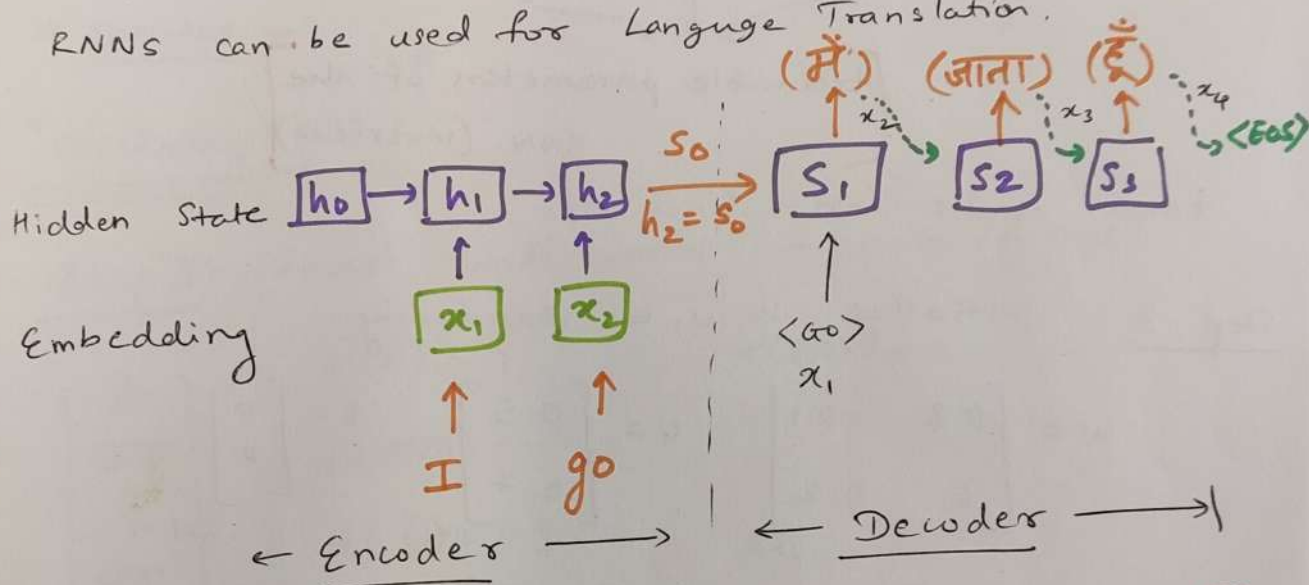


RNNs

Lecture No. 3

"Build RNN from scratch"

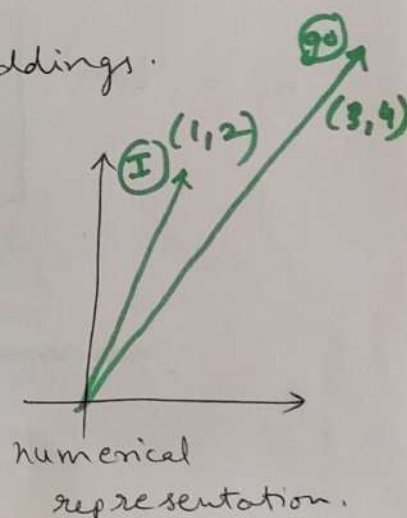
RNNs can be used for Language Translation.



Step 1: Convert I/P words into embeddings.

$$x(I) = x_1 = (1)$$

$$x(go) = x_2 = (2)$$



assume every word is in 1D.

Step 2: Decide hidden state size.

assume hidden state ^{size} = (2)

$$h_0 = \begin{bmatrix} 0 \\ 0 \end{bmatrix}_{2 \times 1}$$

Mathematical relation b/w hidden states;

$$h_1 = \tanh(W \cdot h_0 + U \cdot x_1 + b)$$

(2x1)

$$h_2 = \tanh(W \cdot h_1 + U \cdot x_2 + b)$$

(2x1)

[trainable parameters of the RNN. (matrices)]

\tanh adds non-linearity.

Step 3: initialize w, u, b randomly.

$$W = \begin{bmatrix} 0.3 & -0.1 \\ 0 & 0.2 \end{bmatrix}_{2 \times 2} \quad U = \begin{bmatrix} 0.5 \\ 0.7 \end{bmatrix}_{2 \times 1} \quad b = \begin{bmatrix} 0 \\ 0 \end{bmatrix}_{2 \times 1}$$

$$h_1 = \tanh \begin{bmatrix} 0.5 \\ 0.7 \end{bmatrix} = \begin{bmatrix} 0.46 \\ 0.60 \end{bmatrix}_{2 \times 1}$$

$$h_2 = \tanh \left\{ \begin{bmatrix} 0.3 & -0.1 \\ 0 & 0.2 \end{bmatrix} \begin{bmatrix} 0.46 \\ 0.6 \end{bmatrix} + \begin{bmatrix} 0.5 \\ 0.7 \end{bmatrix} 2 + \begin{bmatrix} 0 \\ 0 \end{bmatrix} \right\}$$

$$= \tanh \begin{bmatrix} 0.078 + 1 \\ 1.52 \end{bmatrix} = \tanh \begin{bmatrix} 1.078 \\ 1.52 \end{bmatrix} = \begin{bmatrix} 0.79 \\ 0.90 \end{bmatrix}$$

memory trans.

"passed to decoder"

⊙ that is why the pressure solely lies on h_2 in RNN.

<Go> → some point of start.

* Decoder -

"Vocabulary"

Step 1: Convert Hindi words into embeddings.

$$\begin{bmatrix} मैं \\ जाता \\ हूँ \\ <Eos> \end{bmatrix} \quad \begin{aligned} x(\text{मैं}) &= 1 \\ x(\text{जाता}) &= 1.1 \\ x(\text{हूँ}) &= 0.9 \end{aligned} \quad x(\text{Eos}) = 0.0$$

$$s_0 = h_2 = \begin{bmatrix} 0.79 \\ 0.90 \end{bmatrix}$$

$$s_t = \tanh \left(W_{\text{dec}} \cdot s_{t-1} + U_{\text{dec}} \cdot \underset{\text{input}}{\underline{x_t}} + b_{\text{dec}} \right)$$

* o/p of s_1 is I/P of s_2 , o/p of s_2 is I/P to $s_3 \dots$

Step 3: Calculate decoder hidden states.

$$s_1 = \tanh \left\{ w_{dec} \begin{bmatrix} 0.7 \\ 0.9 \end{bmatrix} + v_{dec} \cdot 0.5 + b_{dec} \right\}$$

$x < 60 = 0.5$ (chosen random)

w_{dec} , v_{dec} & $b_{dec} \rightarrow$ randomly initialized.

$$s_1 = \begin{bmatrix} 0.27 \\ 0.58 \end{bmatrix}_{2 \times 1} \quad \text{--- ①}$$

$(4 \times 2) (2 \times 1) = 4 \times 1$ matrix.

take this 2×1 matrix and project into 4×1 matrix.
~~st~~ probability matrix. because voc. is 4×1 matrix.

Step 4: Output matrix / Logits matrix.

$$\text{logits}_t = \text{w}_{out} \cdot s_t + \text{b}_{out} \quad \text{① trainable para.}$$

$$w_{out} = \begin{bmatrix} 0.2 & 0.1 \\ 0 & 0.2 \\ -0.1 & -0.2 \\ -0.2 & -0.1 \end{bmatrix}_{4 \times 2} \quad b_{out} = \begin{bmatrix} 0 \\ 0 \\ 0 \\ 0 \end{bmatrix}_{4 \times 1}$$

$$\text{logits}_t = \begin{bmatrix} 0.2 & 0.1 \\ 0 & 0.2 \\ -0.1 & -0.2 \\ -0.2 & -0.1 \end{bmatrix}_{4 \times 2} \begin{bmatrix} 0.27 \\ 0.58 \end{bmatrix}_{2 \times 1} = \begin{bmatrix} 0.11 \\ 0.12 \\ -0.14 \\ -0.03 \end{bmatrix}_{4 \times 1}$$

Logits matrix indicates what translated word is.

~~We~~ But we somehow need to turn these values into
"probability".

• We make use of softmax function, to give higher probabilities to higher values & penalize the lower values. other normalization functions cannot do that. We cannot be completely confident with predictions. Ⓢ

① Positive numbers.

② Sum = 1

③ Get more confidence in our predictions.

} Ⓢ imp.

Step 5:

$$\begin{bmatrix} 0.11 \\ 0.12 \\ -0.14 \\ -0.03 \end{bmatrix} \Rightarrow \text{softmax} \Rightarrow \begin{bmatrix} 0.275 \\ 0.276 \\ 0.212 \\ 0.238 \end{bmatrix} \quad \checkmark \text{ Highest probability (जाता)}$$

Next word prediction : "जाता"

Step 6 : Continue decoding -

$$x_2 = \text{"जाता"} \text{ embedding} = 1.1$$

$$s_2 = \tanh(w_{dec} \cdot s_1 + u_{dec} \cdot x_2 + b_{dec})$$

• Trainable parameters.

$$= \tanh \left(\underline{w_{dec}} \cdot \begin{bmatrix} 0.27 \\ 0.58 \end{bmatrix} + \underline{u_{dec}} * 1.1 + \underline{b_{dec}} \right)$$

$$= \begin{bmatrix} 0.21 \\ 0.49 \end{bmatrix} \rightarrow \text{softmax} \Rightarrow \begin{bmatrix} 0.0927 \\ 0.0978 \\ -0.12 \\ -0.02 \end{bmatrix} \quad \checkmark \underline{\text{जाता}}$$

जाता is predicted again!!

① word 2 vec \rightarrow ~~conv~~ takes in words & converts to vector tokens embeddings.



for. eg: $\boxed{\text{मैं}} \rightarrow \boxed{0.5}$ $\boxed{\text{जाता}} \rightarrow \boxed{1.1}$

or one can train them separately.
in RNNs both can be done.

② RNNs are very specific to tasks.

whereas LMs can do anything & everything together

All the trainable parameters are trained during back propagation using "cross entropy loss" (X).

* Parameters marked in orange are trainable parameters.

"Code RNN from Scratch"

Encoder block → main aim is to find h_1 & h_2 hidden state.

$$h_t = \tanh \left(\underset{\textcircled{1}}{W} \cdot h_{t-1} + \underset{\textcircled{2}}{U} \cdot x_t + \underset{\textcircled{3}}{b} \right)$$

nn. parameter → (trainable) tell pytorch they are trainable
torch.randn → initialize randomly.

(hidden size, hidden size) → size of dim. of W matrix.

self. embedding → easy finding of embedding value for a given token.

For decoder -

$$s_n = \tanh \left(\underset{\textcircled{1}}{W_{\text{dec}}} s_{n-1} + \underset{\textcircled{2}}{U_{\text{dec}}} x_n + \underset{\textcircled{3}}{b_{\text{dec}}} \right)$$

$$\text{logit}_t = \underset{\textcircled{4}}{W_{\text{out}}} s_t + \underset{\textcircled{5}}{b_{\text{out}}}$$

vocabulary size

no. of hidden state size.

⑤ Matrix of probability.

LSTM – Long Short Term Memory

Lecture 5-

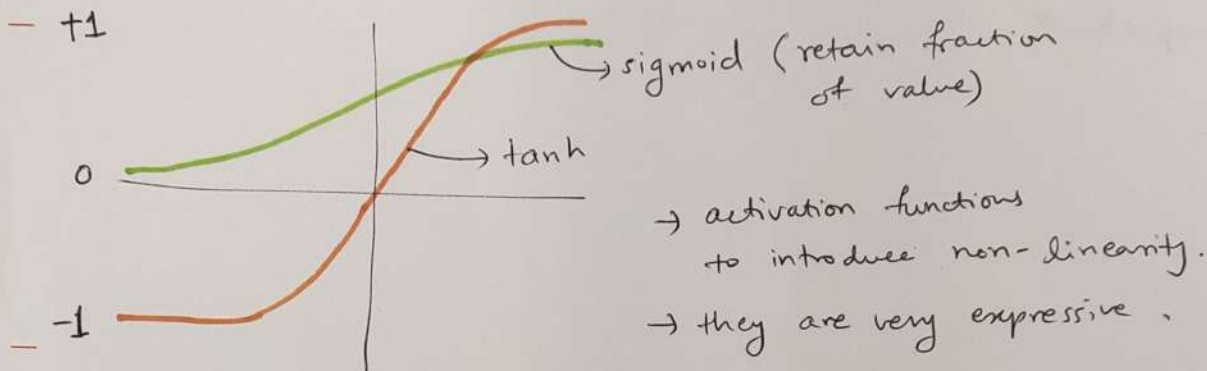
(1997) Long Short term Memory

Neural Networks + Memory = RNN.



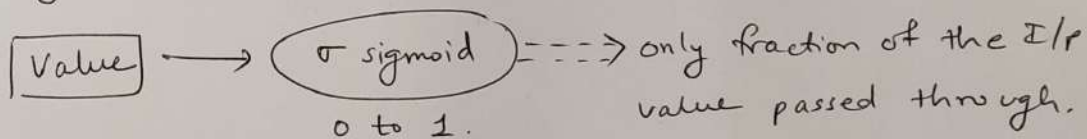
→ Huge info \Rightarrow cannot be taken by RNN.
(long texts cannot be computed by RNNs).

LSTM \rightarrow solved both issues of RNNs.



when values are want to be pressed between 0 & 1
use sigmoid

① Gating mechanism:



② cell state: (missing in RNNs)

Long memory aspect of LSTM.
carried info from way back.

consider it like
a conveyor
belt.
carrying info
from past

we decide what to keep at each stage of cell state & what to remove.

3 Gates

① Forget Gate

what info to throw away from the old cell state.

Till now, input was only Harry (male).

Now, Hermione enters (Female)

so imp to Harry should diminish.

Harry \rightarrow Gate \Rightarrow imp to Harry reduces by some factor
 C_{t-1} f_t

② Input gate $C_{\{t\}}$

Decides how much new info to let inside cell state.

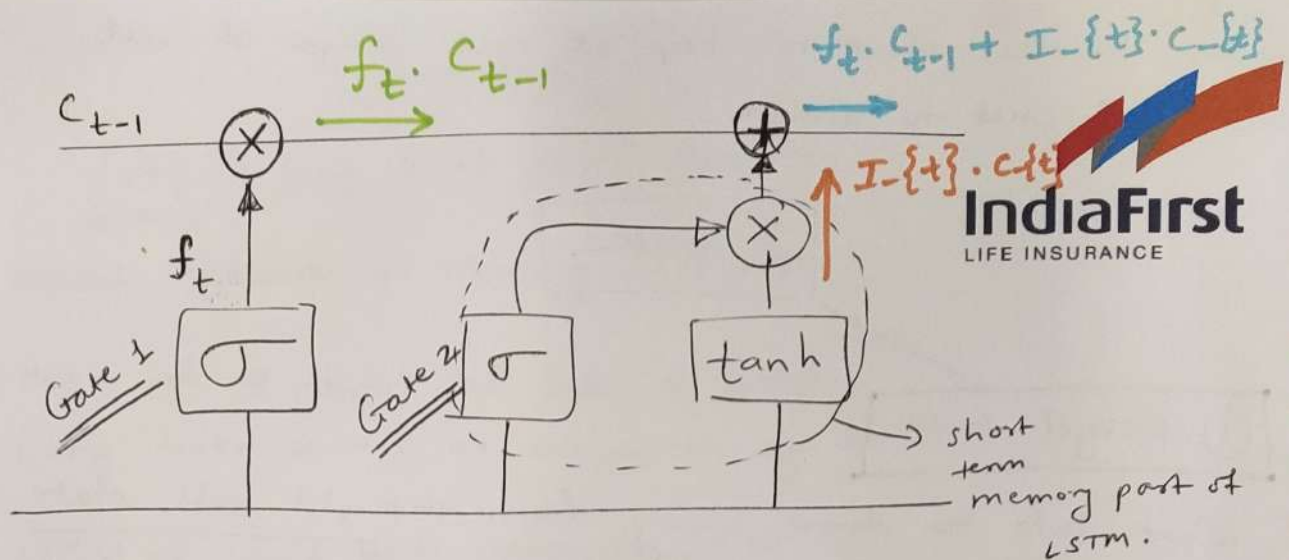
\hookrightarrow let's say about (Hermione)

decided by i/p gate factor $I_{\{t\}}$

③ New cell state $C_t = \underbrace{f_t}_{\text{①}} \cdot C_{t-1} + \underbrace{I_{\{t\}}}_{\text{0-1 (sigmoid)}} \cdot \underbrace{C_{\{t\}}}_{\text{candidate vector}}$

How much of Harry to retain/forget

How much of Hermione to add.



c_{t+1} can lie between -1 & $+1$

③ I_{t+1} decides how much of c_{t+1} should be let in.

④ Update the cell state using the derived formula in ③

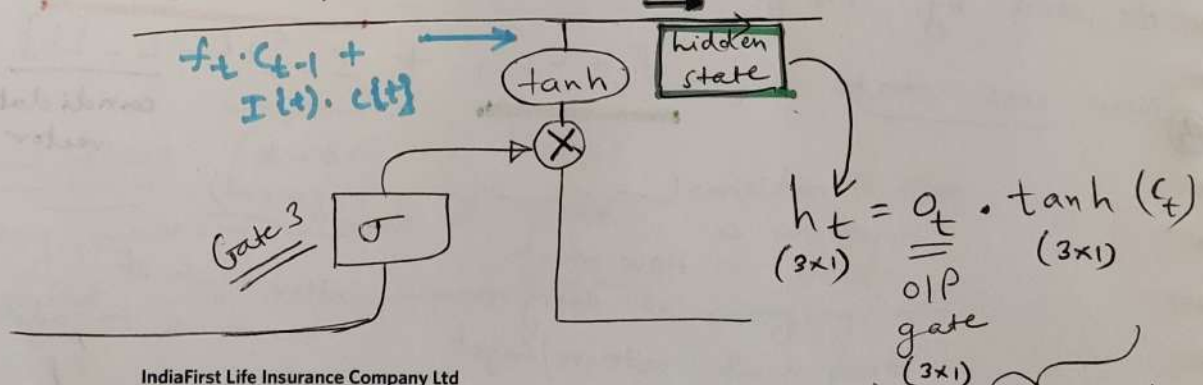
⑤ Output gate — lies btw 0 & 1

O/P = ~~cell state~~ filtered version of cell state.

→ only some part will be displayed and this is controlled through o/p gate.

hidden state ultimately predicts the next word.

↳ what part to retain/remove.



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

(3×1) (3×1) (3×1)

not matrix multi

but element wise multiplication.

$$C_t = f_t * C_{t-1} + I_{\{t\}} * C_{\{t\}}$$

(3×1) (3×1) (3×1) (3×1)

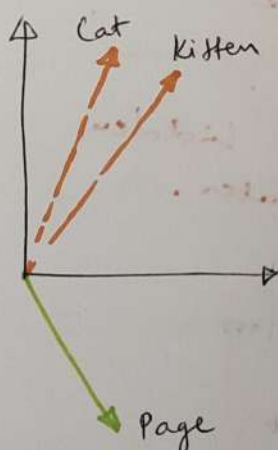
Forget gate prev. cell state I/P gate sigmoid function (value to lie btw 0 & 1) tanh.

W - Weight projection matrix

U - input projection matrix.

takes input and converts it into hidden state dimension.

more dimensions we have for words to embedding, the more context we can capture from them



* For LSTMs,

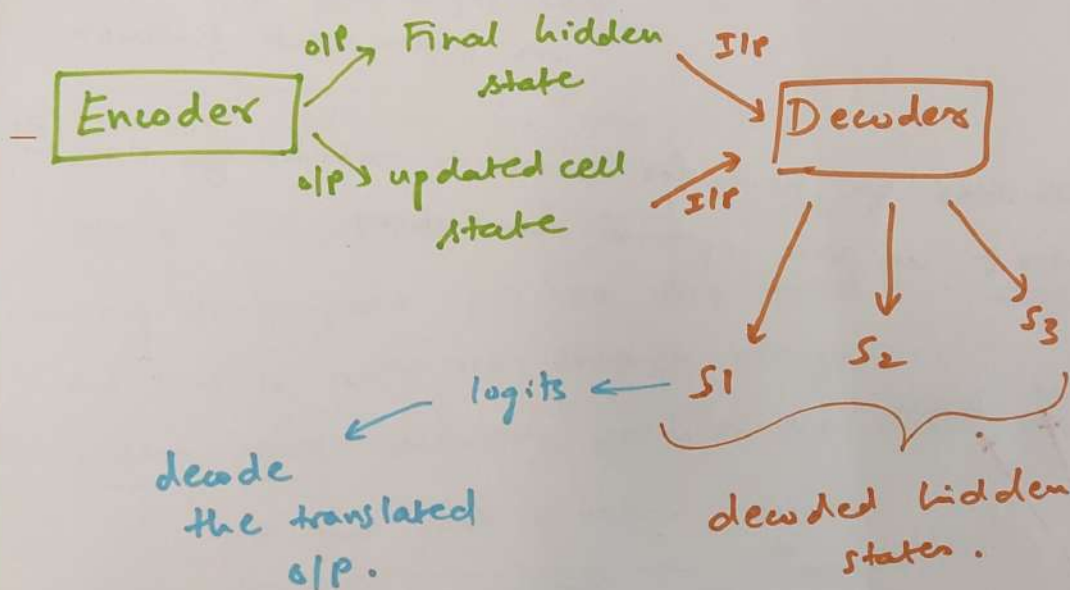
two things are passed to the decoder.
the hidden state and cell state.

whereas in RNNs, only hidden state was passed to decoder.

12 matrices \rightarrow encoder
12 matrices \rightarrow decoder
2: logits

} LSTM.
translator

all parameters optimized through opti back propagation.



Final o/p : मैं मैं मैं
wanted o/p : मैं जाना हूँ

} minimized loss
through
Gradient
Descent.

LSTM for text generation -

We are studying Data Science ?

① IIP (Predictors)

We
[1]

We are
[1 2]

We are studying
[1 2 3]

We are studying
Data

[1 2 3 4]

Target (Labels)

are
[2] → vector

studying
3

Data
[4]

Science
[7]

- ⊙ Autoregressive
- ⊙ model creates it on own

Each IIP - target pair is given to the LSTM.

Step 1: Data Collection.



Step 2: Create vocabulary of words → tokenization. (token IDs)
word level "the" → 1 "Texas" = 3
"dog" → 2

Step 3: IIP target pairs for each headline.

Step 4: LSTM architecture

- convert token into vector embedding.
- Find the hidden state (pass it through architecture)
- logit matrix (Dense NN)

Step 5: Loss for each iteration.

Step 6: Use gradient descent optimization.



padding = makes sure that all inputs are of the same size

if context length = 9.

We are studying Data Science

We are [2] [1]

We are studying [2 11] [100]

We are Data st. [2 11 100] [60]

We are Science st. Data [2 11 100 60] [90]

predictors

0	0	0	0	2	11	100	60	90
(1)	(2)	(3)	(4)	(5)	(6)	(7)	(8)	(9)

Label / Target state

No. of hidden states = Context size.

4/5/20
50/100