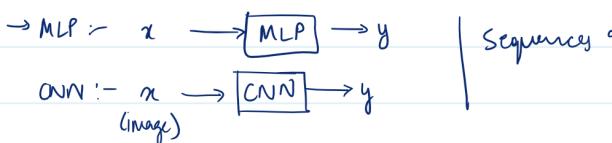


Why RNNs (5.1):-

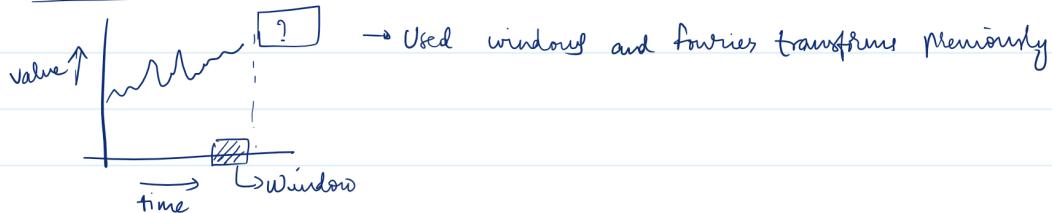


Sequence q

e.g.: This phone is very fast } Both have the same TFIDF/Bow representation. Discarded sequence info.

① Phone this is very fast } but both are not the same.

② Time Series Data :-



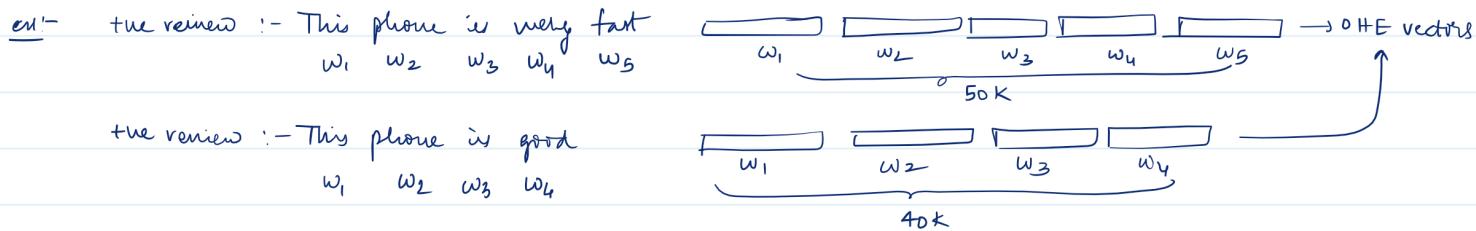
③ Machine Translation :- Language translations. The sentences hold sequential information. Using RNN we work with seq2seq models.

④ Speech Recognition :- Sequence of Audio → English sentence.

⑤ Image captioning (Google Image Search)

→ Our idea is to build a NN that retains & leverages sequential information.

→ Problem with not using RNN & using an MLP directly is inputs are of different sizes



We could do zero padding to make all inputs equal but

(i) Does not work for longer test queries

(ii) Unnecessarily increasing dimensionality

Recurrent Neural Network:- (5.2)

→ Recurrent = repeats.

→ Taking example of Amazon Fine Food Reviews

$$D = \{x_i, y_i\}$$

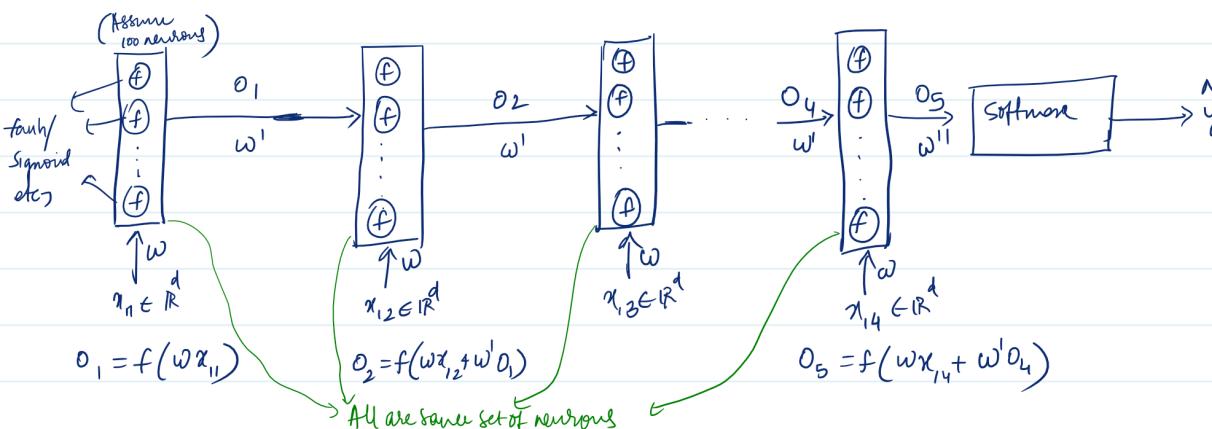
$$x_i \rightarrow x_{i1}, x_{i2}, \dots, x_{im_i} \quad \text{where } m_i = \text{len of } x_i$$

Assume size of vocabulary = d

$x_{ij} \in \mathbb{R}^d$ → Binary vectors (one hot encoding); $x_i = \langle x_{i1}, x_{i2}, x_{i3}, x_{i4}, x_{i5} \rangle, y_i$

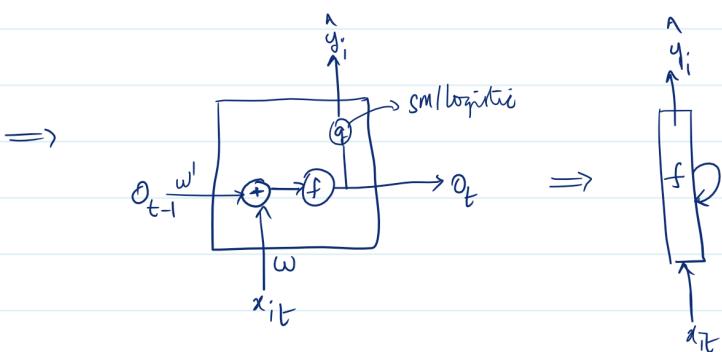
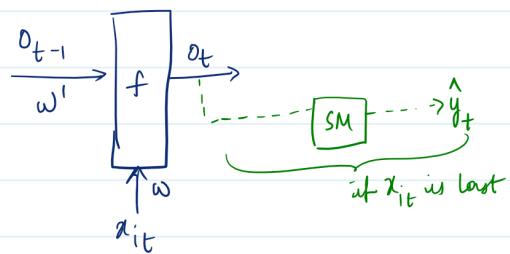
At t=1 (iteration 1) → At t=2 → At t=3 → ... At t=5

w^t is same at all stages



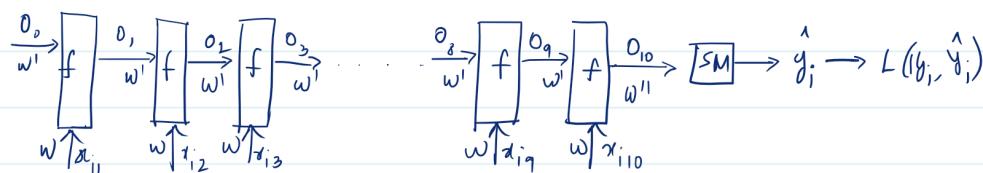
→ For the next sentence the same w are reused.

→ General way of representing RNN



Training RNNs: Backpropagation (5.3) :-

→ Backprop for RNN = backprop over time



$$\begin{aligned} O_1 &= f(wx_{i_1} + w^T O_0) \\ O_2 &= f(wx_{i_2} + w^T O_1) \\ &\vdots \\ O_{10} &= f(wx_{i_{10}} + w^T O_9) \\ \hat{y}_i &= g(w^T O_{10}) \end{aligned}$$

time →

→ Forward propagation = direction of arrows (time)

Back prop : opposite direction.

$$\frac{\partial L}{\partial \hat{y}_i} \quad \left| \quad \frac{\partial L}{\partial \hat{y}_{10}} = \frac{\partial L}{\partial \hat{y}_i} \cdot \frac{\partial \hat{y}_i}{\partial O_{10}} \quad \right| \quad \left| \quad \frac{\partial L}{\partial w^n} = \frac{\partial L}{\partial \hat{y}_i} \cdot \frac{\partial \hat{y}_i}{\partial w^n} \quad \right| \quad \left| \quad \frac{\partial L}{\partial w} = \frac{\partial L}{\partial \hat{y}_i} \cdot \frac{\partial \hat{y}_i}{\partial O_{10}} \cdot \frac{\partial O_{10}}{\partial w} \right.$$

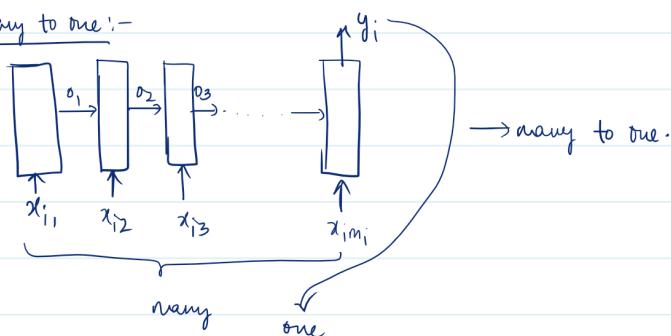
$$\frac{\partial L}{\partial w} = \left(\frac{\partial L}{\partial \hat{y}_i} \cdot \frac{\partial \hat{y}_i}{\partial O_{10}} \cdot \frac{\partial O_{10}}{\partial w} \right) + \left(\frac{\partial L}{\partial \hat{y}_{10}} \cdot \frac{\partial \hat{y}_{10}}{\partial O_{10}} \cdot \frac{\partial O_{10}}{\partial w} \right) + \left(\frac{\partial L}{\partial \hat{y}_i} \cdot \frac{\partial \hat{y}_i}{\partial O_{10}} \cdot \frac{\partial O_{10}}{\partial O_9} \cdot \frac{\partial O_9}{\partial O_8} \cdot \frac{\partial O_8}{\partial w} \right) + \dots \quad ()$$

Lots of multiplications involved

→ Since we are performing backpropagation over time & $\frac{\partial L}{\partial w}$ will have lots of multiplications based on number of words, we run into vanishing gradient & exploding gradient problem. This is NOT because of many layers but because of sequential (sentence) length. LSTMs & GRUs solve this problem.

Types of RNN Architectures (5.4) :-

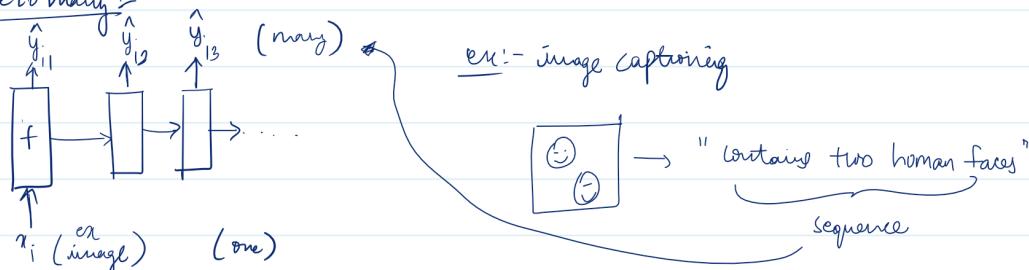
① Many to one:-



Example:- ① sentiment classification (food review)

② Movie review → star rating (multiclass clasifn)

② One to many:-



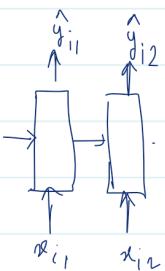
② Many to Many :-

input \rightarrow seq

output \rightarrow seq

type 1 :- input & output same length

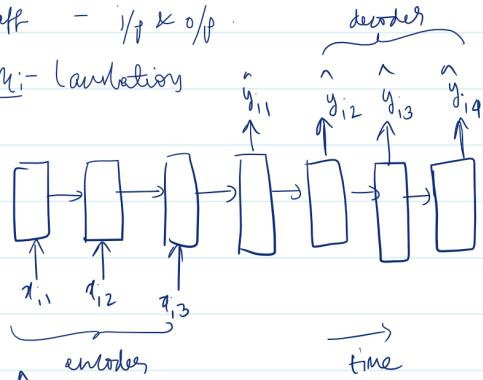
e.g. :- parts of speech of sentence.



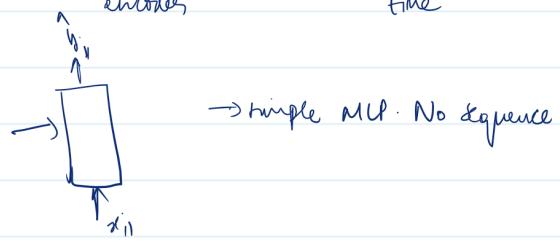
seq2seq models

type 2 :- diff - i/p & o/p

e.g. :- Translation

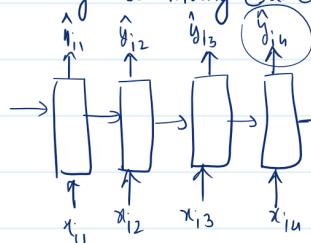


③ one to one :-



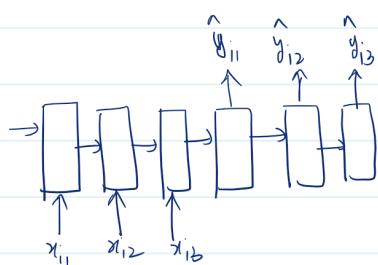
Need for LSTM/GRUs (5.5) :-

\rightarrow Take many - to - many (same length for example)



$\hat{y}_{i,4}$ in this model depends mostly on $x_{i,4}$. Its value does not depend on $x_{i,1}$ because the gradient vanishes by the time it reaches the calculation of $\hat{y}_{i,4}$.
but in real world there may be cases where $\hat{y}_{i,4}$ is more dependent on $x_{i,1}$ than $x_{i,4}$.
There is no long term dependency in this case. (in simple RNN)

But its output depends lot more on earlier inputs.

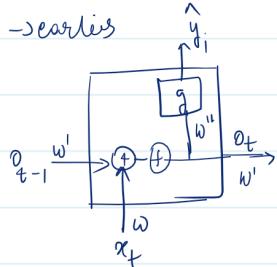


$\hat{y}_{i,3}$ would depend on $x_{i,1}$ in case of translation. Simple RNN can't handle this.

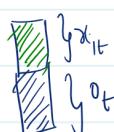
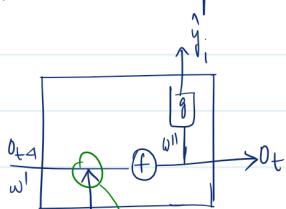
LSTM (5.6) :-

\rightarrow It can take care of both long term & short term dependencies.

\rightarrow earlier



Now



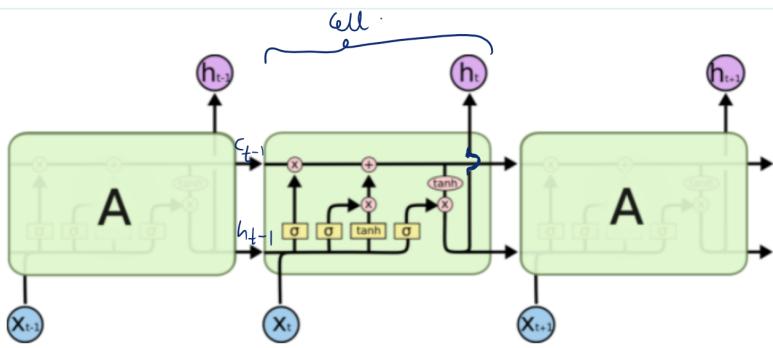
$$o_t = f([w^i, w^o][x_t, o_{t-1}])$$

$$\hat{y}_{i,t} = g(w^o o_t)$$

concatenation operation. Not addition like before -

$$o_t = f(w x_t + w' o_{t-1})$$

$$\hat{y}_{i,t} = g(w' o_t)$$



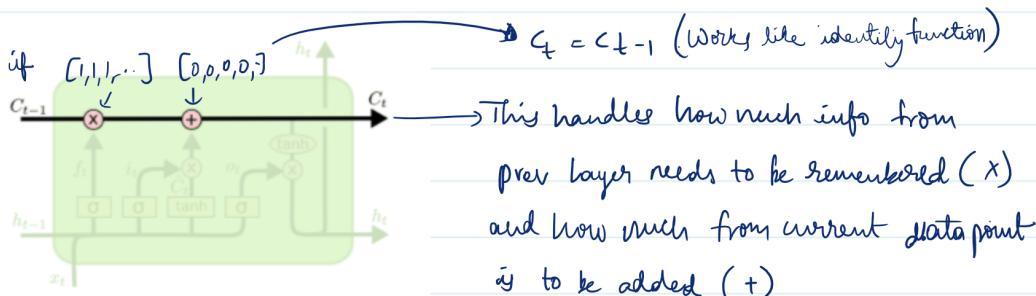
The repeating module in an LSTM contains four interacting layers.

3 outputs.

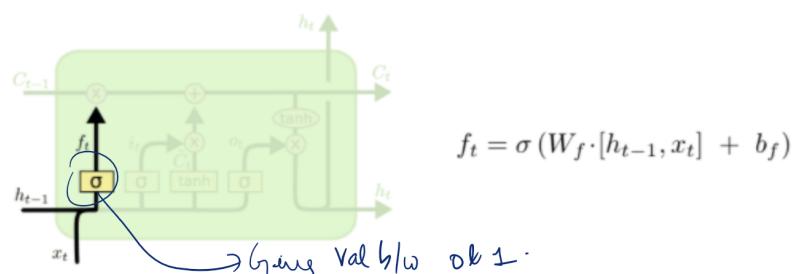
$$h_t = y_t$$

$$x_t = x_t$$

$$c_t = o_t$$

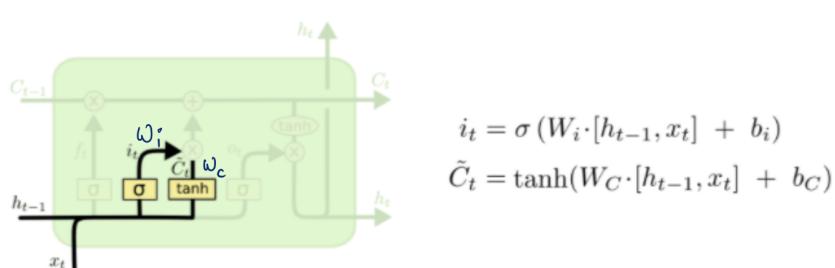


Forget Gate Layer :-



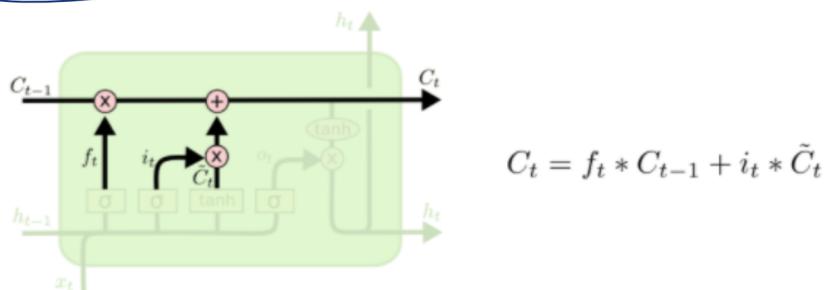
→ handles how much information we are going to throw away from the current layer.

Input Gate Layer :-

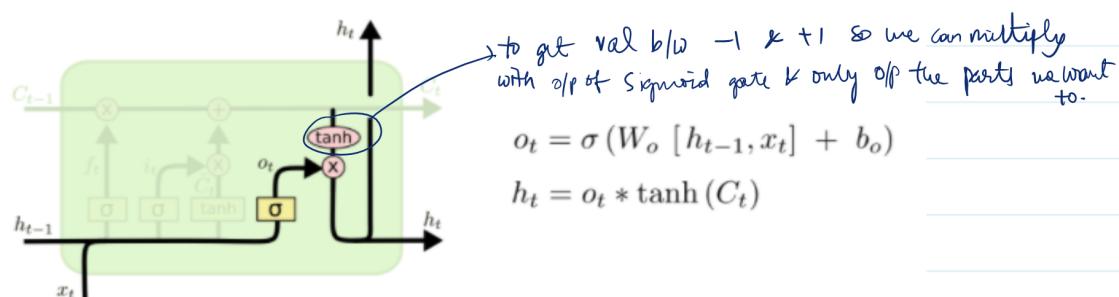


→ handles what new information we are going to store in cell state.

Cell State :-



Output Layer :-



to get val b/w -1 & +1 so we can multiply with o/p of Sigmoid gate & only o/p the parts we want to.

$$o_t = \sigma(W_o \cdot [h_{t-1}, x_t] + b_o)$$

$$h_t = o_t * \tanh(C_t)$$

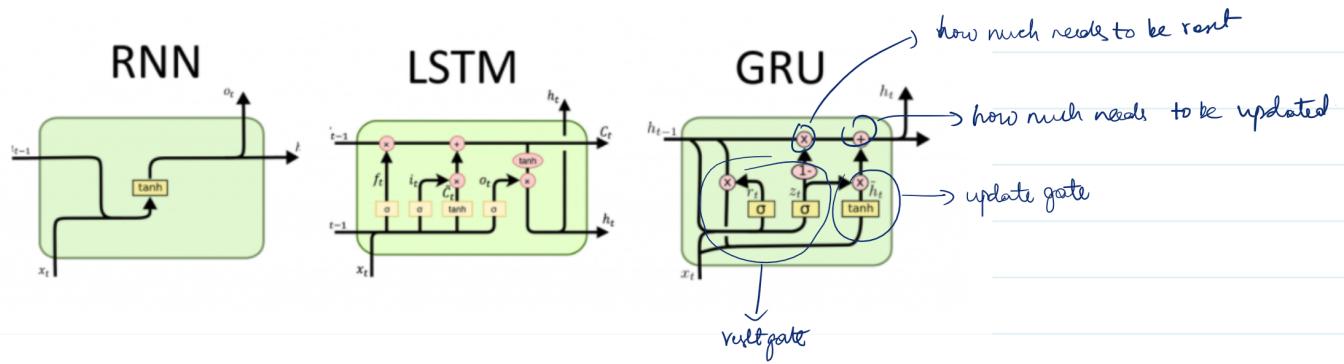
→ n weight. All are concatenation of x_t & h_{t-1} . All learn diff weight though.

GRU (Gated Recurrent Unit) :-

LSTM-1997 , LRU-2014

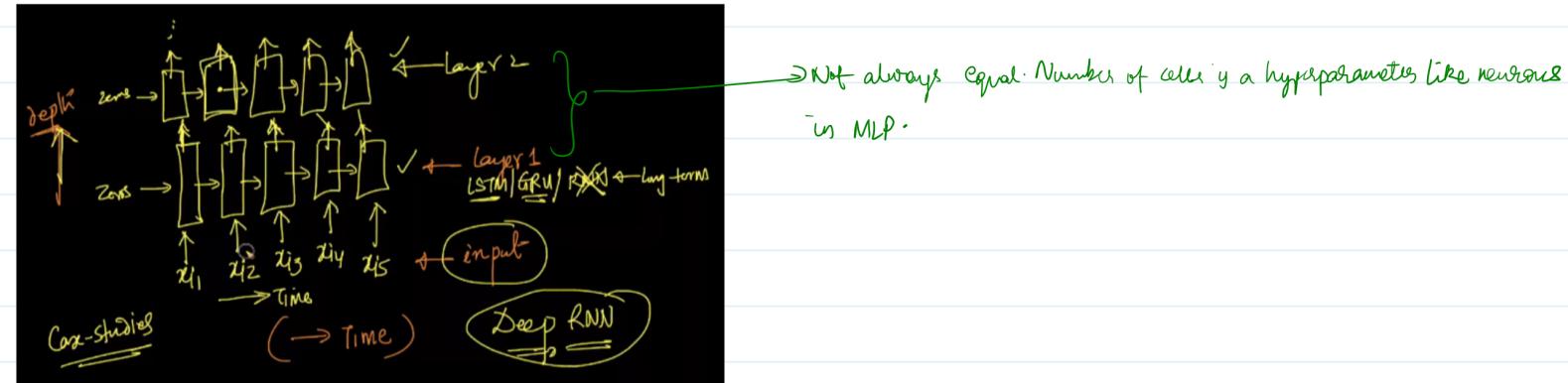
→ simplified & faster to train versions of LSTMs. As powerful as GRU. Merges the cell state & hidden state.

2 gates \rightarrow reset gate, update gate



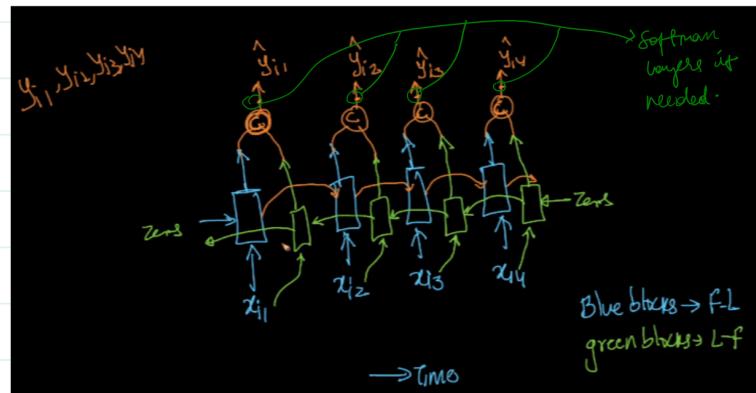
fewer equations \Rightarrow fewer updates \Rightarrow fewer backpropagation \Rightarrow faster.

Deep RNN (5.8):-



Bidirectional RNN(5.9):-

→ What if y_{it} is dependent on x_{it+k} ? Single RNN can't work with such data.



UML Cases :-

→ NLP tasks

→ Speech to text.

IMDB Text Classification (5.10):-

→ IMDb kelas

$$x_1 = \langle x_{11}, x_{12}, \dots \dots \rangle, y_1$$

$$x_2 = \langle x_{21}, x_{22}, \dots \dots \rangle, y_2$$

V = set of all the words in all the reviews.

Step 1:-

word	freq
The	10000
a	20000
:	:

Step 2:- sort by frequency

Step 3:- give words numbers

word	freq
a	20000
an	18000
the	10000
:	:

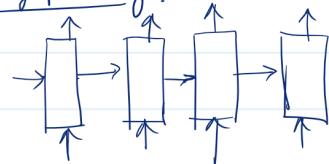
Step 4:- x_1 = This phone is a good phone

$$<81, \dots, 1, \dots >$$

$$x_{11}, \dots, x_{14}, \dots >$$

→ padding :- convert all reviews to the same length by inserting zeros at the beginning.

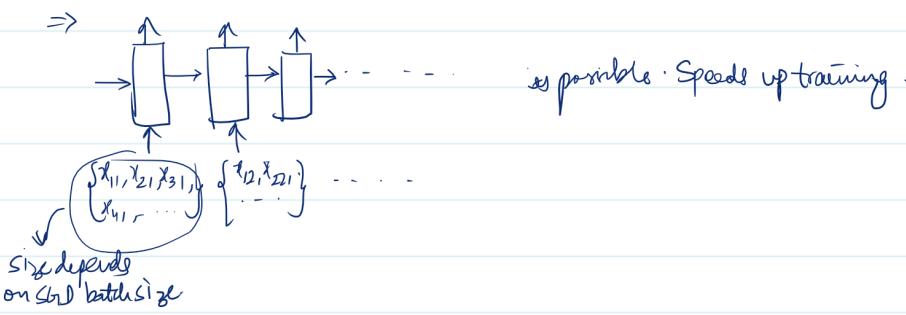
Why padding?



$$\begin{aligned} x_1 &= \langle x_{11}, x_{12}, \dots, x_{1150} \rangle \\ x_2 &= \langle x_{21}, x_{22}, \dots, x_{2100} \rangle \end{aligned} \quad \left. \begin{array}{l} \text{Normally, sent one ip after another.} \\ \downarrow \\ \text{Similar to SGD batch size = 1.} \end{array} \right\}$$

with padding, all elements have the same size.

⇒



is possible. Speeds up training.

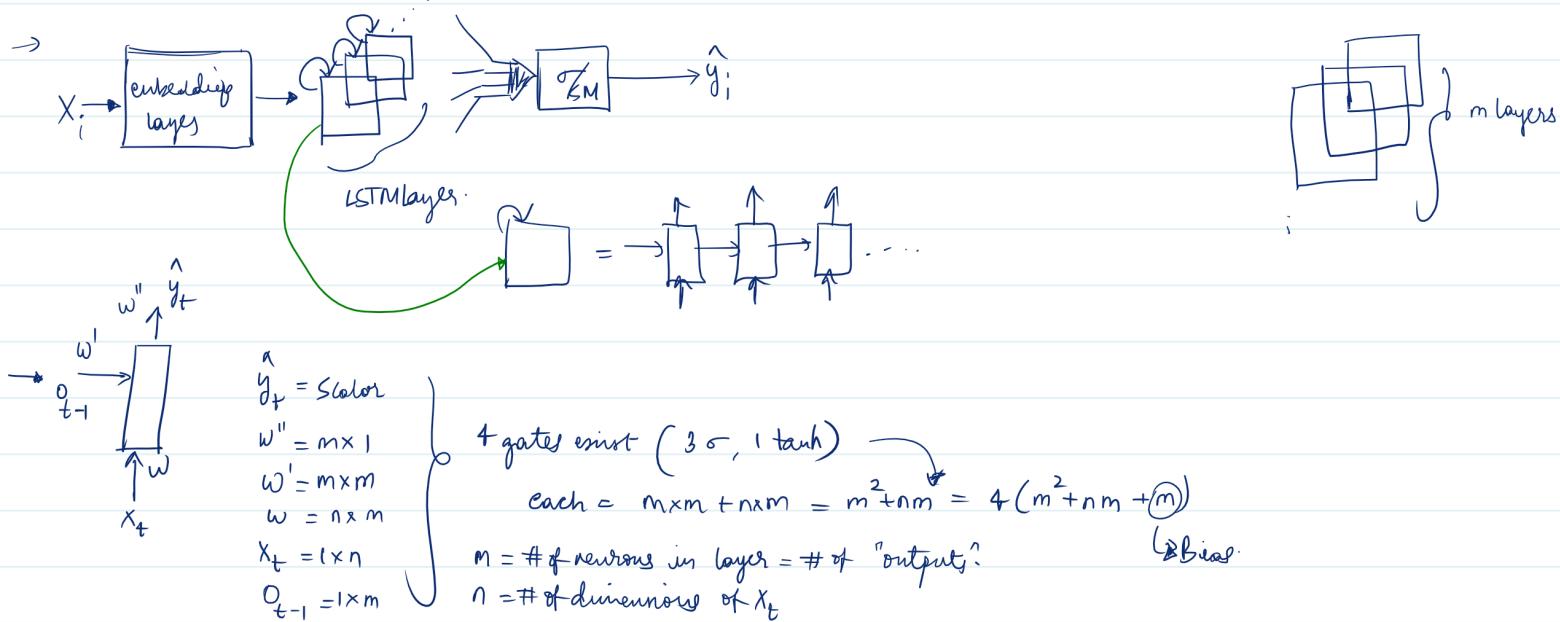
→ embedding takes integers & returns a vector.

e.g. - 1 → embedding → [.....]
 embedding vector length.

→ if there are m inputs & n outputs to LSTM,

$$\# \text{params} = 4(mn + n^2 + n)$$

bias (Keras by default adds bias)



2) Also I'm finding it difficult to understand top level on how the process happens, so here the input we'll pass in $x_{train[0]}$ right? It will be a vector of 600 dimension because of pad sequence. Then first embedding layer converts each feature of 600 features to 32 dimensions right? So it will be like a list of lists. Then LSTM take in each word(of 32 dim) i.e (first list in list of lists) and then takes the second word and so on and performs forward and backward prop and outputs the final binary decision for one sentence right? Is this how it works internally from a top level point of view. Please correct me if I'm mistaken anything.

Reply

Jul 02, 2019 02:46 AM

AppliedAI

yes, you are correct.