

Date : 26/7/24
MON TUE WED THU FRI SAT SUN

COURSE 35

- SEQUENCE MODELS

Sequence data:

Speech recog., music gen^r, sentiment classif; DNA anal; mlc transl; video activity recog; etc; named-entity recogn.

Notation

e.g. $X: \text{Harry Potter and Hermoine Granger}$
 $x^{<1>} \leftarrow \text{invented } \frac{x^{<2>}}{a} \text{ new spell. } \frac{x^{<3>}}{T_x=9}$

$y: \underline{1} \ 1 \ 0 \ 1 \ 1 \ 0 \ 0 \ 0 \ 0 \ T_y=9$
 $y^{<1>} \quad y^{<3>}$

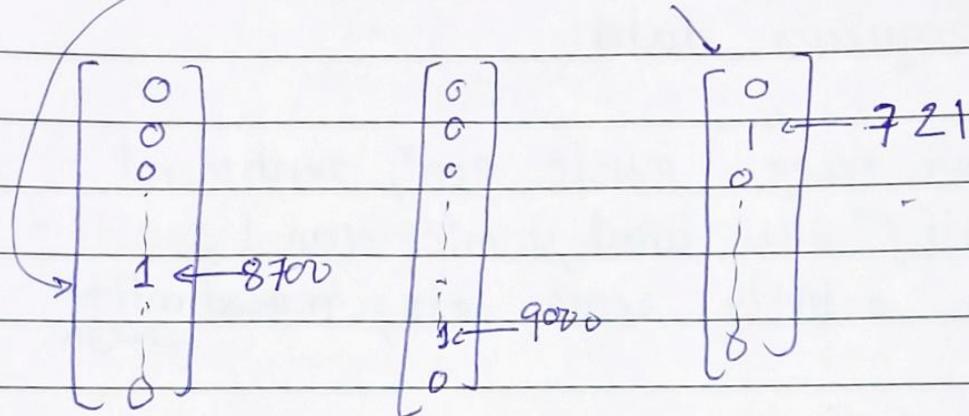
$x^{<t>} \Rightarrow t^{\text{th}}$ element _{of p} $\overline{T_x^{(i)}}$
 $x^{(i)<t>} \Rightarrow t^{\text{th}}$ elem. in sequence of
 i^{th} training example.

→ A dictionary for vocabulary of around 50k words is used.

↓

The sequence words are 1-hot encoded.

eg. Harry Potter and.

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

9 such vectors created if $Tx = a$

Problem with standard network:

- IIP & OIP can be diff. length in diff eg.
- Doesn't share features learned across diff. positions of text
(If "Harry" comes & is recog. as name then further anywhere "Harry" comes it should recognize it as name.)

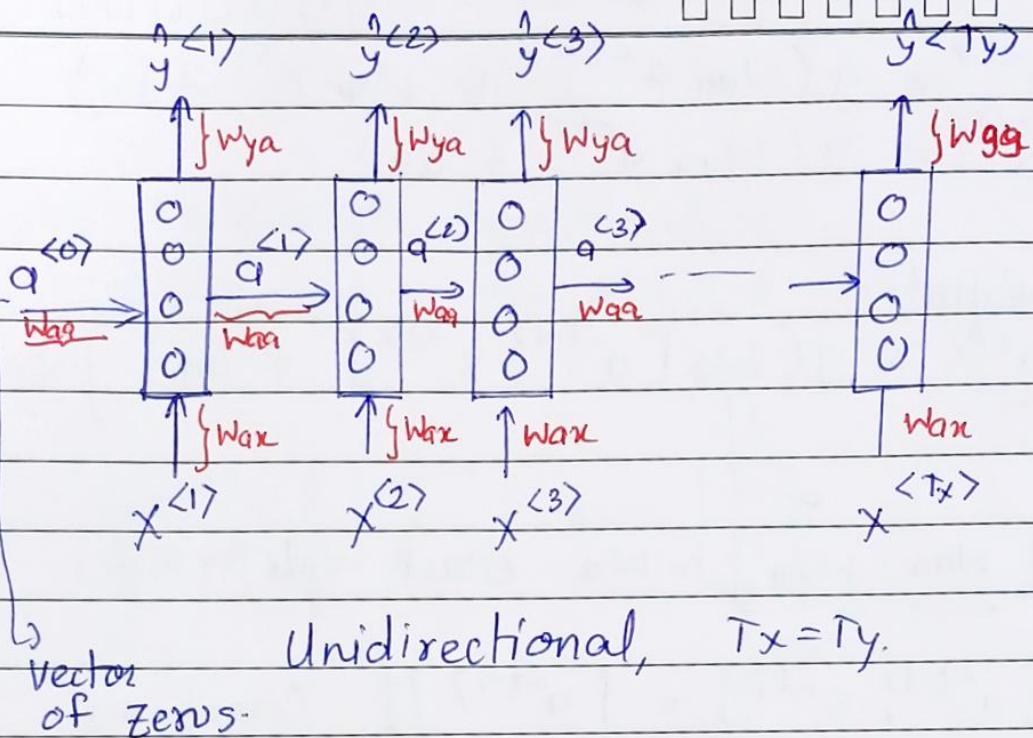
Recurrent Neural Networks.

E.

Date :

MON TUE WED THU FRI SAT SUN

<input type="checkbox"/>						
--------------------------	--------------------------	--------------------------	--------------------------	--------------------------	--------------------------	--------------------------



Prediction at a certain time uses info from inputs only EARLIER in sequence. not LATER.

e.g. He said, "Teddy Roosevelt is President".

He said, "Teddy bears are cute".

→ Forward Propagation ←

$$a^{<0>} = \vec{0}$$

$$a^{<1>} = g_1 \left(W_{aa} a^{<0>} + W_{ax} x^{<1>} + b_a \right)$$

$\underbrace{\quad}_{\text{gen'lly tanh/ReLU}}$

$$y^{<1>} = g_2 \left(W_{ya} a^{<1>} + b_y \right)$$

$\underbrace{\quad}_{\text{gen'lly sigmoid.}}$

Parameters W_a, b_a, W_y, b_y
will remain same for
all elements of sequences.

Date :

MON TUE WED THU FRI SAT SUN

$$a^{<t>} = g(W_a a^{<t-1>} + W_x x^{<t>} + b_a)$$

$$\hat{y}^{<t>} = g(W_y a^{<t>} + b_y)$$

Simplified:

* $a^{<t>} = g(W_a [a^{<t-1>}, x^{<t>}], + b_a)$ *

$$\begin{bmatrix} W_a a^{<t-1>} \\ W_x x^{<t>} \end{bmatrix} = W_a \text{ (stack side by side)}$$

$$[a^{<t-1>}, x^{<t>}] = \begin{bmatrix} a^{<t-1>} \\ x^{<t>} \end{bmatrix} \quad \begin{matrix} \uparrow & \text{(stack on} \\ \downarrow & \text{top of} \\ , & \text{other)} \end{matrix}$$

* $\hat{y}^{<t>} = g(W_y a^{<t>} + b_y)$ *

→ Backprop ← "Backprop through time".

Loss for a single word:

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

Loss for full sequence:

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

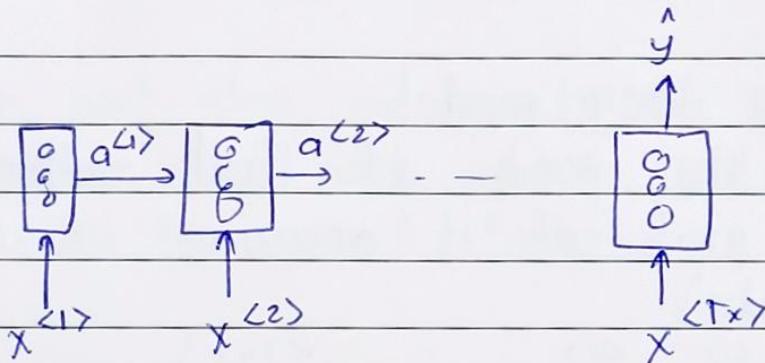
Date :

MON TUE WED THU FRI SAT SUN

$T_x \neq T_y$ RNN

e.g. sentiment classification.

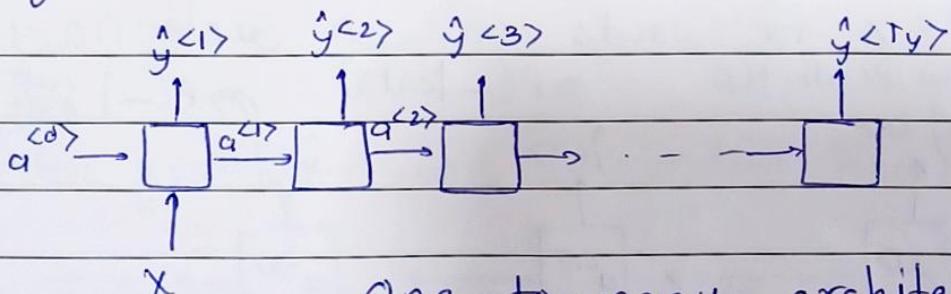
$$x = \text{text} \quad y = 0 \ 1 \ 2 \ 3 \ 4 \ 5$$



Many-to-one architecture.

e.g. Music generation.

$$y^{(0)} \rightarrow y^{(1)} \ y^{(2)} \ \dots \ y^{(T_y)}$$

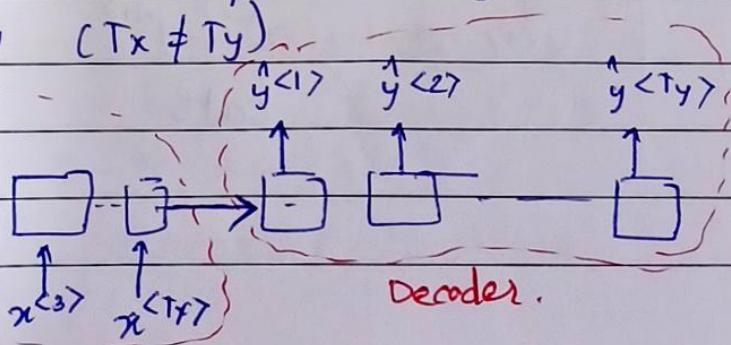
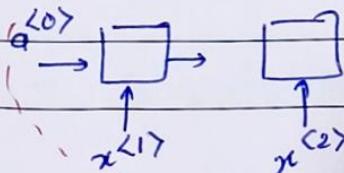


One-to-many architecture.

e.g. Previous page's Many-to-many ($T_x = T_y$)

Many-to-many ($T_x \neq T_y$)

Encoder



Decoder.

⇒ Machine Translation.

Date :

MON TUE WED THU FRI SAT SUN

Language Modelling

e.g. In speech recog:

"The apple and pair salad"

"The apple and pear salad".

* Work of a lang. model:
what is the prob. of that sentence

(meaning prob. of it occurring in some data).

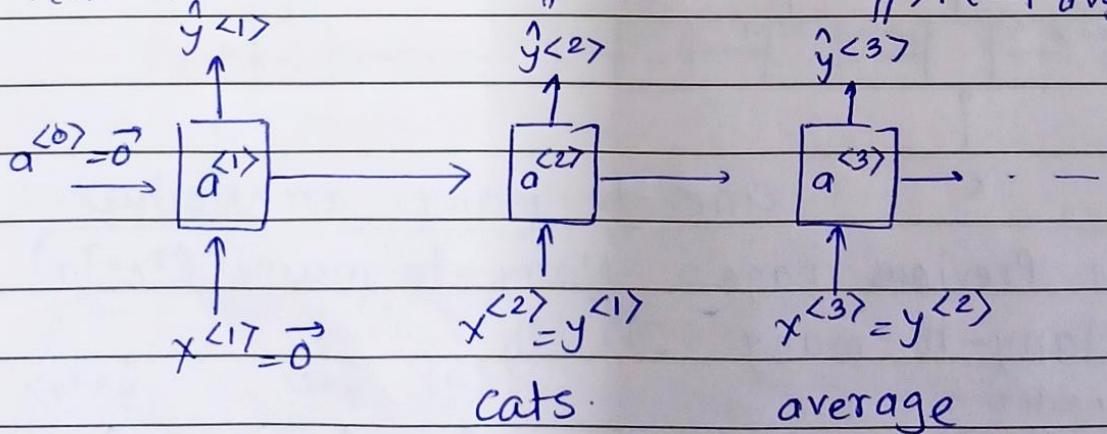
i.e. $P(y^{(1)}, y^{(2)}, \dots, y^{(T_y)})$

→ Tokenize ←

<EOS> added @ end.

<UNK> for words not in dictionary.

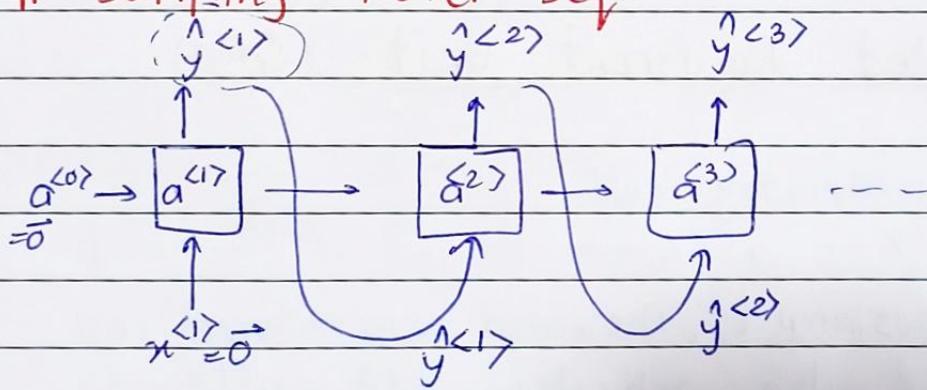
$P(a)P(\text{apple}) \cdot P(\text{cat}) \cdots P(z) \Rightarrow P(\text{-1 cats}) \Rightarrow P(\text{-1 avg})$



$\hat{y}^{<1>}$ will be softmax prediction of
'what is the chance that first word
is some word a--z from dict' or
 $<\text{EOS}>$ or $<\text{UNK}>$!

In next step, we provide the correct word 'cats' and make prediction $\hat{y}^{(2)}$ given the word 'cats'.

Sampling Novel Seq.



→ For first word we np-random-choice.

- when u get <EOS> token, stop.

character-level lang. model

Vocabulary = [a, b, c .. , z, ʌ, ɔ:, ɔ:, ʃ, ɒ:
0-9, A-Z ..]

+ computationally expensive to train.

Vanishing & Exploding gradients

e.g. The cat, —, was sleeping.
 The cats, —, were sleeping

very long-term dependency

weakness of a basic RNN Algo: bcz
 difficult for backprop to update
 earlier weights caused by last layer.

Gated Recurrent Unit: (GRU)

c = Memory cell.

$$c^{<t>} = a^{<t>}$$

@ timestamp t , the word was stored in
 $c^{<t>}$ & also activation o/p = $a^{<t>}$

$$\begin{aligned}\tilde{c}^{<t>} &= \text{candidate for replacing } c^{<t>} \\ &= \tanh(w_c [c^{<t-1>}, x^{<t>}] + b_c)\end{aligned}$$

Γ_u = Update gate

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

↳ assume always either 0 or 1.

bcz for almost majority values $\sigma = 0$ or 1

elem-wise

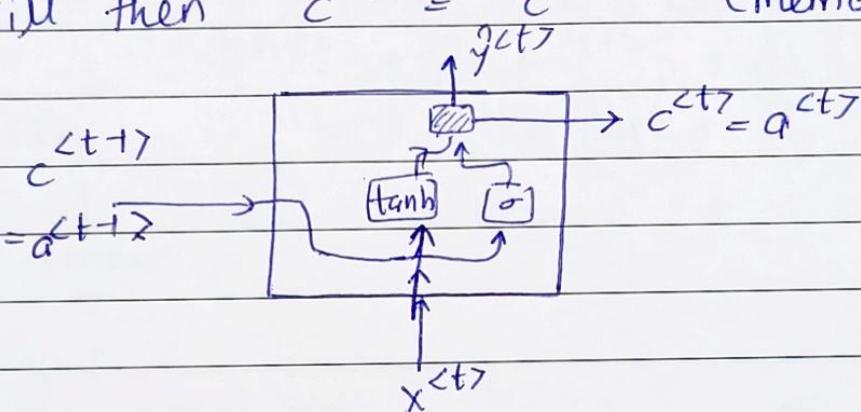
Date :

MON TUE WED THU FRI SAT SUN

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

Job of update gate : to decide whether
to update the $c^{(t)}$ value with $\tilde{c}^{(t)}$
or not.

Till then $c^{(t)} = c^{(t-1)}$ (memorize)



$c^{(t)}$, $\tilde{c}^{(t)}$, Γ_u all same dimensions.

We can assign ~~some~~^{some} bits for
storing singular/plural, food, etc.

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

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

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

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

LSTM (Long Short Term Memory)

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

$$\Gamma_u = \sigma (W_u [a^{<t-1>}, x^{<t>}] + b_u)$$

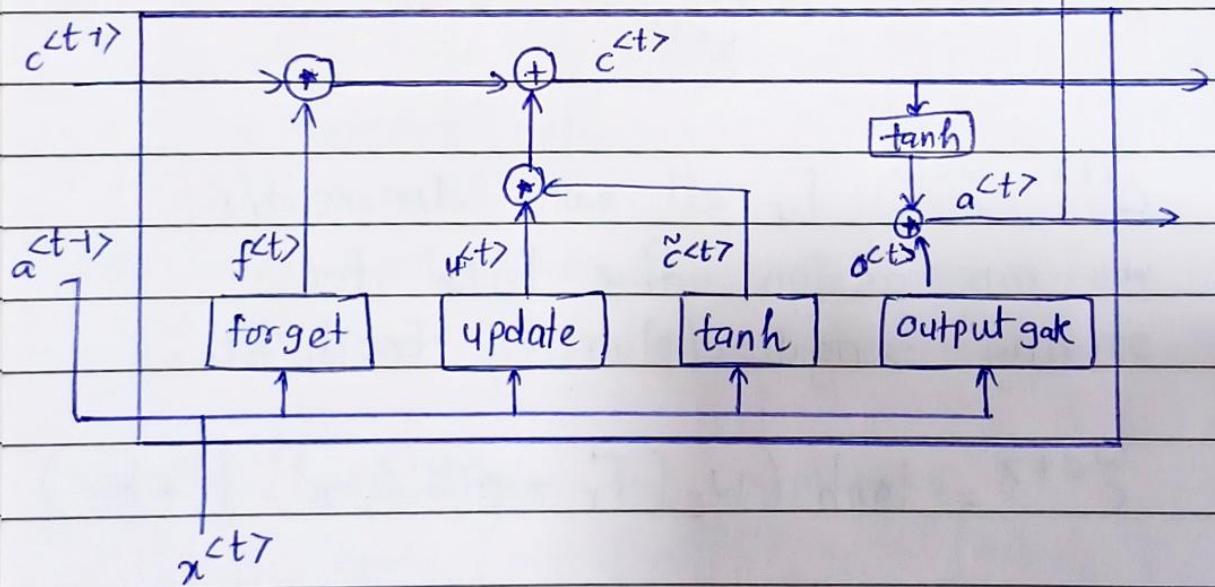
$$\Gamma_f = \sigma (W_f [a^{<t-1>}, x^{<t>}] + b_f)$$

$$\Gamma_o = \sigma (W_o [a^{<t-1>}, x^{<t>}] + b_o)$$

$$c^{<t>} = \Gamma_u + \tilde{c}^{<t>} + \Gamma_f * c^{<t-1>}$$

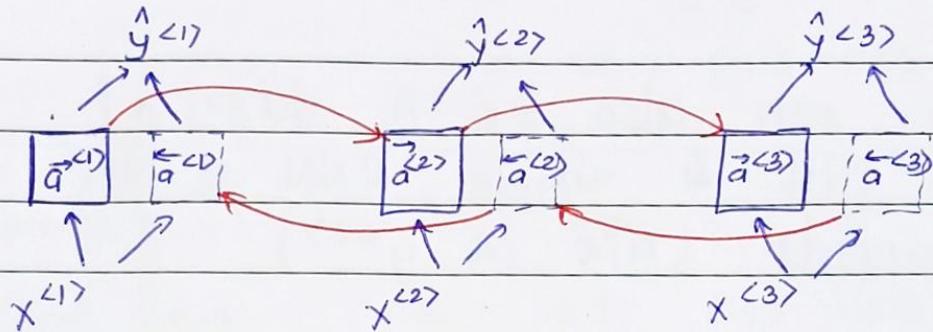
$$a^{<t>} = \Gamma_o + \tanh c^{<t>} \quad y^{<t>}$$

↑
softmax



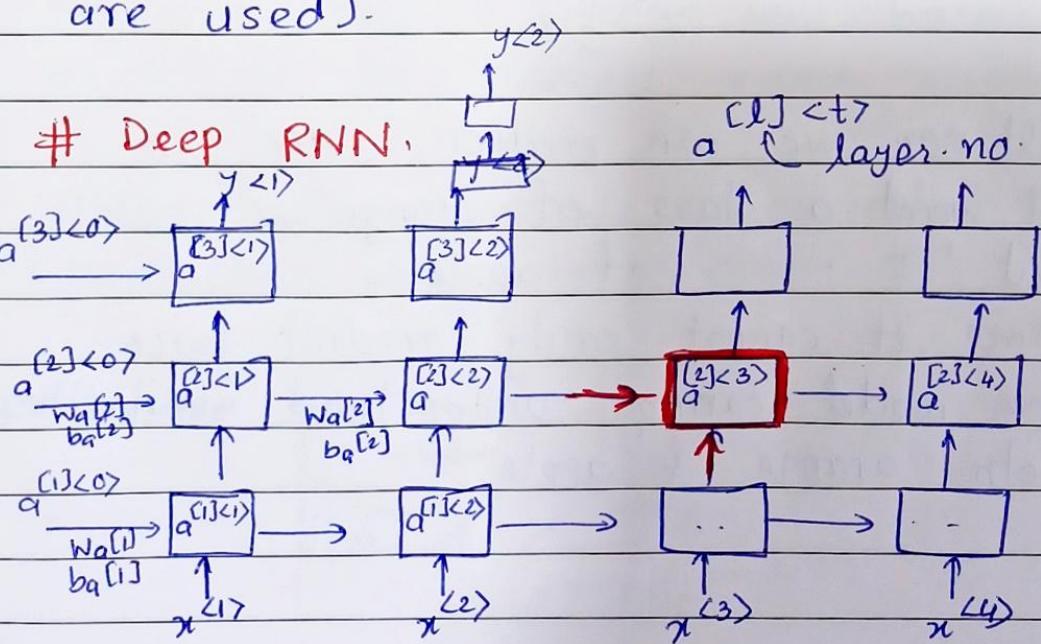
LSTM prev. invented, more computational,
but more powerful than GRU

Bidirectional RNN



All this is
Acyclic Graph just for prop.

These indiv. blocks can be LSTM or GRU
Disadv of BRNN: You need the entire sequence of data before making any predictions. e.g. In speech recogn, u have to wait till a complete sequence. (Actually, much more complex models are used).



Date :

MON TUE WED THU FRI SAT SUN

<input type="checkbox"/>						
--------------------------	--------------------------	--------------------------	--------------------------	--------------------------	--------------------------	--------------------------

$$a^{(2) \times 3} = g \left(W_a^{(2)} [a^{(2) \times 2}, a^{(1) \times 3}] + b_a^{(2)} \right)$$

There can also be a mix of deep RNN & simple RNN or NN afterwards (like in $y^{(2)}$)

Word Representation.

Word embeddings - lets algo. automatically understand analogies like man-woman, King-queen.

Vocabulary

$$0_{5391} = \begin{bmatrix} 0 \\ 0 \\ \vdots \\ 1 \\ 0 \\ \{ \end{bmatrix} \rightarrow \text{Man}(5391)$$

Till now we can predict 'juice' in
'I want a glass of orange - !'

But 'I -- of apple - '

Here, it cannot easily predict juice.

The model cannot understand similarities betn orange & apple.

Date :

MON TUE WED THU FRI SAT SUN

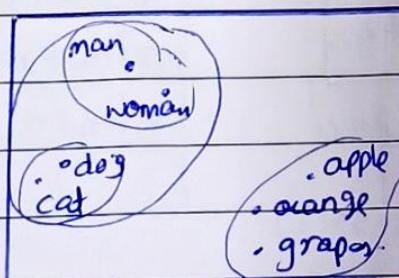
→ Featurized repr - Word embedding ←

	Man (5391)	Woman (1853)	King	Queen	Apple	Orange
Gendr	-1	1	-0.95	0.97	0	0.01
Royal	0.01	0.02	0.93	0.95	-0.01	0.0
Age	,	,	,	,	,	,
Food	,	,	,	,	,	,
size	,	,	,	,	,	,
cost	,	,	,	,	,	,
verb	,	,	,	,	,	,
Noun	,	,	,	,	,	,
etc.	e ⁵³⁹¹	e ⁹⁸⁵³				

eg. 300 such features.

So, now Apple & Orange will have much of the features same

Using t-SNE algo, we can convert this 300D vectors to 2D plots.



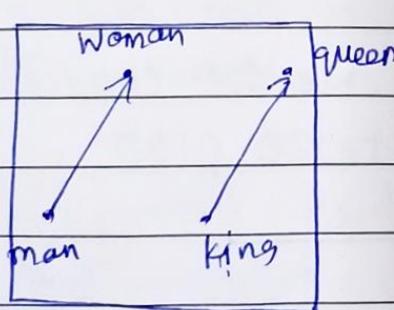
eg. use in named entity recog.

\hookrightarrow Sally Johnson is an orange farmer.
 \hookrightarrow Train.

\hookrightarrow Robert Lin is an apple farmer.
 \hookrightarrow Test

1. Learn word emb. from large text corpus.
2. Transfer emb. to new task with smaller training set.
3. Continue to finetune w.f with new data.

$$\begin{matrix} \text{eman} - \text{ewoman} \approx \\ \text{eking} - \text{equeen} \end{matrix} \left[\begin{array}{c} -1 \\ \approx 0 \\ \approx 0 \\ 0 \end{array} \right] = \left[\begin{array}{c} -2 \\ 0 \\ 0 \end{array} \right]$$



300D space

points

(not on t-SNE)

Find word w such that
maximize
similarity
 $\text{sim}(ew, eking - eman) + ewoman$

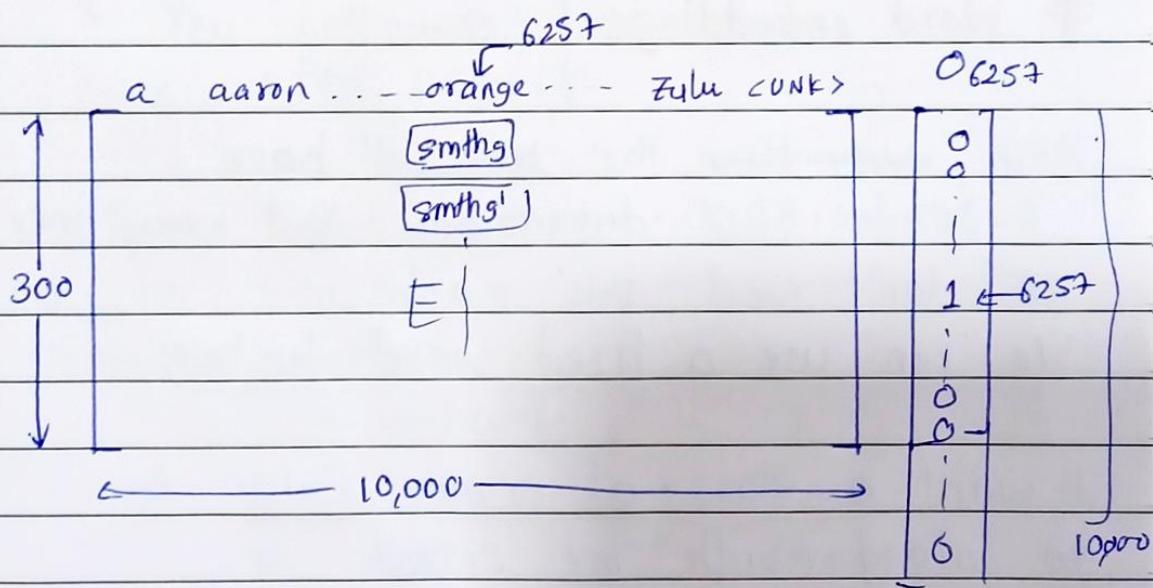
Date :

MON TUE WED THU FRI SAT SUN
 \rightarrow Cosine similarity \leftarrow

↳ commonly used.

 $\text{sim}(\text{ew}, \text{eking}-\text{eman} + \text{ewoman})$

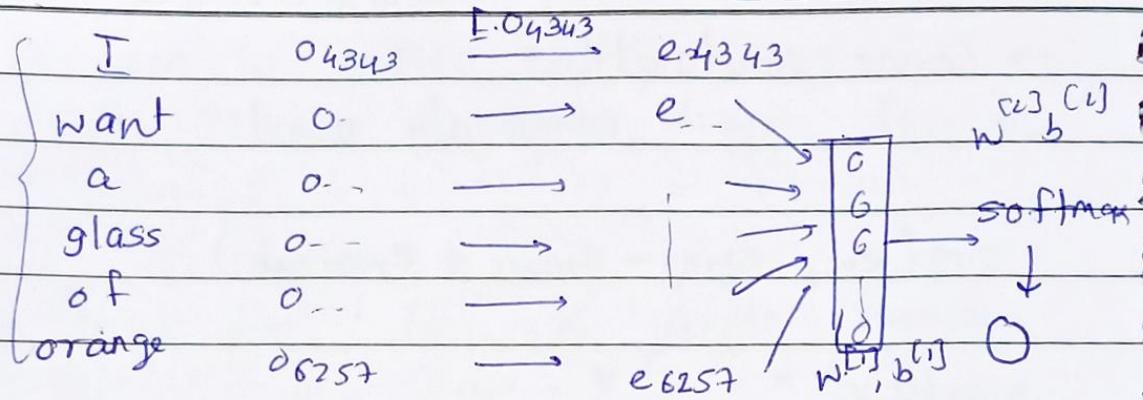
$$\text{sim}(u, v) = \frac{u^T \cdot v}{\|u\|_2 \|v\|_2}$$

 \rightarrow Embedding Matrix. \leftarrow 

$$E \cdot \begin{bmatrix} 0_{6257} \\ (300, 10K) \end{bmatrix} \rightarrow (10K, 1) = \begin{bmatrix} \text{smthg} \\ \text{smthg}' \\ | \\ | \end{bmatrix} = e_{6257} \cdot (300, 1)$$

Date :

MON TUE WED THU FRI SAT SUN



1 ok possibility
for softmax's OLP

Word embeddings

Here, ~~every time~~ the NN will have $6 \times 300 = 1800$ dimensions. But every time we don't need this much.

We can use a fixed word history

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

e.g. of context: Last 4 words,
4 words on L & R, Last 1 word,
nearby 1 word

Date : 28/7/24
MON TUE WED THU FRI SAT SUN

Word2Vec :

- * Embeddings (feature vectors for each & every word) are not hand crafted, instead they are learnt using neural network training.
- * 1. Take a fake problem → like find missing word
2. Solve it using NN.
3. You get word embeddings as a SIDE EFFECT
- * There lived a king called Ashoka.
eg. for step 2.
Training samples:

X (context)	y
lived, a	→ There
a, king	→ lived
king, called	→ a
;	

This is called semi-supervised learning

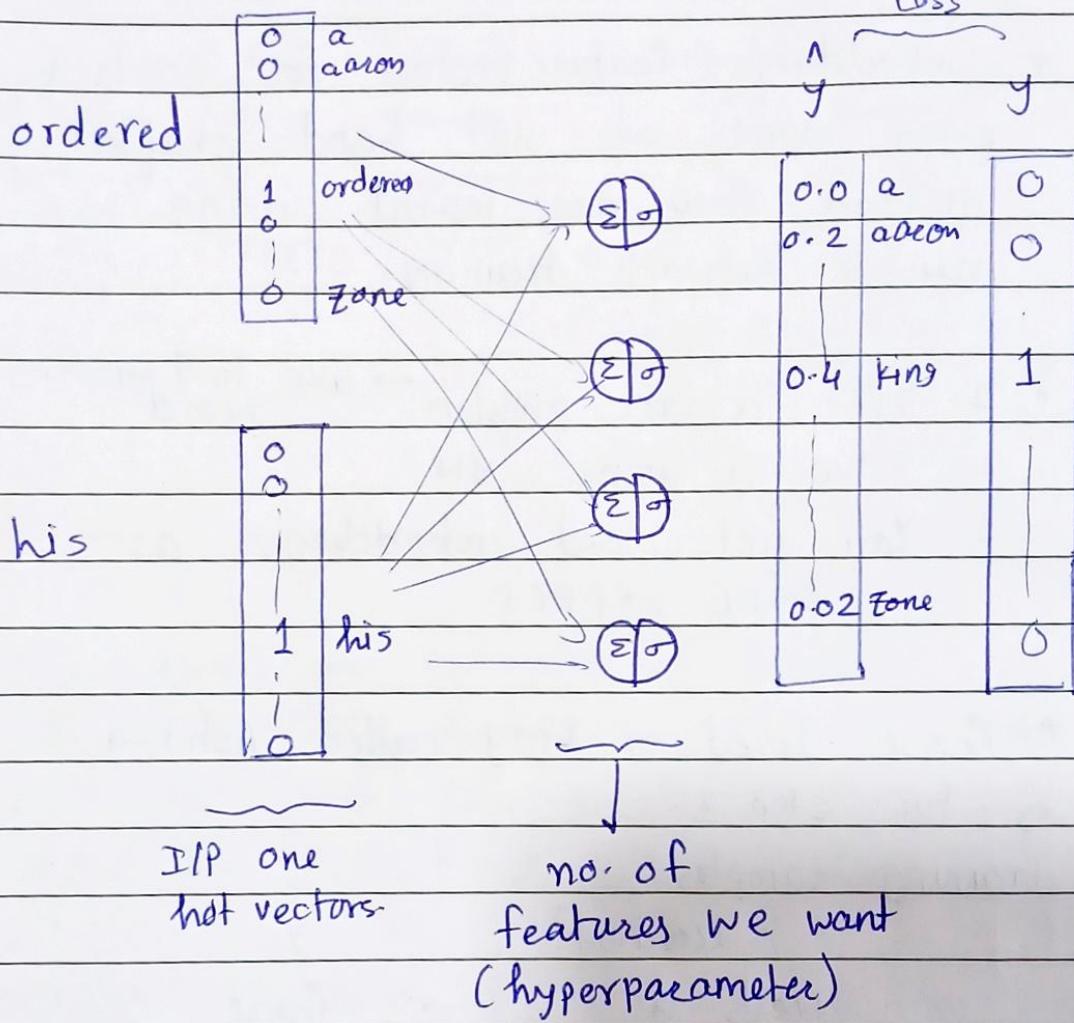
CBOW: Continuous Bag of Words

target ↓ context

Date :

MON TUE WED THU FRI SAT SUN

king ordered his Backprop.



After training the model over all examples, the I/P vector for 'king' would be very similar to that of 'emperor'!

Skip-Gram:

↓
just reverse
the above NN.

King ordered his
 context ↑
 target

For skip-gram:

$$p(t|c) = \frac{e^{\theta_t^T e_c}}{\sum_{j=1}^{10k} e^{\theta_t^T e_c}}$$

→ Negative Sampling ←

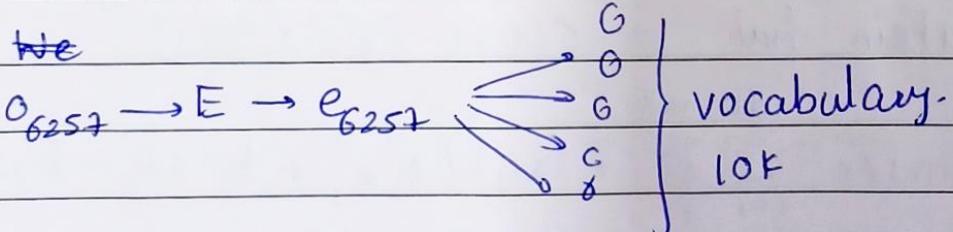
Is there any relation betⁿ context word & another word.

c	x	y	
		word(t)	ty?
	orange	juice	1
K	orange	king	0
	orange	emperor	0

$k = 5-20$ smaller dataset

$k = 2-5$ larger dataset

$$p(y=1 | c, t) = \sigma(\theta_t^T \cdot e_c)$$



Now instead of outputting full vocabulary predictions, we make it to predict on the +ve sample & k -ve samples deliberately.

⇒ like a $10, k$ binary classification prblm.
split into 4-5 binary class everytime.

Date :
MON TUE WED THU FRI SAT SUN

?? How to choose -ve samples ??

- 1) Sample most freq. used words.
- 2) Sample uniformly from vocabulary
- 3) $P(w_i) = \frac{f(w_i)^{3/4}}{\sum_{j=1}^{10k} f(w_j)^{3/4}}$

$f(w_i)$ = observed freq. of a particular word in English.

GLOVe (Global vectors for word repr)

$x_{ij} = \# \text{ times } j \text{ appears in context of } i.$ \uparrow t
 $\uparrow c$

⇒ i & j kitne baar agar-bagal mein or within some $\pm n$ words ke range me aatein hai. ⇒ CLOSE PROXIMITY

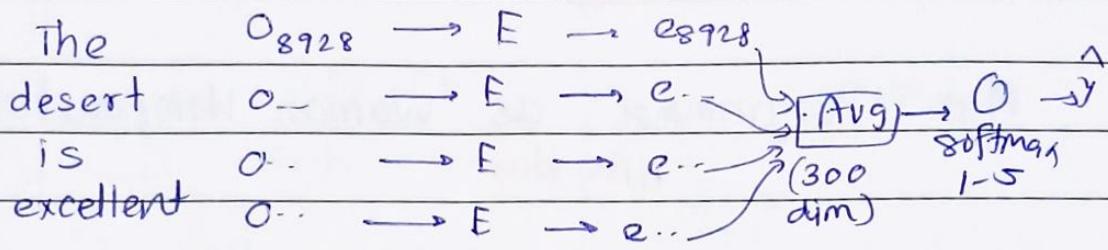
$$\text{Minimize: } \sum_{i=1}^{10k} \sum_{j=1}^{10k} f(x_{ij}) (\theta_i^T e_j + b_i + b_j^T - \log x_{ij})^2$$

Date :

MON TUE WED THU FRI SAT SUN



Sentiment Classification

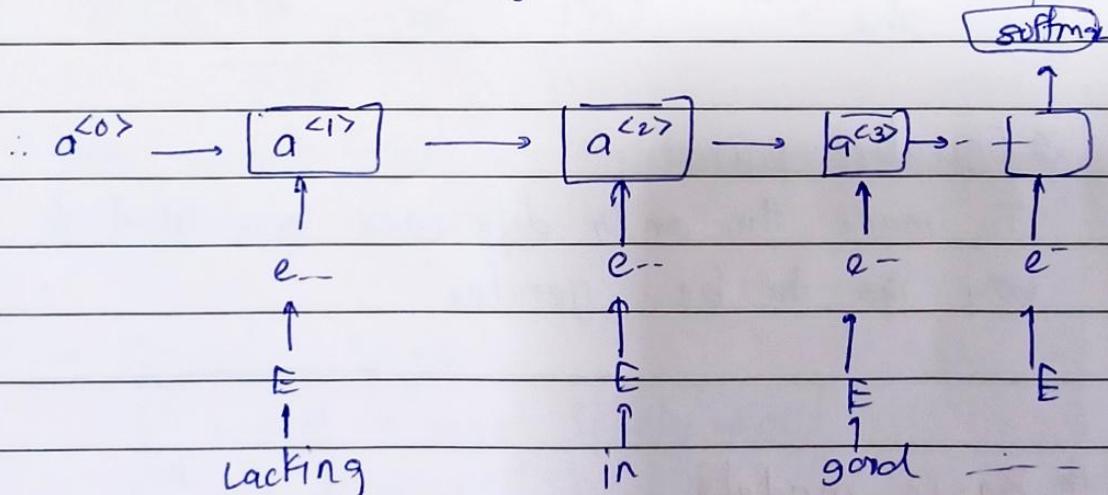


Feature vectors
trained over huge dataset

Instead of just normal summing of vectors, we use RNN.

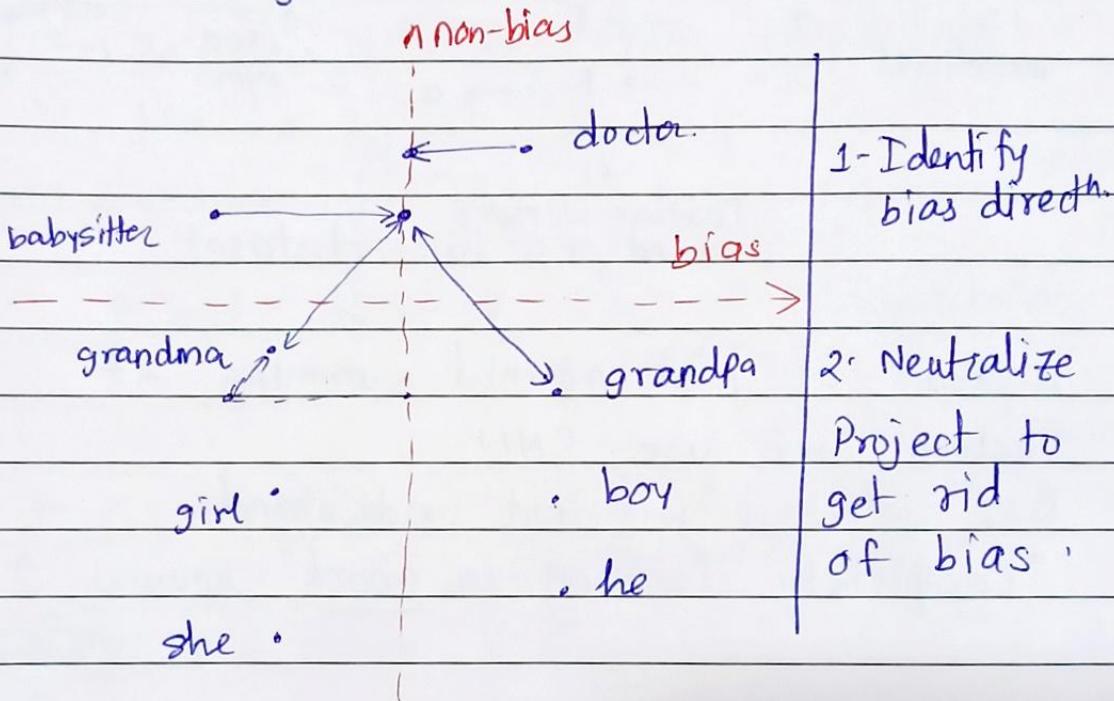
Bcz o/w it will not understand:

"Completely lacking in good" review.



Problem of bias

Man: Programmer as Woman: Homemaker

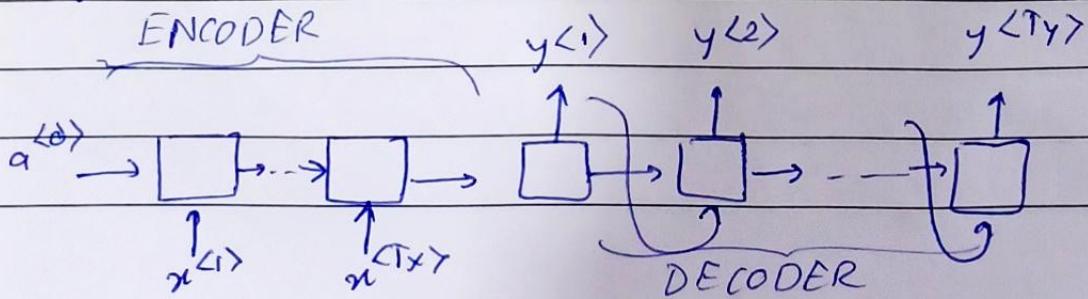


3. Equalize pairs

To make the only difference b/w girl & boy ~~as~~ to be 'gender'!

Basic models

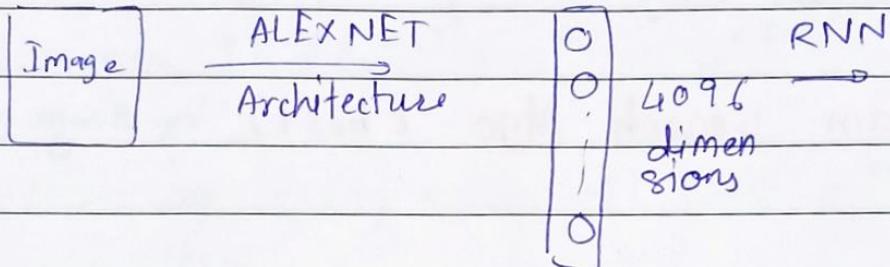
eg. M/c translation.



Date :

MON TUE WED THU FRI SAT SUN

eg. Image captioning



→ Finding Most Likely ←

- Jane is visiting Africa in september
- Jane is going to be visiting . . .
- In sept, Jane will visit Africa

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

$\underbrace{\quad}_{\text{Eng}}$ $\underbrace{\quad}_{\text{French.}}$

We want the y that maximizes
conditional prob

→ Greedy Search:

1st word → most likely word

2nd word → 2nd

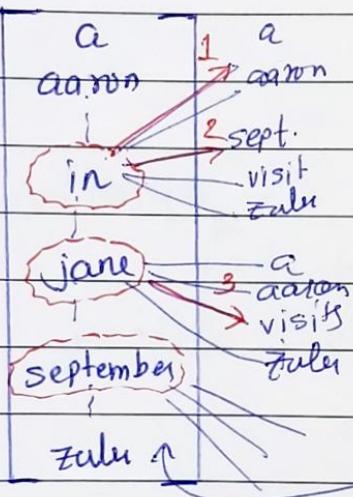
But this approach might output
Jane is going to be - —

→ Approximate Search:

$$\arg \max P(y^{(1)} - y^{(T_y)} | u)$$

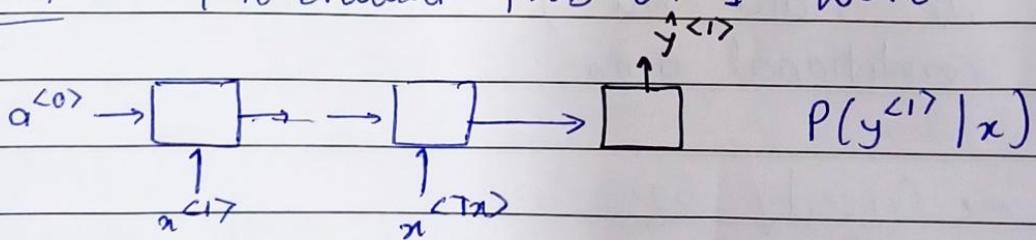
→ Beam Search Algo (best) eg. Beam width $B = 3$

Step-I



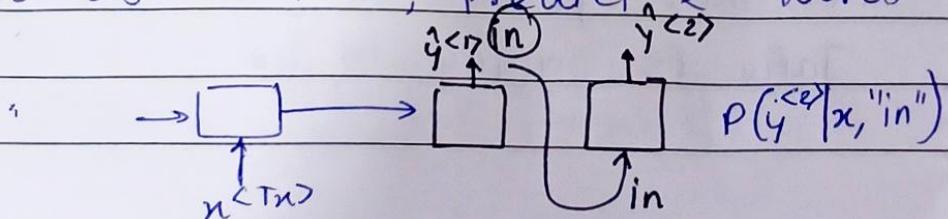
NO top 3 ~~is~~ from September, so it is dropped.

Step-II: Try to evaluate prob. of 1st word.



$B=3 \Rightarrow$ will predict top 3 such words.

Step-III: By hardwiring each of the previous 3 words as 1st word, predict 2nd word.



Date :

MON TUE WED THU FRI SAT SUN

<input type="checkbox"/>						
--------------------------	--------------------------	--------------------------	--------------------------	--------------------------	--------------------------	--------------------------

$$P(y^{(1)}, y^{(2)} | x) = P(y^{(1)} | x) \times P(y^{(2)} | x, y^{(1)})$$

$$\arg \max_y \prod_{t=1}^{T_y} P(y^{(t)} | x, y^{(1)}, \dots, y^{(t-1)}) \quad 29/7/24$$

multiplying many such < 1 nos. = ^{very low} num.
 \Rightarrow underflow

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

But still these probs will be negative overall \Rightarrow Algo will try for shorter translations.

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

$0 < \alpha < 1$ Hyperparam

Bleu Score:

Bilingual Evaluation Understudy

Use:

If there are multiple great answers.

Date : 29/7/24
 MON TUE WED THU FRI SAT SUN

Ref1: The cat is on the mat

Ref2: There is a cat on the mat

MToP: The cat the cat on the mat

<u>Bigrams</u>	Count	Clip count
----------------	-------	------------

① thecat	2	1
② cat the	a	0
③ cat on	1	1
④ on the	1	1
⑤ the mat	1	1

($\frac{5}{6}$)

$$P_n = \frac{\sum_{n\text{-grams}} \text{count clip}(n\text{-gram})}{\sum_{n\text{-grams}} \text{count}(n\text{-gram})}$$

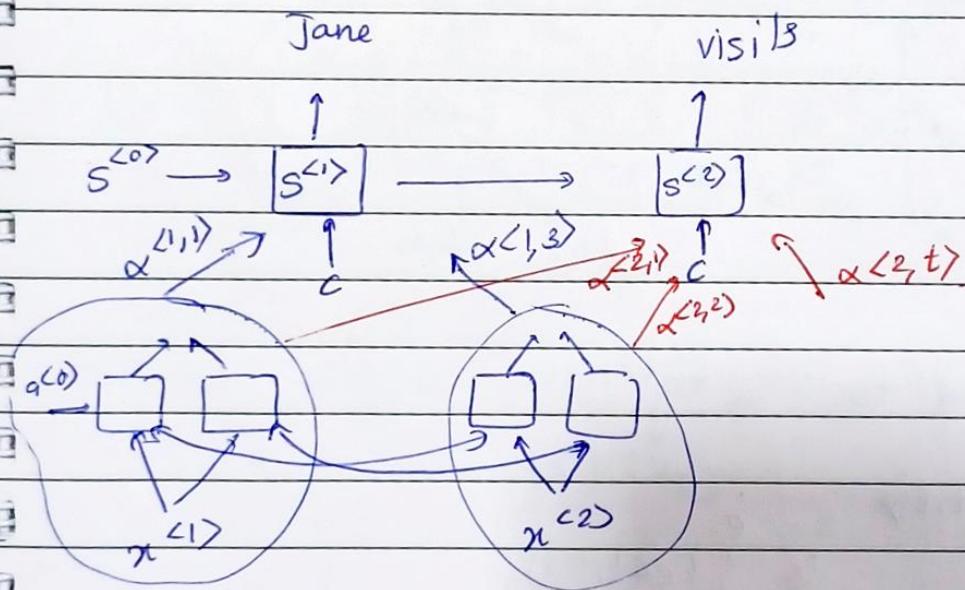
$$\left[\begin{array}{l} \text{Combined BLEU Score} \\ = BP \cdot \exp \left(\frac{1}{4} \sum_{n=1}^4 P_n \right) \end{array} \right]$$

$$\text{Brevity penalty} = \begin{cases} 1 & \text{if MToP length} > \frac{\text{ref. O/P length}}{\text{length}} \\ \exp \left(1 - \frac{\text{ref. O/P length}}{\text{MToP length}} \right) & \text{otherwise} \end{cases}$$

penalizes shorter translations

Attention

The Encoder - Decoder architecture cannot handle long sentence.
 Encoder cannot generate a vector summarizing the full sentence @ a time
 ∵ We need an attention region



$\alpha^{(i,j)}$ are weights/attn wts.

Model: $a^{<t>} = (\vec{a}^{<t>}, \vec{o}^{<t>})$
 concatenation of both.

$\alpha^{<t,t+1>} = \text{amt of attention } y^{<t>} \text{ should pay to } a^{<t+1>}$

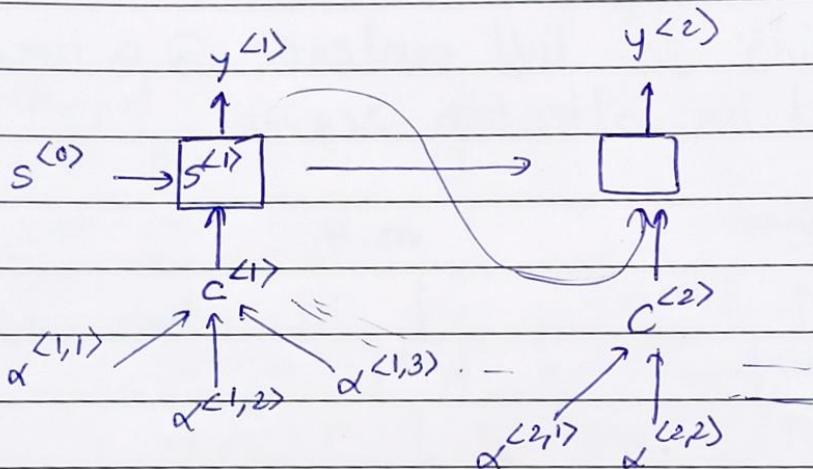
Date :

MON TUE WED THU FRI SAT SUN

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

Sum of all attentions
for a particular $y^{(t)}$
 $= 1$

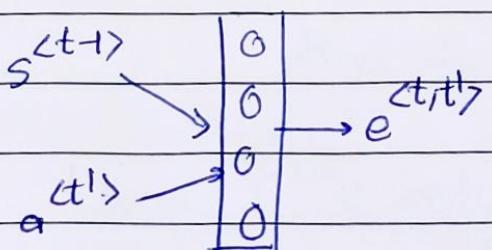
$$c^{(1)} = \sum_{t'} \alpha^{(1,t')} a^{(t')}$$



Formula for $\alpha^{(t,t')}$

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

To compute $e^{(t,t')}$:



Train a very
small NN

→ Takes quadratic time
 $T_x \times T_y$ but acceptable

Previously, phonemes were used
for sp-recog.

Date :

MON TUE WED THU FRI SAT SUN

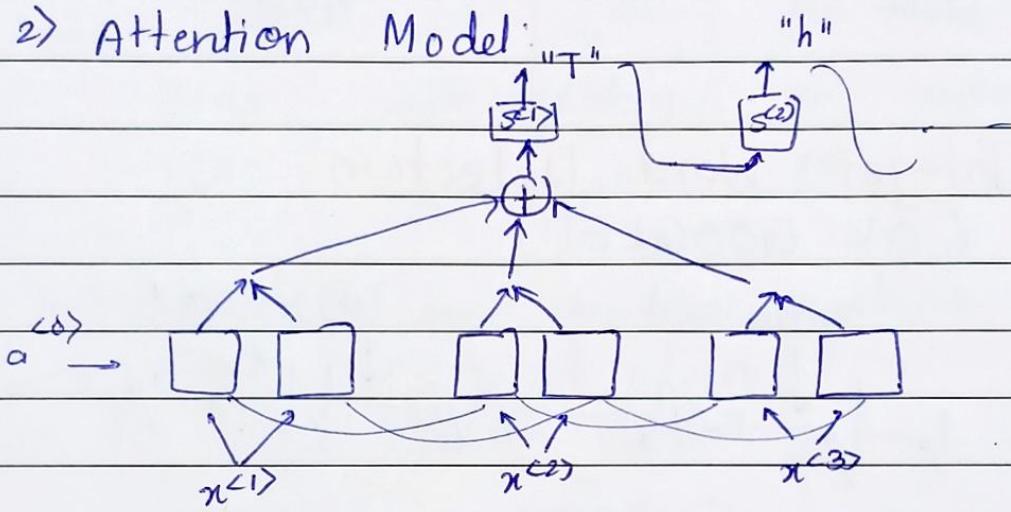
Speech Recognition:

i) CTC cost for speech recogn.

Connectionist Temporal Classif

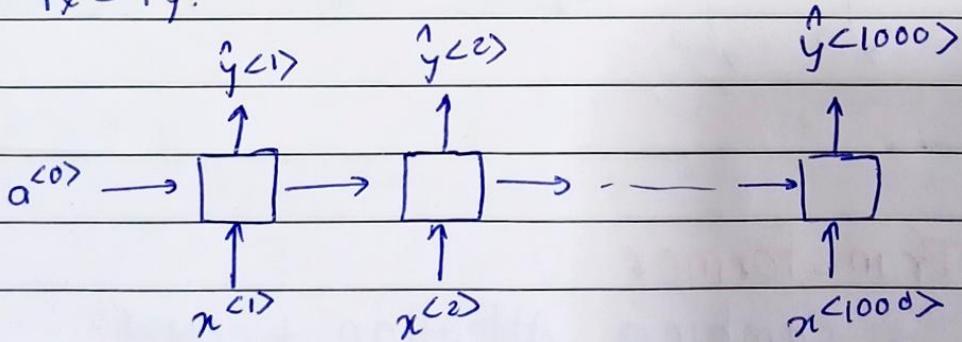
(OR)

2) Attention Model:



i) CTC: Simple RNN ~~or~~ with GRU/LSTM.

$$T_x = T_y.$$



10 sec. audio

Feature rate = 100Hz

⇒ Total of 1000 inputs.

Date :

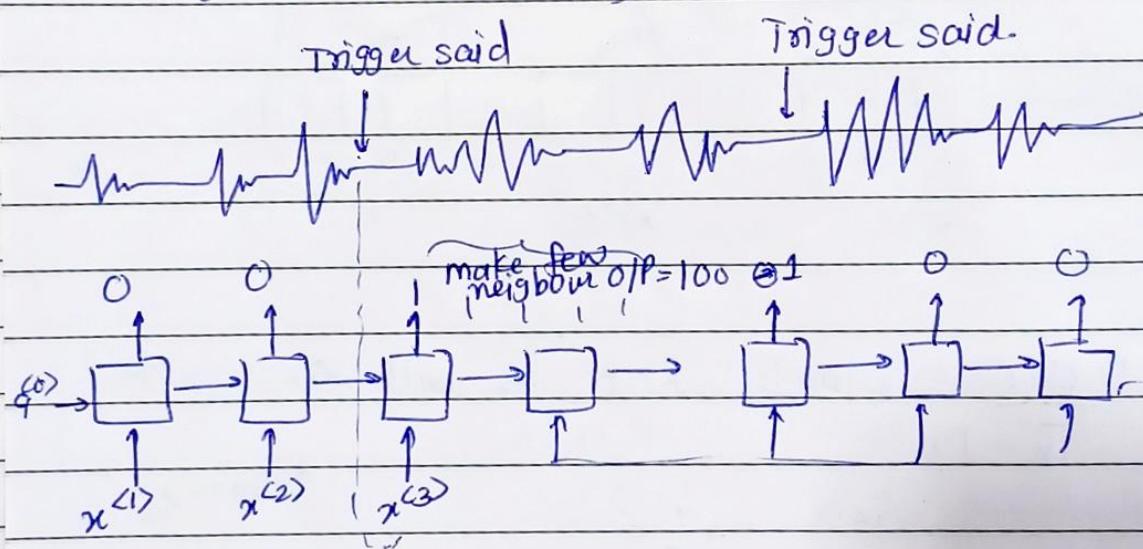
MON TUE WED THU FRI SAT SUN



CTC generates O/P like:

ttt _ -- h - eee - - L --- ggg -
↑ ↑
blank space

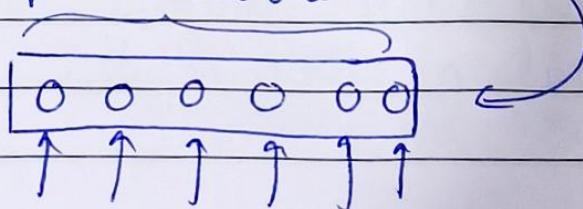
Trigger Word Detection. (OK GOOGLE!)



Transformer

↳ Combines Attention + CNN
style of processing

compute in Parallel



→ Self-attention ←

compute Attn based representation
for each I/P word:

$$A(q, k, v)$$

eg. $x^{(1)}$ $x^{(2)}$ $x^{(3)}$ $x^{(4)}$ $x^{(5)}$
Jane visite l'Afrique en septembre
 $A^{(1)}$ $A^{(2)}$ $A^{(3)}$ $A^{(4)}$ $A^{(5)}$

Depending on how we are thinking
of Africa → historical interest ?,
holiday dest. ?
etc.

We will choose to represent it differently
 $A^{(3)}$ will look at the surrounding
words to get context.

$$A(q, k, v) = \sum_i \frac{\exp(q \cdot k^{(i)})}{\sum_j \exp(q \cdot k^{(j)})} v^{(i)}$$

soft max
—
like

Each input has $q^{(i)}$, $k^{(i)}$ & $v^{(i)}$.

$$\begin{aligned} q^{(3)} &= W^Q \cdot x^{(3)} \\ k^{(3)} &= W^K \cdot x^{(3)} \\ v^{(3)} &= W^V \cdot x^{(3)} \end{aligned} \quad \left\{ \begin{array}{l} x^{(3)} = \text{word emb} \\ \text{of } x^{(3)} \\ \text{l'Afrique} \end{array} \right.$$

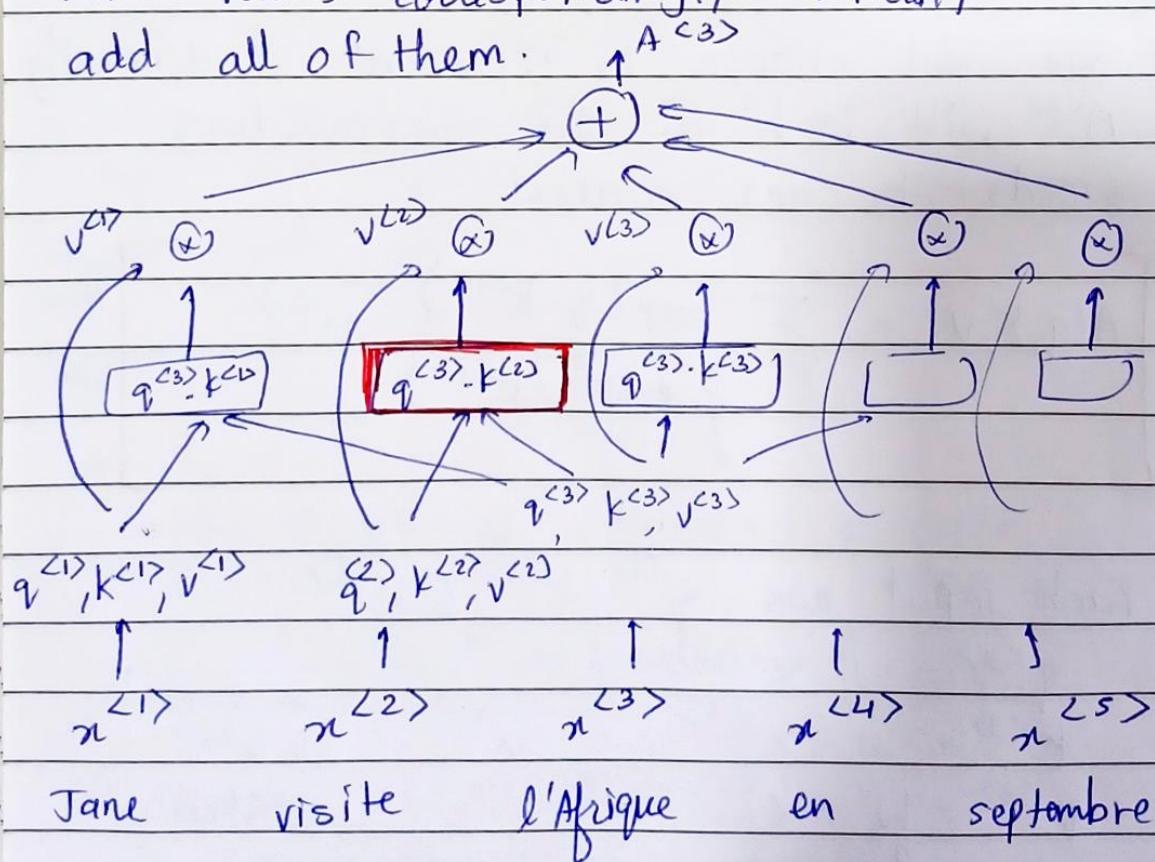
PARAMETER MATRICES

$q^{(3)}$ is a question you can ask to l'Afrique like : "what's happening there?"

↳ Answers are $q^{(3)} \cdot k^{(1)}$, $q^{(3)} \cdot k^{(2)}$
 $q^{(3)} \cdot k^{(3)}$... dot prod

If $k^{(1)}$ was repr. a person & $k^{(2)}$ an action, then $q^{(3)} \cdot k^{(2)} > q^{(3)} \cdot k^{(1)}$

Take softmax of all $q^{(3)} \cdot k^{(i)}$ & multiply them with values correspondingly. Finally add all of them.



Date :

MON TUE WED THU FRI SAT SUN



→ Multi-Head Attention ←

Attention (w_i^Q, w_i^K, w_i^V)

$$W_i^Q, W_i^K, W_i^V$$

$$q^{(1)}, k^{(1)}, v^{(1)}$$

1

$x^{(1)}$

$x^{(2)}$

$x^{(3)}$

$x^{(4)}$

$x^{(5)}$

e.g. For head = 1 → we know what's happening.

For head = 2 → we know when happening?

For head = 3 → who did?

etc.

Multihead(Q, K, V) = concat(head₁, head₂, ..., head_H) $\begin{smallmatrix} \downarrow \\ W_o \end{smallmatrix}$

head_i = Attention (w_i^Q, w_i^K, w_i^V) $\begin{smallmatrix} \downarrow \\ \text{zeros} \end{smallmatrix}$

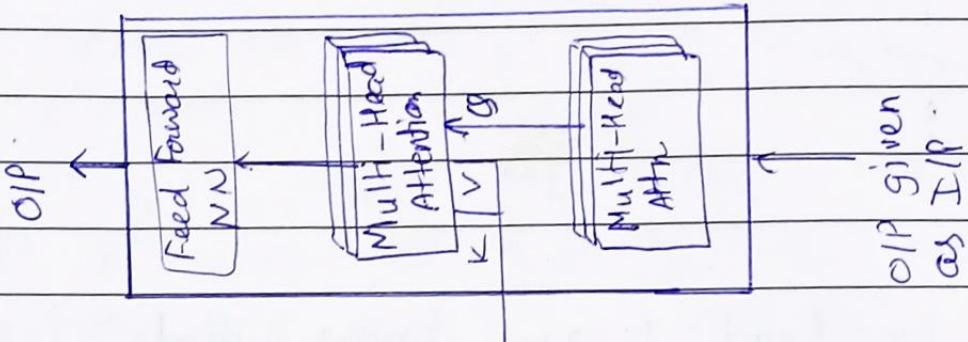
Transformer Network

Date :

MON TUE WED THU FRI SAT SUN



repeat
N times



Positional
Encoding -
<EOS>

Encoder
repeat
N times

<SOS>

