

$$1. \text{ input } x = [0, 1, 0, 1, 1, 0]$$

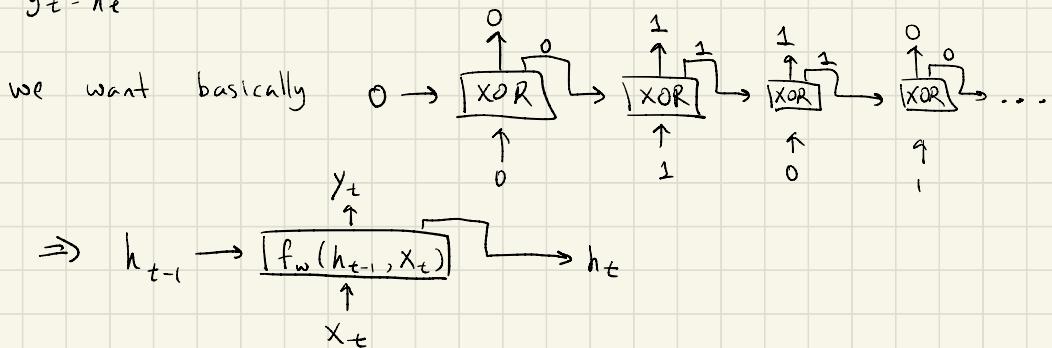
$$\text{output } y = [0, 1, 1, 0, 1, 1]$$

$$h_0 = 0$$

$$h_t = h_{t-1} \text{ XOR } x_t \quad * \text{ can't use XOR, but XOR is the relationship}$$

$$y_t = h_t$$

we want basically



$$\text{let } f_w(h_{t-1}, x_t) = h_{t-1} \cdot \bar{x}_t + \bar{h}_{t-1} \cdot x_t$$

We see that this acts as an XOR gate for binary inputs

h_{t-1}	x_t	$f_w(h_{t-1}, x_t)$
0	0	0
0	1	1
1	0	1
1	1	1

$$2. f_t = \sigma(W_f \cdot [h_{t-1}, x_t] + b_f)$$

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

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

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

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

$$h_t = O^t * \tanh(C_t)$$

Assume $h_0 = 0$ $C_0 = 0$

let $W_f = \begin{bmatrix} 1 \\ 1 \end{bmatrix}$ $b_f = 0$

$$W_i = \begin{bmatrix} -1 \\ 1 \end{bmatrix} \quad b_i = 0$$

$$W_c = \begin{bmatrix} 0 \\ 1 \end{bmatrix} \quad b_c = 0$$

$$W_o = \begin{bmatrix} 0 \\ 0 \end{bmatrix} \quad b_o = 1$$

h_{t-1}	x_t	f_t	i_t	\tilde{C}_t	O^t	C_t	h_t	← note $C_{t-1} = h_{t-1}$
0	0	0	0	0	1	0	0	
0	1	0	1	1	1	1	1	
1	0	1	0	0	1	1	1	
1	1	0	0	1	1	0	0	can see that we get XOR from h_{t-1} and x_t

3. $\text{best}_{\leq t}$ is defined at t

Highest scoring beam in B_t is y_t^1 with score S_t^1

$$S_t^1 \leq \text{best}_{\leq t}$$

probability $\in [0, 1]$ due to $\sum_x p(y_{t+1} | x, y_{\leq t}) = 1$

$$\text{best}_{\leq t+1} = S_t^1 + \log p(y_x | x, y_{\leq t})$$

$S_t^i \leq S_t^1 \leq \text{best}_{\leq t}$ for all other beams in B_t

$$S_{t+1}^1 = S_t^1 + \log p(y_x | x, y_{\leq t}) \quad \text{and since probability } \in [0, 1]$$

All beams in $t' > t$ will be no better than $\text{best}_{\leq t}$

$$4. h_t = W^T h_{t-1}$$

let $\lambda_1, \dots, \lambda_n$ be the eigenvalues of weight matrix $W \in \mathbb{C}^{n \times n}$

$$\rho(W) = \max \{|\lambda_1|, \dots, |\lambda_n|\}$$

Assume initial hidden state h_0 .

$$h_t = W^T h_{t-1}$$

$$h_{t-1} = W^T h_{t-2}$$

:

$$h_t = W^T h_0 \Rightarrow h_t = (W^T)^t h_0$$

take derivative $\frac{dh_t}{dh_0} = (W^T)^t$

decompose $W = P D P^{-1}$ where P is eigenvectors of W
 D is diagonal matrix of W 's eigenvalues

$$= ([P D P^{-1}]^t)^t$$

$$= ([P^{-1}]^t D^t P^t)^t \quad \text{note } D = D^T$$

$$= ([P^{-1}]^t D P^t)^t$$

$$= [P^{-1}]^t D^t P^t$$

we can see as $t \gg 0$:

if $\rho(W) > 1$ then at least one value of D^t

will exponentially grow toward $+\infty$ \leftarrow exploding gradient

if $\rho(W) < 1$ then all eigenvalues will grow towards 0

\downarrow
vanishing gradient

\therefore eigenvalues determine whether exploding/vanishing gradient

$$\begin{aligned}5. \text{ a) } h_i^{t+1} &= q(h_i^t, \text{Agg}(H_{i:t}')) \\&= q(h_i^t, \sum_{j \in N(v_i)} f_j(h_j^t)) \\&= q(h_i^t, \sum_{j \in N(v_i)} v^t h_j^t) \quad \text{where } v_j \in N(v_i)\end{aligned}$$

$$5.b) h_1^t = [1, -1] \quad h_2^t = [-1, 1] \quad h_3^t = [0, -1] \quad h_4^t = [1, 0]$$

$$\text{Agg}(H_{1:t}^t) = [0.6, 0.2, 0.2] \quad \begin{bmatrix} f(h_2^t) \\ f(h_3^t) \\ f(h_4^t) \end{bmatrix}$$

$$f(x) = 2x$$

$$h_1^{t+1} = g(h_1^t, \text{Agg}(H_{1:t}^t)) = W(h_1^t)^T + \max \{ \text{Agg}(H_{1:t}^t), 0 \}$$

where $W = [1, 1]$

$$f(h_2^t) = [-2, 2]$$

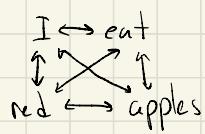
$$f(h_3^t) = [0, -2]$$

$$f(h_4^t) = [2, 0]$$

$$\text{Agg}(H_{1:t}^t) = [0.6, 0.2, 0.2] \quad \begin{bmatrix} [-2, 2] \\ [0, -2] \\ [2, 0] \end{bmatrix} = [-0.8, 0.8]$$

$$\begin{aligned} h_1^{t+1} &= W(h_1^t) + \max \{ \text{Agg}(H_{1:t}^t), 0 \} \\ &= [1, 1] \begin{bmatrix} 1 \\ -1 \end{bmatrix} + \max \{ [-0.8, 0.8], 0 \} \\ &= 0 + [0, 0.8] = [0, 0.8] \end{aligned}$$

5.c)



$$2 \times 6 = \boxed{12}$$

$$5. d) \quad h_i^{l+1} = \sum_{j \in S} (\text{softmax}_j(Q^l h_i^l \cdot K^l h_i^l) V^l h_j^l) \quad (17)$$

$$h_i^{l+1} = q\left(h_i^l, \sum_{j \in N(v_i)} (V^l h_j^l)\right) \quad (15) \quad \text{from part a}$$

$$W_{ij} = \text{softmax}_j(Q^l h_i^l \cdot K^l h_j^l)$$

$$\Rightarrow (17) \text{ is } h_i^{l+1} = \sum_{j \in N(v_i)} W_{ij} (V^l h_j^l) = W(h_i^l)^T \sum_{j \in N(v_i)} (V^l h_j^l)$$

5. e) For fully connected graphs, the number of connections increase exponentially, thus it becomes exponentially more expensive to learn longer term dependencies.