

Neural Topic: LSTM

neural
network
layer

pointwise
operation

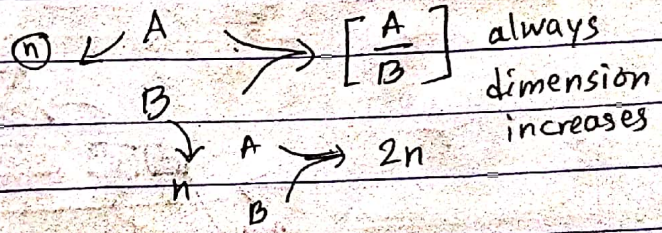
vector
transfer

concatenate

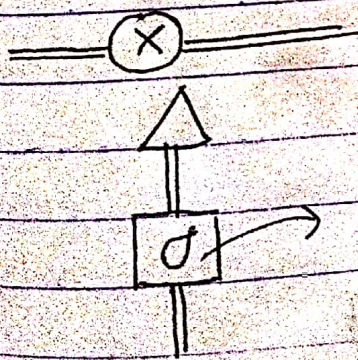
copy

example

$A \oplus B = C$



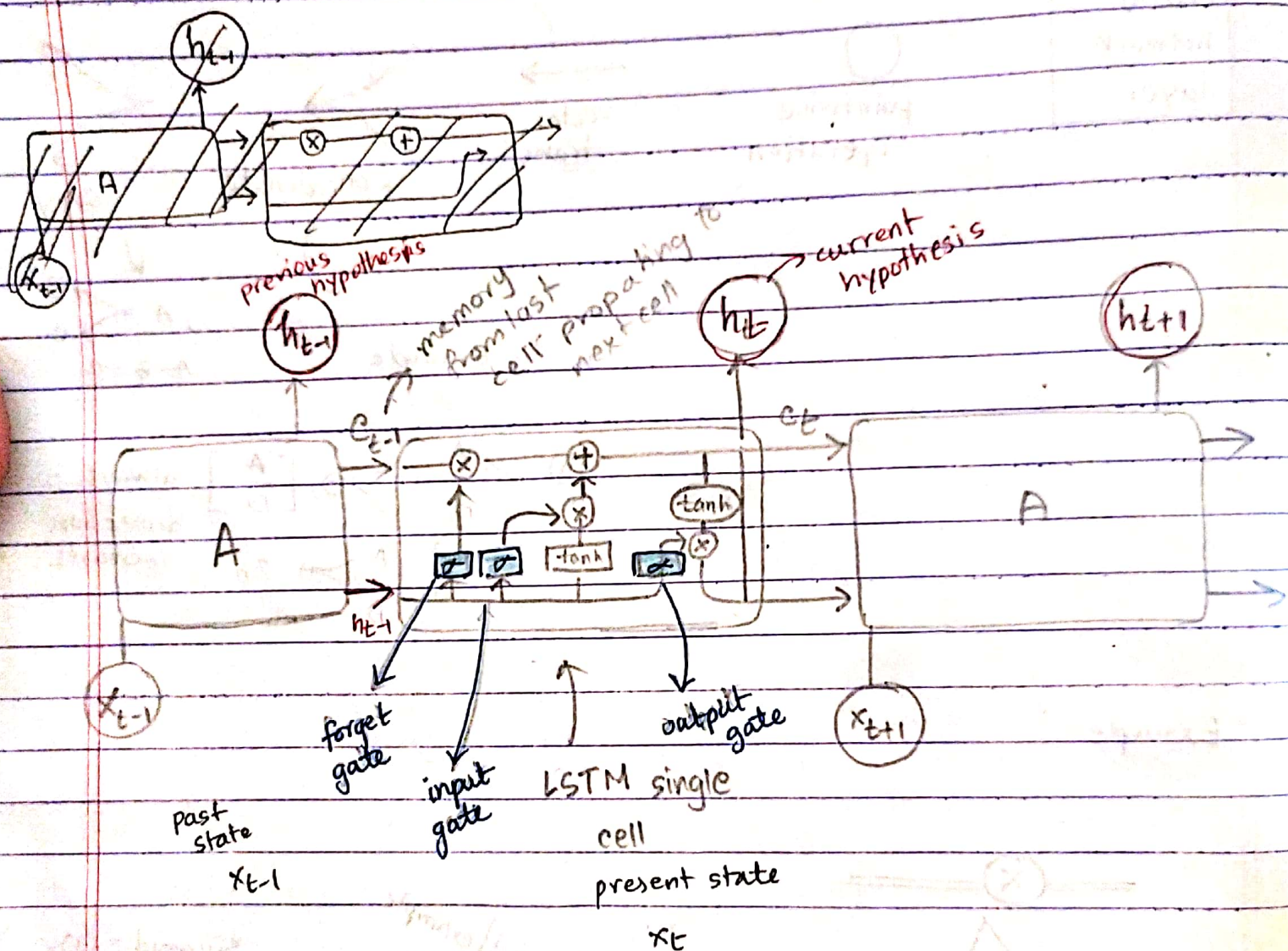
Example



activation
function [using sigmoid] example

sigmoid = (0~1)
acts like a
gate designed
by the activation
function of the
neural network

Example



Here are 4 neural network gates and 3 ~~act~~ ^{act} of them are σ (sigmoid). We will represent them using 3 names.

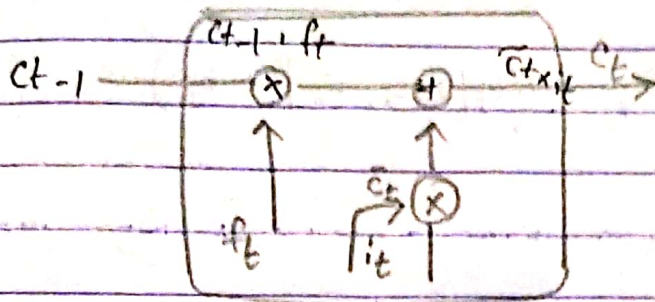
(a) forget gate

(b) input gate (jei info ashtese tar koto percent info ami current cell e hoi korbo)

(c) output gate

koto tuku output and present cell theke next cell e propagate korbo.

~~tanh~~
 $\tanh \rightarrow$ enhances or shrinks current info.



$$c_t = f_t * c_{t-1} + i_t * \tilde{c}_t$$

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

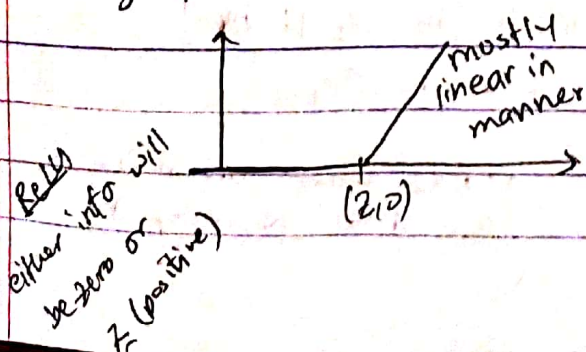
$$i_t = \sigma(w_i[h_{t-1}, x_t] + b_i)$$

$$\tilde{c}_t = \tanh(w_{info}[h_{t-1}, x_t] + b_{info})$$

↓
 updated/enhanced information

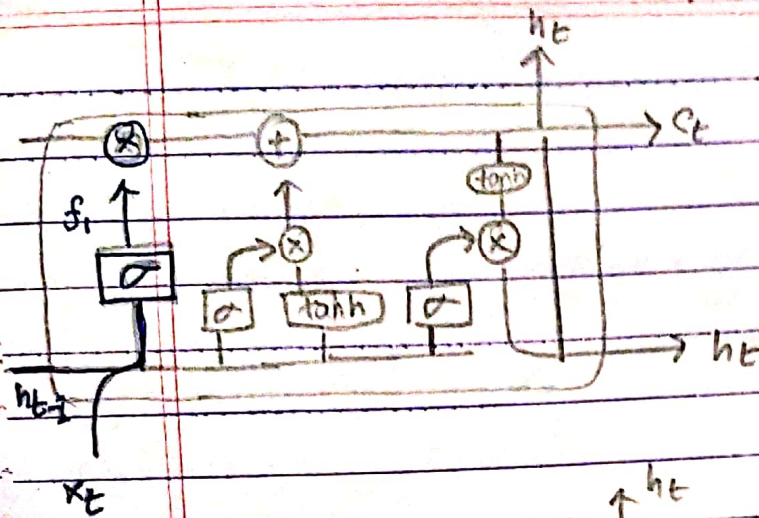
What if we use ~~RE~~ ReLU replaces tanh?

tanh has non linear characteristics as a result it can shrink and enhance. But ReLU has properties such as the graph shown below:

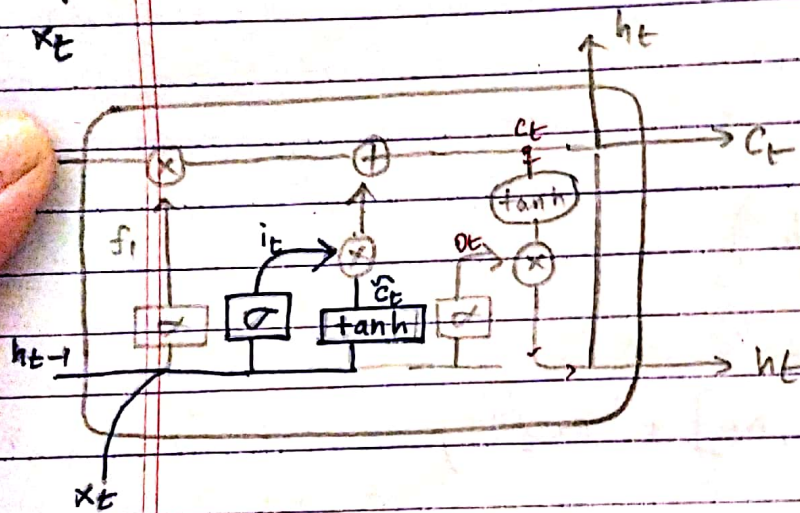


so it cannot inf ~~prop~~ enhance and shrink info as fast so in LSTM structure we need tanh.

Summary



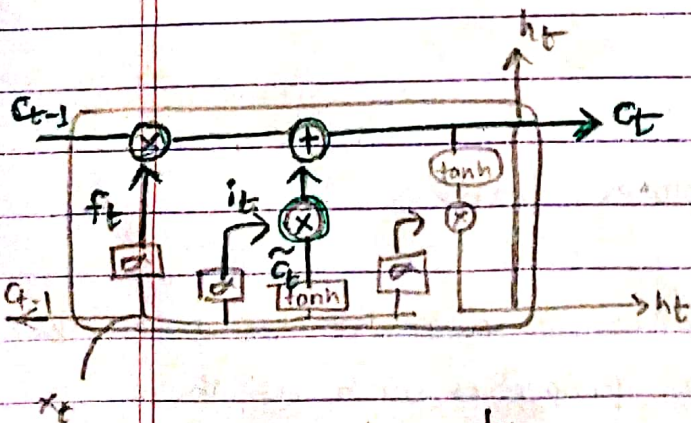
$$f_t = \sigma(w_f[h_{t-1}, x_t] + b_f)$$



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

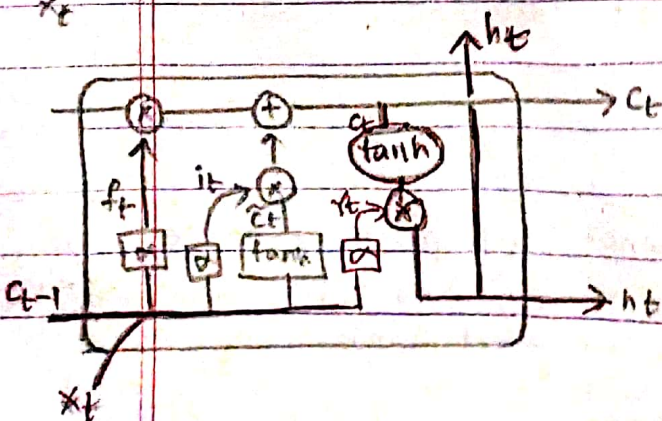
$$\tilde{C}_t = \tanh(w_c[h_{t-1}, x_t] + b_c)$$

i_t decides what component is to be updated. \tilde{C}_t provides change contents.



$$C_t = f_t * C_{t-1} + i_t * \tilde{C}_t$$

updating the cell state



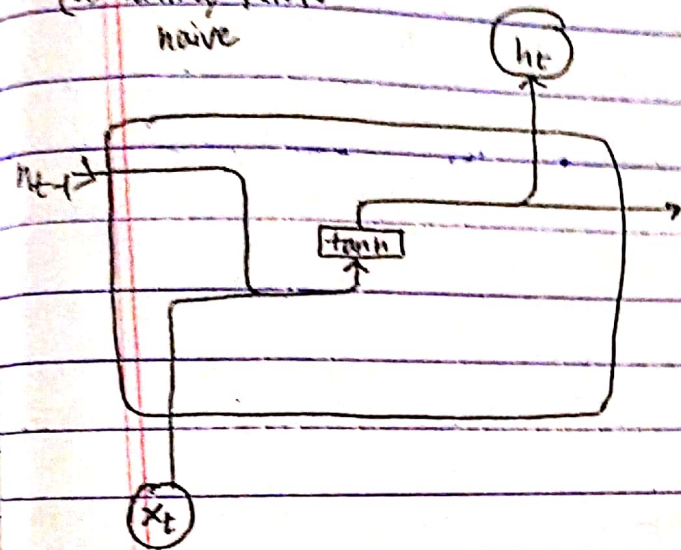
$$y_t = \sigma(w_o[h_{t-1}, x_t] + b_o)$$

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

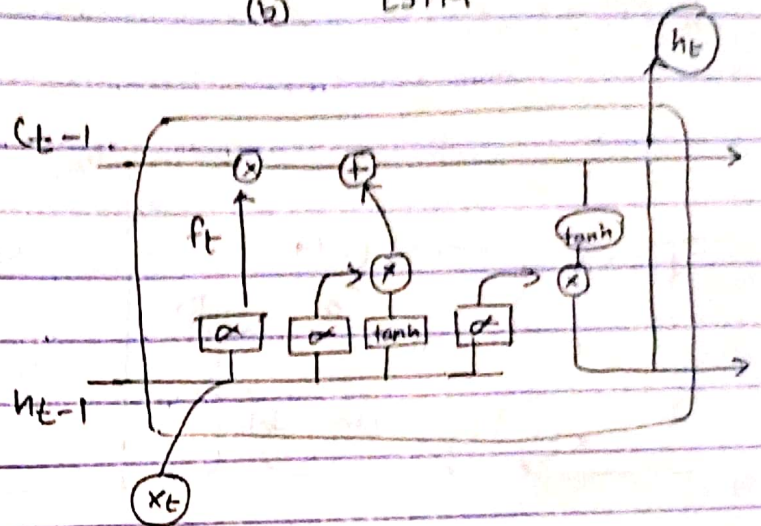
decide what part of the cell state to output

RNN VS LSTM

(a) vanilla/RNN
naïve



(b) LSTM



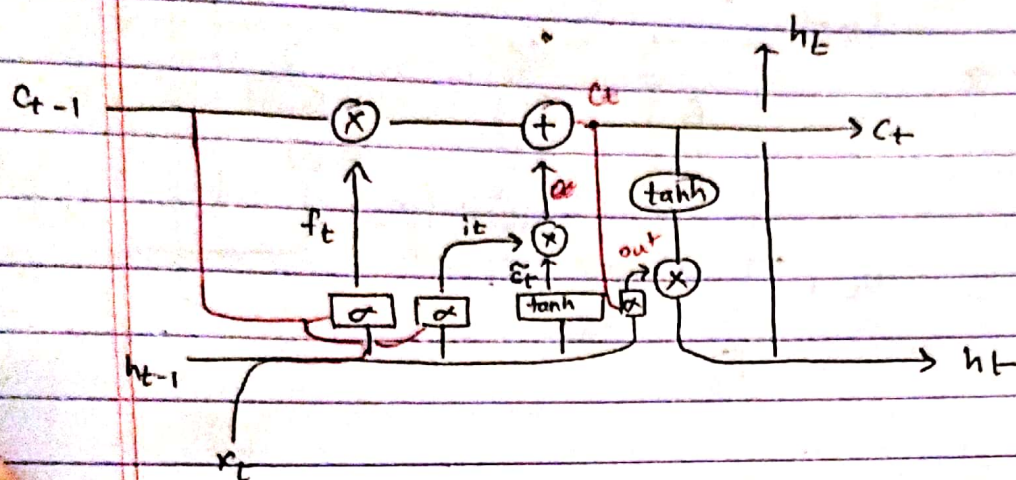
$$h_t = \tanh(W_R [h_{t-1}, x_t] + B_R)$$

Difference

- ① no cell state in RNN.
- ② No gates in RNN, no sigmoid locking key.
LSTM has 3 sigmoid keys \rightarrow i_t, f_t, o_t
- ③ LSTM expensive, can hold longer memory.
prone to vanishing gradient problem.

gates regulates better flow of info

Peephole LSTM



normal LSTM e only (h_{t-1} and x_t) mile sigmoid er box e dhukto but peephole LSTM e tara arekta jinish ke concate kortese that is previous cell state (memory & info)

So (C_{t-1}) copy hole 1st sigmoid then second ~~sig~~ sigmoid e dhukbe.

\therefore basically f_t and i_t er ^{info} C_{t-1} er info dhuktese, but output gate o_t dhuktese
 \downarrow
 current cell info.

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

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

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

$$C_t = f_t * C_{t-1} + i_t * \tilde{C}_t$$

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

c changes slowly $\rightarrow c^t$ is c^{t-1} added by something

jodi cell info onk taratari update hai, taile tar shathe c_{t-1}
and c_t then current cell state onk fast update hobe
so it can hold longer memory.

h h faster $\rightarrow h_t$ and h_{t-1} can be different.