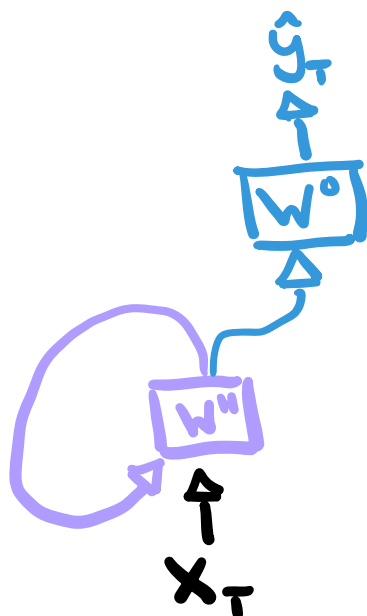
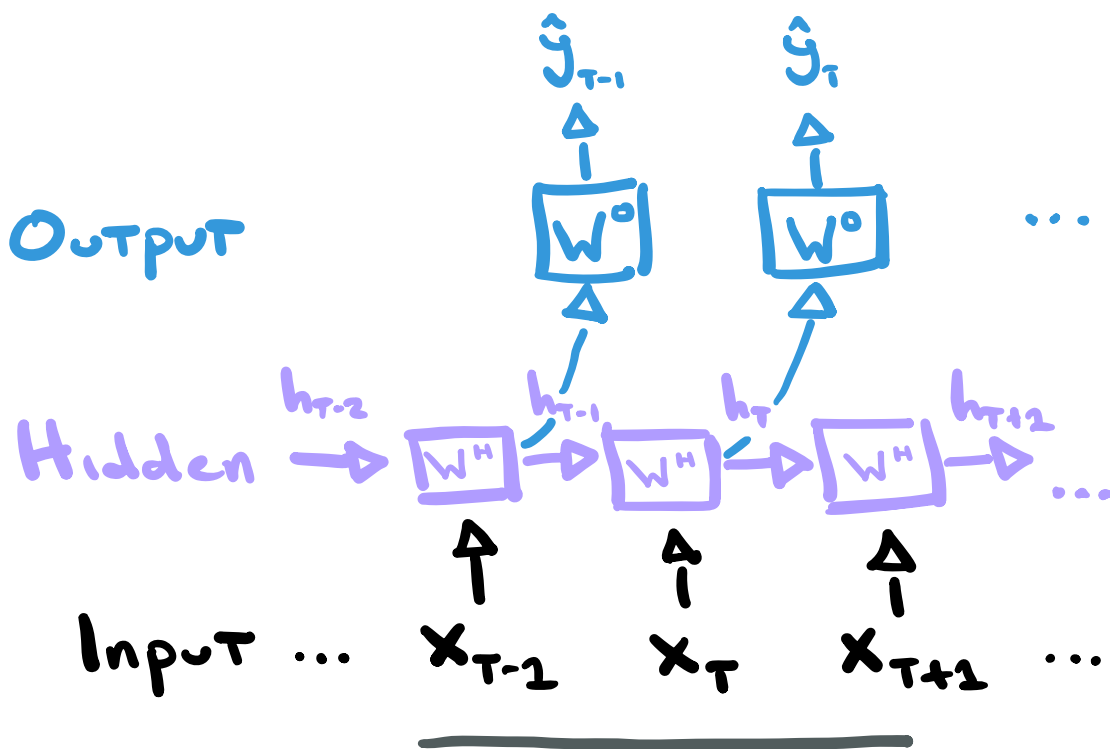


DS 4440

Recurrent Neural Networks 2

Last Time: RNNs for Sequences

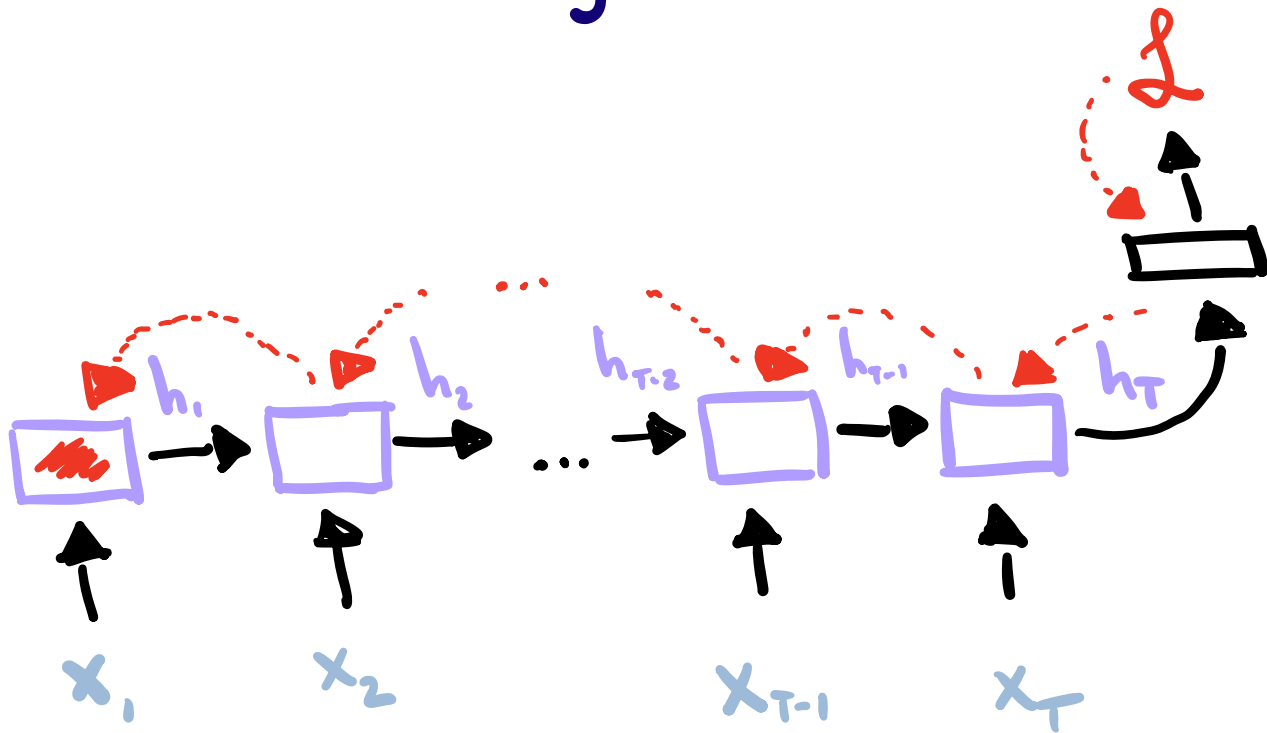


Problem long term dependencies

(0 0 0 0 0)	0
(1 0 0 0 0)	1
...	
(1 0 0 0 0)	1

x y

Signal must flow
Through entire
Sequence!

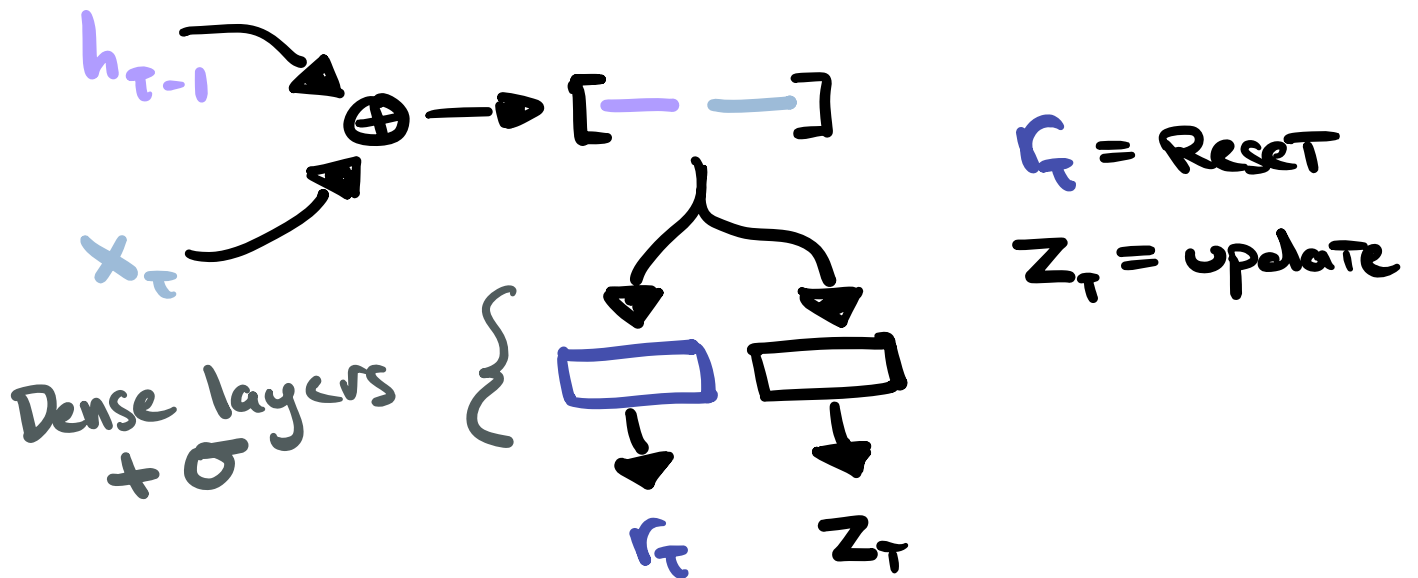


Gated Recurrent Units (GRUs)

Similar to **LSTMs** (next) but
Simpler.

IDEA Introduce gating mechanisms

To allow Model to update h
or opt not to (Skip)



$$r_t = \sigma(x_t W_r^x + h_{t-1} W_r^h + b_r)$$

$$z_t = \sigma(x_t W_z^x + h_{t-1} W_z^h + b_z)$$

Reset

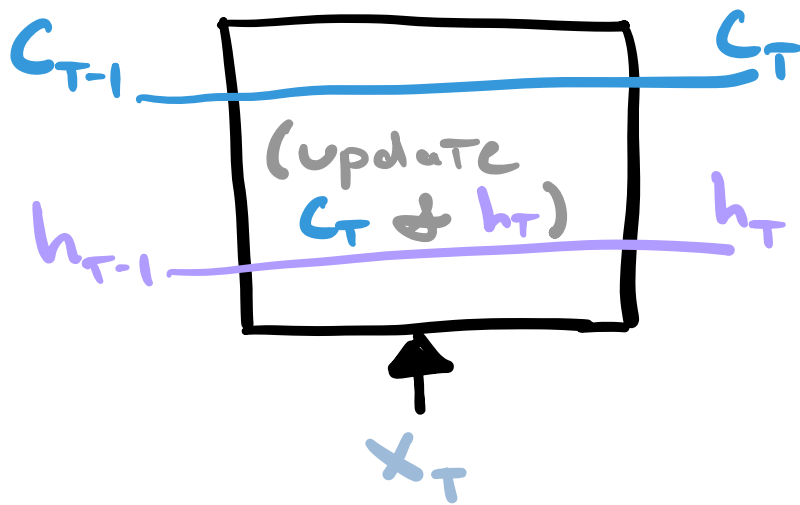
$$\tilde{h}_t \leftarrow \text{Tanh}(x_t W_h^x + \underbrace{(r_t \odot h_{t-1})}_{\text{Permits forgetting}} W_h^h + b_h)$$

Update

$$h_t \leftarrow \underbrace{z_t \odot h_{t-1}}_{\text{Copy}} + \underbrace{(1 - z_t) \odot \tilde{h}_t}_{\text{Update}} \quad \left| \begin{array}{l} \text{See} \\ \text{Notebook} \end{array} \right.$$

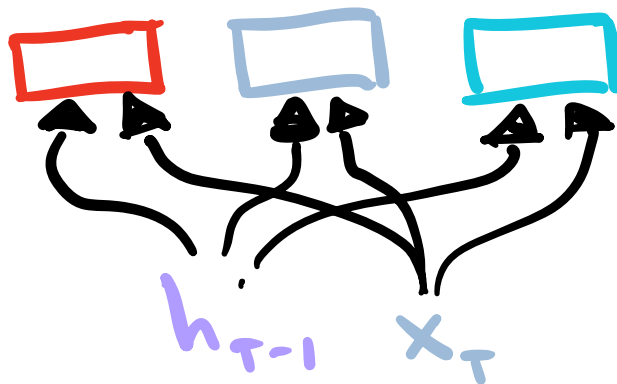
Long Short Term Memory (LSTM)

Like GRUs, use gating for long term dependencies.



Three gates:

forget Input Output



Assume batch size \underline{b} , \underline{d} hidden dims.

$$\underline{x}_T \in \mathbb{R}^{b \times d_{in}} \quad \underline{h}_{T-1} \in \mathbb{R}^{b \times d}$$

$$\underline{f}_T = \sigma(\underline{x}_T \underline{W}_f^x + \underline{h}_{T-1} \underline{W}_f^h + \underline{b}_f)$$

$$\underline{i}_T = \sigma(\underline{x}_T \underline{W}_i^x + \underline{h}_{T-1} \underline{W}_i^h + \underline{b}_i)$$

$$\underline{o}_T = \sigma(\underbrace{\underline{x}_T}_{b \times d} \underbrace{\underline{W}_o^x}_{d_{in} \times d} + \underbrace{\underline{h}_{T-1}}_{b \times d} \underbrace{\underline{W}_o^h}_{d \times d} + \underbrace{\underline{b}_o}_{1 \times d})$$

Apply These gates to a Candidate
"Memory cell" $\tilde{\underline{C}}_T$.

$$\underbrace{\tilde{\underline{C}}_T}_{b \times d} = \text{Tanh}(\underline{x}_T \underline{W}_c^x + \underline{h}_{T-1} \underline{W}_c^h + \underline{b}_c)$$

Next update The Memory

$$\underline{C}_T \leftarrow \underline{f}_T \odot \underline{C}_{T-1} + \underline{i}_T \odot \tilde{\underline{C}}_T$$

Finally, the Output layer, which

is a modified version of C_T .

$$h_T \leftarrow O_T \odot \tanh(C_T)$$

Info can be stored in C_T
w/o affecting h_T if O_T near
0.