

Recurrent Neural Nets (RNN)

L-1- Introduction Why RNN?

MLP $x \rightarrow \square \rightarrow y$
(vector)

CNN $x \rightarrow \square \rightarrow y$
(Image)

RNN $x \rightarrow \square \rightarrow y$
(Sequence)

$x \rightarrow$ Sequence

eg. ① Sequence of words

eg Amazon fine food reviews. +ve/-ve

NLP

This phone is very fast \rightarrow sequence of words.
fast is very this phone \rightarrow this makes no sense

All the techniques seems till now we have discarded the sequence information.
 \hookrightarrow {BOW, TF-idf, }
 \quad Avg W2V

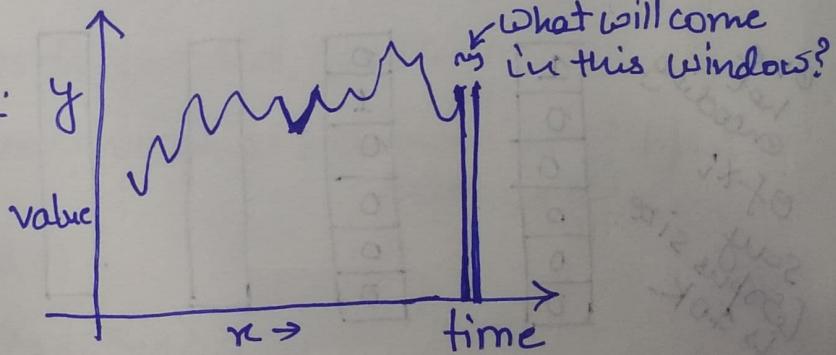
② Time series Data.

\rightarrow window the data

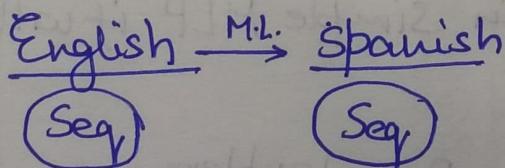
\rightarrow "freq" domain

\rightarrow Stock market/Car cab

prediction.



③ Machine translation



④ Speech Recognition

Sequence in Audio \rightarrow English text
(sequence)

Image \rightarrow Captioning $\xrightarrow{\text{English text}}$ $\xrightarrow{\text{text search}}$
 $\xrightarrow{\text{like google search.}}$

this is useful to google image search

Core Idea: When output depends on the Seq. of input rather than set of inputs.



new type of NN

→ [retains & leverage the step info?]

e.g. → Amazon food reviews

Sentence → $0/1$

$$x_i \quad y_i$$

→ This phone is very fast

$$x_{i1}$$

$$x_{i2}$$

$$x_{i3}$$

$$x_{i4}$$

$$x_{i5}$$

Word 1

Word 2

Word 3

Word 4

Word 5

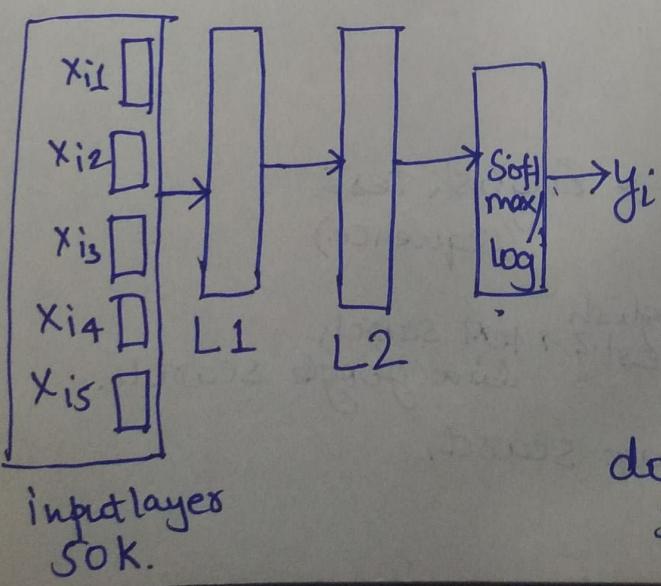
One hot encoding
of x_i
say size
is 10K.

1
0
0
0
0

0
1
0
0
0

Now for 1st point I have input of 5x10K.

So why cannot we train simple MLP if will also have sequence info?



But problem →

$x_i \rightarrow$ This is very good phone

$x_2 \rightarrow$ This phone is good.

① Varying words size.
(All sentence of different length)

Now you can say we can do zero padding for shorter Sentence.

Now if test sentence can be even bigger length than what will be our input dimension.

So, ~~if~~ In short you will be having unnecessarily lot of zeros in input data & v.v. large input dimension. and Billion of weights to train.

L-2. RNN (Recurrent Neural Network)

↓
repeating

(e.g.) Amazon food reviews \rightarrow binary classification {0,1}

$D\{x_i, y_i\}$
 ↓
 sentence \rightarrow binary
 (seq. of words)

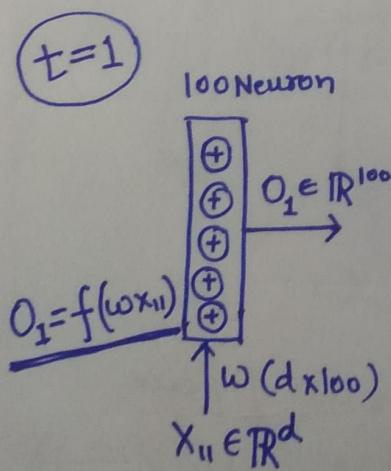
$\{x_{11}, x_{12}, x_{13}, x_{14} \rightarrow x_1 \rightarrow y_i\}$
 $x_{21}, x_{22}, x_{23} \rightarrow x_2 \rightarrow y_i$

words: one-hot encoding

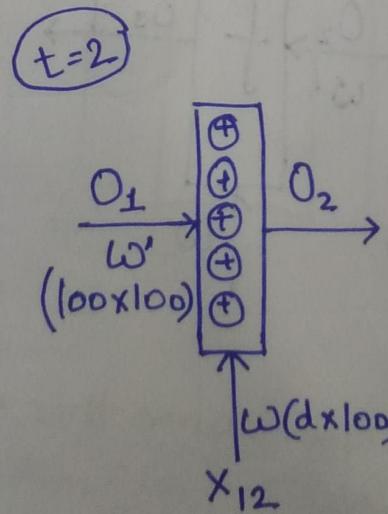
{size of vocabular: dimension of input corpora. \mathbb{R}^d }

Task: seq. of words $\rightarrow \{0,1\}$

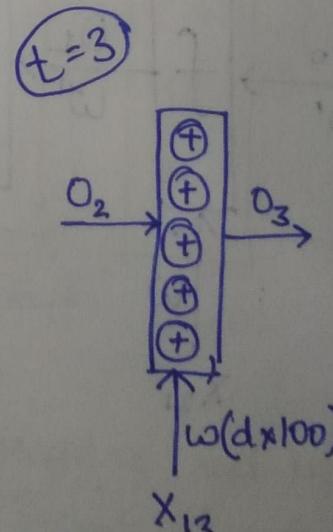
$$x_1 \rightarrow \langle x_{11}, x_{12}, x_{13}, x_{14}, x_{15} \rangle, y_i$$



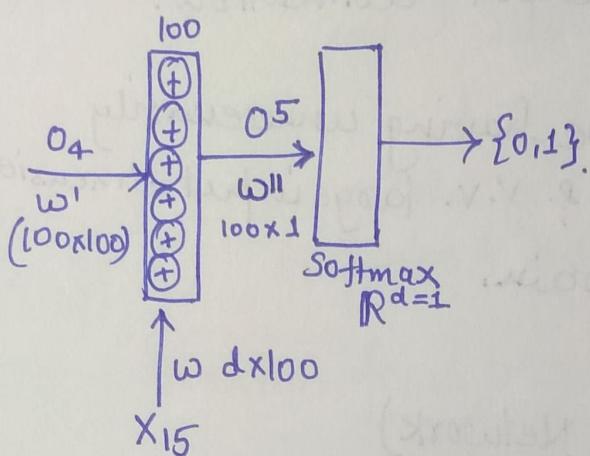
f: tanh, sigmoid,
relu
 w : matrix ($d \times 100$)



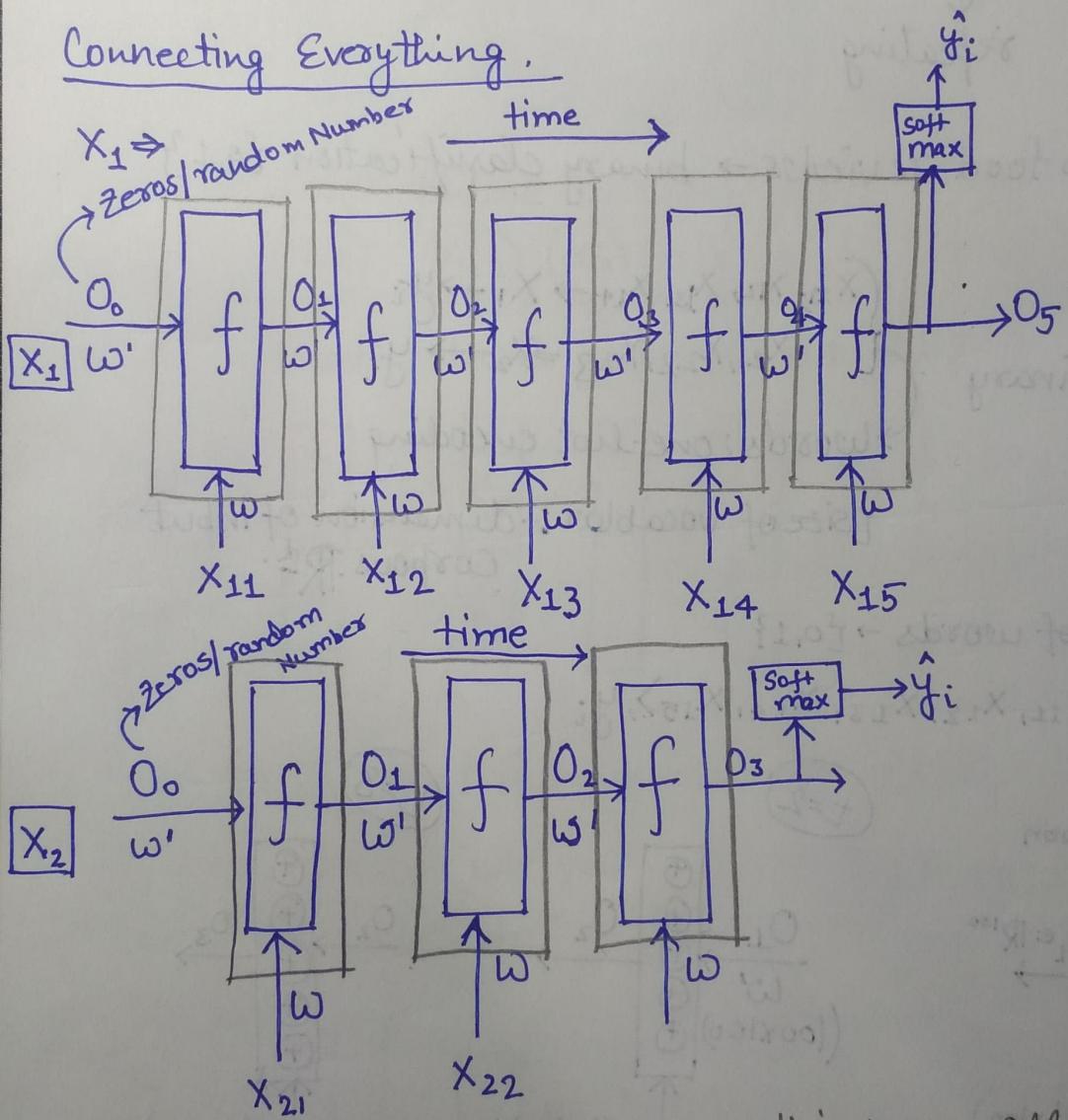
$O_2 = f(wx_{12} + w' O_1)$	$O_3 = f(wx_{13} + w' O_2)$
-----------------------------	-----------------------------



$t = 5$



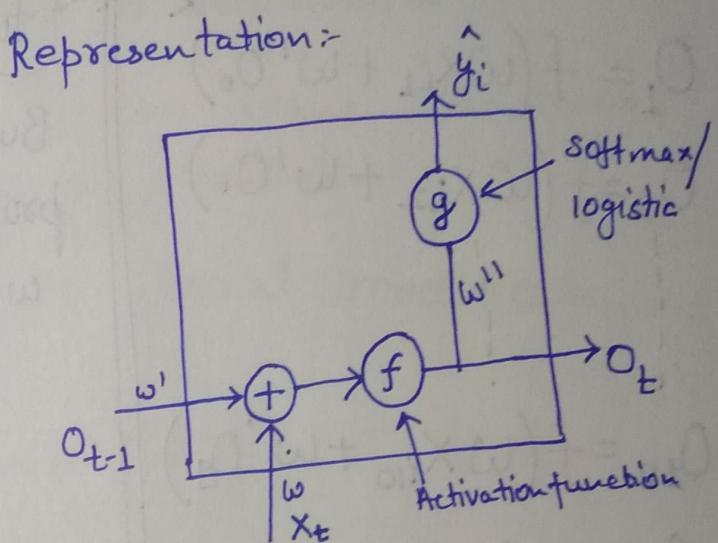
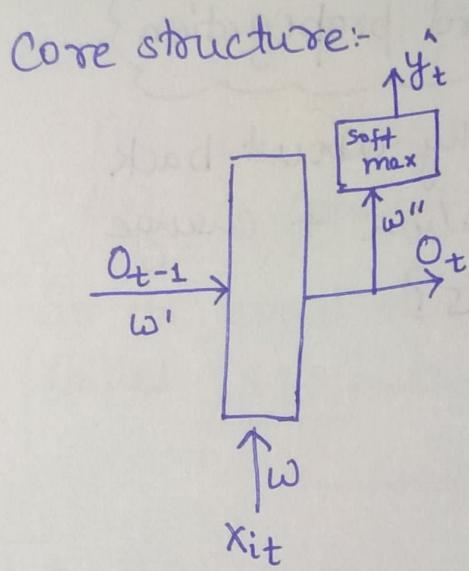
Connecting Everything.



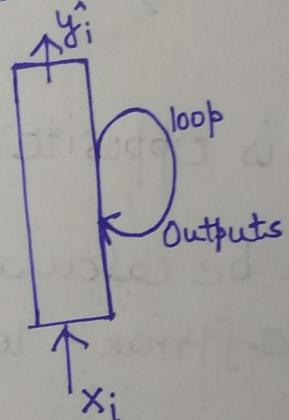
pencil box represent the same this means all \square are pointing to same layer.

$$(0.01 \times 100) + \dots + 0 = 0$$

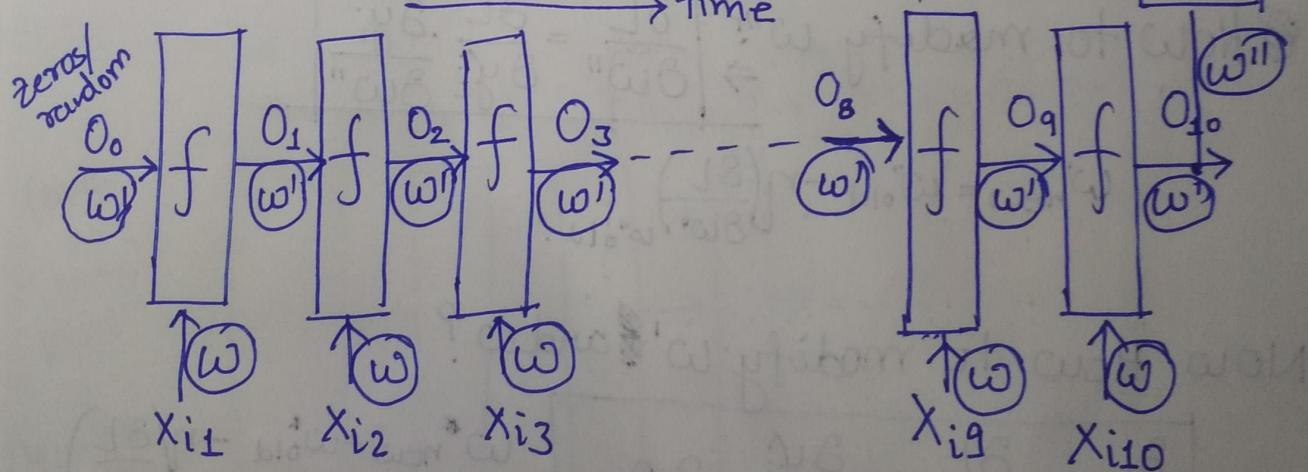
$$(0.01 \times 100) \times 100 = 0$$



Other representation:



L-3 - Training RNN: Backprop over time



$$x_i = \langle x_{i1}, x_{i2}, x_{i3}, x_{i4}, x_{i5}, x_{i6}, \dots, x_{i10} \rangle$$

this is just 1 layer repeating overtime. they are not different layers.

$$O_1 = f(\omega x_{i1} + \omega' O_0)$$

$$O_2 = f(\omega x_{i2} + \omega' O_1)$$

⋮ ⋮

$$O_{10} = f(\omega x_{i10} + \omega' O_9)$$

$$\hat{y}_i = g(\omega'' \cdot O_{10})$$

forward propagation

But why about back propagation & change weights?

Backprop is opposite direction of arrows.

① $\frac{\partial L}{\partial \hat{y}_i}$ can be calculated as loss function of softmax & log-loss is derivative.

② $\frac{\partial L}{\partial O_{10}} = \frac{\partial L}{\partial \hat{y}_i} \cdot \frac{\partial \hat{y}_i}{\partial O_{10}}$, chain rule.

③ How to modify ω'' ? $\frac{\partial L}{\partial \omega''} = \frac{\partial L}{\partial \hat{y}_i} \cdot \frac{\partial \hat{y}_i}{\partial \omega''}$

$$\omega''_{\text{new}} = \omega''_{\text{old}} - \eta \left(\frac{\partial L}{\partial \omega''} \right)_{\omega''_{\text{old}}}$$

Now how to modify ω' and ω ?

$$\frac{\partial L}{\partial \omega} = \frac{\partial L}{\partial \hat{y}_i} \cdot \frac{\partial \hat{y}_i}{\partial O_{10}} \cdot \frac{\partial O_{10}}{\partial \omega}$$

$$\omega_{\text{new}} = \omega_{\text{old}} - \eta \left(\frac{\partial L}{\partial \omega} \right)_{\omega_{\text{old}}}$$

similarly

$$\frac{\partial L}{\partial \omega'} = \frac{\partial L}{\partial \hat{y}_i} \cdot \frac{\partial \hat{y}_i}{\partial O_{10}} \cdot \frac{\partial O_{10}}{\partial \omega'}$$

$$\omega'_{\text{new}} = \omega'_{\text{old}} - \eta \left(\frac{\partial L}{\partial \omega'} \right)_{\omega'_{\text{old}}}$$

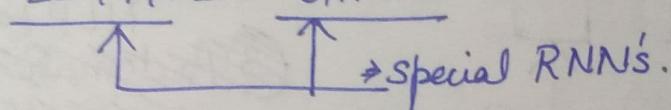
And this has to be done with each x_{ij} i.e dataset j = length of each point

Now if you backprop again for deep network

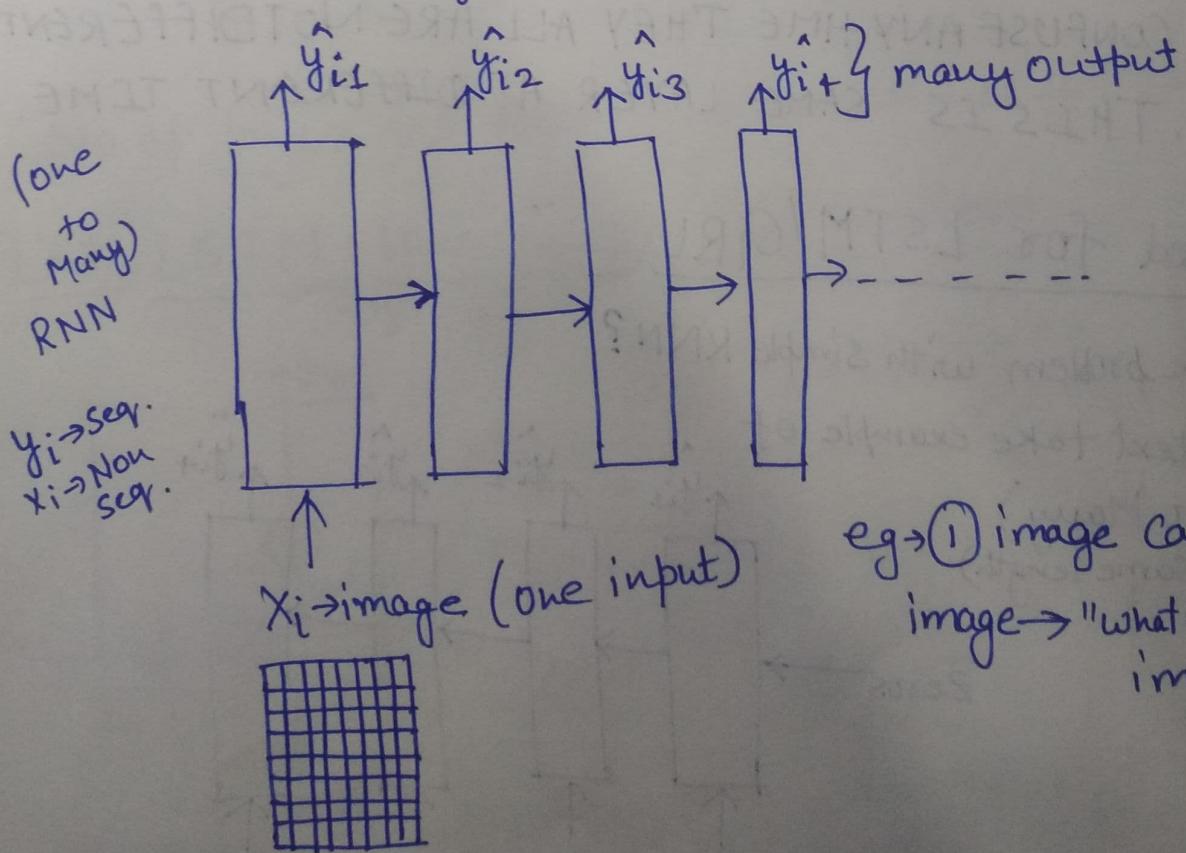
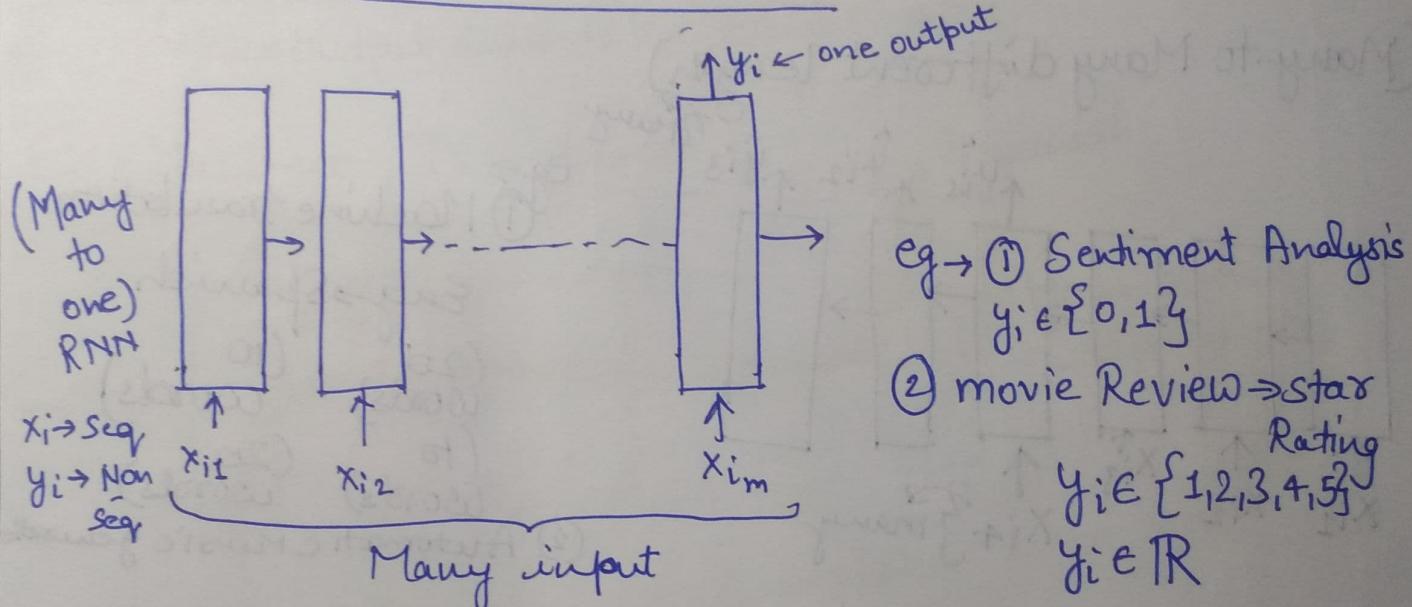
$\frac{\partial L}{\partial w}$ for 1st word in x_i can cause vanishing gradient & exploding gradient.

So to avoid this we have special modification

Called LSTM's and GRU's.

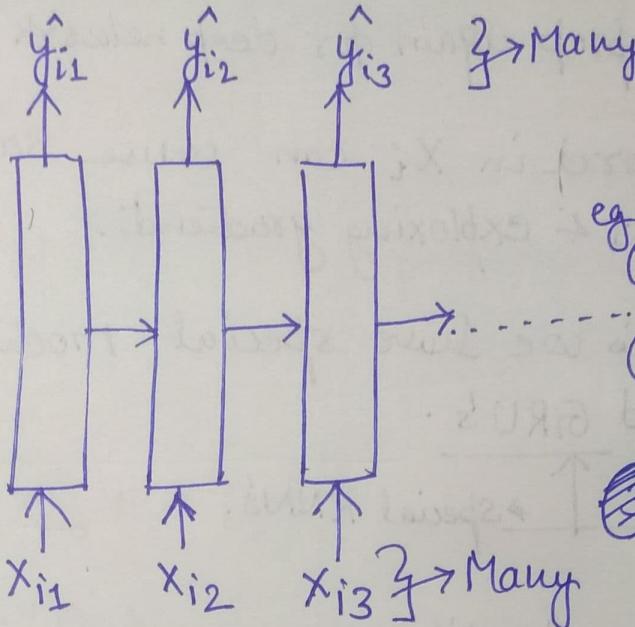


L4 - Types of RNN Architecture



(Many to Many RNN)

input \rightarrow seq
output - seq
Also called (Seq2seq model).

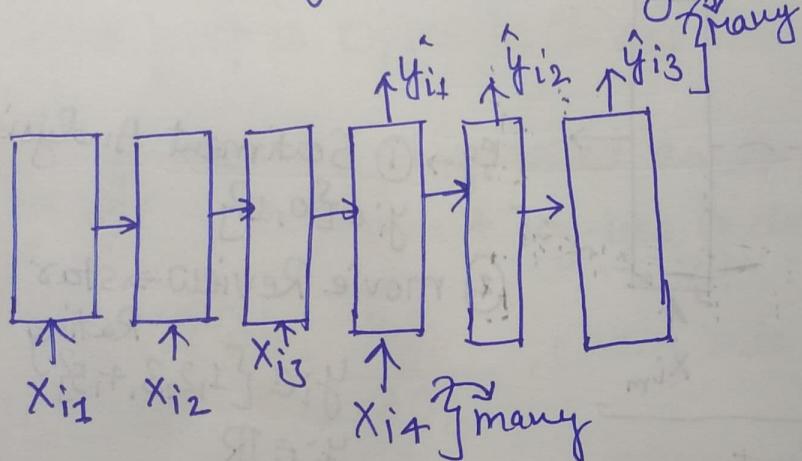


eg → ① input - output same length.

② part of speech (same length) NLP task

Automatic music generation.

(Many-to-Many different length)



eg → ① Machine translation

Eng \rightarrow Spanish

(20 words) (10 words)
(10 words) (20 words)

② Automatic music generation

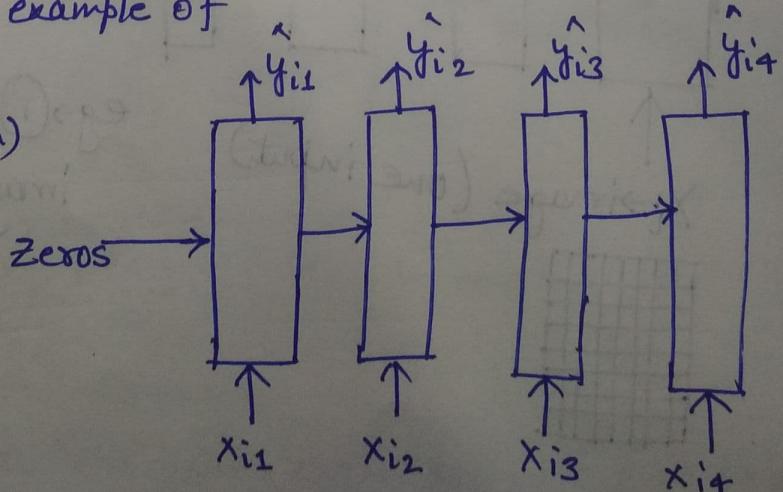
* DO NOT CONFUSE ANY TIME THEY ALL ARE NOT DIFFERENT LAYER. THIS IS SAME LAYER AT DIFFERENT TIME.

L-5- Need for LSTM/GRU

What is the problem with simple RNN?

In this context take example of

(Many to Many)
(Same length)



\hat{y}_{i4} depends on x_{i4} & o_3 more in this case.

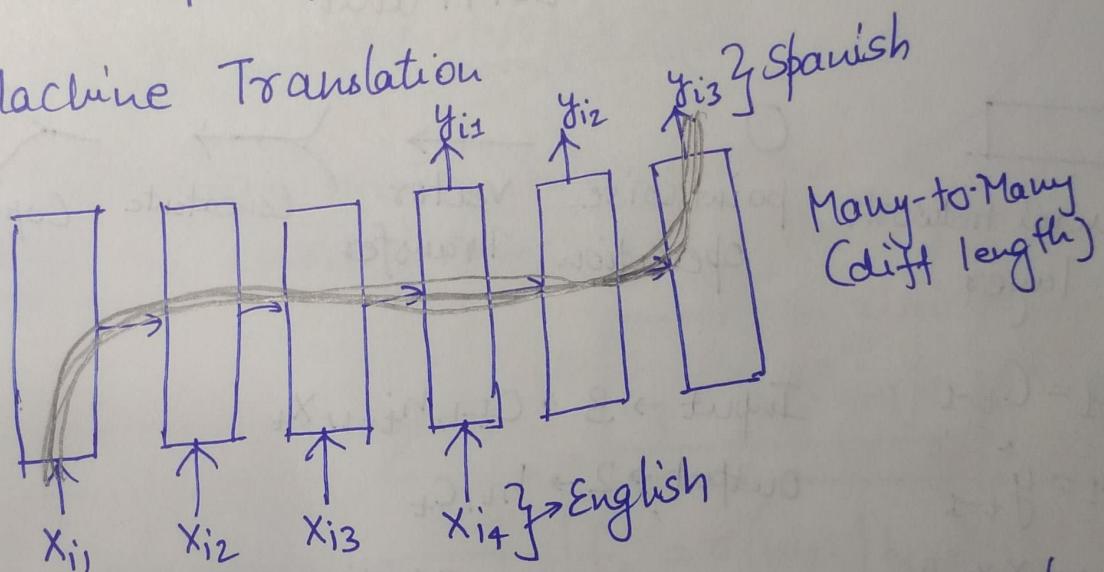
→ There can be real world applications where \hat{y}_{i4} must depend more on x_{i1} rather than x_{i4} .

Ex → part of speech and in that architecture described before \hat{y}_{i4} & x_{i1} are very far away.

So key problem "Simple RNN's cannot handle long term dependencies".

↳ later output depends on earlier input.

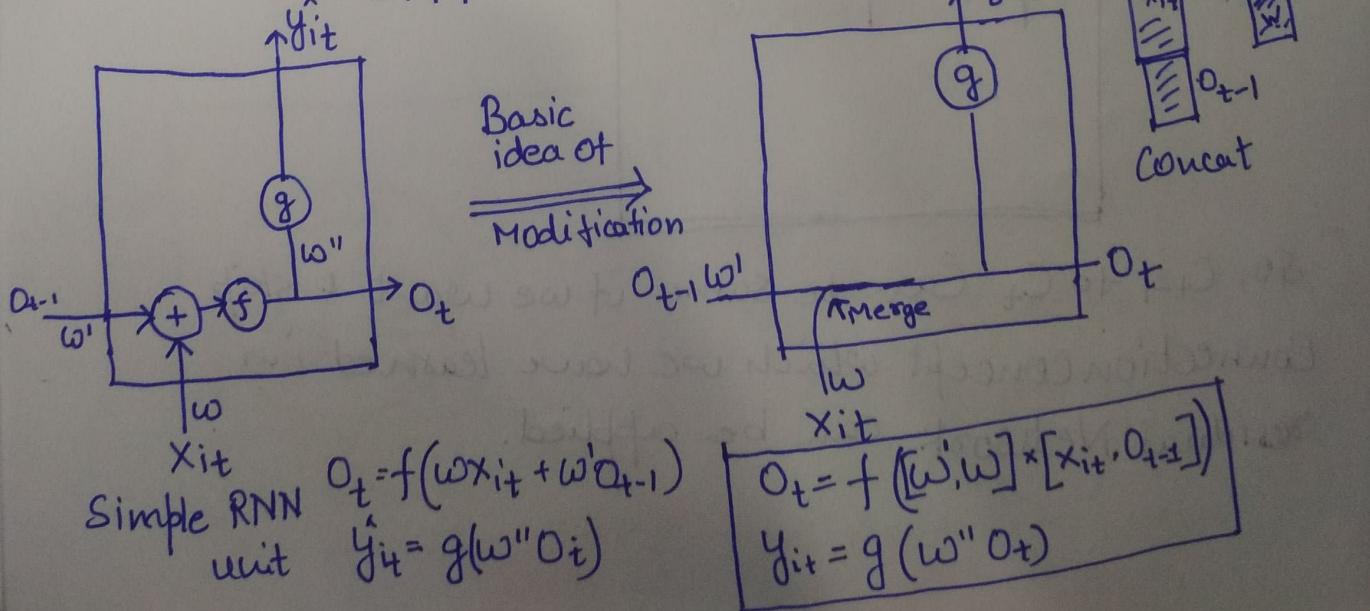
Ex → Machine Translation



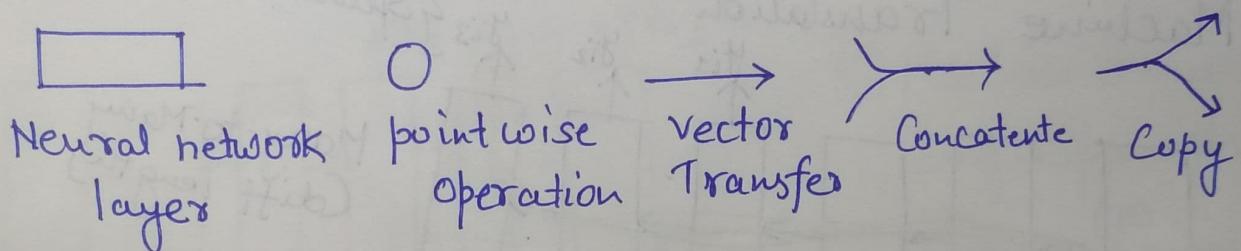
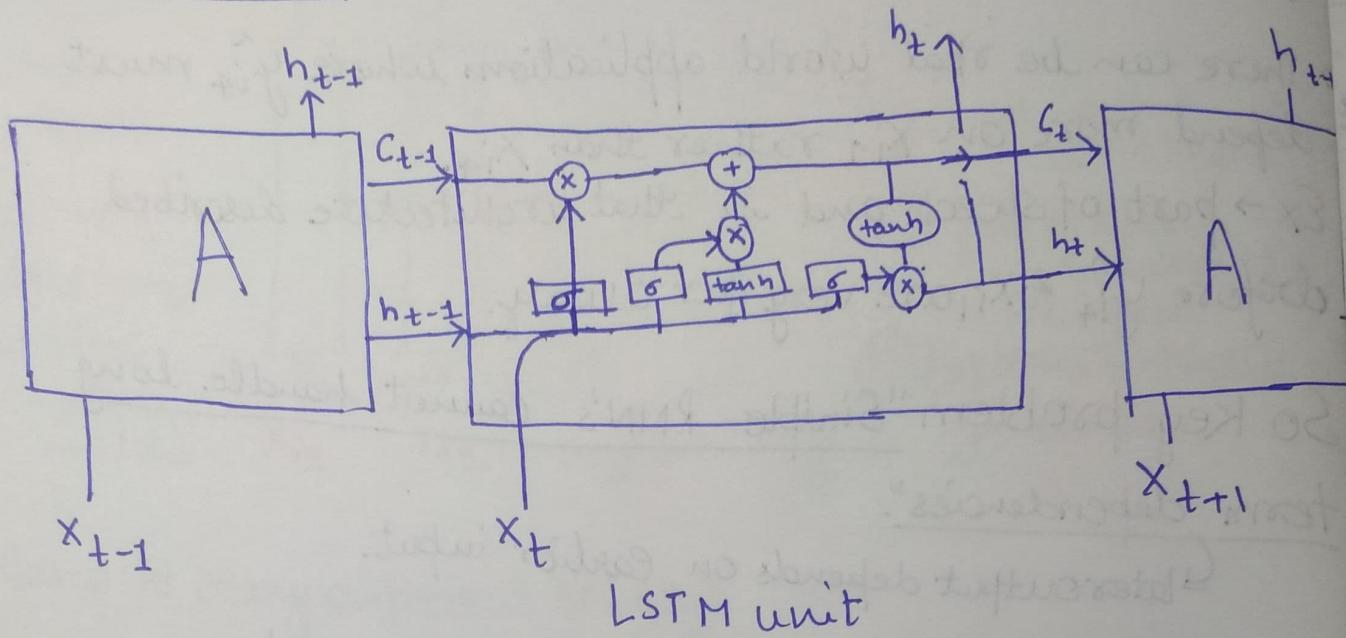
No. of layers far more and dependencies got even far.

L-6 LSTM - Long Short Term Memory

→ 1997



LSTM Unit



$$O_{t-1} = C_{t-1}$$

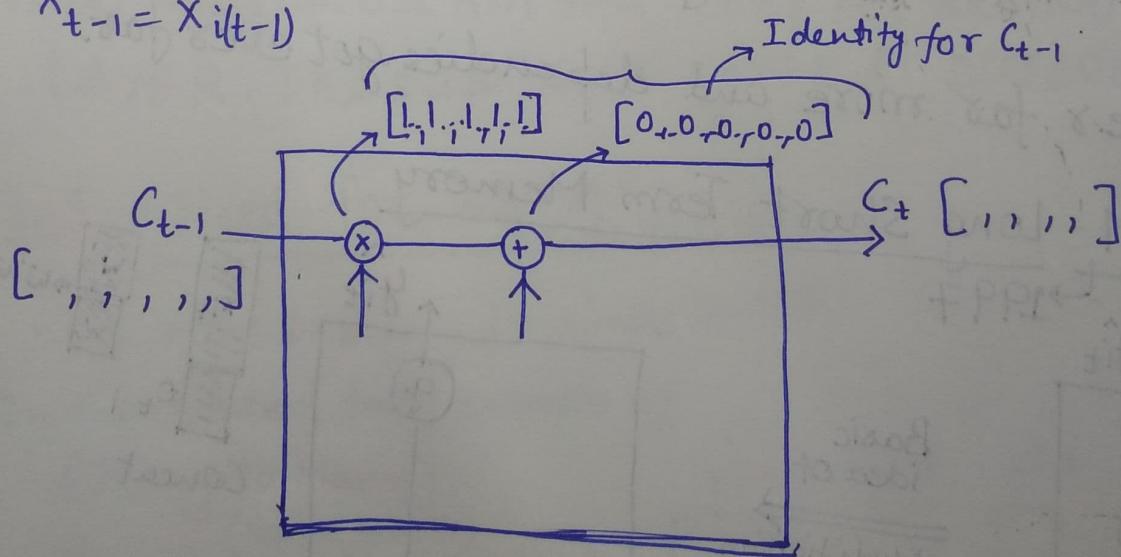
$$h_{t-1} = \hat{y}_{t-1}$$

$$x_{t-1} = x_{i(t-1)}$$

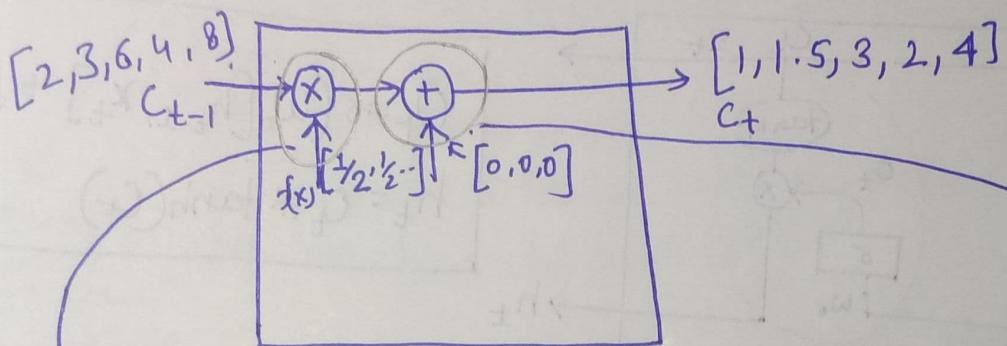
Input $\rightarrow 3 \rightarrow C_{t-1}, h_{t-1}, x_t$

Output $\rightarrow 2 \rightarrow h_t, C_t$

Identity for C_{t-1}



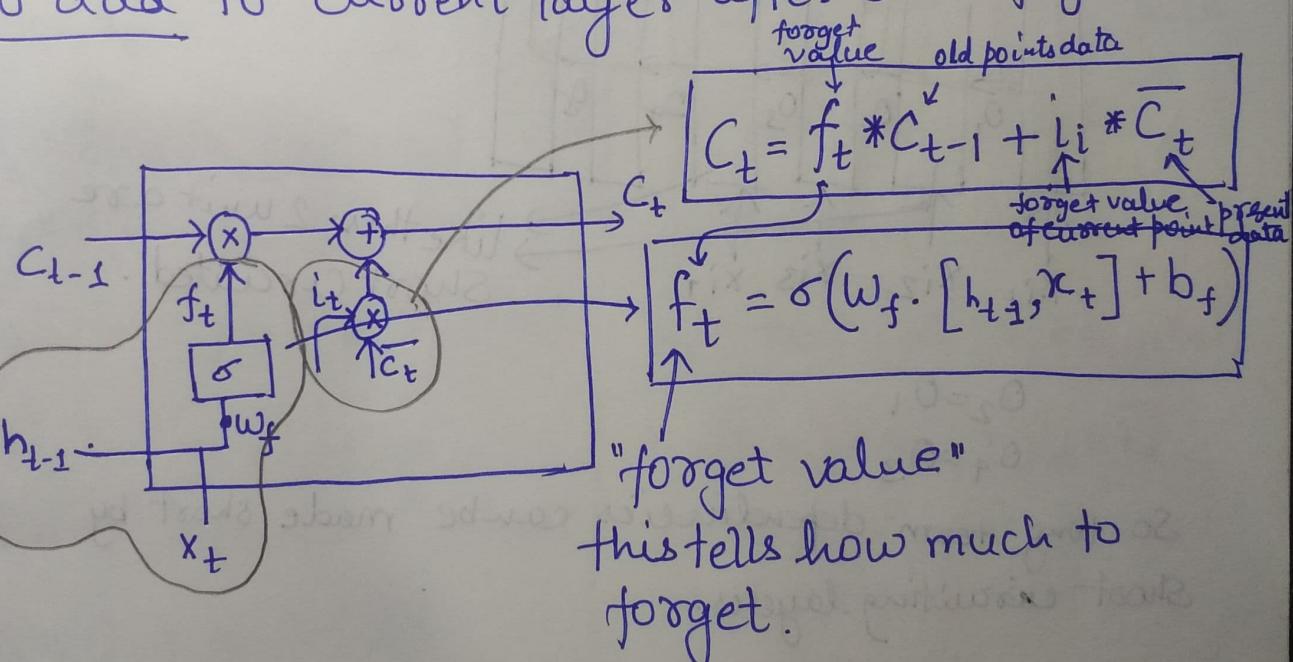
So, C_{t-1} to C_t can be same if we want to skip Connection concept which we have learned in residual Network can be applied.



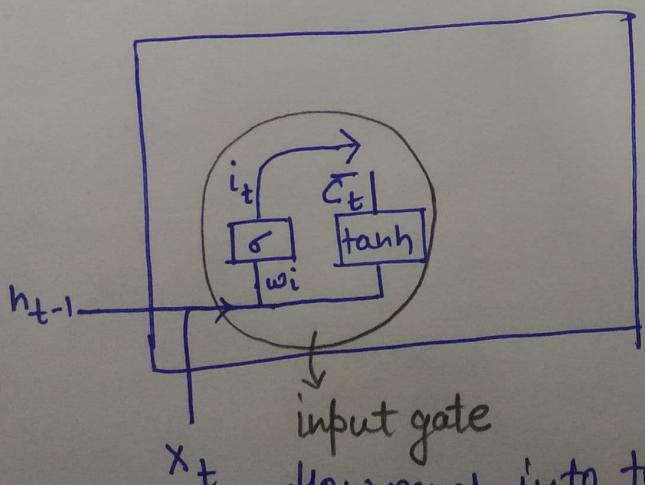
So this tells how much of previous layer's percent you want take forward for each bit.

How much to remember or forget of each bit also called forget gate.

This tells how much information to add to current layer after some forgets.



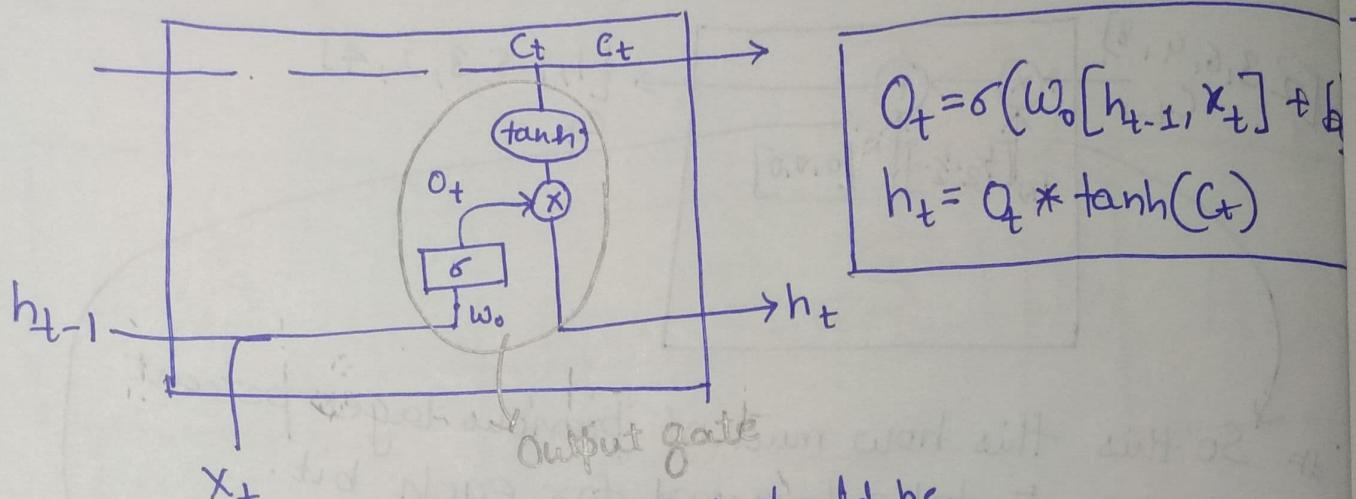
"forget value"
This tells how much to forget.



$$i_t = \sigma(w_i \cdot [h_{t-1}, x_t] + b_i)$$

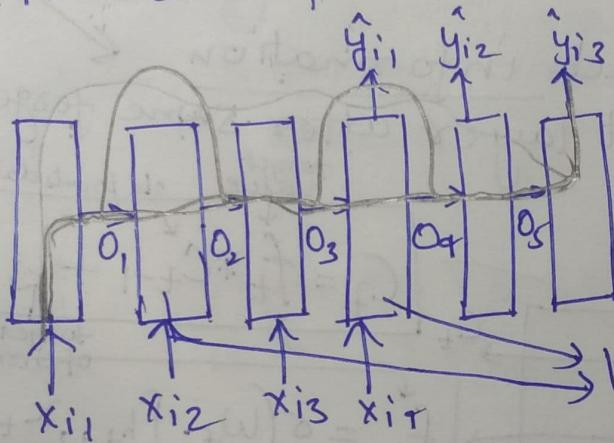
$$c_t = \tanh(w_c \cdot [h_{t-1}, x_t] + b_c)$$

How much info to take forward of present datapoint to future?



* How output should be generated based on processed input.

* LSTM can help in short circuit connections.



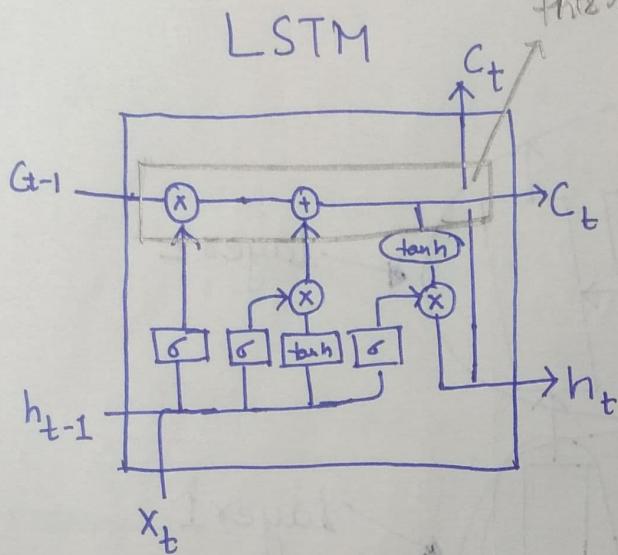
$$O_2 = O_1$$

$$O_4 = O_3$$

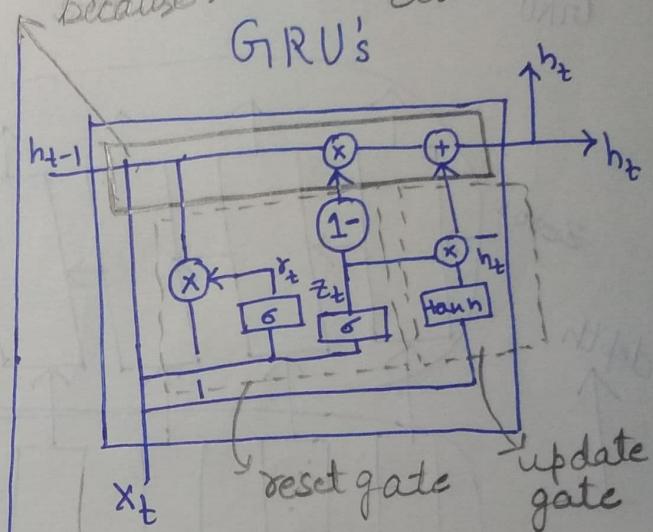
So long term dependencies can be made short by short circuiting layers.

L-7 GRU (Gated Recurrent Unit)

→ more modern than LSTM
 → 2/4 → Simplified version
 → faster to train / lesser equations.
 Only 2 gates → as powerful as LSTM
 - reset
 - update



this kind of structure remain same because it helps to short circuit.



$$\begin{aligned}
 f_t &= \sigma(W_f \cdot [h_{t-1}, x_t] + b_f) \\
 i_t &= \sigma(W_i \cdot [h_{t-1}, x_t] + b_i) \\
 \tilde{C}_t &= \tanh(W_c \cdot [h_{t-1}, x_t] + b_c) \\
 C_t &= f_t * C_{t-1} + i_t * \tilde{C}_t \\
 o_t &= \sigma(W_o \cdot [h_{t-1}, x_t] + b_o) \\
 h_t &= o_t * \tanh(C_t)
 \end{aligned}$$

$$\begin{aligned}
 Z_t &= \sigma(W_z \cdot [h_{t-1}, x_t]) \\
 r_t &= \sigma(W_r \cdot [h_{t-1}, x_t]) \\
 \tilde{h}_t &= \tanh(W_u \cdot [r_t * h_{t-1}, x_t]) \\
 h_t &= (1 - Z_t) * h_{t-1} + Z_t * \tilde{h}_t
 \end{aligned}$$

faster because of lesser equation.

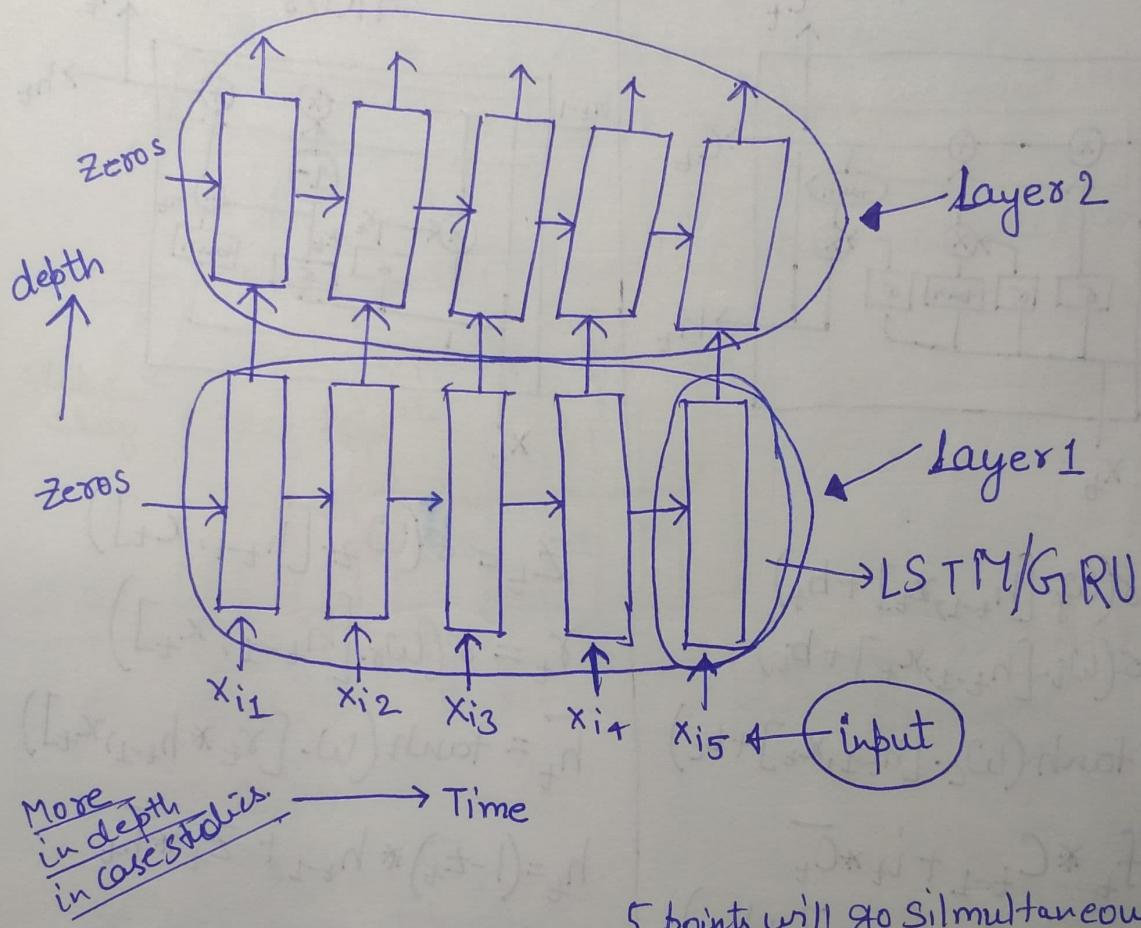
L-8- Deep RNN's

MLP → one layer → Deep MLP

CNN → Conv. Layer & Maxpool → Deep CNN

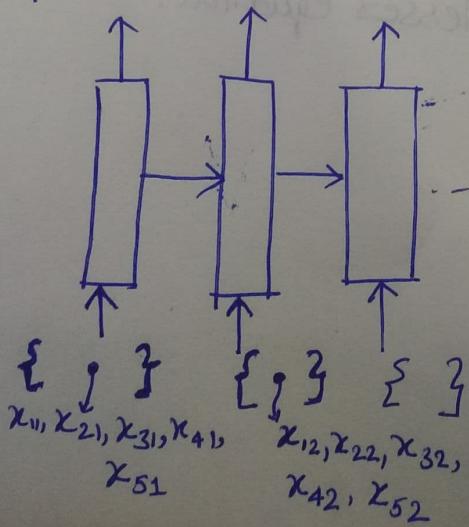
RNN → Deep RNN ... ??

LSTM /
GRU



SGD with batch size = $k(5)$

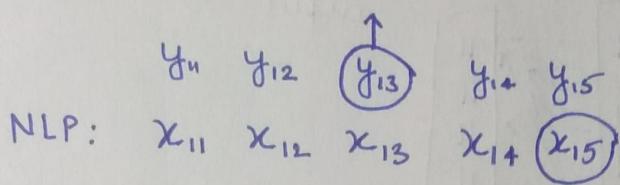
Speed up.



→ But to perform this speed up hack all the points must have same length so for that we can also use the concept of padding.

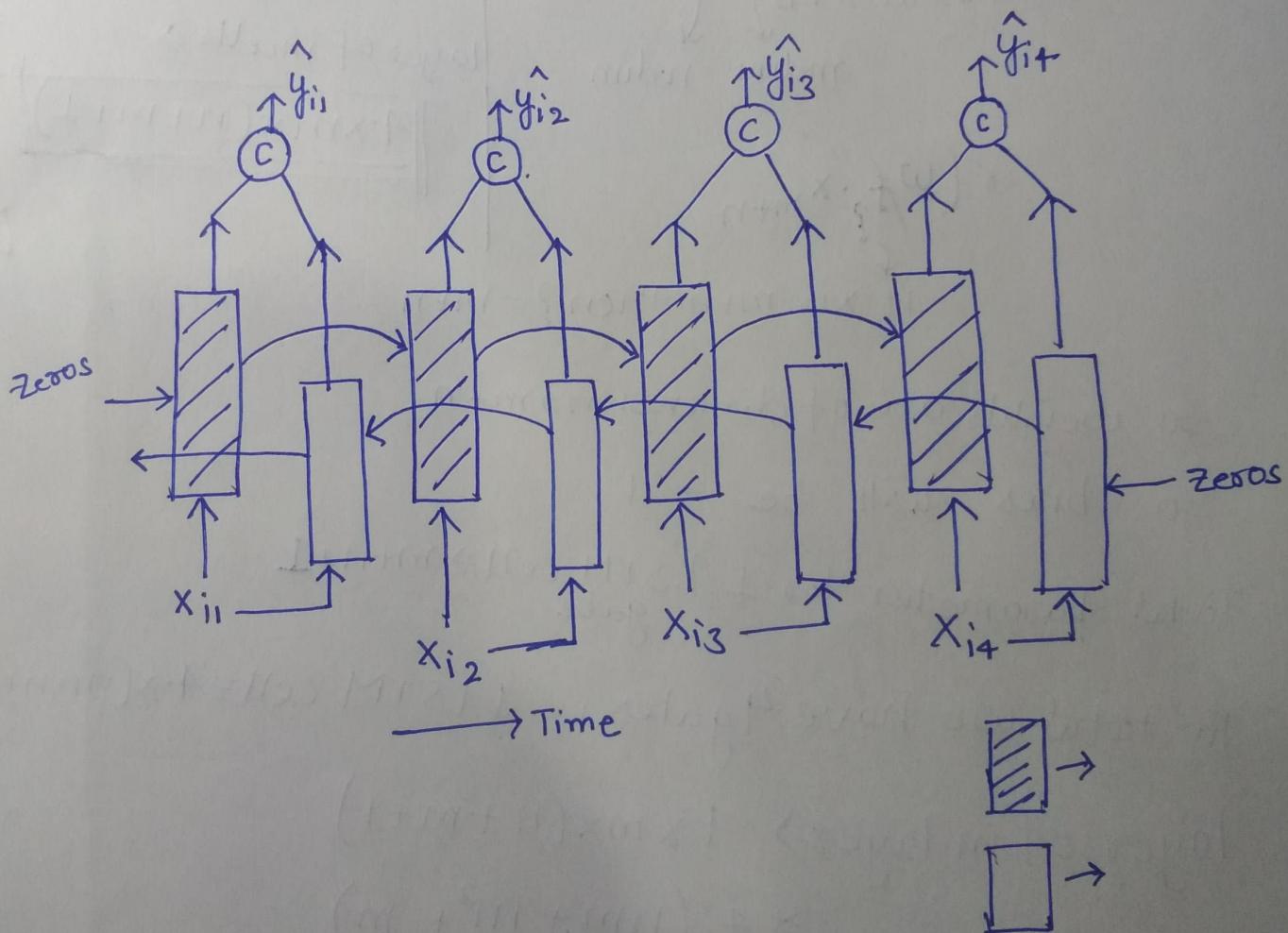
L-9

Bidirectional RNN

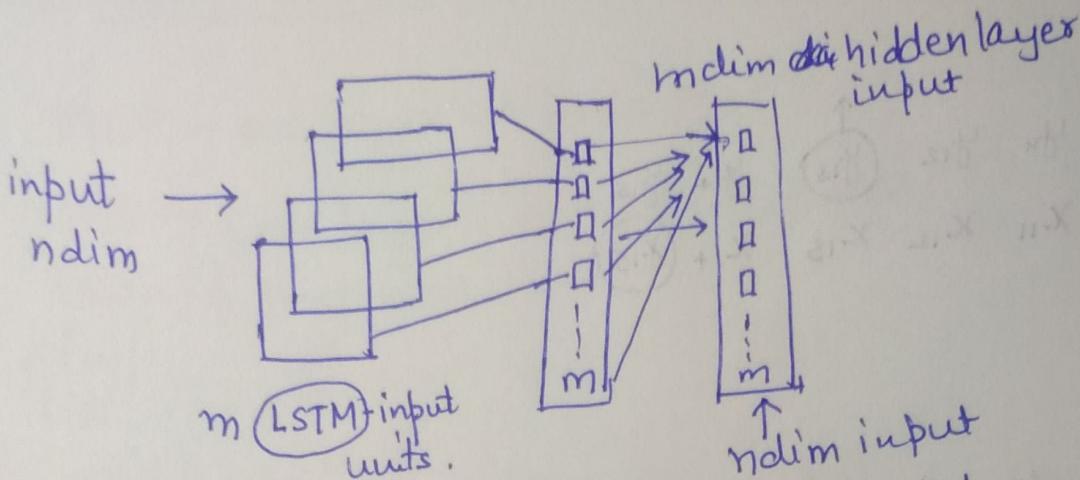


Q. What if y_{i3} is dependent on x_{i5}

\hat{y}_{it} is dependent on x_{it+k}



L-10 Numbers of parameter in LSTM



first lets see number of parameters in first layer 1 unit.

$$f_t = \sigma(W_f \cdot [h_{t-1}, x_t] + b_f)$$

look like $\rightarrow W^T \cdot x + b$

$m \cdot \text{dim}$ $n \cdot \text{dim}$

1 cell \rightarrow 1 gate $\Rightarrow n+m+1$
all 4 gates $\Rightarrow (n+m+1) \times 4$

layer of m cells \Rightarrow

$$\boxed{4 \times m \times (n+m+1)}$$

$$\sigma(W_f \cdot x \cdot m+n)$$

if x is $m+n$ then $? = m+n$

So weights are of dimension $\Rightarrow m+n$

So bias will be $\Rightarrow 1$

Total parameter for 1 LSTM cell $\Rightarrow m+n+1$ gate.

In total we have 4 gates in 1 LSTM cell $= 4 \times (m+n+1)$.

$$\begin{aligned} \text{layer of } m \text{ layers} &\Rightarrow 4 \times m \times (n+m+1) \\ &\Rightarrow 4 (nm + m^2 + m) \end{aligned}$$

if default Keras implementation bias i.e. m is ignored.