

## 24) RNN

### Sequence Processing

- Data ordered sequentially & neighbors are correlated
- RNN leverages it to ① predict the short-term future ② categorize the past



① Time Series data: easy to interpret as a numerical table

② Language: encode each letter/word numerically

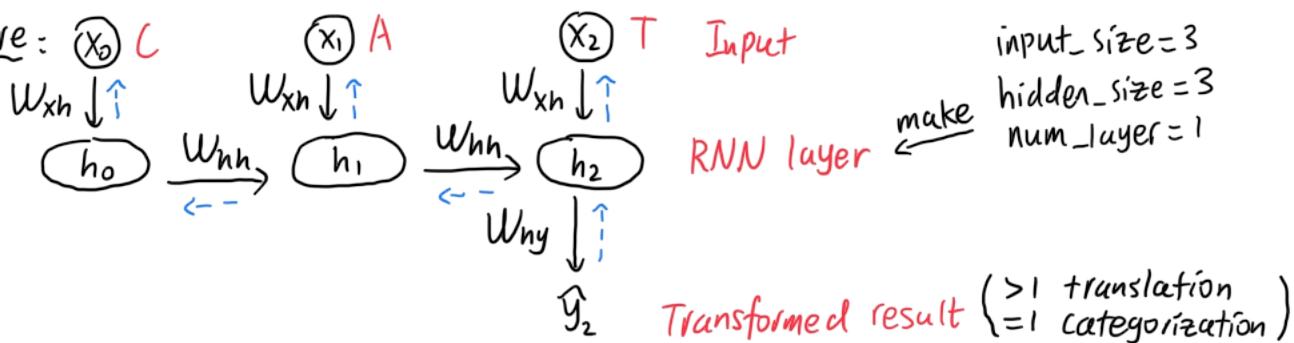
Ex)  $\text{ord}(\text{Hello}) \Rightarrow 8-5-12-12-15$ ,  $\text{vocab\_index}(\text{Entire Words}) \Rightarrow 234-3048$

- Non-RNN Limitation: ① Vanilla CNN: fixed input length, ignores context & order  
② FFN: ignores structures, independent learning params

### Vanilla RNN

- Recurrent network designed to sequence processing

• Architecture:



• Backprop Risks: Due to  $\alpha \frac{\partial J}{\partial W}$ , a deep RNN is oversensitive to recent inputs vanishing to distant inputs

• metaparams:

- ① Input size: # features / data channels (1 for words)
- ② Sequence length: length of overlapping data segment
- ③ Batch size: # data segments processed per epoch

translation data: 0 1 2 3 4 5 6	↓ train	① = 1 ② = 4 ③ = 2
[0 1 2 3] 4		
[2 3 4 5] 6		
x0 y		

### More on RNNs

- ① Hidden State  $h$ : activations of a hidden layer  $\in \mathbb{R}^{\text{hidden\_size}}$   
memory from previous sequence
- Changes at each timestep  $t$  from  $x_t$ :  $h_t = \tanh(W_{xh}x_t + W_{hh}h_{(t-1)})$
- Initializes with 0s so it's not biased with random memory at  $t=0$

② Embedding: represents each letter as  $N$ -dim vector with linear weights

- Interpretation: Over training, embedding layer learns structural patterns b/t data  
 $N$  is arbitrary # and is abstract to interpret  $h_i [ \begin{smallmatrix} 1 & 2 & 3 \\ 2 & 0 & 1 \\ 4 & 1 & 2 \end{smallmatrix} ]$   $\leftarrow$  dim index

- Issue in One-Hot Encoding: makes each encoded letter is orthogonal  
 $(\text{did not capture natural correlation in a word})$  
 $h_i [ \begin{smallmatrix} 1 & \dots & h_i & \dots & 0 \\ 0 & 0 & 1 & 0 & 0 \\ 0 & 0 & 0 & 1 & 0 \\ 0 & 0 & 0 & 0 & 0 \end{smallmatrix} ]$

## GRU & LSTM

- Gate: a sigmoid output that controls the flow of information when multiply

- Gated Recurrent Unit (GRU): useful for longer patterns but need more data

- Process:

① Reset Gate:  $r_t = \sigma(W_r[x_t; h_{t-1}])$

$\hookrightarrow$  model learns how much of past to update

② Update Gate:  $z_t = \sigma(W_z[x_t; h_{t-1}])$

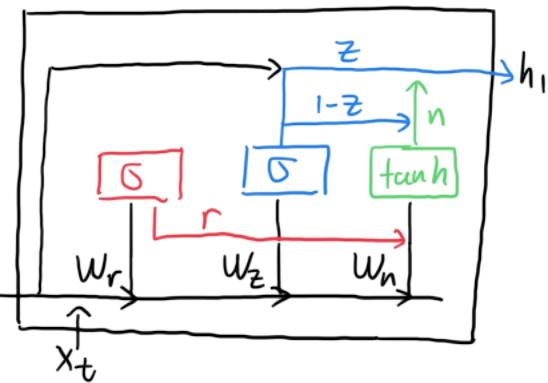
$\hookrightarrow$  model controls what information to update

③ New information:  $n_t = \tanh(W_n[x_t; r_t h_{t-1}])$

$\hookrightarrow$  model learns a new hidden state candidate

④ Hidden state (output):  $h_t = (1 - z_t) n_t + z_t h_{t-1}$

architecture (1 hidden unit)



- Long Short-Term Memory (LSTM): useful for memory-heavy & complex tasks

- Process:

① Forget Gate:  $f_t = \sigma(W_f[x_t; h_{t-1}])$  knocks out unimportant items

② Input Gate:  $i_t = \sigma(W_i[x_t; h_{t-1}])$  filters how much new data to add

③ Cell Gate:  $g_t = \tanh(W_i[x_t; h_{t-1}])$  computes actual new data

④ Output Gate:  $o_t = \sigma(W_o[x_t; h_{t-1}])$  controls how much info to predict

⑤ Cell State:  $c_t = f_t c_{t-1} + i_t g_t$  is what LSTM remembers

⑥ Hidden State:  $h_t = o_t \tanh(c_t)$  is what LSTM predicts