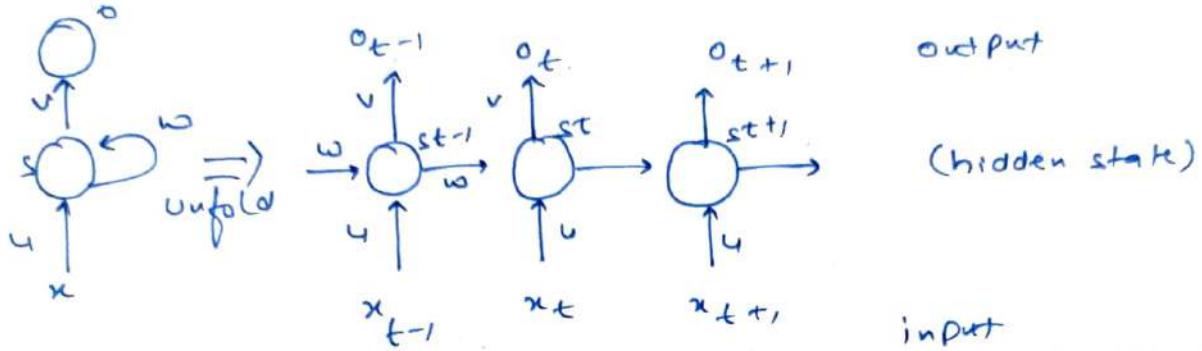


Recurrent neural network :- Difference between feed forward ①
neural network and Recurrent neural network

- RNN :- same task for every element of a sequence with the output being dependent on previous computation
- RNN have a memory which stores information about what has been calculated so far.

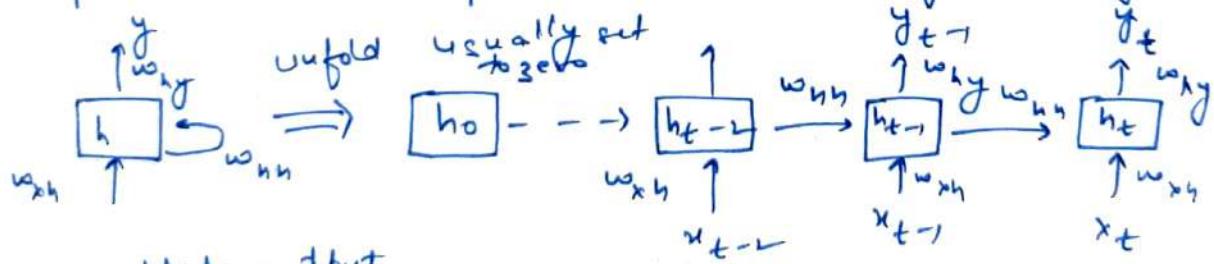


u, v, w :- weights, s, s_{t+1}, s_{t-1}, s_t :- hidden state

same weights used for every hidden layer :: it performs the same task on all the input of hidden layer to produce the output

No. of layer in RNN = no. of words in a sentence sequence
- output from previous step are fed as input to the current step.

- RNN used for sequential learning problem
- RNN save the output of a layer & feed this back to the input in order to predict the output of the layer



- multiple output

- single output

- there is a parameter sharing for

w_{hh} :- shows relation b/w x_t & h_{t-1}

relationship b/w x_t & x_{t-1}

w_{xh} :- identical feature extraction from inputs

Parameter in RNN $x_i, h_{i-1}, w_{xh}, w_{hh}, w_{hy}, b, y$

- Tasks for RNN :- 1. for each timestep of the input sequence & predict output y synchronously.
2. for the input sequence and predict the scalar value of y (like end of the sequence)

3. for the input sequence x of length T_x generate output sequence y of different length T_y
- In RNN outputs are calculated not only by weights applied on input like a regular NN, but also by a "hidden state" vector representation the context based on prior, input, input / output
 - In RNN there are multiple copies of same network with same function and same parameters
 - New / current hidden state h_t can be calculated as

$$h_t = f(h_{t-1}, x_t)$$

$$\text{or } h_t = f(w_{hh} h_{t-1} + w_{xh} x_t + b)$$

$$h_t = \tanh(w_{hh} h_{t-1} + w_{xh} \cdot x_t + b)$$

output vector at time step t calculated as

$$y_t = w_{hy} h_t$$

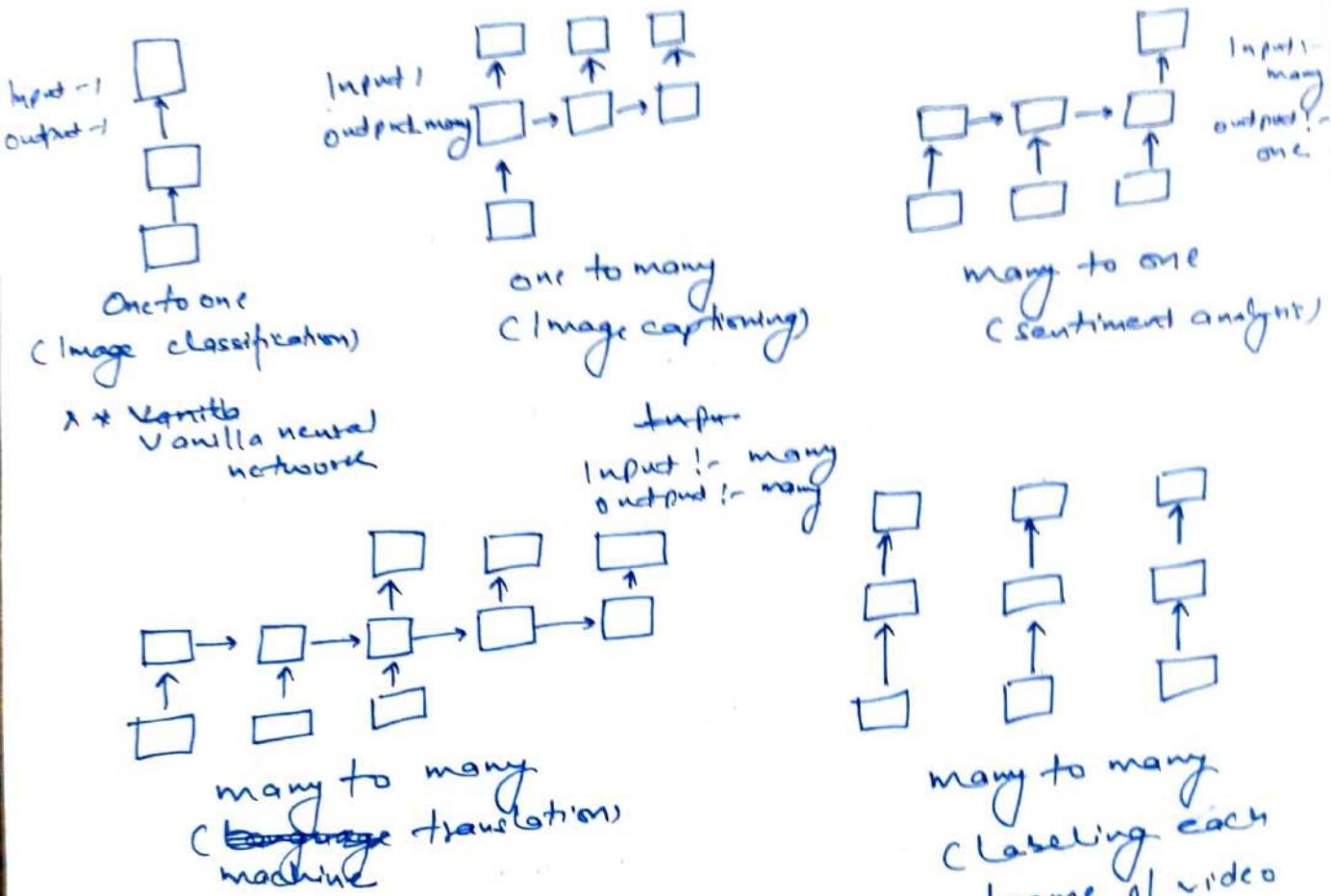
Training through RNN.

1. words are first transformed to machine readable vectors (x_t)
2. RNN process the sequence of vectors one by one that's single time step of the input is supplied to the RNN (x_t is supplied to the RNN)
3. calculated of hidden state h_t (using above formula)
4. The current h_t become h_{t-1} for the next time step.
5. We can go as many time steps the problem demands and combine the information from all previous state
6. finally RNN calculated output y_t (using above formula)
7. The output is then compared to the actual output and the error is generated
8. The error is then backpropagated to the network to update weights (using backpropagation through time) and the network is trained.

Sequence classification

- Application of RNNs :-
- 1. Sentiment classification (video classification)
 - 2. Sequence labeling (part-of-speech tagging, name entity recognition, word prediction)
 - 3. Image captioning
 - 4. Sequence generation (e.g. machine translation)

Various architecture of RNN
 RNN has the flexibility for handling various type
 of data



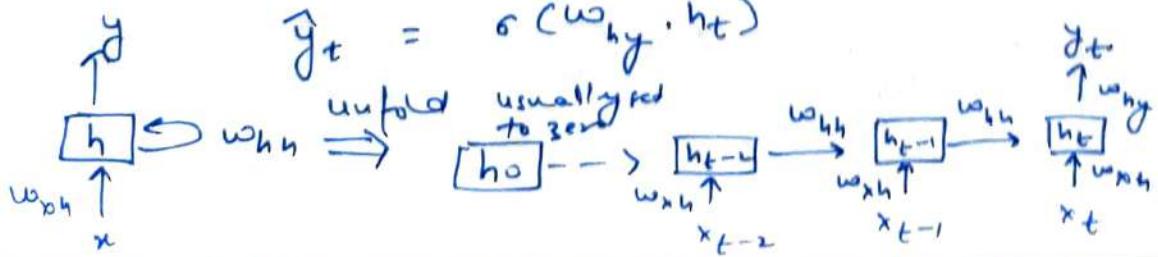
Back propagation through time (BPTT)
 BPTT is the application of Backpropagation training algorithm to RNN applied to sequence of data like a time series or Backpropagation algorithm that works on sequence in time is called BPTT

The main goal of using this algorithm is to obtain the parameters that optimize the cost function.

Here formula used

$$h_t = \sigma(w_{hh} \cdot h_{t-1} + w_{xh} \cdot x_t + b)$$

$$\hat{y}_t = \sigma(w_{hy} \cdot h_t)$$



Let error is calculated using cross entropy loss function

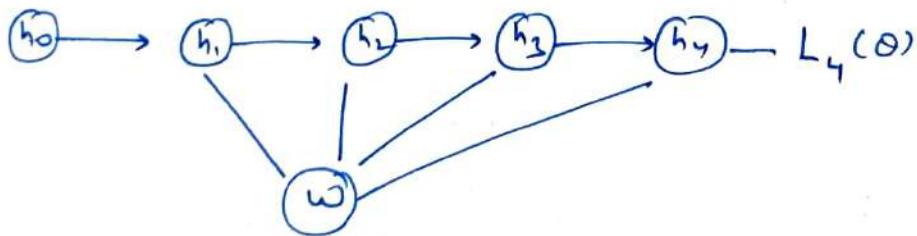
so overall loss of RNN

$$E(y, \hat{y}) = - \sum_t y_t \log \hat{y}_t$$

for each training epoch - start by training on shorter sequences and then train on progressively longer sequences until the length maximum sequence (1, 2, 3, ..., n-1, n)

for each length of sequence unfold the network into a normal feed forward network that has k hidden layers.

- BPTT works by unrolling all input time stamp. Each time stamp has one input time stamp one copy of the network and one output
- Errors are then calculated and accumulated for each time step. stamp.
- Backpropagate errors across the unfolded network and update the weights.



$$\frac{\partial L_4(\theta)}{\partial \omega} = \frac{\partial L_4(\theta)}{\partial h_4} \cdot \frac{\partial h_4}{\partial \omega}$$

$$\frac{\partial h_4}{\partial \omega} = \frac{\partial h_4}{\partial \omega} + \frac{\partial h_4}{\partial h_3} \cdot \frac{\partial h_3}{\partial \omega} + \frac{\partial h_4}{\partial h_3} \cdot \frac{\partial h_3}{\partial h_2} \cdot \frac{\partial h_2}{\partial \omega} + \frac{\partial h_4}{\partial h_3} \cdot \frac{\partial h_3}{\partial h_1} \cdot \frac{\partial h_1}{\partial \omega}$$

$$\frac{\partial h_4}{\partial \omega} = \sum_{k=1}^4 \frac{\partial h_4}{\partial h_k} \cdot \frac{\partial h_k}{\partial \omega}$$

(3)

$$\frac{\partial L_4(\theta)}{\partial w} = \frac{\partial L_4(\theta)}{\partial h_4} \cdot \sum_{k=1}^t \frac{\partial h_k}{\partial h_k} \cdot \frac{\partial h_k}{\partial w}$$

In general at time stamp 't'

$$\frac{\partial L_t(\theta)}{\partial w} = \frac{\partial L_t(\theta)}{\partial h_t} \cdot \sum_{k=1}^t \frac{\partial h_k}{\partial h_k} \cdot \frac{\partial h_k}{\partial w}$$

$$\frac{\partial h_t}{\partial h_k} = \frac{\partial h_t}{\partial h_{t-1}} \times \frac{\partial h_{t-1}}{\partial h_{t-2}} \times \frac{\partial h_{t-2}}{\partial h_{t-3}} \dots \frac{\partial h_{k+1}}{\partial h_k}$$

high - exploding gradient.

low/less - vanishing gradient

Limitation of RNN

if we have long chain of gradient, we can have two type of problems:-

① Vanishing gradients:- is going to be a product of many terms and if all these terms are very small then the gradient will vanish

② Exploding gradients:- is going to be a product of many gradient terms and if all these terms are very large then the gradient will explode

Limitation of RNN:- ① for deep RNN, RNN suffers by vanishing/exploding gradient problem

2. During back propagation RNNs suffer from vanishing gradient problem

- The vanishing gradient problem is when the gradient shrink as it back propagates through time.

- If the gradient value become extremely small, it doesn't contribute to much learning.

- It cannot process very long sequence if using tanh or ReLU as an activation function

- Training an RNN is very difficult task.

(2) RNN suffers from short term memory. If input sequence is long enough, RNN may have lost important information from the beginning.

Eg:- If input is a long paragraph of text to do prediction

RNN can forget what it seen in longer sequences, thus having a short term memory.

Need of long short term memory (LSTM)

LSTM is able to decide whether to keep the existing memory via the introduced gates.

RNN has two limitations:-

1. Due to vanishing gradient problem, RNN's effectiveness is limited. When it needs to go back deep into the context
2. There is no finer control over which part of the context needs to be carried forward and how much of the part needs to be forgotten.

Gradient are values used to $w' = w - \eta \frac{\partial L}{\partial w}$

update a neural network weight $w' = 0.001$ (small change in weight values)

so layers that get a smaller gradient update step learning (usually early layers)

As earlier layers don't learn, RNN's can forget what it seen in longer sequences thus having a short term memory

RNN:- suffers from short term memory. If a sequence / sentence is long enough, it will not carry information from earlier time steps to later ones.

for example:-

If we try to process a paragraph of text consist of many words, to do predictions using RNN,

it may leave out important information from ④
beginning
eg The clouds are in the —

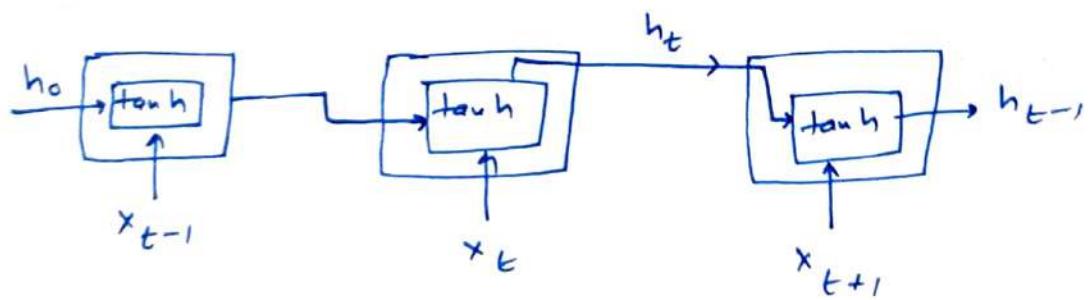
Ans! sky
like consider a language model trying to predict the next word based on the previous one
In such cases where the gap between the relevant information and the place that it is needed is small RNN can learn to use the part of information.

eg:- I grew in france.... I speak fluent —

~~Ans!~~ Ans! - french.

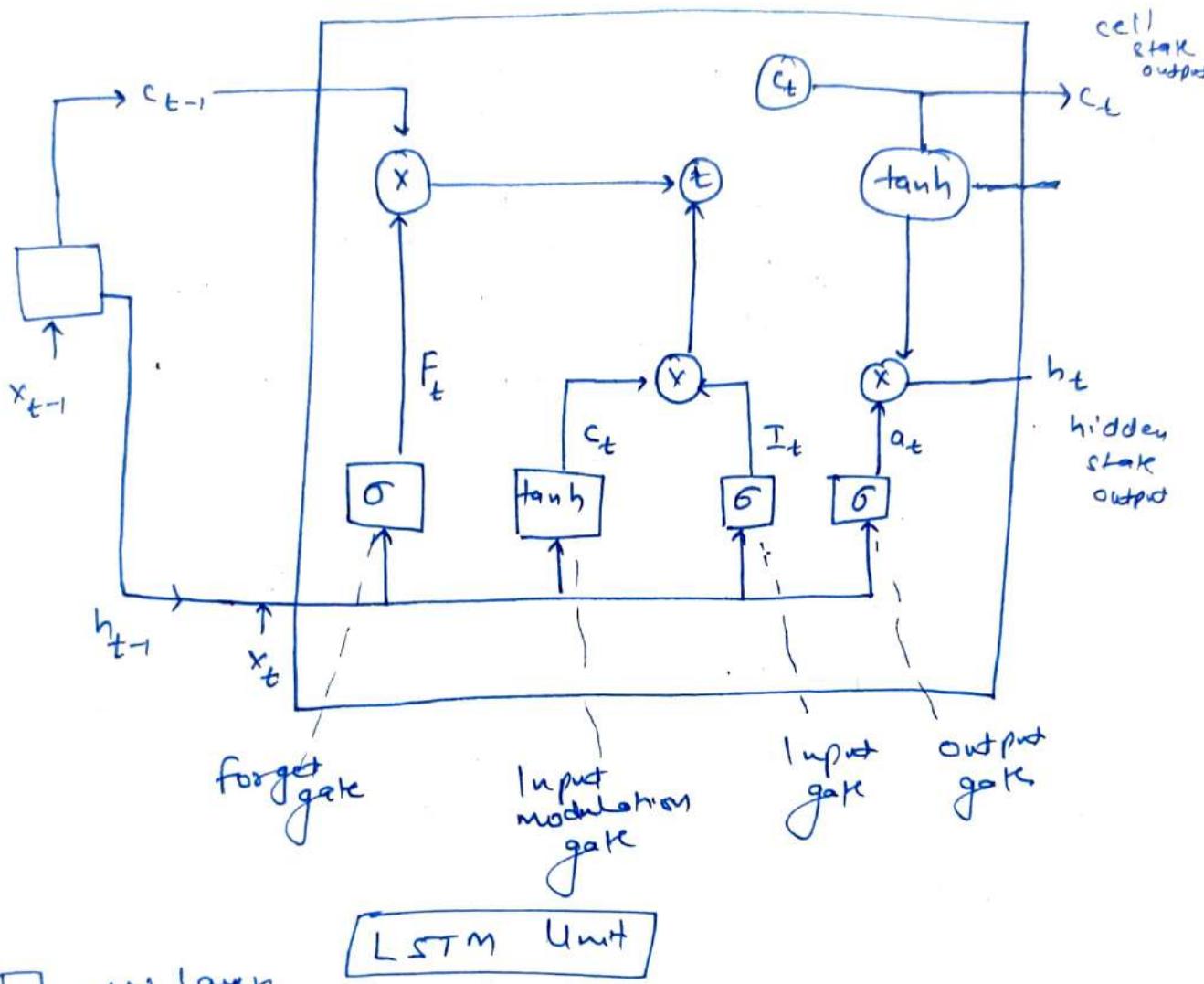
Here the gap between the relevant information and the point where it is needed to become very large. RNN unable to predict

General introduction of LSTM LSTM networks are special kind of RNN, capable of learning long term dependencies



RNN
The repeating module in a standard RNN contains a single layer (\tanh)

LSTM contains four interacting layers in a repeating module and different memory blocks called cells



Terms used in LSTM

cell state: memory of LSTM network. It carries relevant information throughout the processing of the sequence. As the cell state goes on its journey, information gets added / removed to the cell state via gates

Gates → Input gates
→ Forget gates
→ Output gates

The gates are different neural networks that decide which information is allowed on the cell state

During training, the gates can learn what information is relevant to keep forged

Sigmoid activation function: range [0,1]. It is helpful to update or forget data because any number getting multiplied by 0 is zero, causing values to disappear or be forgotten. and if any number multiplied by 1 is the same value, therefore value stay the same or is kept

Element wise / pointwise product \odot

$$\begin{bmatrix} 1 & 2 \end{bmatrix} \odot \begin{bmatrix} 2 & 3 \end{bmatrix}$$

$$\begin{bmatrix} 1 \times 2 & 2 \times 3 \end{bmatrix} = \begin{bmatrix} 2 & 6 \end{bmatrix}$$

Step by step procedure of LSTM

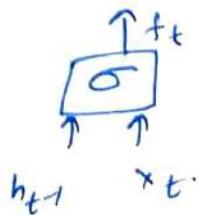
step 1: forget gate layer [What to keep and what to throw from old information]

In first step, it decided what information to throw away from the cell state and this decision is made by a sigmoid layer (or forget gate layer).
Forget gate layer / sigmoid layer takes h_{t-1} and x_t as input and outputs a number between 0 & 1 for each number in the cell state

c_{t-1} Here 1 represents - completely keep this information

0 represent - completely ignore / forget information from this state

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



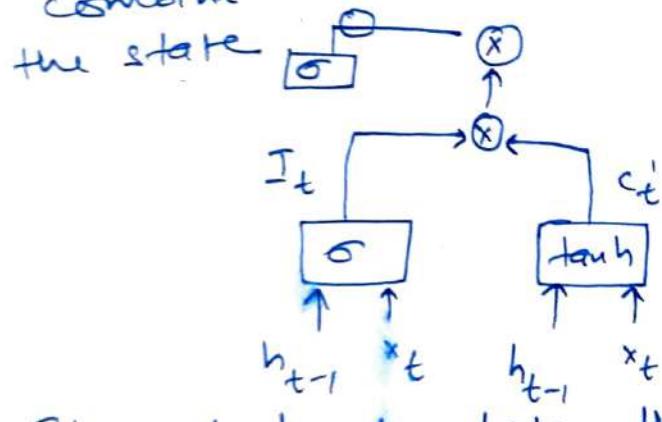
Step 2: Input Modulation gate layer and output gate layer.

In second step, it is decided what new information to store in the cell state. Or addition of useful information to the cell state is done by input gate.

It has two parts

① A sigmoid layer (input gate) decides which values to update layer

② Next a tanh layer (input modulation gate layer) creates a vector of new candidates values (c_t') that could be added to the state. Next combine these two to create an update to the state



Sigmoid function filters the values to be remembered using input h_{t-1} and x_t

tanh function creates a vector that gives output from -1 to +1 which contains all possible values from h_{t-1} and x_t

Next combine these two to create an update to the state

$$I_t = \sigma(w_i [h_{t-1}, x_t] + b_i)$$

$$c_t' = \tanh(w_c [h_{t-1}, x_t] + b_c)$$

How update of old cell state c_{t-1} to the new cell state c_t

① Multiply c_{t-1} by f_t - forgetting the things we decided to forget earlier ($f_t \odot c_{t-1}$) ⑥

ii) $c_t = (f_t \odot c_{t-1}) + (I_t \odot c'_t)$ Memory cell update

$I_t \odot c'_t$ is the new candidate values, scaled by how much we decided to update each state value.

Working of LSTM part 2

Step 3: Calculation of output

In this step, we decide what will be the output of this hidden layer.

This output will be based on cell state (filtered by

Sigmoid layer (output gate layer)

It decides what part of the cell state goes to output [The of gate decides what the next hidden state should be.]

first cell state pass through the tank, it pushes the values between -1 and +1 and then multiply it by output of the sigmoid gate (output layer) so that we only output the parts we decided to.

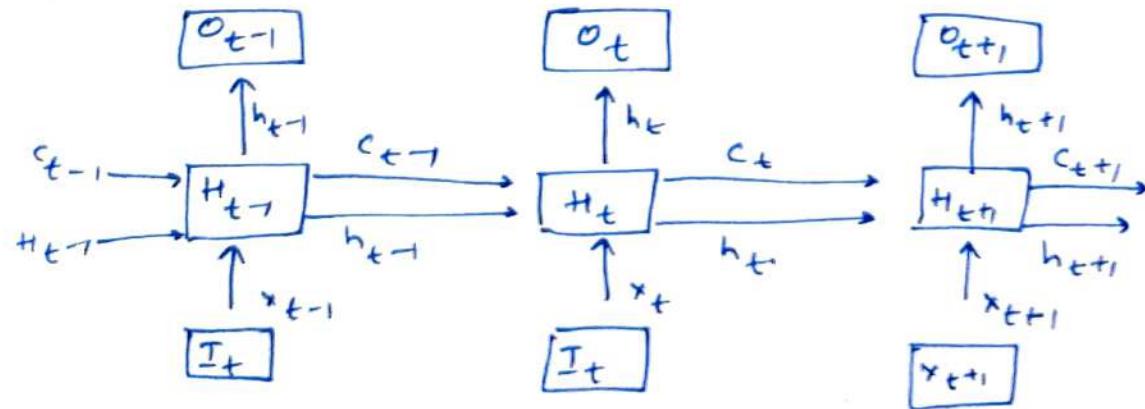
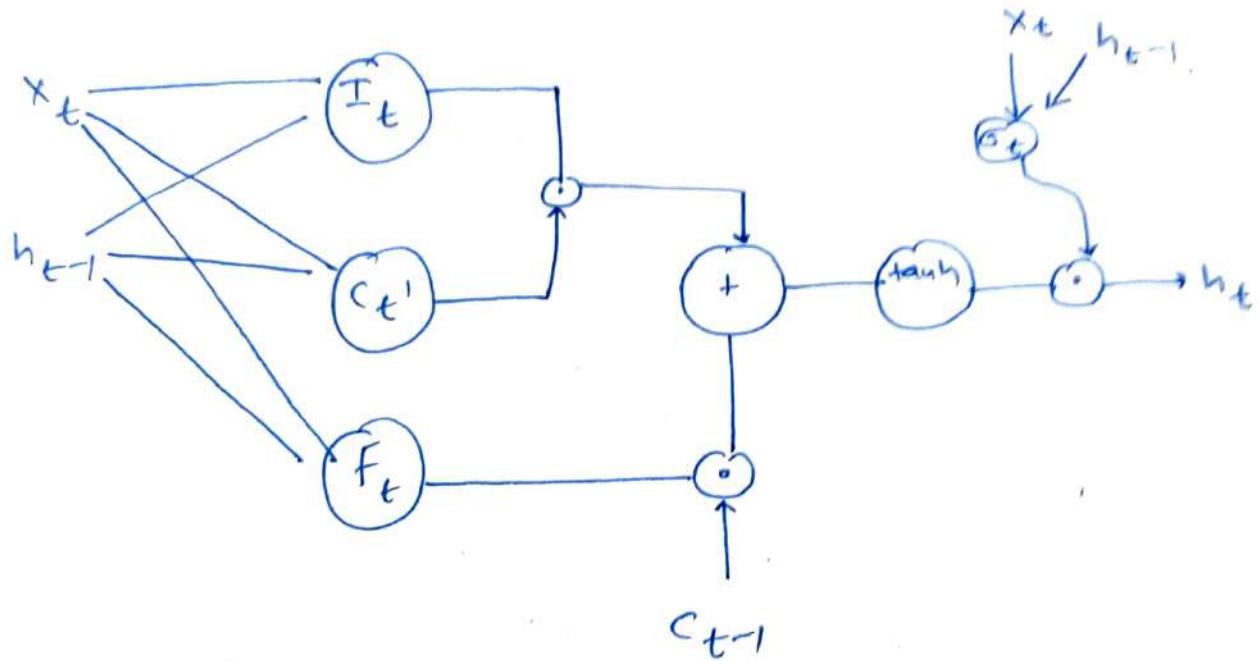


Fig LSTM

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

hidden
Layer output

$$h_t = o_t \otimes \tanh(c_t)$$



Gated Recurrent Recurrent Unit

(GRU)

$c_t \rightarrow$ next state
 h_t



Both solve the vanishing gradient problem of standard RNN

- Similar to LSTM, but there is no memory cell state
- It uses hidden state to transfer information
- I_t only has two gates
 - a reset gate
 - a update gate

update gate

- The update gate acts similar to the forget and input gate of LSTM
- It decides what information to throw away and what new information to add.

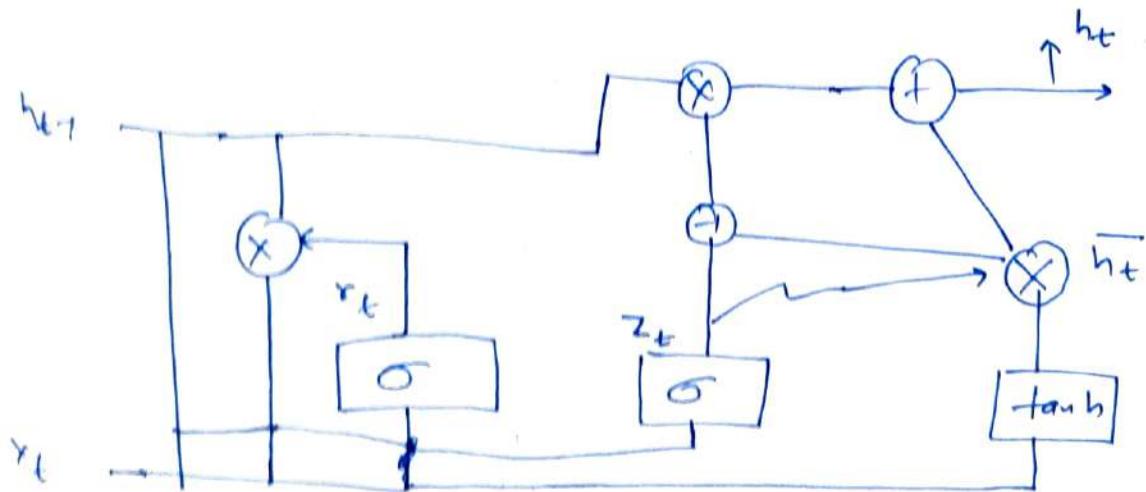
Reset gate: It is used to decide how much past information to forget

Advantages of GRU

- ① It is simple
- ② faster
- ③ optimizes quicker (at least in sentiments)

It is similar to LSTM in using gating functions, but differ from LSTM in that it doesn't have a memory cell / cell state and uses the hidden state to transfer information.

- Each GRU unit consist of update gate and reset gate while LSTM consist of four (W) gates
- LSTM or GRU can learn to keep only relevant information to make prediction and forget non-relevant data.



$$z_t = \sigma(w_2 \cdot [h_{t-1}, x_t])$$

$$r_t = \sigma(w_r \cdot [h_{t-1}, x_t])$$

$$\tilde{h}_t = \tanh(w \cdot [x_t * h_{t-1}, x_t])$$

$$h_t = (1 - z_t) * h_{t-1} + z_t * \tilde{h}_t$$

Machine translation translate a text from one natural language to another
English \rightarrow Hindi.

Sequence to sequence model are used for machine translation which are RNN based model
- finds the most relevant word in sentence for target language

Two type of sequence to sequence models

1. Encoder- Decoder model
2. Attention model

Encoder-Decoder model for machine translation.

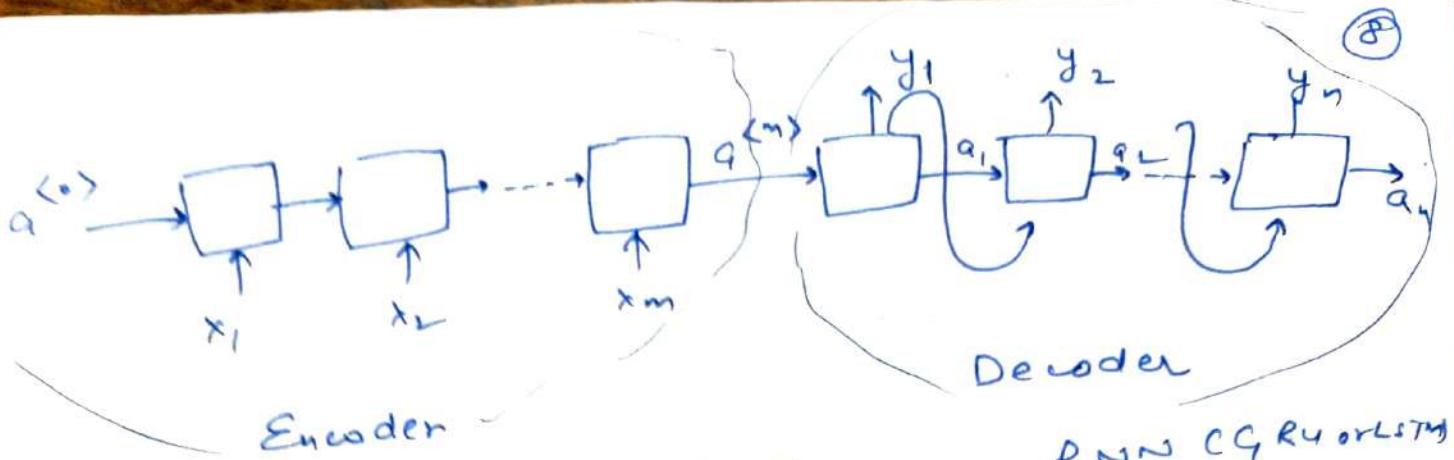
Given (English) Input sequence x_1, x_2, \dots, x_m
where x_i - words in input sequence
 m = number of words in input sequence
model have to find out / predict

(Hindi) output sequence - y_1, y_2, \dots, y_n

where y_i - words in output sequence

n = no. of words in output sequence

where m and n may be equal / different



The encoder network is built as an RNN (GRU or LSTM) where we feed the English words, one word at a time.

After processing the input sequence the output of the RNN is a vector that represents the input sequence.

- After encoder, there is a decoder network which takes as input the encoding output produced by encoder network
- Decoder is then trained to output the translation one word at a time until it outputs the whole o/p sequence
- Decoder acts mainly as a language model probabilistic model that are able to predict the next word in the sequence given the word that precedes it instead of having an input of all zeros here the input of decoder is an output vector of the encoder network

So machine translation model (like encoder-decoder) can be thought of as conditional language model for a system that translates English to Hindi. That means instead of modelling the probability of any ~~sequence~~ sentence it is now shaping the likelihood of the output (Hindi translation) conditioned on some input sentence (English text).

Due to the probabilistic behaviour of model there can be more than one translation to input sentence English sentence

$$\rightarrow \text{main aim} \rightarrow \arg\max P(y^{(1)}, y^{(2)}, \dots, y^{(T_y)} | x)$$

or

$$\arg\max P(y_1, y_2, \dots, y_n | x)$$

L most likely translation

Our aim to find the most suitable translation

so model doesn't sample the outputs at random, but it finds the Hindi sentence that maximize the conditional probability.

- To maximize the conditional probability or to find a most likely translation for a given input sentence the most common algorithm is used that is called Beam search algorithm

Beam search To explain beam search, let we are using a

sequence to sequence model that was an encoder and decoder framework with LSTM or GRU as the basic block

the encoder maps a source sequence encodes the source information (Let here English sentence) and the decoder takes the encoded data from encoder as input to produce an output sequence (Let here Hindi sentence)

- we always want the best and most likely words to be selected for the translation that matches the meaning of the English sentence

- There could be multiple likely translation of the source sentence (English) in the target language (Hindi).

So our goal is to pick the best and most likely translated word, so we choose the target word with the maximum probability based on the source sentence

The Beam search algorithm selects multiple alternatives for an input sequence at each timestamp based on conditional probability

$$J = \hat{f}_1^1, \hat{f}_1^2, \hat{f}_1^3$$

The number of multiple alternatives depends on a parameter called Beam width B

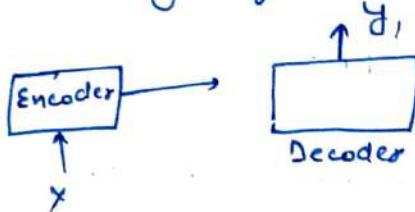
e.g. e.g. $B=2$

so at each time step the beam search selects B number of best alternatives with the highest probability as the most likely possible choices for the timestep.

Example:- Let Beam width = 3

vocabulary of Hindi word = 10,000 words

Step 1

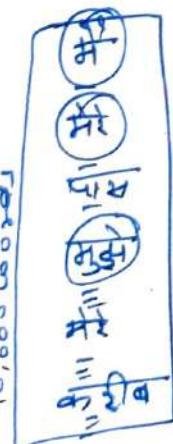


$$P(y_t | x) = [मैं, मेरे, मुझे]$$

(I have some money)

Hindi vocabulary

step 1 finds the 3 words with the highest softmax probability based on the input sequence



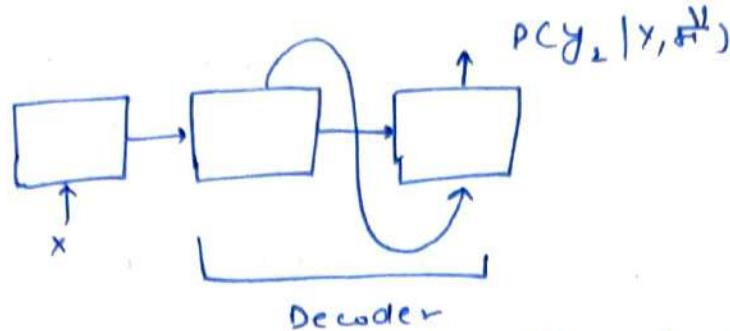
$$B=3$$

Step 2 This step finds top 3 pairs of words for first and second word of the translated sequence (Hindi)

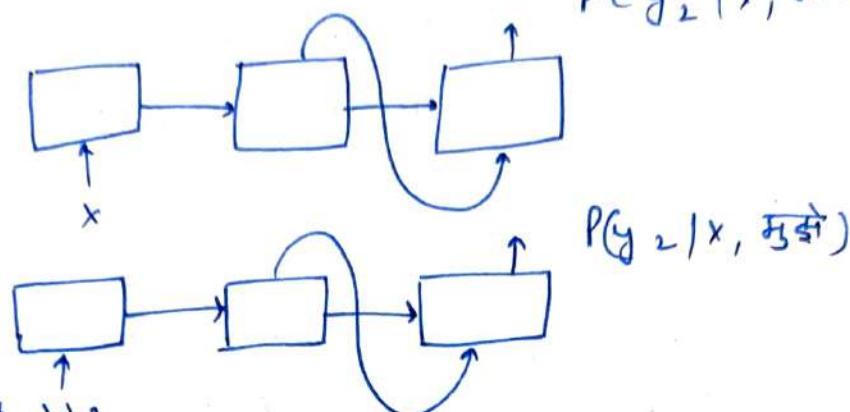
मेरा साथ करीब

मेरा साथ पास करीब

मुझे +



$$P(y_2 | x, \text{मेरा})$$



$$P(y_2 | x, \text{मुझे})$$

$$P(y_1, y_2 | x) = [\text{मेरा गास}, \text{मेरे गास}, \text{मेरे साथ}]$$

- The first three selected words (मेरा, मेरे, मुझे), from step 1 is given as input to second step
- Then softmax function is applied to all the 10,000 words in the Hindi vocabulary to find the 3 best alternatives for the second word that are most likely to form a pair with first word conditional probability
- First, we will take word मेरा and apply the softmax function to all 10,000 words in the vocabulary which calculated the probability with मेरा (10^{1000})
- Similarity probability for other two words, मेरे, मुझे is calculated using softmax functions ($10,000 + 10,000 = 20000$)
- Run 30,000 different combinations to check the top 3 pairs of the first and second words.
- Let the top 3 first and second word pair combinations are मेरा गास, मेरे गास, मेरा साथ

Let us did not find any word pair with 索引 as first word and second word having a high conditional probability so we dropped it. (10)

same process is repeated to find 3 best pair for the first, second and third word based on the input sequence and the chosen first and second word.

$$P(y_1|x) = [\text{ए}, \text{मैं हूँ}, \text{मुझे}]$$

$$P(y_1, y_2|x) = [\text{मैं हूँ}, \text{मैं हूँ तो}, \text{मैं हूँ तो}]$$

$$P(y_1, y_2, y_3|x) = [\text{मैं हूँ तो}, \text{मैं हूँ तो तो}]$$

This process continues and finally we get two output sentence with highest conditional probability and different lengths.

I have some money — [मैं हूँ तो, मैं हूँ तो]

[मैं हूँ तो तो, मैं हूँ तो तो] (0.2)

(first sentence)

- we finally search pick the output of the decoder as the sentence with highest probability.

I have some money — मैं हूँ तो तो

→ Beam search consider multiple best options based on beam width using conditional probability

- If we take higher beamwidth, it will give a better translation but would use a lot of memory and computational power

→ So beam search consider multiple best options based on beam width using conditional probability

- If we take higher beam width, it will give a better translation but would use a lot of memory and computational power

for calculating conditional probability for 30000 words at each time step.

B = 3 100000 words
 vocabulary

If we take lower beam width will result in more low quality translation but will be fast and efficient in terms of memory usage and computation power

BLEU score A method for automatic evaluation of m/c translation

- BLEU is a score for comparing a candidate machine translation of text to one or more reference translations (human translated)
- It can also be used to evaluate text generated for a suite of NLP tasks (like generation, recognition, image caption generation, text summarization, sequence recognition)
- BLEU is a metric for evaluating a generated sentence (often translation) to a reference sentence
 - BLEU metric ranges from 0 to 1
 - score = 1.0 A perfect match result
 - score = 0.0 A perfect mismatch results
- Bleu score can be used to evaluate the performance of a MT system, particularly when there are multiple equally good answers
- Bleu score was developed for evaluating the prediction made by automatic machine translation system, but it is not perfect evaluation method.
- Benefit of Bleu score.
 1. Quick and inexpensive to calculate
 2. Easy to understand
 3. Language independent
 4. It correlates highly with human evaluation
 5. widely adopted method.

Machine translation system requires
① A numerical translation closeness metric

- (1)
- A good quality human reference translation
 - BLEU score is based on "f_n string matches".
 - It is a score that quantifies how good a machine translation is by computing a similarity score based on n-gram precision.
 - Look at each word in the candidate n-gram precision: sentence and assign it a score of 1 if it shows up in any of the reference sentence otherwise 0

example: MT output - The cat ¹ the ² cat ³ on the mat

Reference 1: There ¹ is ² a cat ³ on the mat

Reference 2: The cat ¹ is ² on ³ the mat

Unigram precision:

Unique Unigram

count / score

candidate sentence

the

3

The cat the cat on the mat cat

2

Sequence of n or \leftarrow n-gram
consecutive words

Unigram one words sequence

bigram two words sequence

trigram three words sequence

4-gram four words sequence

- The main task of a BLEU implementor is to compare n-gram of the candidate with Two n-grams of the reference translation and count the number of matches

- The number of matches are position independent
if there are m_1 of more m_2 of matches between candidates and reference translations, then better is the translation

- Higher the Bleu score better is the translation
- More reference human translations result in better and accurate score

Unigram precision

Candidate sentence

The cat the cat on the mat

number of unigrams in
candidate sentence = 7

Unique Unigram

the

cat

on

mat

count / score

2

2

1

1

7

Total counts for
the unique unigram in the
candidate sentence

$$\text{unigram precision} = \frac{10}{7} = 1.0$$

Bigram precision :- candidate sentence - The cat the cat on the mat

Bigram: the cat, in cat the, in cat on, on the, the mat

1 2 3 4 5 6

count / score

Unique Bigram

the cat

2

$$\text{Bigram precision} = \frac{5}{6}$$

cat the

0

$$= 0.833$$

cat on

1

on the

1

the mat

1

Total count

5

Drawback Unigram precision is very easy to be over confident about the quality of MT
eg of poor MT output with high precision.

candidate sentence the the the the the the

$$\text{unigram precision} = \frac{7}{7} = 1.0$$

To overcome the unigram precision problem, clipped count and modified precision were proposed
modified n-gram precision.

(12)

Unigram modified precision

unique unigram	count	clipped count
the	3	2
cat	2	1
on	1	1
mat	1	1
	<u>7</u>	<u>5</u>

unigram modified precision $\frac{5}{7} = 0.714$

Unique Bigram	count	clipped count	
the cat	2	1	
cat the	0	0	
cat on	1	1	
on th	1	1	
the mat	1	1	
	<u>7</u>	<u>4</u>	

Unigram modified precision

unique unigram	count	clipped count
the	7	2
precision	$\frac{2}{7}$	

Problem in modified n-gram precision

Bleu scores tend to favour short translation, which can produce very high precision scores, even using modified precision.

example: candidate sentence :- the cat

Reference sentence :- the cat is on the mat

modified unigram precision $= (1+1)/2 = 1$

solution is to Brevity penalty

$$BP = \begin{cases} 1 & \text{if } c > r \\ e^{(1-\frac{r}{c})} & \text{if } c \leq r \end{cases}$$

r = length of reference sentence
 c = length of candidate sentence

for above example

$$c=2, r=6$$

as $c < r$

$$BP = e^{(1-\frac{r}{c})} = e^{(1-\frac{6}{2})} = e^{-2}$$

Bleu score

$$BP \times \exp \left[\frac{1}{N} \sum_{n=1}^N P_n \right]$$

Geometric mean
of all n-gram
precision

P_n - modified precision
for n-gram

Bleu only utilizes a BP for shorter sentence
The redundant of longer sentences are not penalized properly

Attention Model

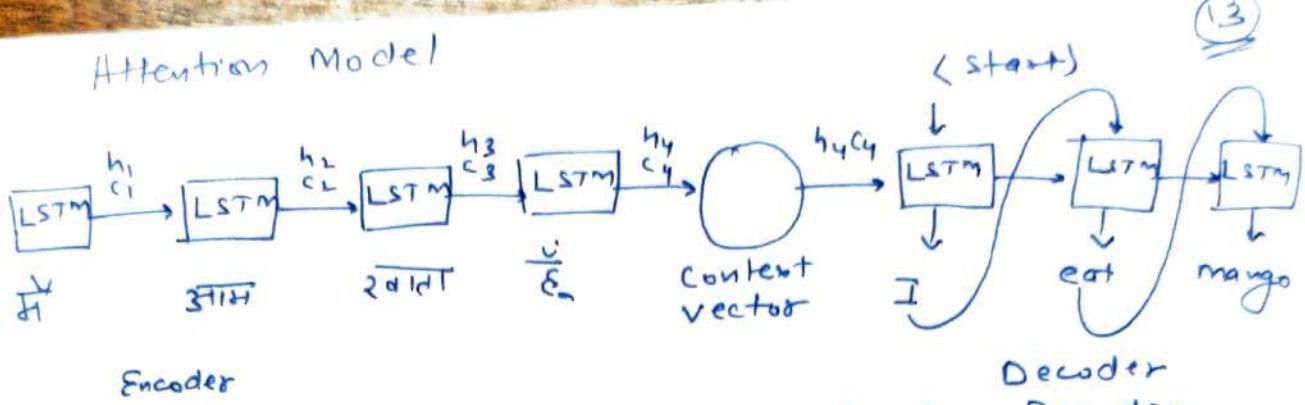


fig translation using Encoder - Decoder

- For long sentences (input) context vector may not have all necessary information
- To solve this problem, attention mechanism is used with encoder - Decoder RNN / LSTM. Context h_4c_4 contains good summary of I , eat, mango.
- ↑
fixed size
Attention means directly your focus at something and taking greater notice
- or concentrating on one or a few things while ignoring others or concentrating on one or a few things while ignoring others
- The attention mechanism is an improvement over the encoder - decoder based on neural based machine translation system in NLP.
- In encoder-decoder RNN / LSTM
 - Encoder process the entire input sequence and encode it into a context vector, which is the last hidden state of the LSTM / RNN.

\overrightarrow{I} $\overrightarrow{\text{eat}}$ $\overrightarrow{\text{mango}}$

I eat mango

t_1	I	0	0	0
t_2	α_{11}	α_{21}	α_{31}	α_{41}
	0	0	0.5	0.5
t_3	α_{12}	α_{22}	α_{32}	α_{42}
	0	1	0	2
	α_{13}	α_{23}	α_{33}	α_{43}

$$t_1 \quad \overrightarrow{I} \quad [1 \ 0 \ 0 \ 0] \ I$$

$$t_2 \quad \overrightarrow{\text{eat}} \quad [0 \ 0 \ 0.5 \ 0.5] \ \text{eat}$$

$$t_3 \quad \overrightarrow{\text{mango}} \quad [0 \ 1 \ 0 \ 0] \ \text{mango}$$

α_{ij}

improvement
 $i = \text{word}$
 $j = \text{timestamp}$.

It is expected that context vector contains good summary of input sequence. All intermediate states of the encoder are ignored and the final state is supposed to be the initial hidden state of the decoder.

- Decoder unit produces the word in a sentence one after another

Drawback of Encoder-Decoder machine translation system:- If the encoder makes a bad summary, the translation will also be bad.

Usually encoder creates a bad summary when it tries to understand long sentences.

- RNN cannot remember long sentence and sequences due to vanishing/exploding gradient problem. It can remember the parts which it has just seen. So the performance of the encoder-decoder network degrades rapidly as the length of the input sentence increases.

LSTM is supposed to be better than RNN, but it becomes forgetful in specific case.

But in both RNN and LSTM there is no mechanism to give more importance to some of the input words compared to others while translating the sentence.

So when model generates a sentence (machine translation) it searches for a set of positions in the encoder hidden states where the most relevant information is available. This idea is called attention.

- Attention is proposed as a solution to the limitation of encoder-decoder model encode the input sequence to one fixed length vector from which to decode each output timestep. This issue is believed to be more of a problem when decoding long sequences.

The attention model allows an RNN to pay attention to specific part of the input that is considered as being important which improves the performance of the resulting model in practice. (19)

In attention model, instead of sending last hidden state output of decoder, we could take a weighted average of the corresponding word representation and feed it to the decoders.

for example; at timestamp $t=2$, we can take a weighted average of the representation \overrightarrow{h}_{T_x} & total output eat from the decoder. By doing this, we are not overloading the decoder with irrelevant information (like ~~the time~~ ~~STH~~)

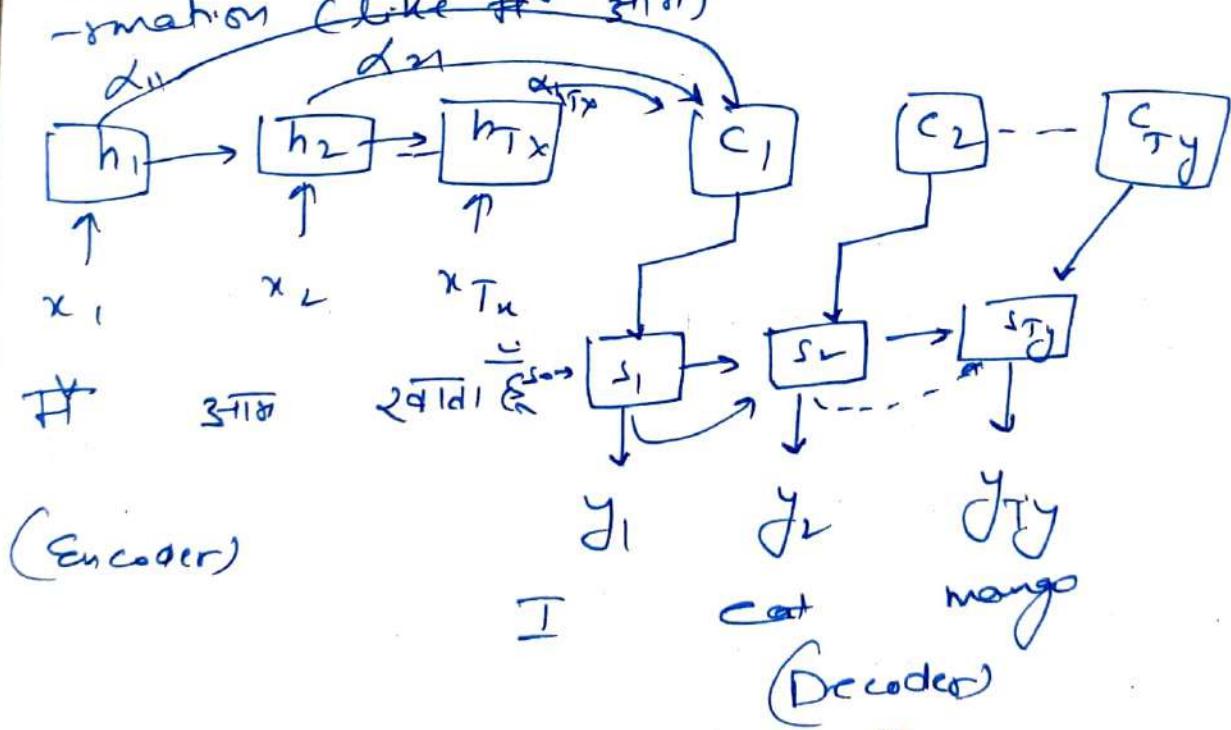


fig. Attention mechanism
in Encoder-decoder model

Prediction of II words

$$t=2 \quad c_2 = \alpha_{12} \cdot h_1 + \alpha_{22} \cdot h_2 + \alpha_{32} \cdot h_3 + \alpha_{42} \cdot h_4$$

$$\alpha_{12} = g(s_2, h_1) \quad \alpha_{12} = \frac{e^{\alpha_{12}}}{e^{\alpha_{12}} + e^{\alpha_{22}} + e^{\alpha_{32}} + e^{\alpha_{42}}}$$

$$Q_{2L} = a(s_1, h_2) \alpha_{2L} = \frac{e^{Q_{2L}}}{e^{Q_{1L}} + e^{Q_{2L}} + e^{Q_{3L}} + e^{Q_{4L}}}$$

$$Q_{3L} = a(s_1, h_3) - - -$$

$$c_t = \sum_{j=1}^T \alpha_{jt} h_j$$

h_j - encoder hidden state
at j^{th} timestamp.

α_{jt} - probability of focusing on
the j^{th} word to produce
the word

$$\alpha_{jt} = \frac{e^{Q_{jt}}}{\sum_{k=1}^T e^{Q_{kt}}}$$

$$Q_{jt} = \underbrace{a(s_{t-1}, h_j)}_{\text{tanh}}$$

$$\alpha_{jt} = \frac{\exp(Q_{jt})}{\sum_{k=1}^T \exp(Q_{kt})}$$

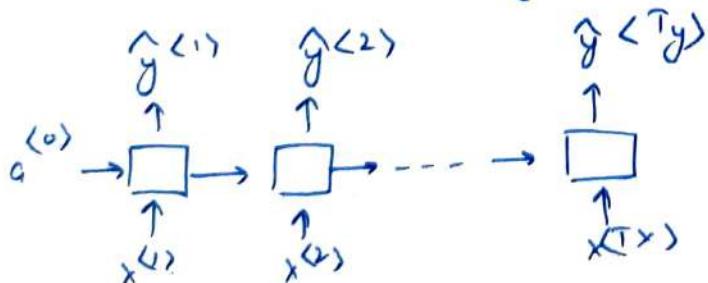
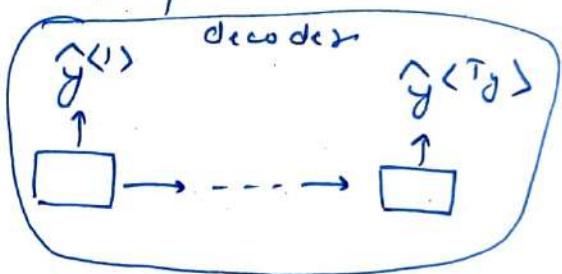
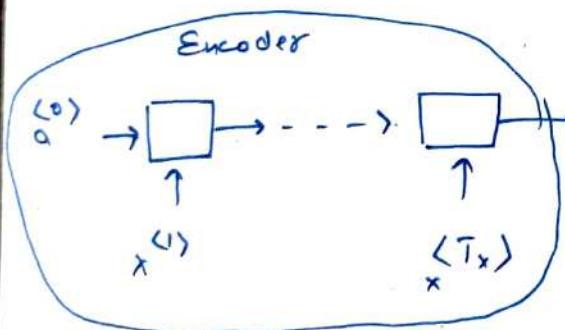
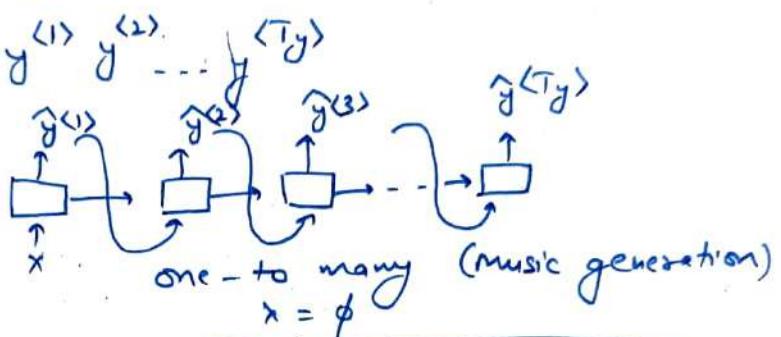
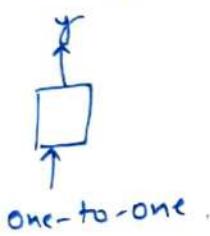
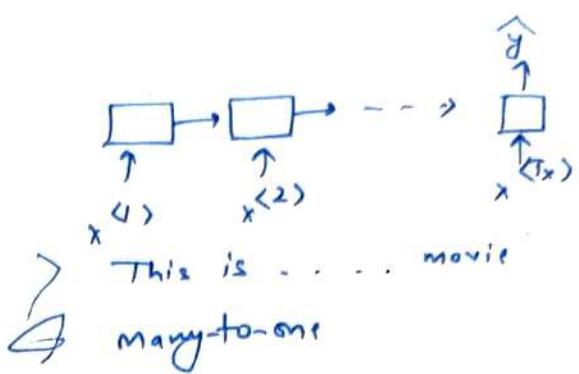
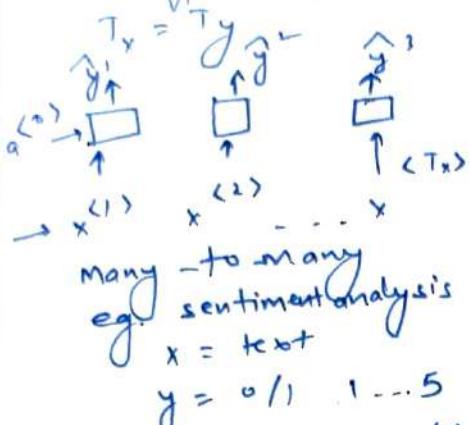
The importance of the j^{th} input word for decoding the t^{th} output word can be calculated by

$$Q_{jt} = F_{ATT}(s_{t-1}, h_j) \quad a F_{ATT} = \begin{matrix} \text{feed forward} \\ \text{neural network} \end{matrix}$$

the alignment model score :- how well each word
input (h) matches the current output of the decoders (1)

$$\overbrace{=}^{=}$$

Different types of RNN



Language model and sequence generation
what is language modelling?

① — Speech recognition

② — The apple and pair salad

③ — The apple and pear salad

$$P(1) = 3 \cdot 2 \cdot 10^{-12}$$

$$P(2) = 5 \cdot 7 \cdot 10^{-10}$$

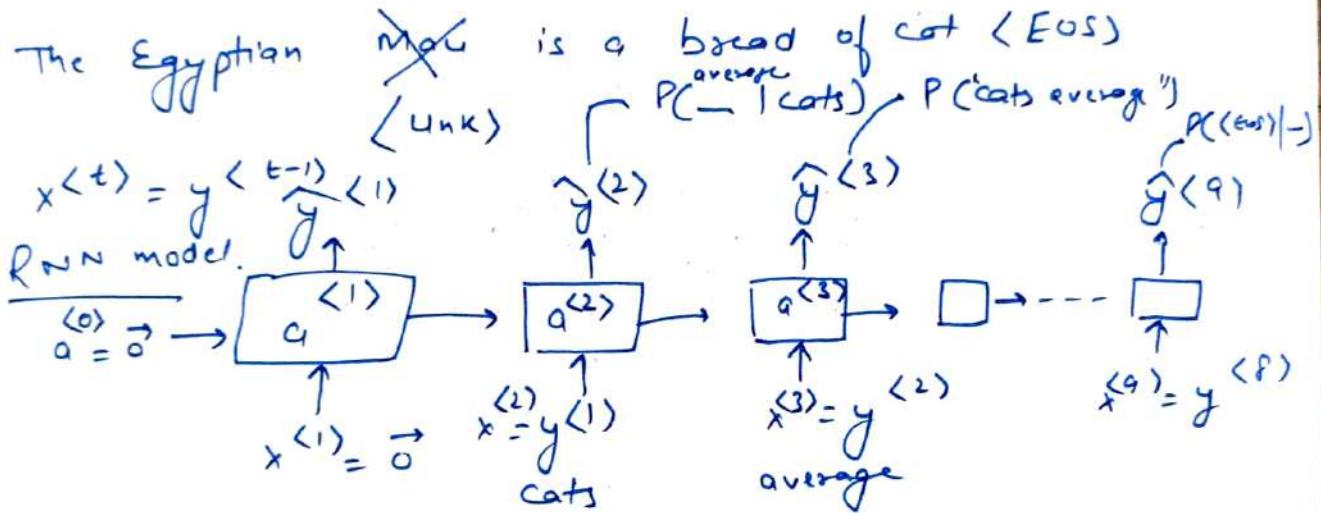
$P(\text{sentence}) = ?$

$$P(y^{(1)}, y^{(2)}, \dots, y^{(T_y)})$$

Training set: Large corpus of English text

Cats average
 $y^{(1)}$ Tokenize
 $y^{(2)}$ $y^{(3)}$ $y^{(4)}$ \dots $y^{(9)}$ $y^{(10)}$ \dots $y^{(T_y)}$

15 hours of sleep a day. $\langle \text{EOS} \rangle$
 end of sentence



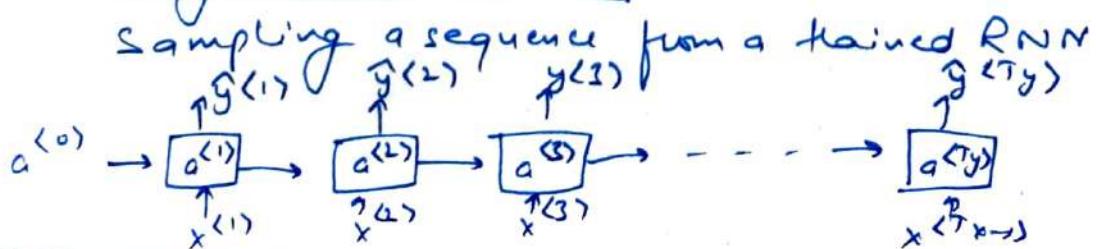
$$P(a) = P(a_{\text{aran}}) \cdots P(\text{cats}) \cdots P(zulu) \cdots P(\text{UNK}) \cdots P(\langle \text{EOS} \rangle)$$

$$L(\hat{y}^{(t)}, y^{(t)}) = - \sum_i y_i^{(t)} \log \hat{y}_i^{(t)}$$

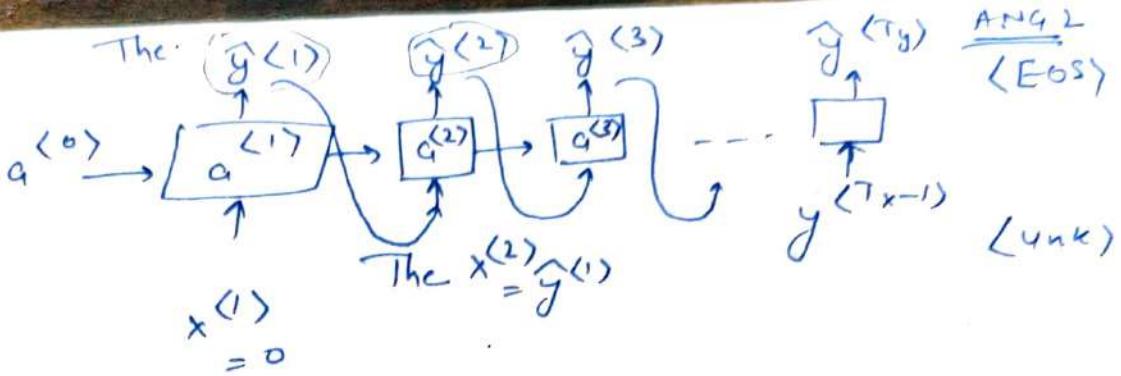
$$L = \sum_t L^{(t)}(\hat{y}^{(t)}, y^{(t)})$$

$$\begin{aligned} & P(y^{(1)}, y^{(2)}, y^{(3)}) \\ &= P(y^{(1)}, P(y^{(2)} | y^{(1)})) \\ &= P(y^{(3)} | y^{(1)}, y^{(2)}) \end{aligned}$$

Sampling novel sequence



sampling



$P(a) P(a|a_{\text{prev}}) \dots P(z_{\text{last}}) P(y_{\text{unk}})$ n.p. random choice

$P(-1 \text{ the})$

Vocabulary - [a, b, c, ..., z] — [A, ..., Z]

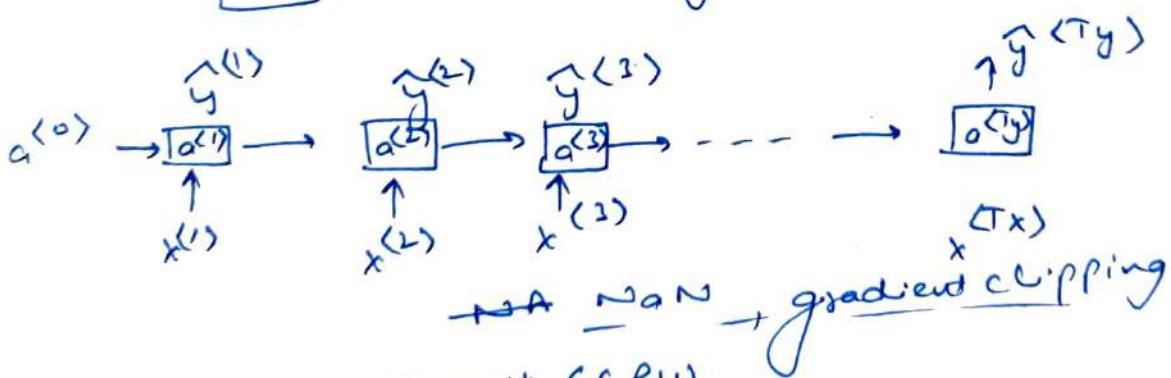
$y^{(1)}, y^{(2)}, y^{(3)}, y^{(4)}$ (cat, average)
c a t space

Sequence generation

Vanishing gradient with RNNs

The cat, which already ate ---, was full

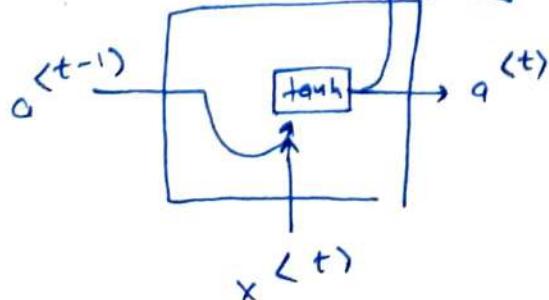
The cats, which already ate ---, were full



Gated Recurrent Unit (GRU)

RNN unit

$$a^{(t)} = g(w_a[a^{(t-1)}, x^{(t)}] + b_a)$$



1. RNN (Simple RNN)

c = memory cell

$$c^{(t)} = \tanh(\omega_c [c^{(t-1)}, x^{(t)}] + b_c)$$



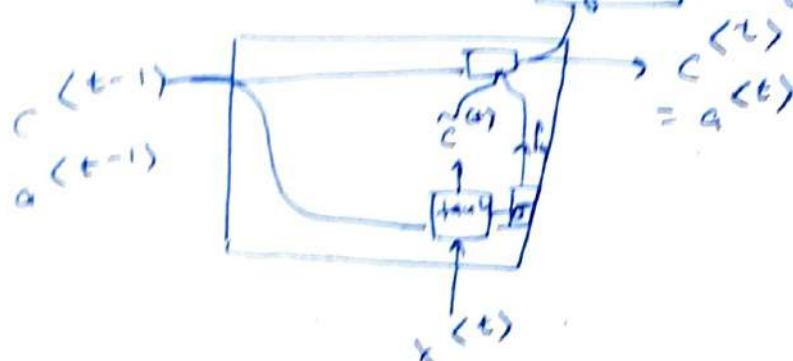
$$\Gamma_u = \sigma(\omega_u [c^{(t-1)}, x^{(t)}] + b_u)$$

$$\rightarrow \Gamma_{u(t)} = 1, \Gamma_u = 1 \quad \Gamma_u = 0, \Gamma_u = 0, \Gamma_u = 1, \Gamma_u = 0$$

The cat, which already ate... $\boxed{\text{was}}$ full

$$c^{(t)} = \Gamma_u * c^{(t)} + (1 - \Gamma_u) * c^{(t-1)}$$

(softmax) $\rightarrow \hat{y}^{(t)}$



$$\Gamma_r = \sigma(\omega_r [c^{(t-1)}, x^{(t)}] + b_r)$$

LSTM Long short term memory

GRU

$$c^{(t)} = \tanh(\omega_c [\Gamma_r * c^{(t-1)}, x^{(t)}] + b_c)$$

$$\Gamma_u = \sigma(\omega_u [c^{(t-1)}, x^{(t)}] + b_u)$$

$$\Gamma_r = \sigma(\omega_r [c^{(t-1)}, x^{(t)}] + b_r)$$

$$c^{(t)} = \Gamma_u * c^{(t)} + (1 - \Gamma_u) * c^{(t-1)}$$

$$a^{(t)} = c^{(t)}$$

L = 19

A. 1184, 1

$$c^{(t)} = \tanh(w_t[a^{(t-1)}, x^{(t)}] + b_t)$$

$$l_u = (w_u[a^{(t-1)}, x^{(t)}] + b_u) \quad \text{update}$$

$$l_p = \sigma(w_p[a^{(t-1)}, x^{(t)}] + b_p) \quad \text{forget}$$

$$l_o = \sigma(w_o[a^{(t-1)}, x^{(t)}] + b_o) \quad \text{output}$$

$$c^{(t)} = l_u * c^{(t-1)} + l_p * c^{(t-1)}$$

[Propagating memory
, t-1)

$$a^{(t)} = l_o * c^{(t)}$$

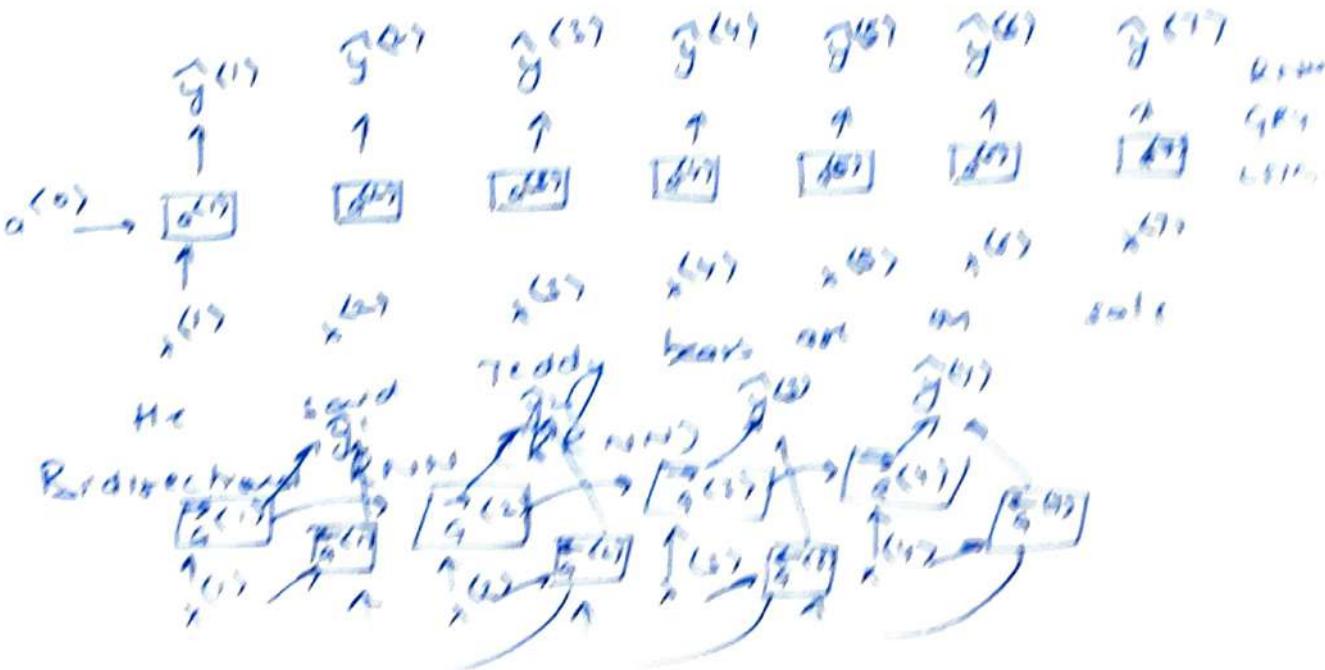
$$a^{(t)} = l_o + \tanh c^{(t)}$$

Bi-directional RNN

Getting information from the future

He said, "Teddy bears are on sale!"

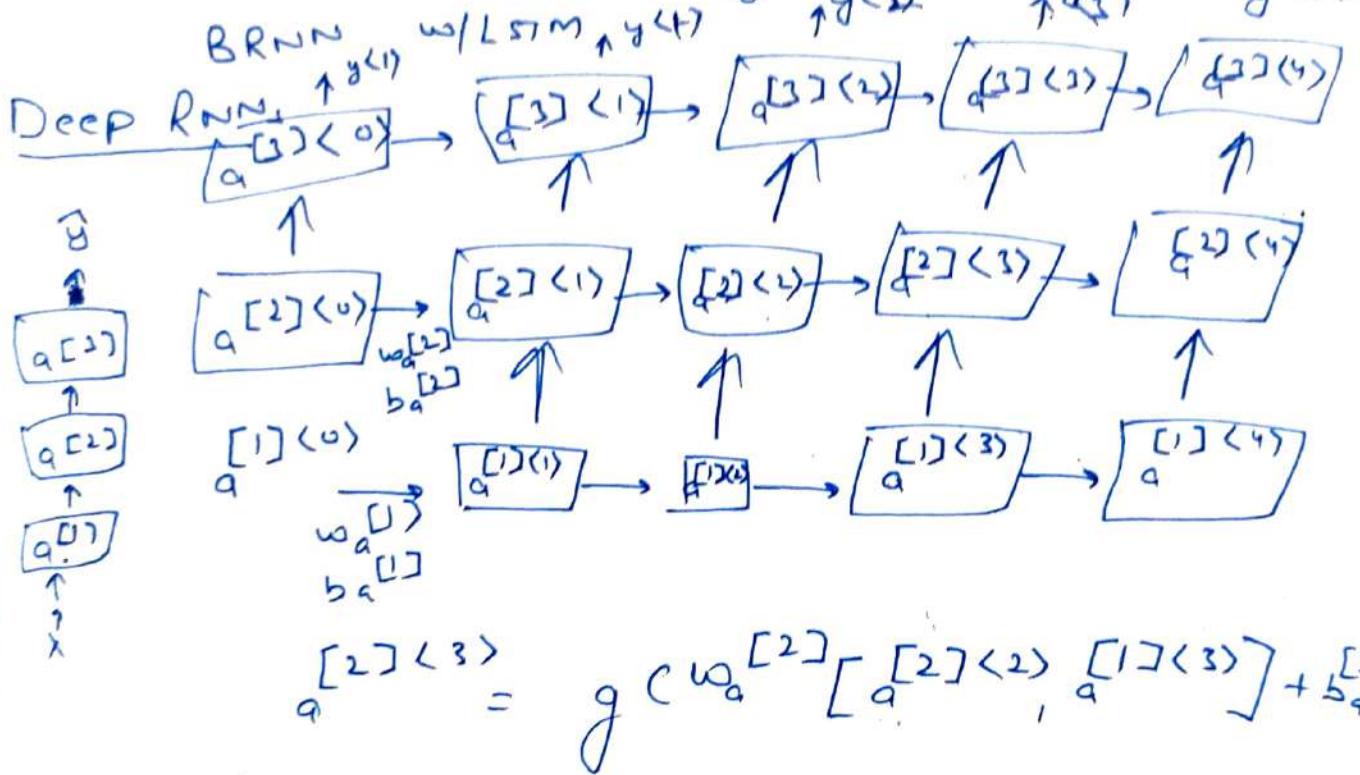
He said, "Teddy Roosevelt was a great President!"



Acyclic graph

$$\vec{y}^{(t)} = g(w_y [\vec{a}^{(t)} \leftarrow \vec{a}^{(t)}] + b_y)$$

He said "teddy Roosevelt" $\vec{y}^{(1)} \vec{y}^{(2)} \vec{y}^{(3)} \vec{y}^{(4)}$



Word Representation

1- hot representation

Man (539)	woman (9853)	King (4914)	Queen (7157)	Apple (451)	Orange (6257)
$\begin{bmatrix} 0 \\ \vdots \\ 0 \\ 1 \\ 0 \\ \vdots \\ 0 \end{bmatrix}$	$\begin{bmatrix} 0 \\ \vdots \\ 0 \\ 1 \\ 0 \\ \vdots \\ 0 \end{bmatrix}$	$\begin{bmatrix} 0 \\ \vdots \\ 0 \\ 1 \\ 0 \\ \vdots \\ 0 \end{bmatrix}$	$\begin{bmatrix} 0 \\ \vdots \\ 0 \\ 1 \\ 0 \\ \vdots \\ 0 \end{bmatrix}$	$\begin{bmatrix} 0 \\ \vdots \\ 0 \\ 1 \\ 0 \\ \vdots \\ 0 \end{bmatrix}$	$\begin{bmatrix} 0 \\ \vdots \\ 0 \\ 1 \\ 0 \\ \vdots \\ 0 \end{bmatrix}$

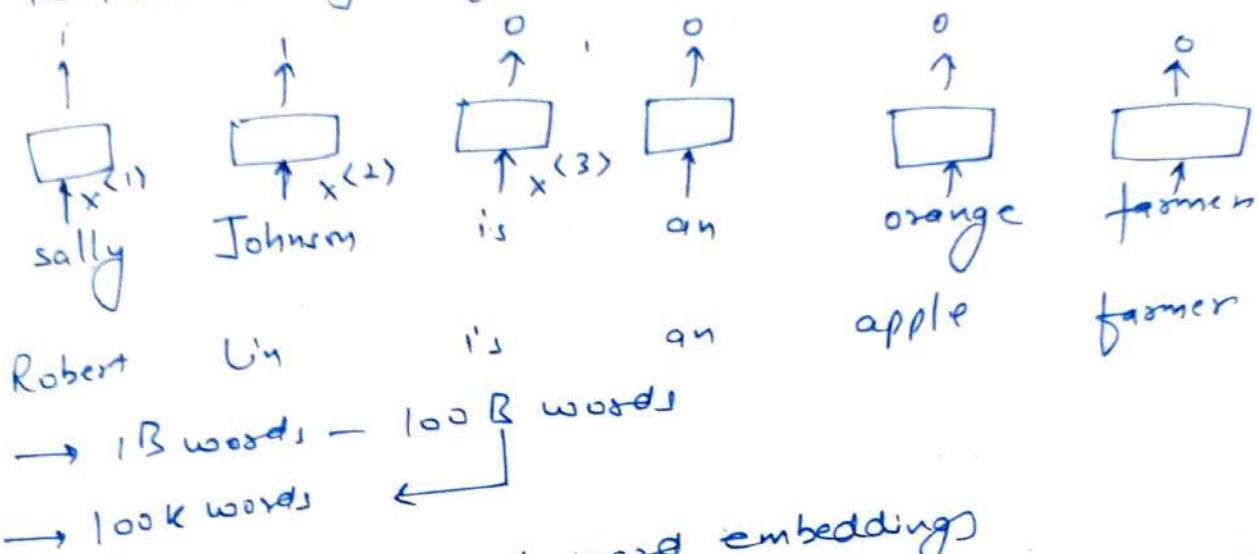
featureized representations : word embedding

Gender	-1	1	0.95	0.97	0.00	0.01
Royal	0.01	0.02	0.93	0.95	-0.01	0.00
King	0.03	0.01	0.7	0.69	0.02	-0.02
Food	0.04	0.01	0.02	0.01	0.95	0.97

Visualizing word embeddings

Using word embedding

Named entity recognition example



Transfer learning and word embeddings

Transfer learning and word embeddings from large text corpus (1-100B words)

- (1). Learn word embedding from large text corpus (1-100B words)
(or download pretrained embedding online.)
- (2). Transfer embedding to new task with smaller training set?
3. optional! continue to finetune the word embeddings with new data.

Relationship to face encoding

Properties of word embedding

man → woman

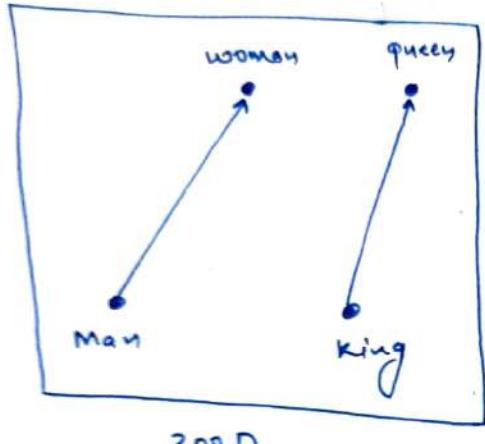
$$e_{\text{man}} - e_{\text{woman}} \approx \begin{bmatrix} -2 \\ 0 \\ 0 \end{bmatrix}$$

king → ? (Queen)

$$e_{\text{king}} - e_{\text{queen}} \approx \begin{bmatrix} -2 \\ 0 \\ 0 \end{bmatrix}$$

$$e_{\text{man}} - e_{\text{woman}} \approx e_{\text{king}} - e_{\text{?}}$$

\hookrightarrow queen

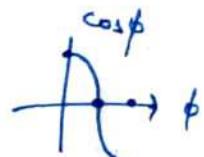
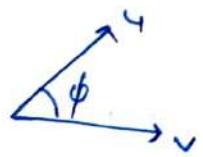


$$e_{\text{man}} - e_{\text{woman}} \approx e_{\text{king}} - e_{?} (e_w)$$

$$\operatorname{argmax}_w \text{Sim}(e_w, e_{\text{king}} - e_{\text{man}} + e_{\text{woman}}) \\ (30 - 75\%)$$

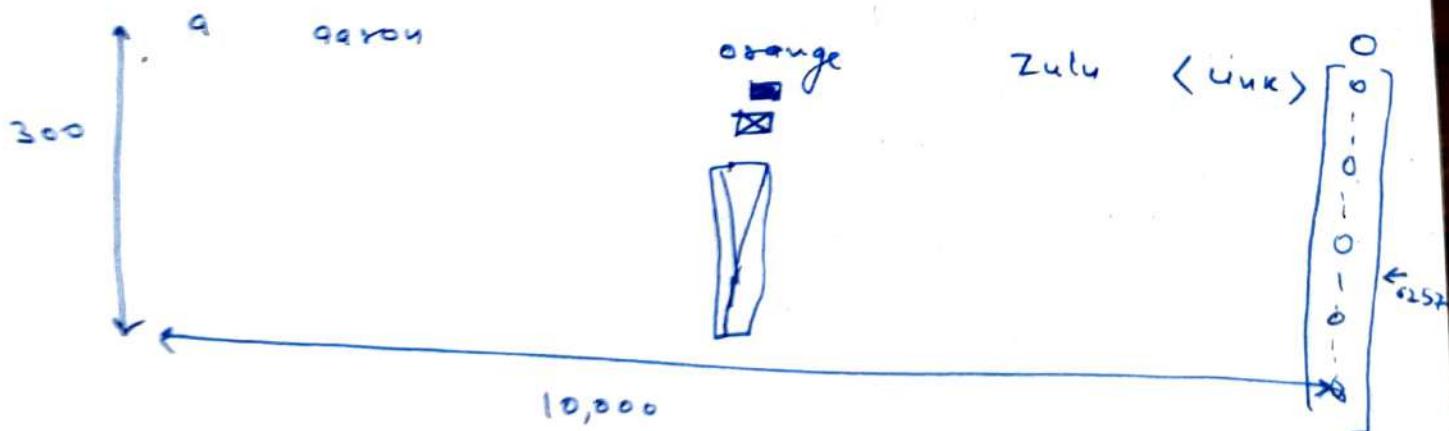
Cosine similarity

$$\text{sim}(u, v) = \frac{u^T v}{\|u\|_2 \|v\|_2}$$



Embedding matrix

$$\|u - v\|^2$$

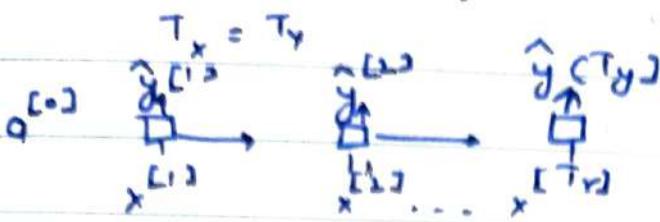


$$E \cdot o_{6257} = \begin{bmatrix} \text{[square]} \\ \text{[square]} \\ \text{[triangle]} \end{bmatrix}_{(300, 1)} = e_{6257} \quad e_w = \text{embedding for } w$$

$E \cdot o_j = e_j$
embedding for word j

Recurrent Neural Network (RNN)

Different types of RNN



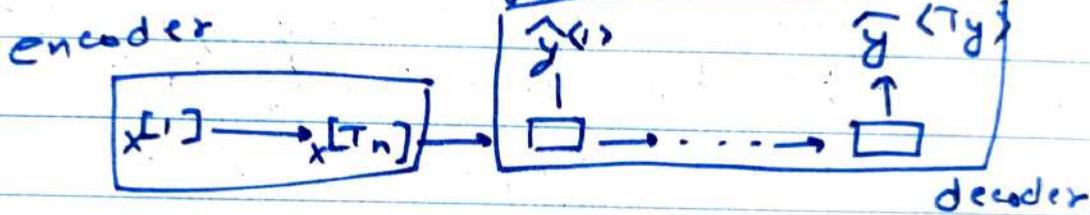
Many - to - many

Many - to - one e.g. Movie review

One - to - one

One - to - many e.g. Music generation

Many - to - many Input and output dimension is different



→ Language Model and sequence operation

What is language modelling?

Speech recognition

Language modelling with an RNN

Training set

- Tokenize

$$x^{(t)} = y^{(t-1)}$$

(unk)

$$a^{(1)} = \vec{o} \rightarrow \boxed{a^{(1)}} \rightarrow$$

$x^{(1)} = \vec{o}$

<EOS> End of sentence
P(cats average) / P(cats average)

$$\begin{aligned} & \hat{y}^{(2)} \quad \hat{y}^{(3)} \quad \hat{y}^{(4)} \\ & \uparrow \quad \uparrow \quad \uparrow \\ & \boxed{a^{(1)}} \rightarrow \boxed{a^{(2)}} \rightarrow \boxed{a^{(3)}} \rightarrow \boxed{a^{(4)}} \rightarrow \boxed{a^{(5)}} \\ & x^{(2)} = y^{(1)} \quad x^{(3)} = y^{(2)} \quad x^{(4)} = y^{(3)} \quad x^{(5)} = y^{(4)} \\ & P(\text{EOS}) = - \end{aligned}$$

$$L(\hat{y}^{(t)}, y^{(t)}) = - \sum_i y_i^{(t)} \log \hat{y}_i^{(t)}$$

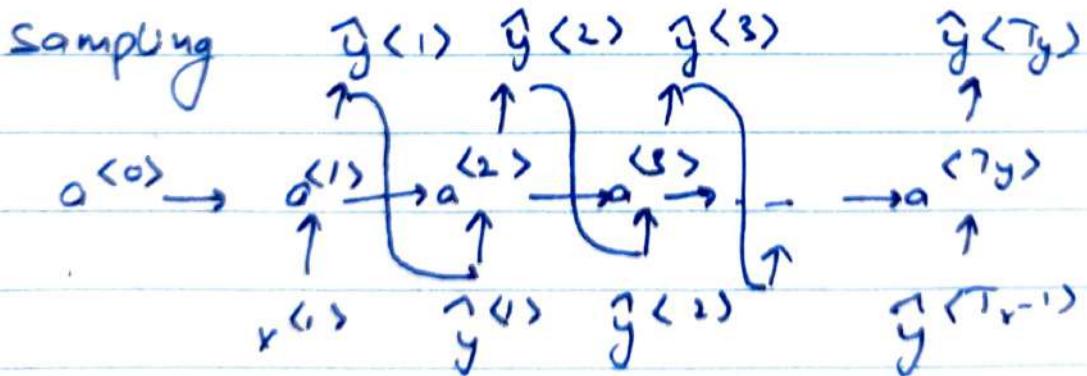
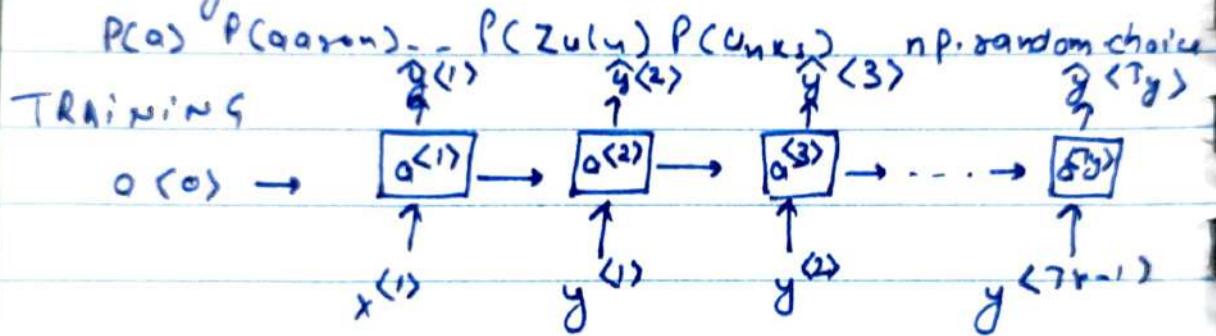
$$L = \sum_t L^{(t)}(\hat{y}^{(t)}, y^{(t)})$$

$$\begin{aligned} P(y^{(1)}, y^{(2)}, y^{(3)}) \\ P(y^{(1)}) P(y^{(2)} | y^{(1)}) \\ P(y^{(3)} | y^{(1)}, y^{(2)}) \end{aligned}$$

→ Sampling novel sequences

Sampling a sequence from a trained RNN

Training
Sampling



→ Vanishing gradients with RNN

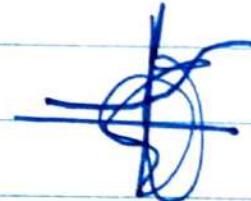
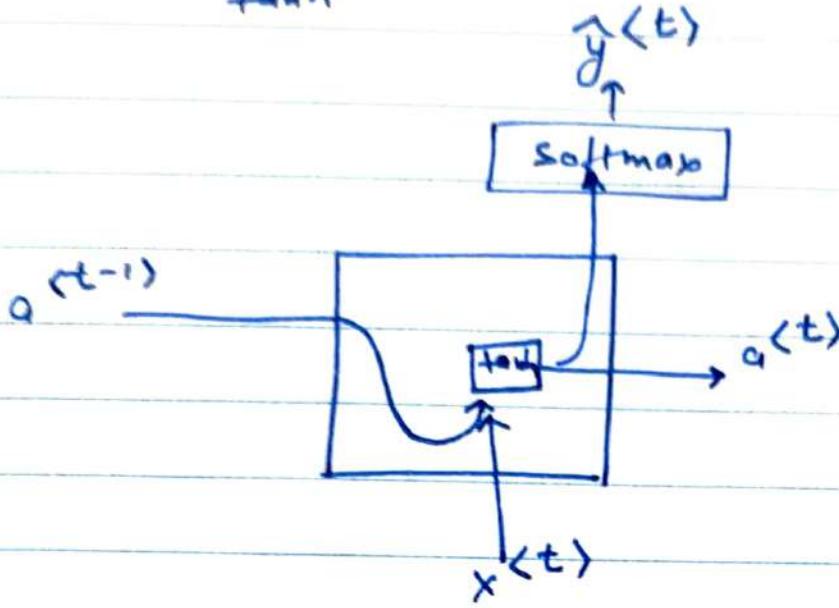
The cells, which already are ---, were full

The cells, which already are ---, were full
Non gradient clipping

→ Gated Recurrent Unit GRU

$$a^{(t)} = g(w_a[a^{(t-1)}, x^{(t)}] + b_a)$$

$\downarrow \tanh$



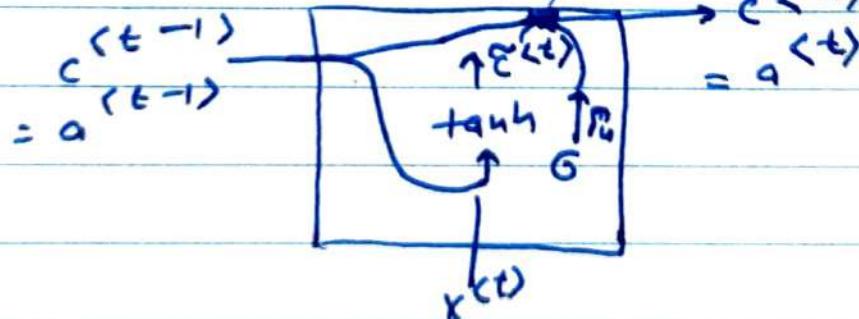
c = memory cell

$$c^{(t)} = \tilde{a}^{(t)}$$

$$\tilde{c}^{(t)} = \tanh(w_c[c^{(t-1)}, x^{(t)}] + b_c)$$

$$\begin{aligned}\Gamma_u &= \sigma(w_u[c^{(t-1)}, x^{(t)}] + b_u) \\ c^{(t)} &= \Gamma_u * \tilde{c}^{(t)} + (1 - \Gamma_u) * c^{(t-1)}\end{aligned}$$

softmax $\rightarrow \hat{y}^{(t)}$



→ LSTM (Long short term memory) unit

$$\tilde{c}^{(t)} = \tanh(w_c[a^{(t-1)}, x^{(t)}] + b_c)$$

$$\text{update } \Gamma_u = \sigma(w_u[a^{(t-1)}, x^{(t)}] + b_u)$$

$$\text{forget } \Gamma_f = \sigma(w_f[a^{(t-1)}, x^{(t)}] + b_f)$$

$$\text{output } \Gamma_o = \sigma(w_o[a^{(t-1)}, x^{(t)}] + b_o)$$

$$c^{(t)} = \Gamma_u * \tilde{c}^{(t)} + \Gamma_f * c^{(t-1)}$$

$$a^{(t)} = \Gamma_o * c^{(t)}$$

(peephole connection)

→ Bidirectional RNN (BRNN)

Getting information from the future.

He said "Teddy bears are on sale!"

He said "Teddy Roosevelt was a great
President!"

→ Deep RNNs

→ Word Representation

Featureized representation: word embeddings

Visualizing word embeddings

→ Using word embeddings

Relation to face encoding (embedding)

→ Properties of word embeddings.

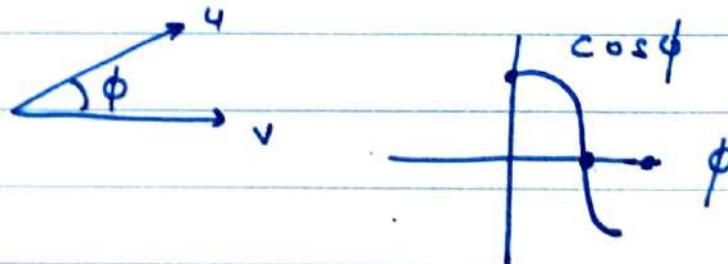
$$e_{\text{man}} - e_{\text{woman}} \approx \begin{bmatrix} -2 \\ 0 \\ 0 \\ 0 \end{bmatrix}$$

$$e_{\text{king}} - e_{\text{queen}} \approx \begin{bmatrix} -2 \\ 0 \\ 0 \\ 0 \end{bmatrix}$$

$$e_{\text{man}} - e_{\text{woman}} \approx e_{\text{king}} - e_{\text{queen}}$$

Linguistic regularities in continuous space word representations. Mikolov et al 2013
final word w: arg max_w sim($e_w, e_{\text{king}} - e_{\text{man}}$)
cosine similarity

$$\text{sim}(u, v) = \frac{u^T v}{\|u\|_2 \|v\|_2}$$

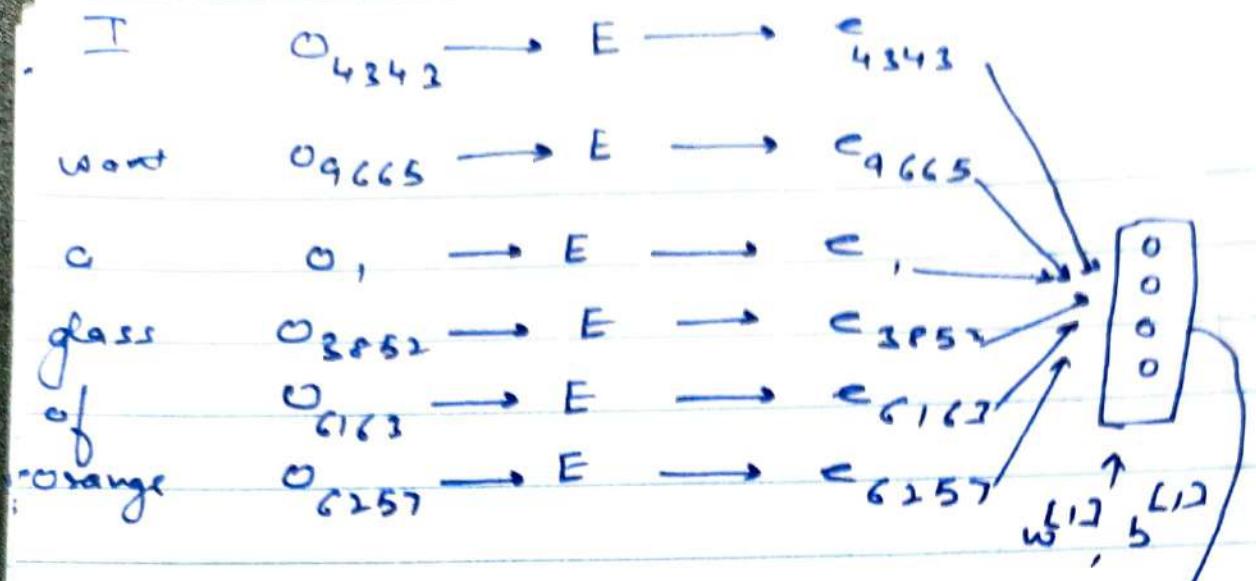


$$\|u - v\|^2$$

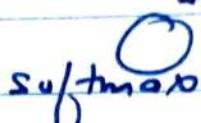
→ Embedding matrix

→ Learning word embedding
neural language model

I want a glass of orange ?
4343 9665 1 3852 6162 6257



use four word history



Bengio et al., 2003, A neural probabilistic language model.

with my cereal

4 word on left & right

Last word orange?

Near by 1 word glass ?

→ Word2Vec

Skip-grams
Model

$$\text{Vocab size} = 10,000 \text{K}$$

$$O_C \rightarrow E \rightarrow e_C \rightarrow O \rightarrow \hat{y}$$

$e_C = E_{O_C}$ soft max

$$\text{softmax: } p(t|c) = \frac{e^{\theta_t^T c_0}}{\sum_{j=1}^{10,000} e^{\theta_j^T c_0}}$$

θ_t = parameter associated with output t

$10,000$

$$L(\hat{y}, y) = - \sum_{i=1}^{10,000} y_i \log \hat{y}_i$$

$$y = \begin{bmatrix} 0 \\ \vdots \\ 1 \\ \vdots \\ 0 \end{bmatrix} \leftarrow \text{4834}$$

Problems with softmax

classification

hierarchical

$$\log |V|$$

→ Negative sampling

Defining a new learning problem
context word target

10,000 binary classification $K+1$

Selecting negative examples

$$P(w_i) = \frac{f(w_i)^{3/4}}{\sum_{j=1}^{10,000} f(w_j)^{3/4}} = \frac{1}{|V|}$$

→ GloVe word vectors

global vector for word representation

$$x_{ij} = x_{ji} \leftarrow$$

Pennington et. al., 2014. GloVe vector for word representation.

Model

$$\text{minimize} \sum_{i=1}^{1000} \sum_{j=1}^{1000} f(x_{ij}) (\theta_i^T e_j + b_i + b_j - \log x_{ij})$$

↑ ↓
t c

↑
o₂
o₁

↑
e_c

↑
weighting " $\theta_i^T e_c$ "

$$f(x_{ij}) = 0$$

$$\text{if } x_{ij} = 0$$

$$\text{" } 0/\log 0 \text{ " } = 0$$

$$e_w^{(\text{final})} = \frac{e_w + \theta_w}{2}$$

A note on the featurization view of word embeddings

$$\text{minimize} \sum_{i=1}^{1000} \sum_{j=1}^{1000} f(x_{ij}) (\theta_i^T e_j + b_i + b_j - \log x_{ij})$$

→ Sentimental classification

Simple sentiment classification model

The dessert is excellent

8928 2468 4694 3180

The $\theta_{8928} \rightarrow E \rightarrow e_{8928}$

desert $\theta_{2468} \rightarrow E \rightarrow e_{2468} \rightarrow \text{Avg}$

is $\theta_{4694} \rightarrow E \rightarrow e_{4694} \rightarrow 300D$

excellent $\theta_{3180} \rightarrow E \rightarrow e_{3180} \rightarrow 300D$

100B words $\hat{y} \leftarrow \text{s. if max } 1-5$

→ Debiasing word embeddings

Bulukbasi et al 2016 man is to computer programmer as woman is to homemaker

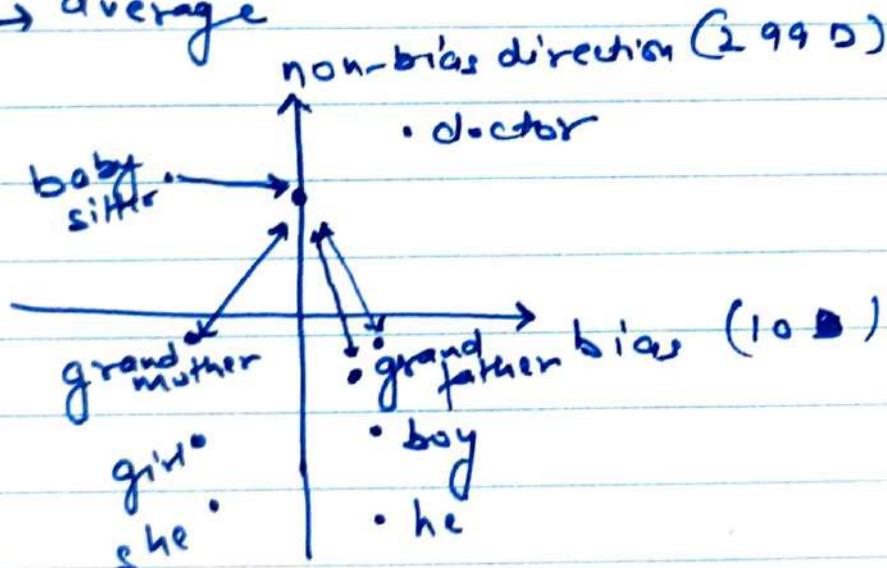
? Debiasing word embeddings

Addressing bias in word embeddings

1. Identify bias directions

$$\left\{ \begin{array}{l} e_{\text{he}} - e_{\text{she}} \\ e_{\text{male}} - e_{\text{female}} \end{array} \right.$$

average



2. Neutralize: for every word that is not definitional, project to get rid of bias

3. Equalize pairs

grandmother
girl

grandfather
boy

Bulukbasi et al., 2016

→ Basic models

Sequence to sequence model

$x^{(1)}$ $x^{(2)}$ $x^{(3)}$ $x^{(4)}$ $x^{(5)}$

Jane visite l'Afrique en septembre

Jane is visiting Africa in september

$y^{(1)}$ $y^{(2)}$ $y^{(3)}$ $y^{(4)}$ $y^{(5)}$ $y^{(6)}$ $y^{(7)}$

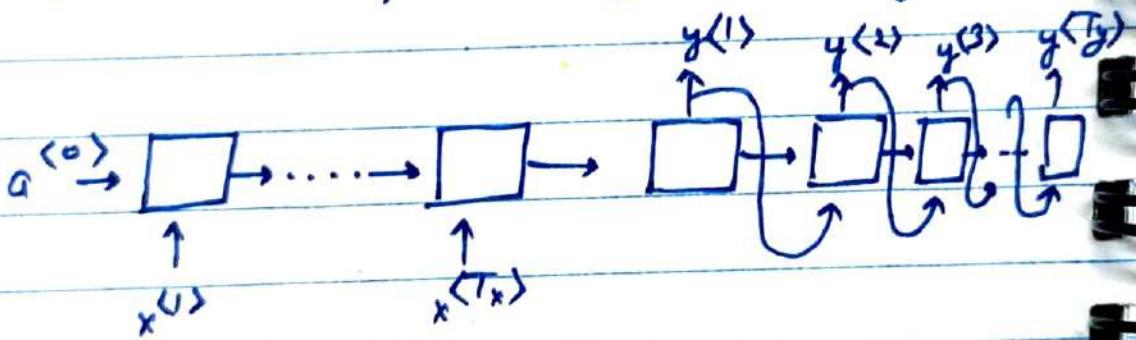


Image captioning

A cat sitting on a chair

→ Picking the most likely sentence

Machine translation as building

a conditional language model

Finding the most likely translation
 $P(y^{(1)}, \dots, y^{(T_y)} | x)$

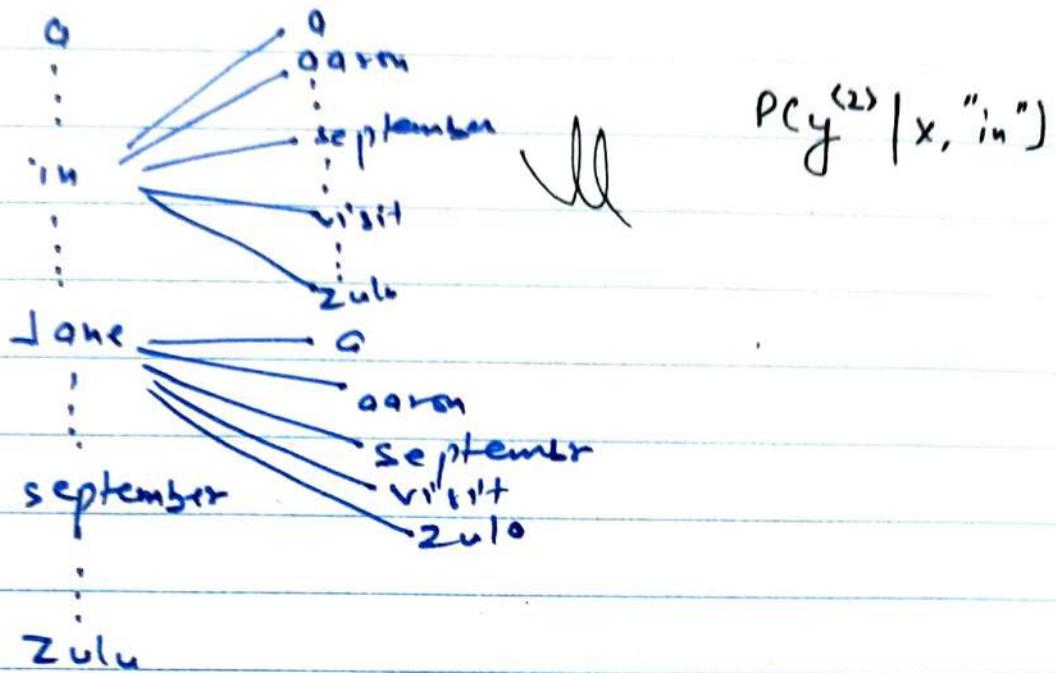
arg max _{$y^{(1)}, \dots, y^{(T_y)}$}

$P(y^{(1)}, \dots, y^{(T_y)} | x)$

why not a greed search?

→ (Basic search) Beam search
 beam width $P(y^{(1)}|x)$

step 1 step 2



$$a^{(0)} \rightarrow \boxed{\quad} \rightarrow \boxed{\quad}$$

Beam search ($B = 3$) (Possibly 3 at a time)

- $P(y^{(1)}, y^{(2)}|x)$
- Refinements to beam search
- Length normalization

$$\arg \max_y \prod_{t=1}^T P(y^{(t)}|x, y^{(1)}, \dots, y^{(t-1)})$$

$$\arg \max_y \sum_{t=1}^T \log P(y^{(t)}|x, y^{(1)}, \dots, y^{(t-1)})$$

$$+ \frac{\log P(y|x)}{\log P(y|x)}$$

$$\frac{1}{Ty^\alpha} \sum_{t=1}^{Ty} \log P(y^{(t)} | x, y^{(1)}, \dots, y^{(t-1)})$$

$\alpha = 0.7 \quad \alpha = 1$
 $\alpha = 0$

Beam width B?

Large B better result, slower
 small B worse result, faster

$1, 3 \rightarrow 10, 100, 1000 \rightarrow 3000$

Unlike exact search algorithm like BFS or DFS, Beam search run faster but is not guaranteed to find exact max. for arg max $P(y|x)$.

→ Error analysis in Beam search

- Bleuscore

Papineni et al 2002, A method for automatic evaluation of machine translation

Bilingual evaluation understanding

Precision modified precision

Blue Bleu score on bigrams

Bleu score on unigrams

$$P_i = \frac{\sum_{\text{unigamey}} \text{count}_{\text{clip}}(\text{unigram})}{\sum_{\text{Unigamey}} \text{count}_{\text{unigamey}}}$$

$$P_1, P_2, = 1.0$$

Bleu details

P_n = Bleu score on n-grams only

combined Bleu score

$$\text{BP cap} \left(\frac{1}{\Psi} \sum_{n=1}^{\Psi} P_n \right)$$

$$\text{BP} = \begin{cases} 1 & \text{if MT_output_length} \geq \text{reference_output_length} \\ \exp \left(\frac{1 - \text{MT_output_length}}{\text{reference_output_length}} \right) & \text{otherwise} \end{cases}$$

→ Attention model intuition

The problem of long sequences

Bahdanau et al 2014 neural machine trans-

-lation by jointly learning to align and
translate

→ Attention model

$$c^{(1)} = \sum_{t'} \alpha^{(1, t')} a^{(t')}$$

$$a^{(t')} = (\gamma^{(t')}, \delta^{(t')})$$

computing attention $\alpha^{(t, t')}$

$\alpha^{(t, t')}$ = amount of attention $y^{(t)}$
should pay to $a^{(t')}$

$$\alpha(t, t') = \frac{\exp \langle t, t' \rangle}{\sum_{t''=1}^T \exp \langle t, t' \rangle}$$

Xu et.al., 2015, show attend and tell:
neural image caption generation with
visual attention

Attention examples

→ speech recognition

x
audio clip

y
transcript

"The quick brown fox"

CTC cost for speech recognition
(Connectionist temporal classification)

Graves et al 2006 Connectionist temporal
classification labeling unsegmented
sequence data with RNN

Basic rule: collapse repeated
characters not separated by "blank"