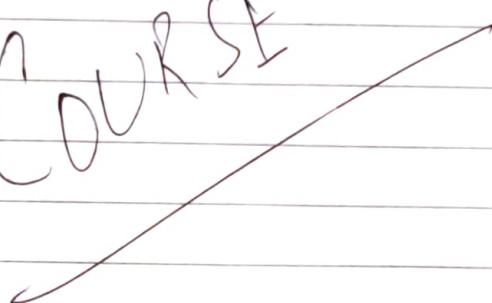


Sequence Modeling

COURSE

05



## SEQUENCED DATA :-

1. Speech Recognition
2. Music Generation
3. Sentiment Classification
4. DNA Sequence Analysis

### Motivating Example :-

Problem: Name entity recognition. Recognize names in below sentences.

$x: \text{Harry Potter \& Hermione Granger invented}$

$x^{(1)} \quad x^{(2)} \quad x^{(3)} \quad x^{(4)} \quad \dots \quad x^{(t)}$

a new spell.  $T_x = 9 = \text{length}$

$x^{(9)}$

$y: 1 \quad 1 \quad 0 \quad 1 \quad 1 \quad 0 \quad 0 \quad 0 \quad 0$

$y^{(1)} \quad y^{(2)} \quad y^{(3)} \quad \dots \quad y^{(9)}$

$$T_y = 9 = \text{length}$$

$x^{(i) \langle t \rangle} = i^{\text{th}} \text{ training sample, } t^{\text{th}} \text{ word}$

$T_x^{(i)} = \text{length of } i^{\text{th}} \text{ training sample.}$

$y^{(i) \langle t \rangle} = t^{\text{th}} \text{ element in } i^{\text{th}} \text{ training sample.}$   
i.e. in OUTPUT SEQUENCE.

$Ty^{(i)}$  = length of  $i^{\text{th}}$  OUTPUT SAMPLE.

Representing Word :=  
   →   

x: Harry Potter & Hermione Granger invented a  
 $x^{(1)} \quad x^{(2)} \quad x^{(3)} \quad x^{(4)} \quad \dots$

Vocabulary:

a	1
aa	2
aaa	3
aaaa	4
aaaaa	5
aaaaaa	6
aaaaaaa	7
aaaaaaa	8
aaaaaaa	9
aaaaaaa	10
aaaaaaa	11
aaaaaaa	12
aaaaaaa	13
aaaaaaa	14
aaaaaaa	15
aaaaaaa	16
aaaaaaa	17
aaaaaaa	18
aaaaaaa	19
aaaaaaa	20
aaaaaaa	21
aaaaaaa	22
aaaaaaa	23
aaaaaaa	24
aaaaaaa	25
aaaaaaa	26
aaaaaaa	27
aaaaaaa	28
aaaaaaa	29
aaaaaaa	30
aaaaaaa	31
aaaaaaa	32
aaaaaaa	33
aaaaaaa	34
aaaaaaa	35
aaaaaaa	36
aaaaaaa	37
aaaaaaa	38
aaaaaaa	39
aaaaaaa	40
aaaaaaa	41
aaaaaaa	42
aaaaaaa	43
aaaaaaa	44
aaaaaaa	45
aaaaaaa	46
aaaaaaa	47
aaaaaaa	48
aaaaaaa	49
aaaaaaa	50

In this example we are using vocabulary of size 50.

$$x^{(1)} \rightarrow \begin{bmatrix} 0 \\ 0 \\ 1 \\ 0 \\ 0 \\ 0 \end{bmatrix} \rightarrow 4075 \quad \text{One-Hot}$$

$$x^{(2)} \rightarrow \begin{bmatrix} 0 \\ 0 \\ 0 \\ 1 \\ 0 \\ 0 \end{bmatrix} \rightarrow 367$$

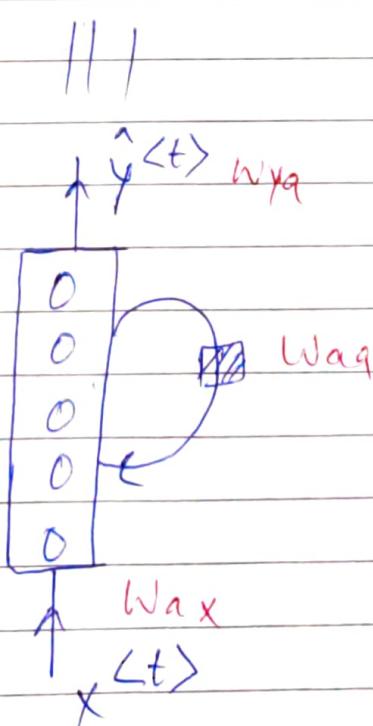
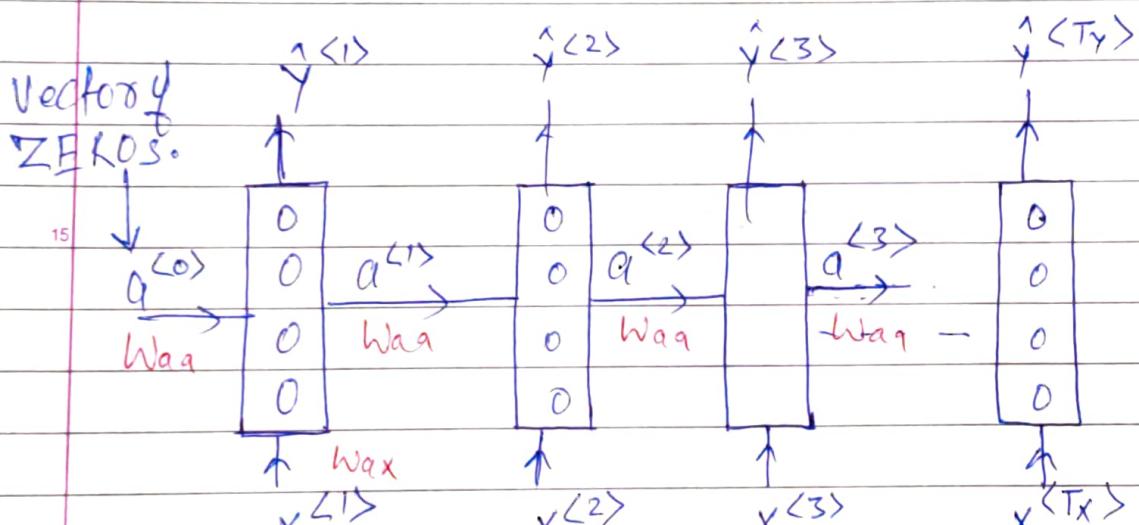
## Why Not Standard Neural Network :-

Problems :-

1. Input/Output can be different lengths in a different examples.

2. Doesn't share features learned across different positions of text.

## Recurrent Neural Network :-



(\*)  $W_{aa}$ ,  $W_{ax}$  &  $W_{ya}$  all are same across the network.

## Let's Apply Activation Functions:-

$$a^{(1)} = g_1(W_{aa} * a^{(0)} + W_{ax} * x^{(1)} + b_a)$$

②  $y^{(1)} = g_2(W_{ya} * a^{(1)} + b_y)$

$g_1 = \tanh/\text{ReLU}$      $g_2 = \text{sigmoid}$

Generalizing:-

$$a^{(t)} = g_1(W_{aa} * a^{(t-1)} + W_{ax} * x^{(t)} + b_a)$$

15     $y^{(t)} = g_2(W_{ya} * a^{(t)} + b_y)$

⊗ Simplifying :-

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

$$W_a = 100 \quad \begin{bmatrix} W_{aa} & W_{ax} \\ 100 & 1000 \end{bmatrix} = W_a \quad (100, 10100)$$

$$[a^{(t-1)}, x^{(t)}] = \begin{bmatrix} a^{(t-1)} \\ x^{(t)} \end{bmatrix} \quad \begin{array}{c} \downarrow 100 \\ \uparrow 1000 \end{array} \quad \begin{array}{c} \uparrow 100 \\ \downarrow 1000 \end{array} = \underline{\underline{10100}}$$

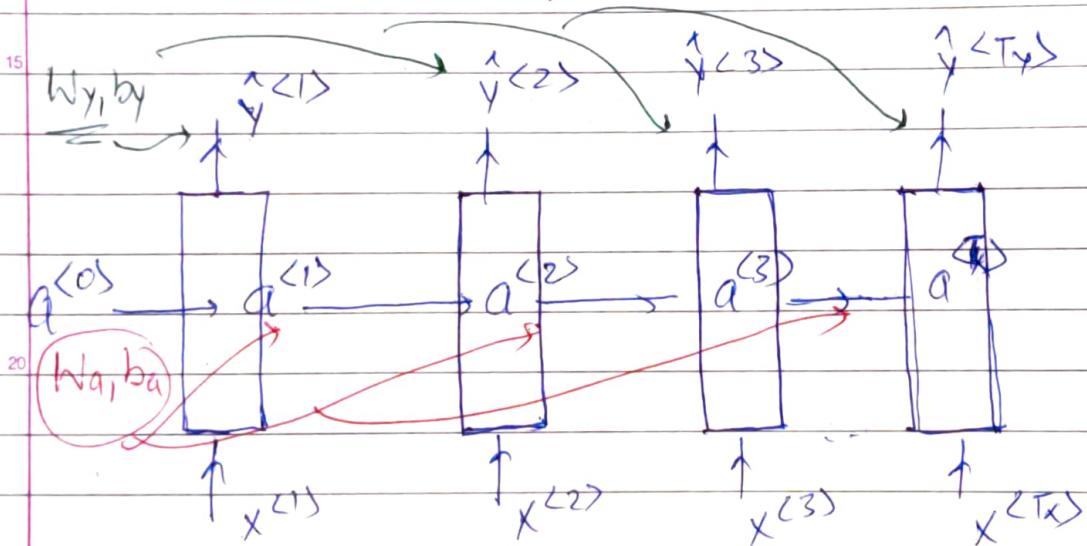
30    We can check

$$[W_{aa} \mid W_{ax}] \begin{bmatrix} a^{(t-1)} \\ x^{(t)} \end{bmatrix} = W_{aa} * a^{(t-1)} + W_{ax} * x^{(t)}$$

In similar way :-

$$\hat{y}^{(t)} = g_2(W_y a^{(t)} + b_y)$$

### Back Propagation RNN :-



$(W_a, b_a) \quad (W_y, b_y)$  weights are shared

Elementwise loss function

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

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

OVERALL LOSS.

## Different Types RNN

$T_x = \text{Input length} | T_y = \text{Output length}$ .

$T_x = T_y = \text{many to many}$

$T_y = 1, T_x > 1 = \text{many to one (Sentiment Class)}$

$T_x = 1, T_y > 1 = \text{One to many}$

(music generation)

## [Language Model]

### [Speech Recognition]

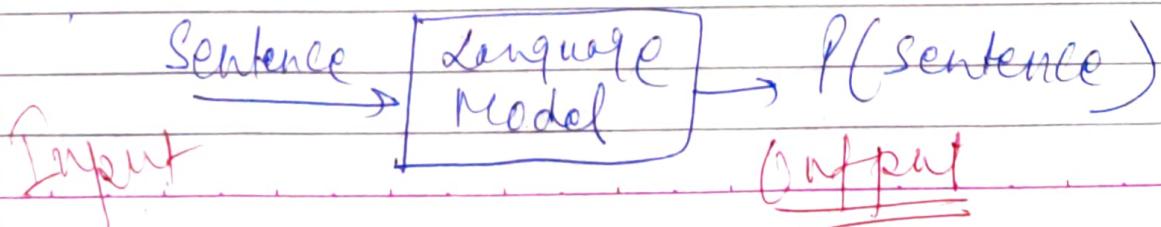
S1: The apple & pear salad.

S2: The apple & pear salad.

$$P(S1) = 3.2 \times 10^{-15}$$

$$P(S2) = 5.7 \times 10^{-10}$$

S2 is much more likely.



END of Sentence.

Cats average 15 hrs of sleep a day. <EOS>

$y^{(1)}$   $y^{(2)}$   $y^{(3)}$   $y^{(6)}$   $y^{(9)}$

The Egyptian May is a board of cat. <EOS>

$\uparrow$   $\text{UNK}$  (Unknown)

$P(a) P(\text{aaron}) \approx P(\text{cats})$

$\uparrow$   $y^{(1)}$   $P(\text{UNK})$   
-P(EOS)

$\uparrow$   $y^{(2)}$   $P(\_ | \text{cats})$

$P(\_ / \text{cats}, \text{avg})$

$a^{(1)}$

$a^{(0)} = 0$

$a^{(2)}$

$y^{(1)}$

$a^{(3)}$

$y^{(2)}$

$\uparrow$   $y^{(9)}$   
 $\uparrow$   
 $x^{(9)}$

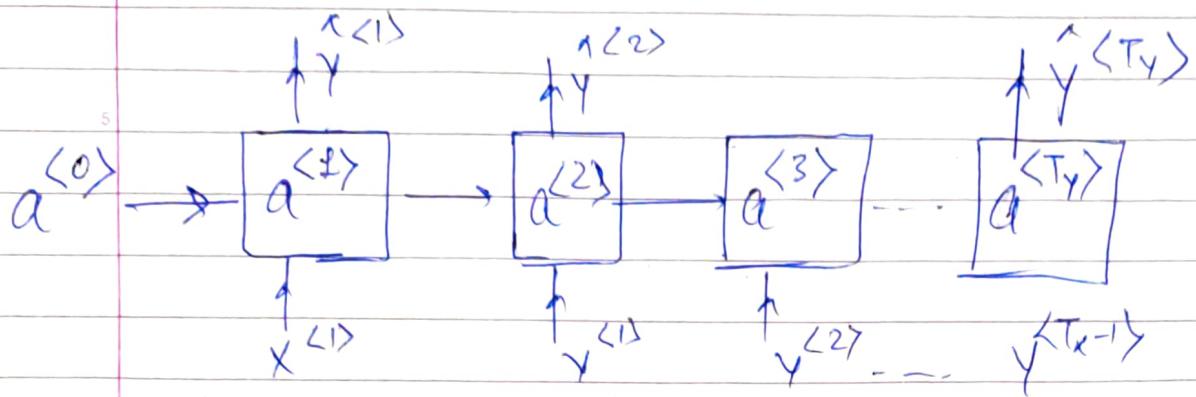
$= y^{(8)} = \text{day}$

Cost function

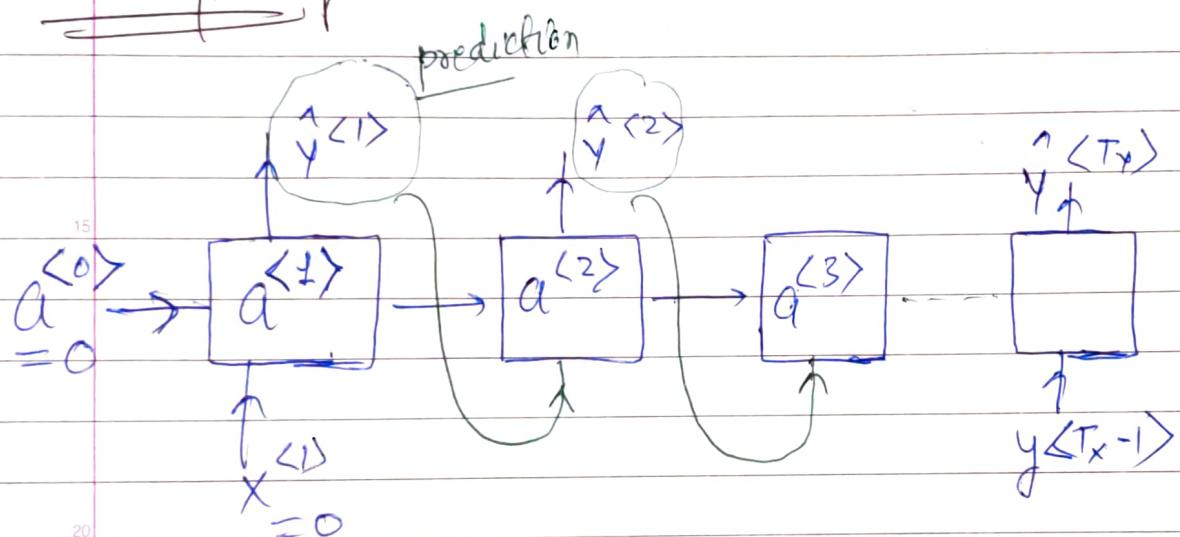
$$\mathcal{L}(\hat{y}, y^{(t)}) = - \sum_i y_i^{(t)} \log \hat{y}_i^{(t)}$$

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

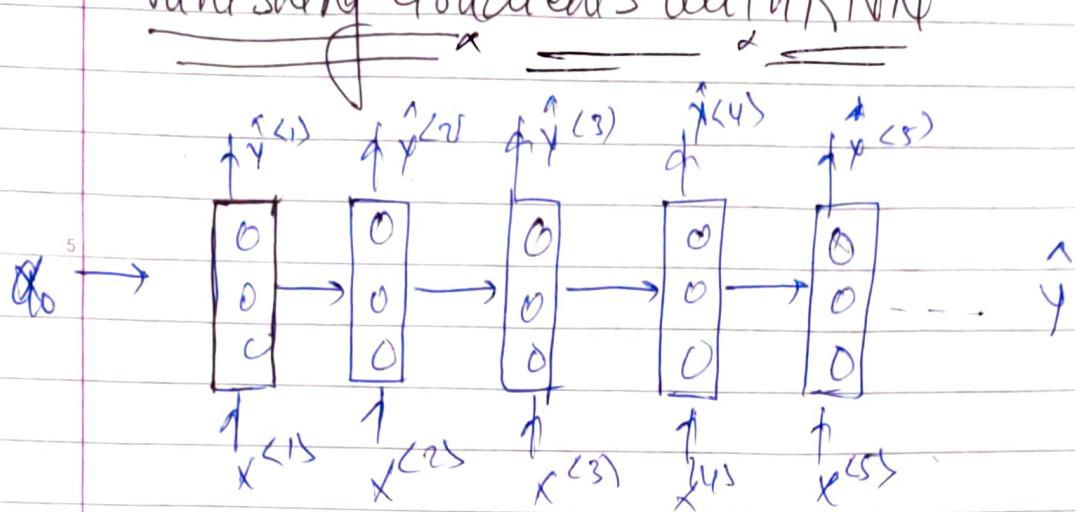
# Sampling a sequence from trained RNN



Sampling :-



## Vanishing Gradients with RNN



Problem :- When the Gradients are being propagated back in time all the way to the initial layer. The gradients coming from the deeper layer have to go through continuous MATRIX MULTIPLICATIONS because of

the ~~chain~~ Chain Rule. If they have small values ( $< 1$ ) they shrink exponentially until they vanish & make it impossible that model can learn.

→ Gradient got Vanish. (Vanishing Gradient)

In the same way if (values  $> 1$ ) multiplying over the time it become NaN.

Gradient Explode (Exploding Gradient).

## Solution :- Gradient Clipping

$\text{g} \leftarrow \begin{cases} \text{g} & \text{if } \|\text{g}\| > \text{threshold} \\ \frac{\text{threshold} \times \text{g}}{\|\text{g}\|} & \text{otherwise} \end{cases}$ 
] Exploding  
Gradient

## Vanishing Gradient Solution :-

### Redesign RNN :

The standard RNN operates in such a way that hidden state activation are influenced by the other local activation closest to them correspond as

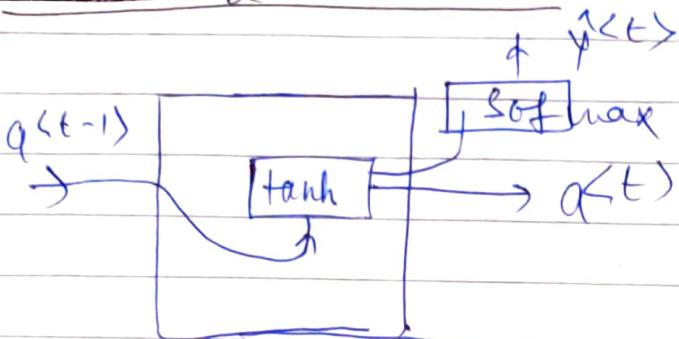
"Short Term Memory", while network weights are influenced by the computations that take over entire

long sequences correspond to "long-term memory". Hence RNN redesign known as "LSTM". "GRU" etc.

Camlin Page  
Date / /

-: GRU :- [ Simplified ]  
[ Version ]

RNN Unit look like :-



GRU :-

$C$  = memory cell

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

$$\tilde{C}^{t-1} = \tanh(W_C [C^{t-1}, X^{t-1}] + b_C)$$

Gate name GAMMA\_U  $\Rightarrow \Gamma_u = \underline{(0 \text{ or } 1)}$

in practice  $\Gamma_u$  = sigmoid function which is most of the time 0 or 1

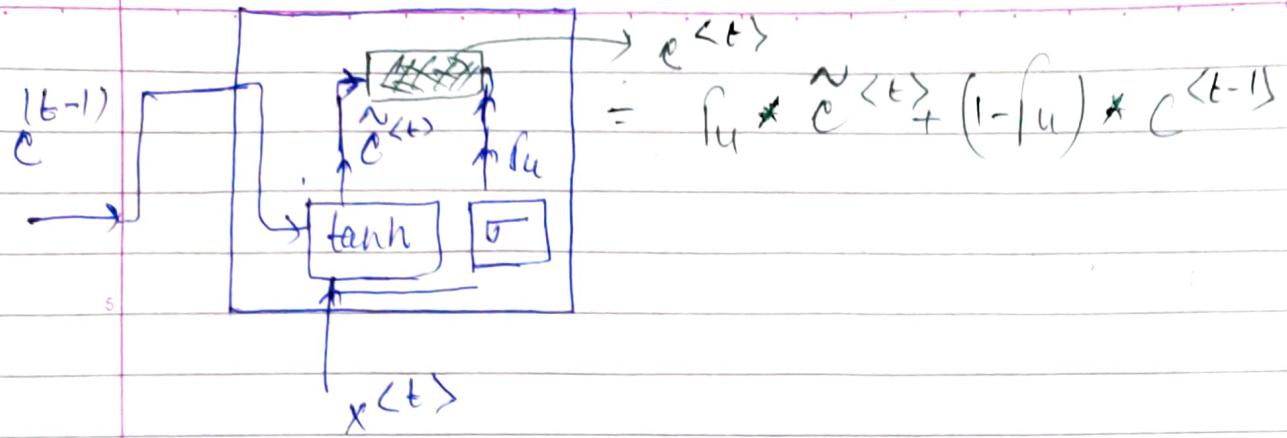
$$\Gamma_u = \sigma(W_u [C^{t-1}, X^{t-1}] + b_u)$$

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

$\hookrightarrow$  element wise multip

it says  $\Gamma_u = 1$  update gate by  $\tilde{C}^{t-1}$

$\Gamma_u = 0$  take old value (Not Update).



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

Full GRU := Full GRU introduce one more SIGMOID GATE.

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

$f_r$  = Intuition of  $f_r$  is tell how much relevance of  $c^{(t-1)}$  earlier it was full.

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

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

↳ Element wise multiply.

= LSTM =

Putting LSTM same as GRU in equation.

$$5 \quad \hat{c}^{(t)} = \tanh(W_c [a^{(t-1)}, x^{(t)}] + b_c)$$

$$10 \quad f_u = \sigma(W_u [a^{(t-1)}, x^{(t)}] + b_u)$$

~~Note~~ Unlike GRU where we ~~had~~ ( $f_u + 1 - f_u$ )

we have two gates more

$$15 \quad f_f = \sigma(W_f [a^{(t-1)}, x^{(t)}] + b_f)$$

↳ forget gate.

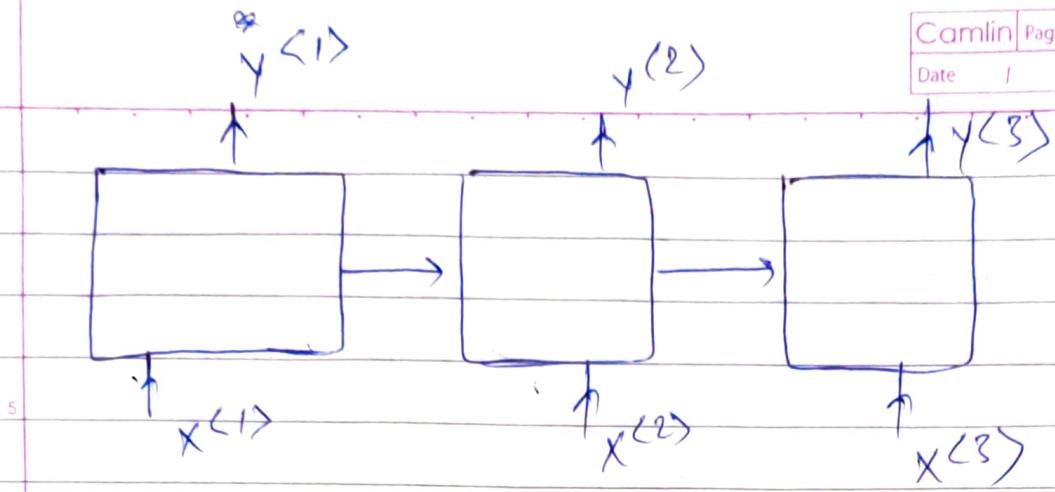
$$20 \quad f_o = \sigma(W_o [a^{(t-1)}, x^{(t)}] + b_o)$$

↳ output gate

$$25 \quad c^{(t)} = f_u * \hat{c}^{(t)} + f_f * c^{(t-1)}$$

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

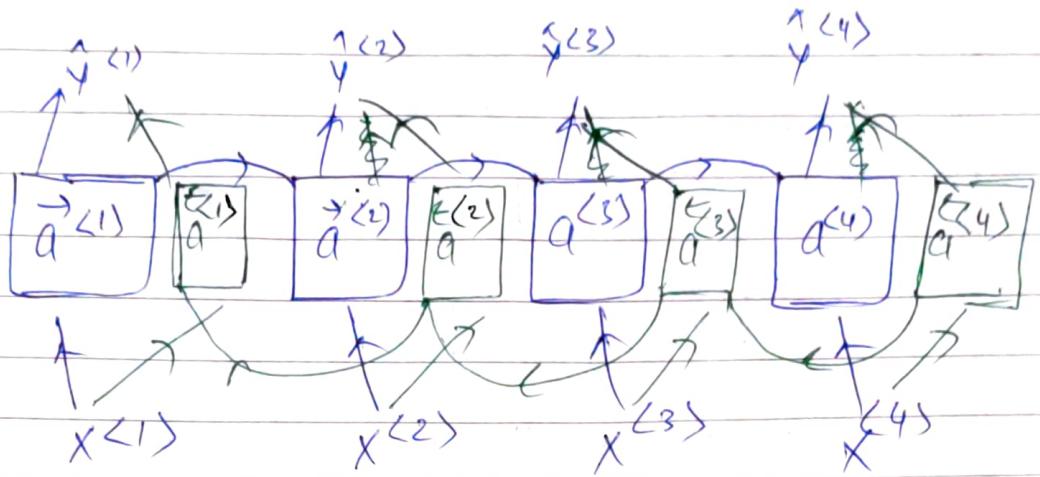
30 All these equations is one LSTM Unit.



① Advantage of GRU simple model so much bigger network can be used. - - ANDREW NG.

But LSTM is morefull. LSTM is more generally default choice,

## Bi-Directional-RNN :=



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

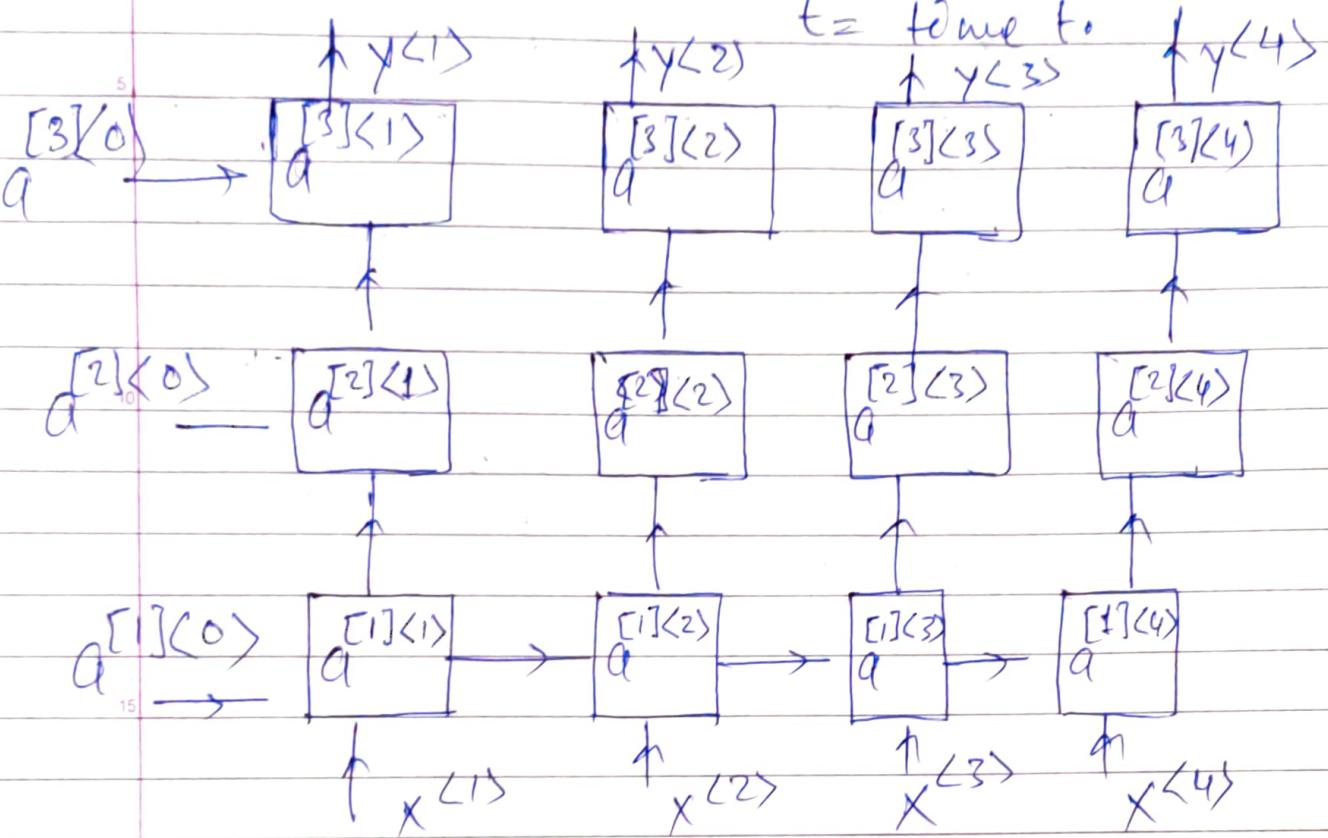
# Deep RNN

~~RECURSIVE~~

$a^{[l]} \leftarrow t$

$l = l^{\text{th}} \text{ layer}$

$t = \text{time step}$



20

25

30

# WEEK-02

Camlin	Page
Date	/ /

## [INTRO WORD EMBEDDINGS]

$$V = [a, \text{aaron}, \dots, \text{zulu}, \text{UNKS}] \quad \|V\| = 10k$$

$\mathbf{t}$  = hot representation ( $\text{Man} = 5391$ ) ( $\text{Woman} = 9853$ )

$$\hookrightarrow \text{Man} = \begin{bmatrix} 0 \\ 0 \\ 0 \\ 0 \\ 0 \\ 5391 \\ 0 \end{bmatrix} \quad \text{Woman} = \begin{bmatrix} 0 \\ 0 \\ 0 \\ 0 \\ 0 \\ 9853 \\ 0 \end{bmatrix}$$

feature representation: Word Embedding

	Man	Woman	King	Queen	Apple	Orange
Gender	-1	1	-0.95	0.97	0.00	0.01
Royal	0.01	0.02	0.93	0.95	-0.01	0.00
Age	0.03	0.02	0.07	0.69	0.03	0.02
Food	0.04	0.01	0.02	0.01	0.95	0.97
Size	!	!	!	!	!	!
lost	300	300	300	300	Orange & Apple feature look pretty & will.	Apple feature
	e5391					

## Transfer Learning Word Embeddings :-

1. Learn word embeddings from large text corpus - (1-100 Billion words)
2. Transfer embedding to new task with smaller training set. (100k words)
3. Optional:- Continue to fine out the word embeddings with new data

### - : Analogies :-

	(5391)	(9853)	(4914)	(7157)	(456)	6757
Gender	-1	1	-0.95	-0.97	0.00	0.01
Royal	0.01	0.02	0.93	0.95	0.01	0.00
Apple	0.03	0.02	0.70	0.69	0.03	-0.02
Food	0.09	0.01	0.02	0.01	0.95	0.97
	e5391	e9853	e4914	e7157	e456	e6757

$$e_{\text{man}} - e_{\text{woman}} \approx \left[ \begin{array}{c} -2 \\ 0 \\ 0 \end{array} \right]$$

$$e_{\text{king}} - e_{\text{queen}} \approx \left[ \begin{array}{c} -2 \\ 0 \\ 0 \end{array} \right]$$

(Mikolov et. al. 2013)

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

Find word  $\arg \max_w \text{sim}(e_w, e_{\text{king}} - e_{\text{man}} + e_{\text{woman}})$

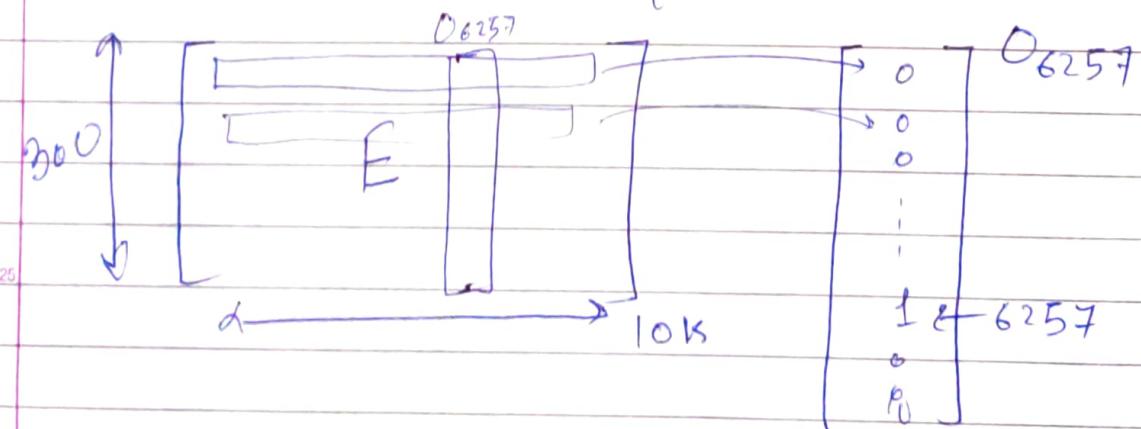
t-SNE ( $\mathbb{R}^N \rightarrow \mathbb{R}^2$ )

$$\text{sim}(u, v) = \frac{u^T v}{\|u\|_2 \|v\|_2} \quad [\text{Cosine Similarity}]$$

[Learning Embedding Matrix]

Learn Word Embedding matrix

a aaron - - orange - - zulu (UNK) <sup>e(6257)</sup>



$$E \cdot O_{6257} = \begin{bmatrix} \boxed{\text{X}} \\ \boxed{\text{X}} \\ \boxed{\text{X}} \\ \boxed{\text{X}} \end{bmatrix} = e_j = \text{embedding for word } j.$$

We initialize  $E$  randomly & use Gradient Descent to learn the parameter of  $E$  ( $300 \times 10000$ ) dimension matrix.

Specific ALGO Learn matrix  $E$ .

[Neural Language Model]

I want a glass of orange — .  
 (4343) (9665) (1) (3852) (6163) (6257)

(One-hot) (Embedding Mat)

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

want  $o_{9665} \rightarrow E \rightarrow e_{9665} \rightarrow \begin{matrix} 0 \\ 0 \\ 0 \\ 0 \\ 0 \\ 0 \end{matrix}$  10K  
 a  $o_1 \rightarrow E \rightarrow e_1 \rightarrow \text{softmax}$

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

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

orange.  $o_{6257} \rightarrow E \rightarrow e_{6257}$

$w^{(1)} b^{(1)}$   
Neural Net

Hyper parameters last k-words  $k=4$

int this case input to neural net

$$4 \times 300 (\text{Stack word}) = \underline{\underline{1200}}$$

know use backprop to maximize the prediction. This way we will learn good E (Embedding matrix).

Hyper param  $k$  = context

- last 4 words → last 1 word
- 4 words on left & right.

Word 2 Vec

Simpler & less computational way to learn (E) embedding word matrix.

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

choose random word say = orange then in  $\pm 5$  vicinity form below pair

context      Target

orange      juice

orange      glass

orange      my

orange      go

## [Skip Gram Model]

Camlin	Page
Date	/ /

Vocab size = 10K

Context  $c$  ("orange")  $\rightarrow$  Target ("juice")

Learn this mapping.

$$o_c \xrightarrow{E} e_c = (E * o_c) \rightarrow \text{Softmax Unit} \rightarrow \hat{y}$$

$$\text{Softmax} = p(t|c) = \frac{e^{\theta_t^T * e_c}}{\sum_{j=1}^{10K} e^{\theta_j^T * e_c}}$$

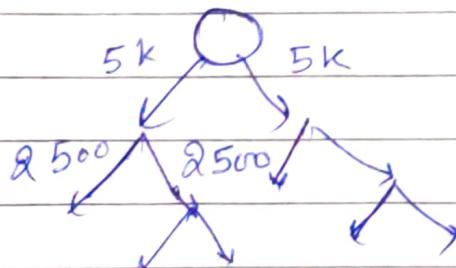
$$L(\hat{y}, y) = - \sum_{i=1}^{10K} y_i \log \hat{y}_i = \text{Loss function}$$

$$Y = \begin{bmatrix} 0 \\ 1 \\ 1 \\ 0 \end{bmatrix} \in \mathbb{R}^4 \quad Y_2 = \begin{bmatrix} 0 \\ 1 \\ 0 \\ 0 \end{bmatrix}$$

10K one hot

④ Problem: Softmax is very slow to compute

⇒ Use Hierarchical Softmax classifier.



## Negative Sampling :=

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

[look around ± 5]

Context	Word	Target
orange	juice	1
orange	king	0

15 dictionary for -ve words.

↑ pick random word from

orange	book	0	k times -ve example.
orange	the	0	$k = 5 \approx 20$
orange	of	0	smaller Data set
X		Y	$k = 2, 5$ larger Data set.

25 ~~softmax~~  $P(Y=1 | c_i, t) = \sigma(\theta_i^T e_c)$

every +ve example we have k -ve example

we train only on  $(k+1)$  samples  
a Neural Network. [  $k$  negative,  $1$  positive ]

## Glove Word Vectors

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

Till now we were choosing  $C_{it}$  in close  
vicinity  $\pm 5$ . But for GLOVE,

$x_{ij} = \text{times } i \text{ appears in context of } j$   
 $\uparrow \quad \uparrow$   
 $C_t \quad T(c)$

$x_{ij} = \text{is count how often words appear}$   
close to each other.