

Courses 5 : Sequence models.

[04.02.19]. v.v.v.y Balme.

WEEK 1 : Recurrent neural networks. (v.v.y week for me).

L-1 why sequence models. [10.20 AM 09.01.19]

Example of sequence data and application.

1. Speech recognition. (consequential)
2. Music generation.
3. Sentiment classification. (v.v.y)
4. DNA sequence analysis.
5. machine translation. ✓
6. video activity recognition. [video picture card 200 09.01.19] (25 for
activity detection)
7. Name entity recognition. —

L-2 Notation

[10.30 AM 09.01.19]

Motivating example.

Name entity recognition :- (where the people's name in a sentence
used in search engines).

→ Indexing different types of names from text (people, name, currency name,
city name, company name etc.)

X: Harry] potter and Hermione Granger invented a new spell.
<1> x <2> x <+>
Y: J <1> - J <2> J <+>

T_x = length of input sequence.

T_f = length of output sequence.

$\{x^{(i)} \langle t \rangle\} = t^{\text{th}}$ element of ^{ith} antitraining example i .

$\{T_x^{(i)}\}$ = input sequence length of in example =

$\{y^{(i)} \in t\}$ = element of i th training example in output

$T_y^{(i)}$ = output sequence length of i th example

How to represent words of a sentence.

X: Harry Potter and Hermoine Granger invented a new spell-
(9) *z*

Vocabulary

$x^{(1)}$

$x^{(2)}$

\dots

$x^{(n)}$

$\left[\begin{matrix} a \\ \text{caron} \\ \vdots \\ \text{and} \\ \vdots \\ \text{harry} \\ \vdots \\ \text{POTEN} \\ \vdots \\ \text{zallu} \end{matrix} \right]$

$\left[\begin{matrix} o \\ \text{o} \\ \vdots \\ \text{o} \\ \text{l} \\ \text{o} \\ \vdots \\ \text{o} \\ \vdots \\ \text{o} \end{matrix} \right] \xleftarrow{\leftarrow 407S}$

$\left[\begin{matrix} e \\ \text{o} \\ \vdots \\ \text{o} \\ \text{l} \\ \text{o} \\ \vdots \\ \text{o} \end{matrix} \right]$

$\left[\begin{matrix} \text{o} \\ \text{o} \\ \vdots \\ \text{o} \\ \text{o} \\ \vdots \\ \text{o} \end{matrix} \right]$

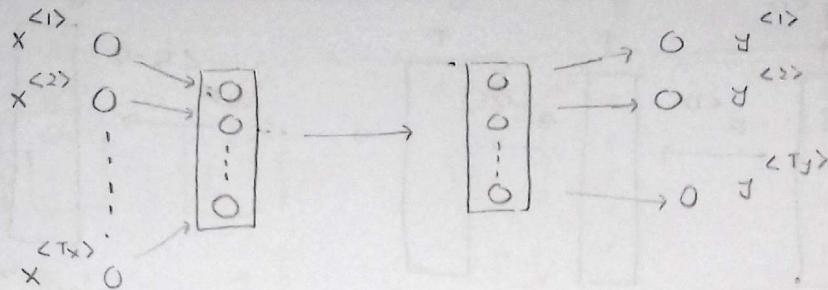
** one hot vector to represent words

[guitar word 24 days] vocabulary size 28 vectors create 210].

** we can create a vocabulary or dictionary by taking
10000 / 20,000 common words in our training set.

L-3 Recurrent neural network model. [10:55 AM 09-02-19]

why not a standard network?



Problems :-

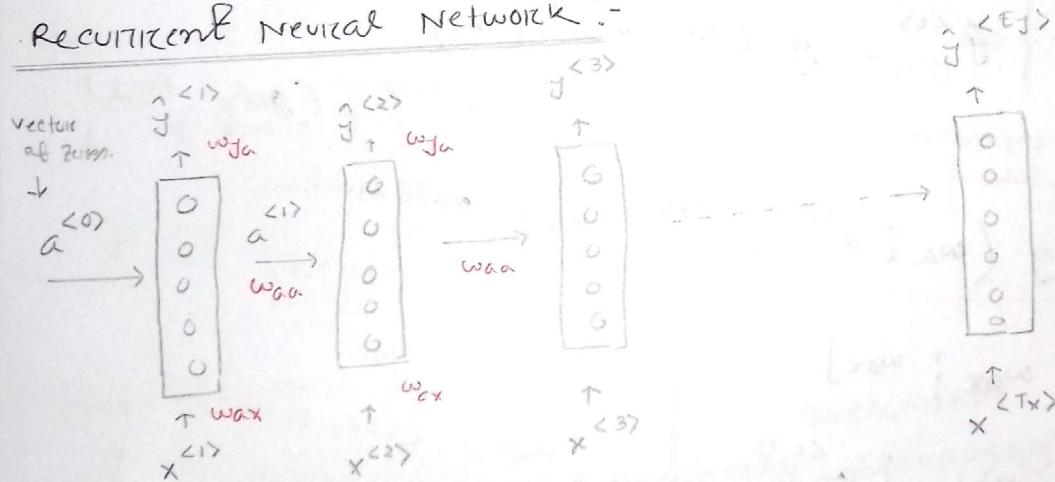
1. Inputs, outputs can be different lengths in different examples.

→ 2. Doesn't chance features across different positions of a text.

(If we have a vocabulary size of 20,000 and max word size = 100
then number of parameters will be very large)

→ This problem can be solved by recurrent neural network //

Recurrent Neural Network :-

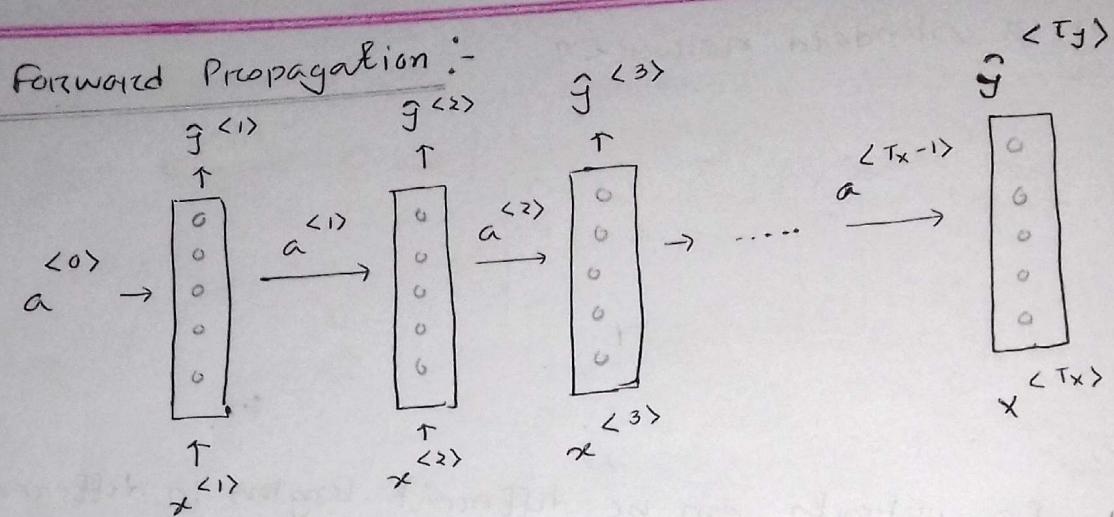


[** one weakness of RNN is that it uses only the earlier information

[Sometimes we need to process from earlier as well as later words in that case performance of RNN falls.]
(slide 9 example chart will be)

(Ears)

Forward Propagation :-



$$\begin{aligned}
 a^{<0>} &= \bar{x} \\
 a^{<1>} &= g(w_{aa} a^{<0>} + w_{ax} x^{<1>} + b_a) \leftarrow \tanh/\text{relu} \\
 g^{<1>} &= g(w_{ya} a^{<1>} + b_y) \leftarrow \text{sigmoid} \\
 &\vdots \\
 a^{<t>} &= g(w_{aa} a^{<t-1>} + w_{ax} x^{<t>} + b_a) \\
 g^{<t>} &= g(w_{ya} a^{<t>} + b_y).
 \end{aligned}$$

This equation can be changed slightly \Rightarrow (General Form)

$$a^{<t>} = g(w_a [a^{<t-1>} \text{, } x^{<t>}] + b_a) \quad \text{main equation.}$$

Let $w_{aa} = (100, 100)$

calculation -

$$w_a = [w_{aa} \text{ ; } w_{ax}]$$

$$a^{<t-1>} = 100$$

$$[a^{<t-1>} \text{, } x^{<t>}] = \begin{bmatrix} a^{<t-1>} \\ x^{<t>} \end{bmatrix}$$

$$w_{ax} = (100, 1000)$$

$$x^{<t>} = 10000$$

$$so \quad w_a = (100, 10100)$$

$$\begin{bmatrix} a^{<t-1>} \\ x^{<t>} \end{bmatrix} \begin{bmatrix} 100 \\ 10000 \end{bmatrix} = 10100$$

$$\therefore [w_{aa} \text{ ; } w_{ax}] \begin{bmatrix} a^{<t-1>} \\ x^{<t>} \end{bmatrix}$$

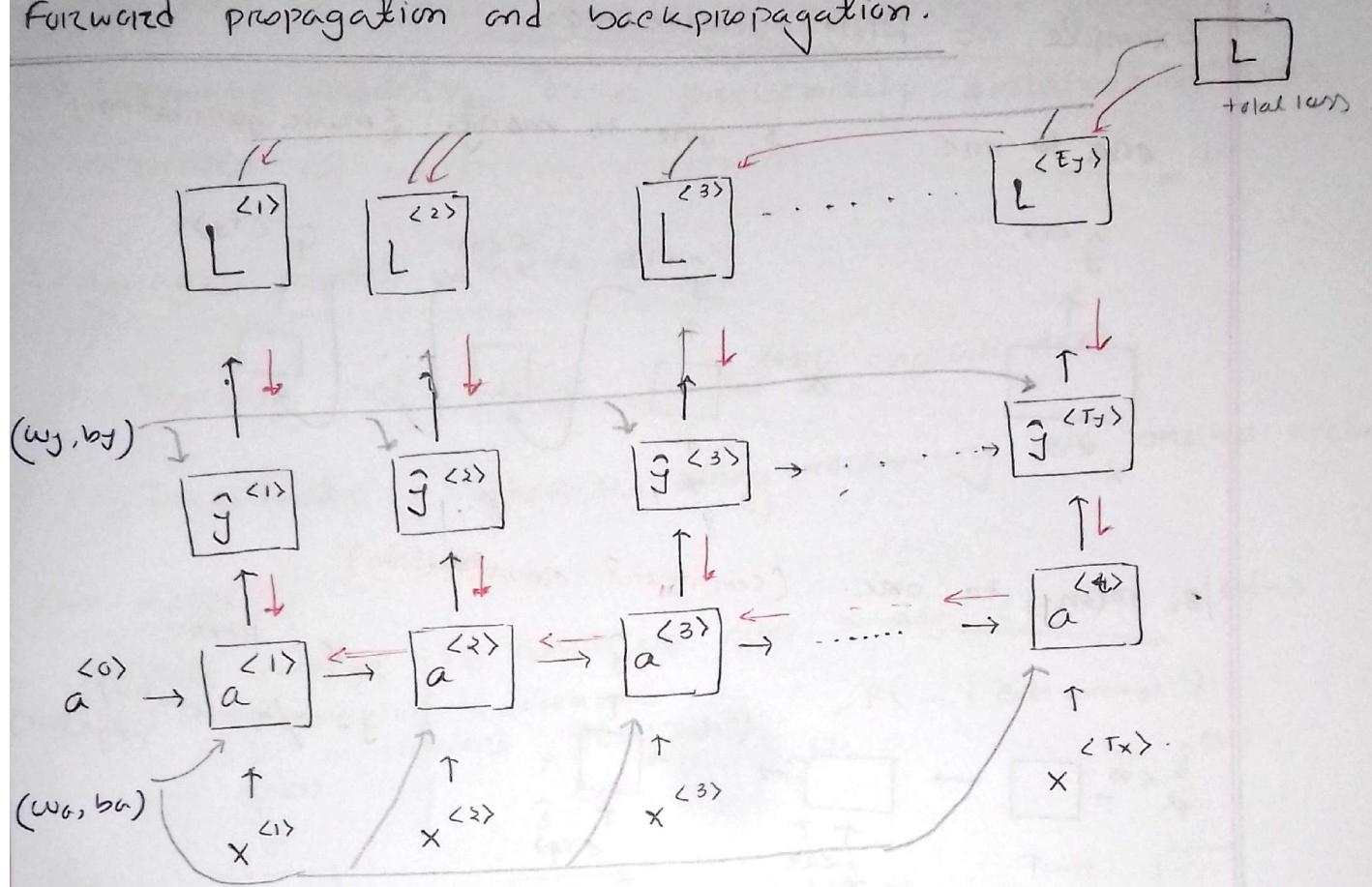
$$= w_{aa} a^{<t-1>} + w_{ax} x^{<t>}$$

$$g^{<t>} = g(w_y a^{<t>} + b_y)$$

m.E.

L-04 Backpropagation through time. [11.35 AM 09.02.19]

Forward propagation and backpropagation.



$$L^{(t)}(\hat{y}^{(t)}, y^{(t)}) = -y^{(t)} \log \hat{y}^{(t)} - (1-y^{(t)}) \log (1-\hat{y}^{(t)})$$

$$\mathcal{L}(\hat{y}, y) = \sum_{t=1}^{T_f} L^{(t)}(\hat{y}^{(t)}, y^{(t)}) \quad [\text{total loss}]$$

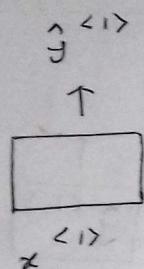
Backpropagation through time.

[Programming framework automatically implements backpropagation no we don't have to worry much about it]

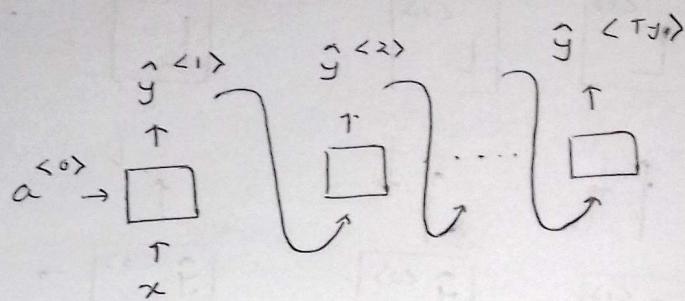
L-5 Different types of RNN [11.50 Am 04-01-17]

Example of RNN architectures.

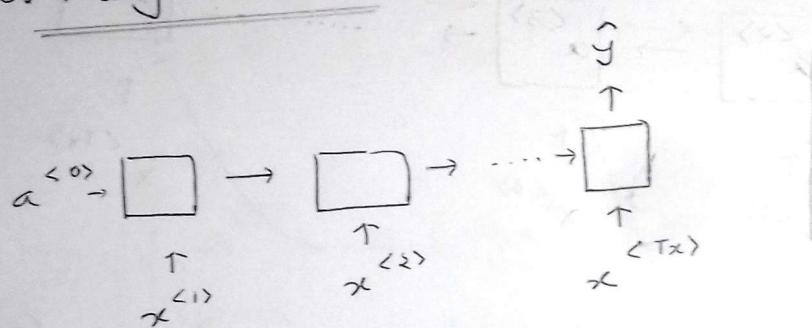
1. one to one.



2. one to many. (music generation)



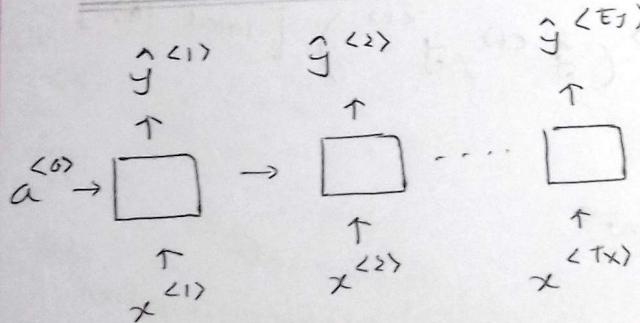
3. Many to one. (sentiment classification)



$$x = \text{text} \\ y = 0/1 \text{ (positive/negative)}$$

4. Many to many

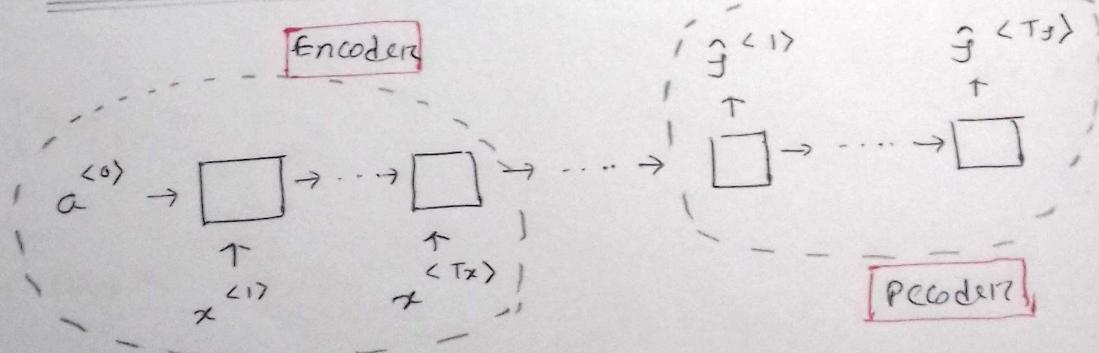
(Name entity recognition) ($T_x = T_y$)



$$T_x = T_y$$

5. Many to many

($T_x \neq T_y$) (machine translation)



L-6 Language Model and sequence. [12.15 PM 04.01.17]

[Most basic and important task in natural language processing]

** Language modeling is a probability distribution over sequences of words or sentences.

Language model with an RNN :-

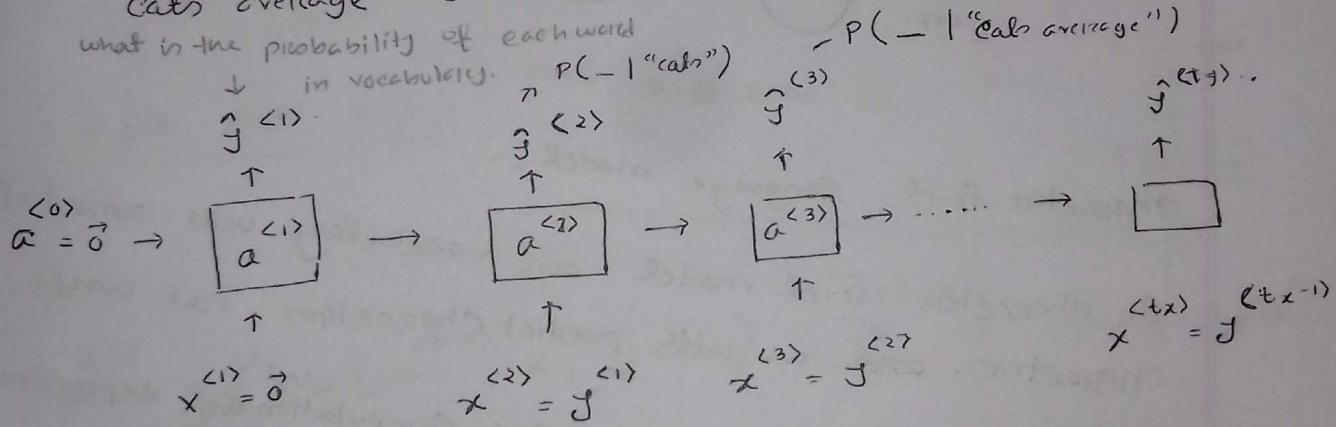
1. Training set : large corpus of english text.

2. Tokenize : Tokenize using vocabulary into one-hot vectors.

RNN model :-

Cats average 15 hours of sleep a day. <EOS> | unk>

what is the probability of each word



** Each step of RNN will look at some set of preceding words and try to predict next word.]

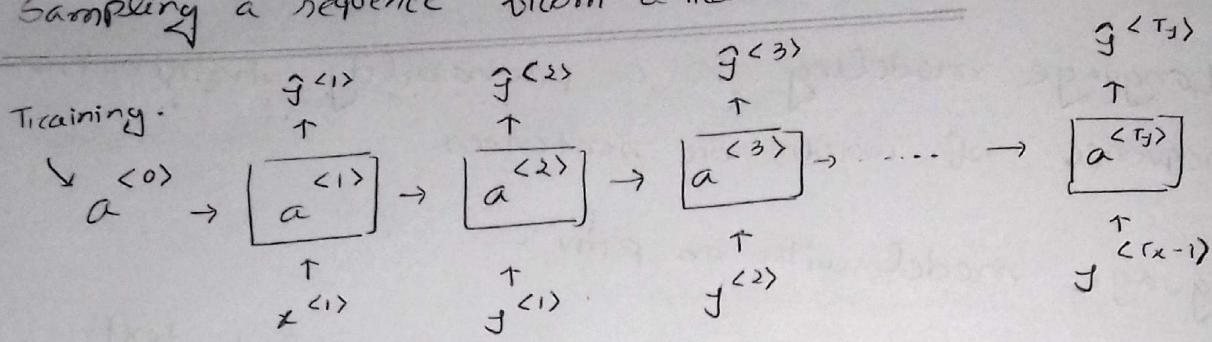
$$\text{Loss function } L = - \sum_i y_i \log \hat{y}_i^{(t)}.$$

$$L = \sum_{t=1} L(\hat{y}^{(t)}, y^{(t)})$$

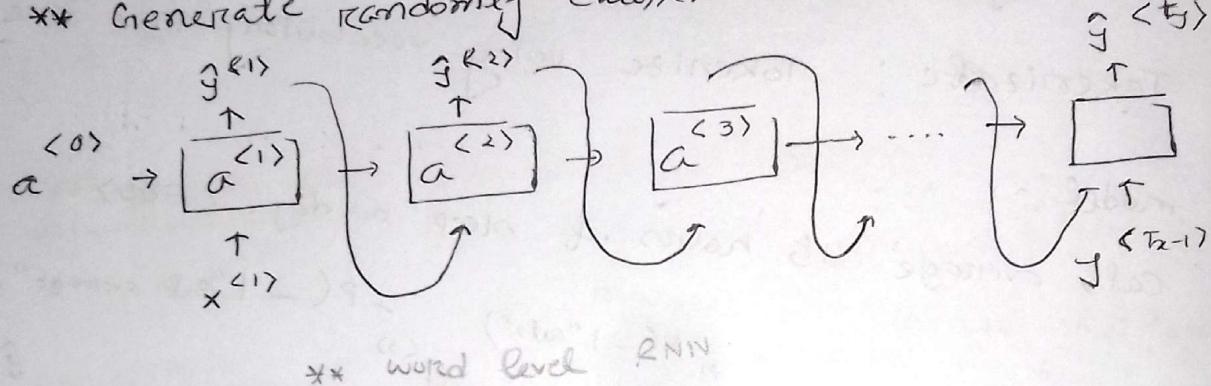
(Programming exercise - J implement
loss)

L-7 Sampling Novel sequences. [4.30 PM 09.02.12]

Sampling a sequence from a trained RNN:-



** Generate randomly chosen sentence from RNN model.



Character-level language model :-

In character-level model our vocabulary will consist of characters. and we will predict characters not words.

** [character-level models are computationally expensive and they are not as good as word level]

[Programming exercise -] Language model implement GAN RNN] (slide 04 2nd part 1st example their code)

L-8 Vanishing gradients with RNNs. [4.50 pm 09.02.19]

[Basic RNN has vanishing gradients problem].

** [In language model there can be very long term dependencies]

(slide 10 RNN-02 vanishing gradient issue) [2:45]

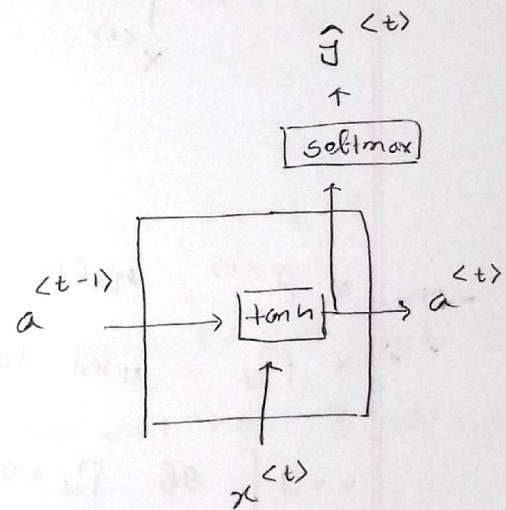
(Q3 Problem 2st part 8th) method concentrate on 2nd

~~L-9~~ Gated Recurrent Unit (GRU) [7.00 pm 09.02.19]

[Much better in capturing long range connections and helps in vanishing gradient problem]

Simple RNN Unit :-

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



2nd paper on it.

on the properties of neural machine translation.

[1] on the properties of recurrent neural networks

[2] Empirical Evaluation of Gated Recurrent Networks

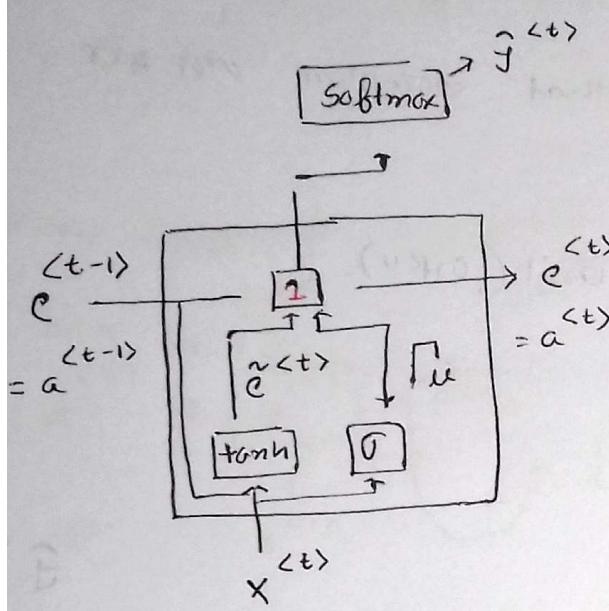
on sequence modeling.

* 2nd paper 2nd part 2nd

GIRU (simplified).

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

c = memory cell.



$$c^{t-1} = a$$

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

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

↓ "update" (gate) elementwise.

$$c^{t-1} = F_u * \tilde{c}^{t-1} + (1 - F_u) * c^{t-1} \quad \text{--- (1)}$$

** gate is update value.

* \tilde{c}^{t-1} holds the value of the memory.

* F_u used to update a particular memory cell.

v.v.y [If $F_u = 0$ then $c^{t-1} = \tilde{c}^{t-1}$ that helps to reduce vanishing gradient problem and learn dependencies through longer network.]

** Dimension of c^{t-1} , \tilde{c}^{t-1} , F_u are same.

Full GRU :-

$$\begin{aligned}\tilde{h} &\leftarrow \tilde{c}^{(t)} = \tanh \left(w_c [\Gamma_{12} * e^{(t-1)}, x^{(t)}] + b_c \right) \quad \begin{array}{l} \text{new candidate} \\ \text{gate memory cell} \end{array} \\ \text{ee} &\leftarrow \Gamma_u = \sigma \left(w_u [e^{(t-1)}, x^{(t)}] + b_u \right) \\ r &\leftarrow \Gamma_r = \sigma \left(w_r [e^{(t-1)}, x^{(t)}] + b_r \right) \\ h &\leftarrow c^{(t)} = \Gamma_u * \tilde{c}^{(t)} + (1 - \Gamma_u) * c^{(t-1)}\end{aligned}$$

** This helps to catch long range dependencies and makes RNN better.

L-10. Long Short Term memory. [7.45 pm 04.02.19]

[LSTM allows the errors to back propagate through unlimited number of time steps. A LSTM unit has 3 gates,

1. input gate.

2. output gate.

3. forget gate.

It calculates the hidden state by taking combination of these 3 gates.]

[It's unique mechanism helps to overcome vanishing and exploding gradient problem]

[LSTM is more powerful than GRU]

GRU OR LSTM
- can use over GRU
- more powerfull

GRU

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

$$r_u = \sigma(w_u [r_u^{(t-1)}, x^{(t)}] + b_u)$$

$$r_{uc} = \sigma(w_{uc} [c^{(t-1)}, x^{(t)}] + b_{uc})$$

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

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

LSTM

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

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

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

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

$$c^{(t)} = r_u * \tilde{c}^{(t)} + r_f * c^{(t-1)}$$

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

LSTM in pictures :- (understanding LSTM blog post code (part 1)).

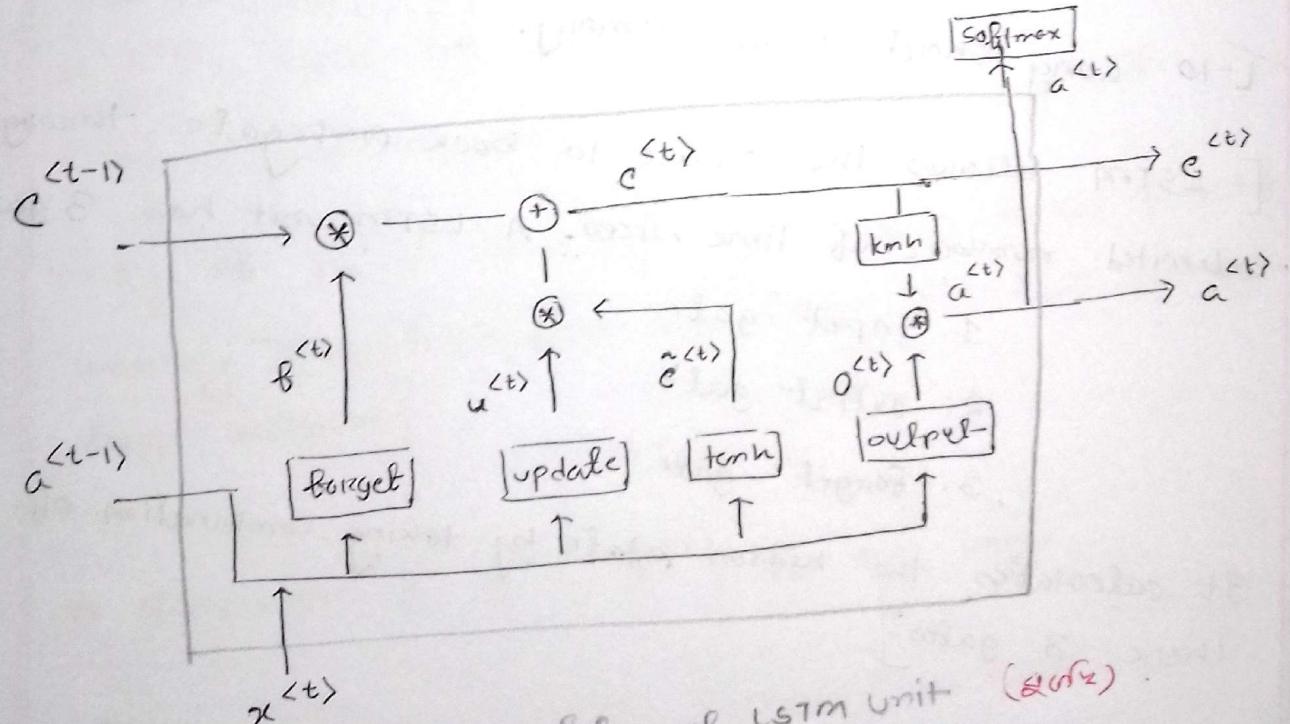


Fig:- Graphical representation of LSTM unit (part 1).

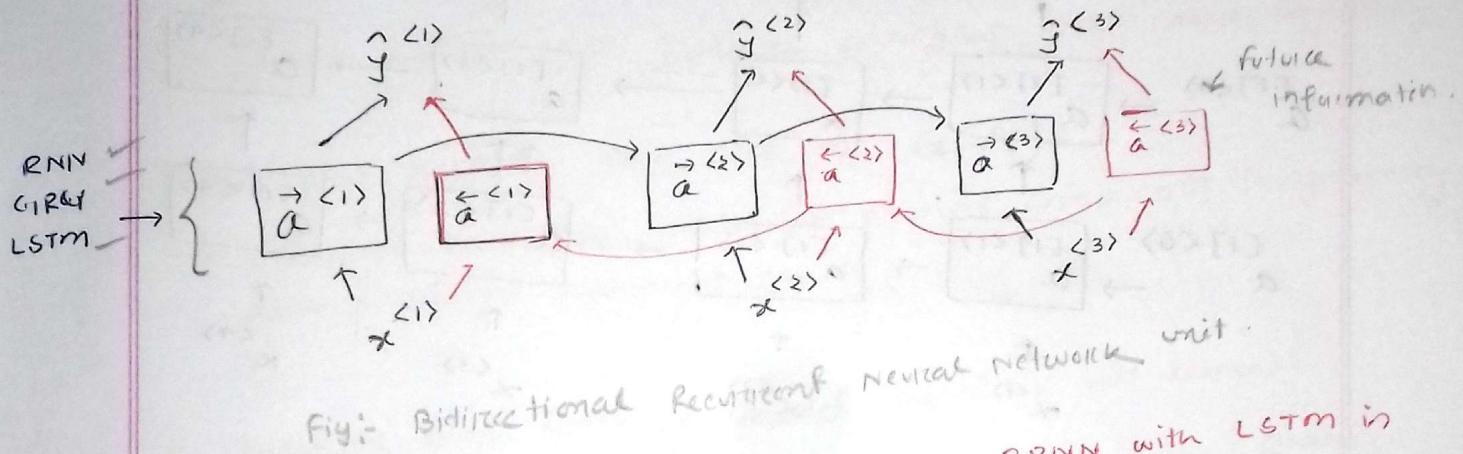
** LSTM unit is good in memorizing previous unit values.

** peephole connection → used to effect gate output.

L-11 Bidirectional RNN [8.13 PM 09.01.19] (BRNN)

[Bidirectional RNN takes information from previous as well as future words]

$$\hat{y}^{(t)} = g(\omega_y [\vec{a}^{(t)}, \hat{a}^{(t)}] + b_y).$$



**** [For natural language processing a pretty reasonable thing to try.]**

Disadvantages:-

- 1. Need entire sequence of data before prediction anywhere.
(For optical time speech recognition it is not very good.)
- But for natural language processing where we have the whole data it is really good.

Deep RNN example.

(may have deep recurrent network).

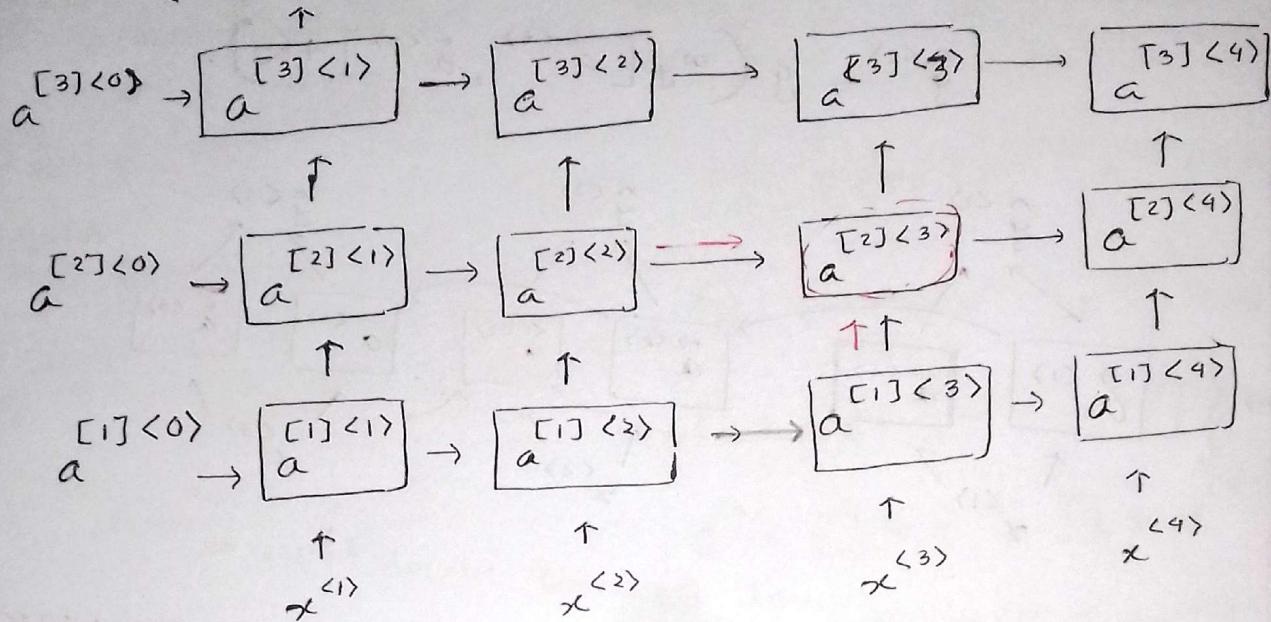


Fig:- A deep RNN architecture.

Lect.

$$a^{[2]<3>} = g \left(w_a \begin{bmatrix} a^{[2]<2>} \\ a^{[1]<3>} \end{bmatrix} + b_a^{[2]} \right).$$

** [This week was very much important, learn about Basic RNN, GRU, LSTM, BRNN, and deepen RNNs].

→ go to related papers after insha-Allah soon.

Programming Assignment :-

1. Building Recurrent Neural Network - step by step.

** यानि का main target कर input preprocessing कर.

उदाहरण से small easily process कर सकते हैं.

** [paper 2012-2015 convolutional neural network का क्या कर]

[Basic RNN & LSTM block का क्या कर]

** (Backward propagation implement करने का क्या कर Programming framework का implement करने)

2. character level language model. (what we will learn).

- How to store text data for processing using an RNN.

- How to synthesize data by sampling predictions at each time step and passing it to the next RNN-cell unit.

- How to build a character-level text generation recurrent neural network.

- why clipping the gradients is important.

1. [LSTM text generation by keras team].

2. [karpathy's blog post].

** यहाँ करना क्या कर
code को किसी implement
करने के लिए क्या करें
2015.

[shakir's द्वारा सेट generation का]

[it's an idea under previous text generation
I use what if]

WEEK-2 Natural language Processing and word embeddings.

[07-02-19]

P1: Introduction to word embeddings.

L-1 Word Representation. [12:30 PM 07.02.19]

$$V = [a, african, \dots, zulu, \text{UNK}] \quad |V| = 10,000$$

1-hot representation :-

Man (5391)	woman (9853)	king (4014)	Queen (7157)	apple (456)	Orange (6257)
$\begin{bmatrix} 0 \\ 0 \\ 0 \\ 0 \\ \vdots \\ 1 \\ 1 \end{bmatrix}$	$\begin{bmatrix} 0 \\ 0 \\ 0 \\ 0 \\ \vdots \\ 1 \end{bmatrix}$	$\begin{bmatrix} 0 \\ 0 \\ 0 \\ 1 \\ \vdots \\ 0 \end{bmatrix}$	$\begin{bmatrix} 0 \\ 0 \\ 1 \\ 1 \\ \vdots \\ 0 \end{bmatrix}$	$\begin{bmatrix} 0 \\ 0 \\ 1 \\ 0 \\ \vdots \\ 0 \end{bmatrix}$	$\begin{bmatrix} 0 \\ 0 \\ 0 \\ 1 \\ \vdots \\ 1 \end{bmatrix}$

If our algorithm learns, it won't learn a glam of apple Juice.

It can not generalize, it won't learn a glam of Orange Juice.

Product

[Because of apple and orange is zero]

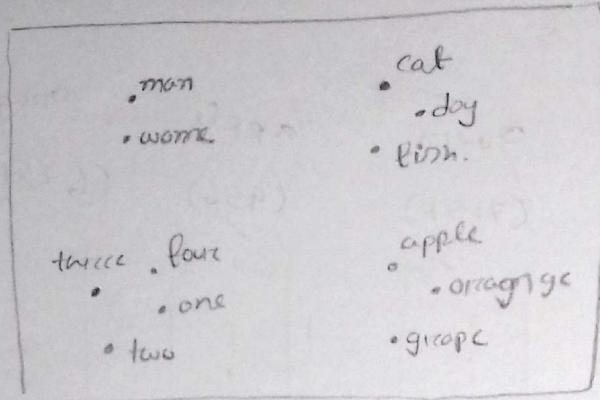
Featureized representation : word embedding :-

features	Man (5391)	woman (9853)	king (4014)	Queen (7157)	apple (456)	Orange (6257)
Gender	-1	1	-0.95	0.95	0.00	0.01
Royal	0.01	0.02	0.95	0.95	-0.001	0.00
Age	0.03	0.02	0.7	0.69	0.03	-0.02
Food	0.09	0.01	0.62	0.01	0.95	0.97
	↑	9853				
	↓	5391				
.....

[It helps us to generalize better over different things]

Visualizing word embeddings.

[] visualizing Data using t-SNE.



300D

↓

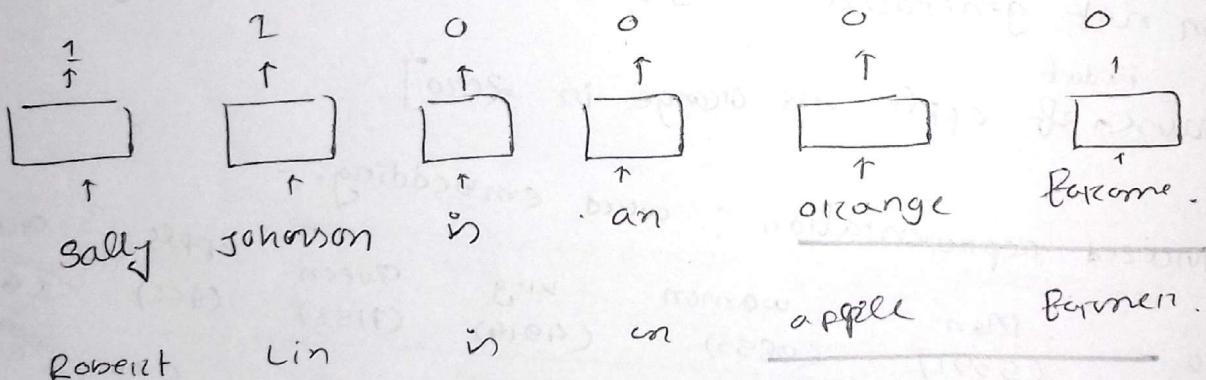
2D

[most important idea in Natural language processing]

L-2 Using word embeddings.

[12:55 pm 07.02.19]

Named entity recognition example.



~~we~~ → we can apply transfer learning here.

(1B - 100B words) → to learn ~~the~~ similar word.

100k → words (small training set)

Transfer Learning and word embeddings.

1. Learn word embeddings from large text corpus. (1-100B words)
(or download pre-trained embedding online).
2. Transfer embedding to new task with smaller training set (say, 100k words) [may use dense vector instead of one-hot vector]
3. optional:- Continue to finetune the word embeddings with new data.

word embedding relation to face encoding.

[It is quite similar which creates encoding for the input].

[one difference is that in face encoding neural network can have unknown face,

But in word embedding we had a vocabulary and has a fixed encoding for the words in vocabulary.]

L-3 Properties of word embeddings. [1:25 pm 07-02-19]

word embeddings are useful for analogy reasoning.

Such as man: woman, king → ?

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

find the word w : array max $\sim \text{sim}(e_w, e_{\text{king}} - e_{\text{man}} + e_{\text{woman}})$

[choose the word which maximizes similarity]

Cosine Similarity (more useful)

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

$$\text{sim}(e_w, e_{\text{king}} - e_{\text{man}} + e_{\text{woman}}).$$

[Slide 17: analogies game, man: woman → slide cursor here]

L-4 Embedding Matrix. [1:40 pm 07-2-19]

[Matrix form. Easy] (uncommented but less code).

E.

Oj.

$\therefore E \cdot O_j$.

$$[300, 10k]$$

$$[]$$

$$[300, 10k] [10k, J] = [300, J]$$

Feature matrix
for a word.

Feature vocab-size.

one-hot vectors for

J-th element in vocabulary.

[In practice, use specialized function to look up word embedding].

P2: Learning Word Embeddings : word2vec & Glove.

L1: Learning word embeddings. [7:05 pm 07.02.15]

I	want	a	glass	of	orange	<u>juice</u>
4343	9665	1	3852	6163	6257	

$$I \quad o_{4343} \rightarrow E \rightarrow e_{4343}$$

$$\text{want} \quad o_{9665} \rightarrow E \rightarrow e_{9665}$$

$$a \quad o_1 \rightarrow E \rightarrow e_1$$

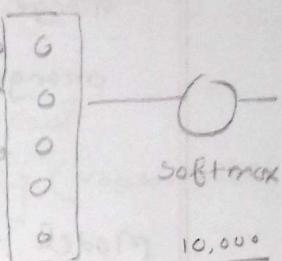
$$\text{glass} \quad o_{3852} \rightarrow E \rightarrow e_{3852}$$

$$\text{of} \quad o_{6163} \rightarrow E \rightarrow e_{6163}$$

$$\text{orange} \quad o_{6257} \rightarrow E \rightarrow e_{6257}$$

one-hot vector

Embedding vector



10,000

w², b²
word embedding

Context / target pairs.

[Context is embedding dimension that helps to learn embedding vector]

[Context is embedding dimension that helps to learn embedding vector]

→ Last 4 words. [Look at the past 4 words and predict the next word]

→ Last 4 words. [Look at the past 4 words and predict the next word]

→ 4 words on left & right [Look at 4 words on both side of Predicting word].

→ Last 1 word.

[Nearly 1 word] - skip gram model.

Here we use ~~one~~ context to predict target word and

eventually learn word embeddings.

Skip-Gram :-

I want a glass of orange juice to along with my cereal.

<u>Context</u>	<u>target</u>
orange	juice
orange	glass
orange	mj.
↑	↑

1.1.1

** Here we randomly chose a word as context word.

** Randomly choose another word as target word within a context window.

Model :-

Let, vocabulary size = ~~10,000~~ 10,000

Context $\in \{ \text{"orange"} \}$ \longrightarrow Target $\in \{ \text{"juicey"} \}$
 6257 - position $\quad\quad\quad$ 4839 - position

$[O_c \longrightarrow E \longrightarrow e_c \longrightarrow o \longrightarrow \hat{o}]$
 [one-hot vector] [embedding] $= E \cdot O_c$ softmax
 [One Context word] [weight matrix] [embedding vector]
 [One Context word]

$$\text{Softmax} : P(t|c) = \frac{e^{\theta_t^T \cdot e_c}}{\sum_{j=1}^{10,000} e^{\theta_j^T \cdot e_c}}$$

θ_t = probability associated with output target word t
 e_c = embedding vector for context word c

$$\text{Loss} : L(\hat{o}, o) = - \sum_{i=1}^{10,000} j_i \cdot \log \hat{j}_i \quad j = \begin{bmatrix} 0 \\ \vdots \\ i \\ \vdots \\ 0 \end{bmatrix} \text{- one-hot vector}$$

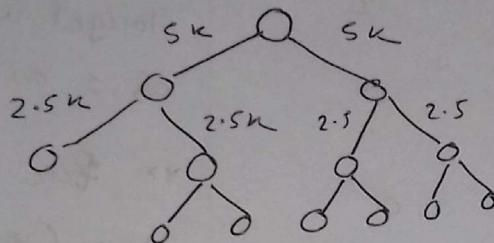
This is called skip-gram model.

Problems with softmax classification.

$$P(t|c) = \frac{e^{\theta_t^T e_c}}{\sum_{j=1}^{10,000} e^{\theta_j^T e_c}}$$

High computational to compute this sum value.

one solution is Hierarchical softmax :- (similar to binary search)



Computation cost

$$\text{becomes} = \log(\text{vocab-size})$$

[To build a hierarchical softmax classifier we can use different heuristics.] (Paper □ right side)

How to sample context c ?

** we can sample context c uniformly randomly but then the frequent words (the, of, a, and, to ...) will be updated and words (orange, juice, ...) will not update much.
So we use some heuristics to sample some common and uncommon words.

[The key problem of this algorithm is that computational cost of softmax function]

L-3. Negative Sampling. [8.10 pm 07.02.19]

[This algorithm is used to solve the problem of skip-gram model]

[paper no. 2(4)].

Defining a new learning problem.

I want a glass of orange juice to go along with my cereal.

	Context	word	target	v.v.y.
	orange	juice	1	** pick a context word, pick a target word and give it a label of 1 as it is positive.
Number of random words from vocabulary (k)	orange	king	0	** fix some number of time (k) for the context word choose random word from vocabulary and label as 0 to negative
	orange	book	0	
	orange	the	0	
	orange	of	0	
	x	y		

Now we have supervised learning algorithm where we input x (pair of words) and predict E_y (target label)

$$\begin{cases} k = 5-20 & \text{smaller dataset} \\ k = 2-5 & \text{large dataset} \end{cases}$$

Model :-

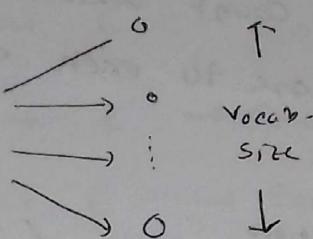
$$\text{softmax} : p(t|c) = \frac{e^{\theta_t^T e_c}}{\sum_{j=1}^{10,000} e^{\theta_j^T e_c}}$$

$$P(j=1 | c, t) = r(\theta_t^T e_c) \leftarrow \text{logistic regression.}$$

orange

G2ST

$$O_{G2ST} \rightarrow E \rightarrow e_{G2ST}$$



10000 binary classification
problem using logistic regression.

Selecting negative examples.

[After selecting positive word for the context word, How do we select negative samples]

→ Probabilities of words can be counted, (unigram distribution)

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

$f(w_i)$ = Frequency of word w_i
in your training set]

[English-GPT Pre-trained word-embeddings vector that are, you can download and use it]

Glove (global vectors for word representation)

A word-embedding model which is "count-based".

$$x_{ij}^t = \underset{c}{\overset{t}{\leftarrow}} \underset{c}{\overset{t}{\leftarrow}} \# \text{ times } i \text{ appears in context of } j$$

** x_{ij} is the count which captures how often the word i, j appears close to each other.

Model :-

$$\text{minimize} \sum_{i=1}^{10,000} \sum_{j=1}^{10,000} f(x_{ij}) \left(\underset{\theta_t^T \cdot e_c}{\downarrow} \underset{\theta_i^T \cdot e_c + b_i + b_j}{\downarrow} -\log x_{ij} \right)$$

weighting term $f(x_{ij}) = 0$ if $x_{ij} = 0$.

[This weighting term ensures that most frequent words do not get much weight and infrequent words less weight]

** (paper no 21) -

[slide 19 part 4 not write]

P3:- Applications using word embeddings.

L-1 Sentiment classification.

[9:35 pm 07.02.19]

[one of the challenges is that we may not have large training dataset. but word embedding can be useful with less data.]

Simple sentiment classification model.

The dessert is excellent

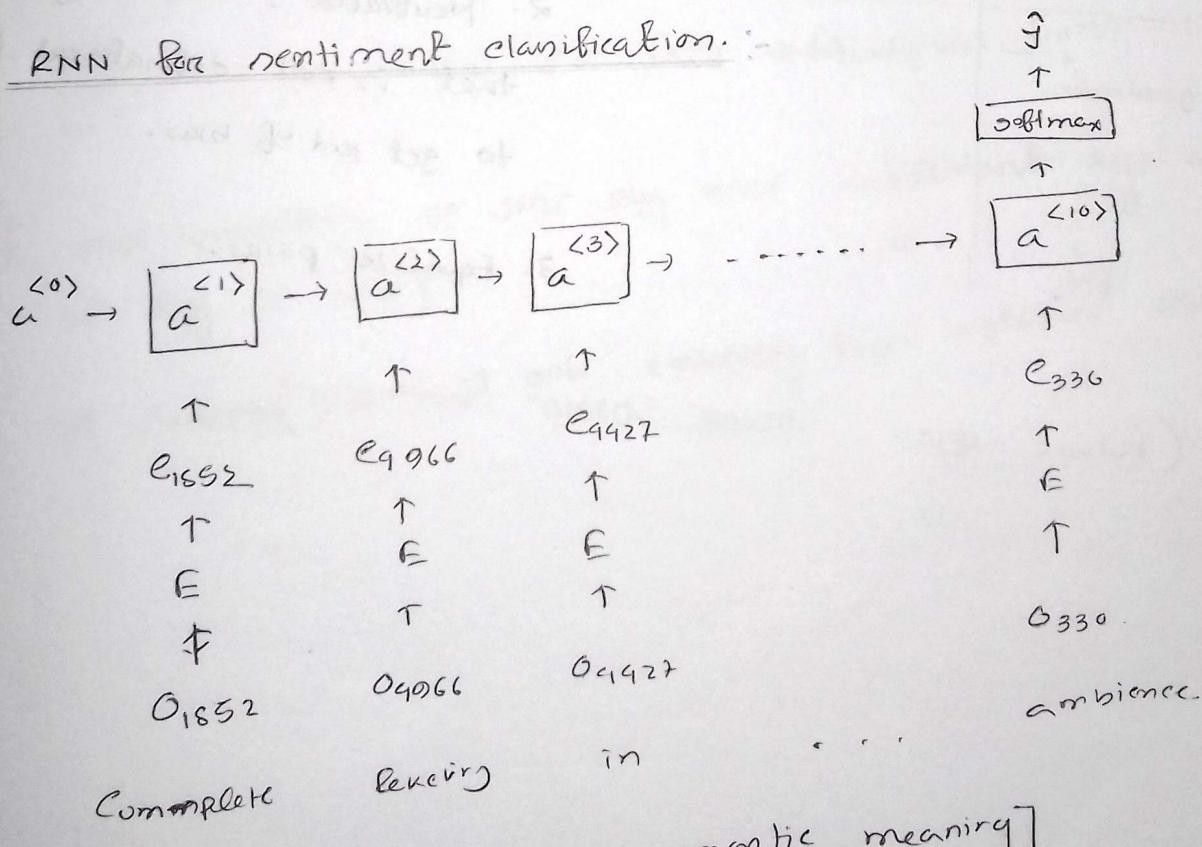
8928 2468 4694 3180

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

dessert $e_{2468} \rightarrow E \rightarrow e_{2468}$ —— $\boxed{\text{Avg}} \rightarrow O \rightarrow \hat{j}$
 is $e_{4694} \rightarrow E \rightarrow e_{4694}$
 excellent $e_{3180} \rightarrow E \rightarrow e_{3180}$
 $\underbrace{e_{3180}}_{300}$

[** The problem is that it does not count order of words.
 It does not take in account semantic meaning.

RNN for sentiment classification.



[This does help in capturing semantic meaning] .

(BLP (o) git for future & one 2(r)

L-2 Debiasing word embeddings

[11:10 pm 07.02.19]

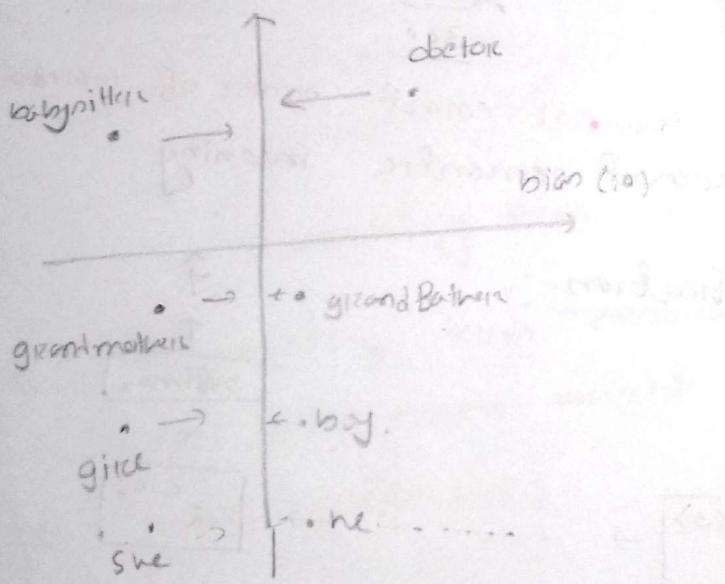
word embeddings can reflect gender, ethnicity, age, sexual orientation and other biases of the text used to train the model.

[we must try to diminish this type of biases as machine learning takes really much important decisions].

[git for paper size] [Balakrishnan et.al. 2016]

Addressing bias in word embeddings.

1. Identify bias direction.



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

3. Equalize pairs.

(future idea 12:07 07.02.19 09:15 2019).

Programming Assignment :-

operations on word vectors.

- Cosine similarity a good way to compare similarity between pairs of word vectors.
- For NLP applications using a pre-trained set of word vectors from the internet is often a good way to start.

Emoji! (Some on sentiment classification)

(Join to sentiment classification idea GPT-3)

[** Join word embedding glove vector API - GPT-3 notes
and for the best generation API]

[Part 2 of GPT-3 paper (2020-2021)]

[word embedding of joint blog once implement GPT-3 API
as blog is that.]

** Consider cosine of word embedding layer implement GPT-3

** Important points.

- In a NLP task where training set is small, word embeddings can be used to help your algorithm.
- word embeddings allow a model to work on words in the test set that may not even appeared in training set.
- To use mini-batches, the sequences need to be padded so that all the examples in mini-batch have same length.
- Design your embedding matrix carefully.
 - Can use pretrained values or initialize in your training set.
- After embedding use your Deep learning architecture.

[V.V.J understand what is happening in week 9 exercise really important]

Week-3: Sequence models & Attention mechanism.

P1: Various Sequence to Sequence Architectures.

L-1 Basic Models. [10:35 AM 09.02.19]

Basic machine translation model.

→ Their paper is references from Google.

Automatic image captioning.

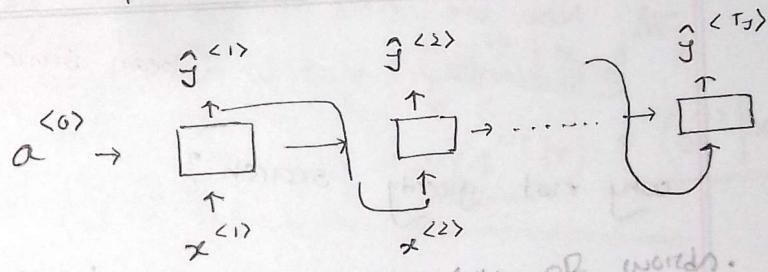
→ AlexNet takes input from Decoder to automatic image caption on it.

→ Decoder circuit - Normally RNN architecture.

L-2 Picking the most likely sentence. [10:50 AM 09.02.19]

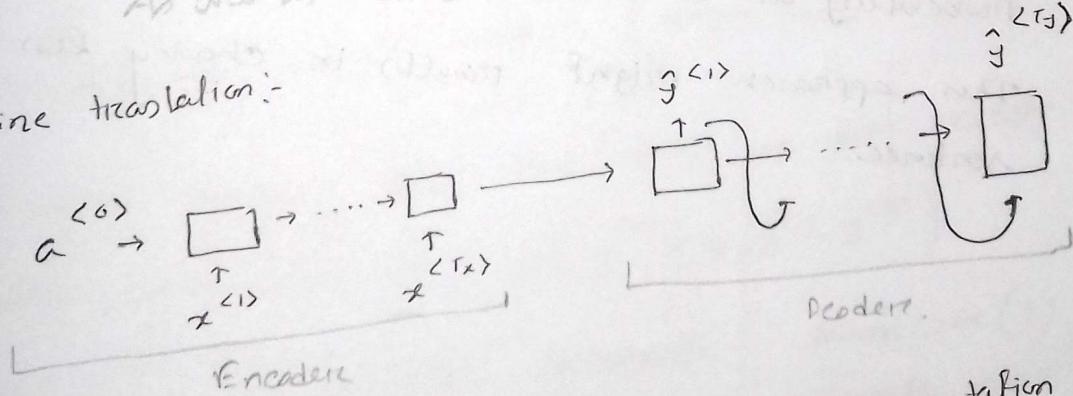
Machine translation as building a Conditional Language Model.

Language model :-



[Language model helps to calculate probability of words.
As well as Encodable sentences.

Machine translation:-



** Encoder circuit figures out some representation of the input sentence and feed it to decoder circuit.

"condition language model"

$$p(y<0>, \dots, y<Tj> | x<0>, \dots, x<Tx>)$$

Probability of an English sentence Conditioned on some French sentence.

finding the most likely translation.

x = a french sentence.

$P(J^{<1>} \dots J^{<T_J>} | x)$ = Probability of corresponding english sentence.

We will not sample randomly. We will take the english sentence that maximizes the probability.

$$\arg \max_{J^{<1>} \dots J^{<T_J>}} P(J^{<1>} \dots J^{<T_J>} | x).$$

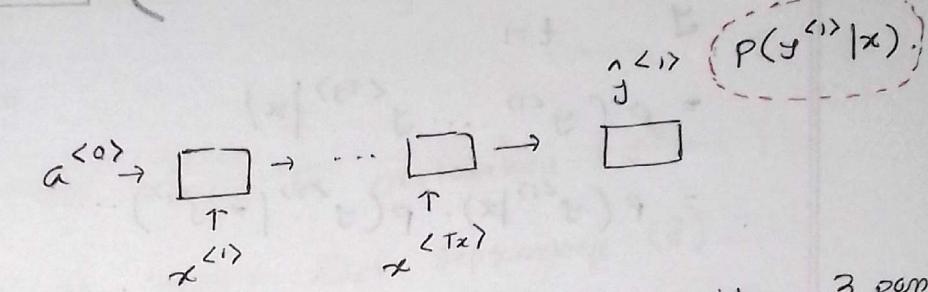
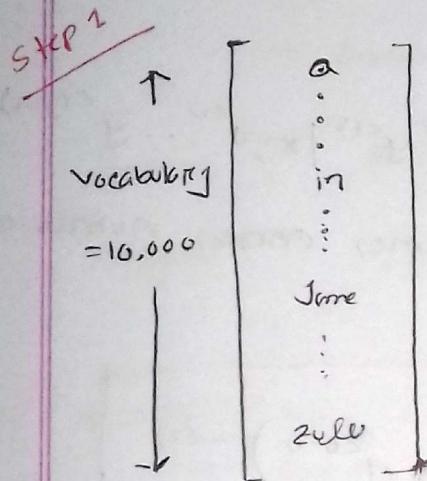
→ Now we need an algorithm which calculates this probability efficiently, that is (Beam Search).

why not greedy search?

→ In greedy search it takes the word which maximizes probability at the current step. This approach might result in choosing less optimal sentence.

[Beam search is an algorithm which try to maximize conditional probability in machine translation]

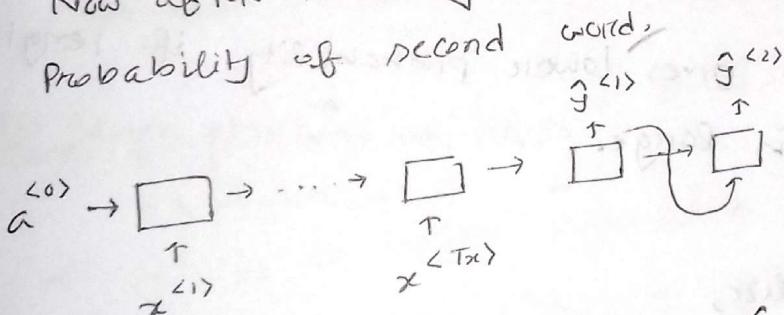
Beam Search algorithm. (slide 1 from google)



$B = 3$ (beam width), it will consider 3 possible output at a time.

Step 2

Now after selecting first 3 words



we will calculate the

$$P(y^{<2>} | x, y^{<1>}) \\ = P(y^{<1>} | x) \cdot P(y^{<2>} | x, y^{<1>})$$

Number of possible words = (Beam size \times vocab-size).

After selecting first two word we will select the 3rd word which maximize $P(y^{<3>} | x, y^{<1>}, y^{<2>})$.

in each step we will memorize output which is calculated used to calculate next step.

beam size (3) possible

[diff - ref - ref
diff out].

[slide 1 was 2nd
from last part 27th]

Length normalization. (Ans)

$$\begin{aligned} \operatorname{argmax}_{\mathcal{Y}} \prod_{t=1}^{T_Y} p(y^{(t)} | x, y^{(1)}, \dots, y^{(t-1)}) & \quad -(1) \\ = p(y^{(1)}, \dots, y^{(T_Y)} | x) \\ = p(y^{(1)} | x) \cdot p(y^{(2)} | x, y^{(1)}) \cdots \cdots \cdots p(y^{(t)} | x, y^{(1)}, \dots, y^{(T_Y-1)}) \end{aligned}$$

As probabilities are small this sometimes causes numerical underflow. we can avoid this by

$$\rightarrow \operatorname{argmax}_{\mathcal{Y}} \sum_{t=1}^{T_Y} \log p(y^{(t)} | x, y^{(1)}, \dots, y^{(t-1)}) \quad -(2)$$

The objective function gives lower probability if length of the sentence (T_Y) is large. (As we multiply more small probabilities).

This can be normalized,

hyper parameter

$$\alpha = 0.7$$

$$\alpha = 0.$$

$$\alpha = 1$$

$$\left[\frac{1}{T_Y^\alpha} \sum_{t=1}^{T_Y} \log p(y^{(t)} | x, y^{(1)}, \dots, y^{(t-1)}) \right]$$

Beam width $\frac{B}{r}$

parameter

large B: better result, slower.
Small B: worst result, faster.

* Beam search runs faster but is not guaranteed to find exact maximum for $\operatorname{argmax}_{\mathcal{Y}} p(y | x)$.

L-5 Error analysis in beam search. [11.50 AM, 09.02.10]

[Slide 9 Beam technique (error analysis contd) (8th).]

In machine translation, we have two main components,

→ RNN. [computes $p(j|x)$]

→ Beam Search.

Human :- Jane visits Africa in September. (j^*)

Algorithm :- Jane visited Africa last September (j).

Case 1 :- $p(j^*|x) > p(j|x)$.

→ Beam search chose j . But j^* attains higher $p(j|x)$.

Conclusion :- Beam search in a fault.

case 2 :- $p(j^*|x) \leq p(j|x)$

→ j^* is a better translation than j . But RNN predicted

$p(j^*|x) < p(j|x)$

Conclusion :- RNN model is a fault.

(Faulty early out) [error analysis for prob].

[After automatic generation of machine translation Bleu score helps to determine how good the translation is]

Bleu score is calculated on machine translation (MT) output and human references.

Example :- Ref 1 :- The cat is on the mat

Ref 2 :- There is a cat on the mat.

MT out : the the the the the the the

[steps after Bleu score calculation method (n-gram) (easy one)]

** Bleu score on bigram. (slide 3 generic & ref)

** Bleu score on unigram :-

$$P_1 = \frac{\sum_{\text{Unigrom} \in \hat{J}} \text{Count}_{\text{ref}}(\text{unigrom})}{\sum_{\text{unigrom} \in \hat{J}} \text{Count}(\text{unigrom})}$$

$$P_n = \frac{\sum_{n\text{-grom} \in \hat{J}} \text{Count}_{\text{ref}}(n\text{-grom})}{\sum_{n\text{-grom} \in \hat{J}} \text{Count}(n\text{-grom})}$$

P_n = Bleu score on n-gram only.

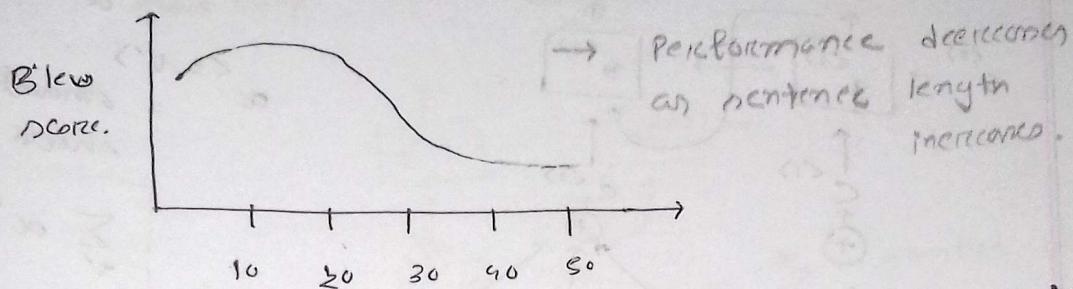
$$\exp \left(\frac{1}{m} \sum_{n=1}^m P_n \right)$$

Combined Bleu score :

$$BP = \begin{cases} 1, & \text{if } \text{MT_output_length} > \text{ref_output_length} \\ \exp \left(1 - \frac{\text{MT_output_length}}{\text{ref_out_length}} \right) & \text{otherwise,} \end{cases}$$

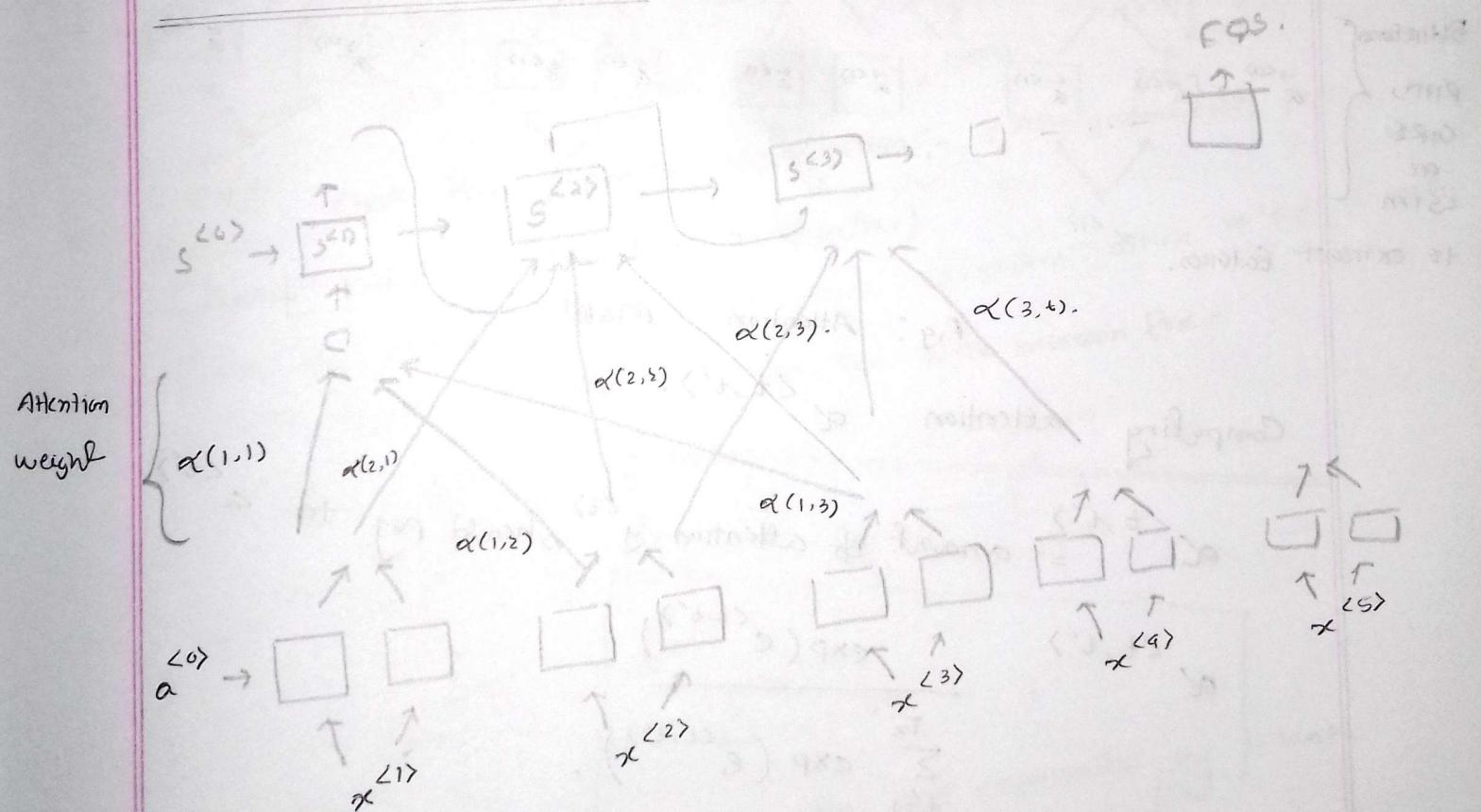
L-7 Attention model intuition [1:00 pm 09.02.19].

[Encoder- Decoder model (Previous model) performs well on short sentences but do not perform much well on long sentence]



** Attention model helps to perform well on long sentence as well.

Attention model intuition.



$$\text{Attention weight} = \alpha(t, t')$$

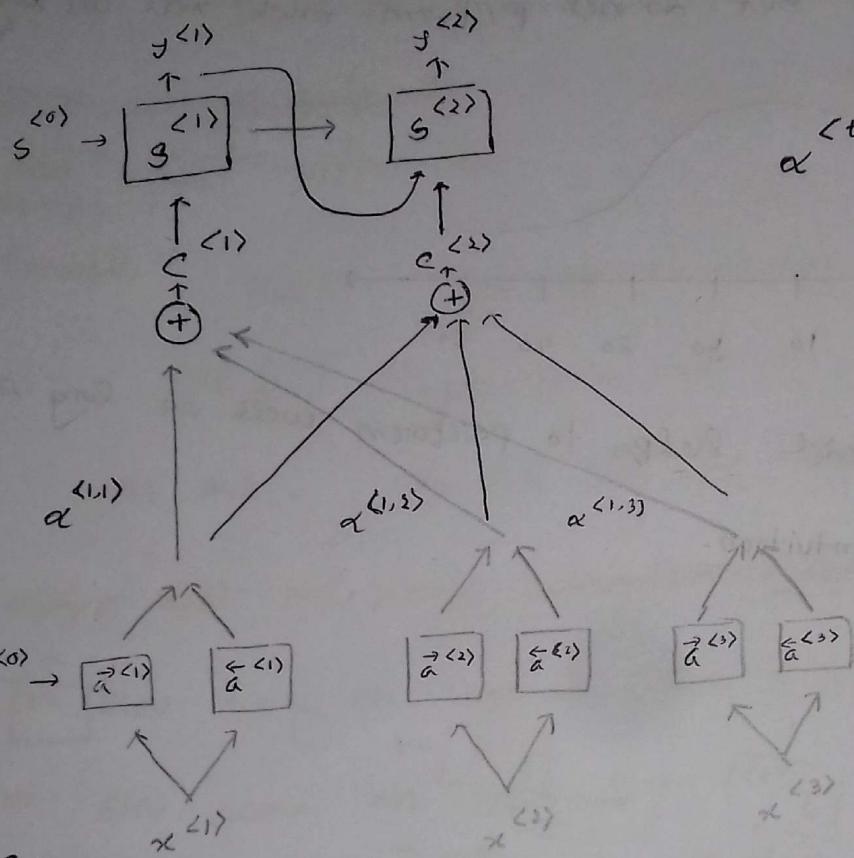
→ When you generating tth english word on how much you should give attention to tth french word.

L-B Attention model.

[1:20 pm 09.02.19]

[Pictorial representation of Attention model of paper by notes 2(t).]

(Slide 9 from 7th Jan 2019)



Bidirectional
RNNs
GRU
or
LSTM

to extract features.

$$\alpha^{(t,t')} = \left(\vec{a}^{(t')} \cdot \vec{a}^{(t)} \right)$$

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

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

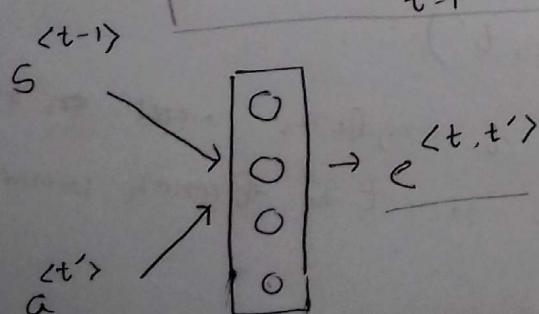
$$c^{(1)} = \sum_{t'} \alpha^{(1,t')} \cdot a^{(t')}$$

Fig: Attention model.

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

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

$$\alpha^{(t,t')} = \frac{\exp(e^{(t,t')})}{\sum_{t'=1}^{T_x} \exp(e^{(t,t')})}$$

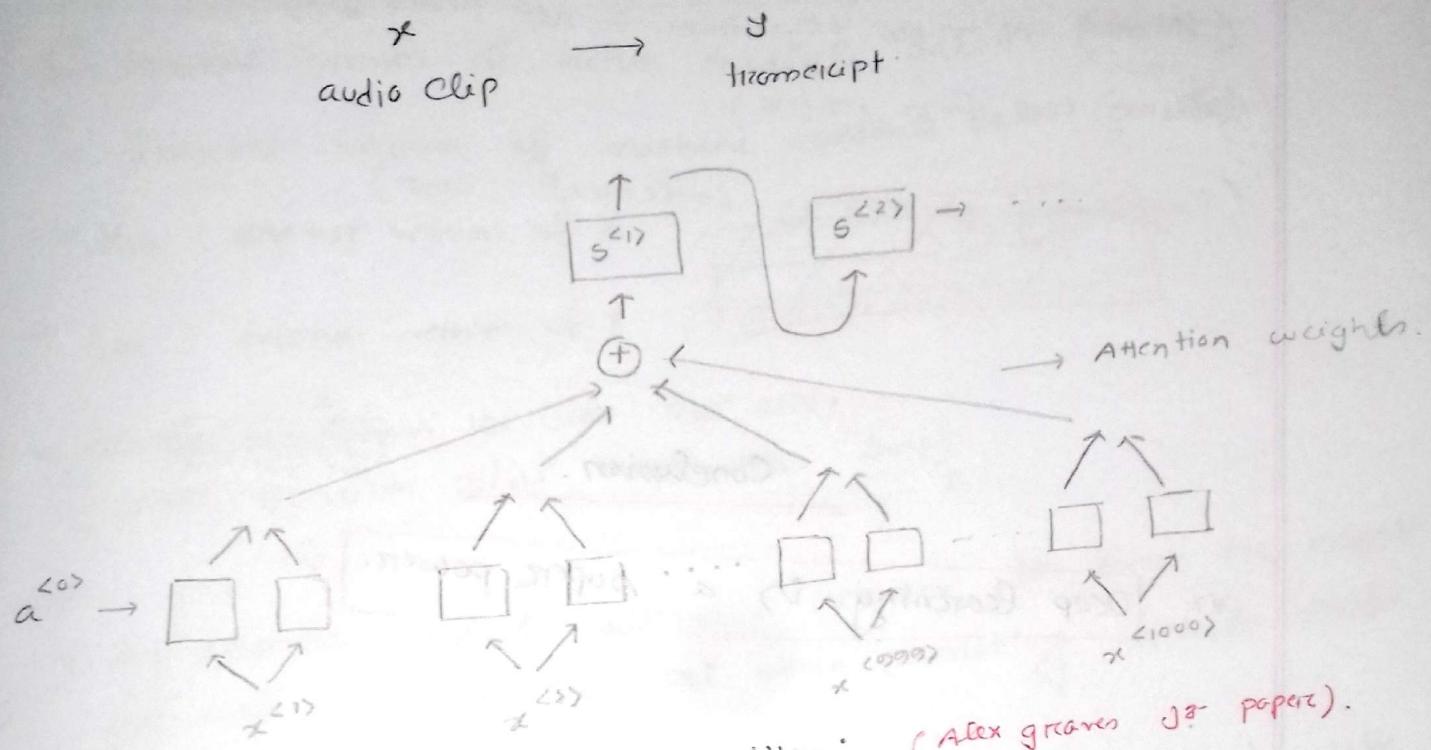


* visualizing $\alpha^{(t,t')}$

[cost is high]

L-1 Speech recognition. [4.40 pm 09.02.19]

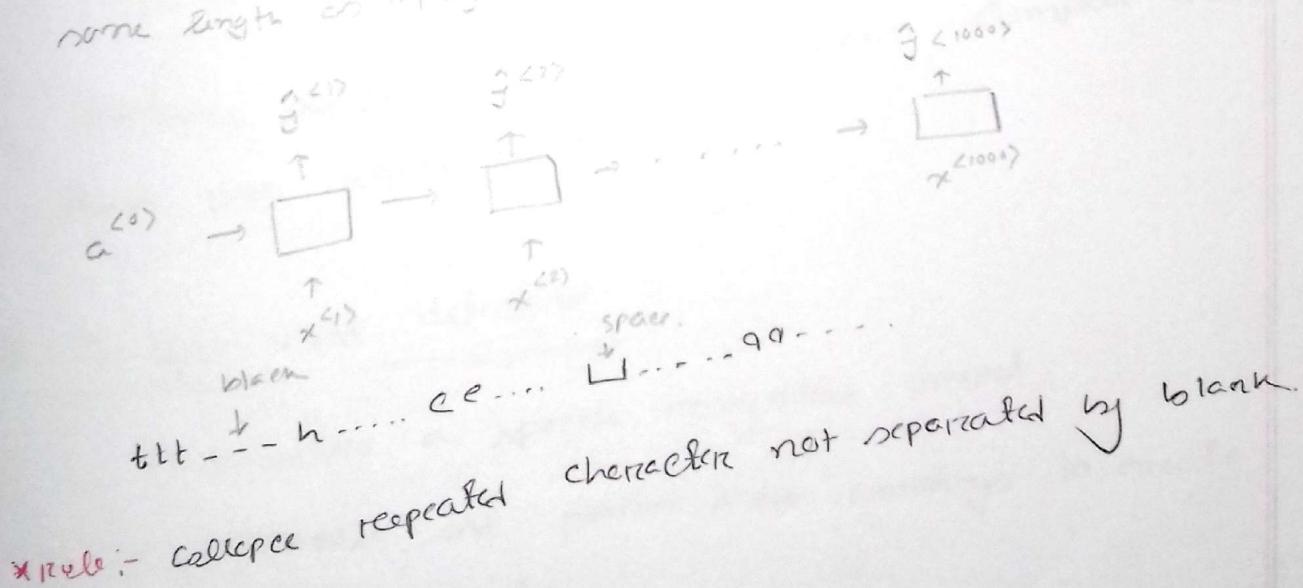
Attention model for speech recognition :-



ETC cost for speech recognition :-

(connectionist temporal classification) where output is not
certain in object recognition

[CTC for predicting correctly the biochemical foot .
name length as input] * The quick brown fox.
 $\text{length} < 1000$



L-2 Trigger word detection. [5.00 PM 00.02.17]

Trigger word detection algorithm :-

(Almost optimized algorithm is not Bound yet)

(slide 0 for more details)

(Programming expense & implement cost)

Conclusion.



** Deep learning is a super power. ✓

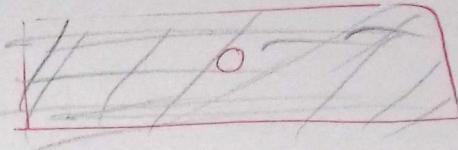
Programming Assignment :-

1. Translating human readable data into machine readable data.

\rightarrow Processed version of human readable data (m, T_x)

\rightarrow Processed version of machine readable data (m, T_y).

$\rightarrow X_{on}$: one-hot version of X



** Attention mechanism is used for this!

(LSTM, Bi-LSTM with their main GMP)

=
An attention model allows a network to focus on most relevant parts of the input when producing specific part of output.

A network ~~can~~ using an attention mechanism can translate from input of length T_x to output length T_y where T_x and T_y can be different.

[Attention model are useful to learn complex mappings
from one sequence to another].

2. Trigger word detection.

\rightarrow Structure a speech recognition project.

\rightarrow Structure a speech recognition project to encode

\rightarrow Synthesize and process audio recordings to encode train/dev datasets.

\rightarrow Train a trigger word detection model and make predictions.