

# Course 5: Sequence Models

## Week 1 : RNNs

1) why RNNs (Supervised learning:  $X \rightarrow Y$ )  
 $G$  (as training set)

eg) Speech recognition

Audio  $X \rightarrow$  Text transcript  $\downarrow$   
sequence data (both i/p & o/p)

eg) Music generation

eg) Sentiment classification

$\rightarrow$  "there is nothing to like in this movie"  $\rightarrow 1\star/5\star$  (rating)

eg) DNA sequence analysis

eg) Machine translation

eg) Video activity recognition

"dip of a race of usain bolt"  $\rightarrow$  Running

eg) Name entity recognition

Yesterday Harry Potter met Hermione Granger  $\rightarrow$  Yesterday Harry Potter met Hermione Granger

## 2) Notation

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

$n$ : Harry Potter and Hermione Granger invented a new spell  $\in \{q\}$

$T_n = 9$

$y$ :  $( \quad | \quad 0 \quad | \quad 1 \quad | \quad 0 \quad 0 \quad 0 )$

$y^{(9)}$

$T_y = 9$

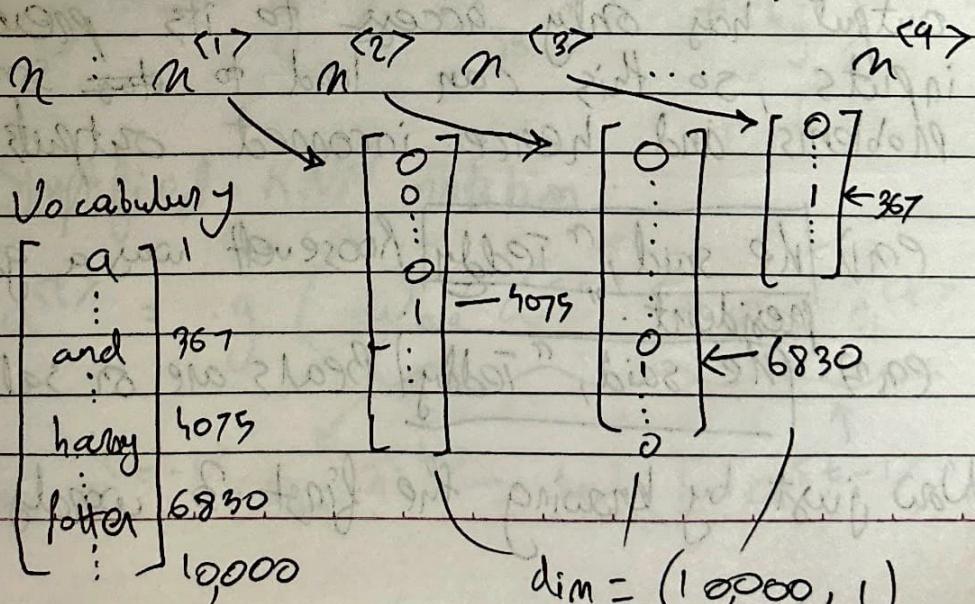
$\boxed{n^{(i)} < t>}$

$T_n^{(i)}$  → input sequence length for training example (i).

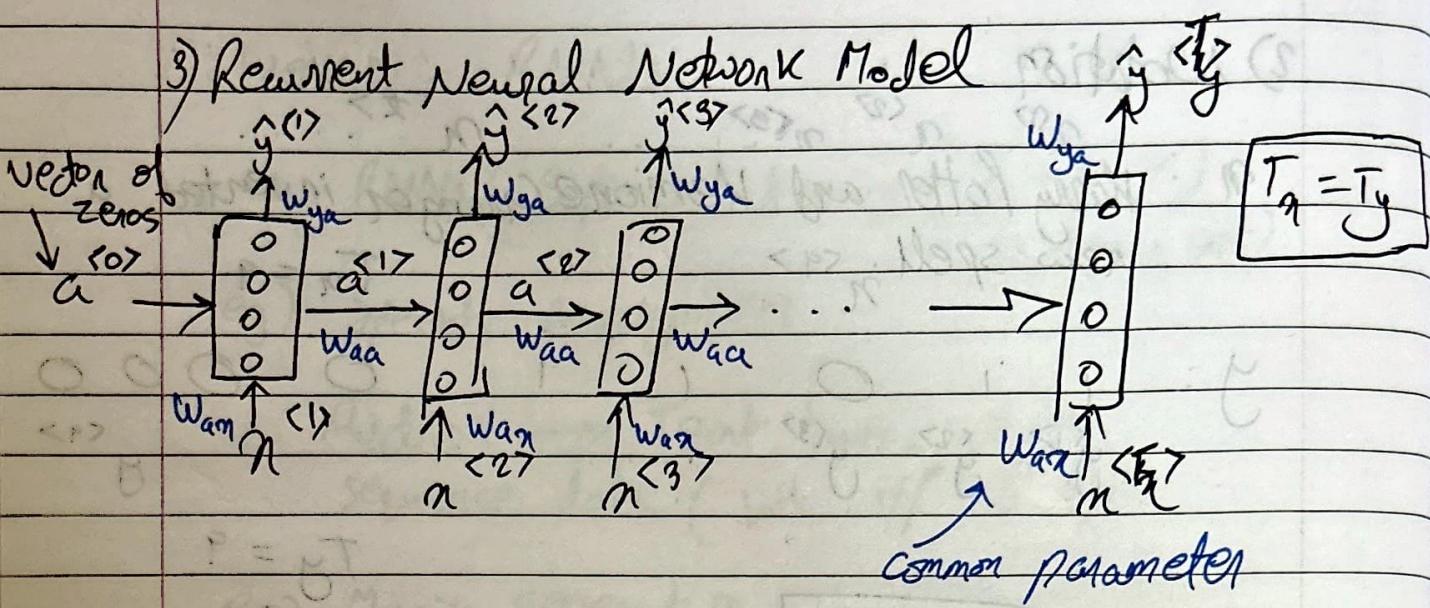
$y^{(i)} \leftarrow$  →  $t$ th element in the output segment for the  $i$ th example

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

## Representing words (One-hot representation)

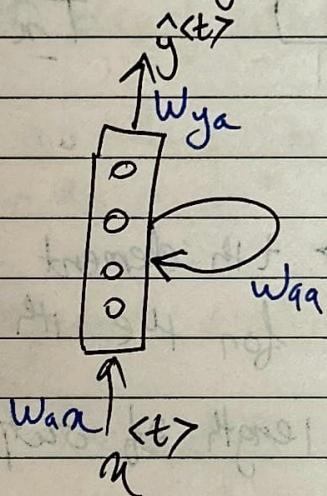


### 3) Recurrent Neural Network Model



\* RNN scans data from left to right

Rolled diag:



Drawback: while predicting  $\hat{y}^{(t)}$ , the output has only access to its previous inputs, so this can lead to ~~stop~~ some problems and hence incorrect outputs. for ex:

(ex 1) He said, "Teddy Roosevelt was a great president."

(ex 2) He said, "Teddy Bears are on sale!"

Now just by knowing the first 3 words of the

Page No. \_\_\_\_\_

$$a \leftarrow w_{an} \cdot n^{<1>} \\ \begin{matrix} \text{used to} \\ \text{compute a like quantity} \end{matrix}$$

above sentences we cannot classify 'Toddy' as a name ~~or~~ of a person or not.

This issue can be resolved using BRNNs  
i.e. Bidirectional RNNs.

Forward Propagation of RNN :

$$a^{<0>} = \vec{o} \quad | \quad a^{<1>} = g(w_{aa} a^{<0>} + w_{an} n^{<1>} + b_a) \\ \hat{y}^{<1>} = g(w_{ya} a^{<1>} + b_y) \\ \text{(maybe some other activation f^n)}$$

$g_1 \rightarrow \text{tanh} / \text{ReLU}$  in  $a^{<1>}$

$g_2 \rightarrow \text{sigmoid}$  in  $\hat{y}^{<1>}$

In general,

hidden state  $a^{<t>} = g(w_{aa} a^{<t-1>} + w_{an} n^{<t>} + b_a)$

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

Simplified RNN notation :

$$a^{<t>} = g(w_{aa} a^{<t-1>} + w_{an} n^{<t>} + b_a)$$

$$w_{aa} a^{<t-1>} + w_{an} n^{<t>} = w_a [a^{<t-1>}, n^{<t>}]$$

$$100 \downarrow [w_{aa}: w_{an}] = w_a$$

$\longleftrightarrow$

$$100 \quad 10000 \quad (100, 10/100)$$

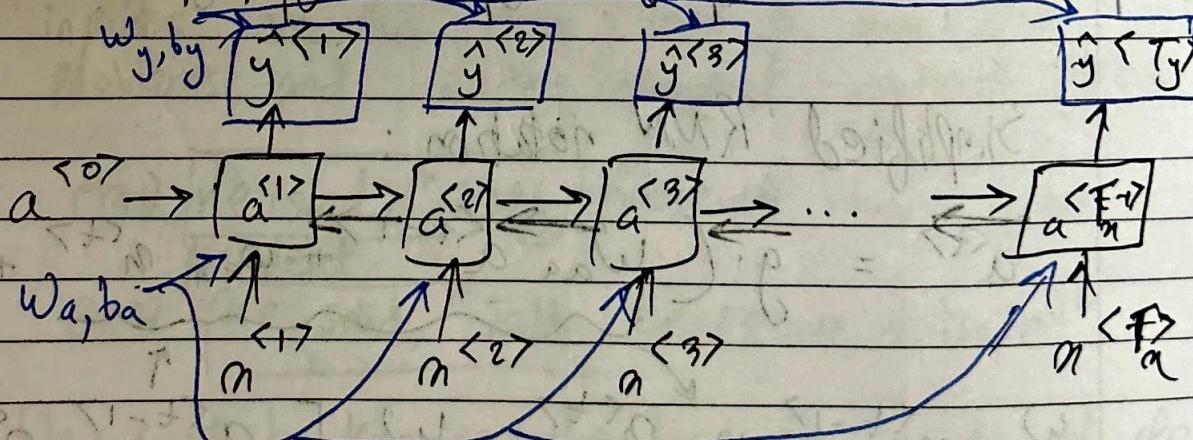
$$[a^{<t-1>}, n^{<t>}] = \begin{bmatrix} a^{<t-1>} \\ n^{<t>} \end{bmatrix} \begin{matrix} \downarrow 100 \\ \uparrow 10/100 \\ \downarrow 10000 \\ \uparrow 10/100 \end{matrix}$$

$$\therefore [w_{aa}: w_{an}] \begin{bmatrix} a^{<t-1>} \\ a^{<t>} \end{bmatrix} = w_{aa} a^{<t-1>} + w_{an} n^{<t>}$$

$$\text{also, } \hat{y}^{<t>} = g(w_y a^{<t>} + b_y)$$

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

4) Backpropagation through time



loss function :

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

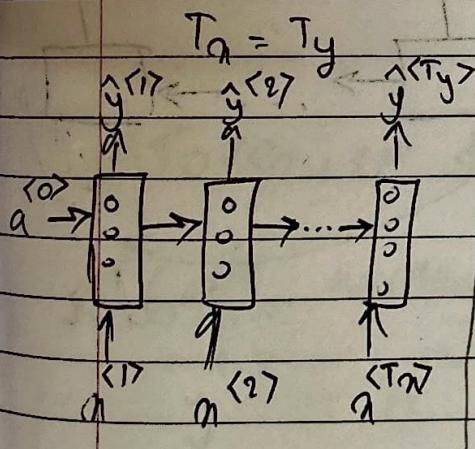
Overall loss,

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

→ sum of individual losses.

\* Backpropagation goes from right to left  
i.e. backwards in time, hence the name  
back propagation through time.

### 5) Different types of RNNs ( $T_a \neq T_y$ )

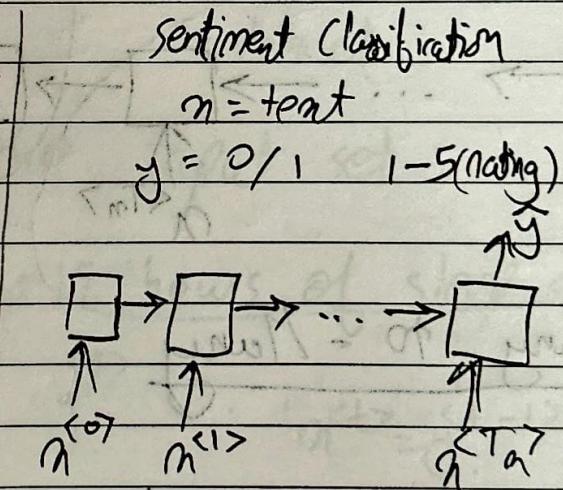


Many -to - Many

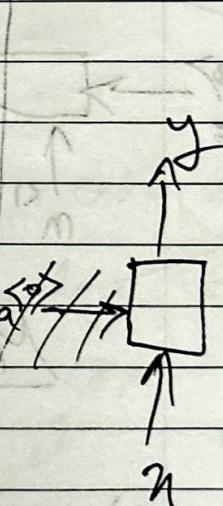
There is ... movie

Many -to - One

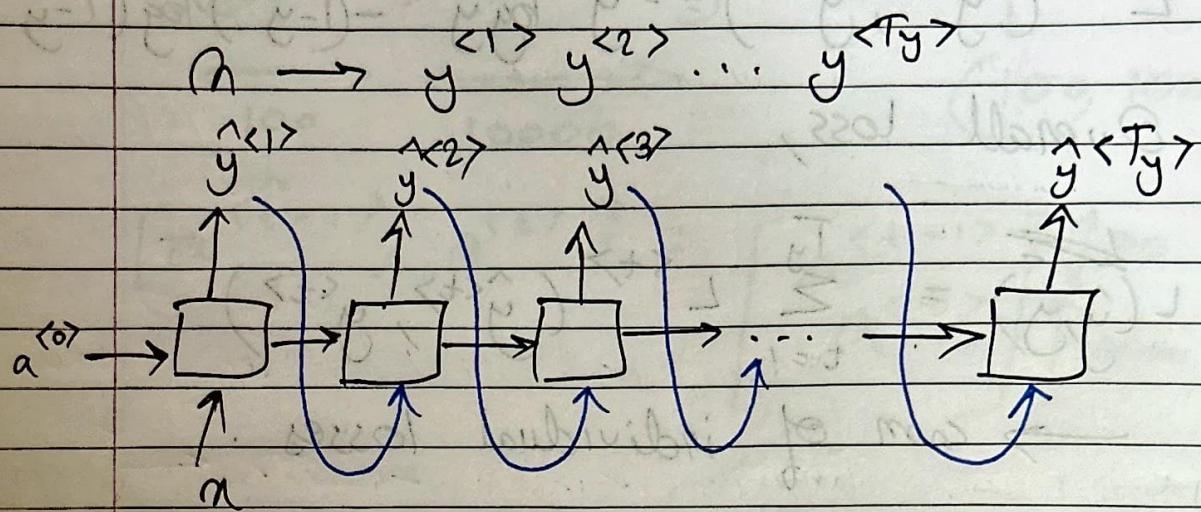
One - to - One



Many -to - One



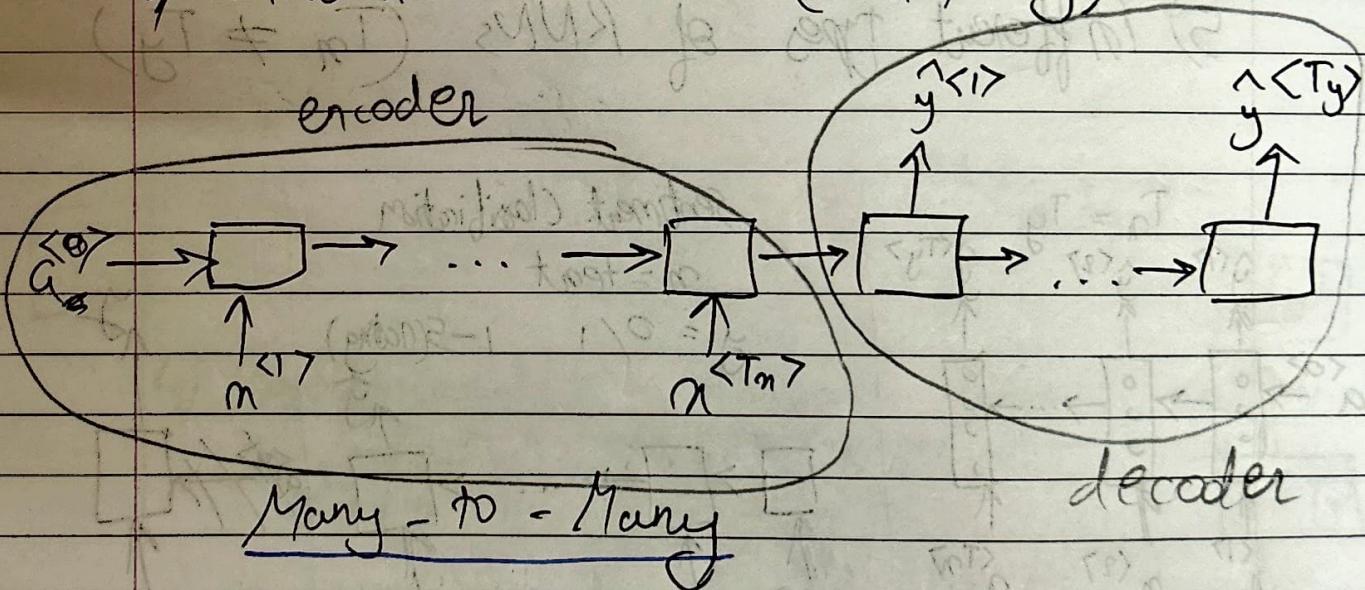
## \* Music generation



## One - to - many

Also note  $n$  can be null set here i.e.  $n = \emptyset$

## \* Machine Translation ( $T_a \neq T_y$ )



- b) Language Model and sequence generation :  
Speech Recognition :

The apple and pain salad.

The apple and pear salad.

$$P(\text{sentence}) = ? \quad P(y^{(1)}, y^{(2)}, \dots, y^{(T_y)}) = ?$$

$$P(\text{The apple and pain salad}) = 3.2 \times 10^{-13}$$

$$P(\text{The apple and pear salad}) = 5.7 \times 10^{-10}$$

This one has more probability hence it will be produced as output.

## Language modelling with an RNN

Training set : large copy of english set.

↑  
dataset

Tokenize the input set

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

$y^{(1)}$   $y^{(2)}$   $y^{(3)}$  ...  $y^{(t)}$   $= y^{(t-1)}$   $y^{(8)}$   $y^{(9)}$

In this case, we ignore the punctuation.

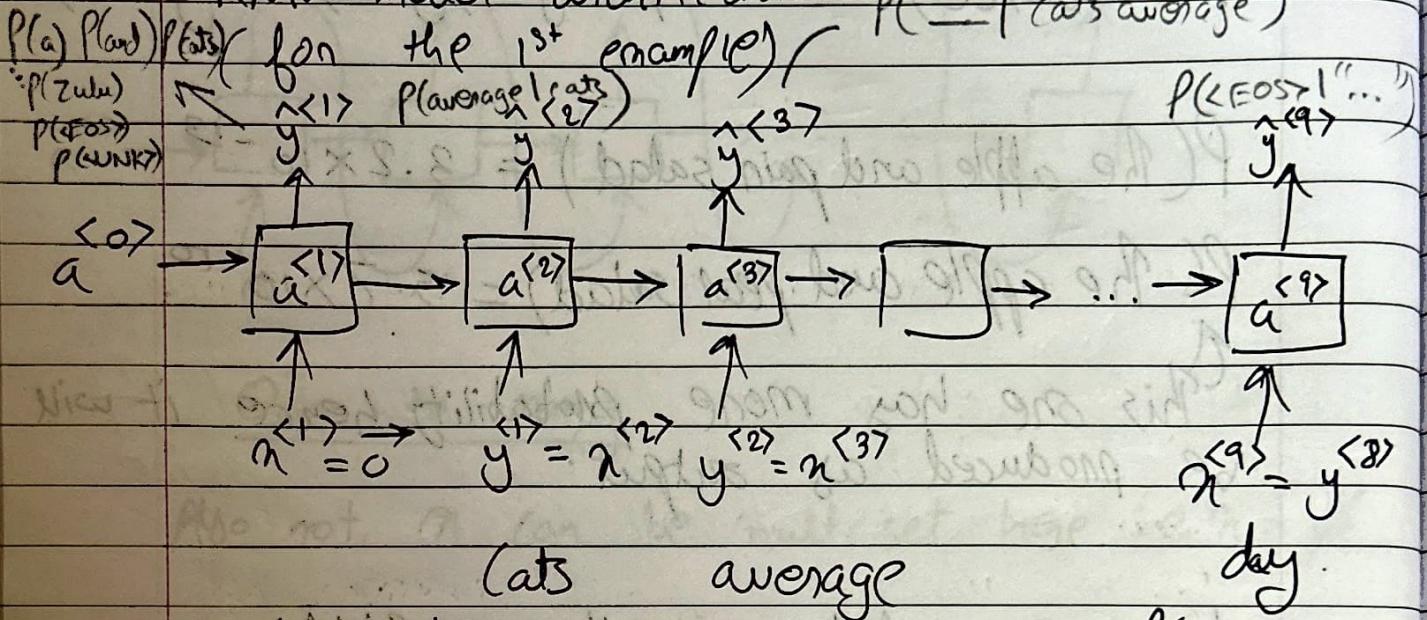
The Egyptian ~~Man~~ is a bread of cat. <EOS>

<UNK>

Unknown word since "Man" won't be in our dictionary of most common 10,000 words

Tokenization basically means taking the input sentence and map to the individual tokens / individual words / individual vocab.

RNN model architecture :  $p(\text{---} | \text{"cats average"})$



i.e. given the first few words what is the  $p(\text{next word})$ ?

Cost function :

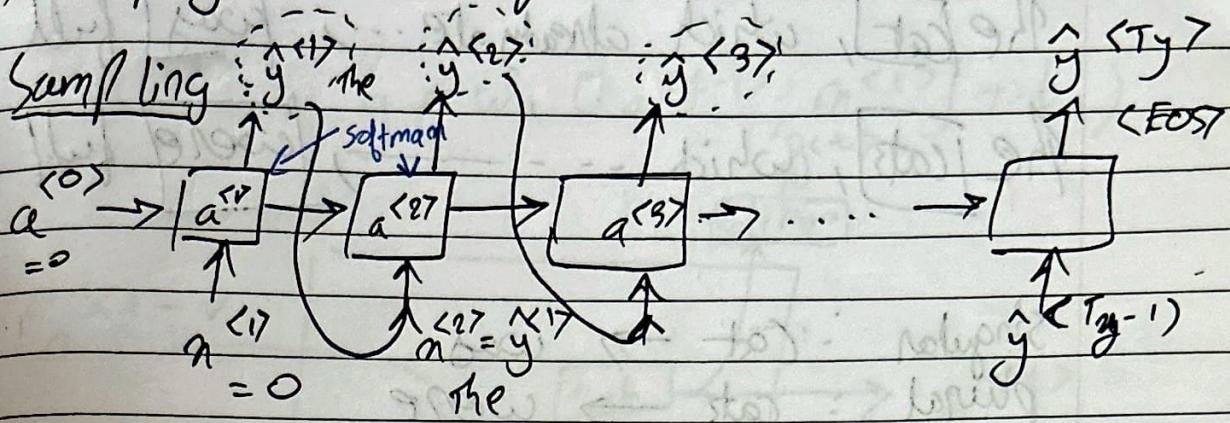
$$L(\hat{y}^{(t)}, y^{(t)}) = - \sum_i y_i^{(t)} \log \hat{y}_i^{(t)}$$

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

Note : Probability of three words in sentence -

$$P(y^{(1)}, y^{(2)}, y^{(3)}) = P(y^{(1)}) * P(y^{(2)} | y^{(1)}) * P(y^{(3)} | y^{(1)}, y^{(2)})$$

## 7) Sampling Novel Sequences



$y^{(1)} \rightarrow p(a), p(aa\text{non}) \dots p(zulu) \ p(<\text{UNK}>)$

$$y^{<2>} \rightarrow p(- | \text{Re})$$

Sampling is generating words based on the probabilities and can be implemented as:

→ np. random. choice

# Character-level language model

Vocabulary =  $[a, b, c, \dots, z, \underline{A}, \dots, \underline{z}]$

eg: (at average)

$y^1, y^2, y^3, y^4, y^5, \dots$

character level produces one character at a time ∴ less effective.

**Pros:** No need to assign `<UNK>` to any word as all the words can be built using vocabulary.

Cons : Computationally expensive to train and also

## 8) Vanishing gradients with RNNs

The [cat], which already ate ..., was full  
 The [cats], which ..., were full

Singular : Cat  $\rightarrow$  was

plural : cats  $\rightarrow$  were

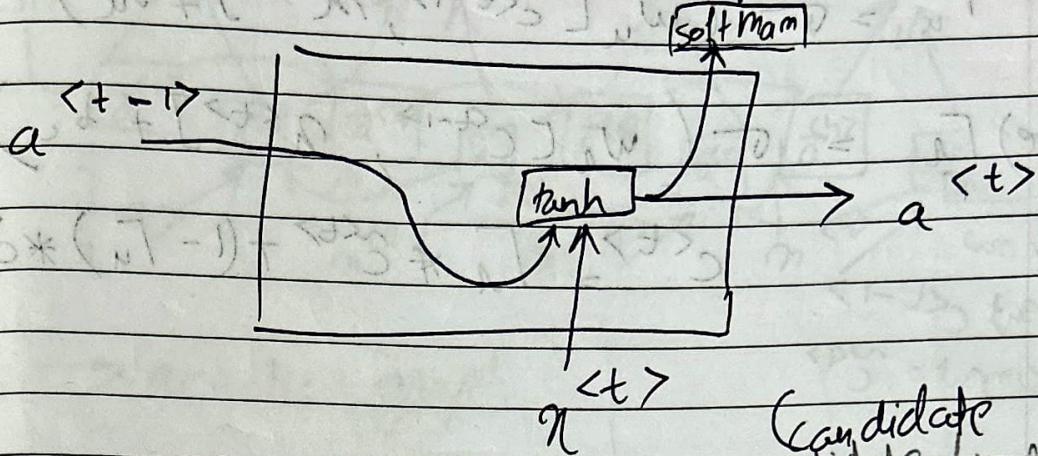
The RNN needs to remember whether the subject was singular or plural but it turns out that even in forward propagation the  $P(\text{was} | \dots)$  or  $P(\text{were} | \dots)$  will depend on some of the previous words more as compared to all the way beginning of the sentence words. Similarly while calculating derivatives in back propagation, the derivatives of the left-most nodes /parameters will less likely depend on those of the right-most ones as it slowly vanishes and becomes less important to compute.

We faced a problem "Exploding gradients" in Deep NNs which can be solved using gradient clipping. NaNs (Not a number) is resulted in case of Exploding gradients

\* Also, it is easier to spot exploding gradients in Deep NNs than vanishing gradients as exploding gradients give us very high numerical values which affect the parameters by a high factor so easier to spot.

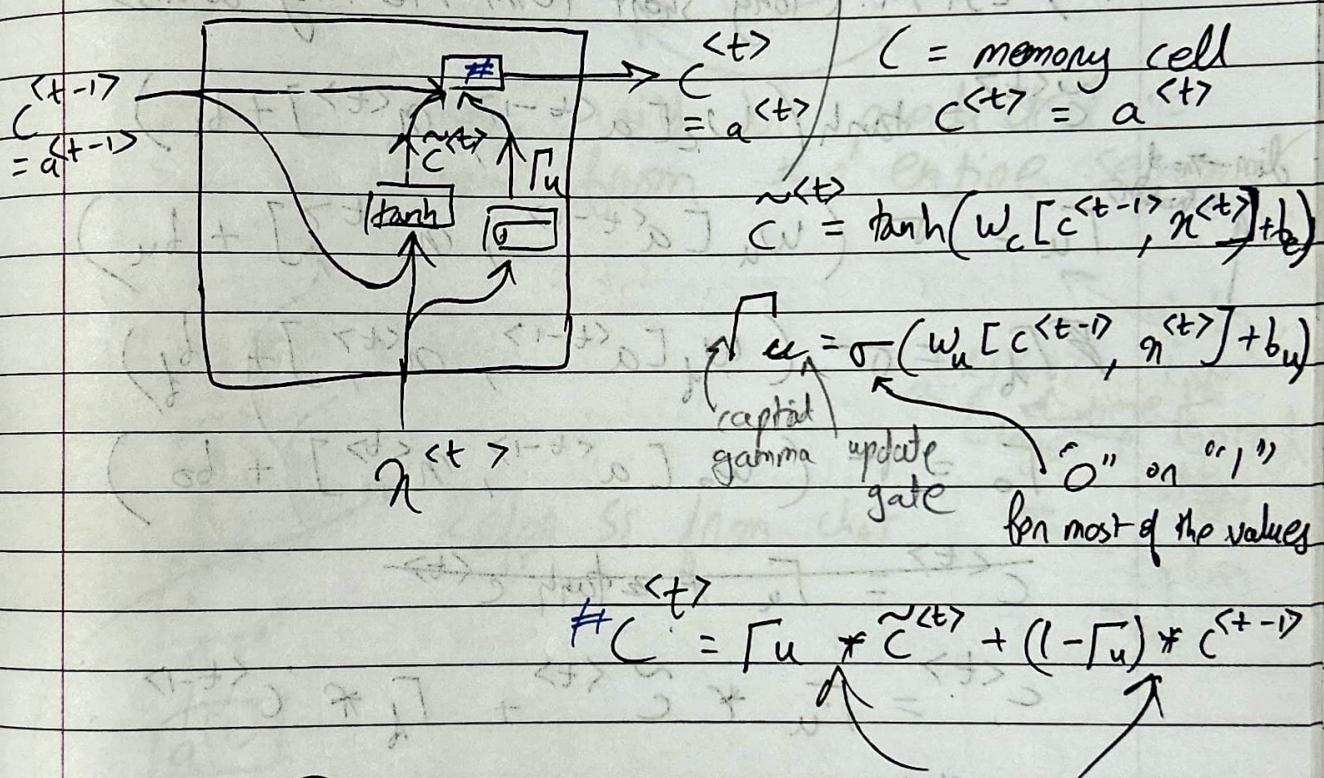
## 9) GRU : Gated Recurrent Unit

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



inRU (Simplified)

(candidate value)  
Candidate for replacing  
memory cell



$\Gamma_u = 1$        $\Gamma_u = 0$        $\Gamma_u = 0$        $\Gamma_u = 0 \dots$  ... element-wise multiplication

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

LSTM  $\rightarrow$  GRU

$c^{(t)}$   $\rightarrow$  cell state  $c$ ,  $\tilde{c}^{(t)}$   $\rightarrow$  candidate state

Page No.

Date

Value  
state

Full GRU

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

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

$$(\text{reset gate}) r_u = \sigma(w_r [c^{(t-1)}, n^{(t)}] + b_r)$$

tells how  
adequate  $c^{(t-1)}$

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

(is to compute  $\tilde{c}^{(t)}$ )

10) LSTM (long short term memory unit)

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

dim  $\rightarrow$  no. of  
activations

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

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

forget  
gate

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

output  
gate

$$\tilde{c}^{(t)} = r_u * \tanh \tilde{c}^{(t)}$$

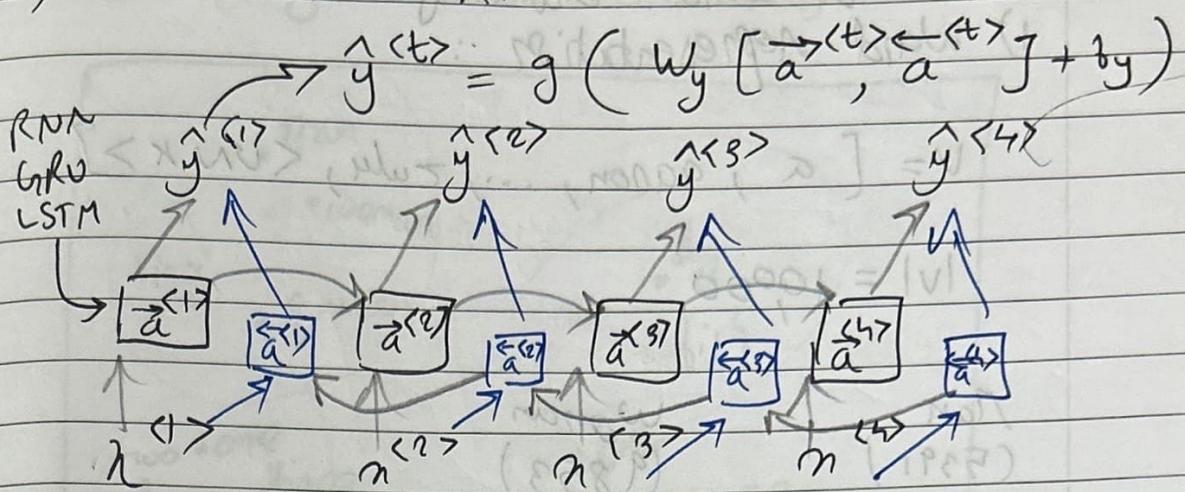
$$c^{(t)} = r_u * \tilde{c}^{(t)} + r_f * c^{(t-1)}$$

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

N/A: feedback connection:  $c^{(t-1)}$

from assignment:  $y_{\text{pred}}^{(t)} = \text{softmax}(w_y a^{(t)} + b_y)$

## 11) Bidirectional RNN

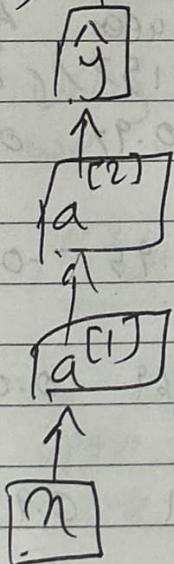


→ acyclic graph

He said "Teddy Resolut :-"

Hence, teddy will be predicted to be a person from the entire sequence.

## 12) Deep RNNs



$a^{[l]} \leftarrow$   
↑ time +  
layer l

Refer ss from chat