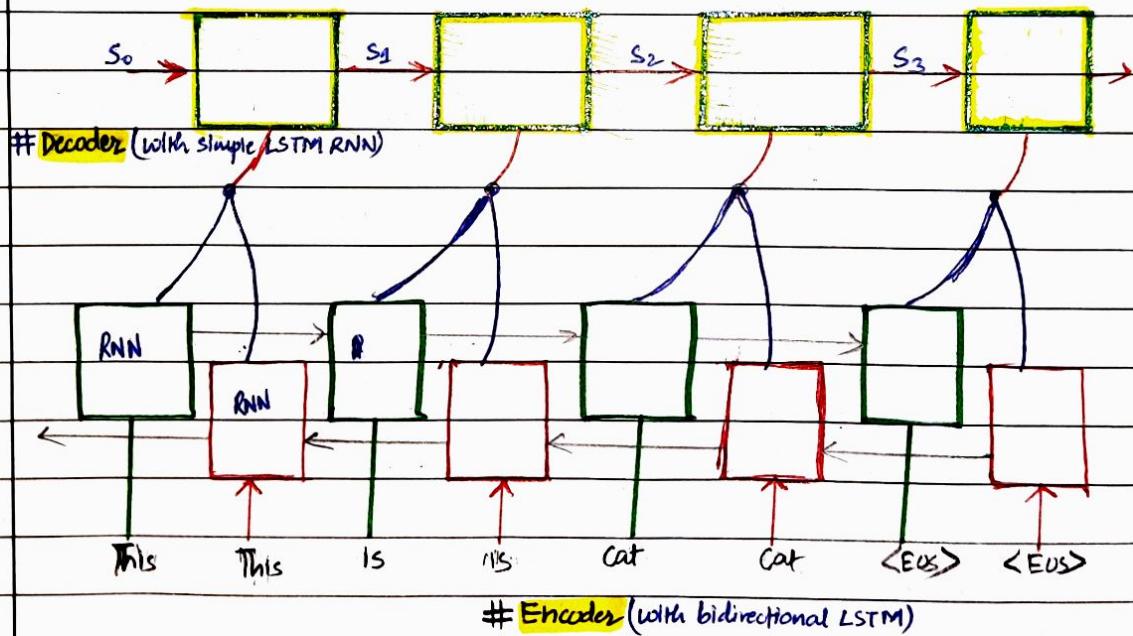


Disadvantages

- Uni-directional RNN (Won't understand sentence in backward direction).
- Collecting content at the end. (QnA, text summarization, text abbreviation, - in these cases all inputs won't contribute)

So we need some kind of a network which will be able to understand that what I need as an i/p to give some output. This is where concept of attention comes into picture.

Encoder-Decoder model with Single LSTM RNN



Encoder (with bidirectional LSTM)

Encoder-Decoder model with Bi-directional LSTM RNN.

In case of translating one language to another, in case of question answer system or fill in the blanks model, System which generate new story by understanding existing story. System that would tell a story about any particular image or any video or video analysis.

In all these cases we should look for a model which will be able to focus on certain set of i/p and based on that focused i/p it should be able to give some o/p.

This model may get failed because we're use one-to-one mapping. ($\text{Input} \rightarrow \text{O/P}$) but we're not focused anywhere.

Deepak Yadav

Teacher's Signature

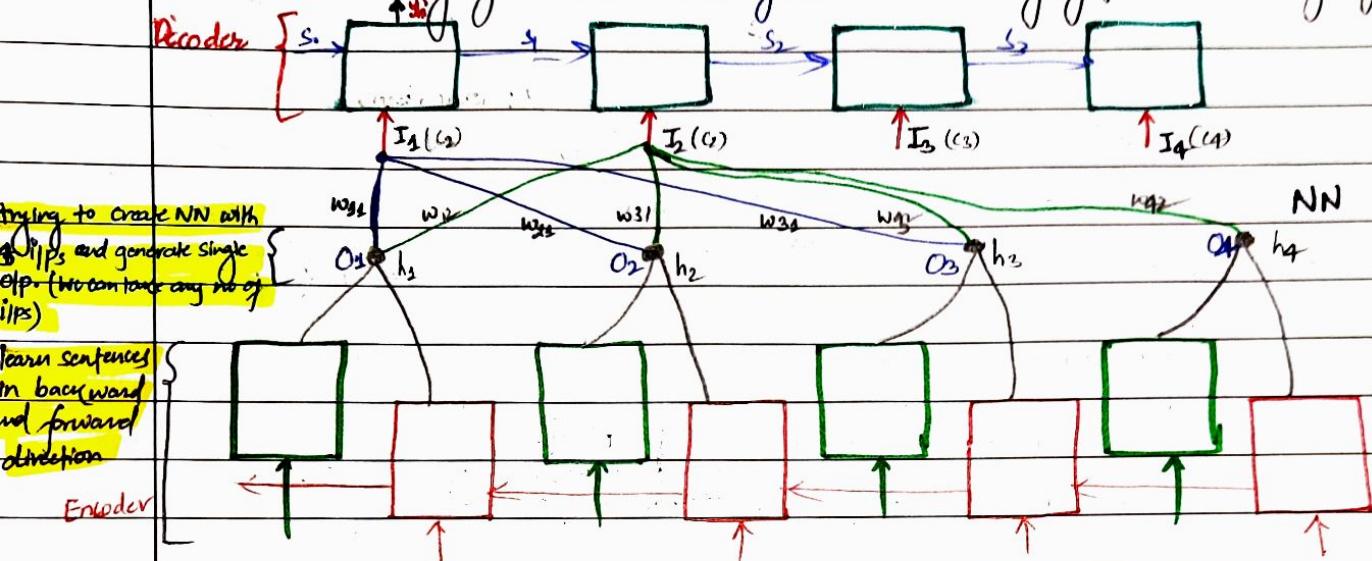


Date: _____

Topic: ATTENTION MODEL

This is where "attention model" comes into picture. In attention model "neural network" added along with the encoder-decoder. "neural network" helps in focusing, which i/p we would like to focus.

Ex:- The ball is running high << Stock Market Rising >> Here "is" is not having significance. "ball" "running" "high" focused.



As we pass "This is Cat" or "The ball is running high" "Here what we want the neural network to do (is that focus on all the words and find out the significance".

"In case of translation" it will focus on entire sentence and finds out significance like noun, pronoun, verb, adjective etc and priority of focusing would be different. This is what NN will do.

Because of applying these neural network researchers have named this concept as "attention model". So that it will be able to focus which word is responsible for giving some kind of o/p we desire.

The Content Vector c_i is given by:

$$c_i = \sum_{j=1}^{T_x} \alpha_{ij} h_j$$

c_i = Content Vector i/p.

T_x = no. of attention i/p that are injected in neural network. (In how many words we're supposed to pay attention)

$C_i = I_i, T_x = 4 \quad \{ h_j = O_i \} \quad \{ \text{above case} \} \quad \alpha_{ij} = w_{ij}$

h_j = combined output (outcome) from encoder side.

$$\alpha_{ij} = (w_{ij}) = \text{Weights} (w_{11}, w_{12}, \dots)$$

Deepak yd

for every i/p to the decoder we have one neural network, which is trying to pick and choose, which is trying to decide that which one is having a high priority and based on that it will be generating i/p to decoder and decoder will generate y_i finally.

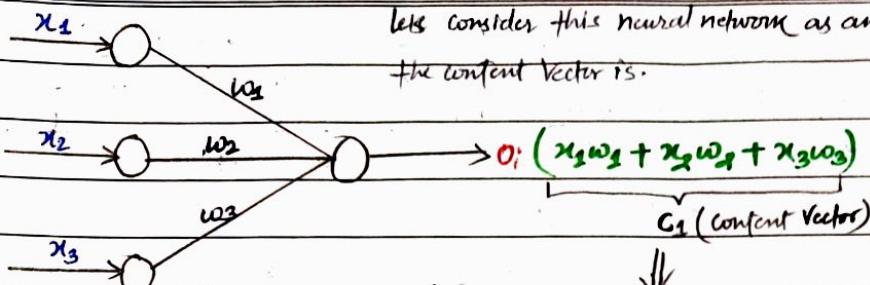
Input inside a decoder cell (one single decoder cell) is nothing but its basically a weight, which we're trying to multiply with weight and o_{ij} (o_{ij} = combined encoder)

Teacher's Signature _____



Date: _____

Topic: _____



The o/p that we're getting from the neural network which will go as i/p to the decoder network, which is responsible for giving final o/p.

$$C_1 = I_1 = O_1 w_{11} + O_2 w_{21} + O_3 w_{31} + O_4 w_{41}, \dots \dots \text{as per previous diagram.}$$

Whenever we're trying to generate weight, we're supposed to choose an equation with activation factor in such a way that it will not be able to overshoot or learn weight by any means.

' α ' is nothing but its weight as we concluded earlier. So the major thing is "How would I be able to control the weight?", because weight should not expand itself or shrink itself. for that we are supposed to put some kind of restrictions although it's a trainable parameter. That means network would try to train it through backward propagation.

Keeping this in mind researcher decides

$$\alpha_{ij} = \frac{\exp(c_{ij})}{\sum_{k=1}^K \exp(c_{ik})}$$

$c_{ij} = a(s_{i-1}, h_j)$ → attention function.

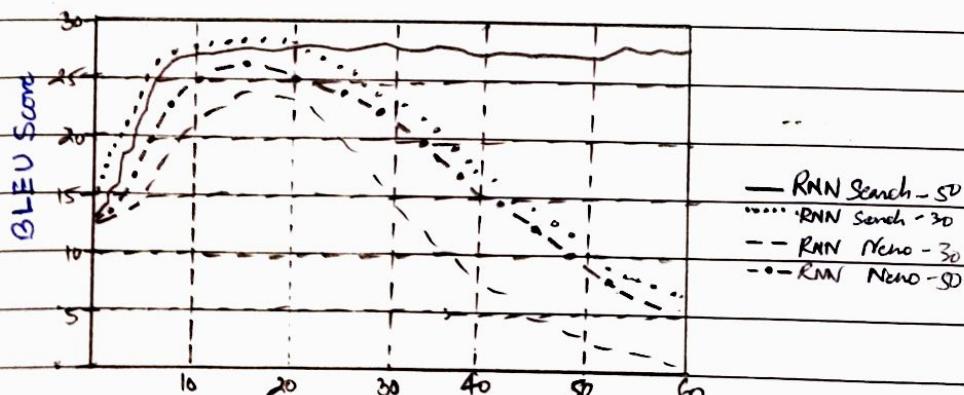
$s \rightarrow$ feedback

$$\sum \alpha_{ij} = 1 ?$$

$$\alpha_{ij} \geq 0$$

outcome of decoder network

outcome from encoder o/p.



Sentence length

Teacher's Signature _____

Deepak ydva



Date: _____

Topic: _____

Steps of using transformer in code:

- (I) Install the transformer \Rightarrow ! Pip install transformers
- (II) Call the pre-trained model (eg. Distilbert)
- (III) Call the tokenizers (for specific pre-trained model, we have specific tokenizers) if you're using pre-trained related to "Distilbert" then use tokenizer relating to same.
 \Rightarrow from transformers import DistilBertTokenizersFast
tokenizers = DistilBertTokenizersFast.from_pretrained('distilbert-base-uncased')
- (IV) Convert encodings into Dataset Objects (eg. tensors in case of tensorflow)

Deepakydu

Teacher's Signature _____



Date: _____

Topic:

TRANSFORMER & BERT

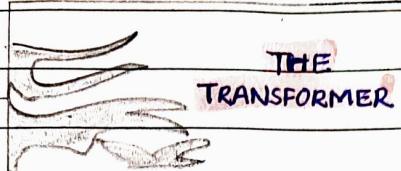
* Jay Naman
A. Krish Naik

29/01/23

TRANSFORMER

INPUT

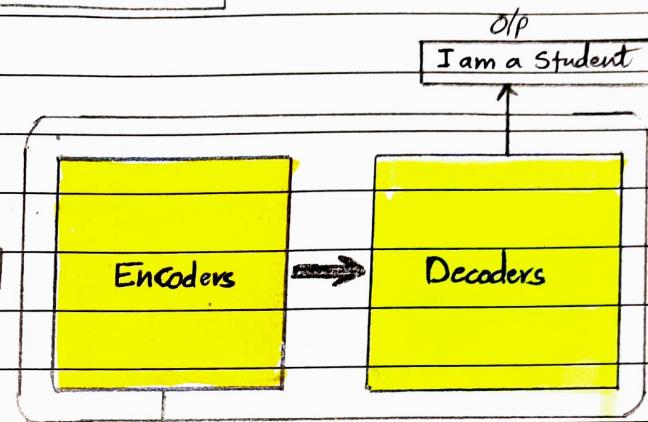
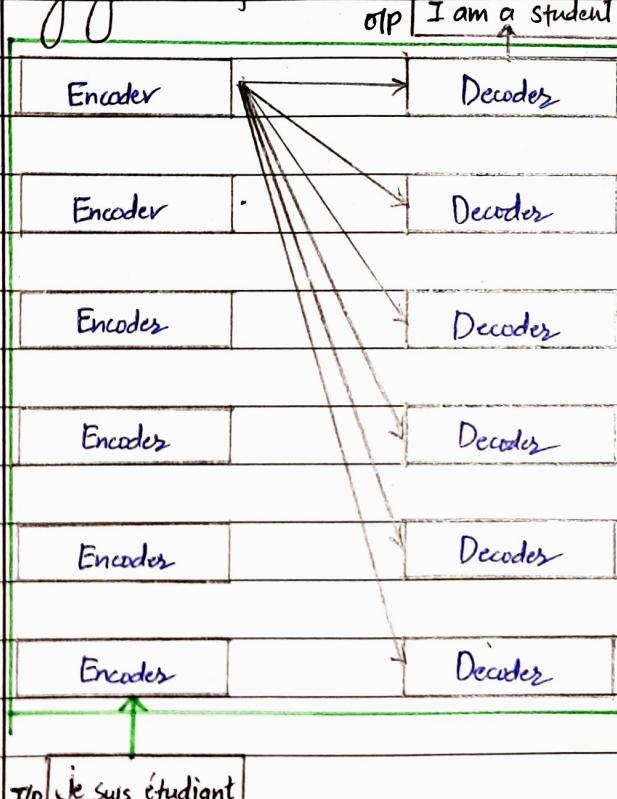
Je suis étudiant



OUTPUT

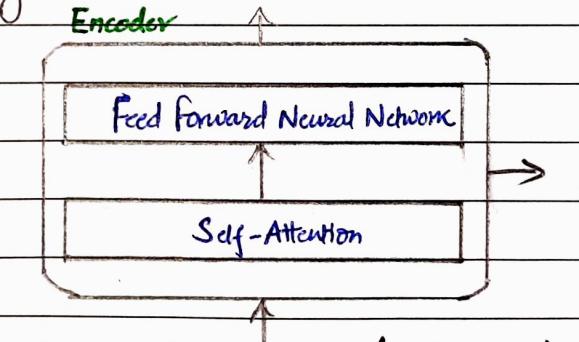
I am a Student

Let's begin by looking at the model as a single black box. In a machine translation application, it would take a sentence in one language, and output its translation in another.



Encoding component, decoding component and connection b/w them.

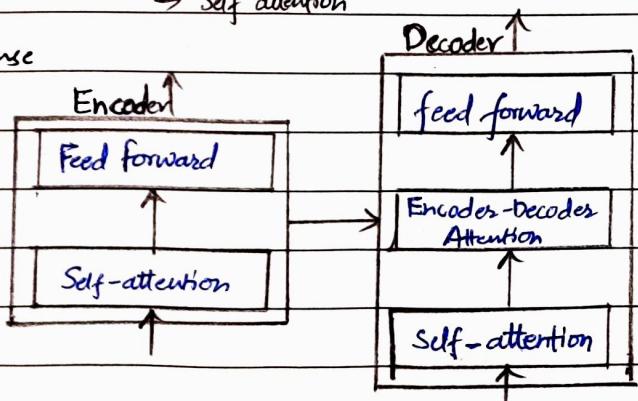
The encoders are all identical in structure (yet they do not share weights). Each one is broken down into two sub-layers.



The stack of 6-encoders and 6-decoders provided good result to the researchers.

- feed forward neural network (similar to ANN)
- Self attention

This is kind of hyper-parameter, you may use them at your convenient.



Deepak ydu

Teacher's Signature _____

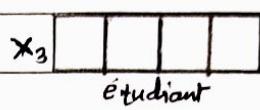
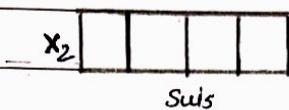
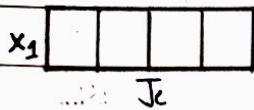


Date: _____

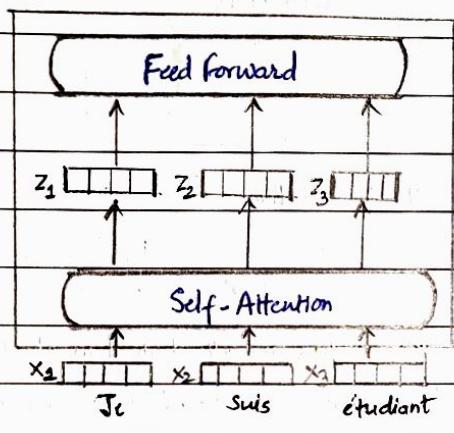
Topic: _____

Bringing The Tensor Into the Picture.

As in the case of NLP applications in general we begin by turning each input word into a vector using an embedding algorithm.



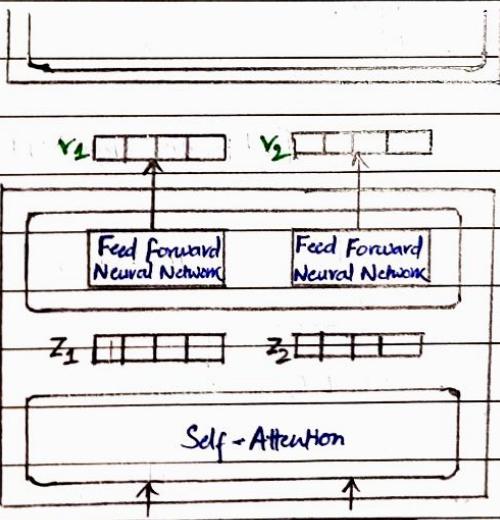
Each word is embedded into a vector of size 512. We'll represent those vectors with these simple boxes.



Here each words like "Je", "Suis", "étudiant" is converted into vector through Word2Vec with 512 dimension and this is basically passed to self-attention. In this process all the words are passed parallelly, all at one time. After that it goes to feed-forward neural network.

Self-Attention at a High level

Encoder 2



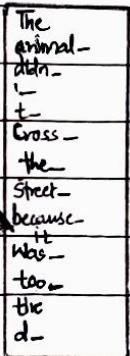
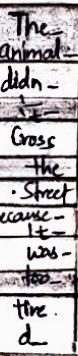
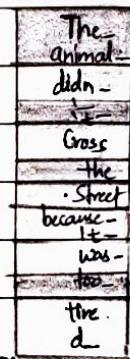
Here lets take an example of below sentence.

"The animal didn't cross the street because it was too tired"

What does it in the sentence refer to? Is it referring to the street or to the animal? It's a simple question to human, but not as simple to an algorithm.

When the model is processing the word "it", self-attention allows it to associate "it" with "animal".

The word at each position passes through a self-attention process. Then, they each pass through a feed-forward neural network... the exact same network with each vector flowing through it separately.



As we are encoding the word "it" in encoder. The animal are strongly correlated and baked a part of its representation into the encoding of "it".

Deepak ydu

Teacher's Signature _____



Date: _____

Topic: _____

Self Attention in Detail

→ Word2Vec Embedding (Converted word to 512 dim)

Step I

Input (512 dim) Thinking Machines

Embedding x_1 [] [] [] [] x_2 [] [] [] []

Queries (64 dim) q_1 [] [] [] [] q_2 [] [] [] []

Keys k_1 [] [] [] k_2 [] [] []

Values v_1 [] [] [] v_2 [] [] []

WQ

WK

WV

Why it is getting multiplied

Because we know that

formula input \times weight + bias

and given activation function.

applied in case of DL model

Concept of ANN.

But in this case we just taking this step and multiplying by WQ we get queries q , parallel which would be 64 dimensions.

Multiplying x_1 by the WQ weight matrix produces q_1 the "query" vector associated with that word.

We end up creating a "query" a "key" and a "value" projection of each word in the input sentence.

Similarly x_1 Vector gets multiplied with WK and then we get k_1 Values. and x_2 with WK get k_2

Values.

Same in case of V_1 and V_2

$x_1 \times W^V \rightarrow V_1$

$x_2 \times W^V \rightarrow V_2$

Step-II

<u>Input</u>	Thinking	Machines
Embedding	x_1 [] [] []	x_2 [] [] []
Queries	q_1 [] [] []	q_2 [] [] []
Keys	k_1 [] [] []	k_2 [] [] []
Values	v_1 [] [] []	v_2 [] [] []
Score	$q_1 \cdot k_1 = 112$	$q_2 \cdot k_2 = 96$

The score is calculated by taking the dot product of the query vector with key vector of the respective word we're scoring. So if we're processing the self attention of the word in position 1, the first score would be the dot product of q_1 and k_1 . The second score would be the dot product of q_1 and k_2 .

In the second step we will be creating score for each word by multiplying query with the keys.

Step-III

<u>Input</u>	Thinking	Machines
Embedding	x_1 [] [] []	x_2 [] [] []
Queries (64 dim)	q_1 [] [] []	q_2 [] [] [] (64 dim)
Keys	k_1 [] [] []	k_2 [] [] []
Values	v_1 [] [] []	v_2 [] [] []
Score	$q_1 \cdot k_1 = 112/8$	$q_2 \cdot k_2 = 96/8$
Divide by $8(\sqrt{64})$	14	12
Software	0.88	0.12

What ever score we get we are dividing by 8 (square root of the dimension of key vector) that is $\sqrt{64} = 8$.

This leads to have more stable gradients. There could be other possible values here, but this default

After that we're passing result to softmax operation. Softmax normalizes the scores so they're all positive and add up to 1.

Teacher's Signature

Deepak yd



Date: _____

Topic: _____

Step-IV As discussed we're applying Softmax activation function in Step-IV. Whenever you provide any I/O to softmax the total value is 1.

Here we got values after dividing Score by 8 are 14 and 12 applying Softmax activation to these values we get 0.88 and 0.12 and summing them we get $(0.88 + 0.12) = 1$. As per Softmax score more importance is given to "Thinking".

<u>Input</u>	<u>Thinking</u>	<u>Machines</u>	In this step whatever softmax value we got in step-IV is multiplied by Values, we get new vector v'_1 and v'_2 .
Embedding Queries	x_1 [] [] []	x_2 [] [] []	
Keys	q_1 [] []	q_2 [] []	
Values	k_1 [] []	k_2 [] []	
Score	$q_1 \cdot k_1$	$q_1 \cdot k_2$	
\times Divide by 8 (\sqrt{dk})	14	12	Finally we add v'_1 and v'_2 to get Z_1 .
Softmax	0.88	0.12	
Softmax Value	v'_1 [] []	v'_2 [] []	
Sum	Z_1 [] []	Z_2 [] []	

Step-VI Z_1 will be the output of the self attention layer when you pass "Thinking" and Z_2 will be the o/p of the machine word when you are passing through the "Self attention".

That concludes the 'self-attention calculation'. The resulting vector is one we send along to feed-forward neural network. In the actual implementation, however this calculation is done in matrix form for faster processing. So let's look at that now.

Matrix Calculation of Self-Attention

The first step is to calculate the Query, Key and Value matrices. We do that by passing our embeddings into a matrix X and multiplying it by weight matrices we've trained (W_Q, W_K, W_V).

$$\begin{array}{ccc}
 \text{Thinking Machine} & \xrightarrow{\text{X } (512 \text{ dim})} & \xrightarrow{\text{W}_Q} Q \text{ (4 dim)} \\
 & \times & = \\
 & &
 \end{array}
 \quad
 \begin{array}{ccc}
 X & \xrightarrow{\text{W}_K} K \\
 \times & = & K^T \\
 & &
 \end{array}
 \quad
 \begin{array}{ccc}
 X & \xrightarrow{\text{W}_V} V \\
 \times & = &
 \end{array}$$

Deepak ydu

Every row in the X matrix corresponds to a Teacher's Signature word in the input sentences.

weights will be initialized randomly.

$$\text{Softmax} \rightarrow \left(\frac{\left(\begin{array}{c} Q \\ K^T \end{array} \right) \cdot \left(\begin{array}{c} V \\ K^T \end{array} \right)}{\sqrt{dk}} \right)$$

o/p of self-attention layer.



Date: _____

Topic: _____

We are considering W^Q , W^K , W^V the same weights for all the words, we will call it as 'single head attention'.

This is further improved when we use "Multi-headed - attention".

If we do the same self-attention calculation we outlined above, just 8-different times with different weight matrices, we end up with eight different Z-matrices.

The Beast with Many Heads:

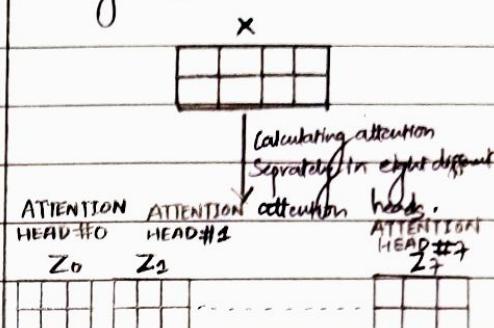
"Multi-headed" Attention

→ By using multi-headed attention, improves the performance of the attention layer in two ways.

(1) It expands the model's ability to focus on different positions.

In the example above "Z₁" contains a little bit of every other encoding, but it could be dominated by the actual word itself.

If we're translating a sentence like "The animal didn't cross the street because it was too tired", it would be useful to know which word "it" refers to.



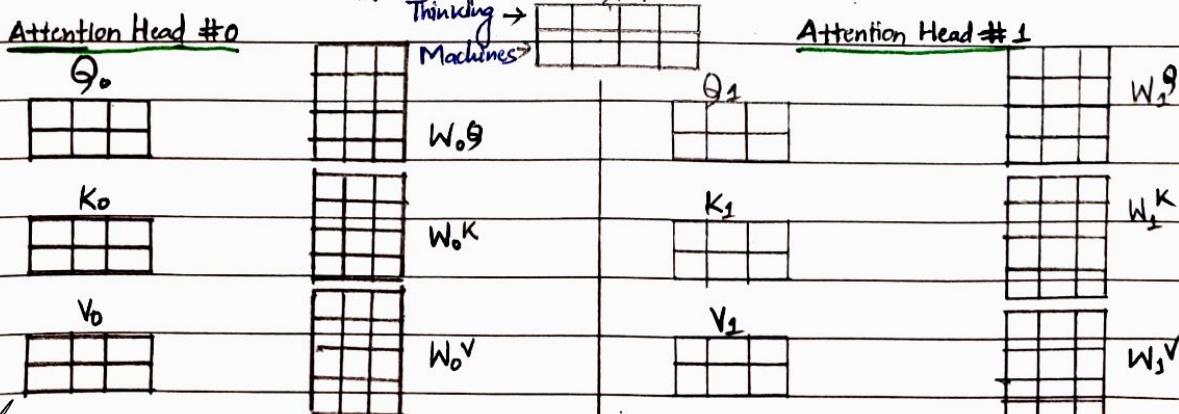
This leaves us with a bit of a challenge. The feed-forward is not expecting 8 matrices - it's expecting a single matrix!

We concat the matrices then multiple them by an additional weight matrix W_O .

(2) It gives the attention layer multiple "representation subspaces".

As we'll see next, with multi-headed attention we have not only one, but multiple sets of Query / Key / Value weight matrices (the transformer uses eight attention heads, so we end up with eight sets for each encoder / decoder). Each of these sets is randomly initialized. Then after training, each set is used to project the input embeddings (or vectors from lower encoder / decoder) into a different representation subspace.

→ P.T.O



Deepak ydu

With multi-headed attention, we maintain separate Teacher's Signature

$Q/K/V$ weight matrices for each head resulting in different $Q/K/V$ matrices. As we did before we multiply X by $W_Q/W_K/W_V$ matrices to produce $Q/K/V$ matrices,



Date: _____

Topic: _____

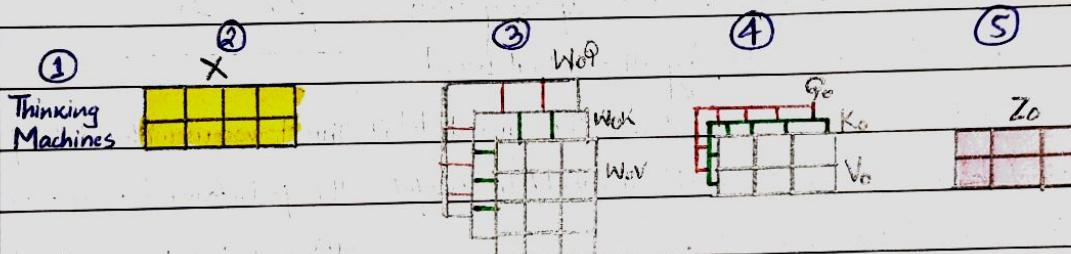
 $Z_0 \ Z_1 \ Z_2 \ Z_3 \ Z_4 \ Z_5 \ Z_6 \ Z_7$ 

→ Concatenate all the attention heads.

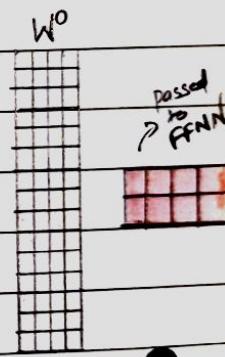
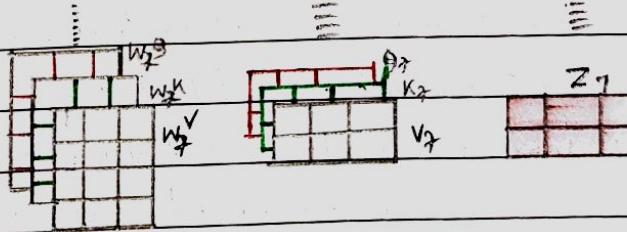
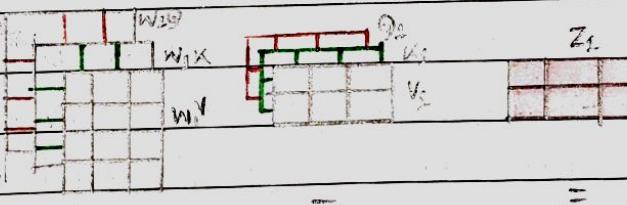
(II) Multiply with a weight matrix W^O that was trained jointly with the model.

 W_O

(III) The result would be the Z -matrix that captures information from all the attention heads. We can send this forward to the FFNN.

 Z 

* In all encoders other than #0 we don't need embedding, we start directly with the OLP of encoder.



① Input Sequence

② We embed each word

③ Split into 8 heads, we multiply X or R with weight matrices (attention head 1 to 7)

④ Calculate attention using the resulting $Q/K/V$ matrices

⑤ Concatenate the resulting Z matrices then multiply with weight matrix W^O to produce the output of the layer.

Deepak ydu

Teacher's Signature _____



Date: _____

Topic: _____

The animal didn't cross the street because it was too tired.

The animal didn't cross the street because it was too tired.

The animal didn't cross the street because it was too tired.

If we add all the attention heads to picture,

The animal didn't cross the street because it was too tired.

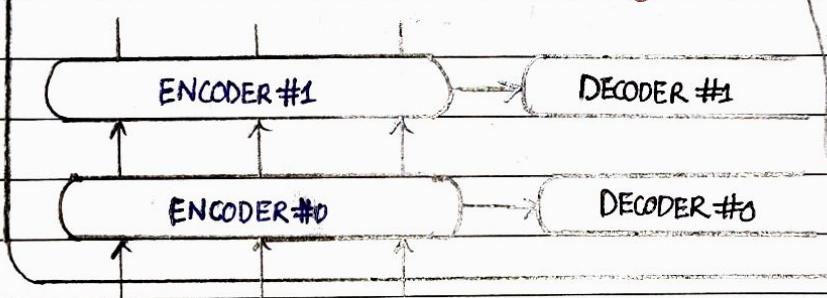
As we encode the word "it" One attention head is focusing most on "the animal", while another is focusing on "tired" - in a sense, the model's representation of the word "(it)" bakes in some of the representation of both "animal" and "tired".

* check blog for more clear picture.

Representing the order of Sequence using Positional Encoding:

Positional encoding deals with calculating distances between the words, lesser the distance the much closer the words are and vice-versa.

→ ordering or word is also important.



Embedding with Time Signal

$$X_1 = \boxed{\text{Yellow}} \quad \boxed{\text{Yellow}} \quad \boxed{\text{Yellow}}$$

$$X_2 = \boxed{\text{Yellow}} \quad \boxed{\text{Yellow}} \quad \boxed{\text{Yellow}}$$

$$X_3 = \boxed{\text{Yellow}} \quad \boxed{\text{Yellow}} \quad \boxed{\text{Yellow}}$$

Positional Encoding

$$t_1 = \boxed{\text{Yellow}} \quad \boxed{\text{Yellow}} \quad \boxed{\text{Yellow}}$$

$$t_2 = \boxed{\text{Yellow}} \quad \boxed{\text{Yellow}} \quad \boxed{\text{Yellow}}$$

$$t_3 = \boxed{\text{Yellow}} \quad \boxed{\text{Yellow}} \quad \boxed{\text{Yellow}}$$

Embeddings

$$X_1 = \boxed{\text{Green}} \quad \boxed{\text{Green}} \quad \boxed{\text{Green}}$$

$$X_2 = \boxed{\text{Green}} \quad \boxed{\text{Green}} \quad \boxed{\text{Green}}$$

$$X_3 = \boxed{\text{Green}} \quad \boxed{\text{Green}} \quad \boxed{\text{Green}}$$

Input

je

Suis

étudiant

To give the model a sense of the order of the words, we add positional encoding vectors. The values of which follow a specific pattern.

If we assumed the embedding has a dimensionality of 4, the actual positional encoding would look like this.

Positional Encoding

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

$$\begin{matrix} 0.94 & 0.001 & 0.91 & 1 \end{matrix}$$

$$\begin{matrix} 0.91 & 0.002 & 0.92 & 1 \end{matrix}$$

Embeddings

$$\boxed{\text{Green}} \quad \boxed{\text{Green}} \quad \boxed{\text{Green}}$$

$$\boxed{\text{Green}} \quad \boxed{\text{Green}} \quad \boxed{\text{Green}}$$

$$\boxed{\text{Green}} \quad \boxed{\text{Green}} \quad \boxed{\text{Green}}$$

Input

je

Suis

étudiant

Deepak ydu

Teacher's Signature _____



Date: _____

Topic: _____

The Residuals

ENCODER #2

Add & Normalize

Feed Forward Feed Forward

Add & Normalize

Self - Attention

 x_1 Thinking x_2 Machines

Add & Normalize

Feed Forward Feed Forward

Add & Normalize

Self Attention

Add & Normalize

Feed Forward Feed Forward

Add & Normalize

Self - Attention

POSITIONAL ENCODING

 x_1 Thinking x_2 Machines

ENCODER #1

Add & Normalize

Feed Forward Feed Forward

 z_1 z_2

Add & Normalize

LayerNorm (H x W)

 z_1 z_2

Self - Attention

 x_1 Thinking x_2 Machines

Add & Normalize

Feed Forward Feed Forward

Softmax

Linear

DECODER #2

Add & Normalize

Feed Forward Feed Forward

Add & Normalize

Encoder-Decoder Attention

Add & Normalize

Self - Attention

Here we have passed all the input at once but we will be generating output one-by-one.

Deepak ydu

Teacher's Signature _____

Date: _____

Topic: _____

BERT

BERT (Bidirectional Encodes Representations from Transformers) is a recent paper published by researchers at Google AI Language.



BERT's key technical innovation is applying the bi-directional training of transformers, a popular attention model, to language modelling.

The paper's results shows that a language model which is bidirectionally trained can have a deeper sense of language context and flow than single-direction language models. In the paper researchers detail a novel technique named Masked LM (MLM) which allows bi-directional training in models in which it was previously impossible.

BERT makes use of transformers, an attention mechanism that learns contextual relations between words (or sub-words) in a text.

In its Vanilla form, Transformer includes two separate mechanisms - an encoder that reads text input and a decoder that produces a prediction for the task. Since BERT's goal is to generate a language model, only the encoder mechanism is necessary.

Application (Uses of BERT)

BERT can help with the following tasks:-

- (I) Text encoding / Similarity retrieval
- (II) Text Summarization / Response Selection.
- (III) Question Answering

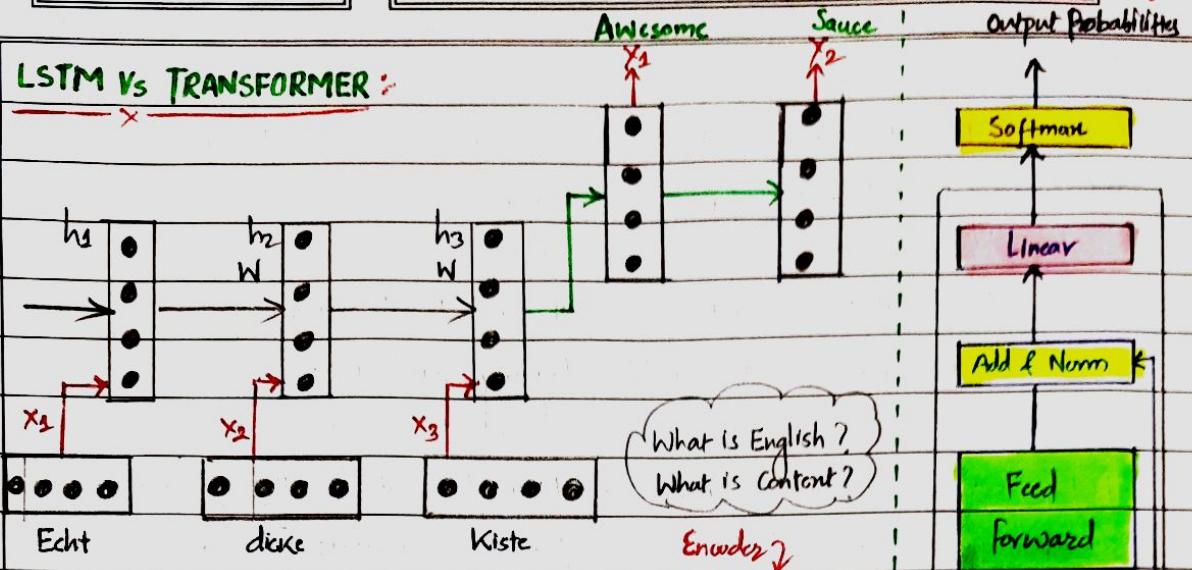
How to Map English Word to French Word?

Date: _____

Topic: _____

Decoder ↘

LSTM VS TRANSFORMER:

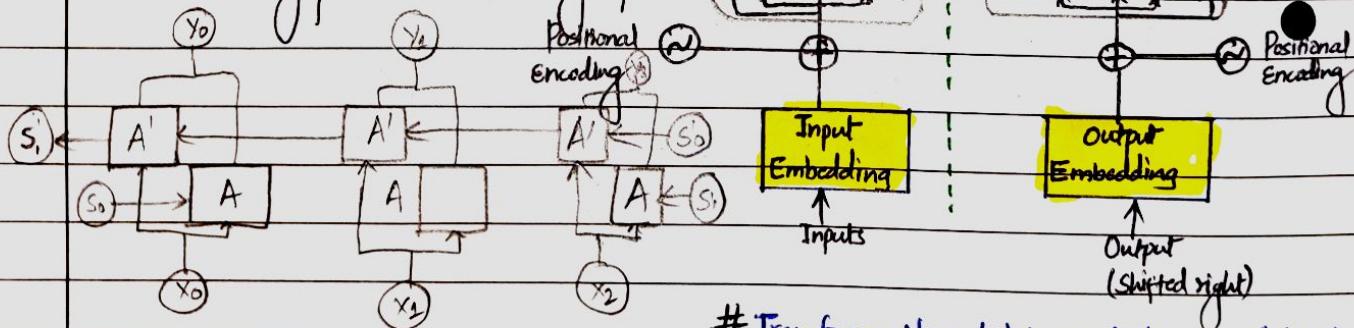


LSTM Network

Before the transformer, LSTM was used for the machine translation process.

~~LSTM Networks has below problems:-~~

- Slow to train
- Words are passed sequentially and are generated sequentially, it can take significant no. of timesteps for neural network to learn (time-consuming)
- True meaning of words are not entirely captured.



Transformer Neural Network Architecture initially

Created to solve problem of Neural Machine Translation

→ first they are faster as words can be processed simultaneously.

→ Second content of word is better learnt

Even bidirectional LSTM "technically learning left to right and right to left." content separately and concatenating them so true content is slightly lost.

But transformer architecture addresses some of these concerns. → As they learns from both direction simultaneously.

Teacher's Signature _____

Deepak ydu



Date: _____

Topic: _____

lets see transformer in action.

lets take an example of language translation "english to french"

The transformer consists of **2 key components** "Encoder" & "Decoder".

"Encoder" takes english words simultaneously and it generates "embeddings for every words", Structurally These "embeddings" are Vectors that encapsulate the meaning of the word.

"Similar words have closer numbers in their vectors."

"The Decoder" takes these embeddings from encoder and the previously generated words of the translated french sentence and then uses them to generate ^{next} a french word, and we keep generating the french translation one word at a time until the end of sentence is reached.

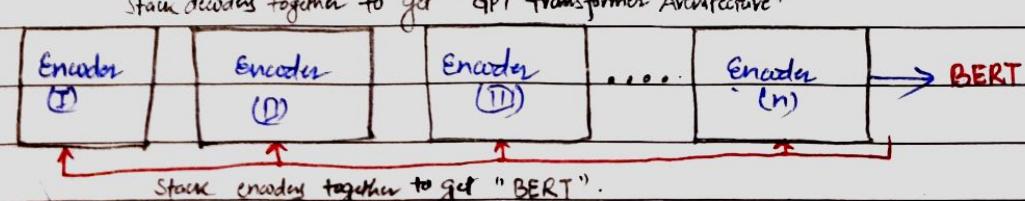
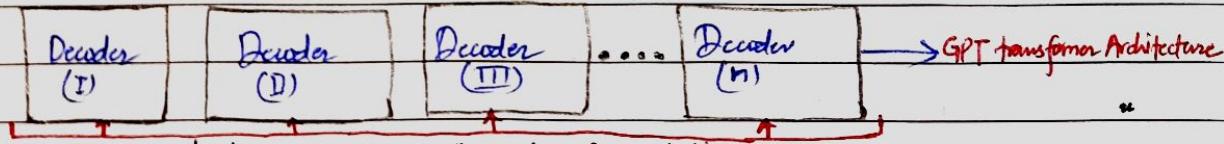
What makes this conceptually so much more appealing than some LSTM cell is that we can physically see a separation in tasks. The encoder learns what is "English" what is "grammar" and more importantly "what is "Content".?

The decoder learns how do "english word relates to french words".

"Both of these encoder and decoder have some underlying understanding of language."

We stack the decoders and we get the GPT transformer architecture.

Conversely if we stack just the encoders we get "BERT" (Bidirectional Encoder Representation from Transformers)



The OG transformer solves problem of "Neural Machine Translation".

We can use BERT for followings :-

(I) Language Translation

(II) Question Answering

(III) Sentiment Analysis

(IV) Text summarization and many more...

Needs language understanding

BERT Training
We can train BERT to understand language and fine tune BERT to learn specific tasks!

> Pre-training ✓
> fine-tuning ✓

Teacher's Signature _____

Deepak ydva



Date: _____

Topic: _____

Pre-training:

The goal of pre-training is to make BERT learn

Pre-training (Pass 1): "What is language? What is content?"

BERT learns language by training on two "Unsupervised tasks"

Simultaneously they are:

Masked Language Model (MLM)

Next Sentence Prediction (NSP)

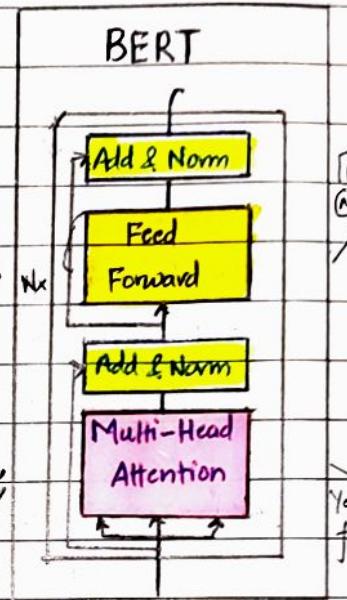
THE [MASK1] brown
for [MASK2] over
the lazy dog →

For MLM BERT takes in a sentence with random words filled

with MASKS. The goal is to output these masked tokens.

It helps BERT understand a bi-directional context within a sentence.

A: Ajay is cool dude
B: He lives in delhi



In the case of Next Sentence Prediction (NSP), BERT takes in

2 sentences and it determines if the second sentence actually follows the first.

In kind of like binary classification problem.

This helps BERT understand content across different sentences themselves, and using both of these together BERT gets a good understanding of language.

Pre-training (Pass 2):

Pre-training (Pass 2): "How to use language for specific task?" (PTO)

The goal of pre-training is to make BERT learn what is language and what is content. BERT learns language by training on 2 unsupervised tasks simultaneously. They are:

Masked Language Model (MLM)

Next Sentence Prediction (NSP)

For MLM BERT takes in a sentence with random words filled with

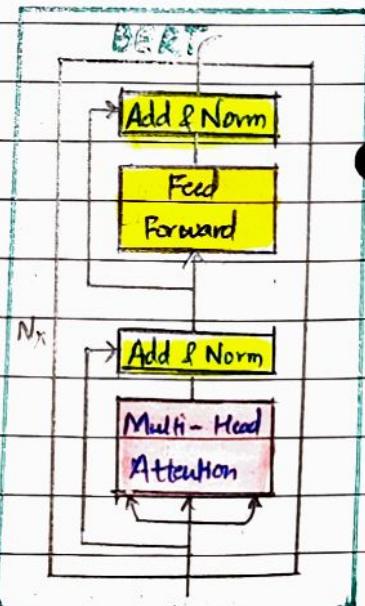
MASKS the goal is to output these MASKED tokens.

The IIP is the set of 2 sentences with some of words being masked each

token is a word. We convert each of these words into embedding using "pre-trained embeddings".

C → binary o/p for next Sentence Prediction

T_n → Word vectors



Teacher's Signature _____

for fig. refers to "BERT
created Pdf research Paper by
Anna Rogers,

Deepak ydwr



Date: _____

Topic: _____

Learned from
Scratch

Fine-Tuning

Fine-tuning (Pass 1): "How to use language for specific task?"

We can further train BERT on very specific NLP tasks for e.g. "Question Answering".

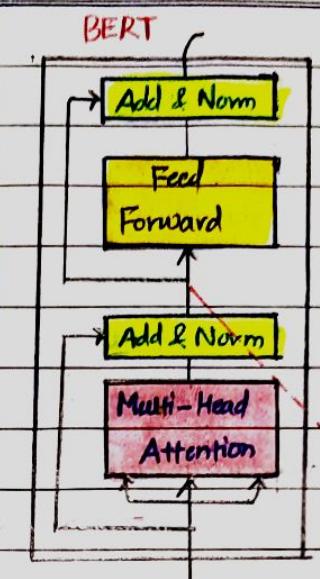
Question Passage

All we need to do is to replace fully connected output layers of the network with fresh set of output layers that can basically output the answer to the question we want.

Then we can perform Supervised using a question answering data set. It won't take long since it's only the output

parameters that are learned from scratch. The rest of the model parameters are just slightly fine-tuned and as a result training time is fast. We can do this for any NLP

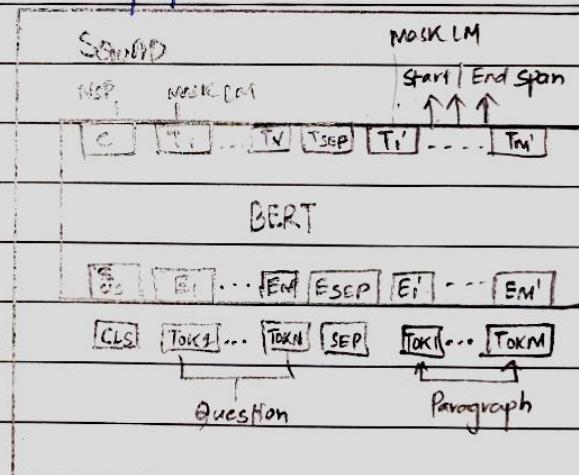
problem. That is replace the old layers and train with the specific dataset.



Fine tuning (Pass 2)

If we wanted to perform Question answering we would train the model by modifying the inputs and the output layer.

We pass in the question followed by a passage containing the answer as inputs and in the output layer we would output these start and end words that encapsulate the answer ensuring that the answer is within the same span of text.





Date: _____

Topic: _____

Pre-training

(Pass-3)

On the left side how we going to generate these embeddings from the word token inputs. Well the initial embedding is constructed from three vectors

Token embeddings are pre-trained embeddings

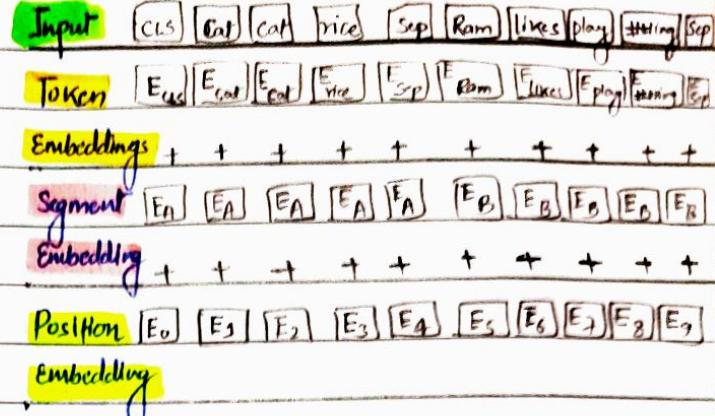
Segment embedding is basically the sentence no. that is encoded into a vector and the

Positional embeddings is the position of the word within that sentence that is encoded into a vector.

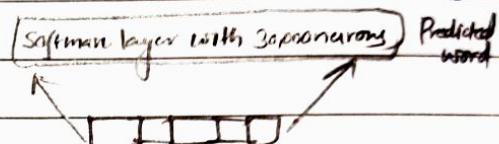
Adding these 3 vectors together we get an embedding vector that we uses as input to BERT

C → O/p of binary and bunch of word vectors, but with training we need to minimize loss.

Two key things to note:

(I) Word vectors T_i have the same size.(II) Word vectors T_i are generated simultaneously.

Segment + Position embeddings add "ordering" for inputs

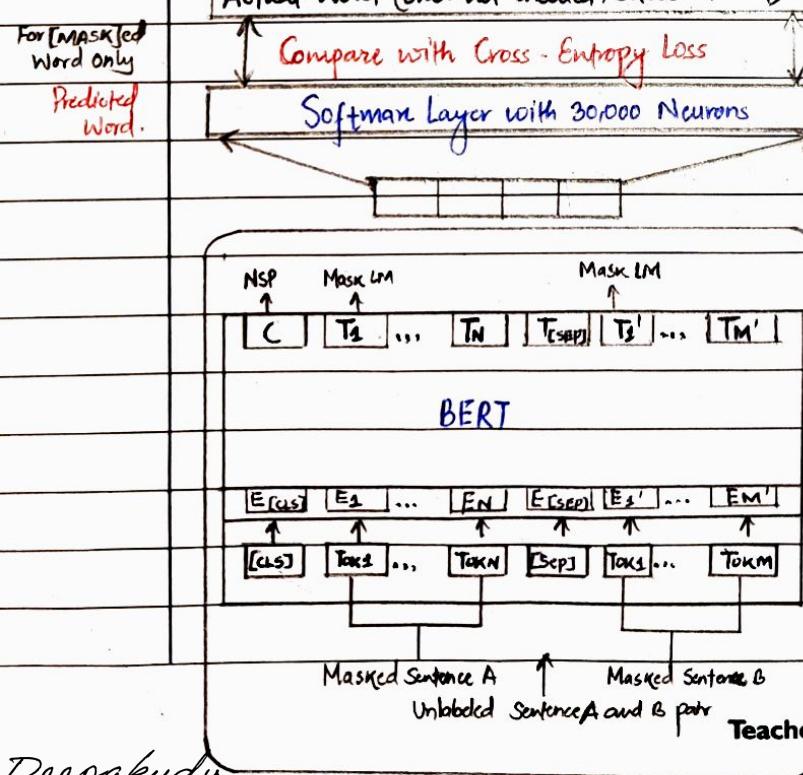


2 We need to take each word vector and pass it into fully connected layer O/p with same no. of neurons equal to the no. of tokens in the vocabulary.

So that would be o/p layer corresponding to 30K neurons in this case. We will apply Softmax activation. This way we would convert word vector to a distribution. The actual label for this distribution would be a one-hot encoding vector for the actual word.

So we compare this two distribution and then pass the network using "Cross-entropy-loss".

But note that the o/p has all the words. Even though those outputs weren't masked at all. The loss though only considers the prediction of the masked words and ignores all the other words that are o/p by the network. This done to make sure more focus is given to predicting the mask value so that it gets them correct and it increases context awareness.



Date: _____

Topic: _____

Key Takeaways:

Pre-training (Summary)

- We pre-train BERT with Mask language modeling and Next Sentence Prediction
- for every word we get the 'Token embedding' from the pre-trained word piece embeddings, add the position and segment embeddings to account for the ordering of the inputs.
- These inputs are then passed into BERT which under hood is a stack of transformer encoders and it outputs a bunch of Word vectors for mask language modeling and a binary value for next sentence prediction.
- The word vectors are then converted into distribution to train using "Cross-entropy-loss".
- One training is complete 'BERT has some notion of language.'

Fine-tuning (Summary)

- In fine-tuning we perform supervised training depending on the task we want to solve and this should happen fast in fact the BERT SQuAD ie Standard Question and Answer model, only takes about 30 mins to fine tune from a language model for a 91% performance.

Performance

Performance depends on how big we want BERT to be. For BERT large model which 340M parameters can achieve way higher accuracy than BERT base model which has only 110 parameters.

nx	BERT Base	24x	BERT Large
110 M Parameters		340 M Parameters	