

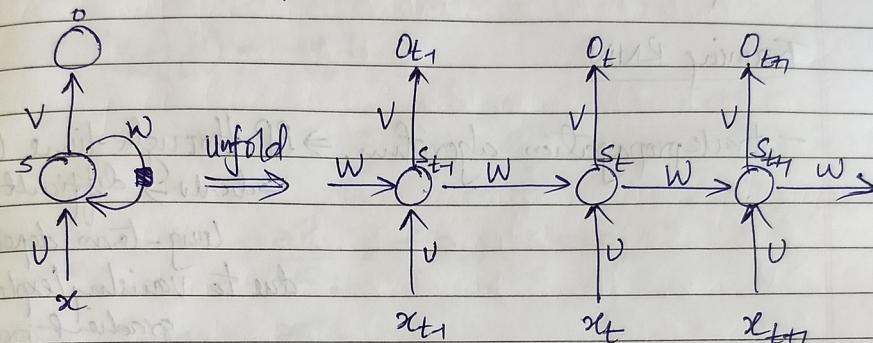
RNN, LSTM, GRU, ConvLSTM

By, Joe Johnson, D17024

Recurrent Neural Networks (RNN)

- make use of sequential information.

recurrent \rightarrow perform the same task for every element of a sequence.



unrolled or unfolded \rightarrow full M/W

\hookrightarrow write out M/W for the complete sequence.
- one layer for each input (each word)

formulas:

- x_t is the i/p at time step t. ($x_t \rightarrow$ one-hot encoding)
- s_t is the hidden state at time step t. (memory)

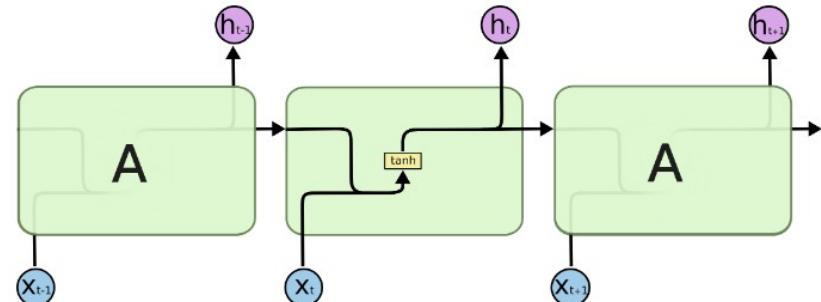
$$s_t = f(v \cdot x_t + W s_{t-1}), \quad f \rightarrow \text{AF (tanh or ReLU)}$$

$s_1 \rightarrow$ first hidden state (initialized to all zeros)

- $o_t = \text{O/p at time step } t. \rightarrow$ vector of prob across O/p
($o_t = \text{softmax}(v \cdot s_t)$)

$s_t \rightarrow$ memory of M/W \rightarrow captures info for all the previous time steps

$o_t \rightarrow$ solely based on the memory at time t's own previous by s_t



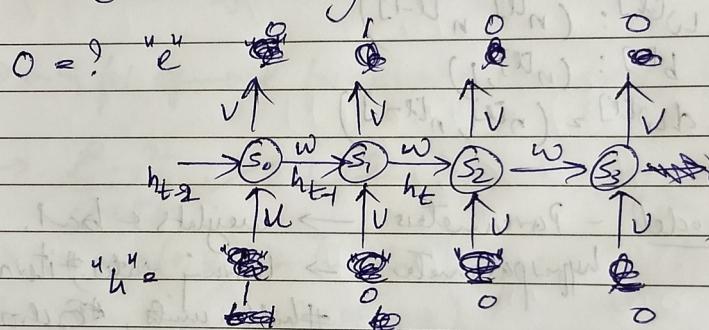
The repeating module in a standard RNN contains a single layer.

Understanding a Recurrent neuron in details.

"hello" $\xrightarrow{\text{task}}$ ~~task~~ $h-e-l-l \rightarrow o$.
 vocabulary \rightarrow 4 letters $\{h, e, l, o\}$

FP in recurrent neuron.

Let. $h = [1 \ 0 \ 0 \ 0]^T$
 $e = [0 \ 1 \ 0 \ 0]^T$
 $l = [0 \ 0 \ 1 \ 0]^T$



$$h_t = f(h_{t-1}, x_t)$$

$$h_t = \tanh(w_{hh} h_{t-1} + w_{xh} x_t)$$

$$y_t = w_{hy} h_t$$

$$x = \begin{bmatrix} "h" & "e" & "l" & "l" \\ 1 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 \\ 0 & 0 & 1 & 1 \\ 0 & 0 & 0 & 0 \end{bmatrix}$$

$x_1 \ x_2 \ x_3 \ x_4$

Step 1:

$$U = w_{xh} = \begin{bmatrix} 0.287 & 0.846 & 0.572 & 0.486 \\ 0.902 & 0.871 & 0.691 & 0.189 \\ 0.5375 & 0.092 & 0.558 & 0.791 \end{bmatrix} \leftarrow \text{randomly generated}$$

$$V \cdot x_t = W_{hh} \cdot x_t = \begin{bmatrix} 0.287 & \dots & \dots \\ 0.902 & \dots & \dots \\ 0.537 & \dots & \dots \end{bmatrix} \begin{bmatrix} 1 \\ 0 \\ 0 \\ 0 \end{bmatrix} = \begin{bmatrix} 0.287 \\ 0.902 \\ 0.537 \end{bmatrix}$$

Step 2 randomly selected.

$$W = W_{hh} = 0.427 \rightarrow, \text{ ab1ign } h_{t-1} = \begin{bmatrix} 0 \\ 0 \\ 0 \end{bmatrix}$$

bias, $b_{hh} = 0.567$

$$W \cdot h_{t-1} + b_{hh} = \begin{bmatrix} 0.567 \\ 0.567 \\ 0.567 \end{bmatrix}$$

Step 3

$$z_t = W \cdot h_{t-1} + V \cdot x_t \rightarrow \begin{bmatrix} 0.287 \\ \vdots \\ 0.287 \end{bmatrix} + \begin{bmatrix} 0.567 \\ \vdots \\ 0.567 \end{bmatrix} = \begin{bmatrix} 0.854 \\ 1.469 \\ 1.04 \end{bmatrix}$$

$$h_t = \tanh(z_t) \rightarrow \begin{bmatrix} 0.693 \\ 0.899 \\ 0.802 \end{bmatrix}$$

Step 4
Repeat from Step 1 to 3 for $X = "e"$ = $\begin{bmatrix} 0 & 1 & 0 & 0 \end{bmatrix}^T$
 $"e" = \begin{bmatrix} 0 & 0 & 1 & 0 \end{bmatrix}^T$

$$h_t = \begin{bmatrix} 0.936 \\ 0.949 \\ 0.762 \end{bmatrix}$$

Step 5

$$y_t = W_{hy} \cdot h_t = V \cdot h_t$$

$$V = W_{hy} = \begin{bmatrix} 0.371 & 0.974 & 0.830 \\ 0.371 & 0.282 & 0.659 \\ 0.649 & 0.098 & 0.334 \\ 0.912 & 0.325 & 0.144 \end{bmatrix} \leftarrow \text{randomly selected}$$

Step 6
Probability for a particular letter \rightarrow Softmax fn.
softmax(y_t)

Unlike traditional DNN
 → RNN share the same parameter (U, V, W) across all steps.
 → Reduce total # of parameters required to learn

hidden state capture some information about a sequence.

Training RNN

- backpropagation algorithm. \rightarrow BP through time (BPTT)
 \hookrightarrow difficulty learning long-term dependencies due to vanishing/exploding gradient problem.

Solutions \rightarrow LSTM.

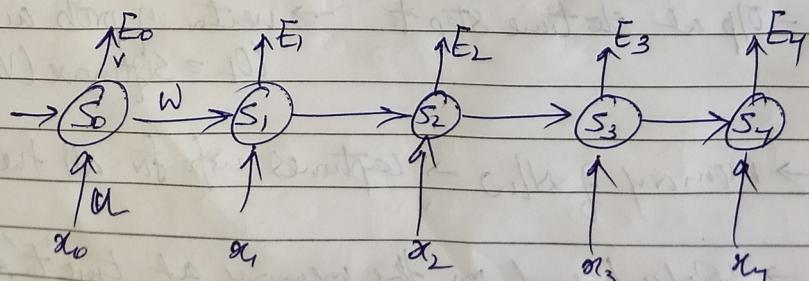
$$s_t = f(Ux_t + Ws_{t-1}) = \tanh(Ux_t + Ws_{t-1})$$

$$y_t = f(Vs_t) = \text{softmax}(V.s_t)$$

loss \Rightarrow Cross entropy loss.

$$E_t(y_t, \hat{y}_t) = -y_t \log \hat{y}_t$$

$$E(y_t, \hat{y}_t) = -\sum_t y_t \log \hat{y}_t$$



BPTT

$$\frac{\partial E_t}{\partial v} = \frac{\partial E_t}{\partial \hat{y}_t} \cdot \frac{\partial \hat{y}_t}{\partial v} = \frac{\partial E_t}{\partial \hat{y}_t} \cdot \frac{\partial \hat{y}_t}{\partial z_t} \cdot \frac{\partial z_t}{\partial v}$$

$$E_t(y_t, \hat{y}_t) = -y_t \log \hat{y}_t$$

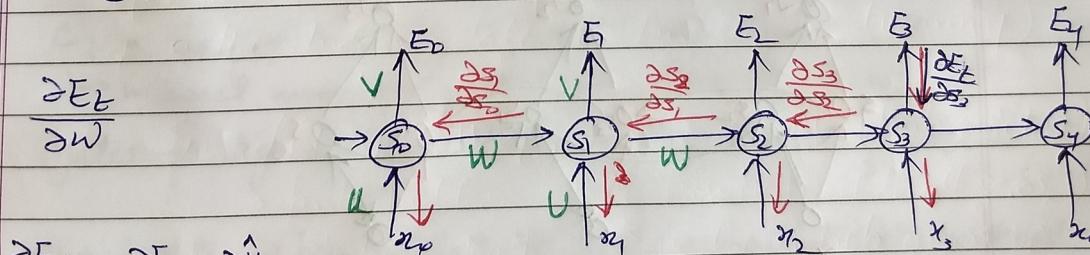
$$\hat{y}_t = \text{softmax}(V.S_t) = \frac{e^{z_t}}{\sum e^{z_t}} \quad \sigma(z_t) \rightarrow \text{softmax}, \quad \sigma'(z_t) = \sigma(z_t)(1 - \sigma(z_t))$$

$$\Rightarrow \frac{\partial z_t}{\partial v} = s_t, \quad \frac{\partial \hat{y}_t}{\partial z_t} = \sigma(z_t)(1 - \sigma(z_t))$$

$$\frac{\partial E_t}{\partial \hat{y}_t} = -\frac{1}{\hat{y}_t}$$

$$\begin{aligned} \frac{\partial E_t}{\partial z_t} &= \frac{\partial (-y_t \log \sigma(z_t))}{\partial z_t} = -y_t \frac{1}{\hat{y}_t} \frac{\partial \hat{y}_t}{\partial z_t} \\ &= \hat{y}_t - y_t \end{aligned}$$

$$\Rightarrow \frac{\partial E_t}{\partial v} = (\hat{y}_t - y_t)(s_t)^T$$



$$\frac{\partial E_t}{\partial w} = \frac{\partial E_t}{\partial \hat{y}_t} \cdot \frac{\partial \hat{y}_t}{\partial s_{t-1}} \cdot \frac{\partial s_{t-1}}{\partial w}$$

$s_t = \tanh(u x_t + w s_{t-1}) \Rightarrow s_{t-1} \text{ depend on } s_{t-2} \text{ etc.}$

$$\Rightarrow \frac{\partial E_t}{\partial w} = \sum_{k=0}^t \frac{\partial E_t}{\partial \hat{y}_t} \cdot \frac{\partial \hat{y}_t}{\partial s_t} \cdot \frac{\partial s_t}{\partial s_k} \cdot \frac{\partial s_k}{\partial w}$$

$$\frac{\partial E_t}{\partial v} = \frac{\partial E_t}{\partial y_t} \cdot \frac{\partial y_t}{\partial s_t} \cdot \frac{\partial s_t}{\partial v}$$

$$s_t = \tanh(v \cdot a_t + w s_{t-1})$$

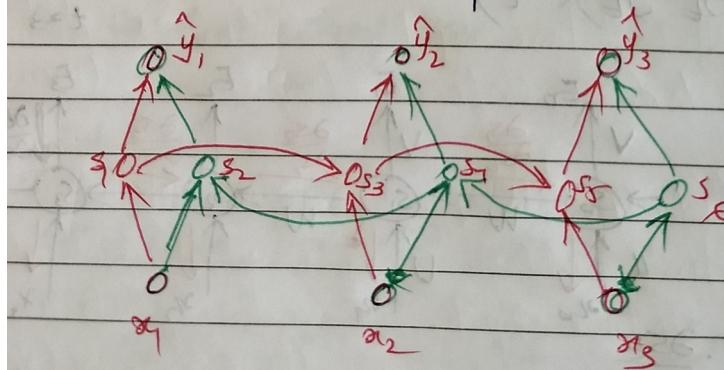
$$\tanh(A+B) = \frac{\tanh(A) + \tanh(B)}{1 + \tanh(A)\tanh(B)}$$

$$\frac{\partial \tanh(x)}{\partial x} = \operatorname{sech}^2 x$$

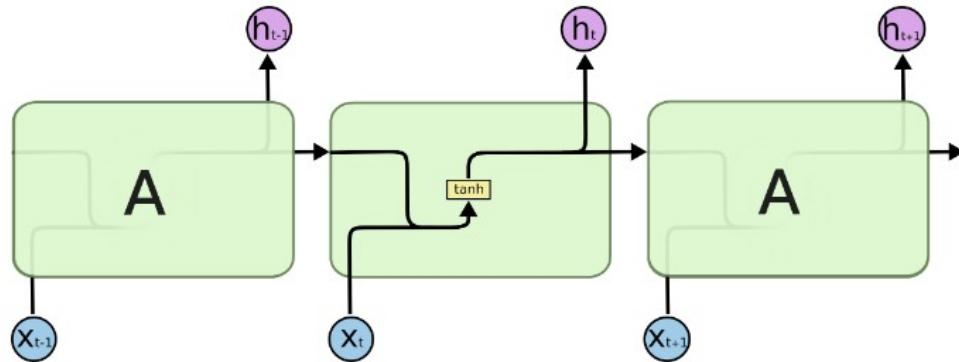
$$\frac{\partial E_t}{\partial v} = \sum_{k=0}^t \frac{\partial E_t}{\partial y_t} \cdot \frac{\partial y_t}{\partial s_t} \cdot \frac{\partial s_t}{\partial s_k} \cdot \frac{\partial s_k}{\partial v}$$

Bidirectional RNN

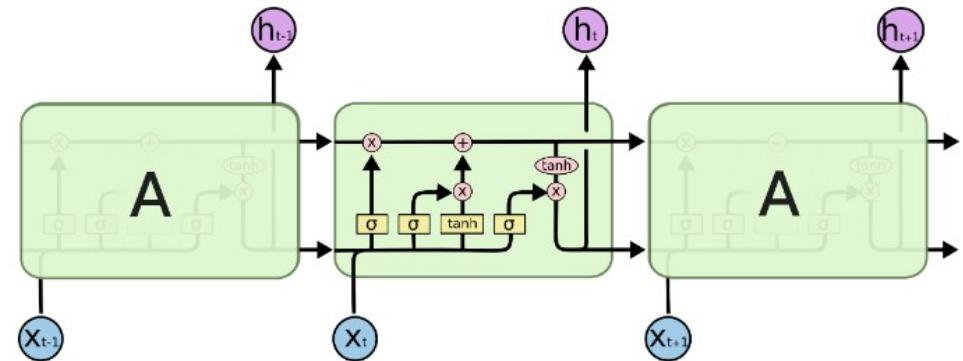
- Output at time 't' depend on both previous and future elements in the sequence.



LSTM



The repeating module in a standard RNN contains a single layer.



The repeating module in an LSTM contains four interacting layers.

$$i_t = \sigma(W_i x_t + b_{ii} + U_i h_{t-1} + b_{hi})$$

$$f_t = \sigma(W_f x_t + b_{if} + U_f h_{t-1} + b_{hf})$$

$$g_t = \tanh(W_g x_t + b_{ig} + U_g h_{t-1} + b_{hg})$$

$$o_t = \sigma(W_o x_t + b_{io} + U_o h_{t-1} + b_{ho})$$

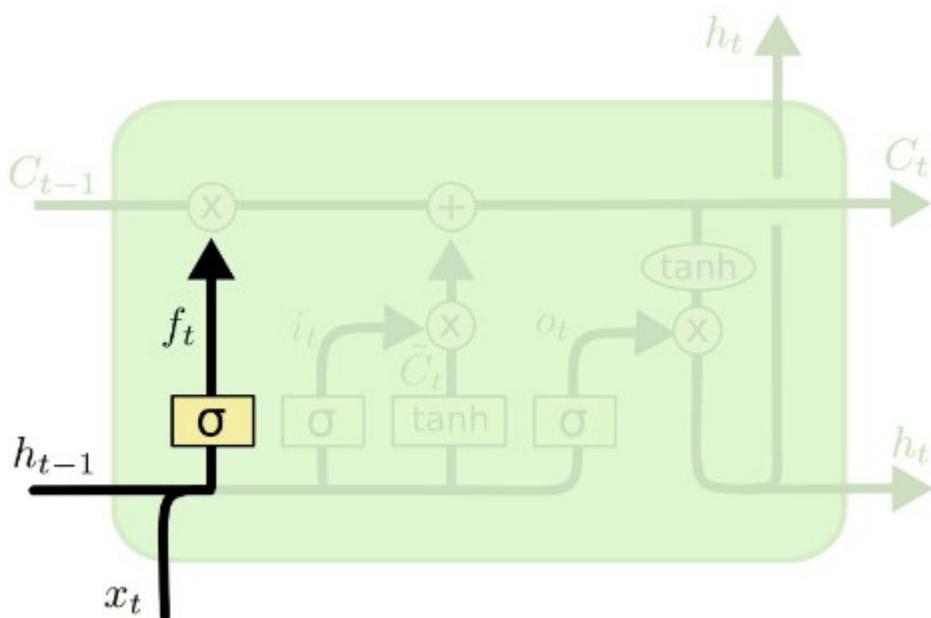
$$c_t = f_t \circ c_{t-1} + i_t \circ g_t$$

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

○ denotes the Hadamard or entry-wise product

Forget gate layer

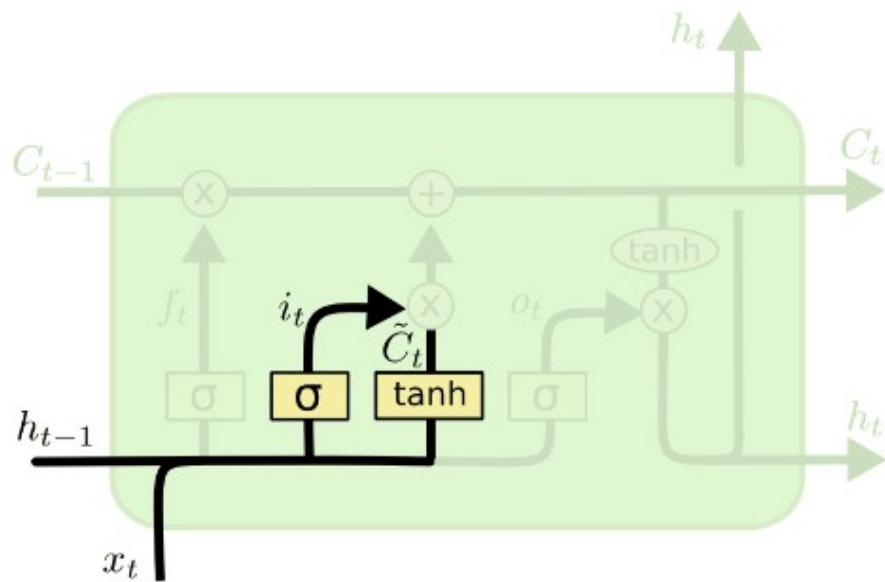
- Decides what information we're going to throw away from the cell state
- This decision is made by a sigmoid layer called the “forget gate layer”
- It looks at h_{t-1} and x_t , and outputs a number between 0 and 1 for each number in the cell state C_{t-1}
- 1 represents “completely keep this” while a 0 represents “completely get rid of this”



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

Input gate layer

- Decide what new information we're going to store in the cell state
- Sigmoid layer called the “input gate layer” decides which values we'll update
- tanh layer creates a vector of new candidate values, $C_{\hat{t}}$, that could be added to the state
- Combine these two to create an update to the state

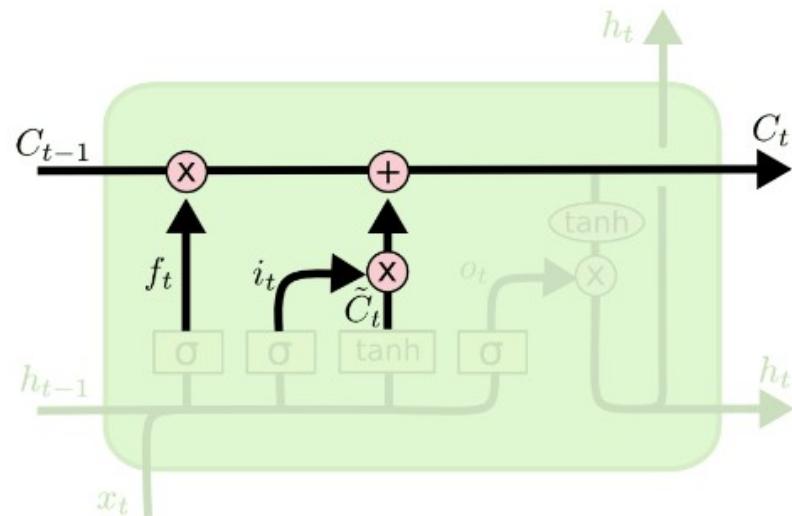


$$i_t = \sigma(W_i \cdot [h_{t-1}, x_t] + b_i)$$

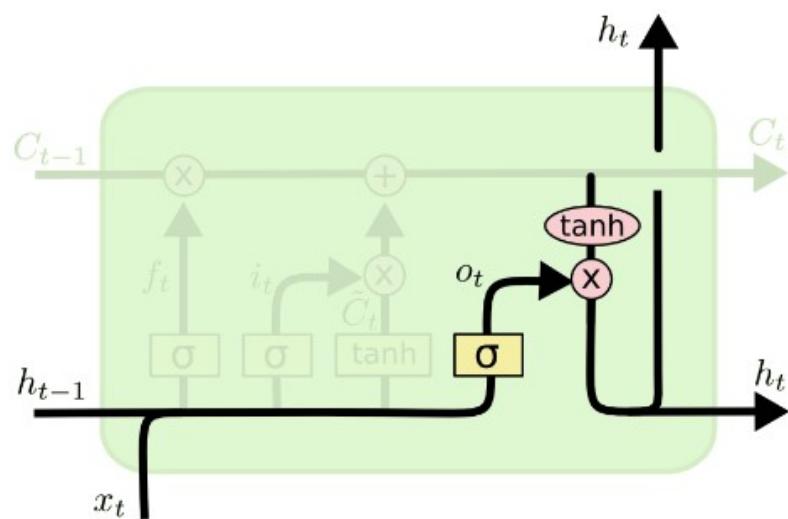
$$\tilde{C}_t = \tanh(W_C \cdot [h_{t-1}, x_t] + b_C)$$

Update cell state, Hidden state, Output

- Update the old cell state, C_{t-1} , into the new cell state C_t



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



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

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

GRU - Gated Recurrent Units \rightarrow learn long term dependence
 2014

$$z = \sigma(x_t V^z + s_{t-1} W^z) \rightarrow \text{update gate}$$

$$r = \sigma(x_t V^r + s_{t-1} W^r) \rightarrow \text{reset gate}$$

$$h = \tanh(x_t V^h + (s_{t-1} \cdot r) W^h)$$

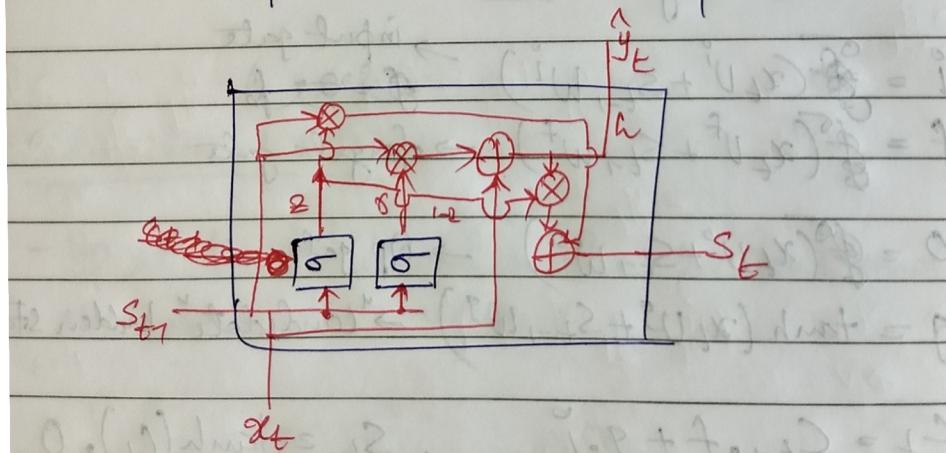
$$s_t = (1-z) \cdot h + z \cdot s_{t-1}$$

reset gate determiner - how to combine the new i/p with the previous memory.

update gate \rightarrow how much of the previous memory to keep around.

If set ~~all~~ reset all '1' & update all '0' \rightarrow plain RNN.

GRU \rightarrow don't possess internal memory ' C_t '.



adv
 - few parameters \rightarrow train faster \rightarrow less data to generalize

Conv-LSTM

FC-LSTM

$$\begin{aligned} i_t &= \sigma(W_{xi}x_t + W_{hi}h_{t-1} + W_{ci} \circ c_{t-1} + b_i) \\ f_t &= \sigma(W_{xf}x_t + W_{hf}h_{t-1} + W_{cf} \circ c_{t-1} + b_f) \\ c_t &= f_t \circ c_{t-1} + i_t \circ \tanh(W_{xc}x_t + W_{hc}h_{t-1} + b_c) \\ o_t &= \sigma(W_{xo}x_t + W_{ho}h_{t-1} + W_{co} \circ c_t + b_o) \\ h_t &= o_t \circ \tanh(c_t) \end{aligned}$$

Input & state at a timestamp are **1D vectors**. Dimensions of the state can be permuted without affecting the overall structure.

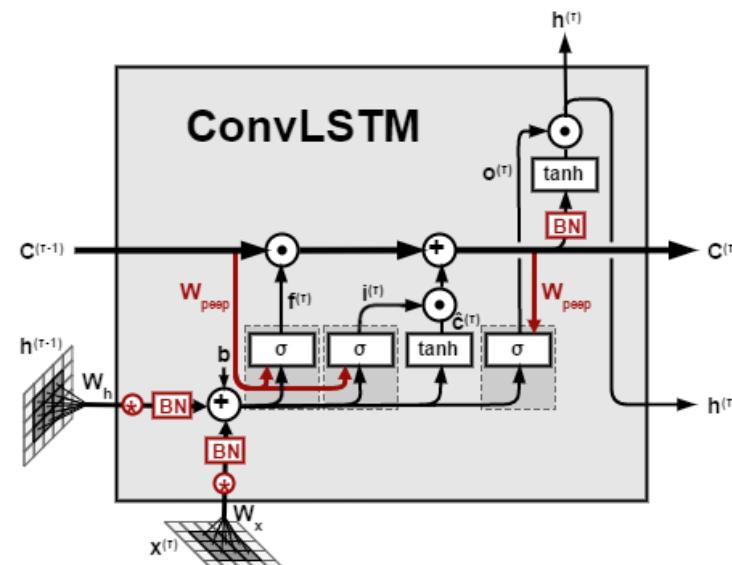
LSTM

$$\begin{aligned} i_t &= \sigma(W_i x_t + b_{ii} + U_i h_{t-1} + b_{hi}) \\ f_t &= \sigma(W_f x_t + b_{if} + U_f h_{t-1} + b_{hf}) \\ g_t &= \tanh(W_g x_t + b_{ig} + U_g h_{t-1} + b_{hg}) \\ o_t &= \sigma(W_o x_t + b_{io} + U_o h_{t-1} + b_{ho}) \\ c_t &= f_t \circ c_{t-1} + i_t \circ g_t \\ h_t &= o_t \circ \tanh(c_t) \end{aligned}$$

ConvLSTM

$$\begin{aligned} i_t &= \sigma(W_{xi} * \mathcal{X}_t + W_{hi} * \mathcal{H}_{t-1} + W_{ci} \circ \mathcal{C}_{t-1} + b_i) \\ f_t &= \sigma(W_{xf} * \mathcal{X}_t + W_{hf} * \mathcal{H}_{t-1} + W_{cf} \circ \mathcal{C}_{t-1} + b_f) \\ \mathcal{C}_t &= f_t \circ \mathcal{C}_{t-1} + i_t \circ \tanh(W_{xc} * \mathcal{X}_t + W_{hc} * \mathcal{H}_{t-1} + b_c) \\ o_t &= \sigma(W_{xo} * \mathcal{X}_t + W_{ho} * \mathcal{H}_{t-1} + W_{co} \circ \mathcal{C}_t + b_o) \\ \mathcal{H}_t &= o_t \circ \tanh(\mathcal{C}_t) \end{aligned}$$

Input & state at a timestamp are **3D tensors**. **Convolution** is used for both **input-to-state** and **state-to-state** connection.



Thank you