

# LSTM

RNN input at time  $t$  is soon forgotten.

$$\vec{y}^{t+100}(\vec{x}^{t+100}, \vec{h}^{t+99}(\vec{x}^{t+99}))$$

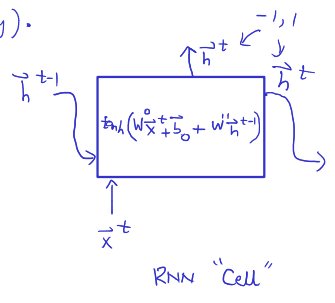
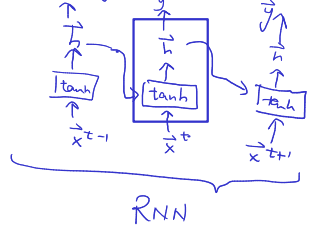
Theoretically,  $\vec{y}^{t+100}$  should depend on  $\vec{x}^1$ .

Practically, it does not.

There is no long-term dependency of outputs on previous inputs.

Because of  $W''$ .

Solved by LSTM (Long Short-term Memory).



Identity

$$\text{tanh}(W''x^t + b) + W''h^{t-1}$$

$$\vec{h}^t = W''h^{t-1} + W''x^t$$

$$= W''(W''h^{t-2} + W''x^{t-1}) + W''x^t$$

$$\vdots$$

$$= (W'')^t \dots + (W'')^1 W''x^0$$

vector

Magnitude depends on largest eval. of  $W''$ .

Affine transformation  $M\vec{x} + \vec{b}$

linear transformation  $M\vec{x}$

## LSTM Cell

$$\vec{v}^t = \begin{bmatrix} \vec{h}^{t-1} \\ \vec{x}^t \end{bmatrix}$$

$(M+D) \times 1$

①  $\sigma(W_f \vec{v}^t + \vec{b}_f)$

$M \times (M+D) \times 1$   $(M+D) \times 1$   $M \times 1$

$M \times 1$

$\vec{f}^t$

$\sigma_f$

$M \times 1$  vector with value between 0 and 1

②  $\sigma(W_i \vec{v}^t + \vec{b}_i)$

$M \times 1$   $(M+D) \times 1$   $M \times 1$

$\vec{i}^t$

$0-1$

③  $\sigma(W_o \vec{v}^t + \vec{b}_o)$

$M \times 1$   $(M+D) \times 1$   $M \times 1$

$\vec{o}^t$

$0-1$

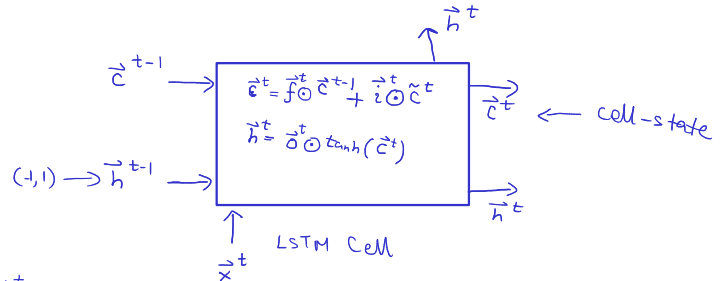
④  $\text{tanh}(W_c \vec{v}^t + \vec{b}_c)$

$M \times 1$   $(M+D) \times 1$   $M \times 1$

$\vec{c}^t$

$-1, 1$   $\text{tanh}_c$

Potentially new cell state



$\odot$  = Hadamard product = element-wise multiplication.

Roles of  $\vec{f}^t, \vec{i}^t, \vec{o}^t$  and  $\vec{c}^t$

Between 0 and 1  $\vec{f}_j^t = \begin{cases} 0 \Rightarrow \vec{c}_j^{t-1} \text{ will be removed/forgotten in the new cell state } \vec{c}_j^t. \\ 1 \Rightarrow \text{retained fully} \end{cases}$

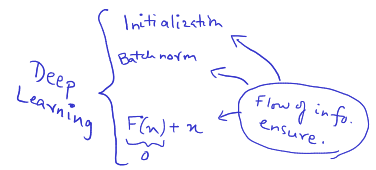
Forget gate on  $\vec{c}^{t-1}$

$\vec{i}_j^t = \begin{cases} 0 \Rightarrow \text{no new information will be added to cell state } \vec{c}_j^t. \\ 1 \Rightarrow \text{add the potential cell state } \vec{c}_j^t \text{ completely (irrespective of forgetness level).} \end{cases}$

Input gate on  $\vec{c}^t$

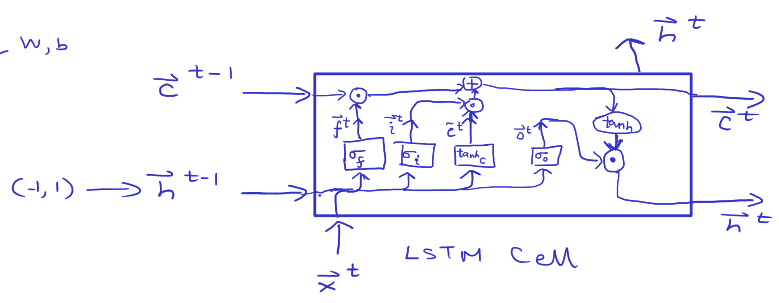
$\vec{o}_j^t = \begin{cases} 0 \Rightarrow \text{completely hide cell state } \vec{c}_j^t. \\ 1 \Rightarrow \text{set } \vec{h}_j^t = \text{tanh}(\vec{c}_j^t). \text{ Cell state } \vec{c}_j^t \text{ is fully exposed in space and time.} \end{cases}$

Output gate on  $\vec{c}^t$



$\square$  = parameterised transformation  $W, b$

$\odot$  = operation



$\vec{c}^t = \vec{f}^t \odot \vec{c}^{t-1} + \vec{i}^t \odot \vec{c}^t$

$\vec{h}^t = \vec{o}^t \odot \text{tanh}(\vec{c}^t)$

Notice that dependency on  $\vec{f}^t$  and  $\vec{i}^t$ ,  $\vec{c}^t$  can be an exact copy of previous cell state  $\vec{c}^{t-1}$ . "Flowing cell state".

If  $\vec{c}^t = 1$ , it can imply that a bracket is currently open. It can remain 1 for a long time until a closing bracket is observed at input.

This gives long term memory.