

RNN

PAGE No.

DATE

* Application

- * RNN works very well with Sequence of Data as your Input.

e.g.: NLP (SPAM, Time goes Day, Cells forecasting) etc.

e.g.

(NLP)

↳ Bow, TF-IDF, word2vec

convert the sentence into vectors

(NLP)

e.g.: Google Assistant, Amazon Alexa, Siri, etc.

e.g.: Time Series

RNN will help to predict [Hence RNN works very well]

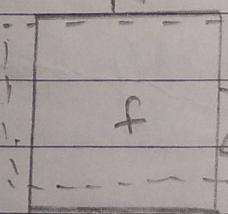
e.g.: Google Image Search

e.g.: Google Translate

e.g.: Image Captioning

* RNN forward Propagation with Time, i.e., $t = 0$ -

O/P \rightarrow $t = 0$.



↳ Hidden layer here we have only no. of hidden neurons $t = 0$ -

we also get O/P with respect to time $(w_c o + w_{u,h})t = 0$ -

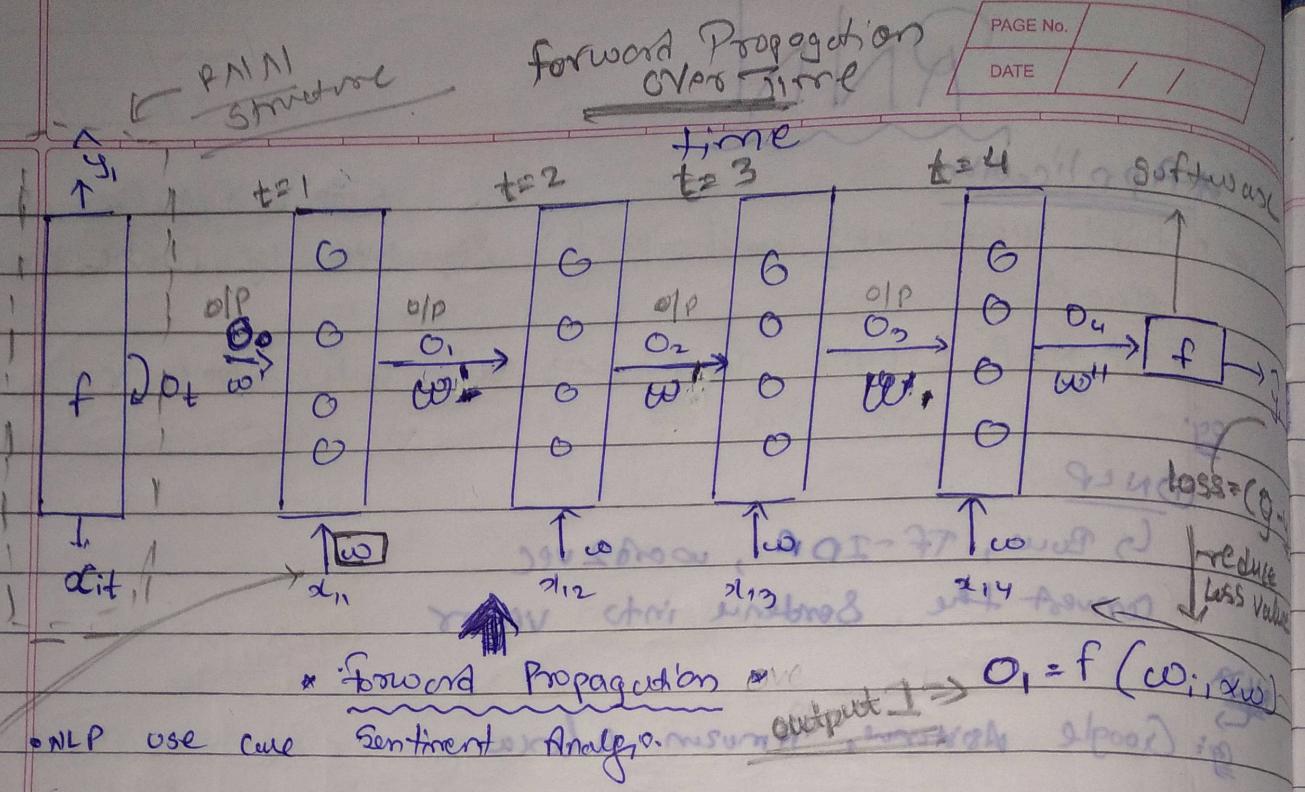
→ x_t input (Any no. of Dimension, i.e. any no. of features)

$$(w_c o + w_{u,h})t = 0$$

• NLP use for Sentiment Analysis

$$\Rightarrow x_t = \langle x_{11}, x_{12}, x_{13}, x_{14} \rangle \quad y; \text{ output}$$

Our RNN act each and every time at $t=1$, will process 1 word, then at $t=2$... process second word...



Input we give to our Hidden Layers $x_i \in \{x_{11}, x_{12}, x_{13}, x_{14}\}$ y_1 global emitting O/P

At time $t=1$ my first word goes inside this particular hidden layer O_{11}

Note one representing the word in Vector. principle name $\rightarrow O_1 = f(x_1 w + O_1 w')$ \leftarrow function of O_1

$\rightarrow O_1 = f(x_1 * w)$ \leftarrow function (sigmoid, tanh, etc.)

$\rightarrow O_2 = f(x_{12} w + O_1 w')$ $\# O_2$ is completely dependent on x_{12} $\#$ O_1

$\rightarrow O_3 = f(x_{13} w + O_2 w')$

$\rightarrow O_4 = f(x_{14} w + O_3 w')$

Note Sequence information is always maintained, it is not discarded.

* Back Propagation (weights gets updated continuously w.r.t time)

$$\rightarrow \frac{\partial L}{\partial y}, \frac{\partial L}{\partial w''} = \frac{\partial L}{\partial y} \cdot \frac{\partial y}{\partial w''} \dots \text{(chain rule)}$$

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

bias or forget weight is not random weight

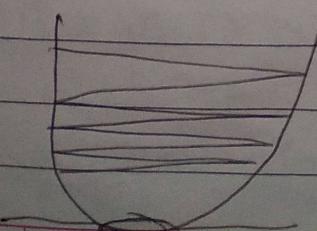
$$\rightarrow \frac{\partial L}{\partial y_i}, \frac{\partial L}{\partial w} = \frac{\partial L}{\partial y_i} \cdot \frac{\partial y_i}{\partial y} \cdot \frac{\partial y}{\partial w} \dots \text{(chain rule)}$$

* Problem in Simple Recurrent Neural Network

Von Neumann
problem

The derivative of Sigmoid is between 0 to 1, now suppose one here we have calculated derivative function it will be ~~small~~ between 0 to 1, but as we go towards this w , this value becomes very small value itself as the derivative is continuously happening, derivative of cell ~~has~~ these particular weights this will become a very very small value (w), when w becomes a very very small value whenever we do an update of this particular weights, it will be negligible, then we ~~will never~~ converge to global minima.

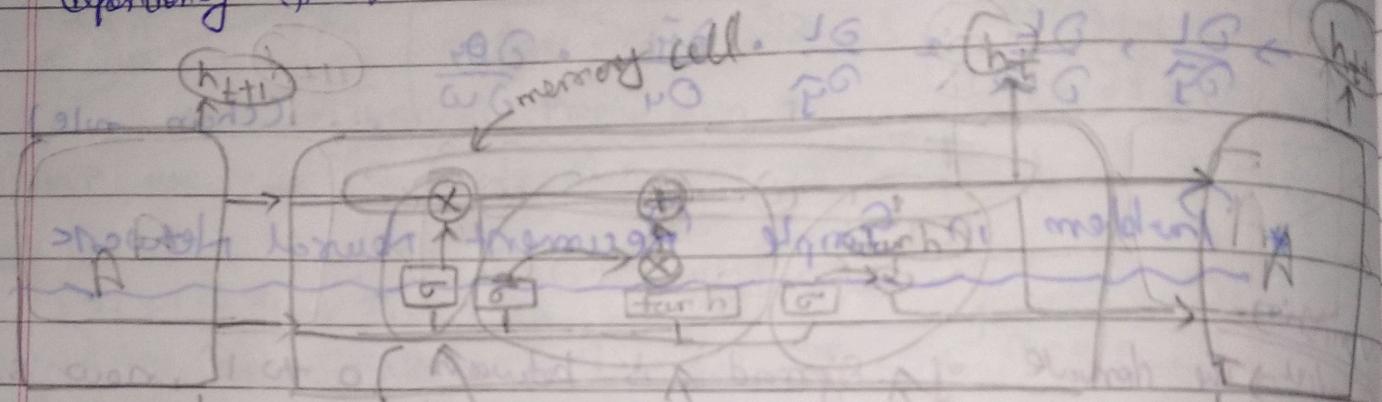
Exponential) If we are using other than sigmoid like ReLU, $1 > 0$ (suppose my derivative is ~~greater than 1~~), so we go ahead towards right to left, it creates exploding gradient problem, when there's ~~no~~ changes will happen so big that it will never reach the global minimum point.



→ To solve this problem we use LSTM RNN
(Long short term memory recurrent neural network)

* LSTM

- LSTM core explicitly designed to avoid the long-term dependency problem.



Different Components

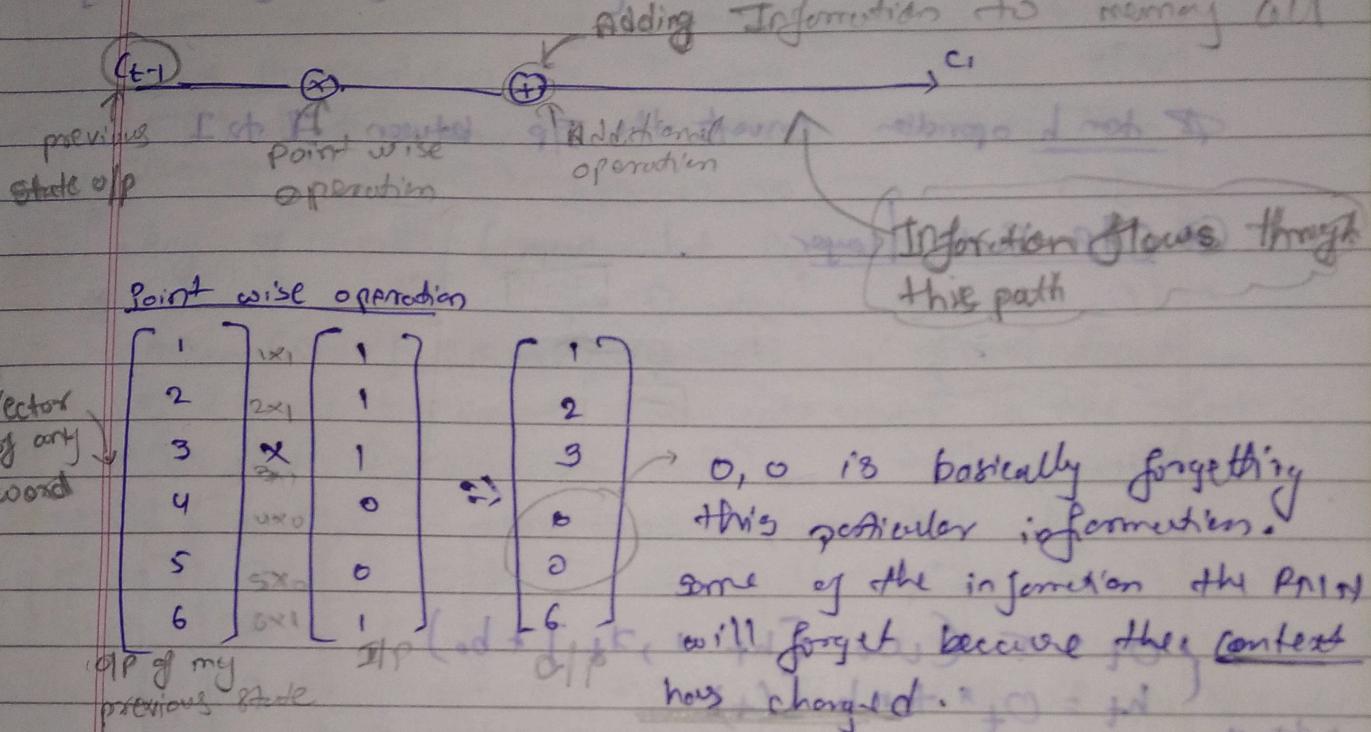
- ① Memory cell
- ② Forget gate
- ③ Input gate
- ④ Output gate

* Notation

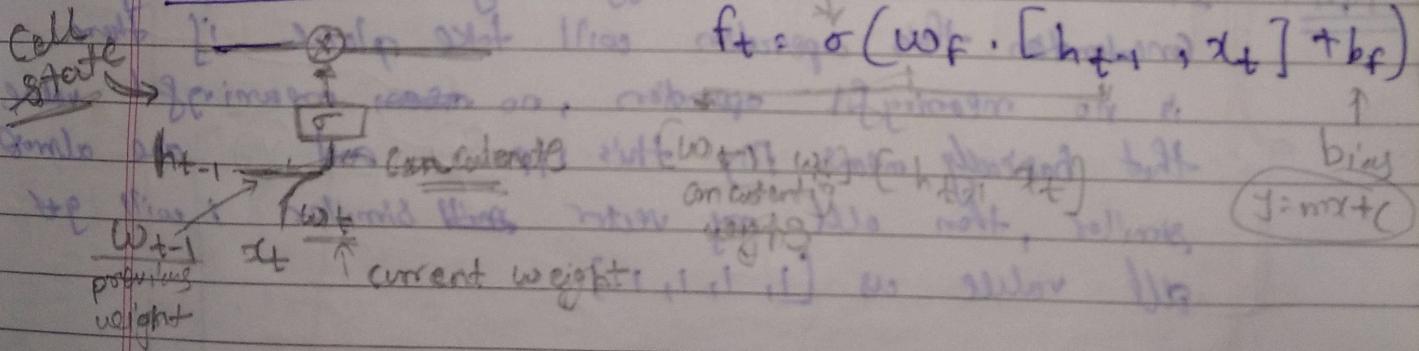
- $\boxed{\quad}$ = Neural Network layer
- \circ = Pointwise operation
- \rightarrow = vector Transfer
- $\nearrow \searrow$ = concatenate
- $\overleftarrow{\overrightarrow{\quad}}$ = copy

① Memory cell

- Memory cell is used to remember and forget things, how do
- How No we remember and forget is based on the context of the input



② Forget Gate



* σ transforms the I/p between 0 to 1

→ Forget Gate → It either lets something in or doesn't let it in, which means it helps us the forget gate, it helps in defining what information to forget what information to remember

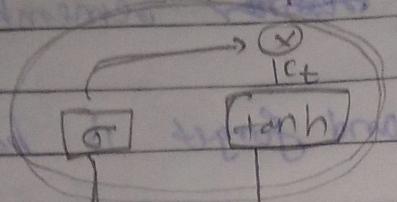
3) IP Gate Layer

Adding Information
to memory cell

PAGE No.

DATE

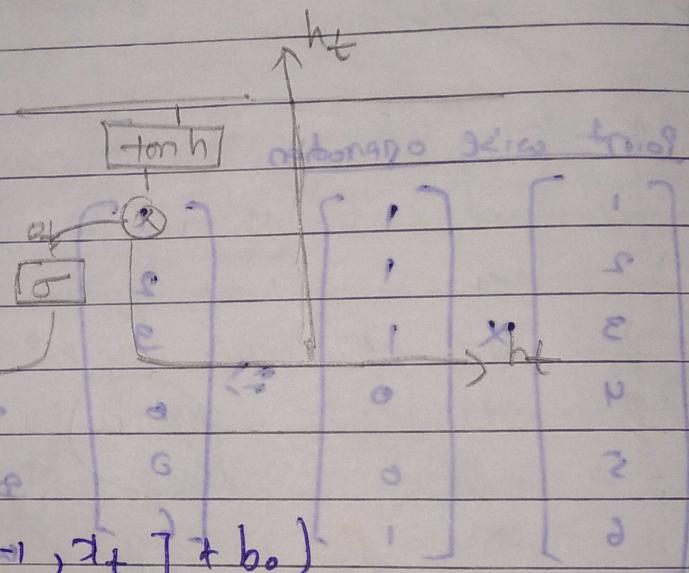
10



$$\begin{cases} i_t = \sigma(w_i \cdot [h_{t-1}, x_t] + b_i) \\ C_t = \tanh(w_C \cdot [h_{t-1}, x_t] + b_C) \end{cases}$$

~~IP~~ tanh operation converts the IP between \mathbb{E} to \mathbb{I}

④ OP Gate Layer



$$\begin{cases} O_t = \sigma(w_O \cdot [h_{t-1}, x_t] + b_O) \\ h_t = O_t \cdot \tanh(C_t) \end{cases}$$

Note:- If there is New meaningful information then only this concatenation this operation will take place if there is no meaningful operation, no ~~mean~~ meaningful value that basically means my this $i|p$ and x are almost similar, then all my vector ~~will~~ similar; will get all value as [1, 1, 1, 1]

$i|p$ & x are almost similar

Word Embedding (Feature Representation)

Two technique in word Embedding

(i) word2vec

(ii) Glove

word representation

$|V| = 10,000$ (10,000 words), Dictionary having 10,000 words

size (Dict) • Bag of words (one hot representation)
 eg Man [5000] among presentation in 5000 location
 If i want to convert man into vector in the form of

(One hot representation), (0001) (0002) (0003)

$\begin{bmatrix} 0 \\ 0 \\ 0 \\ \vdots \\ 1 \\ 0 \end{bmatrix}$ index of 5000. It is difficult to execute
 at index 5000 the value will be 1

woman [9000]

$\begin{bmatrix} 0 \\ 0 \\ \vdots \\ 1 \\ 0 \\ 0 \end{bmatrix}$ index of 9000. It is difficult to execute
 at index 9000 the value will be 1

→ This is very sparse matrix, High Dimension.

→ Sparse matrix \Rightarrow is basically means that you have so many number of 0s and you have less number of 1s

~~Cosine Similarity :- used in recommendation system~~

heavily used

PAGE No.	/ /
DATE	/ /

disadvantage :- when my ML or DL algorithm is getting applied it is very difficult to generalize these vector, we cannot understand like if i want to find out the similar word, it is very very difficult to understand from this, because some of the other index in every word you know that only one vector will be 1 and remaining all vector will be 0, Hence there is no semantic information of the words.

- To overcome this problem we use Word Embedding (Feature Representation)

	(2000)	(5000)	(6000)	(9000)	(1000)	(4000)
Gender	-1	0	-0.92	0.93	0.0	0.1
Royal	0.01	0.02	0.75	0.96	-0.02	0.01
age	0.03	0.02	0.7	0.6	0.95	0.92
food	:	:	:	:	0	:
:	:	:	:	:	0	:

→ Converting vector based on some features (Gender, Royal, age) [Feature 300×10000]

~~Lower Dimension and Dense matrix.~~

→ Cosine is used to find similarity between two vectors

Analogy (Based on Gender)

Boy $\xrightarrow{\text{represent}}$ Girl

$$\begin{aligned} \mathbf{g}_1 &= \begin{bmatrix} -1 \\ 0 \\ 0 \\ 0 \end{bmatrix} & \mathbf{g}_2 &= \begin{bmatrix} 1 \\ 0 \\ 0 \\ 0 \end{bmatrix} & \leftarrow \mathbf{g}_1 \text{ and } \mathbf{g}_2 \text{ are vectors} \\ \mathbf{g}_1 - \mathbf{g}_2 &= \begin{bmatrix} -2 \\ 0 \\ 0 \\ 0 \end{bmatrix} & (-1 - (+1)) = -2 & \text{of Boy and Girl} \end{aligned}$$

↳ cosine similarity

King $\xrightarrow{\text{represent}}$ Queen

$$\begin{aligned} \mathbf{k} &= \begin{bmatrix} \approx -2 \\ 0 \\ 0 \\ 0 \end{bmatrix} & (\mathbf{q} - \mathbf{k}) &= \begin{bmatrix} \approx -0.92 \\ 0.93 \\ 0 \\ 0 \end{bmatrix} & \approx -2 \end{aligned}$$

King + Queen

Apple

(Grouping) based
on Similarity

Ribbedmeat

• Boy • Man

• Girl
• women

• meat

* word Embedding is used to find out the most similar vectors because through the vector is operated in that particular way only

$$|V| = 10000$$



Ans Similarity

Vocabulary size tells us size of the dict
dimension = features [Gender, Royal, age...]

* Bidirectional RNN

$$\|v\| = 10,000$$

Boy is good

0
0
0
0
0
1 ← 2000
0
0
0
0

Boy → 2000 ← index

at index

Embedding Layer \Rightarrow Converts some kind of vector representation

using keras

converted Boy into one-hot. After applying one-hot [2000, 4000, 5500] we will get 1000 values of these values. we will pass these values with Embedding layer.

Embedding layer

1000 dimensions \rightarrow 10

- In Keras we have function called One-hot, which helps us to convert this word(Boy) into one-hot representation.

- Embedding layer is basically creating a feature matrix. Boy

- for 4000 index another 10 vectors of number will come

Boy

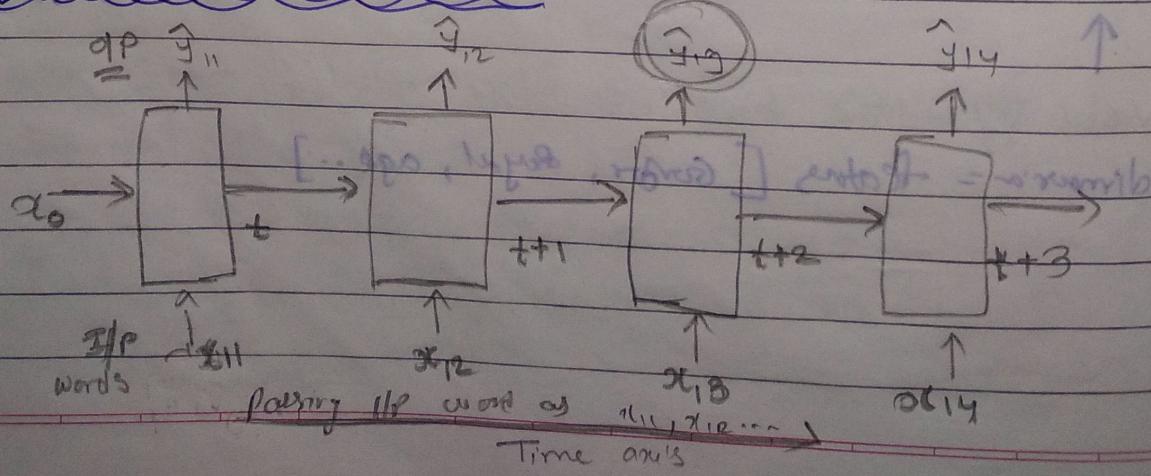
-1
0.001
0.002
0.003

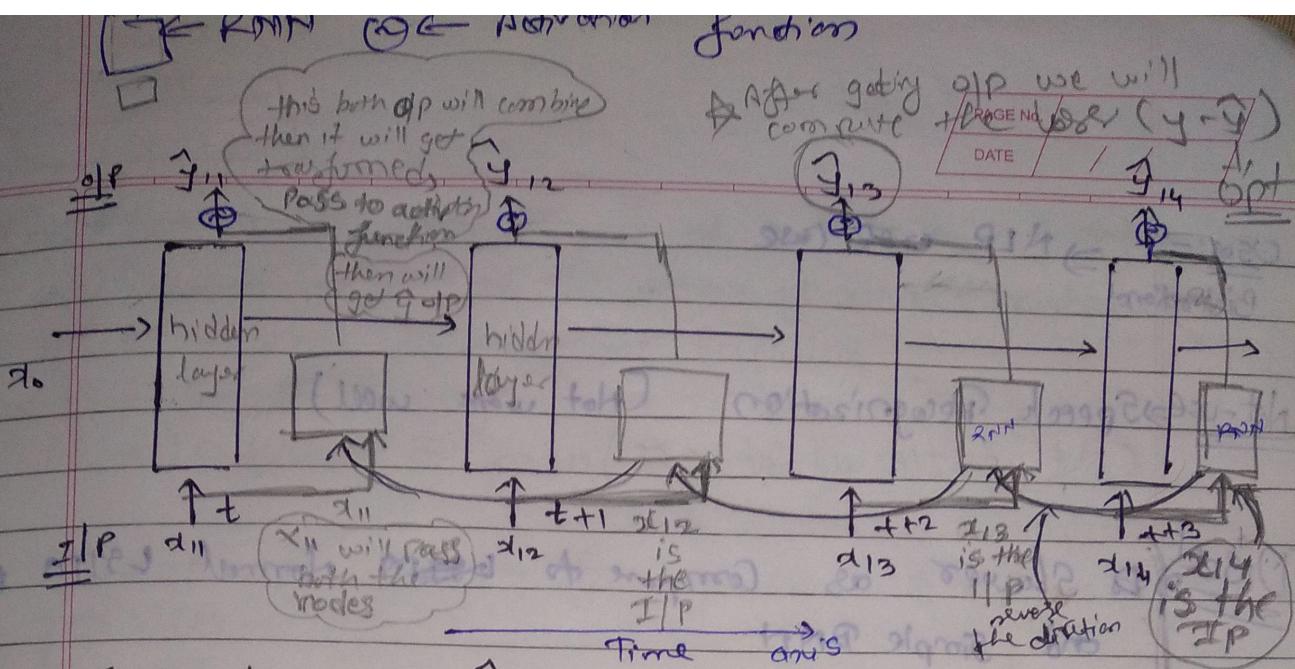
With two half of been in middle now after forward second hidden unit from your reading text it

* Bidirectional RNN

Explanation

$$0000 = \|v\|$$





Explanation \rightarrow To compute $o/p \Rightarrow \hat{y}_{13}$ we need to have information of x_{13}, x_{12} and x_{11} . This information is getting passed along with the $o/p \hat{y}_{11}, \hat{y}_{12}, \dots$

Note \rightarrow But, what if \hat{y}_{13} this particular o/p is dependent on future words.

e.g. He likes eating $\frac{?}{}$ (\Rightarrow can be apple, banana...)

RNN Information passes in Unidirection, it cannot get the information of future words.

- To overcome this \uparrow problem we use Bidirectional.

\rightarrow We will be using Backward layer PNN \square and it will just change the direction of the information

\rightarrow Now, we will just reverse the direction, \Rightarrow

Imp \rightarrow Only thing we are doing is that we are adding the extra LSTM Layer and information is passing to the reverse direction.

Now:- \hat{y}_{13} is able to get a previous words and the upcoming words

Used in NLP task

Uses:- \Rightarrow NLP Use Case
Bidirectional

Not-use \Rightarrow Speech Recognition (Not work well)
(Not get I/p well at ones)

Disadv Slower as compare to ~~normal~~ Normal LSTM RNN,
and Simple RNN

Because we are taking exact replica of some layer
and we are passing information in the bidirectional.

Imp

Note:- In the longer sentence, i will get the information of the
future words

↳ Availability you can和睦 \uparrow edit tomorrow at 8

↳ How is it very helpful because you can edit tomorrow \uparrow

↳ Availability you can edit tomorrow \uparrow

↳ Availability you can edit tomorrow \uparrow

↳ Availability you can edit tomorrow \uparrow

Poem RNN, LSTM, GRU

O/P not consider in Encoder
we are more focus on getting the context vector

I am taking up after the last RNN layer

PAGE No.

DATE

then, context vector will pass to the Decoder

* Sequence To Sequence Learning with NN

Encoder And Decoder

A B C D
are words

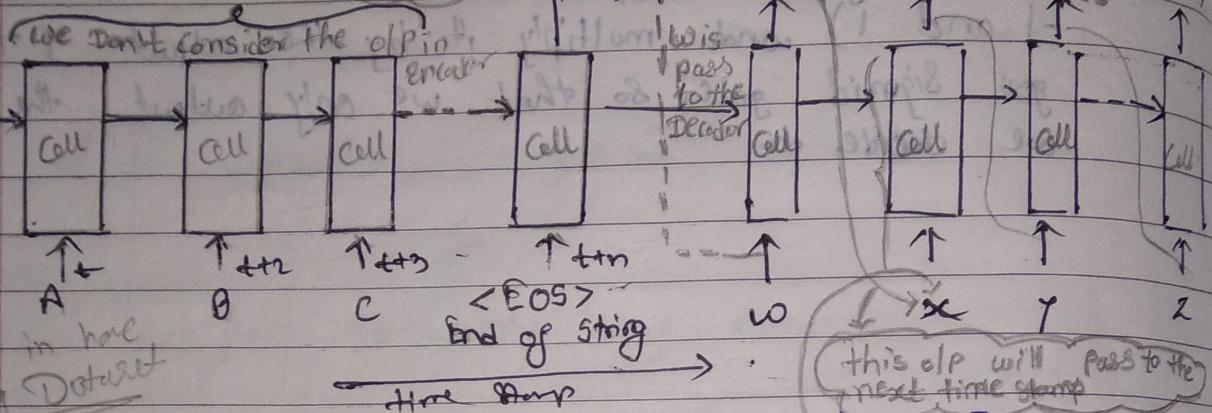
Encoder

Generate vectors

vector

Decoder

O/P word in sequence end of sentence



e.g. Language translation [English \rightarrow French] this is sequence to sequence learning.

(ad + [sd, sd] now is y so)

eg

eg:- Image captioning

eg:- Google search [whenever we are searching Dog image we will get images on Dog]

eg:- LinkedIn [we get some automated reply in the chatbox]

* Two main Component in Sequence to Sequence

① Encoder

② Decoder

Encoder \Rightarrow Encoder it can be any NN like RNN, LSTM, GRU
But usually we mostly select LSTM because it works well

In Encoder the o/p will not be present the o/p will present in the last time step $t+n$

Encoder and Decoder \Rightarrow are in unidirectional

One book $\frac{w}{w+1} \begin{matrix} w & 0 & r & ds \\ 0 & 0 & 0 & 00 \end{matrix}$ ~~not~~ (odd can be any NN) PAGE No.

PAGE No.	
DATE	/ /

eg • In language translation we pass text, each word in every text is passed in the vectors

Vectors like \Rightarrow { Embedding layers
Onhot Representation
Word2vec }

process. Convert the word to vector then pass it to ANN

$x = \langle x_1, x_2, x_3 \rangle \dots$ It can be any no. of i/p we consider as A, B, C

$y = \langle y_1, y_2, y_3 \rangle \dots$ Objekt,

* for prediction $\hat{Y}_1, \hat{Y}_2, \hat{Y}_3$ (Training Data-Set)
LSTM works on the probability

LSTM works on the probability

$$\Pr(\vec{y} | (x_1, x_2, x_3, \dots, x_n))$$

Sequence of α/ρ will be taken up and we will find the mutability of $\beta_1 \beta_2 \dots \beta_n$ and their relation to each other.

\hat{y} is the combination of $\hat{y} < \hat{g}_1, \hat{g}_2, \hat{g}_3 >$

After getting \hat{y} , then i will try to find out the loss (y, \hat{y}), for this i will use Optimizers. Optimizers will reduce the loss function, through the

Buck propagation, we gets updated.

Disadvantages \Rightarrow As the sentence length increase w.r.t the Encoder and Decoder, the accuracy decreases. (Sentence length e.g = 100 words)

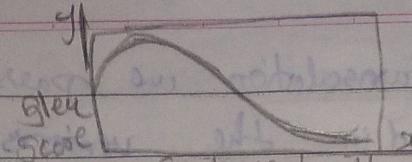
To avoid this problem we will be using Attention

Cell = NN

PAGE No.

DATE

Bleu Score

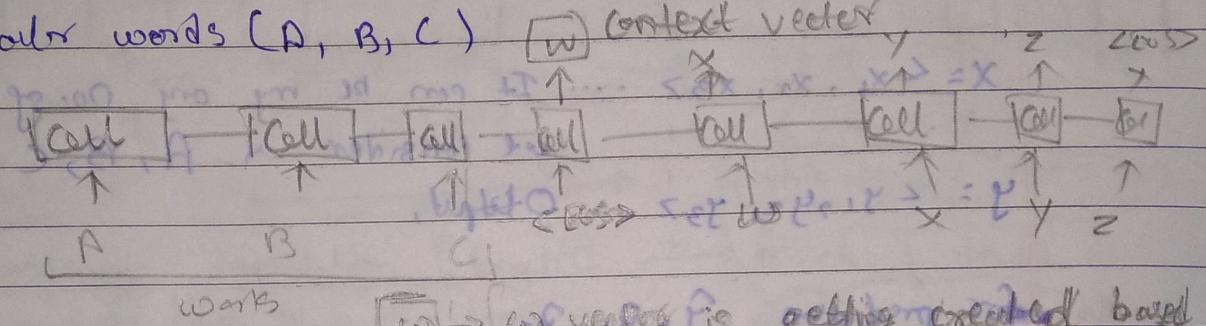


Imp

Dipadu
Note:-

As the sentence length increasing, the BLEU Score (decoding) coming down, hence this was the problem researchers found

whenever the information is long this vector w generated is not able to capture all the information of these particular words (A, B, C) w context vector



$w \rightarrow$ w vector is getting created based on the input sentence, help words that i have given

* Attention Model

Initially we will just take the window size words in our sentences may be 3 or 4 depends on hyperparameter that we select.

If we select 4 to 5 words in the sentence then i pass that particular information to the decoder, the decoder will convert it. It goes on like this. Hence, to solve the problem of Encoder and Decoder this is the solution. It is called as attention model.

Dipadu
Decoder

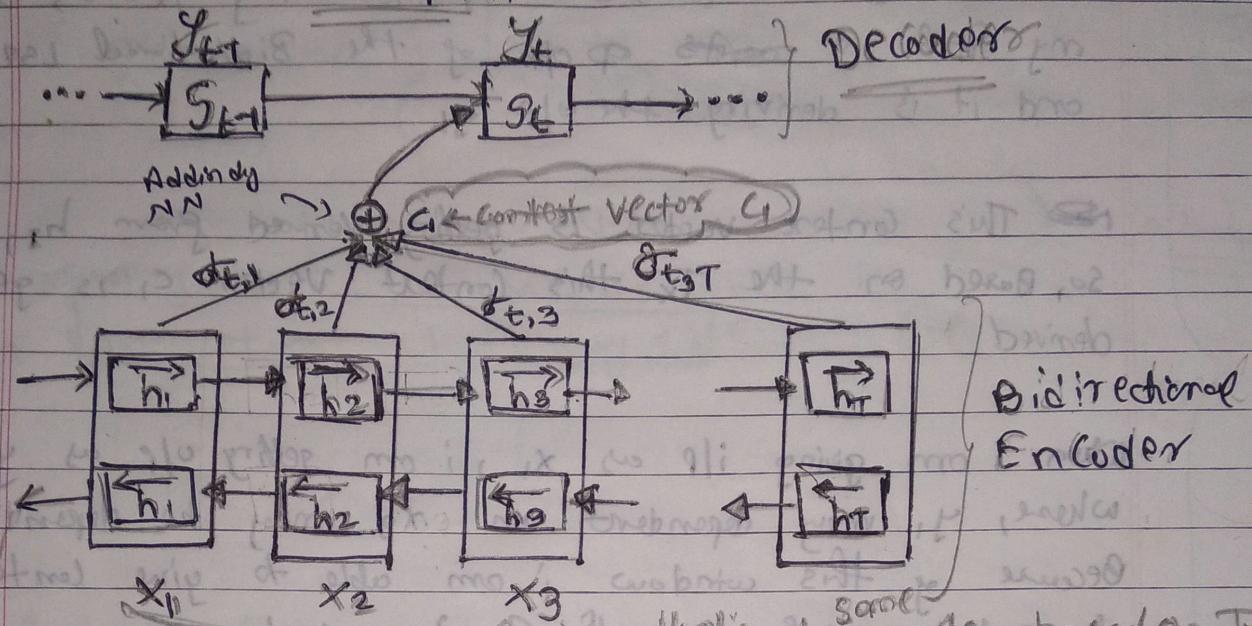
Encoder and Decoder usually do not perform well longer sentences

Researchers experienced and found out lower Bleu Score

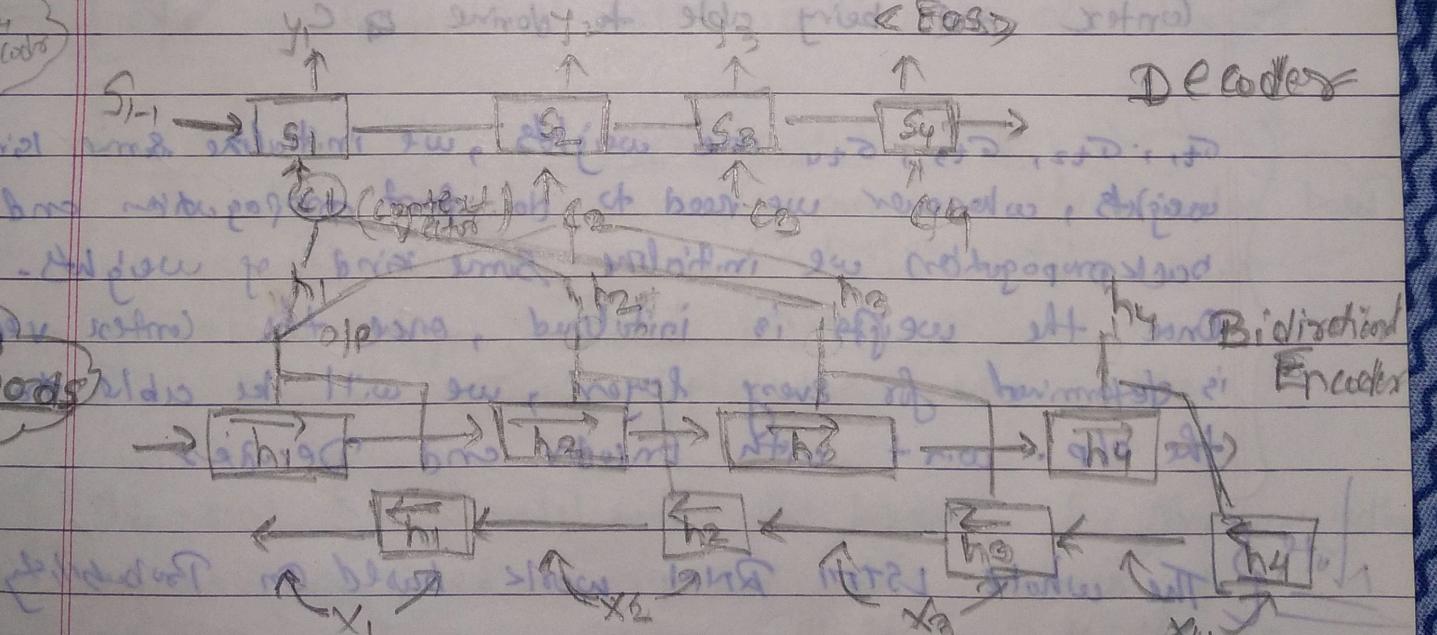
* Understanding Attention Models and Architecture

Refer \Rightarrow RNN, LSTM, Bidirectional LSTM, Sequence 2 seq

LSTM RNN



In Encoder we will be using Bidirectional RNN
In Decoder we will be using LSTM RNN



We have to consider this window, and this window is represented as T_x , here $T_x = 3$

This window can be decided as a hyperparameter. Like initially how many words should I give, then the algorithm will able to translate

~~→ Context vector are getting derived from h_1, h_2, h_3~~
 (Taking window size = 3 of T_x)

Now this h_1 will combine this, since my $T_x = 3$, so
 my first 3 steps of off of the Bidirectional LSTM
 and it is deriving the o/p y_1

Explanation ~~(This Context vector c_1 is getting derived from h_1, h_2, h_3~~
 so, Based on the T_x this Context Vector c_1 is getting
 derived)

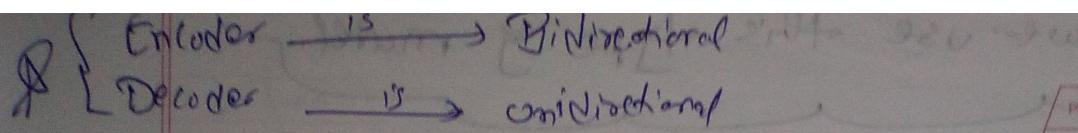
When I am giving i/p as x_1 , I am getting o/p as y_1 ,
 where, y_1 may dependent on x_2 may be dependent on x_3 .
 Because of this window, I am able to give context &
 context information through h_1 and h_3 .

Even this y_1 is dependent on x_2 and x_3 , it's getting
 context and being able to derive c_1

$G_{t1}, G_{t2}, G_{t3}, \sigma_{tw}$ are weights, we initialize some kind of
 weights, whenever we need to do front propagation and
 backpropagation we initialize some kind of weights.
 Once the weights is initialized, once the context vector
 is determined for every layer, we will be able to see
 the o/p. work with Encoder and Decoder

Note → The whole LSTM RNN work based on Probability

- Applications
- Language Translation
 - Text Summarization
 - Q, A



PAGE NO.	
DATE	1/1

Maths

• Condition 1

when $\Rightarrow T = 3$

$$\begin{matrix} \alpha_1, \alpha_2, \alpha_3 \\ \downarrow \\ \alpha_1 \end{matrix}$$

{ window size $\# 3$, it will be getting 3 α values
for the first word x_1 , and this
will be passing to my decoder
and it will create an o/p y_1

• sum of the α_1, α_2 and α_3 will be,

Imp $\sum \alpha_i = 1$ (Always 1) $\quad \textcircled{1}$

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

• Condition 2

Imp $C_f \rightarrow$ Context vector, context vector bio generation,
now it is getting generated, it is just multiplying,
o/p of Bidirectional LSTM with the α values

? Label to sent out word $\xrightarrow{\text{forward}}$
 $\sum_{i=1}^{Tx} \alpha_i h_i$

we can change the α value based
on back propagation

• This is the property of feed forward NN
It is basically multiply, I/P features and weights

• The feed forward NN is basically for Backpropagation
Because, we need to update the values of
 α_1, α_2 and α_3 . If we will not update the α we
will not get good accuracy of getting words

e.g. Some like RNN (feed forward NN)

Q Why we use this type of model.

PAGE NO.	1/1
DATE	/ /

Notes:

Softmax helps to make α value $\alpha = 1$ always
for any very window size, it will make $\alpha = 1$ always

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

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

↑
feed forward NN

BLEU
Score

Even the sentences are long
we are good accuracy

Sentence length

Eg! In Google Translator, even for longer sentences, it is able to give good accuracy. Because of the TimeStamp probably, timestamp acting as the parameter the sentences are changing.

Q Why we use this type of Model?

\Rightarrow In case of Encoder and Decoder, when the researcher were tried with different - different sentence length, does not perform well with the longer sentences.

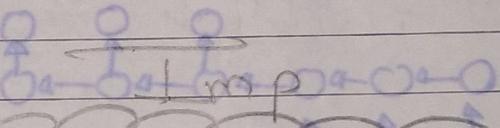
Because of this they have included Bidirectional LSTM encoder, they have made some architecture changes between basically in the center they have use Attention model which is nothing but feed forward NN w.r.t. some Time stamp which is basically indicating ~~the~~ window size.

Encoder, Decoder, Attention Models, Transformers, Bert - I

⊕ CNN always try to train in a way to understand a relationship of all static data, but it will not be able to memorize, if it will not be able to understand a context, if it is not able to understand a context, eg if i am talking about a Question system you will not be able to build it. even like a sentimental system you will not be able to build it. Language translation you will not be able to build it, text summarization you will not be able to do it.

⊕ To overcome this problem (CNN) \rightarrow introduce LSTM

Seq to Seq
 Vec to Seq
 Seq to Vec
 Vec to Vec



Note:- $\text{O} \leftarrow \text{I}$ can be any type of RNN cell

① O_{dpp} used in = Image classification
 O_{h} Bounding box regression
 Cell can be RNN, LSTM, GRU

$12\text{IP} \rightarrow 10\text{P}$
 \uparrow
 $f(x)$ regression (approx)

O_{dpp}

vector \rightarrow vector

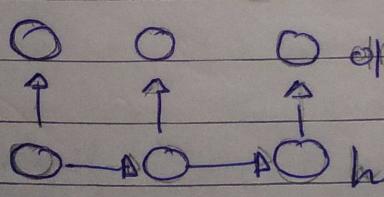
$x \rightarrow \hat{y}$... we go from vector x to another vector \hat{y}

②

$[+]\hat{y} \leftarrow [+]x$

which then produce
Seq of dpp

sequence of
internal state



used in = Image captioning

$x \rightarrow \hat{y}[t]$

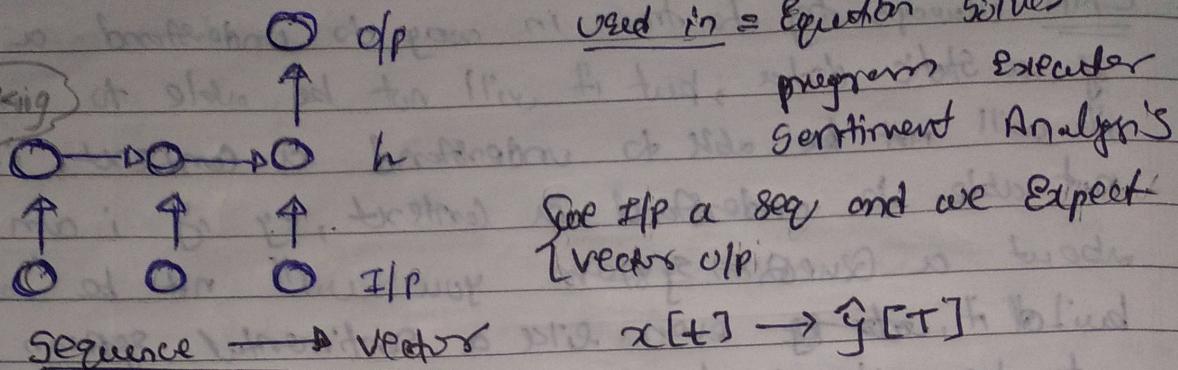
eg: feed back

O_{IP} (single IP)

vector \rightarrow Sequence

③

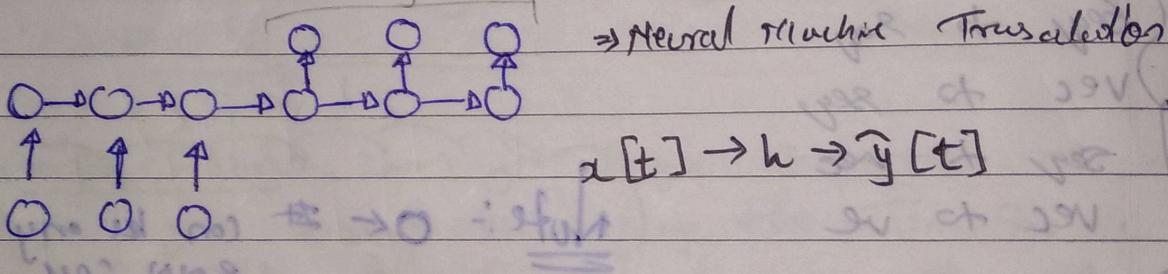
How we looking
only 10 IP



Sequence of I/P feed into the Hidden layer then we produce an O/p Only at the End of the

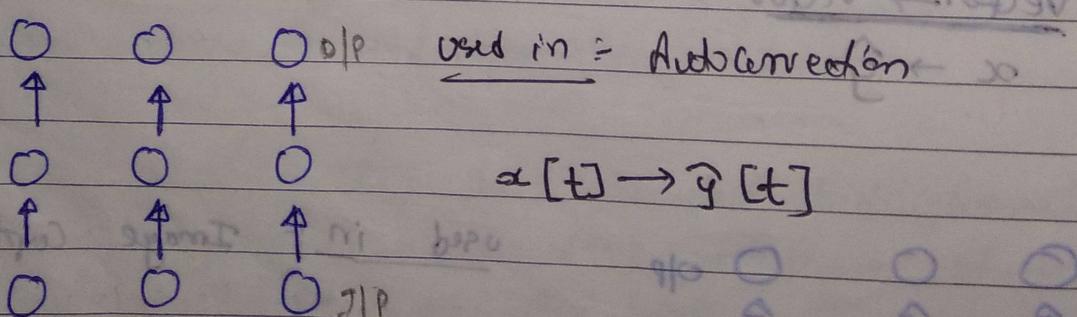
④

it can be only length

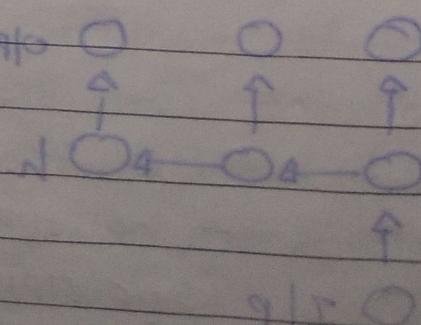


Seq to Seq, but we also come through a Condensed representation (which is actually a vector, so in this case we can call this scheme also encoder-decoder)

⑤



* Seq to Seq



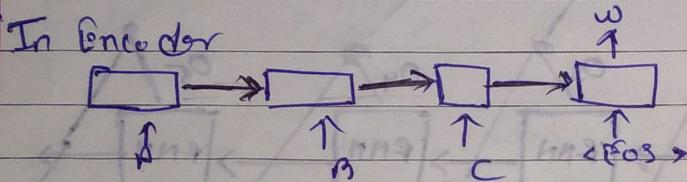
Issue in Encoder

PAGE No.	/ /
DATE	/ /

* Separated I/P and O/P network named as

- I/P Network named as \rightarrow ENCODER
- O/P Network named as \rightarrow Decoder

Context vector $\{ \begin{matrix} \text{Context Vector} \\ \text{w} \end{matrix} \}$ \Rightarrow we are trying to give multiple I/P (e.g. I am fine) it is trying to understand, it is trying to summarise each and every word and then it is trying to send this (w) input to a different network. Giving I/P from the previous network \rightarrow Encoder \rightarrow Decoder (O/P Network)



- ISSUE in Encoder Note
- In Encoder we are trying to propagate our data only in one direction. The word that you are taking before and the word that you are talking after that matters a lot.
 - But in Encoder the input data is propagating only in one direction.

$$h_t = \text{sigm} (W^{hx}x_t + W^{hh}h_{t-1})$$

It is a kind of sigmoidal function that you expect to get current output y_t from previous feedback from particular cell

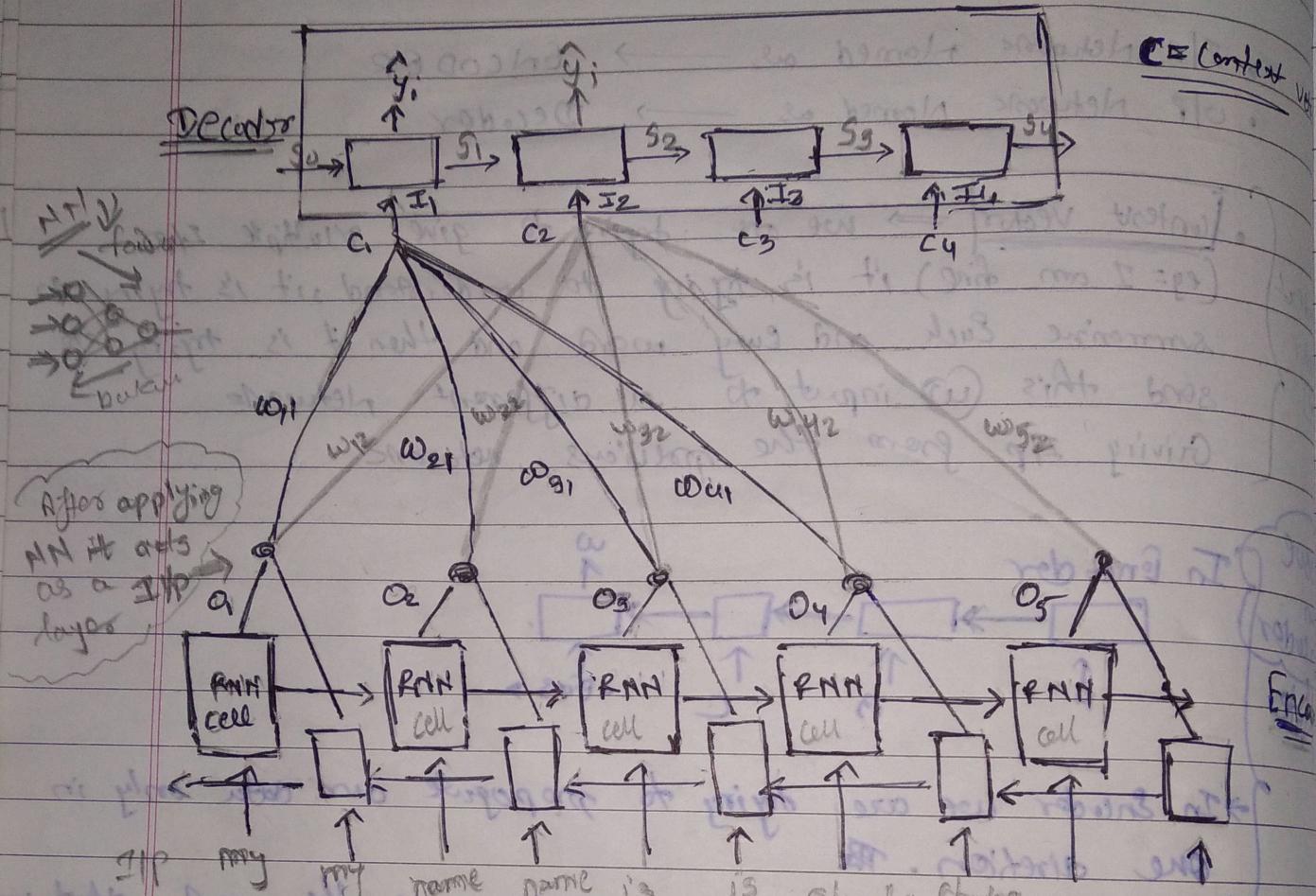
$$y_t = W^{yb}h_t + b$$

- * \rightarrow Internal hidden layer calculation

$$P(y_1, \dots, y_T | x_1, \dots, x_T) = \prod_{t=1}^T P(y_t | y_1, \dots, y_{t-1})$$

refer NLP Book

Attention Models, Transformers - Part 2



Note
In between we have introduced NN to give attention for some of RNN words.

Eg:- Bull is gaining high attention

$O_1, O_2, O_3, O_4, O_5 \Rightarrow$ O/P of Bi-directional

$w_{11}, w_{12}, w_{21}, w_{31}, w_{41} \Rightarrow$ weights (weights will be random)

Implementation
After introducing NN in between the Encoder and Decoder
The O/P of Encoder (Bi-directional) will act as IIP.
Then after apply weights (i.e., $w_{11}, w_{12}, w_{21}, w_{31}, w_{41}$)
then it is given to IIP of Decoder, after that get O/P predicted (g_1) O/P.

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

where,

C_i ← Context vector I/P which will get inside a decoder

$C_i \in I_i$, T_x = hyperparameter (no. of I/P into my NN)

$h_j \in O_i$; (output) combined output of RNN network

α_{ij} = weights (w_{ij}) (learnable parameter)

$$\alpha_{ij} = w_{ij}$$

Content Vector

$$\Rightarrow C_i = \sum_{j=1}^{T_x} \alpha_{ij} w_{ij} + \dots + \alpha_{T_x} w_{T_x}$$

Exponent eg: $\frac{w_1}{w_2} \cdot \frac{w_1}{w_2} \cdot \dots \cdot \frac{w_1}{w_5}$

$$C = \alpha_1 w_1 + \alpha_2 w_2 + \alpha_3 w_3 + \dots + \alpha_T w_T$$

Calculating weight

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

Exponent

Calculating weights
 $i_k \rightarrow h_1, h_2, h_3, h_4, h_5$

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

Attention function

s_i = feed back

s_{i-1} = the o/p that we are getting from previous one, current
 h_j = o/p from the Encoder

Condition $\sum \alpha_{ij} = 1$ (always +ve) — (1)

$$\begin{array}{r} 3000 \\ \times 15 \\ \hline 15000 \end{array}$$

s_i I/P (or)
 α o/p of decoder

$\alpha_{ij} \geq 0$ (α_{ij} will be always +ve)

Application

- (i) Language Translation
- (ii) Q A System
- (iii) Speech Recognition
- (iv) Image Captioning
- (v) Text summarization