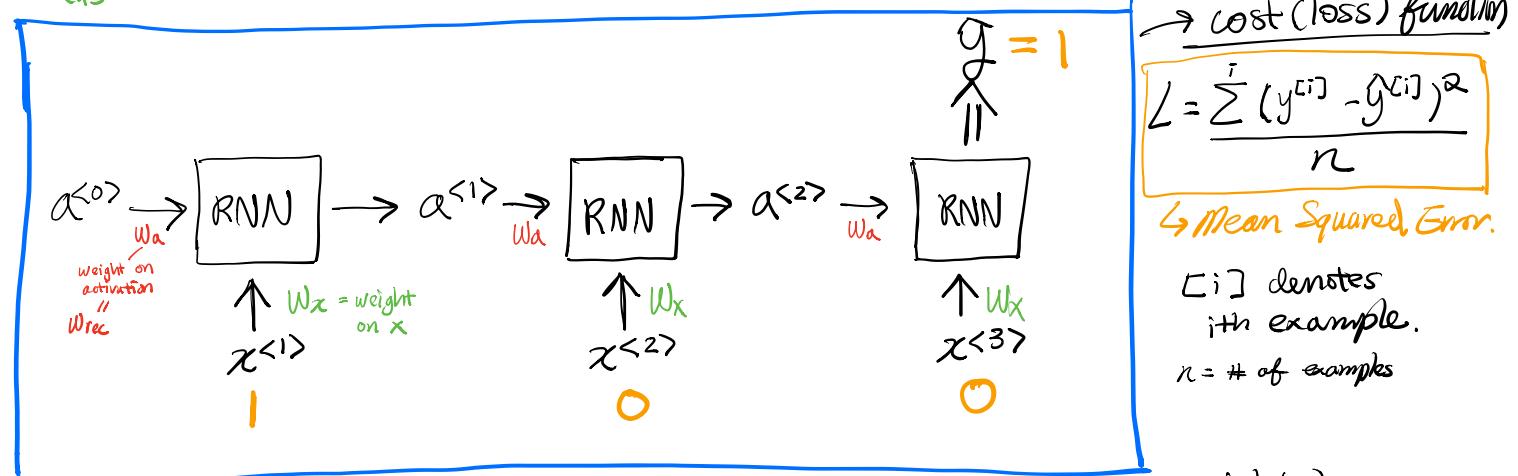


< RNN by Hand >

Data = $\begin{bmatrix} [1 \ 0 \ 0] \\ [0 \ 1 \ 1] \\ [1 \ 1 \ 1] \end{bmatrix} \rightarrow \begin{bmatrix} 1 \\ 2 \\ 3 \end{bmatrix}$

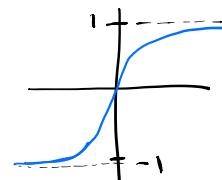
3 examples, length of 3
target (sum of input)

Goal: have RNN sum the input sequence as the output.



others refer it as $S^{[t]}$ for state

$$\tanh(x) = \frac{\sinh(x)}{\cosh(x)}$$



we will ignore bias for

Step 1: initialize variables. & first forward prop,

$w_a = 1$ $w_x = 1$ $\alpha^{<0>} = 0$	$\alpha^{<1>} = \tanh(1 \cdot 0 + 1 \cdot 1) = 0.76$ $\alpha^{<2>} = \tanh(1 \cdot 0.76 + 1 \cdot 0) = 0.64$ $\alpha^{<3>} = \tanh(1 \cdot 0.64 + 1 \cdot 0) = 0.57 = \hat{y}^{[0]} \rightarrow y^{[0]} = 1$
--	--

ex2 \rightarrow

$\alpha^{<1>} = \tanh(1 \cdot 0 + 1 \cdot 0) = 0$ $\alpha^{<2>} = \tanh(1 \cdot 0 + 1 \cdot 1) = 0.76$ $\alpha^{<3>} = \tanh(1 \cdot 0.76 + 1 \cdot 1) = 0.94 = \hat{y}^{[2]} \rightarrow y^{[2]} = 2$
--

ex3 \rightarrow

$\alpha^{<1>} = \tanh(1 \cdot 0 + 1 \cdot 1) = 0.76$ $\alpha^{<2>} = \tanh(1 \cdot 0.76 + 1 \cdot 1) = 0.94$ $\alpha^{<3>} = \tanh(1 \cdot 0.94 + 1 \cdot 1) = 0.96 = \hat{y}^{[3]} \rightarrow y^{[3]} = 3$
--

1st Epoch

Loss = $L = \frac{(1-0.57)^2 + (2-0.94)^2 + (3-0.96)^2}{3} = 1.82$

Step 2: Back propagation

$$L(\theta) = \frac{1}{n} \sum_i^n (h_{\theta}(x^{(i)}) - y^{(i)})^2 = \frac{1}{n} \sum_i^n (\alpha^{<3>(i)} - y^{(i)})^2$$

$$\frac{\partial}{\partial \theta_0} L(\theta_0, \theta_1) = \frac{\partial}{\partial \theta_0} \left(\frac{1}{n} \sum_i^n (h_{\theta}(x^{(i)}) - y^{(i)})^2 \right)$$

$$= \frac{1}{n} \sum_i^n \frac{\partial}{\partial \theta_0} (h_{\theta}(x^{(i)}) - y^{(i)})^2$$

$$= \frac{1}{n} \sum_i^n 2(h_{\theta}(x^{(i)}) - y^{(i)}) \cdot \frac{\partial}{\partial \theta_0} (h_{\theta}(x^{(i)}) - y^{(i)})$$

$$= \frac{\alpha^{<1>}}{n} \sum_i^n (h_{\theta}(x^{(i)}) - y^{(i)})$$

To keep things simple, we simplify the activation function to be:

$$\alpha^{<t>} = W_a \alpha^{<t-1>} + W_x x^{<t>} \quad \begin{matrix} \downarrow & \downarrow & \downarrow \\ h_{\theta} & \theta_0 & \theta_1 \end{matrix}$$

7

State 3 (Last State)

$$\frac{\partial L}{\partial W_x} = \frac{\partial L}{\partial \alpha^{<3>}} \cdot \frac{\partial \alpha^{<3>}}{\partial W_x} = \frac{\partial}{\partial} \sum_i^n (\alpha^{<3>(i)} - y^{(i)}) \cdot x^{<3>(i)} \quad - (a)$$

$$\frac{\partial L}{\partial W_a} = \frac{\partial L}{\partial \alpha^{<3>}} \cdot \frac{\partial \alpha^{<3>}}{\partial W_a} = \frac{\partial}{\partial} \sum_i^n (\alpha^{<3>(i)} - y^{(i)}) \cdot \alpha^{<3>(i)} \quad - (b)$$

$$\begin{cases} \frac{\partial \alpha^{<2>}}{\partial W_x} = x^{<2>} & , \frac{\partial \alpha^{<2>}}{\partial W_a} = \alpha^{<1>} \\ \frac{\partial \alpha^{<3>}}{\partial \alpha^{<2>}} = \frac{\partial}{\partial} (W_a \alpha^{<2>} + W_x x^{<2>}) = W_a \\ \frac{\partial L}{\partial \alpha^{<3>}} = \frac{\partial}{\partial} \sum_i^n (\alpha^{<3>(i)} - y^{(i)}) \end{cases}$$

State 2 ($\alpha^{<2>}$)

$$\frac{\partial L}{\partial W_x} = \frac{\partial L}{\partial \alpha^{<3>}} \cdot \frac{\partial \alpha^{<3>}}{\partial \alpha^{<2>}} \cdot \frac{\partial \alpha^{<2>}}{\partial W_x} = \frac{\partial}{\partial} \sum_i^n (\alpha^{<3>(i)} - y^{(i)}) (W_a) (x^{<2>(i)}) \quad - (c)$$

$$\frac{\partial L}{\partial W_a} = \frac{\partial L}{\partial \alpha^{<3>}} \cdot \frac{\partial \alpha^{<3>}}{\partial \alpha^{<2>}} \cdot \frac{\partial \alpha^{<2>}}{\partial W_a} = \frac{\partial}{\partial} \sum_i^n (\alpha^{<3>(i)} - y^{(i)}) (W_a) (\alpha^{<2>(i)}) \quad - (d)$$

$$\begin{cases} \frac{\partial \alpha^{<2>}}{\partial \alpha^{<1>}} = W_a \\ \frac{\partial \alpha^{<1>}}{\partial W_x} = x^{<1>} & , \frac{\partial \alpha^{<1>}}{\partial W_a} = \alpha^{<0>} \end{cases}$$

State 1 ($\alpha^{<0>}$)

$$\frac{\partial L}{\partial W_x} = \frac{\partial L}{\partial \alpha^{<3>}} \cdot \frac{\partial \alpha^{<3>}}{\partial \alpha^{<2>}} \cdot \frac{\partial \alpha^{<2>}}{\partial \alpha^{<1>}} \cdot \frac{\partial \alpha^{<1>}}{\partial W_x} = \frac{\partial}{\partial} \sum_i^n (\alpha^{<3>(i)} - y^{(i)}) \cdot (W_a) \alpha^{<2>} \cdot (x^{<1>(i)}) \quad - (e)$$

$$\frac{\partial L}{\partial W_a} = \dots \cdot \frac{\partial \alpha^{<1>}}{\partial W_a} = \frac{\partial}{\partial} \sum_i^n (\alpha^{<3>(i)} - y^{(i)}) (W_a) \alpha^{<2>} \cdot (\alpha^{<0>(i)}) \quad - (f)$$

* Remember that I did not use consistent activation function in forward prop + backprop. But the flow of calculations should be same

$$(a) \frac{2}{3} [(0.57-1)(0) + (0.94-2)(1) + (0.96-3)(1)] = -1.36 \quad \leftarrow \text{probably need to multiply by learning rate since } -1.36 \text{ is too big}$$

$$(b) \frac{2}{3} [(0.57-1)(0.64) + (0.94-2)(0.76) + (0.96-3)(0.94)] = \dots$$

(c) (d) (e) (f) → do the same for these equations

Step 3: Update Weights

$$W_x = W_x - (a) + (c) + (e)$$

$$W_a = W_a - (b) + (d) + (f)$$

Step 4: Repeat!

Repeat Step 1 until the model converges to minimum loss.

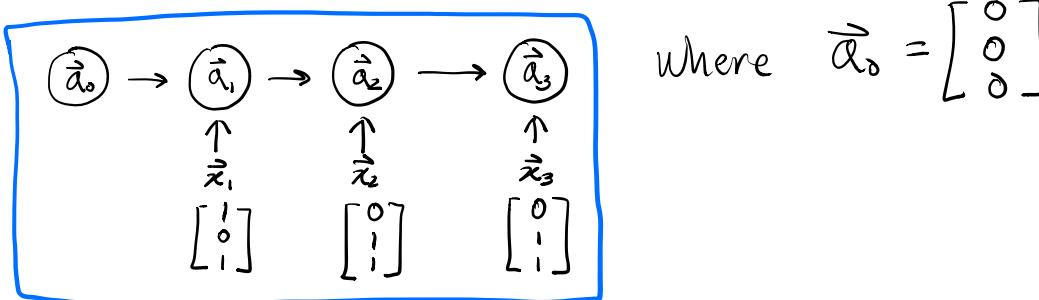
<Vectorization>

vector implementation of above

$$x = \begin{bmatrix} [1, 0, 0] \\ [0, 1, 1] \\ [1, 1, 1] \end{bmatrix} \quad \vec{x}^{(1)} \\ \vec{x}^{(2)} \\ \vec{x}^{(3)}$$

$\vec{x}_1 \quad \vec{x}_2 \quad \vec{x}_3$

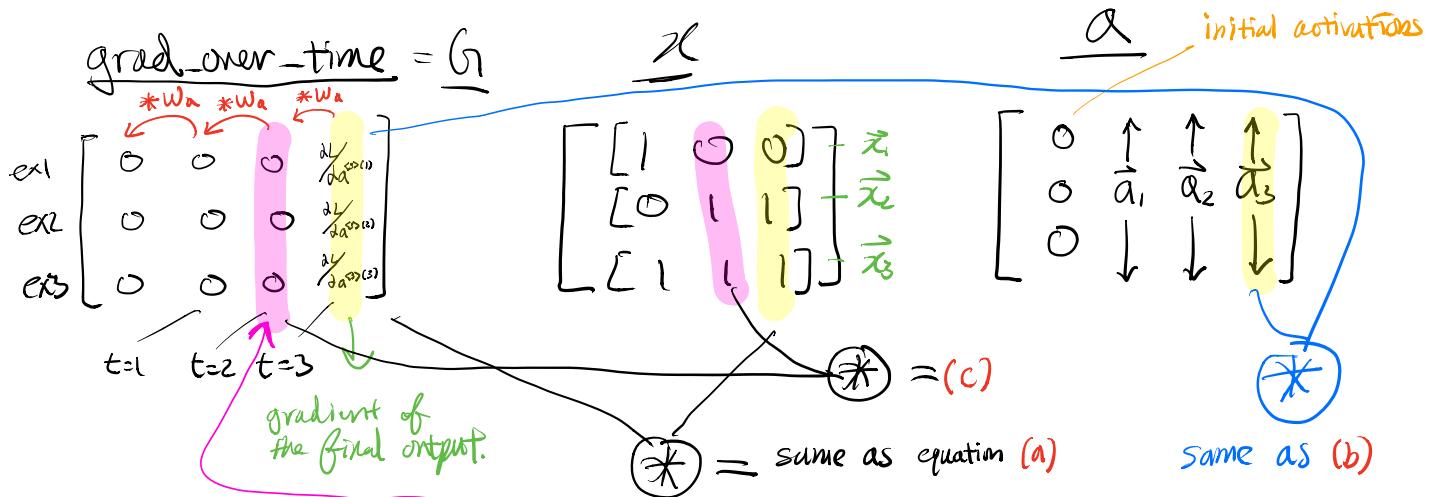
$$\vec{y} = [1, 2, 3]$$



$$\begin{cases} \vec{a}_1 = \vec{a}_0 \cdot W_a + \vec{x}_1 \cdot W_x = \begin{bmatrix} 0 \\ 0 \\ 0 \end{bmatrix} W_a + \begin{bmatrix} 1 \\ 0 \\ 1 \end{bmatrix} W_x = \begin{bmatrix} w_{x1} \\ 0 \\ w_{x3} \end{bmatrix} \\ \vec{a}_2 = \vec{a}_1 \cdot W_a + \vec{x}_2 \cdot W_x = \begin{bmatrix} w_{x1} \\ 0 \\ w_{x3} \end{bmatrix} W_a + \begin{bmatrix} 0 \\ 1 \\ 1 \end{bmatrix} W_x = \begin{bmatrix} w_{x1}w_{a1} + w_{x3} \\ 0 \\ w_{x1}w_{a1} + w_{x3} \end{bmatrix} = \begin{bmatrix} w_{x1}w_{a1} \\ w_{x1} \\ w_{x1}w_{a1} + w_{x3} \end{bmatrix} \\ \vec{a}_3 = \vec{a}_2 \cdot W_a + \vec{x}_3 \cdot W_x = \dots \end{cases}$$

$$L = \frac{1}{n} (\vec{a}_3 - \vec{y})^2 = \frac{1}{n} \left(\begin{bmatrix} a_{31} \\ a_{32} \\ a_{33} \end{bmatrix} - \begin{bmatrix} 1 \\ 2 \\ 3 \end{bmatrix} \right)^2 = \boxed{\quad}$$

Below is in reference to the 'RNN from Scratch' Jupyter Notebook.



$$\frac{\partial L}{\partial a^{>(t)}} = \frac{2}{n} \sum_{i=1}^n (a^{>(i)} - y^{(i)})$$

update value

$$\frac{\partial L}{\partial W} = \left(\frac{\partial L}{\partial a^{>(t)}} \cdot \frac{\partial a^{>(t)}}{\partial a^{<(t)}} \right) \cdot \frac{\partial a^{<(t)}}{\partial W}$$

$$= \boxed{\frac{2}{n} \sum_{i=1}^n (a^{>(i)} - y^{(i)}) \cdot W_a}$$