

COURSE 5: SEQUENCE MODELS

When the points in the dataset are dependent on the other points in the dataset, the data is called sequential.

Week 1:

Video 1: Why sequence models?

Applications :

(i) speech recognition (sequence to sequence)

X: input eg. audio (wave sequence)

Y: output eg. text (text sequence)

(ii) music generation

X: ♪

Y: wave sequence \rightarrow music

(iii) sentiment classification

X: text

Y: integer rating (stars) (1 to 5)

(iv) DNA sequence analysis

X: DNA sequence

Y: DNA sequence identification

etc.

Video 2: Notation

Named Entity Recognition Example

X: "Harry Potter and Hermione Granger invented a new spell."
Y: H I O I I O O O O O

1 - name same size of i/p & o/p sequence
0 - not name

* n^{th} element of ^{input} sequence is represented as $x^{<n>}$.
example: 1st by $x^{<1>} = \text{Harry}$.
 2nd by $x^{<2>} = \text{Potter}$.

* n^{th} element of output sequence (y) is represented as $y^{<n>}$.
example: $y^{<1>} = 1$
 $y^{<2>} = 1$

* $T_x = \text{size of input sequence}$
 $T_y = \text{size of output sequence}$

* $i \rightarrow$ training example

$x^{(i)<t>}$: element t of i^{th} training example in input
 $y^{(i)<t>}$: element t of i^{th} training example in output

$T_x^{(i)}$: input sequence length of i^{th} training example
 $T_y^{(i)}$: output sequence length of i^{th} training example

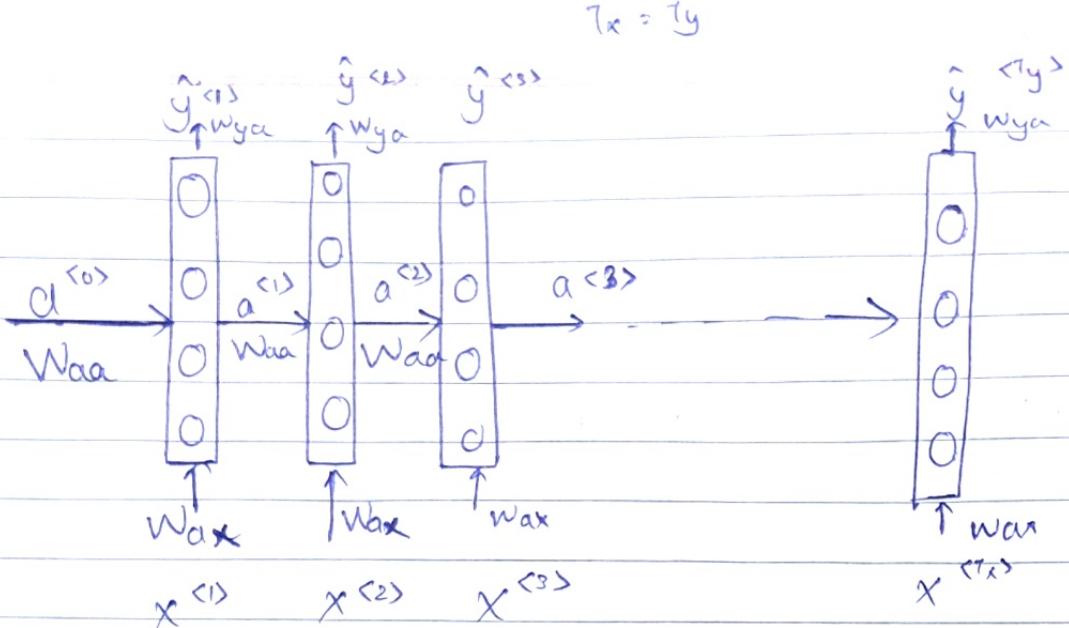
Representing words:

- (i) Build a dictionary. It can have 10,000 / 50,000 / 1,00,000 words.
- (ii) dictionary = vocabulary list
 - each word has unique index representation
 - Sorting is in alphabetical order.
- (iii) One hot encoding system - digit '1' is used to represent a word in the dictionary at its particular index. Others are 0.
- (iv) If there is an unknown word we can add it to our dictionary.

Video 3: Recurrent Neural Networks

Problems faced by standard neural networks in the above problem:

- inputs & outputs are of different length
- doesn't share features ~~across~~ learned across different positions of text.



$a^{<0>}$ - it is initialised with zeros

W_{ax} : (No. of hidden neurons, N_x)

W_{aa} : (No. of hidden neurons, No. of hidden neurons)

W_{ya} : (No. of hidden neurons)

$L \rightarrow R$

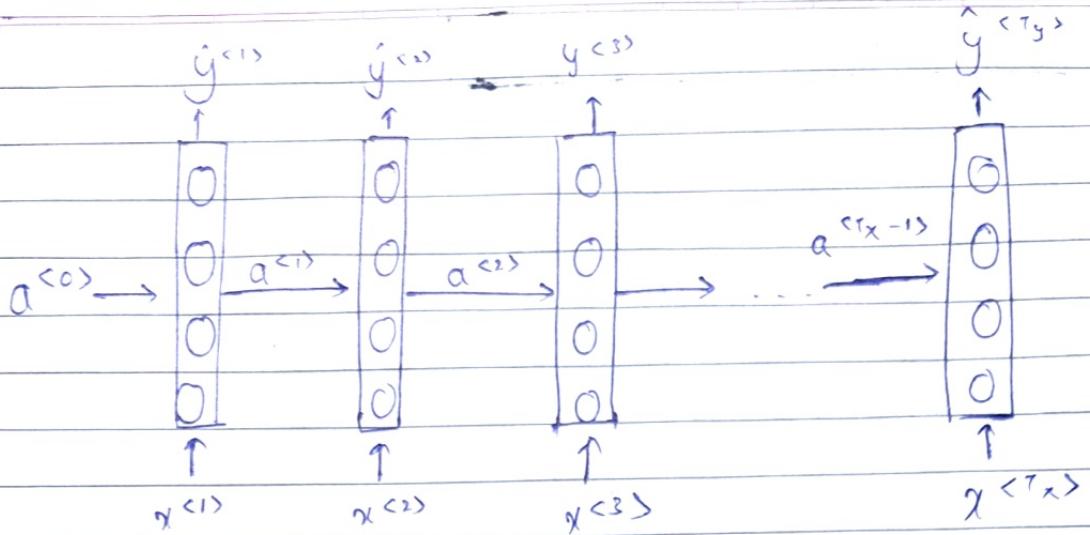
W_{ax} - parameters used in each step

W_{aa} - activation parameters

W_{ya} - governs output

$\hat{y}^{<t>}$ depends on previous input and activations

RNN's use information that is earlier in the sequence to predict. They cannot access future after the word in the sequence.



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

$$a^{<1>} = g_1(W_{aa} a^{<0>} + W_{ax} x^{<1>} + b_a) \leftarrow \tanh / \text{ReLU}$$

$$\hat{y}^{<1>} = g_2(W_{ya} a^{<1>} + b_y) \leftarrow \text{sigmoid}$$

$$\hat{y}^{<2>} = g(W_{aa} a^{<1>} + W_{ax} x^{<2>} + b_a)$$

$$\hat{y}^{<3>} = g(W_{ya} a^{<2>} + b_y)$$

$$\cancel{a^{<t>} = g(W_{aa}[a^{<t-1>}] \times [x^{<t>}] + b_a)}$$

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

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

$$W_a = \begin{pmatrix} W_{aa} & W_{ax} \end{pmatrix} \\ \longleftrightarrow (100 \times 10000) \\ (100, 10100)$$

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

Video 4: Backpropagation Through Time

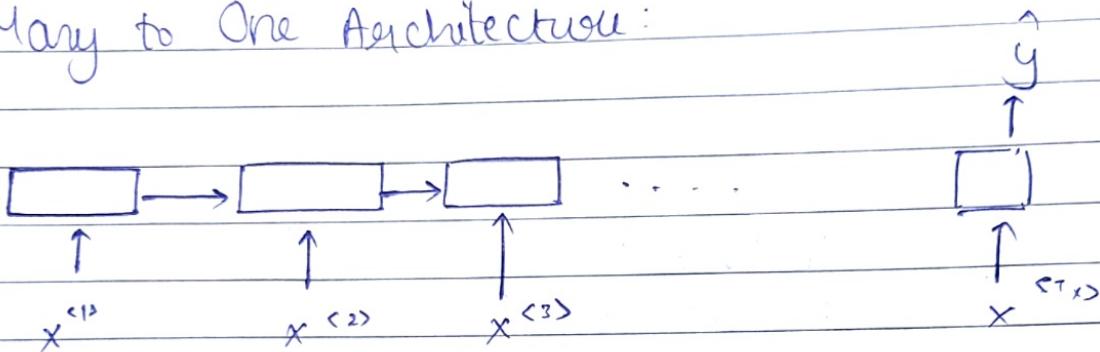
$$L^{<t>}(\hat{y}^{<t>}, y^{<t>}) = -y^{<t>} \log \hat{y}^{<t>} - (1-y^{<t>}) \log (1-\hat{y}^{<t>})$$

$$L(\hat{y}, y) = \sum L^{<t>}(\hat{y}^{<t>}, y^{<t>})$$

Video 5: Different types of RNN

(i) Many to Many Architecture : Input sequence has many inputs as sequence and the outputs sequence also has many outputs.
eg. language translation.

(ii) Many to One Architecture :

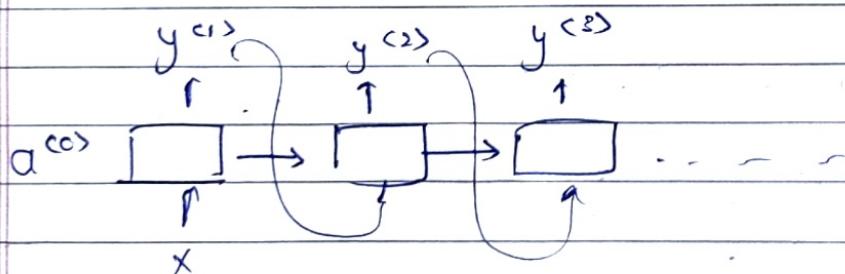


eg. sentiment analysis in movie rating .

(iii) One to many Architecture

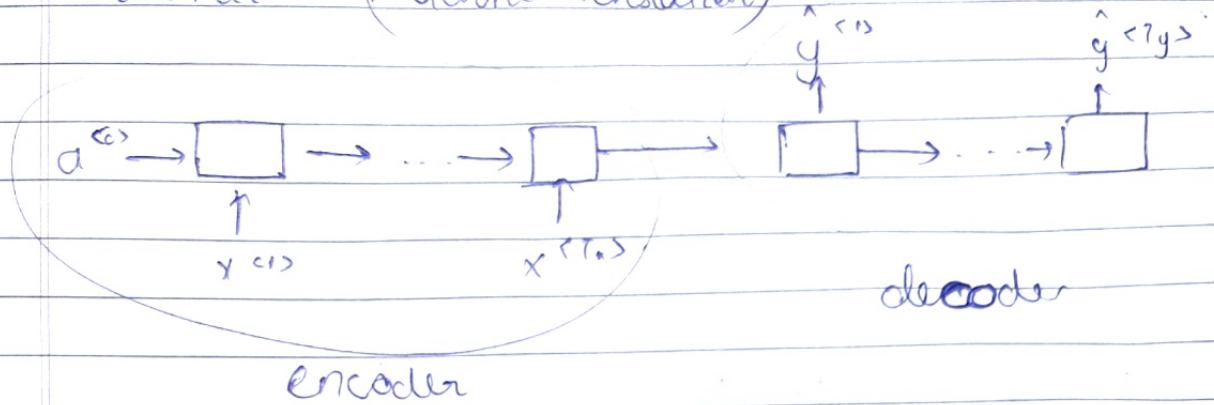
x : input

$x = \phi / \text{integer like "genre"}$



eg. music generation

Alternative (Machine Translation)



Video 6: Language Model & Sequence Generation.

- * Language model gives the probability of a given sequence of words.

Example:

The apple & pear salad
The apple & pears salad

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

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

- * To build Language Model:

(i) Training set: large corpus of english text
 ↑
 huge body

(ii) Tokenise

(iii) One hot vectors

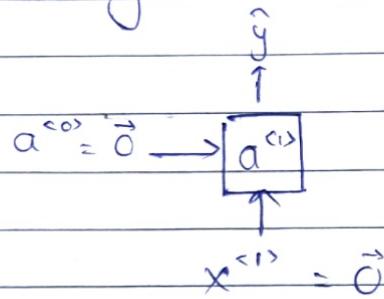
(iv) Extra token <EOS> for end of sentence

* you can tokenise punctuation too by adding to vocabulary.

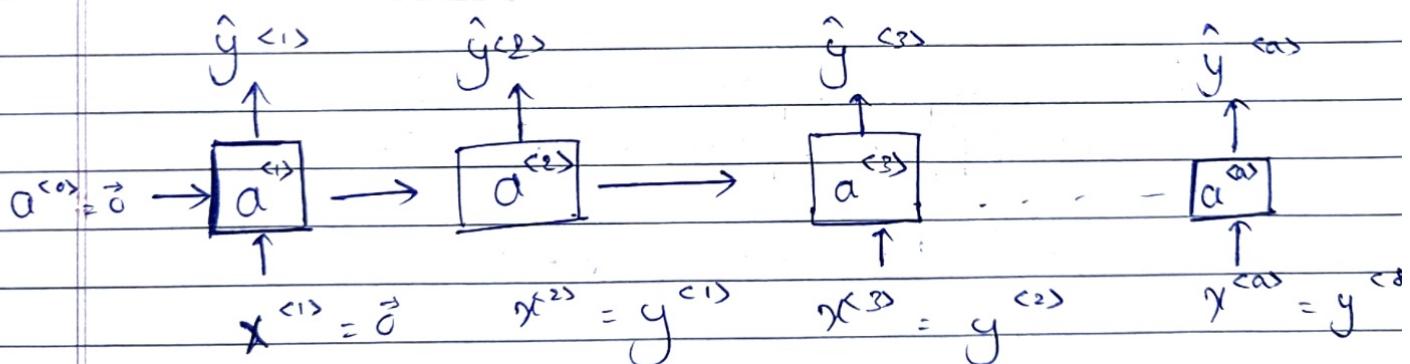
* The Egyptian Mau is a breed of cat. <EOS>

↑
<UNK>
↑
unknown word

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



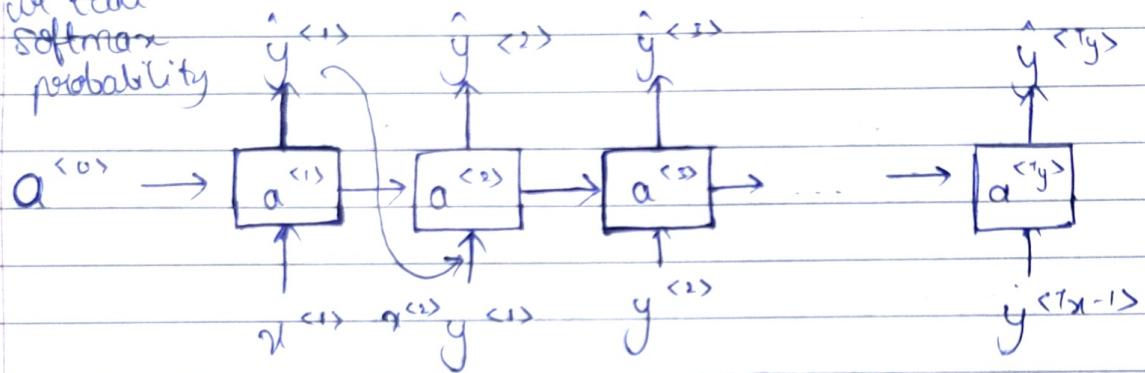
$a^{<1>}$ - makes softmax prediction about probability of first word y



$$\frac{P(y^{<1>} | y^{<0>})}{P(y^{<0>})} \frac{P(y^{<2>} | y^{<1>}, y^{<0>})}{P(y^{<1>})} \frac{P(y^{<3>} | y^{<2>}, y^{<1>})}{P(y^{<2>})} \leftarrow$$

Video 2: Sampling Novel Sequences

we take softmax probability



$a^{<0>} - \text{zeros vector}$

$x^{<t>} - \text{zeros vector}$

Prediction y^t is chosen from distribution obtained by $\hat{y}^{<t>}$. eg. The
In numpy - np.random.choice

Continue till you get <EOS> token
you can ignore <UNK> token

- * we have to build a word level language model.
- * for character level language model -
vocabulary : (a-z A-Z 0-9), punctuation, special characters / token.

In character level model there will be no tokens taken.

However you end up with longer sequences. Unlike word level language models they cannot capture long range dependencies between two words in earlier parts affect words in later part. They are harder to train.

Video: Vanishing Gradients With RNN's

Example:

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

The cats, which , were full.

Cat - was
cats - were } RNN can't

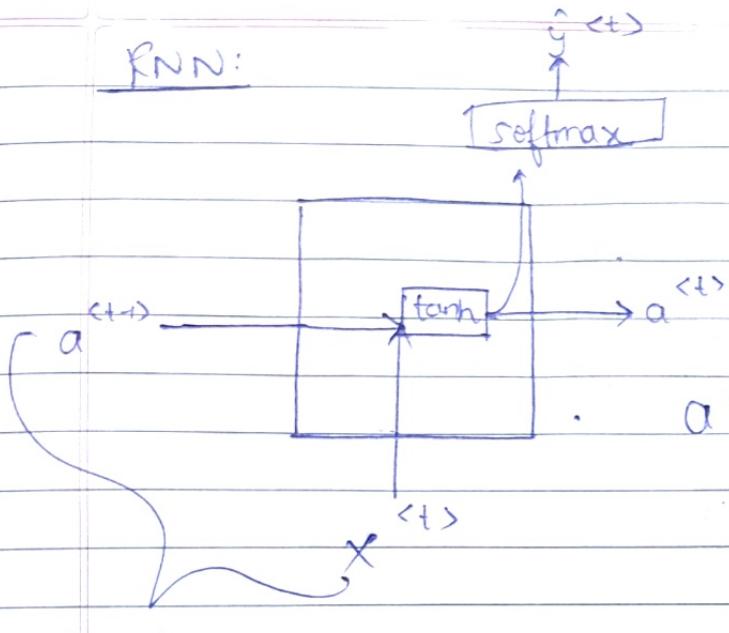
} i.e vanishing gradient

For computing the word "was" we need to compute the gradient for everything behind. Multiplying fractions tends to vanish the gradient while multiplication of large numbers cause it to explode.

Exploding gradients can be seen when weight values become NaN

Solution - gradient clipping

rescale gradient so that it isn't too big

RNN:

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

input

$$C = \text{memory cell}$$

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

↑
at time t

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

 F_u = value between 0 & 1

$$\uparrow = \alpha (W_u[C^{<t-1>}, x^{<t>}] + b_u)$$

update
gate C is updated using \tilde{C}

gate decides whether we update or not

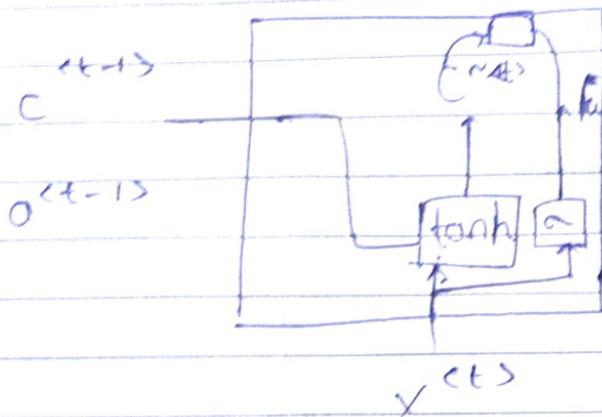
 α is set to 0/1

$$C^{<t>} = F_u * \tilde{C}^{<t>} + (1 - F_u) * C^{<t-1>}$$

Example:

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

GRU:



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

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

~~gates~~
gates

$$f_g = \alpha(W_g[c^{t-1}, x^{t-1}] + b_g)$$

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

Shape of $o^{t-1} = c^{t-1} = \tilde{c}^{t-1} = u^{t-1}$

Video 9: Long Short Term Memory (LSTM)

(GRU):

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

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

$$f_g = \alpha(W_g[c^{t-1}, x^{t-1}] + b_g)$$

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

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

update
forget
about

$$\begin{aligned}\tilde{c}^{<t>} &= \tanh(W_c[a^{<t-1>}], x^{<t>}]) + b_c \\ f_u &= \sigma(W_u[a^{<t-1>}], x^{<t>}]) + b_u \\ f_f &= \sigma(W_f[a^{<t-1>}], x^{<t>}]) + b_f \\ f_o &= \sigma(W_o[a^{<t-1>}], x^{<t>}]) + b_o \\ c^{<t>} &= f_u * \tilde{c}^{<t>} + f_f * c^{<t-1>} \\ a^{<t>} &= f_o * \tanh(c^{<t>})\end{aligned}$$

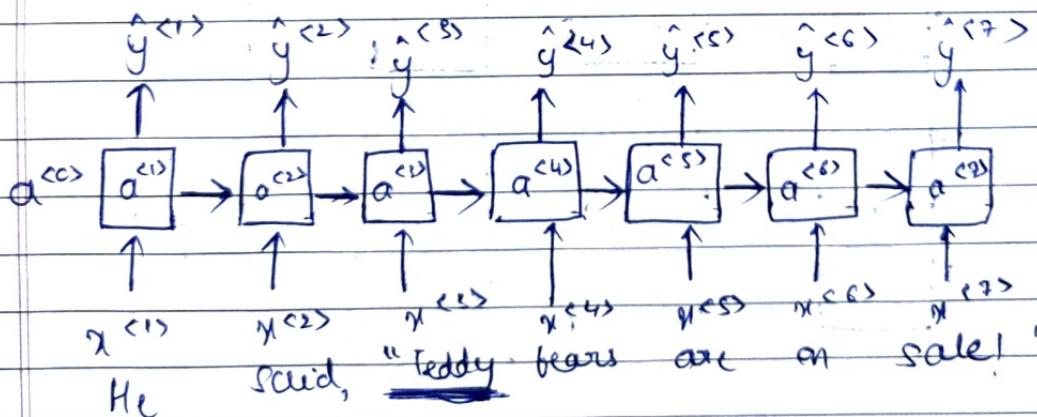
$c^{<t-1>}$ forgetable connection

LSTM >>> GRU
 $c^{<t>} \neq a^{<t>}$

LSTM & GRU's are good at memorising values. Some equal values stored in the memory cell can be used for many many time steps.

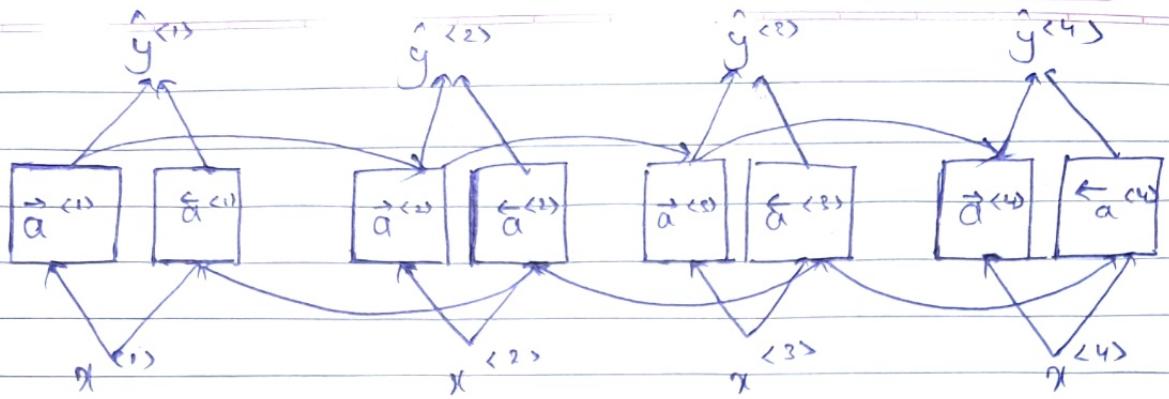
Video 10: Bidirectional RNN (BRNN)

- * Bidirectional RNN helps you take information from both earlier & later in sequence



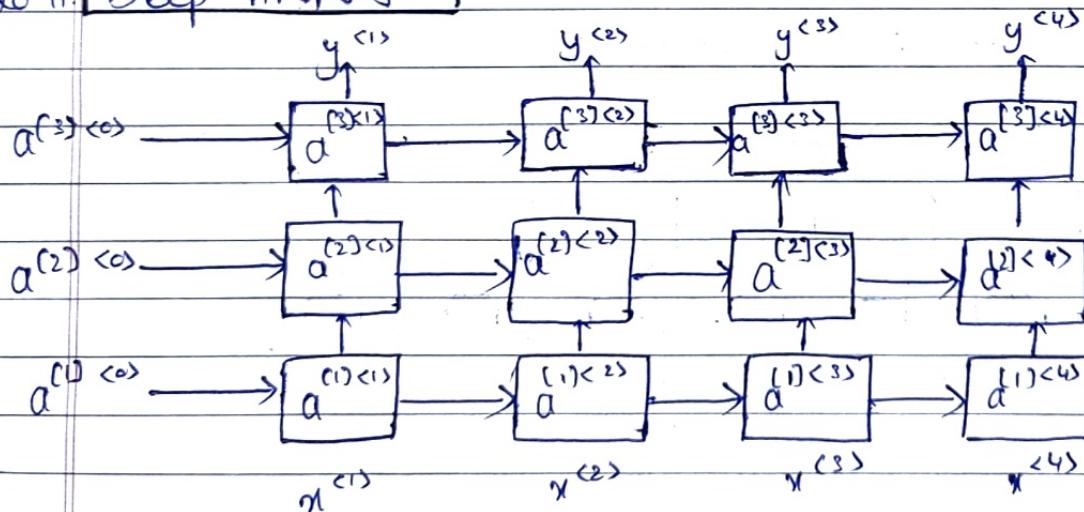
We know what
"teddy bears" and not "said"

BRNN fixes this issue.



- * BiRNN is an acyclic graph
- * Part of the forward propagation goes from left to right & part from right to ~~left~~ left.
- * Blocks can be RNN/LSTM/GRU
- * Disadvantage: whole sequence is required before processing it.

Video III: Deep RNN's



In RNN's layers is considered deep.