

Long short-term Memory (LSTM)

Mathematical formulation

- Suppose that there are h hidden units, a batch size n , and d inputs
- $\mathbf{X}_t \in \mathbb{R}^{n \times d}$ and $\mathbf{H}_{t-1} \in \mathbb{R}^{n \times h}$
- We get the following gate values

$$\mathbf{I}_t = \sigma(\mathbf{X}_t \mathbf{W}_{xi} + \mathbf{H}_{t-1} \mathbf{W}_{hi} + \mathbf{b}_i)$$

$$\mathbf{F}_t = \sigma(\mathbf{X}_t \mathbf{W}_{xf} + \mathbf{H}_{t-1} \mathbf{W}_{hf} + \mathbf{b}_f)$$

$$\mathbf{O}_t = \sigma(\mathbf{X}_t \mathbf{W}_{xo} + \mathbf{H}_{t-1} \mathbf{W}_{ho} + \mathbf{b}_o)$$

- $\mathbf{W}_x \in \mathbb{R}^{d \times h}$
- $\mathbf{W}_h \in \mathbb{R}^{h \times h}$
- $\mathbf{I}_t, \mathbf{F}_t, \mathbf{O}_t \in \mathbb{R}^{n \times h}$

σ = Sigmoid

\odot = element-wise product

Memory cell state

- Let us define the memory cell state $\mathbf{C}_t \in \mathbb{R}^{n \times h}$ for timestep t

$$\mathbf{C}_t = \mathbf{F}_t \odot \mathbf{C}_{t-1} + \mathbf{I}_t \odot \tilde{\mathbf{C}}_t$$

- \mathbf{F}_t addresses how much of the old cell internal state \mathbf{C}_{t-1} we retain
- \mathbf{I}_t governs how much we take new data into account via $\tilde{\mathbf{C}}_t$
- If $\mathbf{F}_t = \mathbf{1}$ and $\mathbf{I}_t = \mathbf{0}$ the memory cell remains constant ($\mathbf{C}_t = \mathbf{C}_{t-1}$)

Input node

- $\tilde{\mathbf{C}}_t \in \mathbb{R}^{n \times h}$
 - Its computation is equivalent to the one of three gates but with a tanh activation
- $$\tilde{\mathbf{C}}_t = \tanh(\mathbf{X}_t \mathbf{W}_{xc} + \mathbf{H}_{t-1} \mathbf{W}_{hc} + \mathbf{b}_c)$$

Link with the input gate

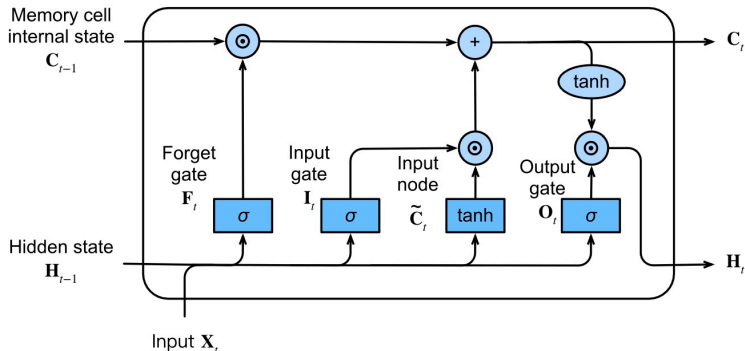
- The value of the input node interacts with the input gate to decide what should be added to the current **internal state**

Output gate and internal state

- Finally, we have to define the output \mathbf{H}_t of the memory cell using both \mathbf{O}_t and \mathbf{C}_t

$$\mathbf{H}_t = \mathbf{O}_t \odot \tanh(\mathbf{C}_t)$$

- When \mathbf{O}_t is close to 0, current memory does not impact the subsequent layer of the network
- When \mathbf{O}_t is close to 1, current memory adds information to the next layer



Gated Recurrent Units (GRU)

$$R_t = \sigma(X_t W_{xr} + H_{t-1} W_{hr} + b_r)$$

$$Z_t = \sigma(X_t W_{xz} + H_{t-1} W_{hz} + b_z)$$

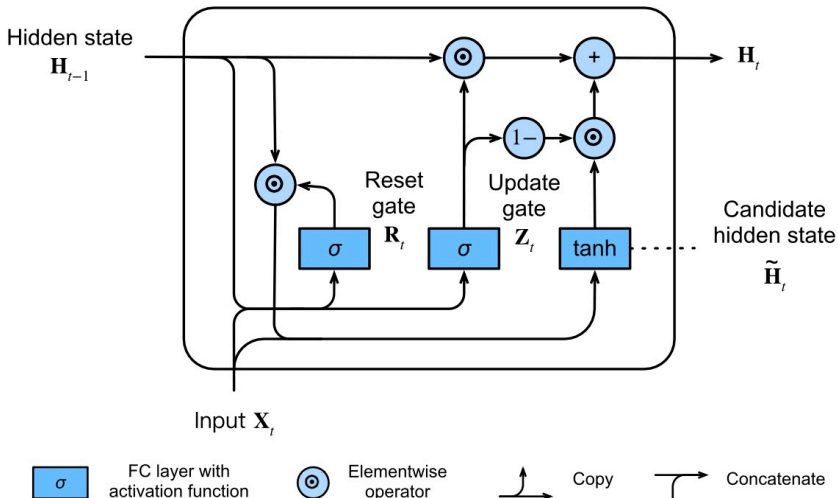
Candidate Hidden State

$$\tilde{H}_t = \tanh(X_t W_{xh} + (R_t \odot H_{t-1}) W_{hh} + b_h)$$

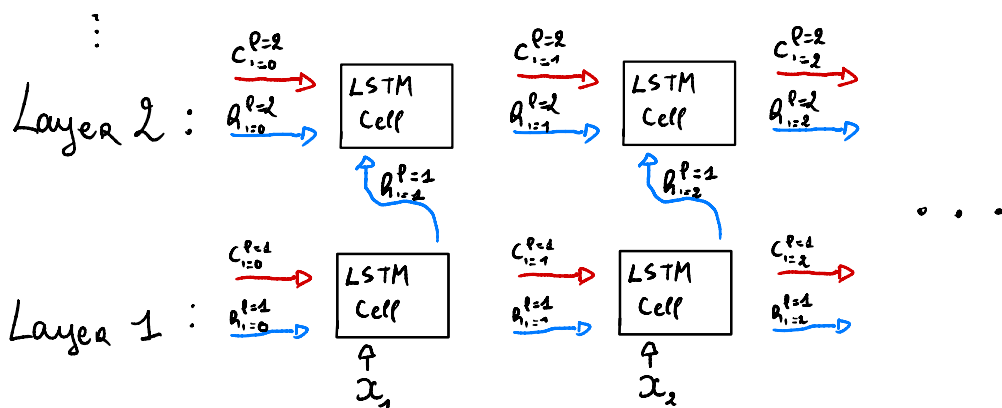
Hidden State

- Finally, we incorporate the effect of the update gate Z_t on the hidden state H_t

$$H_t = Z_t \odot H_{t-1} + (1 - Z_t) \odot \tilde{H}_t$$

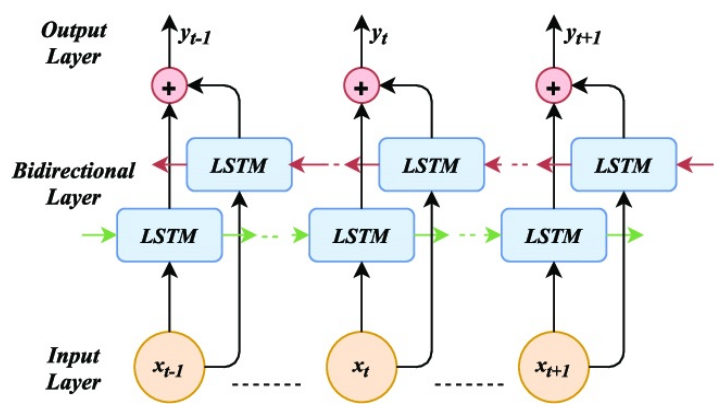


Stacking recurrent cells



(Same thing for GRU)

Bi-directional RNN :



Convolutional RNN :

All matrix products are replaced by convolutional operations. Inputs, hidden states, memory states are not of shape $n \times d$ but $n \times \text{channels} \times \text{height} \times \text{width}$.