

## Why RNN's are needed?

09 October 2024 15:53

### Sequential data

ANN → tabular data

CNN → image data

RNN → Recurrent Neural Network

6 type of sequential model to work on  
sequential data

#### Non-sequential data

Eg: Suppose we have data

ig | marks | gender | placement

ig . O  
marks O O  
gender' O

↑  
sequence of input  
doesn't matter

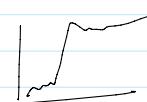
#### Sequential data

Text

Hi my name is Sakalya

→ sequence matters for grammatical correctness

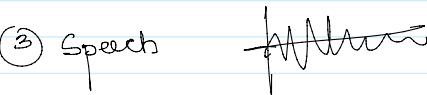
① Text



② Time Series data



③ Speech



④ DNA sequence

RNNs are mostly used in NLP.

### The why of RNN?

Input                      Output

Hi my name is Sakalya

O

12

O

↳ Input size  
varies

I love campus

O

12

O

India won the match

1

12

O

↓ ANN cannot  
handle  
varying input  
size

vectorize (OHE)

unique words

[1, 0, 0, ...] [0, 1, ..., 0]

H?

my

12

O

∴  $12 \times 5 \times 4 < 240$  weights

But now for the next sentence input size is 36248

- Textual data is of varying size

One walkaround

possible

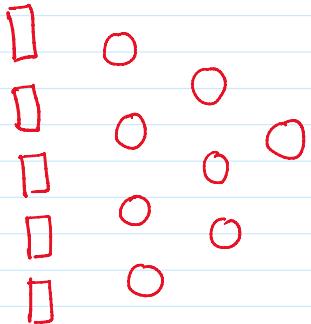
→ zero padding

Select the sentence having max words

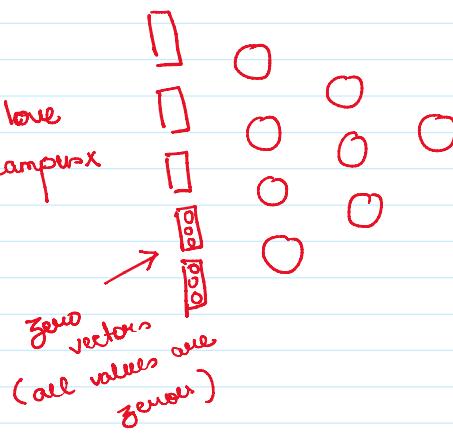
For smaller sentences use zero vectors

Example

Hi my name is  
Satyayya



I love  
campus



Zero  
vectors  
(all values are  
zeroes)

Problems

- this is not a good solution
- leads to lot of unnecessary computation
- Prediction problem → if predicting sentence has 200 words
- completely discarding semantic meaning → ANN has no memory to remember sequence

# Forward Propagation in RNN

09 October 2024 16:21

## Data for RNN

$(t, \text{timesteps}, \text{input-features})$

$t=1$   
review-1

[10000] [01000] [00100]

(3,5)  
↓  
no. of features  
no. of timesteps

	Review	Sentiment
(3,5)	movie was good	1
(3,5)	movie was bad	0
(4,5)	movie was not good	0

vocab  $\rightarrow$  6 words

movie was good bad  
[10000] [01000] [00100] [00010]  
[00001] not

Keras  $\rightarrow$  SimpleRNN  $\rightarrow$  batchsize, timesteps, input-features

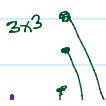
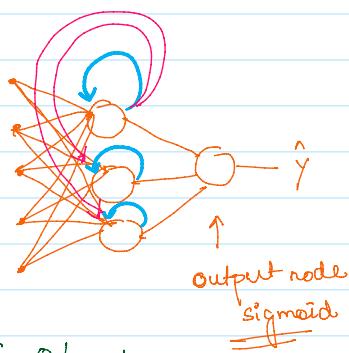
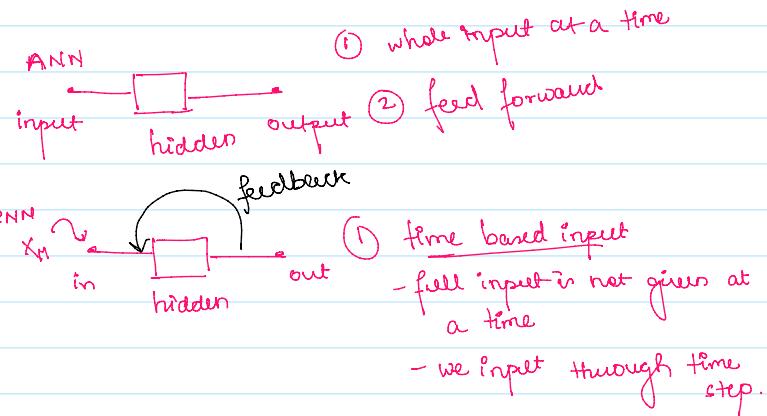
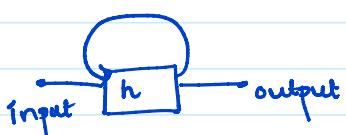
(3, 4, 5)

longest timestep is taken

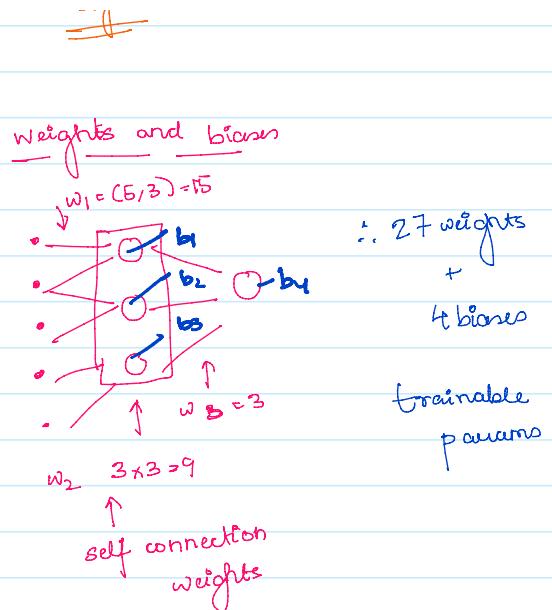
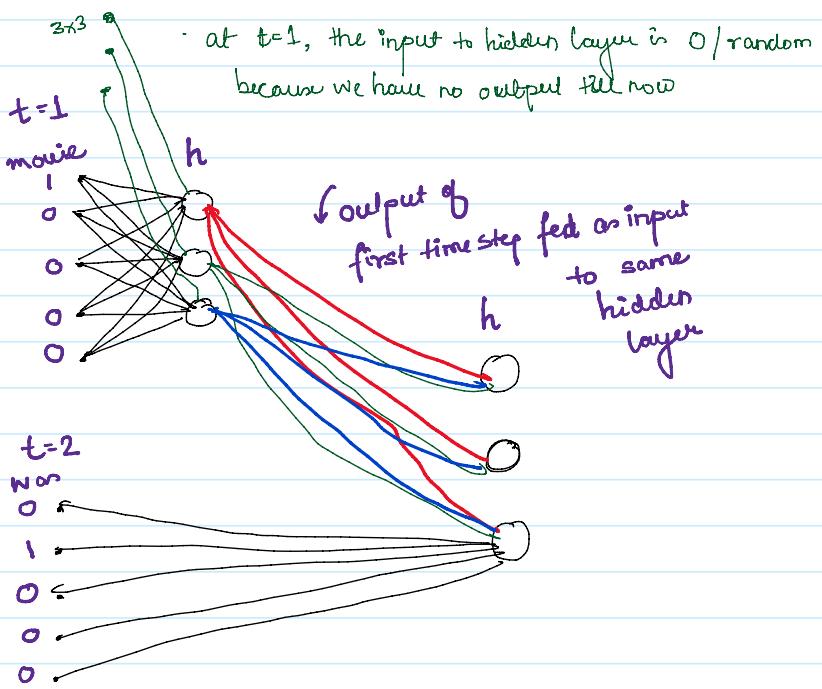
## How RNN works?

	Review	Sentiment
$x_1$	$x_{11}$ movie was good	1
$x_2$	$x_{21}$ movie was bad	0
$x_3$	$x_{31}$ movie was not good	0

movie was good bad  
[10000] [01000] [00100] [00010]  
[00001] not



at  $t=1$ , the input to hidden layer is 0/random  
because we have no output till now



### Code example

```
from keras import Sequential
from keras.layers import Dense, SimpleRNN

model = Sequential()
model.add(SimpleRNN(3, input_shape=(4,5)))
model.add(Dense(1, activation='sigmoid'))
model.summary()

final dense layer for output
```

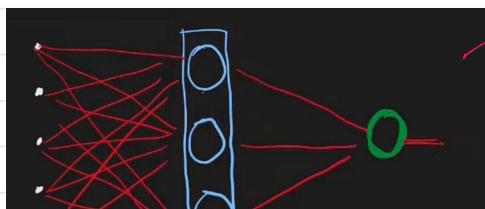
3 hidden layer RNN  
more 4 time steps  
5 input features  
final dense layer for output

Layer (type)	Output Shape	Param #
simple_rnn (SimpleRNN)	(None, 3)	27
dense (Dense)	(None, 1)	4

Total params: 31  
Trainable params: 31 ← trainable params are 31 as we calculated  
( $27+4$ )  
Non-trainable params: 0

### Forward Propagation

reviews	Sentiment
$x_{11}, x_{12}, x_{13}$	1
$x_{21}, x_{22}, x_{23}$	0
...	0

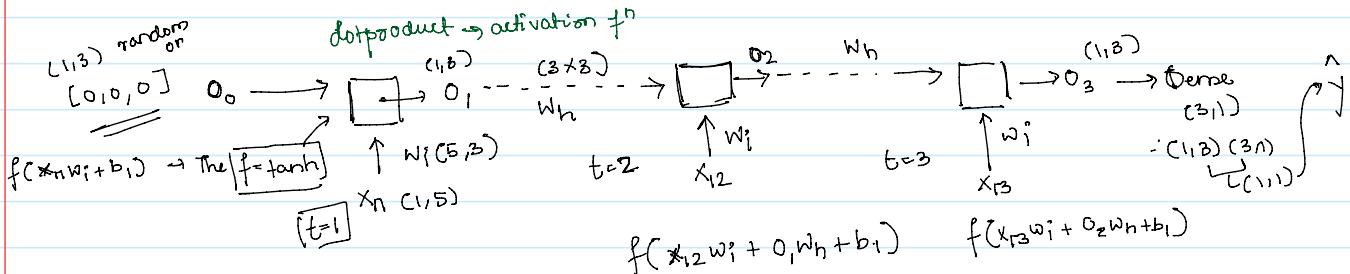


forward propagation  
unfolding through time  
works as a loop

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

vectors

5-dimension



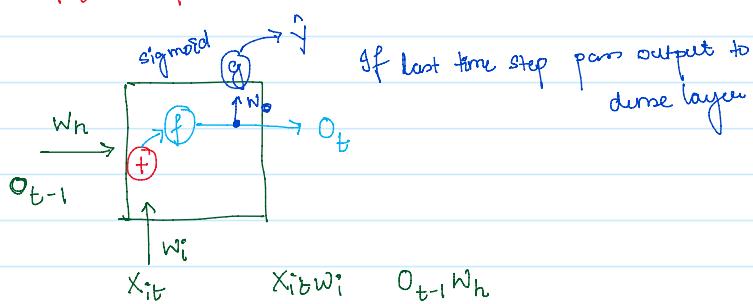
$$\begin{matrix} x_{12} & w_i \\ (1 \times 5) & (5 \times 3) \end{matrix} \quad \begin{matrix} o_1 & w_h \\ (1 \times 3) & (3 \times 3) \end{matrix}$$

$$(1 \times 3) + (1 \times 3) \rightarrow (1 \times 3)$$

### Advantages

1. Parameter / Weight sharing  $\rightarrow$  same weights are used for every input
2. RNN can store sequential information upto 10 time steps

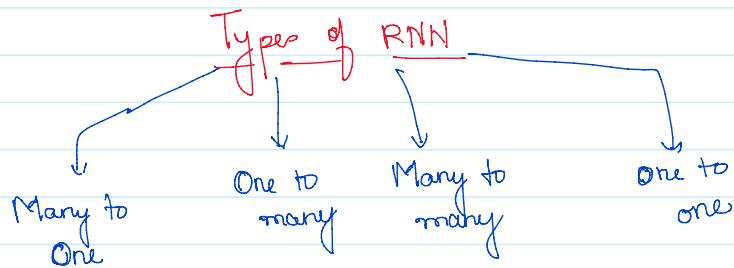
### Simplified Representation



$t = \text{time step}$

# Types of RNN

11 October 2024 11:33



## Many to One

input → sequence

output → non-sequence → int/numbers → scalars (1D)

example → sentiment analysis

$$\boxed{\text{E}} \rightarrow (1D)$$

→ review prediction

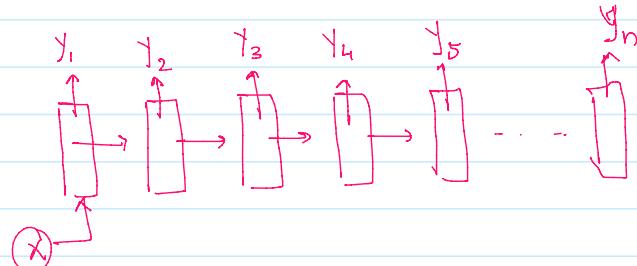
$$\boxed{\text{E}} \rightarrow \begin{matrix} 1 & \longrightarrow & 5 \\ \text{x} & \text{x} & \text{x} & \text{x} & \text{x} \end{matrix}$$

multiclass  
classification

## One to many

input → non-sequential → images, numbers

output → sequence → words etc.



ex. image captioning  
music generation

## Many to many

Input length = output length  
same length

Input → seq

Output → seq

variable length

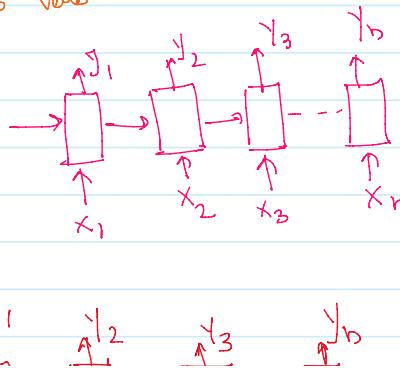
seq2seq model

Ex: machine translation

POS tagging

my name is Sakalya  
pronoun noun verb noun

NER

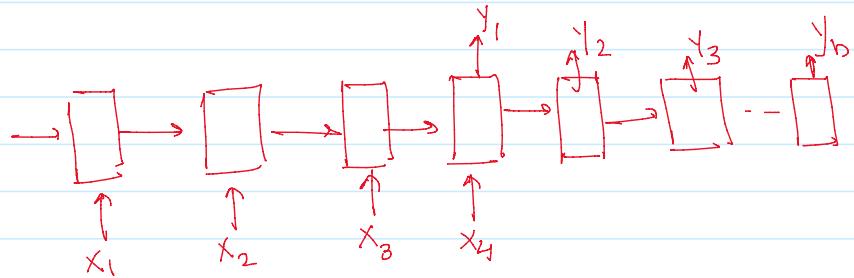


Ex: machine translation

Hi my name is Sakalya  $\rightarrow$  5

$\downarrow$   
H E I M Y S A K A L Y A  $\xrightarrow{\text{word embedding}}$   $\begin{matrix} \text{H} \\ \text{E} \\ \text{I} \\ \text{M} \\ \text{Y} \\ \text{S} \\ \text{A} \\ \text{K} \\ \text{A} \\ \text{L} \\ \text{Y} \end{matrix}$   $\rightarrow$  4

Length cannot be  
determined



first study input and then get output  
also called encoder-decoder arch

One to one (Not typically a RNN because no sequence)

image classification



simple  
neural network

no feedback

# Backpropagation in RNN

11 October 2024 11:50

It is called Backpropagation through time Step

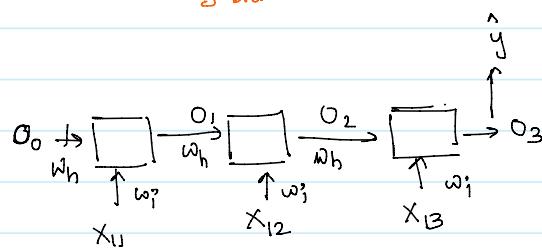
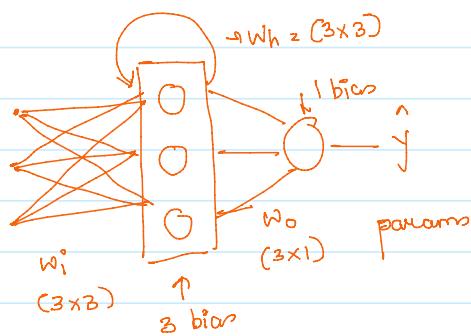
BPTT

Many to One RNN

text  $\rightarrow$  I/O

		X	y
cat mat rat	1	$x_1 [1 \ 0 \ 0] [0 \ 1 \ 0] [0 \ 0 \ 1]$	1
rat mat mat	1	$x_2 [0 \ 0 \ 1] [0 \ 1 \ 0] [0 \ 1 \ 0]$	1
mat mat cat	0	$x_3 [0 \ 1 \ 0] [0 \ 1 \ 0] [1 \ 0 \ 0]$	0

$$v_0 \text{ cab} = \begin{matrix} \text{eat mat rat} \\ [0 \ 0 \ 1] \\ [1 \ 0 \ 0] \end{matrix}$$



$$o_1 = f(x_{11}w_1 + o_0w_h)$$

$$o_2 = f(x_{12}w_2 + o_1w_h)$$

$$o_3 = f(x_{13}w_3 + o_2w_h)$$

$$\hat{y} = \sigma(o_3 \cdot w_o)$$

$$L = -y_i \log(\hat{y}_i) - (1-y_i) \log(1-\hat{y}_i)$$

minimize  
Gradient descent

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

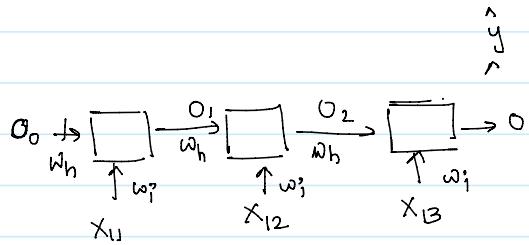
$$w_o = w_o - \eta \frac{\partial L}{\partial w_o}$$

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

$$\omega_h = \omega_i - \frac{\partial L}{\partial \omega_h}$$

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

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



depends depends  
 $\frac{\partial L}{\partial \omega_o} \xrightarrow{\hat{y}} \xrightarrow{o_3} \omega_o$

$$\boxed{\frac{\partial L}{\partial \omega_o} = \frac{\partial L}{\partial \hat{y}} \times \frac{\partial \hat{y}}{\partial \omega_o}}$$

$$o_1 = f(x_{11}\omega_1 + \omega_0\omega_h)$$

$$o_2 = f(x_{12}\omega_2 + \omega_0\omega_h)$$

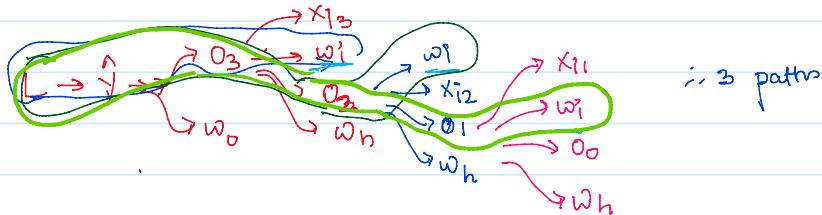
$$o_3 = f(x_{13}\omega_3 + \omega_0\omega_h)$$

$$\hat{y} = \sigma(o_3 \cdot \omega_o)$$

$$L = -y_i \log(\hat{y}_i) - (1-y_i) \log(1-\hat{y}_i)$$

$$\begin{aligned} \frac{\partial L}{\partial \hat{y}_i} &= -\frac{y_i}{\hat{y}_i} - \frac{(1-y_i)}{(1-\hat{y}_i)} \times (-1) = \frac{-y_i}{\hat{y}_i} + \frac{(1-y_i)}{(1-\hat{y}_i)} = \frac{-y_i(1-\hat{y}_i) + (1-y_i)\hat{y}_i}{\hat{y}_i(1-\hat{y}_i)} \\ &= -y_i + y_i \hat{y}_i + \hat{y}_i - \hat{y}_i \hat{y}_i \end{aligned}$$

$$\boxed{\frac{\partial L}{\partial \omega_i}}$$



$$\frac{\partial L}{\partial \omega_i} = \frac{\partial L}{\partial \hat{y}} \frac{\partial \hat{y}}{\partial o_3} \frac{\partial o_3}{\partial \omega_i} + \frac{\partial L}{\partial \hat{y}} \frac{\partial \hat{y}}{\partial o_2} \frac{\partial o_2}{\partial \omega_i} + \frac{\partial L}{\partial \hat{y}} \frac{\partial \hat{y}}{\partial o_1} \frac{\partial o_1}{\partial \omega_i}$$

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

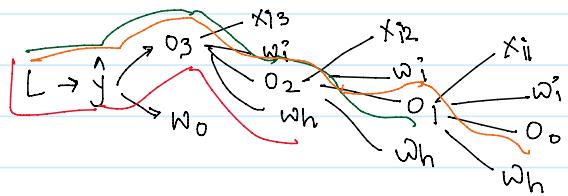
$n = \text{timestep}$

$$\delta_1 = \frac{\partial L}{\partial \hat{y}} \frac{\partial \hat{y}}{\partial o_1} \frac{\partial o_1}{\partial \omega_i} = \frac{\partial L}{\partial \hat{y}} \frac{\partial \hat{y}}{\partial o_3} \frac{\partial o_3}{\partial o_2} \frac{\partial o_2}{\partial o_1} \frac{\partial o_1}{\partial \omega_i}$$

$$\delta_2 = \frac{\partial L}{\partial \hat{y}} \frac{\partial \hat{y}}{\partial o_2} \frac{\partial o_2}{\partial \omega_i} = \frac{\partial L}{\partial \hat{y}} \frac{\partial \hat{y}}{\partial o_3} \frac{\partial o_3}{\partial o_2} \frac{\partial o_2}{\partial \omega_i} \quad \oplus \rightarrow \frac{\partial L}{\partial \omega_i}$$

$$\delta_3 = \frac{\partial L}{\partial \hat{y}} \frac{\partial \hat{y}}{\partial o_3} \frac{\partial o_3}{\partial \omega_i} = \frac{\partial L}{\partial \hat{y}} \frac{\partial \hat{y}}{\partial o_3} \frac{\partial o_3}{\partial \omega_i}$$

$$\frac{\partial L}{\partial \bar{w}_h}$$



$$\frac{\partial L}{\partial \bar{w}_h} = \frac{\partial L}{\partial \hat{y}} \frac{\partial \hat{y}}{\partial o_3} \frac{\partial o_3}{\partial \bar{w}_h} + \frac{\partial L}{\partial \hat{y}} \frac{\partial \hat{y}}{\partial o_3} \frac{\partial o_3}{\partial \bar{w}_2} \frac{\partial o_2}{\partial \bar{w}_h} + \frac{\partial L}{\partial \hat{y}} \frac{\partial \hat{y}}{\partial o_3} \frac{\partial o_3}{\partial \bar{w}_2} \frac{\partial o_2}{\partial \bar{w}_1} \frac{\partial o_1}{\partial \bar{w}_h}$$

$$\boxed{\frac{\partial L}{\partial \bar{w}_h} = \sum_{j=1}^n \frac{\partial L}{\partial \hat{y}} \frac{\partial \hat{y}}{\partial o_j} \frac{\partial o_j}{\partial \bar{w}_h}}$$

$n = \text{timesteps}$

# Problems with RNN

15 October 2024 21:07

2 major problems: *lead to need for LSTMs*

- Problem of Long Term dependencies
  - Unstable Training
- } caused by Unstable gradients

## ① Long term dependency problem

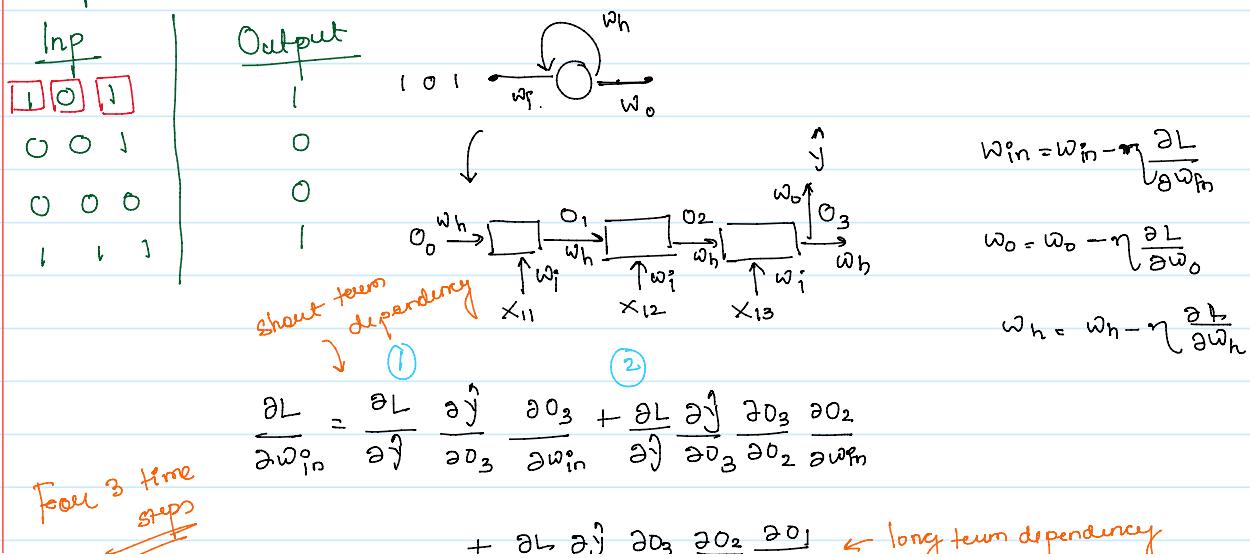
If time steps are very large, it forgets previous data.

Next word prediction  
Maurathi is spoken in Maharashtra *→ short term dependency*

Maharashtra is a beautiful place. I went there last year, but I could not enjoy because I do not understand Maurathi. *→ long term dependency → RNN fails here*

*many time steps*  
causes due to vanishing gradients

### Example



For 3 time steps

$$\frac{\partial L}{\partial w_{in}} = \frac{\partial L}{\partial o_3} \frac{\partial o_3}{\partial o_2} \frac{\partial o_2}{\partial o_1} \frac{\partial o_1}{\partial w_{in}} \quad \leftarrow \text{long term dependency}$$

For 100 time steps

$$\frac{\partial L}{\partial w_{in}} = \frac{\partial L}{\partial o_1} \frac{\partial o_1}{\partial o_2} \frac{\partial o_2}{\partial o_3} \dots \frac{\partial o_{100}}{\partial o_{99}} \dots \frac{\partial o_1}{\partial w_{in}}$$

The long term contribution reduces  
compared

$$\frac{\partial L}{\partial w_{in}} = \frac{\partial L}{\partial o_{100}} \prod_{t=2}^{100} \left( \frac{\partial o_t}{\partial o_{t-1}} \right) \frac{\partial o_1}{\partial w_{in}}$$

$$o_t = \tanh(x_t w_{in} + o_{t-1} w_h)$$

$$\frac{\partial o_t}{\partial o_{t-1}} = \tanh' \underbrace{(x_t w_{in} + o_{t-1} w_h)}_{[0,1]} \cdot w_h$$

$$= \frac{\partial L}{\partial o_{100}} \prod_{t=2}^{100} \left( \tanh' \underbrace{(x_t w_{in} + o_{t-1} w_h)}_{[0,1]} \cdot w_h \right) \frac{\partial o_1}{\partial w_{in}}$$

when we multiply 100 or more numbers in the range 0-1  
it is close to 0.

$\therefore$  Contribution by Long term dependency is almost negligible  
(Vanishing gradient problem)

### Solutions

- different activation functions (relu, leaky relu)
- better weight initialization
- skip RNNs
- LSTMs

### ② Unstable Training

$$\frac{\partial o_t}{\partial o_{t-1}} = \tanh' \underbrace{(x_t w_{in} + o_{t-1} w_h)}_{P} \cdot w_h$$

Suppose we use

Suppose we  
use  
ReLU

→ gradient increases → multiplication  
increases gradient → exploding gradient

### Solution

- gradient clipping - don't let gradient to increase beyond a value
- controlled learning rate
- LSTMs