

Transfer learning

Why:-

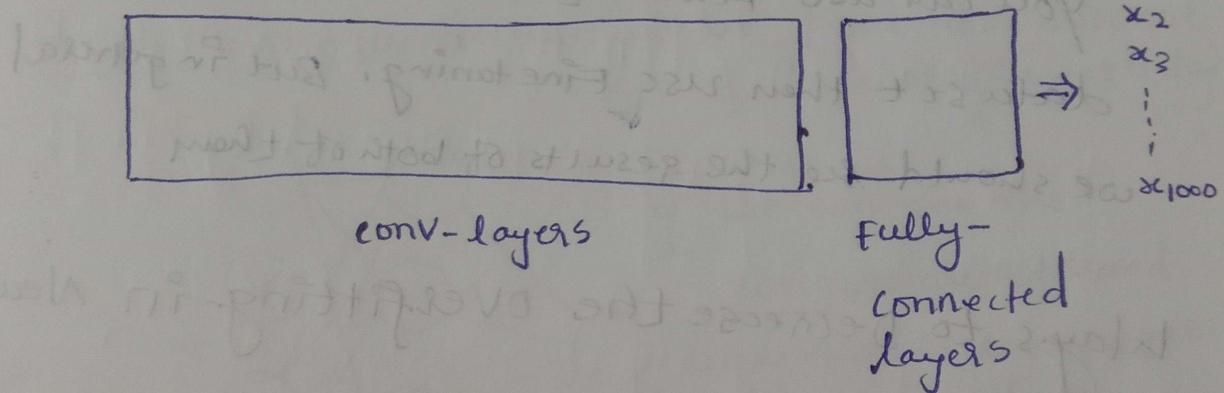
what if your all doing a project, where the pre trained models are not trained on that dataset.

Defⁿ: Transfer learning is a research problem in ML that focuses on storing knowledge gained while solving one problem and applying it to different, but related problems.

* How it works:-

we use the pretrained models \Rightarrow

e.g.: VGG16



Step-1: we freeze the conv-base or conv-layers

then we remove the FC layers.

Step-2: with the dataset u have we train a another Fully Connected layer.

Step 3: now we attach the trained dense Fully connected model with conv-layers that we freezed before.

Ways of Doing Transfer learning :-

- i) Feature extraction (discussed before)
- ii) Finetuning.

Fine tuning:-

In fine tuning we also train the last convolutional layers with Fully connected dense layers.

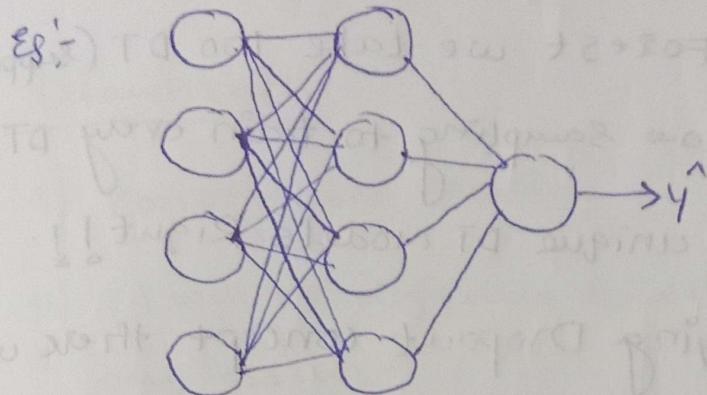
When you should use fine tuning:-

* Generally if you have a small dataset then you can use Feature extraction. If large dataset then use Finetuning. But in general we should see the results of both of them

Ways to Decrease the overfitting in Neural Networks.

- 1) Add more data (using data augmentation)
- 2) Reduce the complexity of Neural Network
- 3) Early stopping
- 4) Regularization L_1 L_2
- 5) Dropout
- 6) Batch Normalization

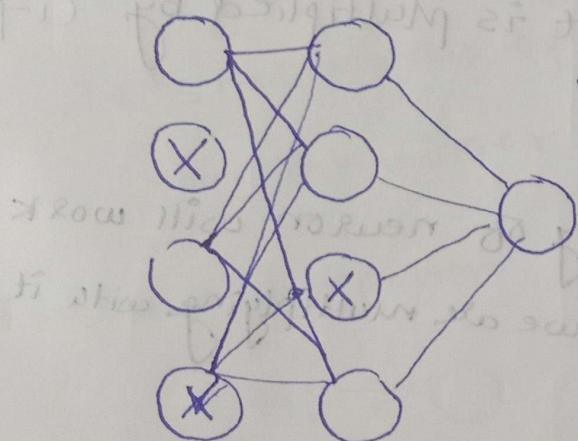
2) Dropouts:- It is technique where we turn off the nodes randomly from hidden & input layer in every epoch. (P)



initial ANN model with $p=0.25$ Then

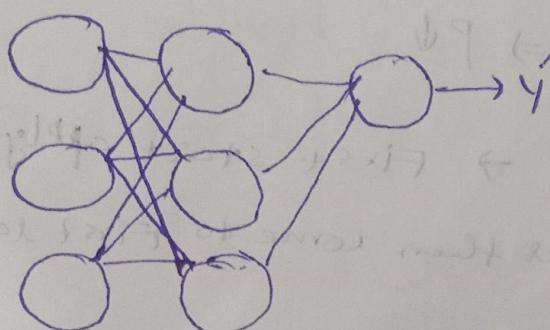
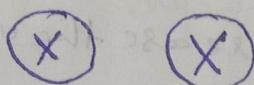
$$p=0.25 \quad p=0.25$$

1-epoch: The architecture will be: (one node will not work in each lay)



{ we can see the weights are reduced now there is less probability for overfitting }

(2-epoch



* Therefore we can say that there will be
[different] architecture.

* Random Forest Analogy:-

As in random forest we take 100 DT (suppose)
and will do col / row sampling to train every DT.

Now we have 100 unique DT models Right!!.

Similarly on applying Dropout concept there will
unique architecture in every epoch.

How prediction works:-

During testing All the neurons will work.

And every weight is multiplied by " $(1-p)$ "

value. Bec'-

$(1-p)$ = probability of neuron will work
right. so we are multiplying with it

Practical Tips:-

1) overfitting $\Rightarrow P \uparrow$ (increase the value of P)

underfitting $\Rightarrow P \downarrow$

2) last layer \rightarrow First start applying dropout
for last layer then come to first layers.

3) In case of RNN : p range \Rightarrow 0.1 to 0.5

CNN : $p \rightarrow 0.4 - 0.5$

RNN : $p \rightarrow 0.2 - 0.3$

Drawbacks

- i) convergence \rightarrow delay
- ii) It will be difficult to debug back when had any issues.

Why RNN's

Eg: hi my name is Jayath (5) 1
It is not my cup of tea (5) 0
you're the winner (3) 1

Let's take a preposition

*hi \rightarrow O

my \rightarrow O

name \rightarrow O

is \rightarrow O

Jay \rightarrow O

unique words (+ve/-ve) output

13

O

O

O

O

O

O

O

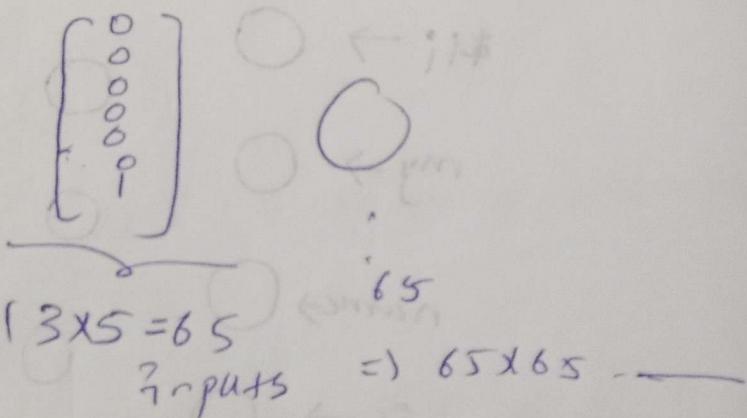
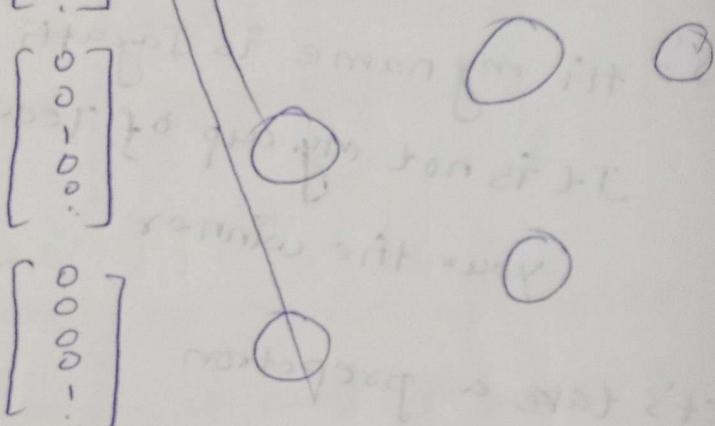
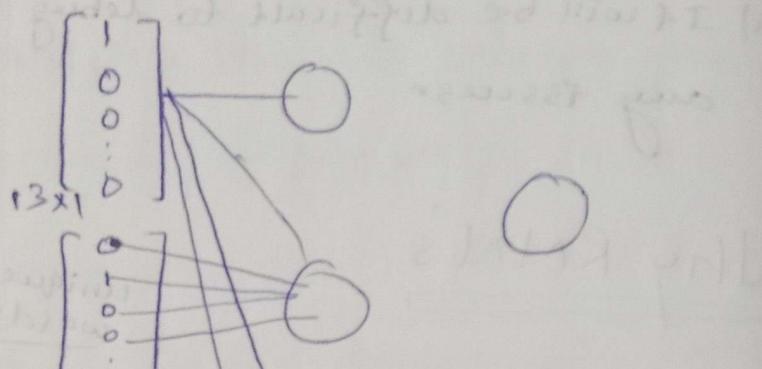
* on doing one hot encoding, we convert into vectors.

$$tr = [1, 0, 0, 0, \dots, 0] \underbrace{0}_{13}$$

$$my = [0, 1, 0, \dots, 0] \underbrace{0}_{13}$$

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

how number of nodes in input and output layer



problems with having large vocabulary

- 1) Text input varying size and shared embeddings
- 2) zero padding \Rightarrow unnecessary computation
- 3) If you did zero padding also, what if the test or new text has different no. of unique words.
Then we get prediction problem

4) Totally disregarding the sequence information.

Then they thought of different architecture called "RNN".

Applications of RNN's - $[0,1,0,1,0] = \text{book}$

1. NLP:

- i) language modelling
- ii) machine translation
- iii) sentiment analysis
- iv) named entity recognition
- v) text generation

2. Speech Recognition (speech to text)

3. Time series Analysis

4. Image captioning

5. Anomaly detection in text (like behaviour)
So...on.

RNN's :- Rnn's are the special type of Neural networks which has memory feature; that is why they work great on "sequential data".

Data for RNN :-

format \Rightarrow (Timesteps, input-features)

	Review	Sentiment
movie =	[movie] was [food]	1
	movie was [bad]	0
was =	movie was [not] good	0
	No. of unique words = 5	

good = [0, 0, 1, 1, 0] - so on.

Now Review-1 can be written as

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

$R-1 \Rightarrow$ [Timestamps, input features]

next t=1 sec for movie to go into archi

t=2 sec for was, rest for word.

t=3 sec for word all.

$\therefore R-1 \Rightarrow [3, 5]$ no. of unique features

$\parallel R-2 \Rightarrow [3, 5]$ no. of unique elements

$$R-3 \Rightarrow [4, 5]$$

* In keras \rightarrow simple RNN \rightarrow
 (batch size, timestamps, input-features)

For above ex:- $(3, 4, 5)$
 ↓ ↓ ↓
 3 4 5
 As we have max.no.of timestamps among the
 3 reviews reviews

How it works?

Above example: Reviews sentiment

$R_1 \Rightarrow$ movie was good 1
 x₁₁ x₁₂ x₁₃

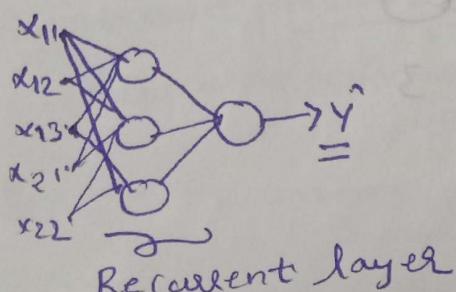
$R_2 \Rightarrow$ movie was bad 0
 x₂₁ x₂₂ x₂₃

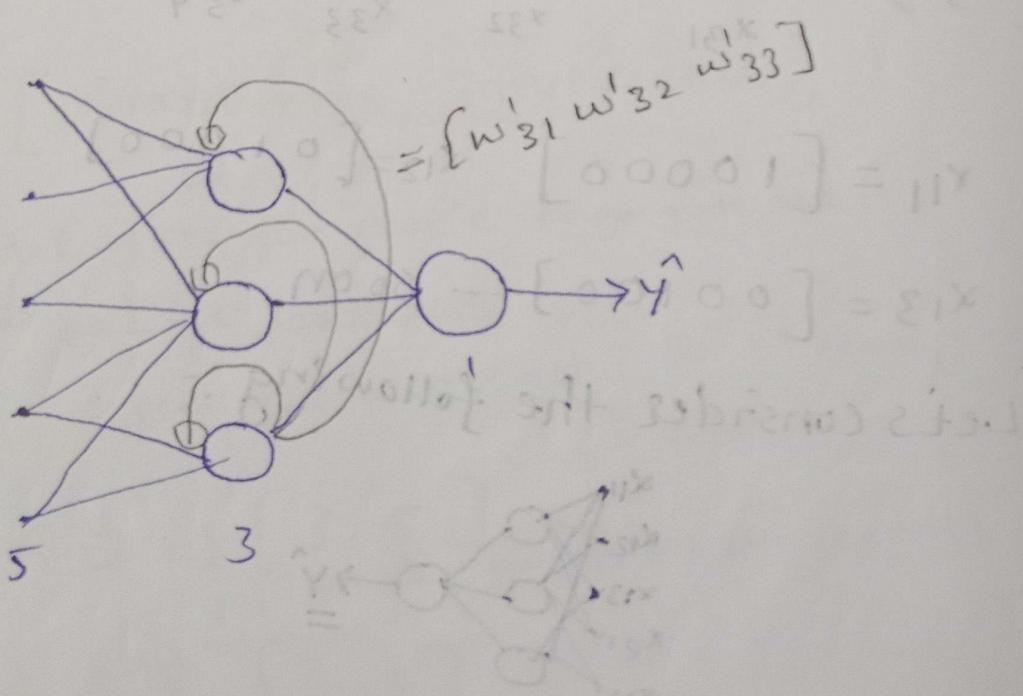
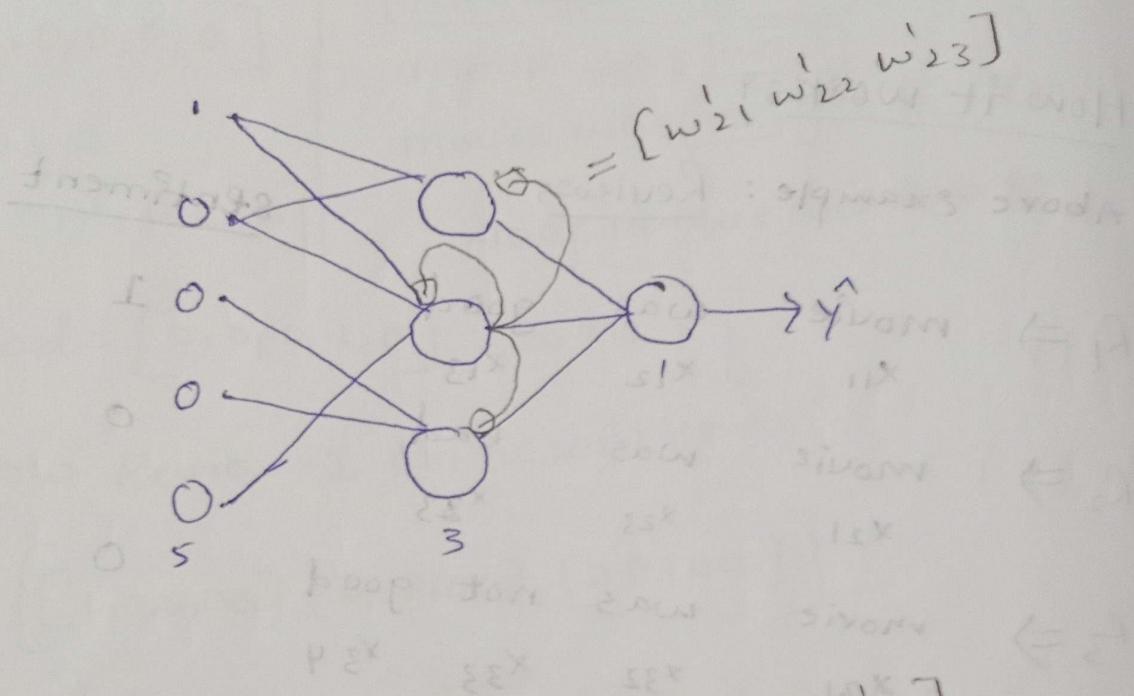
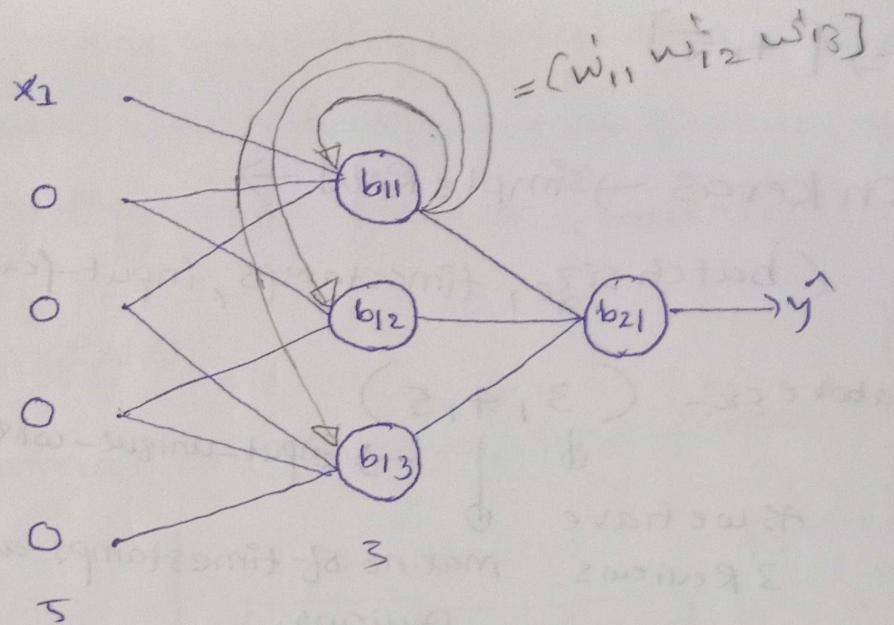
$R_3 \Rightarrow$ movie was not good 0
 x₃₁ x₃₂ x₃₃ x₃₄

$$x_{11} = [1 0 0 0 0] \quad x_{12} = [0 1 0 0 0]$$

$$x_{13} = [0 0 1 0 0] \dots \text{so on.}$$

Let's consider the following:-





Total trainable parameters = ~~add prior seen now~~

$$w_1 = 5 \times 3 = 15 \text{ weights}$$

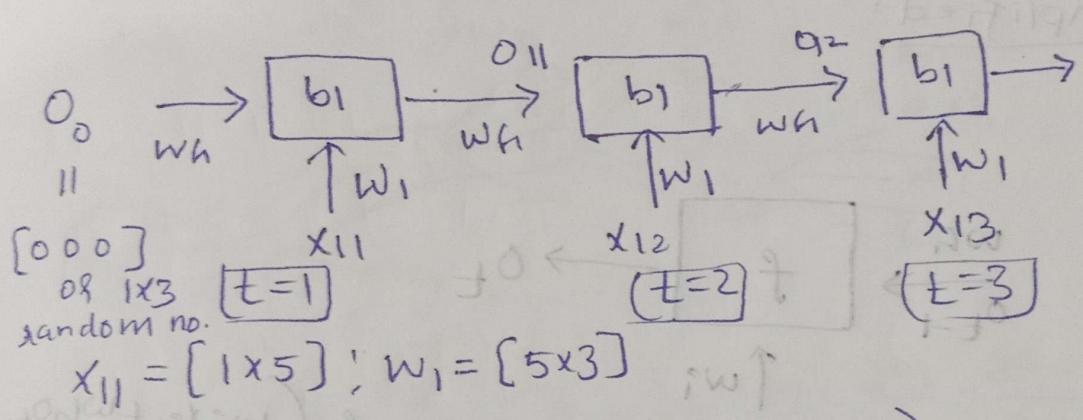
$$w_h = 3 \times 3 = 9 \text{ weights hidden (Recurrent layer)}$$

$$w_2 = 3 \times 1 = 3 \text{ weights}$$

27 weights

$$\text{biases} = 4 \therefore \boxed{T = 31} \text{ trainable parameters}$$

Forward propagation:-



$$x_{11} = [1 \times 5], w_1 = [5 \times 3]$$

$$o_{11} = f(w_1 x_{11} + b_1 + w_h o_0) = [1 \times 3]$$

$$o_{12} = f(w_1 x_{12} + o_{11} w_h + b_1) = [1 \times 3]$$

$$o_{13} = f(w_1 x_{13} + o_{12} w_h + b_1) = [1 \times 3]$$

$$\text{Now } w_2 = [3 \times 1]$$

$$y_1 = f(w_2 o_{11} + b_2) = [1 \times 1] = \text{scalar}$$

In RNN \Rightarrow At every time stamp we send the

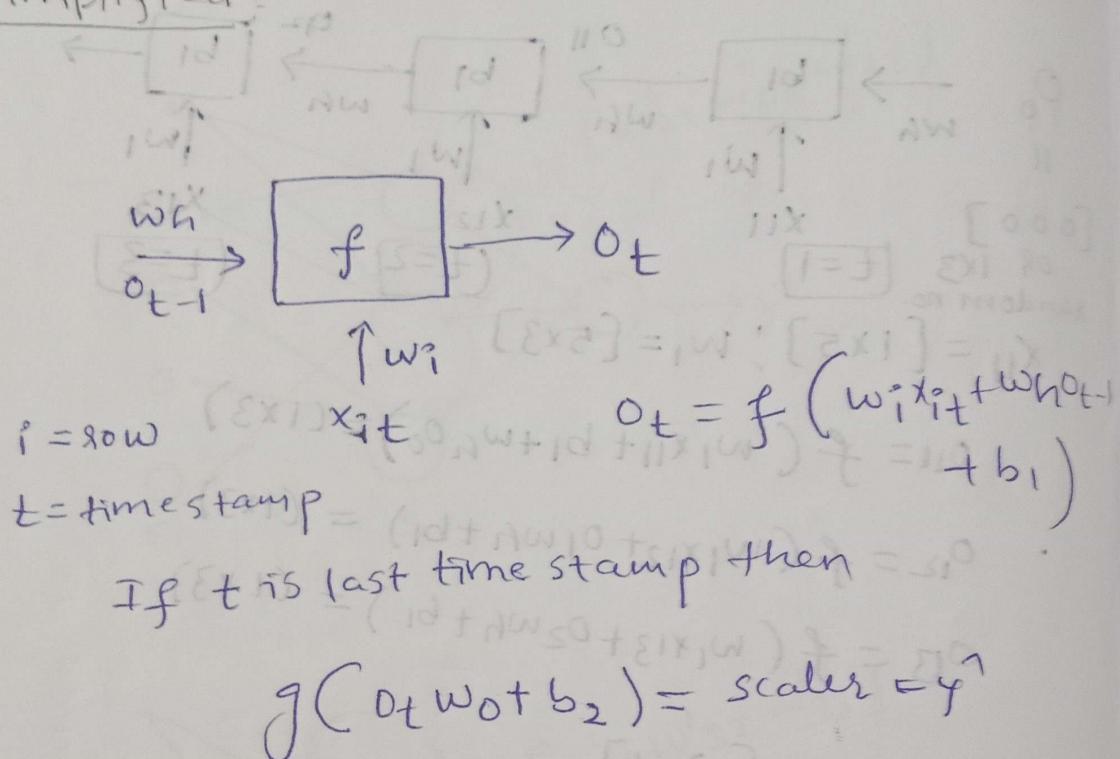
current input and previous output for the new network

Reoccurring(OR) Reusing it as feedback for entire sentence

That's why we call Recurrent NN.

- * we are using the same weights and bias for complete sentence forward prediction; That implies we are using the concept of parameter sharing or weight sharing.
- * RNN is actually remembering the previous words so that the sequence or the context of sentence is maintained.

Simplified:



Backpropagation:

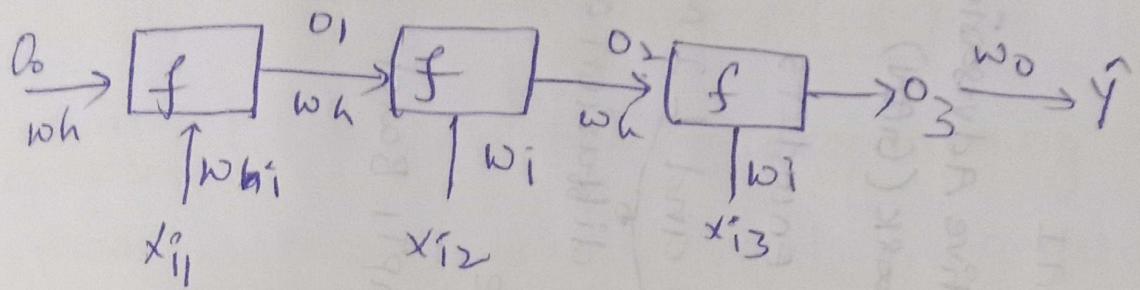
$$o_1 = f(x_{i1} w_i + o_0 w_h)$$

$$o_2 = f(x_{i2} w_i + o_1 w_h)$$

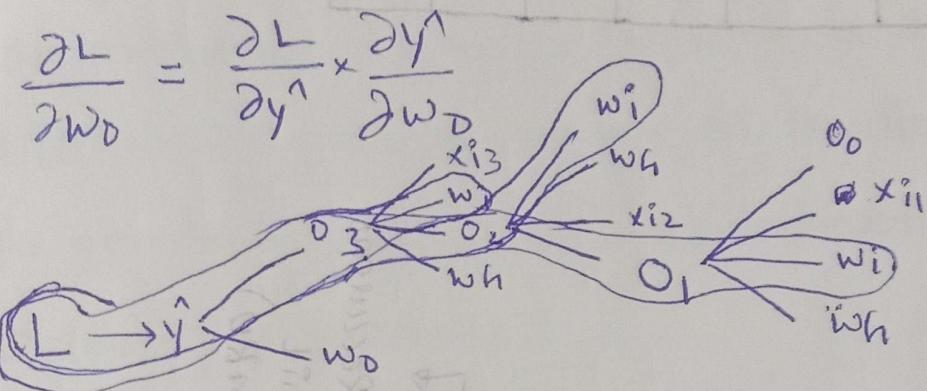
$$o_3 = f(x_{i3} w_i + o_2 w_h)$$

AIM: $\frac{\partial L}{\partial w_0}, \frac{\partial L}{\partial w_h}, \frac{\partial L}{\partial w}$

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



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



$$\frac{\partial L}{\partial w_o} = \frac{\partial L}{\partial y^{\hat{}}} \cdot \frac{\partial y^{\hat{}}}{\partial o_3} \cdot \frac{\partial o_3}{\partial w_o} + \frac{\partial L}{\partial y^{\hat{}}} \cdot \frac{\partial y^{\hat{}}}{\partial o_3} \cdot \frac{\partial o_3}{\partial o_2} \cdot \frac{\partial o_2}{\partial w_i}$$

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

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

(As we took 3 weights)

$$\text{If } \frac{\partial L}{\partial w_h} = \sum_{j=1}^3 \frac{\partial L}{\partial y^{\hat{}}} \cdot \frac{\partial y^{\hat{}}}{\partial o_j} \cdot \frac{\partial o_j}{\partial w_h}$$

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

$n = \text{no. of time stamps}$

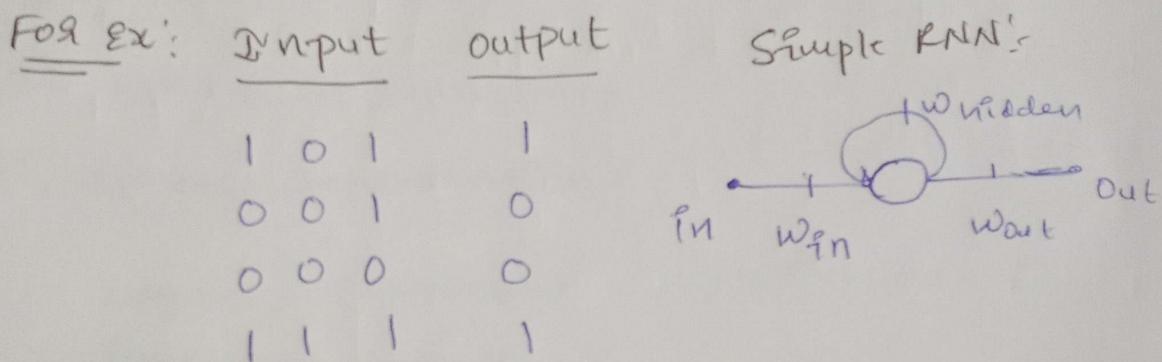
Problems with RNN

i) Long term dependencies caused by

Vanishing gradient problem

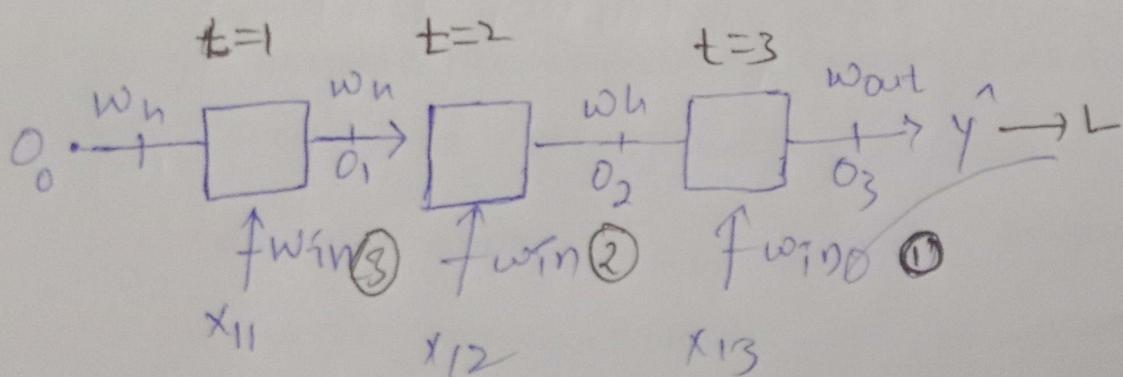
ii) Stagnated Training \Rightarrow reason \Rightarrow exploding gradient problem

#problem 1:- Vanishing gradient descent



$$x = [1, 0, 1]_{1 \times 3}$$

remember we have only one neuron in rnn, but
3 input vectors are there, so we write 3 times
at different time stamps.



$$w_{in} = w_{in} - \left(\frac{\partial L}{\partial w_{in}} \right)$$

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

$$w_{out} = w_{out} - \frac{\partial L}{\partial w_{out}}$$

we apply chain rule; short term dependency

$$\frac{\partial L}{\partial w_{in}} = \frac{\partial L}{\partial y^1} \frac{\partial y^1}{\partial o_3} \frac{\partial o_3}{\partial w_{in}} + \frac{\partial L}{\partial y^1} \frac{\partial y^1}{\partial o_3} \frac{\partial o_3}{\partial o_2} \frac{\partial o_2}{\partial w_{in}} +$$

$$\boxed{\frac{\partial L}{\partial y^1} \frac{\partial y^1}{\partial o_3} \frac{\partial o_3}{\partial o_2} \frac{\partial o_2}{\partial o_1} \frac{\partial o_1}{\partial w_{in}}}$$

[see diagram
in front page]

Long term dependency

we got it for only

3 time stamps.

what if 100 time stamps!, It will look

like this

$$\Rightarrow \frac{\partial L}{\partial y^1} \frac{\partial y^1}{\partial o_{100}} \frac{\partial o_{100}}{\partial o_{99}} \dots \frac{\partial o_2}{\partial o_1} \frac{\partial o_1}{\partial w_{in}} \quad A$$

Imagine how much small it become!

$$\textcircled{A} \quad \frac{\partial L}{\partial y} \frac{\partial y}{\partial o_{100}} \left[\begin{array}{c|c} \frac{\partial o_{100}}{\partial a_{100}} & \frac{\partial a_{100}}{\partial w_{100}} \\ \hline \end{array} \right] = \frac{\partial L}{\partial w_{100}}$$

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

$$o_1 = \tanh \left(x_{11} w_{1n} + o_0 w_n \right)$$

$$o_t = \tanh f(x_{it} w_{in} + o_{t-1} w_n)$$

$$\frac{\partial o_t}{\partial o_{t-1}} = \tanh' (x_{it} w_{in} + o_{t-1} w_n) \quad (w_n)$$

Range: $0-1$ Depends on o_n

$$\textcircled{A} = \frac{\partial L}{\partial y} \frac{\partial y}{\partial o_{100}} \boxed{\prod_{t=2}^{100} \left(\tanh' (x_{it} w_{in} + o_{t-1} w_n) w_n \right)}$$

Assume you have given $w_n : 0$ to 1 ↗

Image range
how small it will be

Gradient is Vanishing ≈ 0

$$\therefore \textcircled{A} \approx 0$$

\therefore The change in L w.r.t w_{in} will be
only contributed by short term dependencies

Solution:

- 1) Using different activation function
instead of \tanh (Range: 0 to 1) like ReLU
/ leaky ReLU
- 2) Better weight initialization (not from
 $0-1$)
- 3) Skip RNNs
- 4) LSTMs

#problem 2: Exploding Gradient descent

If we are using " \tanh " and initializing weights
from $0-1$ then we are getting vanishing
gradient problem.

Assume you have used ReLU and initialized
weights greater than 1

Then ① will dominate the short term dependencies. I mean

$$\frac{\partial L}{\partial w_{in}} = \underbrace{\frac{\partial L}{\partial y^1} \frac{\partial y^1}{\partial o_3} \frac{\partial o_3}{\partial w_{in}}}_{\text{initial}} + \underbrace{\frac{\partial L}{\partial y^1} \frac{\partial y^1}{\partial o_3} \frac{\partial o_3}{\partial o_2} \frac{\partial o_2}{\partial w_{in}}}_{\text{short term}}$$

$$+ \underbrace{\frac{\partial L}{\partial y^1} \frac{\partial y^1}{\partial o_3} \frac{\partial o_3}{\partial o_2} \frac{\partial o_2}{\partial w_{in}}}_{\text{long term}} \quad \textcircled{1}$$

These will become negligible compared to ①
become very very large number if we take 100 time
Stamps!

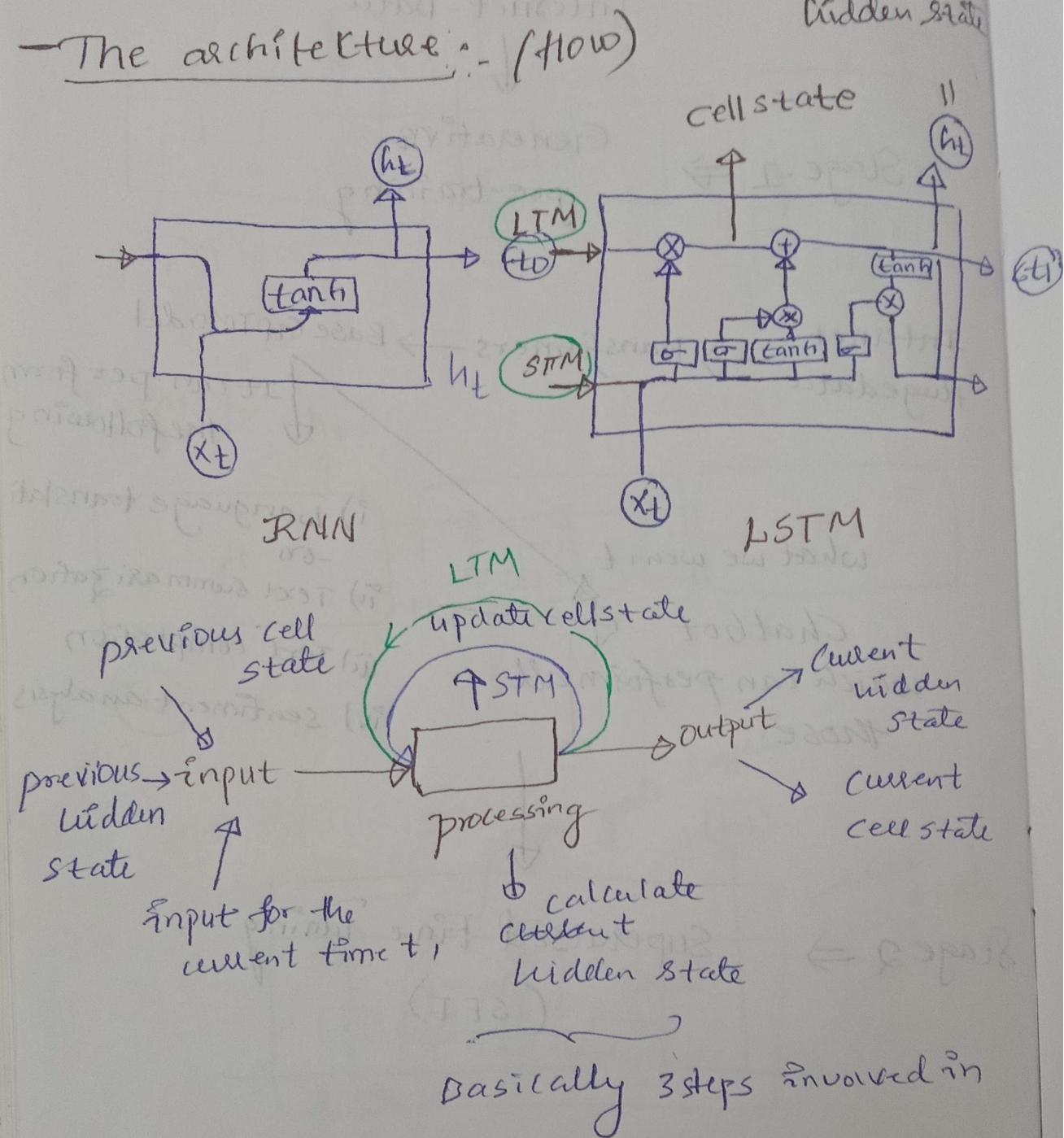
Solutions:

- 1) Gradient clipping
- 2) Controlled learning rate
- 3) LSTMs

LSTM's : Long short Term Memory

- The what

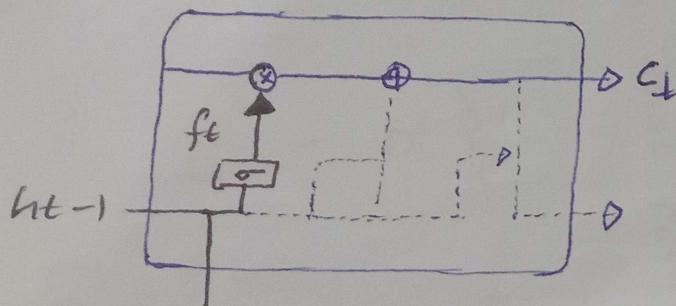
- The architecture :- (flow)



Core Idea: In RNNs we have only one path (h_t) which captures the short term dependencies but in LSTM scientists added another path cell state which also stores long term dependencies (LTM)

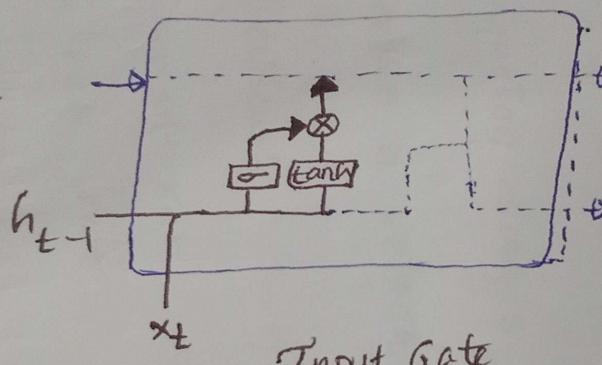
(2)

Gates



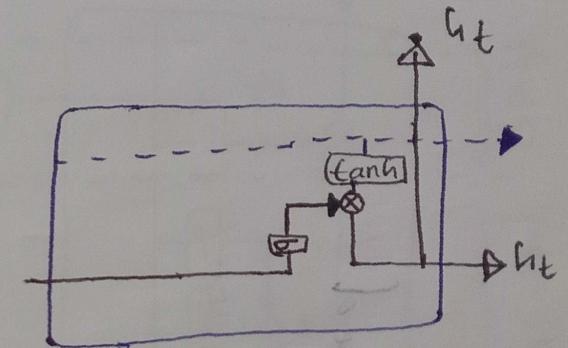
x_t forget gate

To remove something from
 c_t

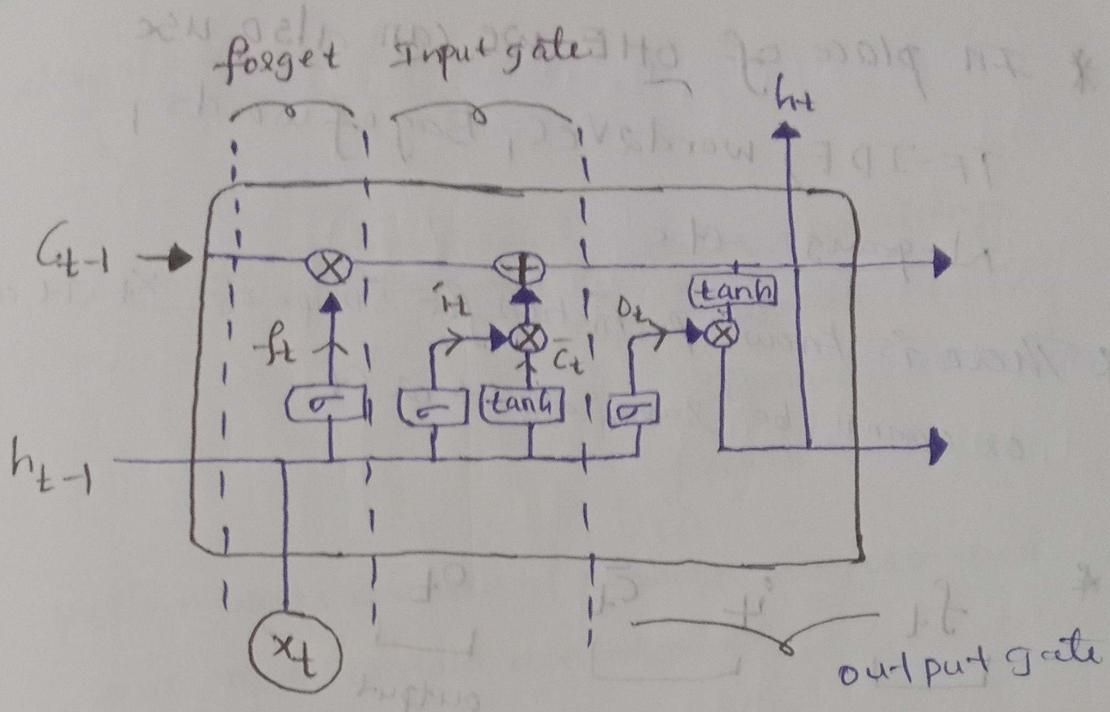


x_t Input Gate

To add something to
 c_t



output gate
to calculate hidden
state h_t



- c_{t-1} & h_{t-1} both are vectors of same dimension always. If c_t & h_t
 - x_t is also a vector like;

E

Reviews Sentiment

envelope sairthing
o

cat mat rat ← o x

cat rat rat ← 1

mat mat cat ← +

the [cat mat rat] ← + 0 -

cat mat rat ← + 0 -

1 0 0 ← +

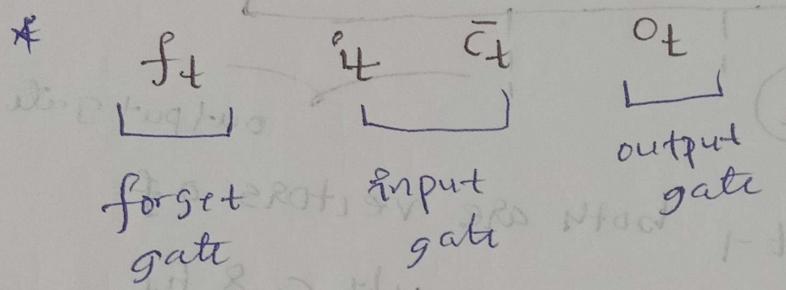
[0 0 0] ← + 1 1 1

0 0 0 ← + 1 1 1

⇒ The word 'cat' = [1 0 0]

* In place of OHE we can also use
TF-IDF, word2vec, Bag of words,
N-grams, etc.

* There is known restriction for shape of
or cannot be same as x_t or c_t .
 x_t It can



All the above are vectors
∴ The shape is same as $[c_t \times m]$

pointwise operations

$$\textcircled{X} \rightarrow$$

Ex:- $c_{t-1} = \begin{bmatrix} 4 & 5 & 6 \\ 7 & 8 & 9 \end{bmatrix}$
 $f_t = \begin{bmatrix} 1 & 2 & 3 \end{bmatrix}$

$$c_{t-1} \otimes f_t \Rightarrow \text{vector}$$

$$f_{oy} = \begin{bmatrix} 4 & 10 & 18 \end{bmatrix}$$

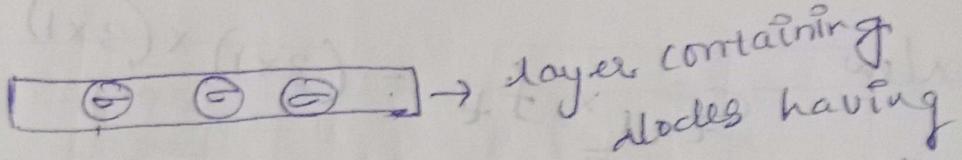
$$||/y \oplus \rightarrow$$

$$\text{tanh} \rightarrow [\tanh(4), \tanh(5), \tanh(6)]$$

$$c_{t-1} \text{ say } [0.12 \ 0.36 \ 0.21]$$

$$[0.01 \ 0.1] = \text{Ans}$$

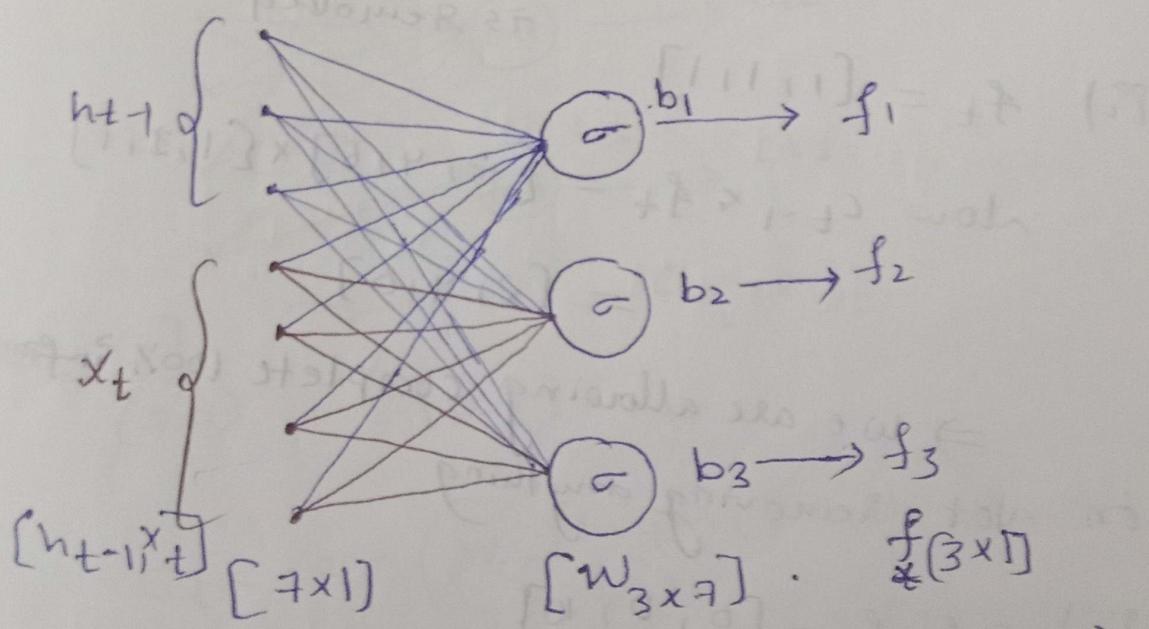
Neural network layers



- * Yellow boxes are neural network layers and each node contain activation fn.

The Forgate :-

Let (suppose) x_t is a 4 dimensional vector
and h_{t-1} is a 3 dimensional vector



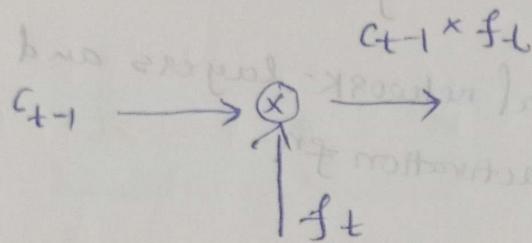
$$(3 \times 7)(7 \times 1) = (3 \times 1)$$

$$f_t = [f_1 \ f_2 \ f_3]$$

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

Now forward $c_{t-1} \times f_t = (3 \times 1) \times (3 \times 1)$

Intuition behind word "Removal" used for \otimes



i) Let's say $c_{t-1} = (2, 4, 6)$

and $f_t = (\frac{1}{2}, \frac{1}{2}, \frac{1}{2})$

Now $c_{t-1} \times f_t = (2, 4, 6) \times (\frac{1}{2}, \frac{1}{2}, \frac{1}{2})$
 $= 2(1, 2, 3) \Rightarrow$ only 50% is removed

ii) $f_t = [1, 1, 1]$

Now $c_{t-1} \times f_t = [2, 4, 6] \times [1, 1, 1]$

$= [2, 4, 6]$

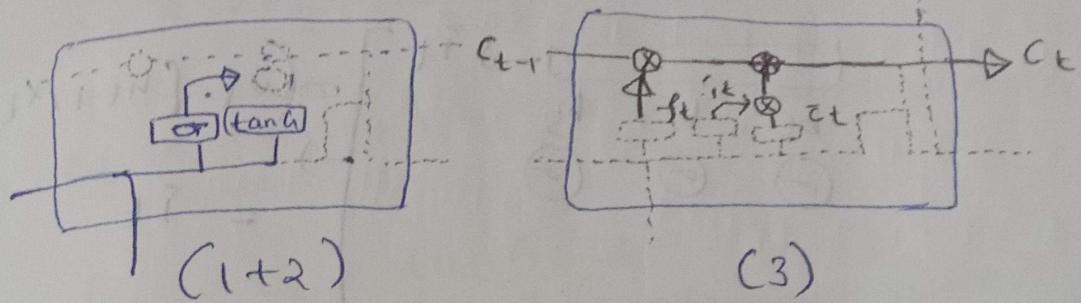
\Rightarrow we are allowing complete 100% information

- in Not removing anything

iii) If $f_t = [0, 0, 10]$

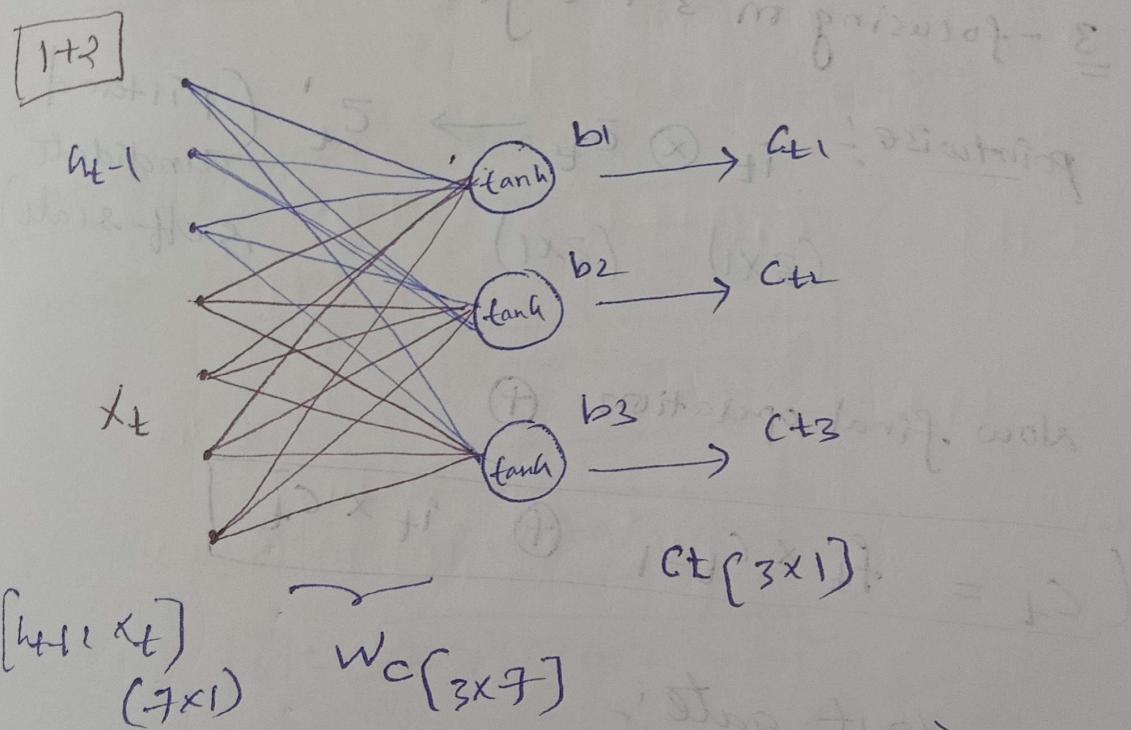
Now $f_{t-1} \times f_t = [0, 0, 10] \Rightarrow 100\%$
 information is removed.

The Input gate :- To Add new input info to c_t



stages cell

- i) \bar{c}_t candidate state
- ii) $i_t \rightarrow$ it is basically filtering
- iii) updating cell state



$$\therefore \bar{c}_t = \tanh \left(W_C [h_{t-1}, x_t] + b_C \right)$$

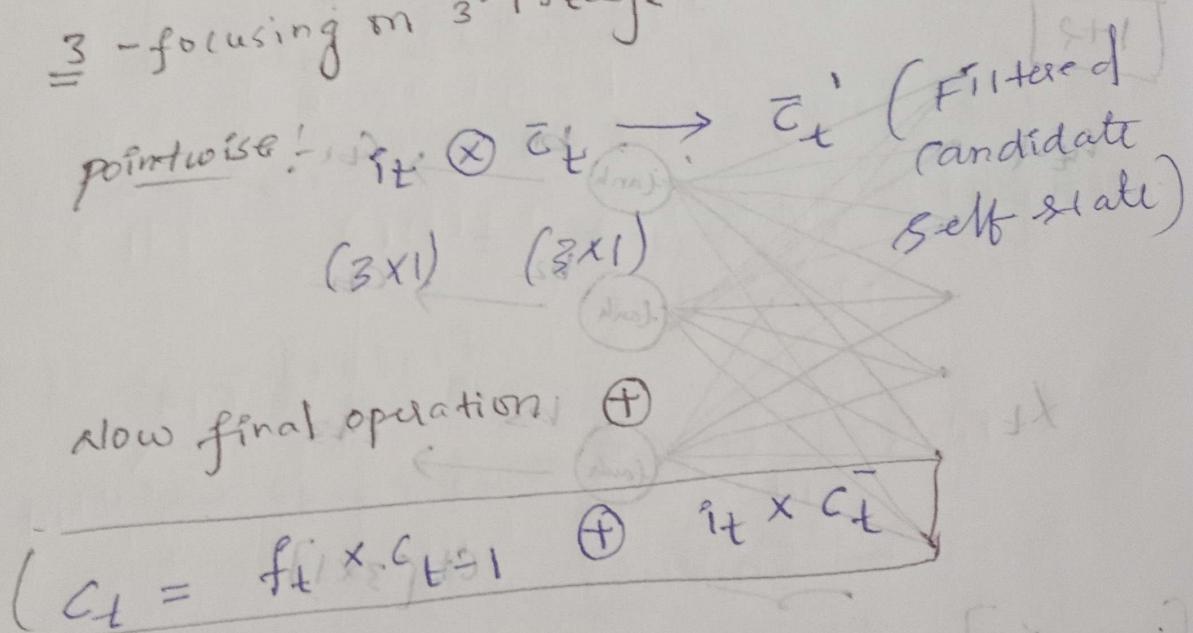
$$\therefore \bar{c}_t \quad 3 \times 1$$

Same for it also

$$\begin{array}{c} \text{Input } x_t \rightarrow i_t \\ \begin{bmatrix} b_1 & b_2 & b_3 \end{bmatrix} \\ \ominus \quad \ominus \quad \ominus \\ \text{(i)} \end{array} \left. \begin{array}{l} \alpha(\text{vec. } [h_{t-1}, x_t] + b_i) \\ = s_t \\ \text{(ii)} \end{array} \right\}$$

h_{t-1} state, it's bias β (i)
 future forward pass \rightarrow filter for c_t (ii)
 Here it acts as filter for i_t (iii)

3 - focusing in 3rd stage



The output gate :-

goal: $c_t \rightarrow h_t$

i) $\tan(c_t) \rightarrow [-1, 1] \Rightarrow$ now we apply
filtration

i) For filtration we require a vector called o_t

$$o_t = \sigma [w_o[h_{t-1}, x_t] + b_o]$$

$o_t \times \tanh(c_t)$ - pointwise multiplication

ii)

$$h_t = o_t \times \tanh(c_t)$$

$(3 \times 1) \quad (3 \times 1)$ - element wise

+ multiplication

Word2Vec

queen - girl + boy

*

Latent semantic Analysis

Cosine similarity :-

$$\cos\theta = \frac{\vec{a} \cdot \vec{b}}{\|\vec{a}\| \|\vec{b}\|}$$

$$\|\vec{a}\| = \sqrt{a_1^2 + a_2^2 + a_3^2 + \dots + a_n^2}$$

$$\|\vec{b}\| = \sqrt{b_1^2 + b_2^2 + b_3^2 + \dots + b_n^2}$$

Topic modelling :-

* The process to find the topics from documents in an unsupervised manner.

* Text mining approach to find recurring patterns

the text documents. (nothing but a row)

* Algorithms or techniques for topic modelling:

1. LDA - Latent Dirichlet Allocation

2. NMF - Non-Negative Matrix Factorization

3. LSA - Latent Semantic Allocation.

Aim of LDA:

i) what are the most important topics

ii) what are the topics assigned to each document

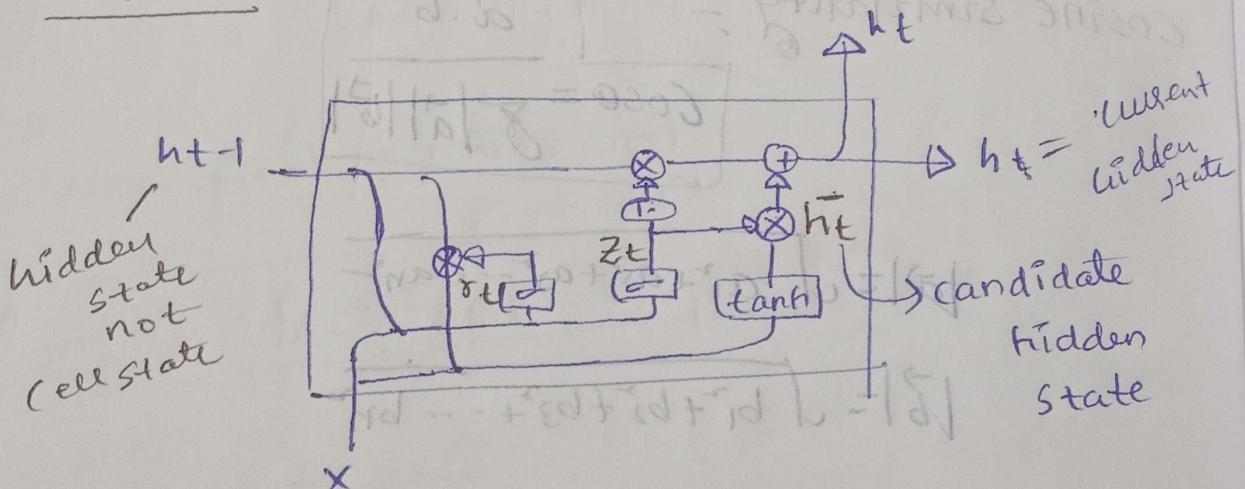
OR

iii) what is the topic to the term distribution (are there like food, clothing but annotating with)

iv) Document to topic distribution. [clothing but annotating with]

words in document with topics].

GIRU's: GATED RECURRENT UNIT



Why?

* It is less complex than LSTM's

* ⇒ Its training time will be less than LSTM's

* GRU ≈ LSTM's

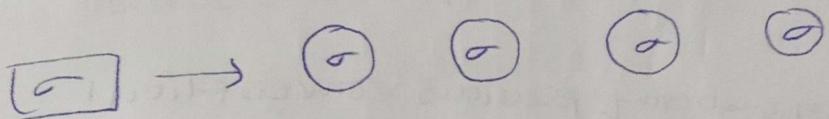
* only two gates (Reset gate & update gate)

* we manipulate the hidden state so that it can carry both long term memory & short term memory.

Goal: $h_{t-1} \xrightarrow{x_t} [h_t]$ current hidden state (?)

$\rightarrow h_{t-1}, h_t, r_t, [x_t], z_t, \tilde{h}_t \rightarrow$ All are vectors
and except x_t all their dimension is same.

$\rightarrow [\sigma] \quad [\tanh] \Rightarrow$ These are neural networks with activation fn sigmoid & tanh



rem: No. of neurons will same in $[\sigma] + [\tanh]$

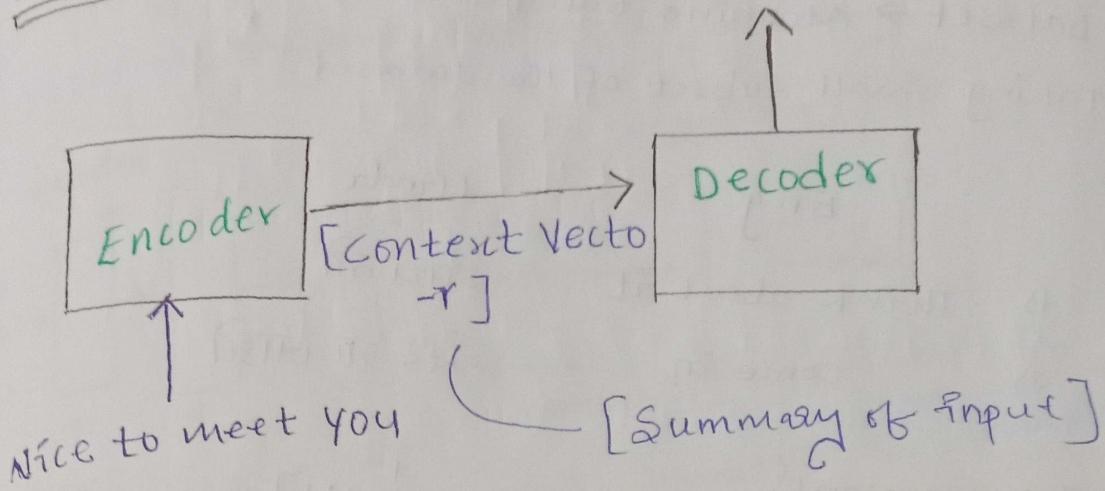
Number of training steps =

$$\left(\frac{\text{No. of rows (or examples)}}{\text{Batch size}} \right) \times \text{Number of epochs}$$

EDA - Encoder Decoder Architecture

Overview:-

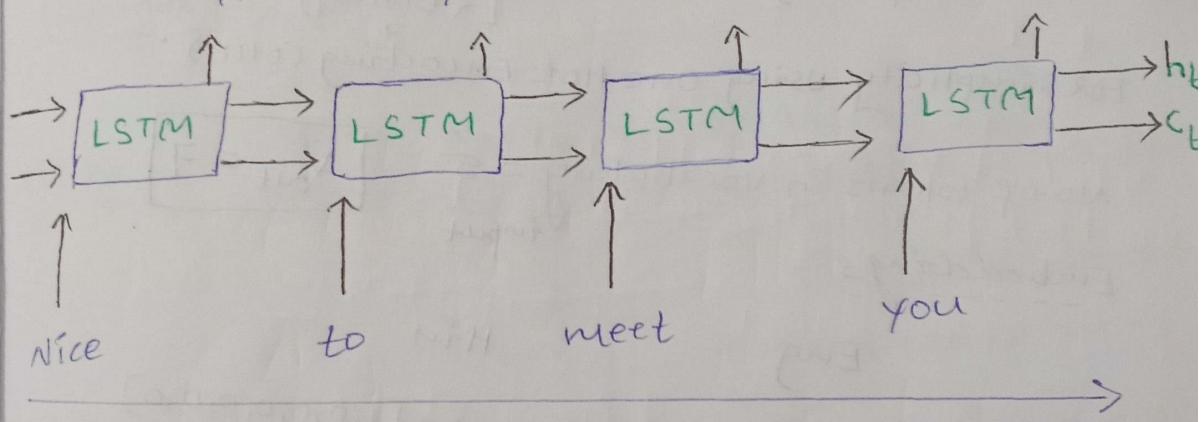
आप से मिल कर अच्छा लगा।



What's under the hood?

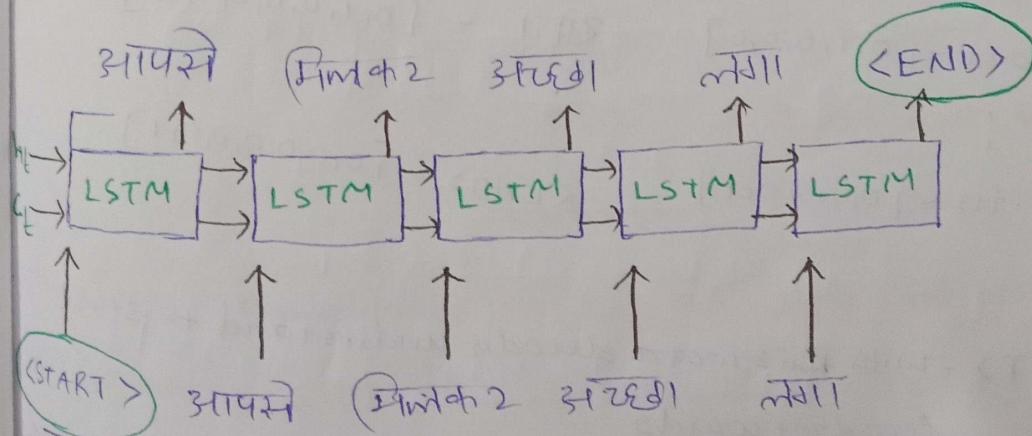
Encoder:

These outputs are ignored



Decoder:

timesteps



Time steps.

Considering Machine translation task - Part 1

Training the architecture using Backpropagation

Dataset → Machine translation

- * Taking small subset of the dataset

Eng	Hindi
1) Think about it	सोच लो
2) come in	आइ आ उत्तमो

Tokenization:-

[¹ think, ² about, ³ it]	[<start>]
[come, in]	[सोच, ³ लो]
4 5	[आइ, ⁵ आ, ⁶ उत्तमो]
	[<END>]

For simplicity using one hot encoding (OHE):

No. of tokens in vocabulary = 5

$$V_{out} = 7$$

Embeddings :-

Eng	Hin
think - [1, 0, 1, 0, 0]	<START> - [1, 0, 1, 0, 0, 0, 0]
about - [0, 1, 1, 0, 0]	सोच - [0, 1, 1, 0, 1, 0, 0]
:	:
.in - [0, 1, 0, 1, 0, 1]	<END> - [0, 1, 0, 1, 0, 0, 1]

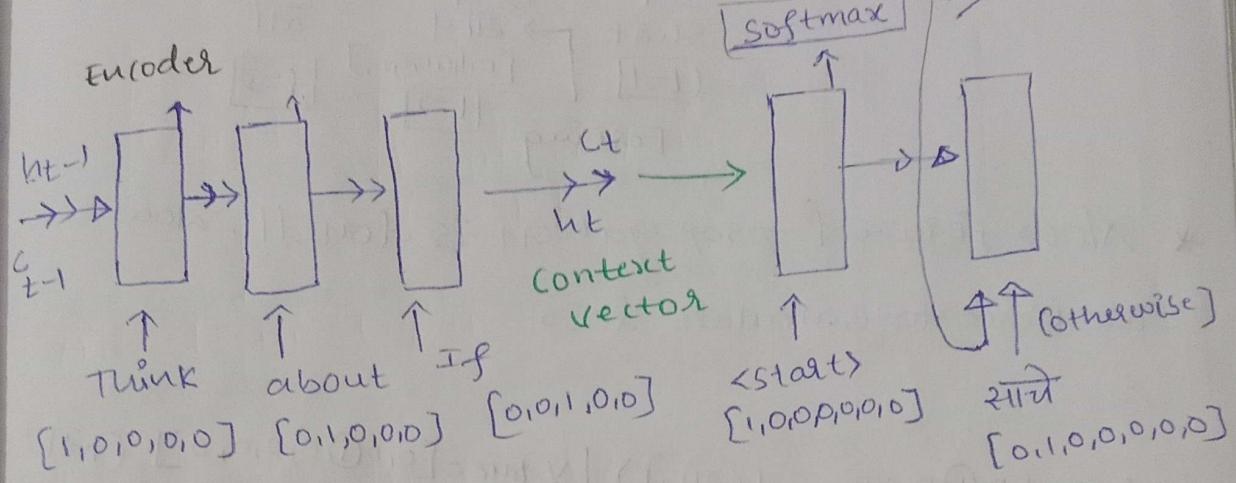
<START> - with this word decoder understand to start decoding words

<END> - After decoding all the words, decoder will through 'END' word

Training:-

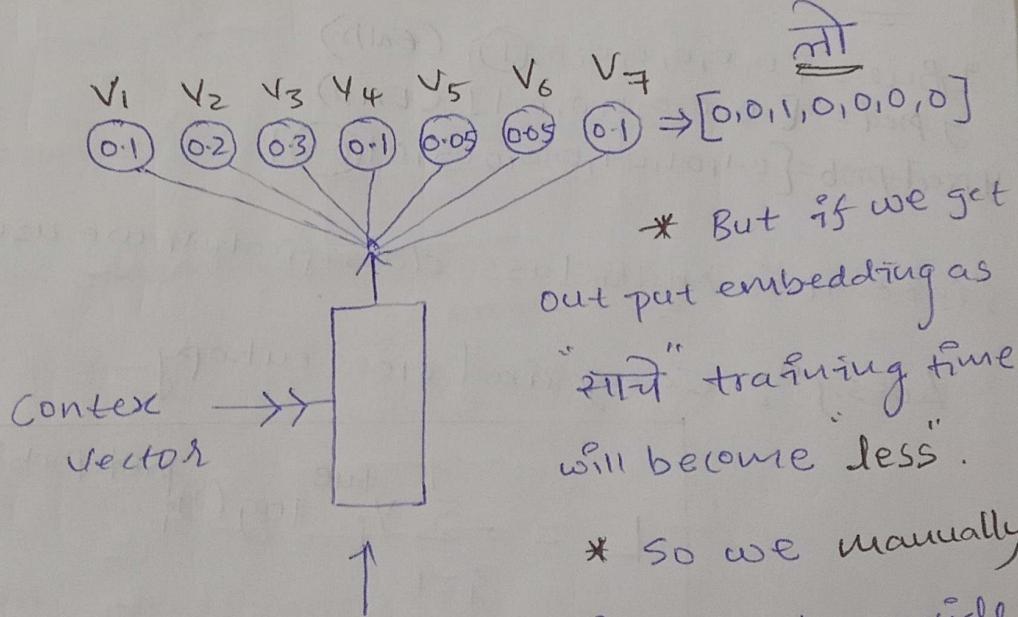
Taking row-2: [Think about it → सोच लो]

* Initially as we know weights and biases are of LSTM are randomly initialized [embedding]



(softmax) → It contains neurons, no. of neurons = V_{out}

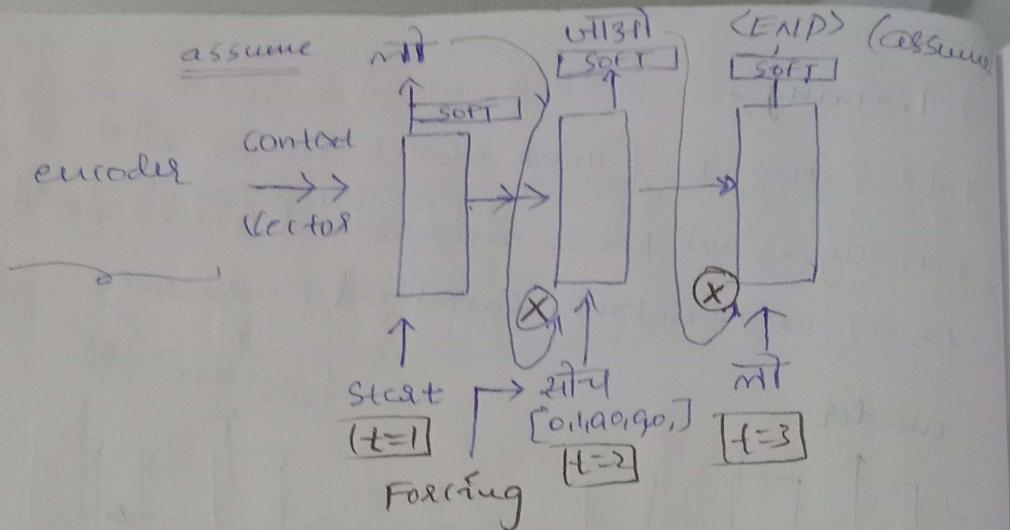
For eg:-



* So we manually force and provide

the input embedding vector as सोच instead मी; for next lstm, this is known as **Teacher forcing**.

⇒ This makes the convergence fast.



* Now forward propagation is done!! , now we have to calculate "loss".

$$\begin{array}{l}
 y_{\text{true}} = [0, 1, 0, 0, 0, 0, 0] \quad (\text{रोचे}) \\
 y_{\text{pred}} = [0, 0, 1, 0, 0, 0, 0] \quad (\text{मी}) \\
 y_{\text{pred-prob}} = [0, 1, 0.2, 0.3, 0.1, 0.05, 0.05, 0.1]
 \end{array}
 \quad \left| \quad \begin{array}{l}
 y_{\text{true}} = [0, 0, 1, 0, 0, 0, 0] \quad (\text{मी}) \\
 y_{\text{pred}} = [0, 0, 0, 0, 0, 0, 1, 0] \quad (\text{जांडा}) \\
 y_{\text{pred-prob}} = [0, 1, 0.05, 0.05, 0.1, 0.2, 0.4, 0.1]
 \end{array} \right.$$

Since it is multiclass classification we use

Loss function: Categorical cross entropy

$$L = - \sum_{i=1}^7 y_i^{\text{true}} \log(y_i^{\text{pred_prob}})$$

$$L_{t=1} = (-1) \log(0.2) = 0.69$$

$$L_{t=3} = (-)(1)(0.05) = (-)(1)\log(0.05) = 1.3$$

$$L_{t=3} = (-1)(0.5) = -(-1)(1) \log(0.5) = 0.3$$

$$L_T = 2.29$$

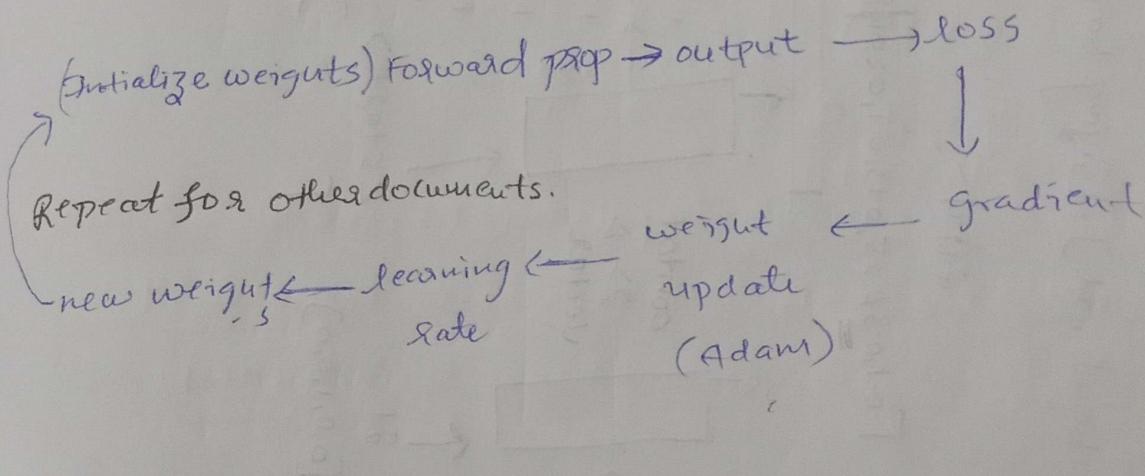
* Now we know loss, what next?

Back propagation:

- Two steps:-
- 1) calculating gradient
 - 2) update parameters

We have to take the gradient of all the weights & biases of LSTM & SOFT MAX Layer w.r.t LOSS
Then we update them using adam or rms-prop
etc optimizers.

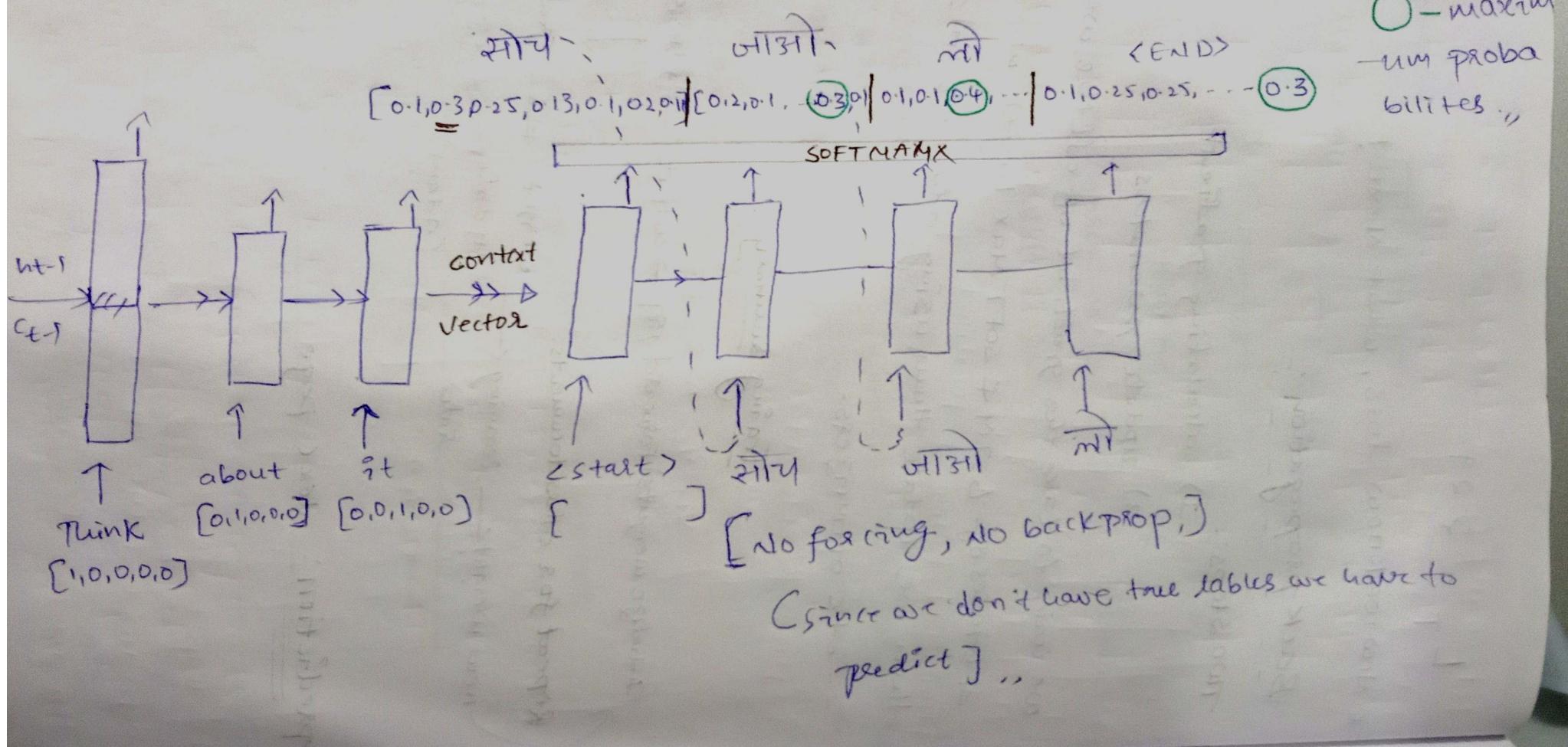
Training summary:



Prediction: next page.

Prediction: <start> → ① शोध ② ले → ③, संदर्भ → ④, तांत्रिक ⑤, जानकी → ⑥, <end> → ⑦

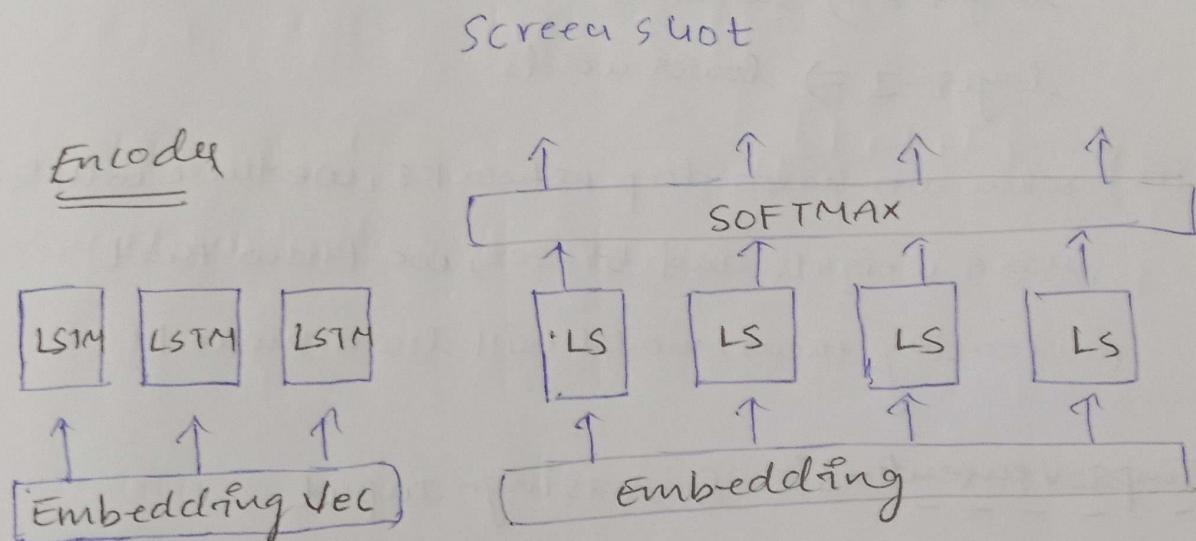
Input: [Think about it] - {giving back data document itself}



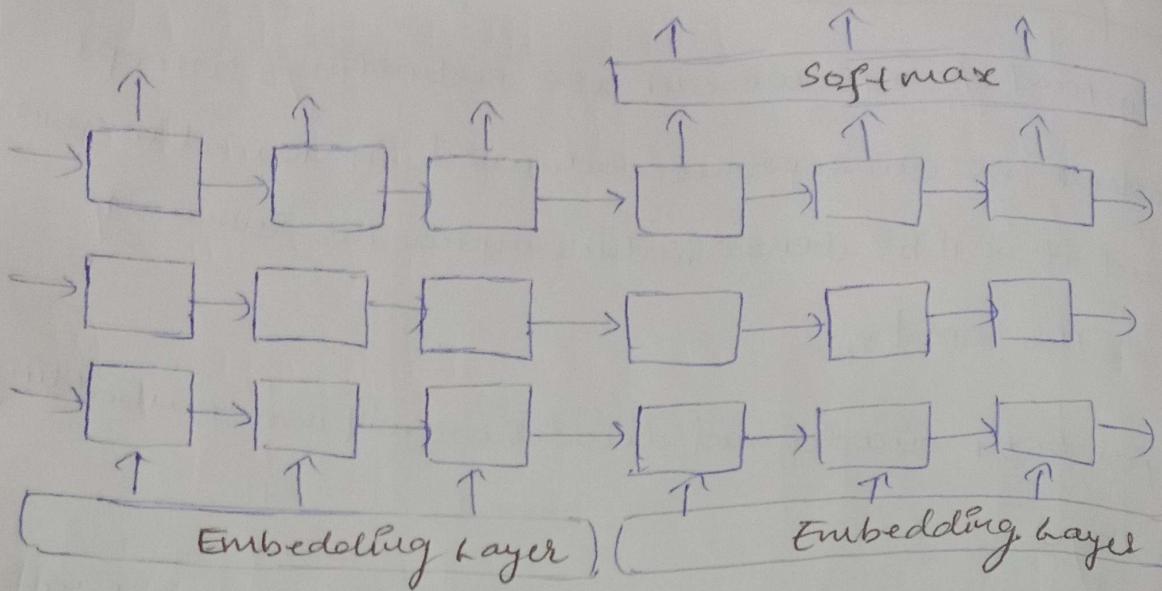
Improvement-1: Embeddings

- * Instead of OHE we can use embeddings concept where we can represent each word in low-dimension and it will be dense, which capture the summary of the word.
- * For both encoder and decoder we will use embedding layer.
- * We can use pretrained embedding techniques like word2vec, glove or custom embedding on training the data

Final architecture :-



Improvement-2 - Deep LSTMs



3 Benefits

i) Long term dependencies

ii) layered representation:

layer-3 :- paragraph

layer-2 \Rightarrow sentences

layer-1 \Rightarrow lower words

iii) when we have deep networks, we will have more weights and biases, we know that increasing them, model will learn more !!

Improvement-3:- Reversing Input

* when we reverse the input the starting words will be close to the corresponding output

* It will not work always, there are some language pairs, like english-french, etc, it will not

where initial words should have more context
than last words!