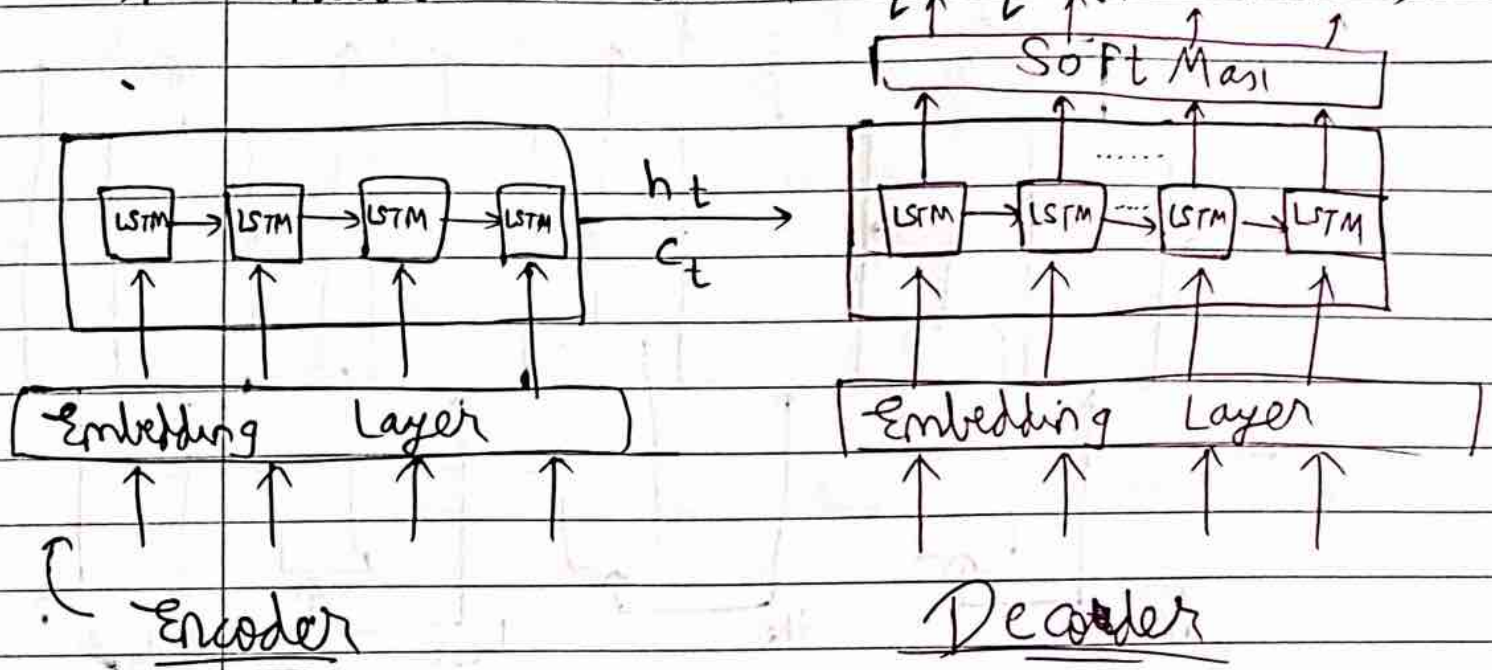


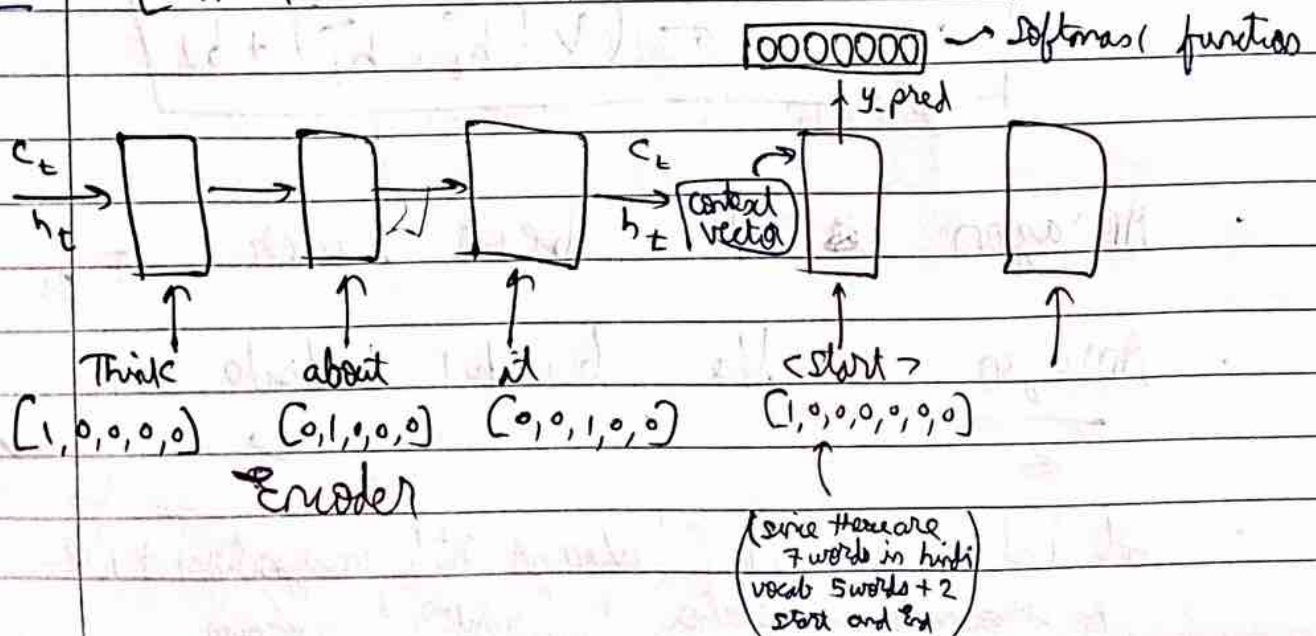
* Encoder-Decoder (Seq 2 seq Architecture)



* Training the Architecture using Back propagation with an example.

1) English 2) Hindi 3) <start>
 1) (think about it) (सोच लो) 7) <end>
 2) (come in) (अंदर आ जाओ)

row 1 - [Think about it → सोच लो]



highest \rightarrow (ली) \rightarrow word

Let's say $y_{pred} = [0.2, 0.1, 0.3, 0.2, 0.1, 0.1, 0]$

$y_{true} = [0, 1, 0, 0, 0, 0, 0]$
 \hookrightarrow सोच

~~X~~

Note :-

- Our output is (ली) word instead of (सोच) word
- In Decoder, the output of previous cell is given as input to its next cell at time for timestamp $(t+1)$.
- But, while training process, even if the previous cell calculates its output (correct or incorrect), we would still send 'True labelled value' to the current cell.
- Eg:- We got out as (ली) but correct word was (सोच), so we will give (सोच) as our new input to the next cell and not (ली) as our input. This is called teacher forcing method.
- This helps in faster training and fast convergence to the true value.

*

Calculating Loss and performing back propagation

Eg:-

	$\begin{matrix} \text{padding} \\ \text{to match} \\ \text{output dim} \end{matrix}$	Loss	$\begin{matrix} \text{padding} \\ \text{to match} \\ \text{output dim} \end{matrix}$
y_{true}	$[0, 1, 0, 0, 0, 0, 0]$ सोच	$[0, 0, 1, 0, 0, 0, 0]$ (ली)	$[0, 0, 0, 0, 0, 0, 1]$ <end>
y_{pred}	$[0.2, 0.1, 0.3, 0.2, 0.1, 0.1, 0]$ ली	$[0.1, 0.0, 0.1, 0.15, 0.2, 0.4, 0.05]$ जीत	$[0.1, 0.1, 0.1, 0.2, 0.2, 0.3]$ <end>

we would use Crossentropy Loss () function

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

$$L_{t=1} = -1 \times \log 0.1 - 0 \times \log 0.2 - 0 \times \log 0.3 \dots 0 \times \log 0$$

$$L_{t=1} = -1 \times \log 0.1 = 1 \quad \therefore L_{t=1} = 1$$

$$L_{t=2} = -1 \times \log 0.1 = 1$$

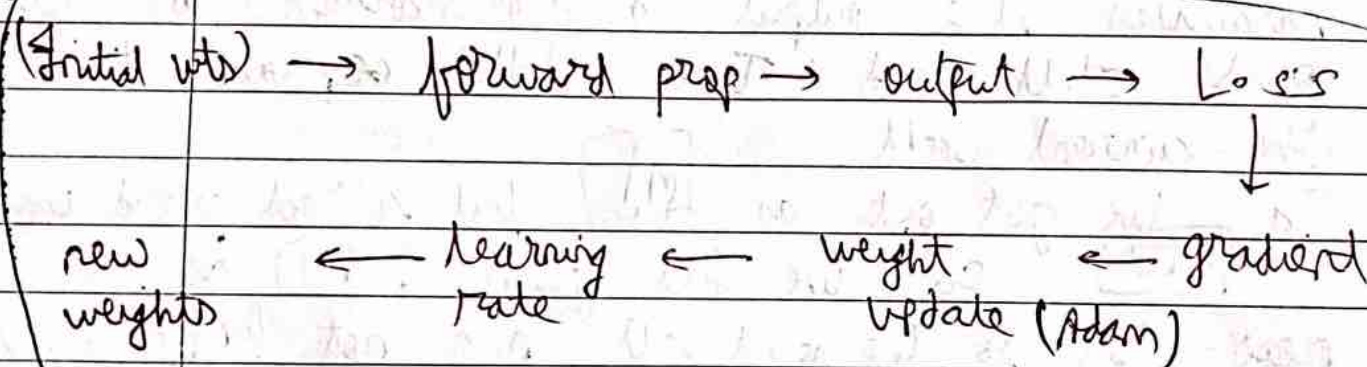
$$L_{t=3} = -1 \times \log 0.3 = 0.31 \text{ (When output is correct Loss } \downarrow \text{)}$$

$$\therefore \text{Total Loss} = \boxed{2.31}$$

$$\therefore \text{average Loss} = \boxed{0.75}$$

Now we will calculate gradients based on Loss function and update weights

✕



✕ Note:-

1) While making prediction we will provide input to current cell is "decoder" as output of previous cell of decoder unlike is training phase (which is obvious).

* Improvements

- 1) Embedding - We won't use one-hot encoding, instead use embedding layer which is a low dimensional dense vector layer, using pretrained (~~word~~ word2vec / Glove)
- 2) Use Deep LSTMs :- Instead of using single LSTM layers use deep LSTM layers, stacked upon each other. It helps in :-
 - i) long term dependency (as it has many context vector unlike using single LSTM layer)
 - ii) Layered represent (It helps in ~~understanding~~ hierarchical understanding of paragraphs i.e. from words \rightarrow sentences \rightarrow paragraphs)
 - iii) If parameters are more (when used deep LSTM), it ensures good accuracy
- 3) Reversing the input :- In encoder, we can provide input in reverse direction too, to improve model performance (but it works in very few cases only.)

Input \rightarrow Think about it
Reverse \rightarrow it about Think

* Encoder

- Purpose :- Understand the input sequence and compress its meaning
- How :-
 - 1) Takes input (e.g. an English sentence)
 - 2) Processes it word by word (using RNN/LSTM)
 - 3) Converts the entire input sequence into a context vector (a fixed size representation)

* Decoder

- Purpose :- Generate the output sequence based on the encoder's context
- How :-
 - 1) Takes the context vector from the encoder as ^{input}
 - 2) Predicts the output sequence one token at a time (e.g. in French)
 - 3) Each prediction is based on the context and previously generated tokens

* Example :-

Input :- "I am happy"

Encoder :- Understands and encodes this into a context vector

Decoder :- Uses this vector to generate output like
"Je suis heureux" (French translation)