

# RNN (Recurrent Neural Network)

→ A type of sequential model to work on sequential data.

Q. What is sequential and Non-sequential data.

Non-sequential data

dataset
iq   munks   gender   pluicmonl

→ In this dataset, while feeding data to ANN, sequence does not matter, output doesn't depend on sequence.

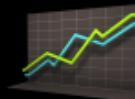
Sequential data

↳ Ex- text dataset

correct = "Hyy, my name is panth...!" } Sequence  
wrong = "Panth...! Hyy my is name" } matters...!!

More :-

- Time series based data



- Speech



- DNA sequence



→ Hence RNNs are highly useful in NLP (Natural Language Processing) tasks.



Q. Why RNN....?

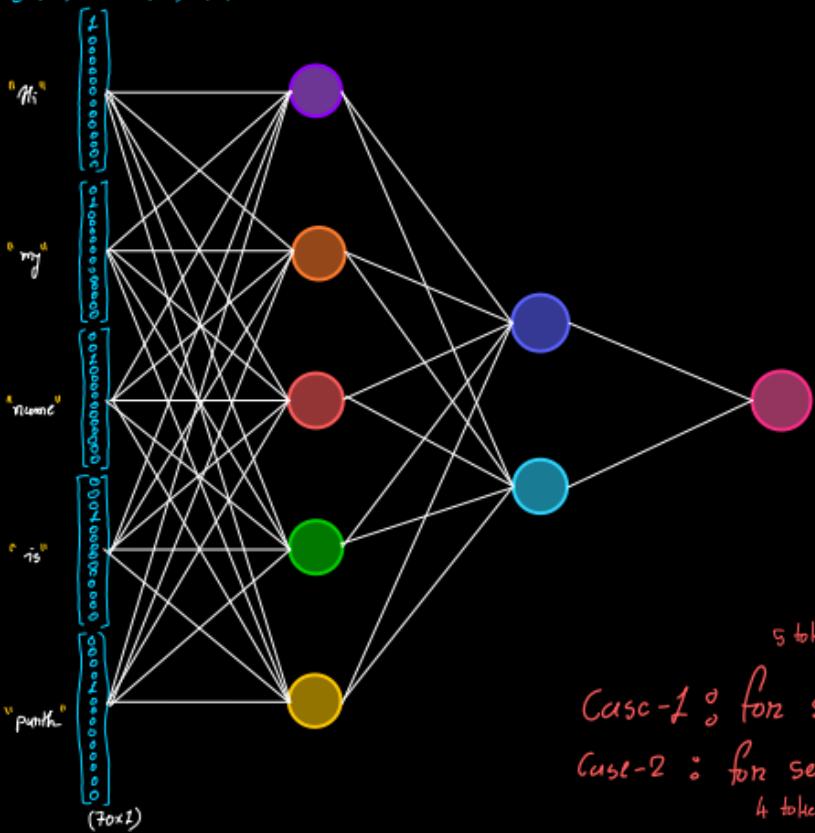
Input	Output	
Hi my name is purna	O	(5)
I love to learn AI/mL	O	(5)
India won the match	I	(4)

total 14 tokens

Vectorize formate

$$[1, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0] \rightarrow "Hi"$$

$$[0, 0, 0, 0, 0, 0, 0, 1, 0, 0, 0, 0, 0, 0, 0] \rightarrow "Love"$$



Case-1 : for sentence-1

Case-2 : for sentence-3  
4 tokens = 56 i/p

$$70 \times 5 = 350 \text{ weights}$$

$\triangle$  I/p is changing which is not

Due to this mismatch of sentences, we will use zero padding technique.

### Sentence-1

"Hi"  
"my"  
"name"  
"is"  
"punct."

### Sentence-2

"I"  
"love"  
"to"  
"learn"  
"AI/ML"

### Sentence-3

"India"  
"work"  
"the"  
"match"  
  
zero padding to  
maintain seq struc-  
ture.

### Problems

- Varying text input size
- Unnecessary computation
- Prediction problem.
- Doesn't follow sequence during input.

ANNs are not capable to remember the sequence of incoming data. hence we need RNN. because it works on sequence based data.

[text2data.com/Demo](http://text2data.com/Demo)

- ANN works on fixed i/p size while RNN can handle any length of sequential data.
- Sequence of data contains some meaning, in ANN while feeding data, i/p data is treated as separate entities, which means there is no sequence will be followed.
- RNNs have **memory feature** which means RNNs can remember the past i/p.
- I/p data has specific format while feeding it to the network.

format = (timesteps, i/p-features)

Review	Sentiment			unique words (vocabulary)
movie was <u>good</u>		1		[ movie, was, good, bad, not ]
movie was <u>bad</u>		0		
movie was <u>not</u> good		0		

Review\_1 (Movie was good)

$$\left[ [1, 0, 0, 0, 0], [0, 1, 0, 0, 0], [0, 0, 1, 0, 0] \right]$$

t=1 t=2 t=3

$$[1, 0, 0, 0, 0] = "movie"$$

$$[0, 1, 0, 0, 0] = "was"$$

$$[0, 0, 1, 0, 0] = "good"$$

$$[0, 0, 0, 1, 0] = "bad"$$

$$[0, 0, 0, 0, 1] = "not"$$

timesteps = Words

(3, 5)

No. of timesteps

No. of i/p features

(each word can be represented as 5 numbers)

Review\_2 (Movie was bad)

$$\left[ [1, 0, 0, 0, 0], [0, 1, 0, 0, 0], [0, 0, 0, 1, 0] \right] \rightarrow (3, 5)$$

Review\_3 (Movie was not good)

$$\left[ [1, 0, 0, 0, 0], [0, 1, 0, 0, 0], [0, 0, 0, 0, 1], [0, 0, 1, 0, 0] \right] \rightarrow (4, 5)$$

data input → format = (timesteps, i/p-features)

while using keras (SimpleRNN class) ↴

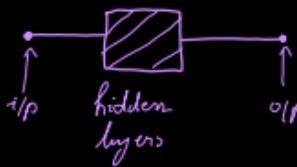
format = (Batch size, timesteps, i/p-features) - 0 user reviews and 1 ground truth

Input shape =  $(3, 4, 5) = 3D$  tensor

Number of reviews      Max size of review      No. of input features (OneHotEncoded representation)  
each word is represented as 5 numbers.

Q. How RNN works....?

data			Sentiment	Numeric representation	
Review				movie	[1, 0, 0, 0, 0]
$x_1$	movie	$x_{12}$ was	$x_{13}$ good	1	was [0, 1, 0, 0, 0]
$x_2$	movie	$x_{21}$	$x_{22}$	0	good [0, 0, 1, 0, 0]
$x_3$	movie	$x_{31}$	$x_{32}$	0	bad [0, 0, 0, 1, 0]
		$x_{33}$ not	$x_{34}$ good		not [0, 0, 0, 0, 1]



→ I/p will be passed to RNN based on time

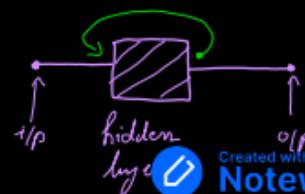
time(t)	input
0	$x_{11}$
1	$x_{21}$
2	$x_{31}$

ANN

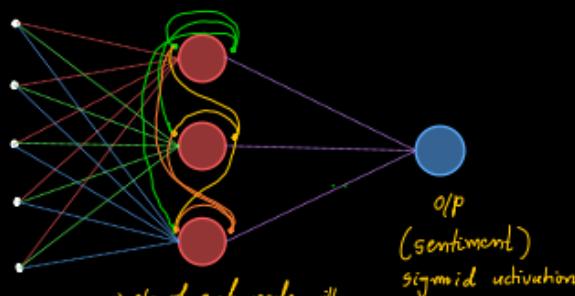
- like ANN, we don't feed entire data set once.
- A feed forward NN. (data flows from I/p to O/p direction).

RNN

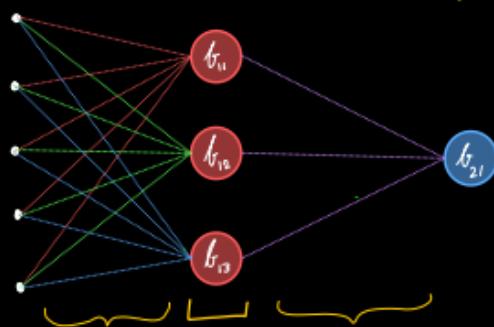
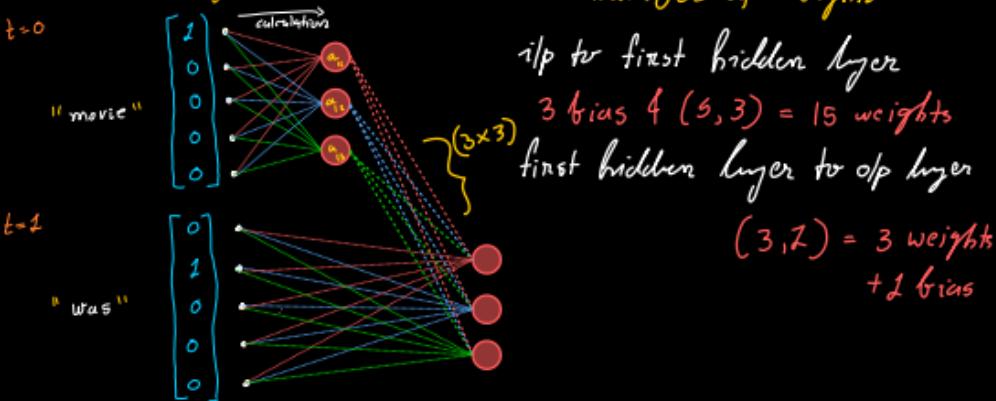
- I/p will be given to RNN based on time steps.
- Hidden layer sends feedback on backward direction, which is the major difference between ANN & RNN.



## Structure of RNN



(because, each token is represented as total 5 training values)



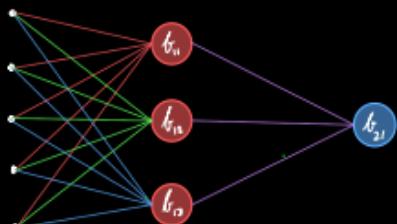
$$5 \times 3 + 3 \times 3 + 3 = 27 \text{ weights}$$

# Q. forward Propagation

review

			sentiment
$x_{11}$	$x_{12}$	$x_{13}$	2
$x_{21}$	$x_{22}$	$x_{23}$	0
$x_{31}$	$x_{32}$	$x_{33}$	0

Where each of  $x_{ij}$  is 5D vector  
 $[1, 0, 0, 0, 0]$



Default activation of each neurons will be zero.

- dot-product  $(x_{11}, w_i)$  where  $x_{11}$  is  $1 \times 5$  and  $w_i$  is  $5 \times 3$ , hence op will be  $1 \times 3$ .

- the op of this dot product will be given to an activation function  $f(x_{11} \cdot w_i + b_2) = O_2$ .

- till now, the process was for  $t=0$ .

- Now for  $t=1$

- i/p =  $x_{12}$

- in between  $O_2 \cdot w_k$  will be work as response for hidden layer.

- final calculation would be  $f(x_{12} \cdot w_i + O_2 w_k + b_2) = O_2$

$$\begin{matrix} & & & \\ & 2 \times 5 & 5 \times 3 & 2 \times 3 & 3 \times 3 \\ \diagdown & & \diagdown & & \diagdown \\ 2 \times 3 & & 2 \times 3 & & 3 \times 3 \\ & & & & \\ & (2 \times 3 + 2 \times 3) = 2 \times 3 & & & \end{matrix}$$

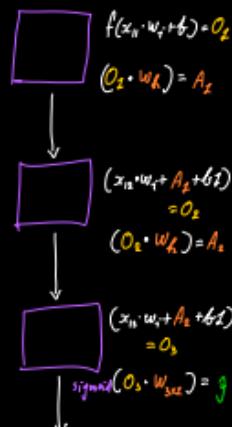
- Now for  $t=2$ .

- i/p =  $x_{13}$

- in between  $O_2 \cdot w_k$

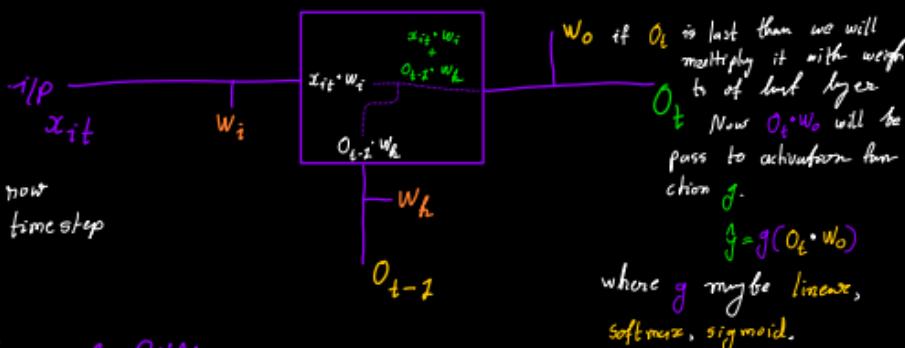
- final calculation  $f(x_{13} \cdot w_i + O_2 w_k + b_2) = O_2$

- Now we will move forward to next layer & multiply  $O_2$  with weights of next layer  $w_k$  hence  $O_2 \cdot w_k = g$  now  $\text{sigmoid}(g)$  will give final result.



## Simplified Representation



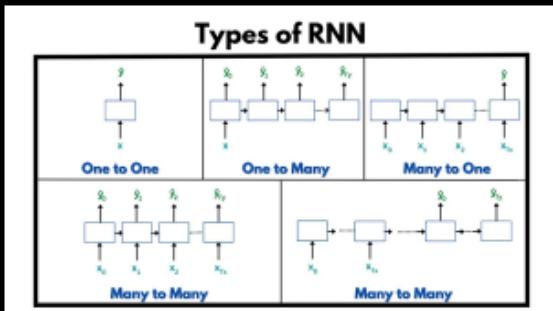


## Types of RNNs.

- 1. One-to-One
- 2. Many-to-one
- 3. Many-to-many
- 4. One-to-Many

### 1. Many-to-one

- Input : sequence (sentence, characters, time series).
- Output : Non-sequential output (Number/Scalar).



### 2. One-to-many

- Input : Non-sequential (Numbers/Image)
- Output : Sequential (sentence, characters, time series)

Input = once, Output = multiplied times

Example : Image captioning

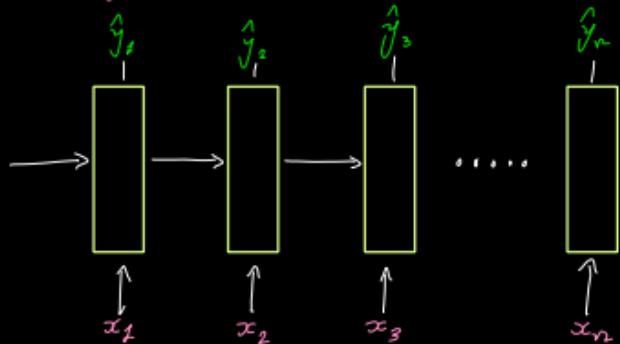
### 3. Many-to-Many

- Input : Sequential data
- Output : Sequential data

Seq 2 Seq.

Seq 2 Seq has 2 types.

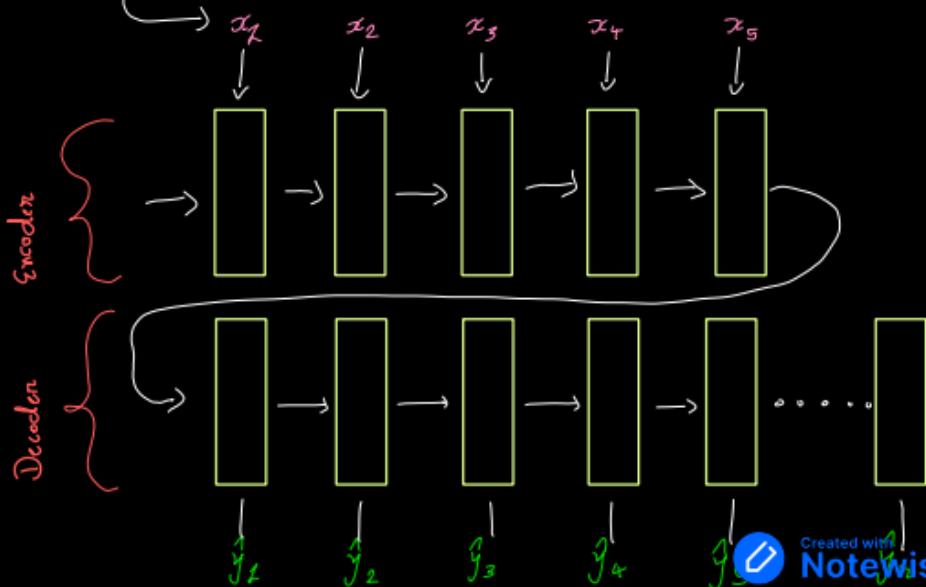
→ Same length [zip seq. size = o/p seq. size] O/p Numeral quantity never



→ Variable length [Machine translation]  
• language translation models.

### Example

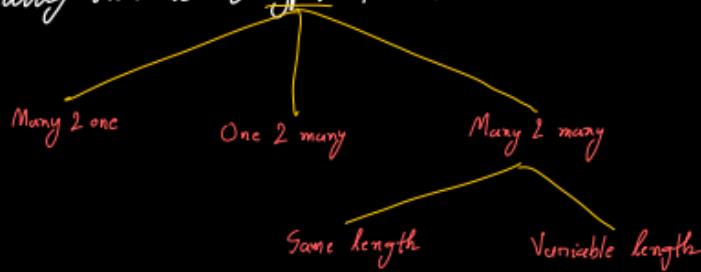
i/p : Hi, my name is paul.  
o/p : 嗨，我的名字是 Paul！



### 3 One-to-One

- input : Non-sequential
- output : Non-sequential      ex: image classification
- It looks like ANN or CNN.

Hence, technically there is 3 types of RNN.



### Backpropagation in RNN

→ Backpropagation in RNN → called as BPTT (Back Propagation Through Time).

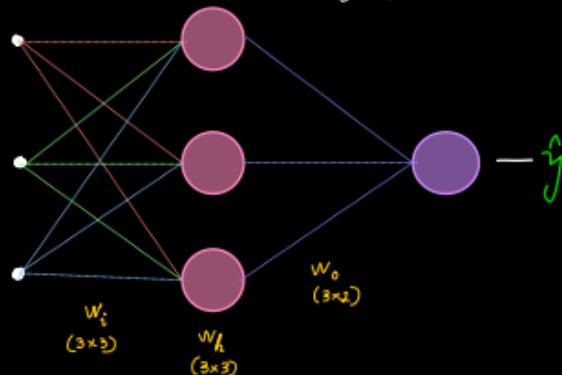
Example text

cat	must	rat	1
rat	rat	must	2
must	must	cut	0

Vocabulary

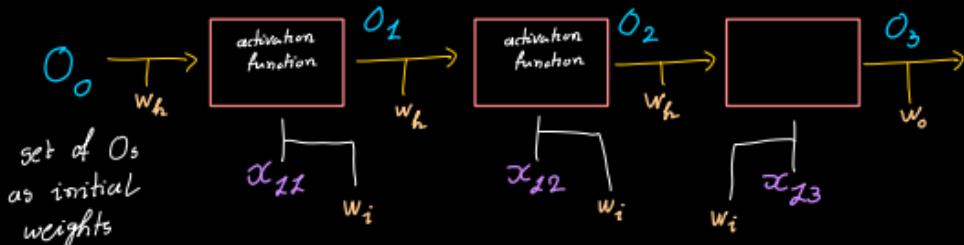
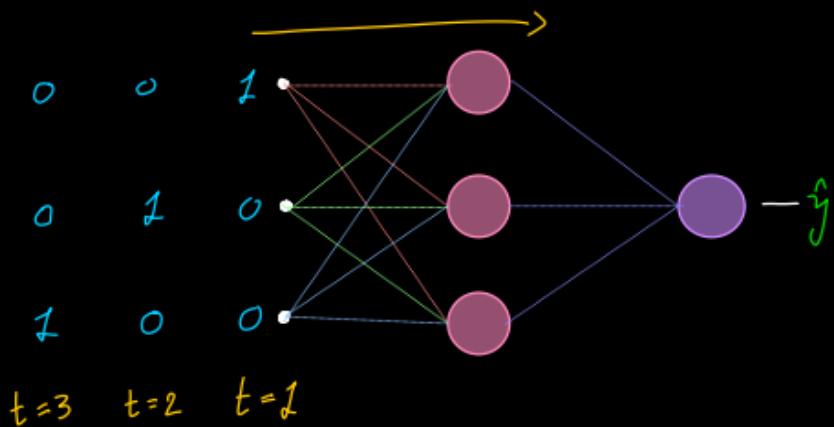
cat must rat  
[200] [020] [001]

x	y
$x_1$ [100] [010] [001]	2
$x_2$ [001] [001] [010]	1
$x_3$ [010] [010] [100]	0



forward propagation for....

$x_1 [100] [010] [001]$



where:

$$O_1 = f(x_{21}w_i + O_0w_h) \quad \text{loss} = -y_i \log \hat{y}_i - (1-y_i) \log (1-\hat{y}_i)$$

$$O_2 = f(x_{22}w_i + O_1w_h)$$

$$O_3 = f(x_{23}w_i + O_2w_h)$$

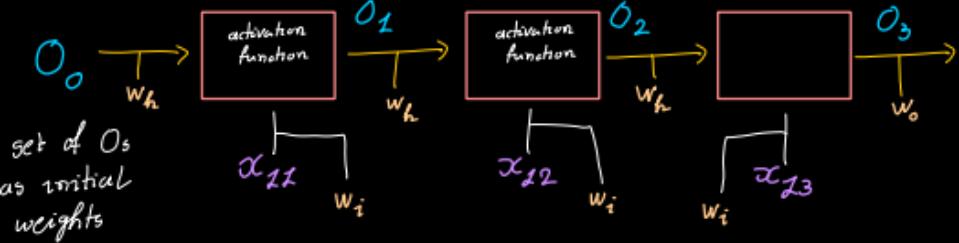
Afterwards we have done forward propagation, now we need to find best optimal value of  $w_i$ ,  $w_h$  and  $w_o$ .

$$\hat{y} = \sigma(O_3 w_o)$$

$$w_i = w_i - \alpha \frac{\partial L}{\partial w_i} \quad | \quad w_h = w_h - \alpha \frac{\partial L}{\partial w_h}$$

$$w_o = w_o - \alpha \frac{\partial L}{\partial w_o}$$





where:

$$O_1 = f(x_{11}w_i + O_0w_h)$$

$$O_2 = f(x_{21}w_i + O_1w_h)$$

$$O_3 = f(x_{31}w_i + O_2w_h)$$

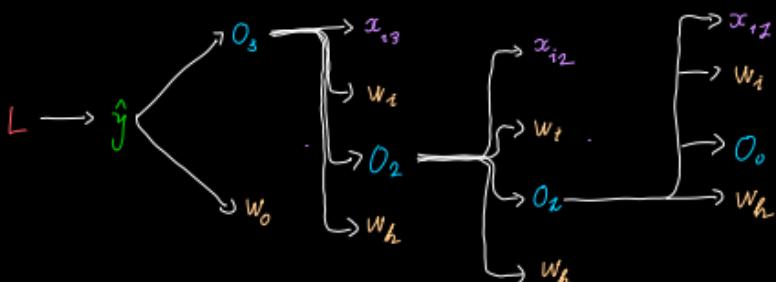
$$\hat{y} = \sigma(O_3w_o)$$

$$\frac{\partial L}{\partial W_o}, \frac{\partial L}{\partial W_i}, \frac{\partial L}{\partial W_h}$$

$$\frac{\partial L}{\partial W_o} = L \rightarrow \hat{y} \begin{cases} \nearrow O_3 \\ \searrow w_o \end{cases}$$

$$= \frac{\partial L}{\partial \hat{y}} \frac{\partial \hat{y}}{\partial W_o}$$

### full path



for  $\frac{\partial L}{\partial W_i}$  we have 3 different paths.

$$L \rightarrow \hat{y} \rightarrow O_3 \rightarrow w_i \quad \frac{\partial L}{\partial w_i} = \frac{\partial L}{\partial \hat{y}} \frac{\partial \hat{y}}{\partial O_3} \frac{\partial O_3}{\partial w_i}$$

$$L \rightarrow \hat{y} \rightarrow O_3 \rightarrow O_2 \rightarrow w_i \quad \frac{\partial L}{\partial w_i} = \frac{\partial L}{\partial \hat{y}} \frac{\partial \hat{y}}{\partial O_3} \frac{\partial O_3}{\partial O_2} \frac{\partial O_2}{\partial w_i}$$

$$L \rightarrow \hat{y} \rightarrow O_3 \rightarrow O_2 \rightarrow O_1 \rightarrow w_i \quad \frac{\partial L}{\partial w_i} = \frac{\partial L}{\partial \hat{y}} \frac{\partial \hat{y}}{\partial O_3} \frac{\partial O_3}{\partial O_2} \frac{\partial O_2}{\partial O_1} \frac{\partial O_1}{\partial w_i}$$

hence, final path is ...

$$\frac{\partial L}{\partial W_i} = \frac{\partial L}{\partial \hat{y}} \frac{\partial \hat{y}}{\partial O_3} \frac{\partial O_3}{\partial w_i} + \frac{\partial L}{\partial \hat{y}} \frac{\partial \hat{y}}{\partial O_3} \frac{\partial O_3}{\partial O_2} \frac{\partial O_2}{\partial w_i} + \frac{\partial L}{\partial \hat{y}} \frac{\partial \hat{y}}{\partial O_3} \frac{\partial O_3}{\partial O_1} \frac{\partial O_1}{\partial O_2} \frac{\partial O_2}{\partial w_i}$$



$$\frac{\partial L}{\partial w_i} = \sum_{j=1}^3 \frac{\partial L}{\partial \hat{y}} \frac{\partial \hat{y}}{\partial o_j} \frac{\partial o_j}{\partial w_i}$$

: where 3 is number of words

$$\text{for } j=1 : \frac{\partial L}{\partial w_i} = \frac{\partial L}{\partial \hat{y}} \frac{\partial \hat{y}}{\partial o_1} \frac{\partial o_1}{\partial w_i}$$

$$\hookrightarrow \frac{\partial \hat{y}}{\partial o_1} = \frac{\partial \hat{y}}{\partial o_3} \frac{\partial o_3}{\partial o_2} \frac{\partial o_2}{\partial o_1}$$

for  $j=2$

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

$$\hookrightarrow \frac{\partial \hat{y}}{\partial o_2} = \frac{\partial \hat{y}}{\partial o_3} \frac{\partial o_3}{\partial o_2}$$

for  $j=3$

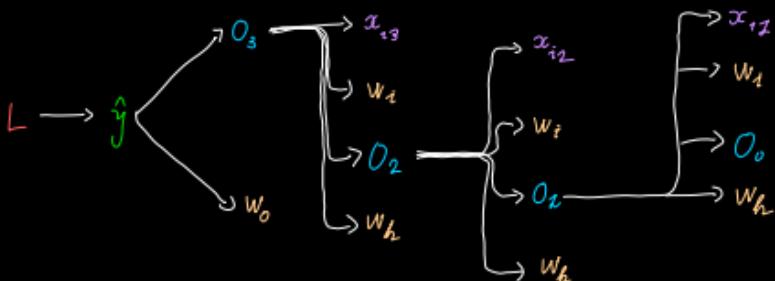
$$\frac{\partial L}{\partial w_i} = \frac{\partial L}{\partial \hat{y}} \frac{\partial \hat{y}}{\partial o_3} \frac{\partial o_3}{\partial w_i}$$

$$\hookrightarrow \frac{\partial \hat{y}}{\partial o_3}$$

hence

$$\frac{\partial L}{\partial w_i} = \sum_{j=1}^n \frac{\partial L}{\partial \hat{y}} \frac{\partial \hat{y}}{\partial o_j} \frac{\partial o_j}{\partial w_i}$$

## full path



For  $\frac{\partial L}{\partial w_h}$  we have two different paths.

$$L \rightarrow \hat{y} \rightarrow O_3 \rightarrow W_h$$

$$L \rightarrow \hat{y} \rightarrow O_3 \rightarrow O_2 \rightarrow W_h$$

$$L \rightarrow \hat{y} \rightarrow O_3 \rightarrow O_2 \rightarrow O_1 \rightarrow W_h$$

hence.....

$$\frac{\partial L}{\partial w_h} = \frac{\partial L}{\partial \hat{y}} \frac{\partial \hat{y}}{\partial O_3} \frac{\partial O_3}{\partial w_h} + \frac{\partial L}{\partial \hat{y}} \frac{\partial \hat{y}}{\partial O_3} \frac{\partial O_3}{\partial O_2} \frac{\partial O_2}{\partial w_h} + \frac{\partial L}{\partial \hat{y}} \frac{\partial \hat{y}}{\partial O_3} \frac{\partial O_3}{\partial O_2} \frac{\partial O_2}{\partial O_1} \frac{\partial O_1}{\partial w_h}$$

$$\frac{\partial L}{\partial w_h} = \sum_{j=1}^n \frac{\partial L}{\partial \hat{y}} \frac{\partial \hat{y}}{\partial O_j} \frac{\partial O_j}{\partial w_h}$$

Where  $n = \text{number of timesteps / words.}$



## Problems with RNNs.

- o RNN suffers from few problems, hence they are not widely used at that level.
- o Hence we uses some other architectures like LSTMs.
- o Problems

1. Problem of long term dependency.
2. Stagnated training

} Unstable gradients

- o RNN is useful for sequential data where each term is depend on its previous term.
- o If the sequence/sentence gets too longer than the last term can't remember the initial term. which is like short term memory loss.

### Example

1. ગરજાતિ રિ સ્પોકન ઇન ગુજરાત.

→ Here word ગરજાતિ રિ depend on another previous word ગરજાતિ, which is not far from it. so the word ગરજાતિ રિ predictable based on its near by word ગરજાતિ રિ.

2. Germany is beautiful place. I went there last year. But i couldn't enjoyed because i don't understand german.

→ As we can see above example of long term dependency where german word depends on Germany.

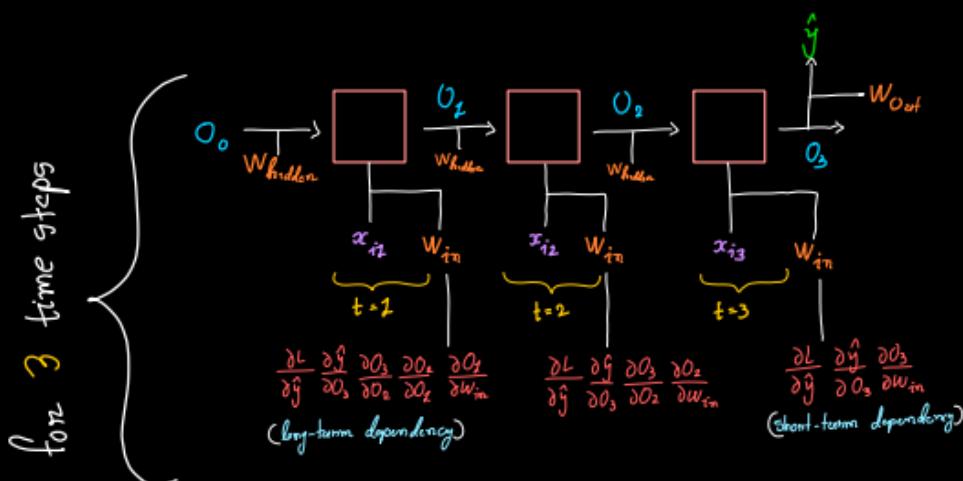
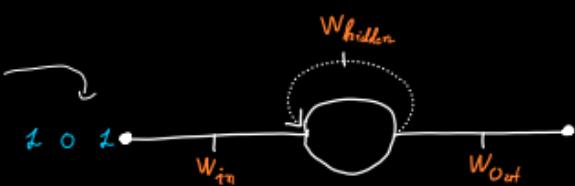
→ This kind of situations RNN can't work properly.

→ This problem is called as long term dependency problem. which happens due to Vanishing Gradient Problem.



## Problem #1 Long term dependency

input	output
1 0 1	1
0 0 1	0
0 0 0	0
1 1 1	1



Lets imagine for 100 time-steps.

$$\frac{\partial L}{\partial w_{in}} = \frac{\partial L}{\partial \hat{y}} \frac{\partial \hat{y}}{\partial O_{100}} \frac{\partial O_{100}}{\partial O_{99}} \frac{\partial O_{99}}{\partial O_{98}} \dots \frac{\partial O_2}{\partial O_1} \frac{\partial O_2}{\partial w_{in}}$$

→ as we go to calculate gradient, these long term dependencies gets smaller. in this situation short-term dependencies gives more contribution for calculating  $\frac{\partial L}{\partial w_{in}}$ .

→ as time steps increases (as the term gets bigger) long-term dependencies will get more small so that's why the responsibility of derivative calculation will be on short-term dependencies.



$$\frac{\partial L}{\partial W_{in}} = \frac{\partial L}{\partial \hat{y}} \frac{\partial \hat{y}}{\partial O_{100}} \frac{\partial O_{100}}{\partial O_{99}} \frac{\partial O_{99}}{\partial O_{98}} \frac{\partial O_{98}}{\partial O_{97}} \dots \dots \dots \frac{\partial O_2}{\partial O_1} \frac{\partial O_1}{\partial W_{in}}$$

Let's generalize this equation.

$$\frac{\partial L}{\partial \hat{y}} \frac{\partial \hat{y}}{\partial O_{100}} \prod_{t=2}^{100} \left( \frac{\partial O_t}{\partial O_{t-1}} \right) \frac{\partial O_1}{\partial W_{in}}$$

Where:

derivative of tanh

$$O_t = \tanh(x_{it} w_{in} + O_{t-1} w_h)$$

$$\frac{\partial O_t}{\partial O_{t-1}} = \tanh'(x_{it} w_{in} + O_{t-1} w_h) w_h$$

derivative of tanh will be between 0 to 1.

Now run long-term dependency.

$$\frac{\partial L}{\partial \hat{y}} \frac{\partial \hat{y}}{\partial O_{100}} \prod_{t=2}^{100} \left( \tanh'(x_{it} w_{in} + O_{t-1} w_h) w_h \right) \frac{\partial O_1}{\partial W_{in}}$$

As timestamp increases, gradients which are at very long distance gets more smaller. This is called as vanishing gradient problem.

Solution

- o Use different activation function.
- o Better weight initialization techniques.
- o Use skip RNNs.
- o LSTM



## Problem #2 Unstable training (Exploding Gradient)

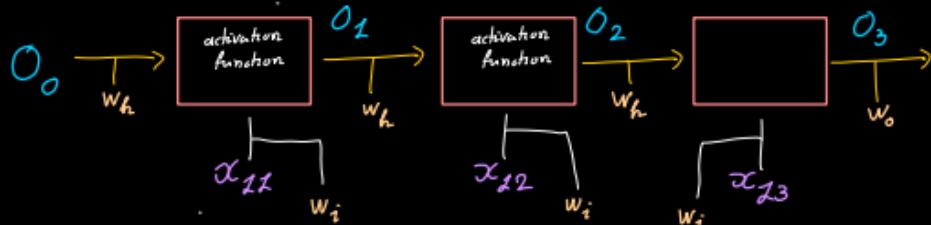
- What if, over long term dependency get such biggen that it dominate short-term dependency or become infinite.
- Hence gradient update will become infinite & learning won't occur
- Suppose we have initialized our weight as 1. with each multiplication this number will become more larger.

### Solution

- Gradient Clipping (set maximum limit)
- Controlled learning-rate
- LSTM



# LSTM (Long Short Term Memory)



Above shown is normal workflow of RNN which can only hold short term context.

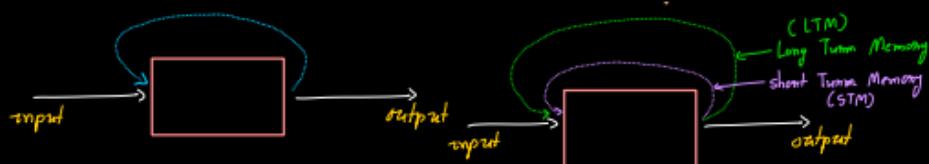
- To store long term context, RNN also maintains a long term memory called as cell state.
- RNN has simple architecture, while LSTM has a little bit complex architecture.

↙ cell state



→ Simple RNN architecture.

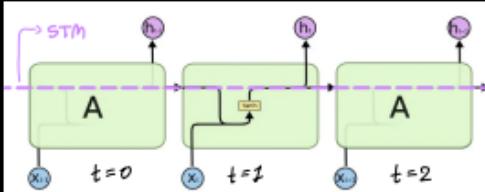
→ LSTM architecture



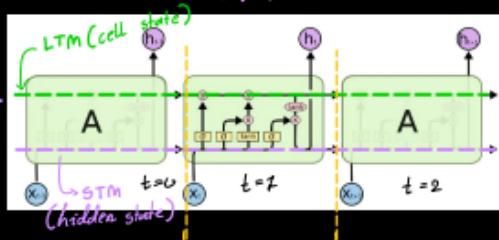
- To maintain state we need a complex architecture which can maintain communication between LTM and STM.



RNN



LSTM



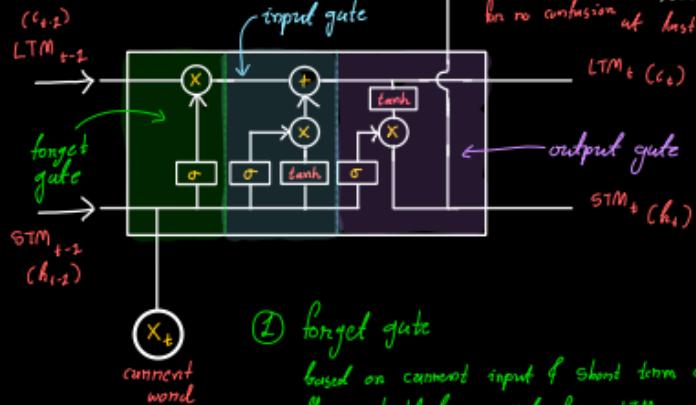
for  $t=t$

$$h_t \rightarrow$$

same  
(STM<sub>t</sub>)

be no confusion  
at last node.

This circuit is specially designed for establishing communication between STM and LTM.



### ② forget gate

based on current input of short term context, it decides what things should be removed from LTM.

### ③ input gate

based on current word input, it decides that which things should be added in LTM.

### ④ Output gate

Based on current input, it decides to return LTM as output after removing something unnecessary things from LTM.

→ Now at last stage & based on current cell state, it will decide the output.

(content in LTM)

→ At each given stage, it creates new STM.

#### Task-1

Update LTM by removing unnecessary thing & adding important things.

#### Task-2

Create STM from next stage.

#### Input

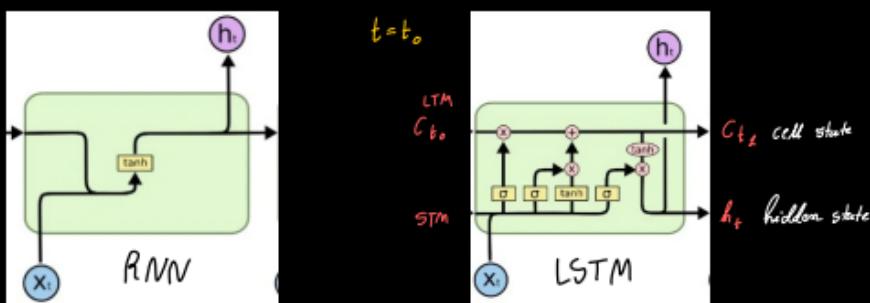
- $C_{t-1}$  LTM of previous state
- $h_{t-1}$  previous hidden state
- $x_t$  current word

#### Output

- $C_t$  current cell state
- $h_t$  current hidden state



# LSTM Architecture



→ LSTM Maintains 2 different / interconnected state / Memory at a time called as

1. LTM (Long Term Memory) also called as Cell state.

2. STM (Short Term Memory) also called as Hidden State.

→ LSTM keeps interaction alive between LTM and STM.

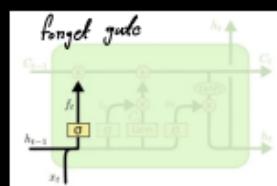
→ Within this process, 2 task is happening.

1. Update cell state.  $\overbrace{\text{forget}}^{\text{remove}}$

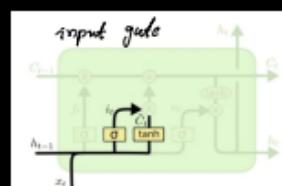
2. Calculating hidden state ( $h_t$ )

→ based on  $x_t$ , it decides that which thing should / should-not be in cell state.

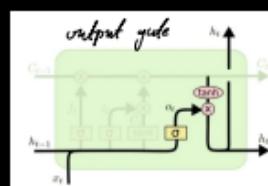
## The Gates in LSTM



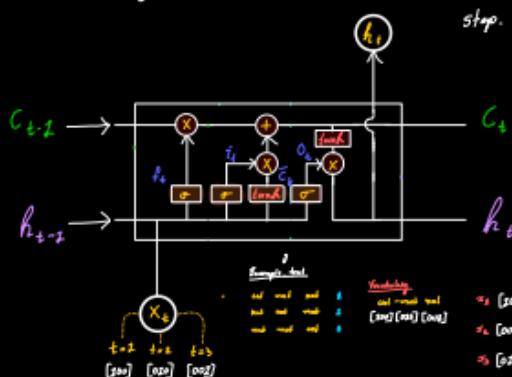
→ To remove something from cell state.



→ To add something to cell state.



→ To calculate hidden-state ( $h_t$ ) for next time step.



$\tilde{C}_t$  = potential important info that should be added in long term context.

→ Mathematically  $h_t$  and  $C_t$  are vectors of 3 dimension [0.2, 0.3, 0.2].

→  $h_t$  and  $C_t$  will always have same dimension but may contain different numbers.

→  $X_t$  is set of input words (already seen in RNN).

→  $f_t$ ,  $i_t$ ,  $C_t$  and  $O_t$  are part of LSTM cell.

↳ candidate cell state.

all of them are vectors. and dimension of all of 4 them will be same. this is a rule.

Pointwise operation

→ Major 3 operations are happening in cell

1.  $\otimes$

↳ example

$$C_{t-1} = [1, 2, 3] \quad f_t = [4, 5, 6]$$

$$o/p = [4, 20, 24]$$

2.  $\oplus$

↳ example

$$C_{t-1} = [1, 2, 3] \quad f_t = [4, 5, 6]$$

$$o/p = [5, 4, 9]$$

3. tanh

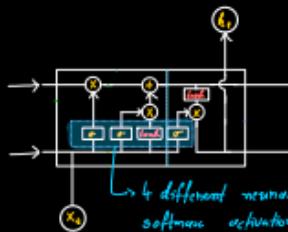
↳ example

$$C_{t-1} = [1, 2, 3] \quad f_t = [4, 5, 6]$$

$$\text{tanh} \quad \text{tanh}$$

$$[0.26159, 0.73642, 0.997505] \quad [0.99112, 0.49440, 0.99448]$$

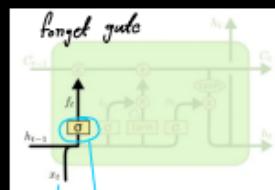
## Neural Network Layers



→ 4 different neural network architectures where 3 ANN with sigmoid activation & 1 with softmax activation. Number of neurons is hyper parameter, but once we finalize the number of neurons, each ANN will have same number of neurons.

Now let's understand each gate.

### #1 The Forget Gate



input

1.  $C_{t-1}$  : Previous cell state
2.  $h_{t-1}$  : Previous hidden state
3.  $x_t$  : Current word

Output

→ updated cell-state after removing unnecessary context.

→ suppose there are 3 units/neurons & activation = sigmoid

→ At current time-step  $x_t$  is 4 dimensional vector.  $[x_{t1}, x_{t2}, x_{t3}, x_{t4}]$

→  $h_{t-1}$  and  $C_{t-1}$  will be of 3 dimension. because the dimension of  $h_{t-1}$  &  $C_{t-1}$  will always be equal to the number of neurons in ANN layer.

→ This gate will work on 2 stages.

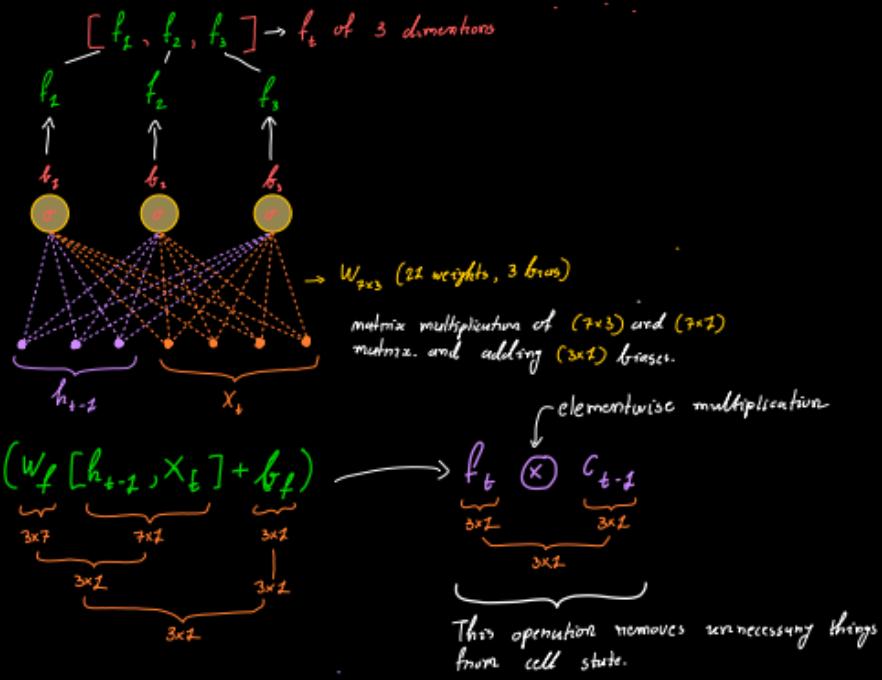
Stage-1 : calculate  $f_t$

Stage-2 :  $C_{t-1} \times f_t$  (element wise multiplication)

removing something from cell state.



Now let's see the ANN architecture of forget gate.  
as we discussed that ANN has 3 neurons with each bias sigmoid as activation function.



example  $f_t \quad c_{t-1} \quad \text{O/P}$

$[1, 1, 1] [4, 5, 6] [2, 2, 5, 3] \rightarrow$  cell state is reduced by half

$[1, 1, 1] [4, 5, 6] [4, 5, 6] \rightarrow$  No changes occurred.

$[0, 0, 0] [4, 5, 6] [0, 0, 0] \rightarrow$  cell state is empty now

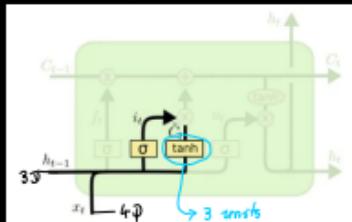
## #2 The Input Gate

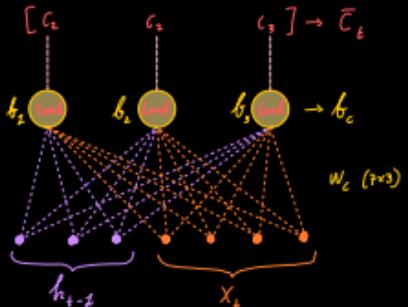
→ This gate works on 3 stages.

Stage-1 calculate  $i_t$  candidate cell state.

Stage-2 calculate  $s_t$ , where  $i_t$  decides that which value will be added in cell state from candidate cell state.

Stage-3 calculate  $c_t$  from previous incoming state  $[f_t \odot c_{t-1}]$ .

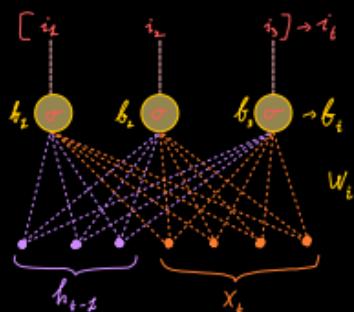




$$\text{tanh}\left(W_f \cdot [h_{t-1}, x_t] + b_f\right)$$

$\underbrace{3 \times 3}_{3 \times 1}$        $\underbrace{3 \times 2}_{3 \times 1}$        $\underbrace{3 \times 1}_{3 \times 1}$   
 $\underbrace{3 \times 2}_{3 \times 1}$        $\underbrace{3 \times 2}_{3 \times 1}$

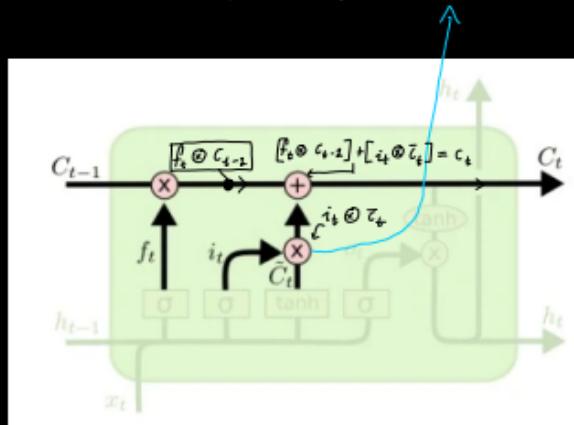
Now let's calculate it which filters the content from  $C_{t-1}$  & adds final filtered app in cell state.  
 → ANN structure for calculate will be same as all previous architecture with sigmoid activation.



$$\text{sigmoid}\left(W_i \cdot [h_{t-1}, x_t] + b_i\right)$$

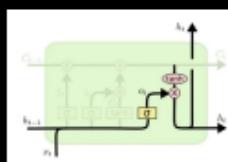
$\underbrace{3 \times 3}_{3 \times 1}$        $\underbrace{3 \times 2}_{3 \times 1}$        $\underbrace{3 \times 1}_{3 \times 1}$   
 $\underbrace{3 \times 2}_{3 \times 1}$

$i_t \otimes \bar{C}_t \rightarrow \text{filtered } \bar{C}_t$



$$\text{final Cell State } (C_t) = [f_t \otimes C_{t-1}] + [i_t \otimes \bar{C}_t]$$

## #3 The Output Gate



Task : to decide the o/p or also we can say "decide the value of hidden state ( $h_t$ )" at current time stamp.

Hence, main point is  $h_t$  is calculated by applying some mathematical operation on final cell state.

→ There are 2 steps to calculate hidden state of current timestamp by using current cell-state.

Step - 1 :

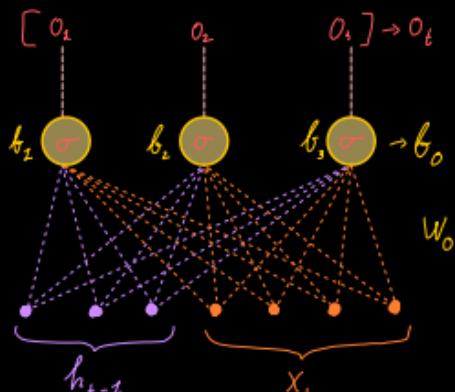
- o apply tanh operation on each element of  $c_t$  that bring them between -1 to +1.

$$\text{tanh}(c_t)$$

Step - 2 :

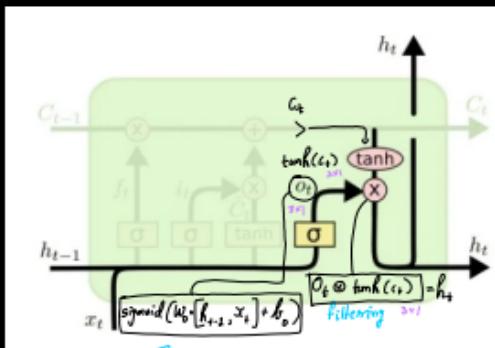
- o apply filtering using  $O_t$  and  $\text{tanh}(c_t)$ .

$$h_t = [O_t \odot \text{tanh}(c_t)]$$



$$\text{sigmoid}(\sigma) \left( W_o \cdot \left[ h_{t-1}, x_t \right] + b_o \right)$$

$\underbrace{W_o}_{3 \times 7} \quad \underbrace{\left[ h_{t-1}, x_t \right]}_{7 \times 2} \quad \underbrace{b_o}_{3 \times 1}$   
 $\underbrace{\sigma}_{\sigma(3 \times 1)} \quad \underbrace{3 \times 2}_{3 \times 2}$



And that is the output of current state.....

