

Backpropagation in RNN.

→ It is called as backpropagation through time

Eg: Take Many to one RNN for sentiment analysis.

Eg:  $\begin{cases} R_1 = \text{cat mat rat} \\ R_2 = \text{Rat rat mat} \\ R_3 = \text{mat mat cat} \end{cases}$

Step 1: Encoding

corpus = [cat mat rat]

cat  $\Rightarrow$  [100] [010] [001]

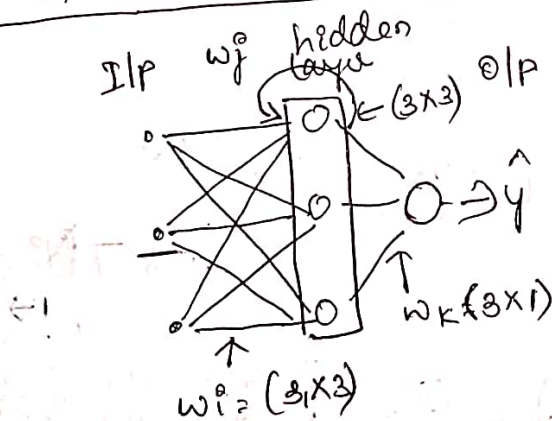
Eg: in vector representation

$R_1 = \begin{bmatrix} 1 & 0 & 0 \end{bmatrix} \begin{bmatrix} 0 & 1 & 0 \end{bmatrix} \begin{bmatrix} 0 & 0 & 1 \end{bmatrix}$

$R_2 = \begin{bmatrix} 0 & 0 & 1 \end{bmatrix} \begin{bmatrix} 0 & 0 & 1 \end{bmatrix} \begin{bmatrix} 0 & 1 & 0 \end{bmatrix}$

$R_3 = \begin{bmatrix} 0 & 1 & 0 \end{bmatrix} \begin{bmatrix} 0 & 1 & 0 \end{bmatrix} \begin{bmatrix} 1 & 0 & 0 \end{bmatrix}$

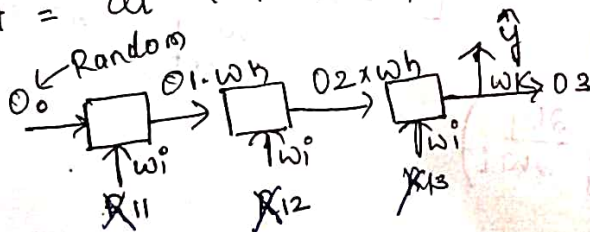
RNN in forward propagation.



bias = 4

we need to send words to RNN in each review.

$R_1 = \text{at } t=1 \text{ send } [100] \Rightarrow \text{cat}$



## Equations for forward propagation

$$O_1 = F(x_{11} \cdot w_i^0 + O_0 w_h) + b_1$$

$$O_2 = F(x_{12} \cdot w_i^1 + O_1 w_h) + b_2$$

$$O_3 = F(x_{13} \cdot w_i^2 + O_2 w_h) + b_3$$

$$\hat{y} = F(O_3 \cdot w_k) + b_4$$

$$\text{loss} = (y - \hat{y})$$

\* we need to reduce loss function using gradient descent.

\* we need to update  $w_i^0, w_h, w_k$

## Weight updation formula

$$w_i^0 = w_i^0 - \eta \frac{\partial L}{\partial w_i^0} \quad \left\{ \begin{array}{l} \text{calculate this} \\ \text{using chain rule.} \end{array} \right.$$

$$w_h = w_h - \eta \frac{\partial L}{\partial w_h}$$

$$w_k = w_k - \eta \frac{\partial L}{\partial w_k}$$

[loss depends on  $w_k$  and  $O_3$ ]

$$\frac{\partial L}{\partial w_k} = \frac{\partial L}{\partial \hat{y}} \cdot \frac{\partial \hat{y}}{\partial w_k}$$

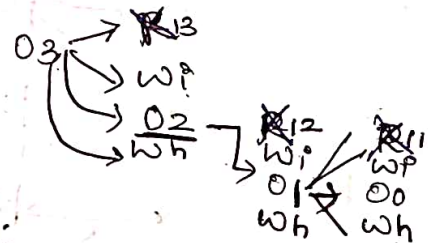
$$\boxed{\frac{\partial L}{\partial w_k} = \frac{\partial L}{\partial \hat{y}} \cdot \frac{\partial \hat{y}}{\partial O_3} \cdot \frac{\partial O_3}{\partial w_k}}$$

To update  $w_i^0$

$L \rightarrow O_3 \rightarrow w_0$

$$\frac{\partial L}{\partial w_i^0} = \left[ \left( \frac{\partial L}{\partial \hat{y}} \cdot \frac{\partial \hat{y}}{\partial O_3} \cdot \frac{\partial O_3}{\partial w_i^0} \right) + \right.$$

$$\left. \left( \frac{\partial L}{\partial \hat{y}} \cdot \frac{\partial \hat{y}}{\partial O_3} \cdot \frac{\partial O_3}{\partial O_2} \cdot \frac{\partial O_2}{\partial w_i^0} \right) + \left( \frac{\partial L}{\partial \hat{y}} \cdot \frac{\partial \hat{y}}{\partial O_3} \cdot \frac{\partial O_3}{\partial O_2} \cdot \frac{\partial O_2}{\partial O_1} \cdot \frac{\partial O_1}{\partial w_i^0} \right) \right]$$



[↑ in Review, there are 3 words hence 3 set]

In simple,

$$\frac{\partial L}{\partial w_i^0} = \sum_{j=1}^3 \left( \frac{\partial L}{\partial \hat{y}} \cdot \frac{\partial \hat{y}}{\partial O_j} \cdot \frac{\partial O_j}{\partial w_i^0} \right)$$

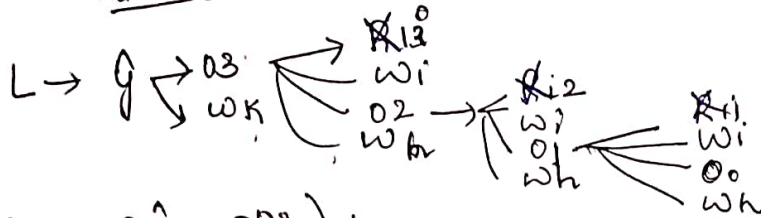
with expansion for  $j=3$

$$= \frac{\partial L}{\partial \hat{y}} \left( \frac{\partial \hat{y}}{\partial o_1} \right) \frac{\partial o_1}{\partial w_i}$$

$$\hat{y} \rightarrow o_3 \rightarrow o_2 \rightarrow o_1$$

$$= \frac{\partial L}{\partial \hat{y}} \cdot \frac{\partial \hat{y}}{\partial o_3} \cdot \frac{\partial o_3}{\partial o_2} \cdot \frac{\partial o_2}{\partial o_1} \cdot \frac{\partial o_1}{\partial w_i}$$

To update  $\left[ \frac{\partial L}{\partial w_h} \right]$



$$\frac{\partial L}{\partial w_h} = \left[ \left( \frac{\partial L}{\partial \hat{y}} \cdot \frac{\partial \hat{y}}{\partial o_3} \cdot \frac{\partial o_3}{\partial w_h} \right) + \left( \frac{\partial L}{\partial \hat{y}} \cdot \frac{\partial \hat{y}}{\partial o_3} \cdot \frac{\partial o_3}{\partial o_2} \cdot \frac{\partial o_2}{\partial w_h} \right) + \left( \frac{\partial L}{\partial \hat{y}} \cdot \frac{\partial \hat{y}}{\partial o_3} \cdot \frac{\partial o_3}{\partial o_2} \cdot \frac{\partial o_2}{\partial o_1} \cdot \frac{\partial o_1}{\partial w_h} \right) \right]$$

$$\frac{\partial L}{\partial w_h} = \sum_{j=1}^n \left( \frac{\partial L}{\partial \hat{y}} \cdot \frac{\partial \hat{y}}{\partial o_j} \cdot \frac{\partial o_j}{\partial w_h} \right)$$

Then we try to reach global minima

### Problems with RNN.

- ① Problem of long term dependency
- ② Problem of ~~unstable~~ stagnated training. } Because of unstable gradients.

#### ① Problem of long term dependency

→ If the sequence (number of timesteps) is too long, it will not remember beginning states.

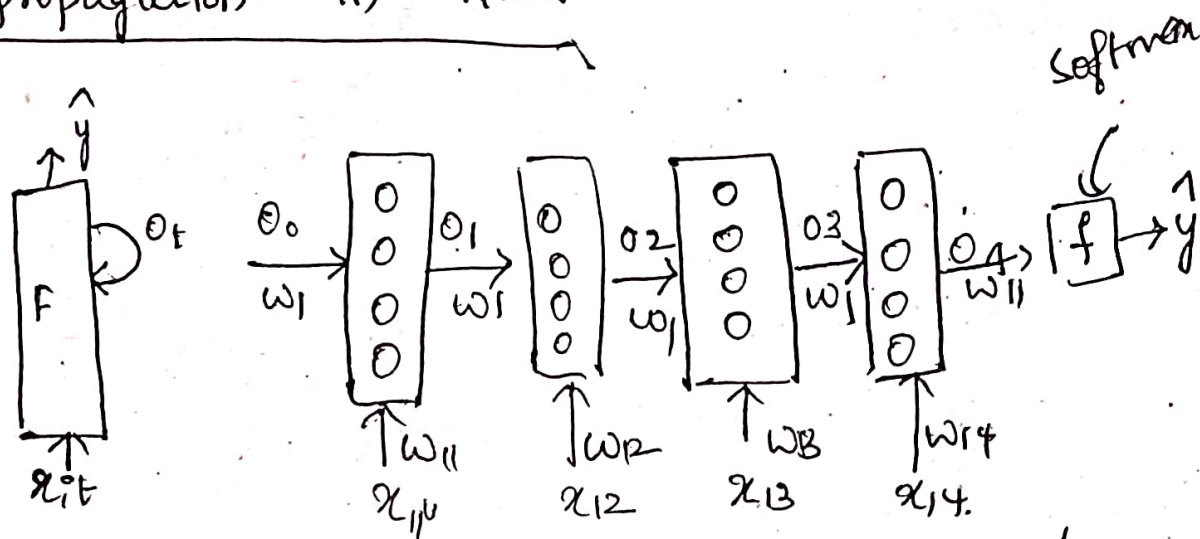
Eg. Maharashtra is a beautiful place. I went there last year But I could not enjoy properly. Because I don't understand long term malati short malathi is spoken in Maharashtra

→ RNN is not having that much memory  
→ It results in vanishing gradient problem.

② \* Sometimes we cannot train properly.  
\* Exploding gradient problem



# Back propagation in RNN



$$O_1 = f(x_{11}w + O_0w_1)$$

$$O_2 = f(x_{12}w + O_1w_1)$$

$$O_3 = f(x_{13}w + O_2w_1)$$

$$O_4 = f(x_{14}w + O_3w_1)$$

\* Reduce loss  
\*

$$\frac{\partial L}{\partial \hat{y}} = \frac{\partial L}{\partial O_4} = \frac{\partial L}{\partial \hat{y}} \cdot \frac{\partial \hat{y}}{\partial O_4}$$

$$\frac{\partial L}{\partial w_{14}} = \frac{\partial L}{\partial \hat{y}} \cdot \frac{\partial \hat{y}}{\partial O_4} \cdot \frac{\partial O_4}{\partial w_{14}}$$

→ updation of weights in each sequence.

## Problems

- ① sigmoid → vanishing gradient
- ② ReLU → if derivative is greater than 1 then ~~exp~~ exploding

It can be solved by  
LSTM.

RNN  $\rightarrow$  sequential data,

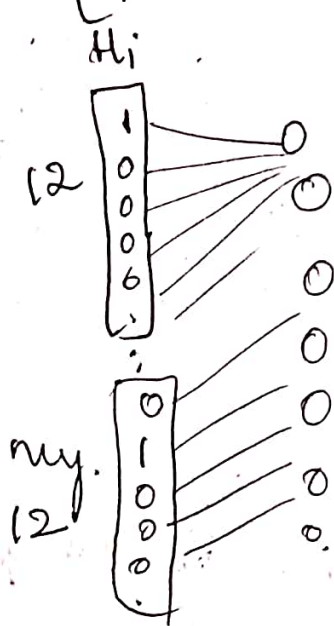
NLP  $\rightarrow$  RNN.

Text

we need to vectorize it.

Hi | my | name one hot Encoder

$[1, 0, 0, 0, 0], [0, 1, 0, \dots], [0, 0, 1, \dots]$



2nd  
36  $\subset$  Input size

3rd  
48

\* Input size we cannot be vary in ANN

$\rightarrow$  Identify which has the highest words

Then put zero padding

$\Rightarrow$  All will be zero.

- ① text i/p is varying
- ② zero padding (unnecessary computation)
- ③ Prediction will be wrong
- ④ sequence information is lost

# RNN forward propagation.

⇒ RNN accepts the data

in C timesteps, input = feeling

① movie was good

② movie was bad

③ movie was not good

Review	Sentiment
movie was good	1
movie was bad	0
movie was not good	0

① we need to convert words to vectors.

Unique words forms a vocabulary

↓ corpus

⇒ Review 1

[[10000], [01000], [00100]]

Each word will be sent to

⇒ I can represent movie with 5 numbers

[10000] [01000]  
good was

[00100] [00010]  
bad

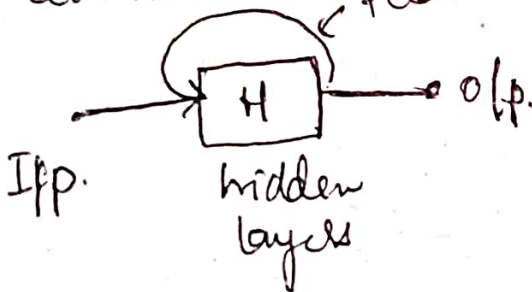
[00001]  
not

RNN separately

at time = 1 [10000]

## RNN architecture

at time = 1





Review	Sentiment
R <sub>1</sub> movie was good	1
R <sub>2</sub> movie was bad	0
R <sub>3</sub> movie was not good	0

Input for RNN will be in the form  
(~~time~~ steps, input-features)

Dictionary/  
corpus. { movie, was good, bad,  
not }

Hence

movie =	[10000]	} using one-hot encoding
was =	[01000]	
good =	[00100]	
bad =	[00010]	
not =	[00001]	

Hence

Review1 = movie was good  
 $\begin{bmatrix} [10000] & [01000] & [00100] \end{bmatrix}$

Review2 = movie was bad  
 $\begin{bmatrix} [10000] & [01000] & [00010] \end{bmatrix}$

shape (3, 5) (Number of ip features)

all 3 words

Number of timesteps

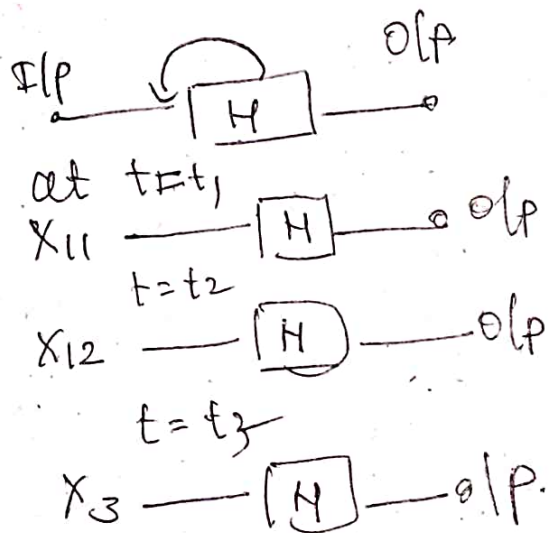


In RNN, each word by word we will send data to RNN.

In keras, (batchsize, timestep,  $q(p)$  features)

(3, 4, 5)  
 ↑                      ↑                      ↑  
 3 reviews are sent together    max length of words    corpus length.

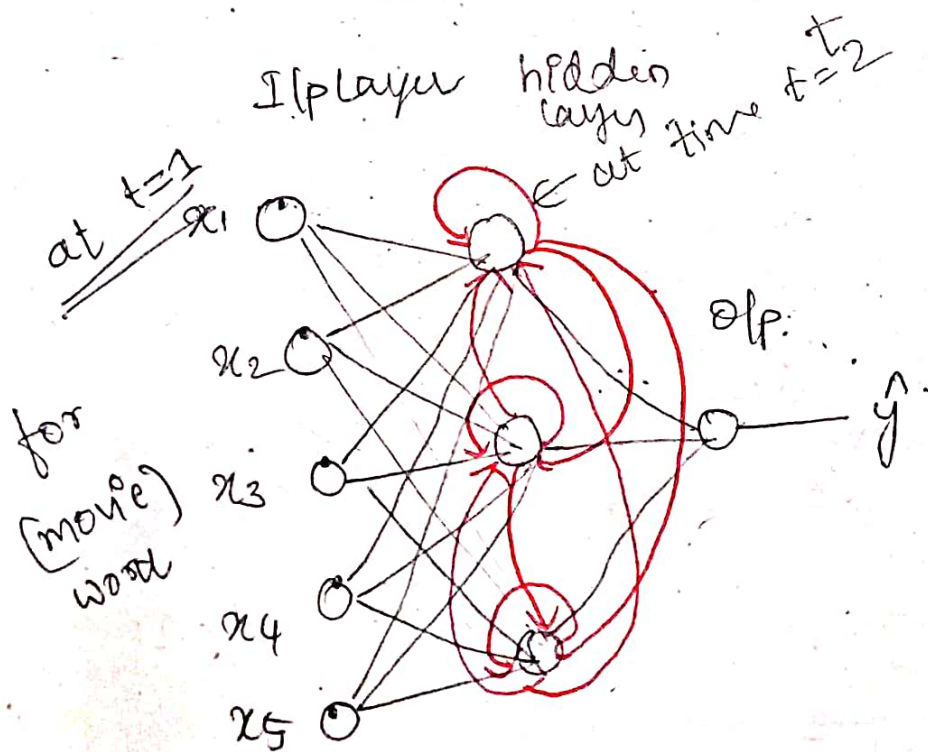
### RNN architecture



\* RNNs are not feedforward.

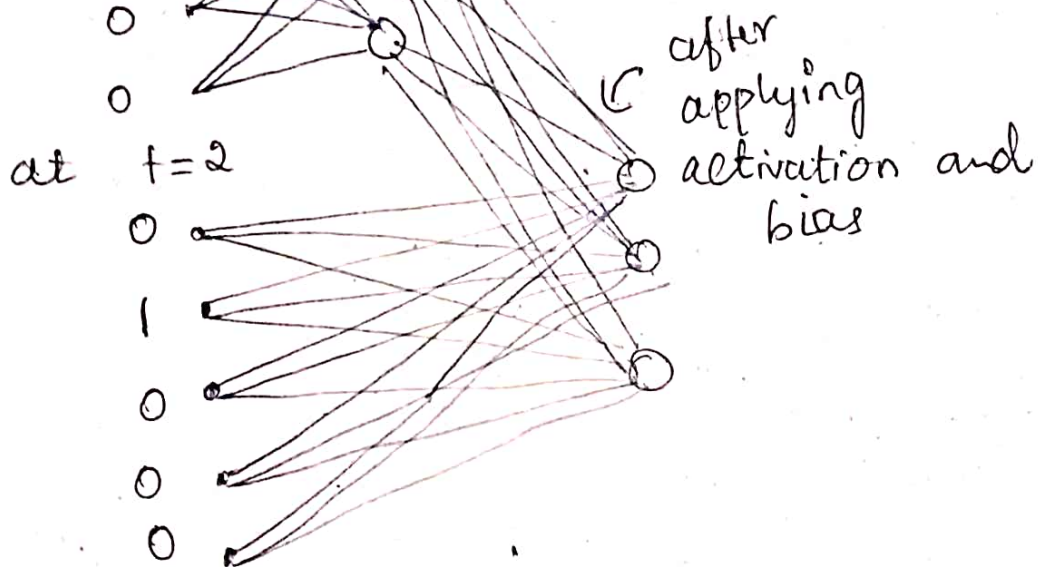
\* There is a concept of state.

\* hidden layer sends a feedback back..



at  $t = 1$

3x3 (3x0)  
or  
Random  
numbers.



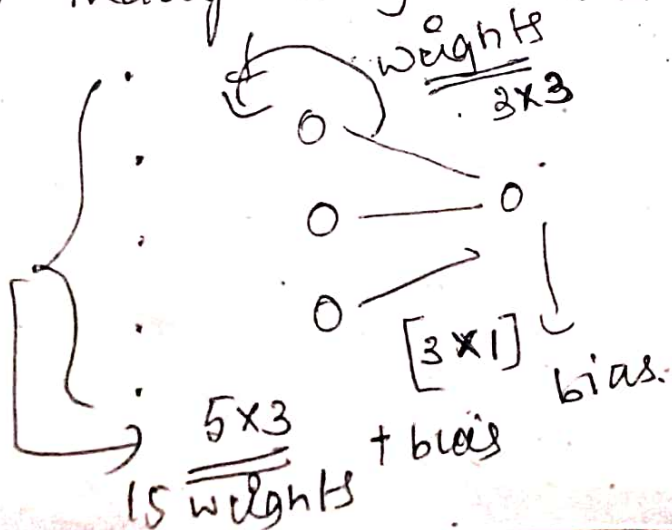
at  $t = 3$ .

0

0

0

How many weights and bias)

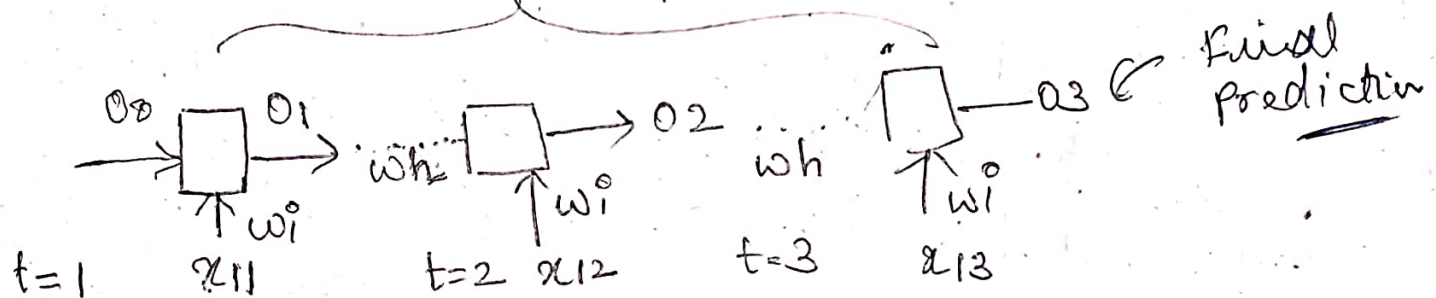


ie  $w_1 = [15 \text{ weights}]$   
 $w_2 = [3 \times 3] = 9$   
 $w_3 = 3$   
 $\text{bias} = 4$   
 $\Rightarrow 32 \text{ Trainable parameters}$

## Forward propagation

Review				Sentiment
$x_{11}$	$x_{12}$	$x_{13}$		1
$x_{21}$	$x_{22}$	$x_{23}$		0
$x_{31}$	$x_{32}$	$x_{33}$	$x_{34}$	0

completes Review 1.



## RNN Forward propagation.

Step 1:

at  $t=1$

① send 1st word  $x_{11}$

② multiply with all weights  $w_1^0$

③ Then activation function will act upon it.  
 $o_1 = f(x_{11}w_1^0) + b_1$

④ Then add bias to it

at  $t=2$

① we use the same network with same weights.



② previous  $O_1$  is given as input with weighted connection.

$$O_2 = f(x_{12}w_i^o + O_1w_h + b_2)$$

at  $t_3$

$$O_3 = f(x_{13}w_i^o + O_2w_h + b_3)$$

Note: Here weight sharing will be done.

---

1) A bidirectional RNN is a type of neural network architecture that allows for processing of i/p sequence in both forward and backward directions.

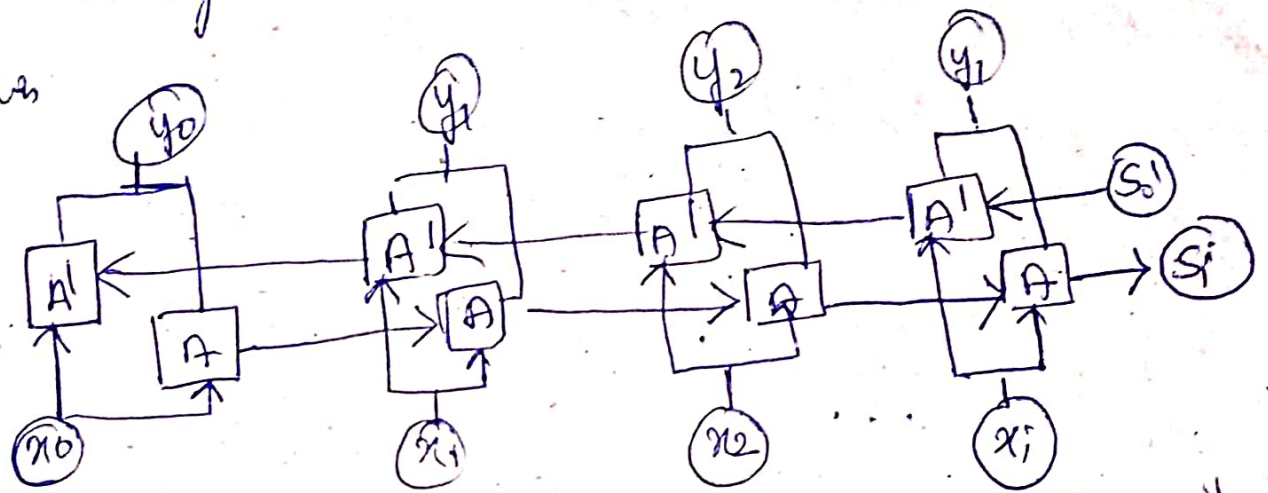
2)  $\rightarrow$  This is achieved by having two separate hidden layers for each time step: one for processing the sequence from left to right (forward) and one for processing the sequence from right to left (backward).

$\rightarrow$  The o/p of these two hidden layers are then combined to generate a final o/p which takes into account information from both directions of the i/p sequence.

This approach is particularly useful for tasks where context from both past and future time steps is important.



→ During training, the weights of both forward and backward RNNs are updated using backpropagation through time (BPTT) to minimize the loss function.



① → Bidirectional recurrent neural networks (BRNN) are really just putting two independent RNNs together.

→ The inp sequence is fed in normal time order for one network, and in reverse time order for another. The o/p of the two networks are usually concatenated at each time.

→ This structure allows the networks to have both backward and forward information about the sequence at every time step.

## Problems of RNN

① vanishing Gradient problem: RNNs

can have difficulty learning long-term dependencies in sequences due to the vanishing gradient problem.

② Exploding Gradient Problem: In addition to the vanishing gradient problem, RNN can also suffer from the opposite problem. This occurs when the gradients become too large during backpropagation.

③ Memory Limitation;

④ Computationally Expensive.

⑤ Overfitting.



## Bidirectional RNN

- A typical state in an RNN relies on the past and present events..
- There can be situations where a prediction depends on the past present and future events.
- Take speech recognition, when you use a voice assistant, you initially utter a few words after which the assistant interprets and responds.
- This interpretation may not entirely depend on the preceding words; the whole sequence of words can make sense only when succeeding words are analyzed.
- To enable the past and ~~future~~ traversal of i/p bidirectional RNNs are used.
- Bidirectional RNN → combination of 2 RNNs  
one RNN moves forward beging from the start of the data sequence and other moves backward. beging from end of data sequence.

→ A BRNN has an additional hidden layer to accomodate the backward training process.

At any time  $t$ .

$$O_f = \phi(x_t \times \overset{\text{Forward}}{w_{ij}} + O_b \times \overset{\text{Forward}}{w_{kb}} + b_i) \rightarrow \text{Forward.}$$

$$O_n = \phi(x_t \times \overset{\text{backward}}{w_{ij}} + O_n \times \overset{\text{Backward}}{w_{hn}} + b_n) \rightarrow \text{Backward.}$$

where  $\phi$  = Activation function  
 $b$  = bias.

The hidden state at  $t$  is given by a combination of  $O_f$  (Forward) and  $O_n$  (Backward).

The O/p at any hidden state is defined as

$$O_h = H_t \times w_{ij} + b_j$$

$\uparrow$  new weight

In BRNN, since there is forward and backward passes happening simultaneously, updating the weights for two processes could happen at the same point of time. This leads to erroneous results.



