

learning long-range dependencies.

Depending on activation and parameters we get explode / vanish if the gradients are large that's called exploding gradient problem.

⇒ The basic idea of using a gating mechanism to learn long-term dependencies.

## \* Implementation of Vanilla RNN :-

Let suppose we have two sentences

1. Vocabulary = [<sup>0</sup>'Mango', <sup>1</sup>'Banana', <sup>2</sup>'color',  
<sup>3</sup>'is', <sup>4</sup>'red', <sup>5</sup>'yellow']

(i) Mango is red color.

(ii) Banana is yellow color,

(iii) hair has black color.

Now, how we represent each word?

For instance let represent them with one-hot-encoding:

Player	red	red	colours.
1	0	0	0
0	0	0	0
0	0	0	1
0	1	0	0
0	0	1	0
0	0	0	0

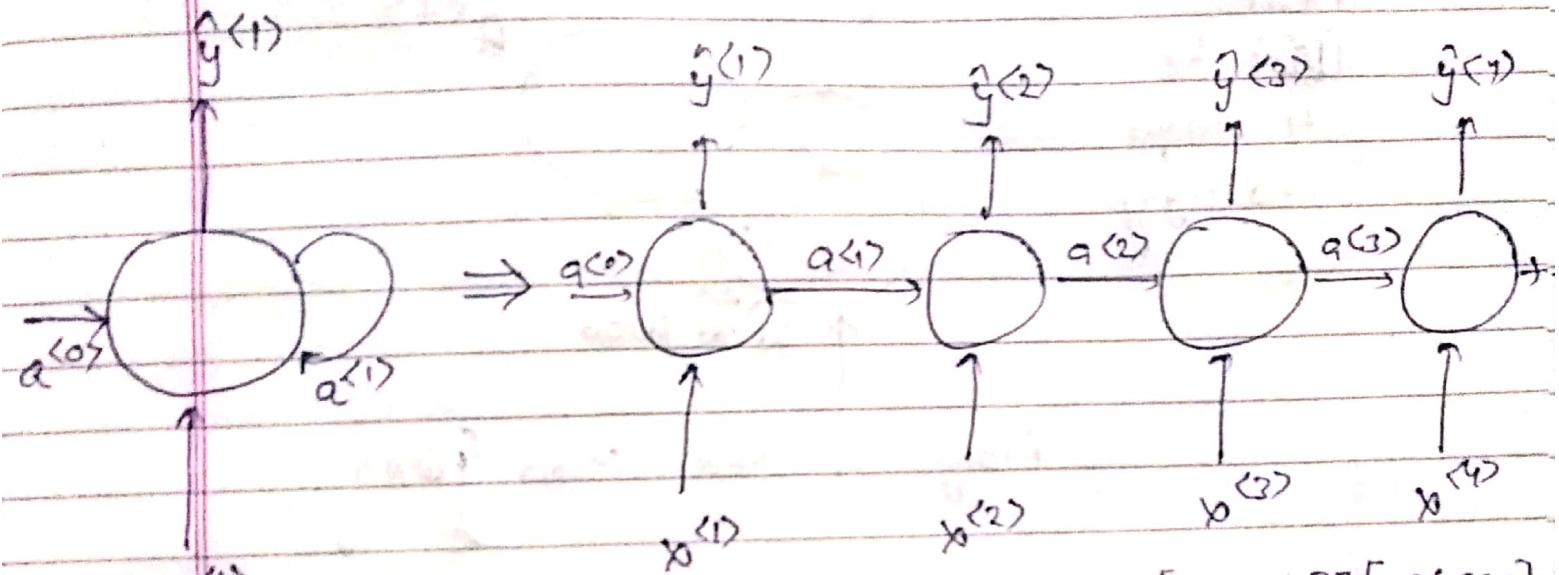
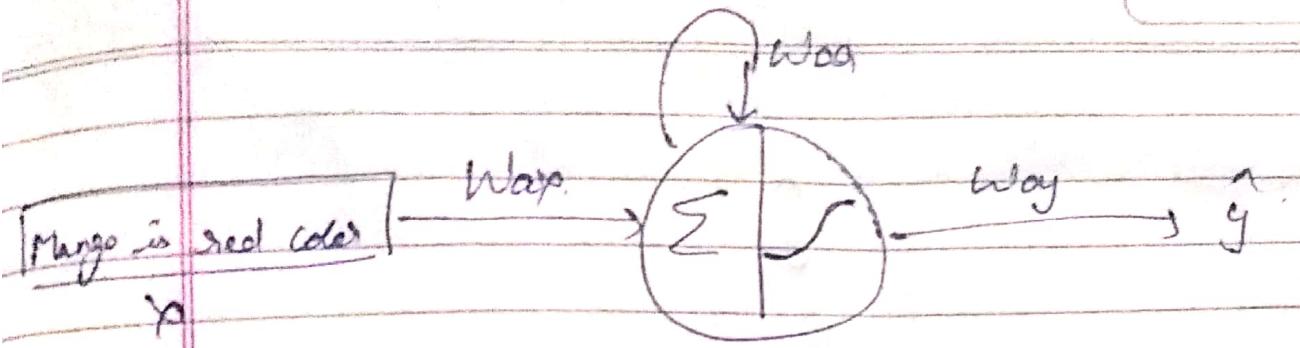
Input size = [vocabulary.len x 1].

For instance, we don't train these sentence that's why there are OHE but after we use pre-trained model which are trained on 100B corpora so then this matrix is not represent by OHE. It is consist some positive values which represent some dependencies.

Remember that weights are shared with all time steps.

Parameter representation:

1. Idea : weight of activation  $\times$  activation
2. ~~target~~
2. Way : weight of input  $\times$  activation
3. Way : weight of activation  $\times$  output.
4. bias : Activation bias -
5. Output bias



Mango [1 0 0 0 0 0]  
is [0 0 0 1 0 0]  
red [0 0 0 0 1 0]  
color [0 0 1 0 0 0]

[1 0 0 0 0 0].7 [0 0 0 1 0 0].7 [0 0 0 0 1 0]7 [0 0 1 0 0 0]7  
Mango is red colour.

(111100)<sup>T</sup>  
1st row  
Anterior n-th size.

How we represent Input:

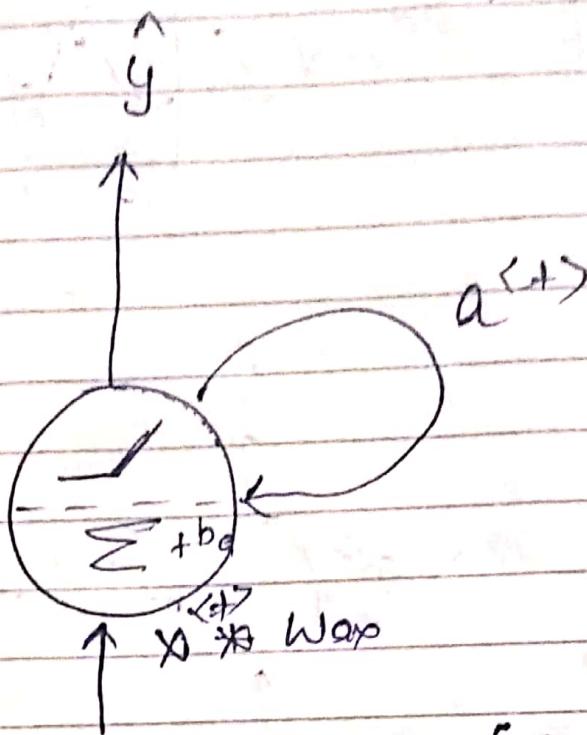
Mango is red colour	1	0	0	0
	0	1	0	0
	0	0	1	0
	0	1	0	0
	0	0	1	0

Vocabulary-Ter

10 x 4 → Sentence-LEN

Here.

(16) Ke  
4 loops  
challengy.



Mango is red colour.  $[10 \times 4]$ .

$\downarrow$        $\downarrow$   
vocab      timestep

1.  $X\_Input = [batch\_size, vocab\_size, sequence/timestep]$

2.  $W_{aa} = [hidden\_unit, hidden\_unit]$

3.  $W_{ap} = [hidden\_unit, hidden\_unit]$

4.  $W_{gy} = [hidden\_unit, output\_unit]$

5.  $b_a = [hidden\_unit, 1]$

6.  $b_y = [output\_unit, 1]$

You can say that number of vocab  
= no. of hidden units.

7.  $y_t\text{-pred} = [ \text{Output-unit}, \text{vocab-size}, \text{timestep} ]$

8.  $a_{\text{prev}} = [ \text{hidden-unit}, \text{vocab-size} ]$

In our example we have 3 sentence.

1. Mango is red color.

2. Banana is yellow color.

3. Cat is black pool.

Vocab-size = [10] # Unique words.

$b_a = (2 \times 1)$   $b_y = (1 \times 1)$

$x_{\text{input}} = (1 \times 10 \times 4)$ .

$w_{ax} = (10 \times 2)$

$a_{\text{prev}} \Rightarrow (2 \times 10) \Leftarrow \text{anext}$

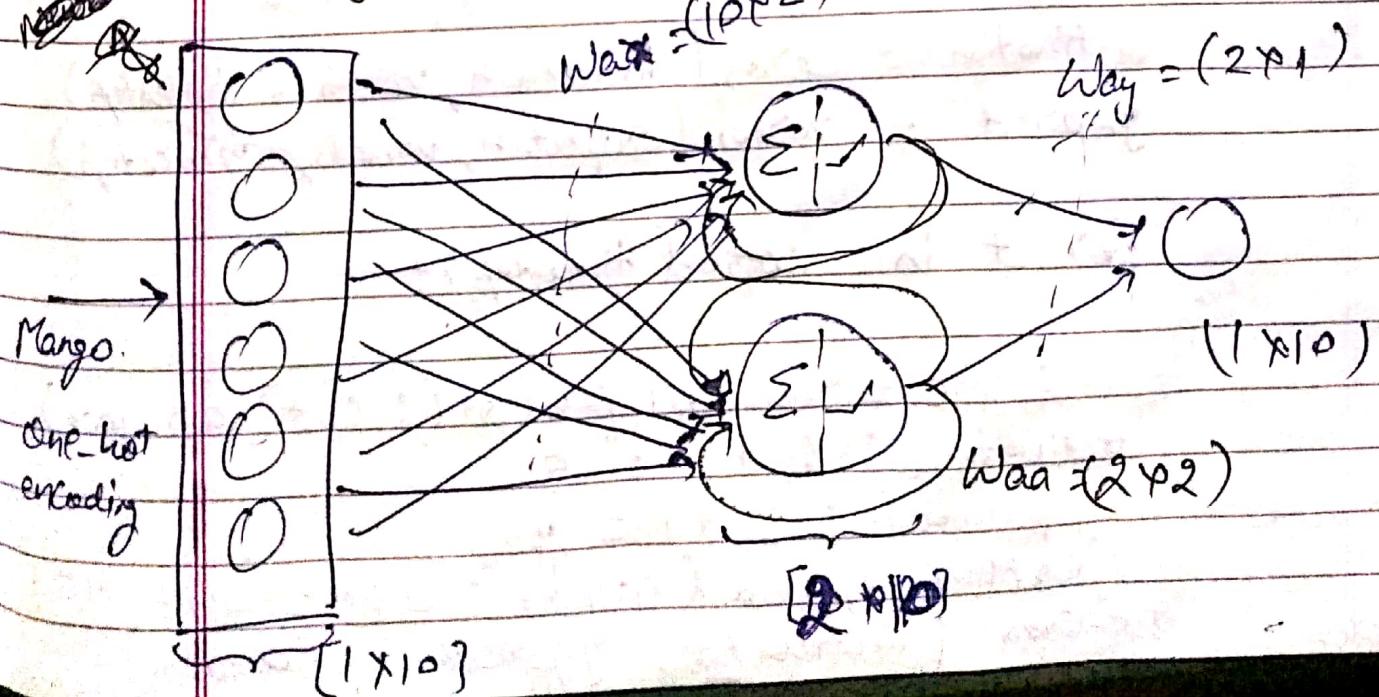
$w_{ay} = (2 \times 1)$

$y_t\text{-pred} = (1, 10 \times 4)$

$l_{\text{at-ay}} = (2 \times 1)$ .

$\text{Waa} = (10 \times 2)$

$\text{Way} = (2 \times 1)$



$$a_{\text{next}} = \tanh(\text{dot}(W_{a,T}, a_{\text{prev}}) + \text{dot}(W_{y,T}, x_t) + b_a)$$

$$y_{\text{op}} = \text{softmax}(\text{dot}(W_{y,T}, a_{\text{next}}) + b_y)$$

def runn-cell-forward(x\_t, a\_prev, parameters):

$$a_{\text{next}} = \tanh(x_t, a_{\text{prev}}, \text{parameters})$$

$$y_{\text{op}} = \text{softmax}(a_{\text{next}})$$

$$\text{cache} = (a_{\text{prev}}, x_t, a_{\text{next}}, \text{parameters})$$

return a\_{\text{next}}, y\_{\text{op}}, \text{cache}

def runn-forward(x\_t, a\_prev, parameters):

$$\text{caches} = []$$

$$\text{activation} = \text{zeros}(\text{hidden\_n}, \text{vocab}, \text{timestep})$$

$$y_{\text{pred}} = \text{zeros}(\text{output\_n}, \text{vocab}, \text{timestep})$$

for t in range(timestep):

$$q, y, c = \text{runn-cell-forw}(x_t[:, :, t], a_p, p_w)$$

$$\text{activation}[:, :, t] = q$$

$$y_{\text{pred}}[:, :, t] = y$$

$$\text{caches.append}(\text{cache}); \text{caches} = (\text{caches}, x_t)$$

return (activation, y\_{\text{pred}}, \text{caches})

Steps in forward propagation we check

We have,  $\mathbf{x}$ , and  $a_{\text{prev}}$ .

$$a_{\text{next}} = \tanh(\text{dot}(\mathbf{Waa}, a_{\text{prev}}) + \text{dot}(\mathbf{Wax}, \mathbf{x}) + \mathbf{ba})$$

$$\hat{y}_{\text{cap}} = \text{softmax}(\text{dot}(\mathbf{Way}, a_{\text{next}}) + \mathbf{by}).$$

## \* Backpropagation through time

height are same for all timesteps in forward propagation.

$$\text{Cost function } J^t(\theta) = - \sum_{j=1}^{|\mathcal{M}|} y_{t,j} \log \hat{y}_{t,j}$$

$$J(\theta) = -\frac{1}{T} \sum_{t=1}^T \sum_{j=1}^{|\mathcal{M}|} y_{t,j} \log \hat{y}_{t,j}$$

where  $M$  = Vocabulary  $J(\theta)$  = cost function

In our example we have 4 time steps.

$$J(\theta) = -\frac{1}{4} \left[ \sum_{j=1}^{10} y_{0,j} \log \hat{y}_{0,j} + \sum_{j=1}^{10} y_{1,j} \log \hat{y}_{1,j} \right. \\ \left. + \sum_{j=1}^{10} y_{2,j} \log \hat{y}_{2,j} + \sum_{j=1}^{10} y_{3,j} \log \hat{y}_{3,j} \right]$$

Actually, we are considering average cost in all time steps.

Let first write all parameters

$\Rightarrow b_a, b_y, W_{ay}, W_{aq}, W_{ay}$

Cache = ( $a^{<+>} \downarrow, a^{<+1>}, y^{<+>}, \text{parameters}$ )

$$\frac{\partial J}{\partial a^{<+1>}} = \frac{\partial J}{\partial a} \cdot \frac{\partial a}{\partial a^{<+1>}}$$

$\frac{\partial J}{\partial a^{<+>}}$	$\frac{\partial a}{\partial w_p}$	$\frac{\partial a}{\partial w_q}$	$\frac{\partial a}{\partial b}$
---------------------------------------	-----------------------------------	-----------------------------------	---------------------------------

$$\frac{\partial J}{\partial a^{<+>}}$$

$$\frac{\partial J}{\partial a^{<+>}} = \frac{\partial J}{\partial a} \cdot \frac{\partial a}{\partial a^{<+>}}$$

$\tanh(-)$

We have  $a$ .

$\tanh(a) \uparrow$

$$a^{<+>} = \text{act} = \tanh(W_p \cdot x + W_q \cdot a_{\text{prev}} + b_a)$$

$$y_{-op} = \text{softmax}(W_y \cdot a_{\text{next}} + b_y),$$

$$J = -\frac{1}{T} \sum_{i=1}^{|M|} \sum_{j=1}^T y_{j,i} \log(y_{-op}_{j,i})$$

and we have to find

$$\frac{\partial J}{\partial \theta} = \frac{\partial J}{\partial (W_y, b_y, W_p, W_q, b_a)} = \frac{\partial J}{\partial y_{-op}} \cdot \frac{\partial y_{-op}}{\partial W_y}$$

$$\frac{\partial J}{\partial y_{-op}} \cdot \frac{\partial J}{\partial (W_y, b_y, W_p, W_q, b_a)} = \frac{\partial J}{\partial y_{-op}} \cdot \frac{\partial y_{-op}}{\partial a_{\text{next}}} \cdot \frac{\partial a_{\text{next}}}{\partial W_p, W_q, b_a}$$

I

$$Q^{(t)} = \tanh(W_{xp}x^{(t)} + W_{aa}a^{(t-1)} + b)$$

$$\frac{\partial \tanh(u)}{\partial u} = \frac{1}{\cosh^2(u)} = \frac{e^u - e^{-u}}{e^u + e^{-u}}$$

$$= \frac{[(e^u + e^{-u})^2 - (e^u - e^{-u})^2]}{(e^u + e^{-u})^2}$$

$$= 1 - \left( \frac{e^u - e^{-u}}{e^u + e^{-u}} \right)^2$$

$$\frac{\partial \tanh(u)}{\partial u} = 1 - \tanh(u)^2$$

$$1. \frac{\partial Q^{(t)}}{\partial W_{xp}} = [1 - \tanh(W_{xp}x^{(t)} + W_{aa}a^{(t-1)} + b)] \cdot x^{(t)}$$

$$2. \frac{\partial Q^{(t)}}{\partial W_{aa}} = [1 - \tanh(W_{xp}x^{(t)} + W_{aa}a^{(t-1)} + b)] \cdot a^{(t-1)}$$

$$3. \frac{\partial Q^{(t)}}{\partial b} = \sum_{\text{Batch}} (1 - \tanh(W_{xp}x^{(t)} + W_{aa}a^{(t-1)} + b))^2$$

Is it over? no! Why because  $a_{\text{prev}}$  is also contain some words from previous activations.

$$4. \frac{\partial Q^{(t)}}{\partial u^{(t)}} = W_{xp} \cdot T (1 - \tanh(W_{xp}x^{(t)} + W_{aa}a^{(t-1)} + b))^2$$

$$5. \frac{\partial e^{(t)}}{\partial a^{(t-1)}} = W_{aa}^T \cdot (1 - \tanh(W_{aa}^{(t)} + W_{ba}^{(t-1)} + b)^2)$$

$\partial \text{tanh} = [\text{hidden-unit}, \text{vocab-size}]$ .

$\partial W_{ap} = [\text{batch-unit}, \text{hidden-unit}]$

$\partial W_{aa} = [\text{hidden-unit}, \text{hidden-unit}]$

$\partial b_a = [\text{hidden-unit}, 1]$

$\partial b_y = [\text{output-unit}, 1]$

$\partial x = [\text{batch-unit}, \text{vocab-size}]$

$\partial a_{\text{prev}} = [\text{hidden-unit}, \text{vocab-size}]$

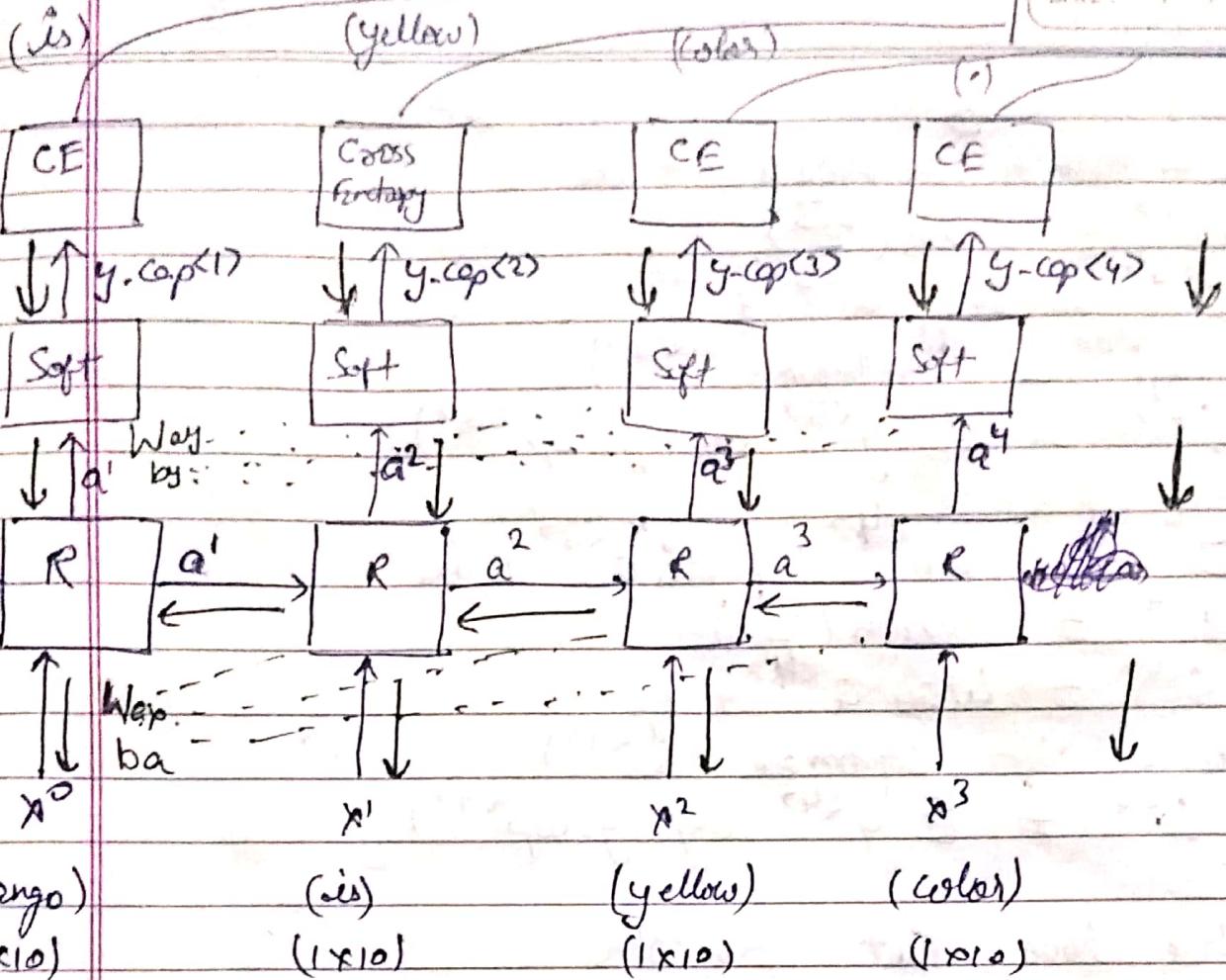
loss for single time step. [for BiClassification]

$$L^{(t)}(\hat{y}^{(t)}, y^{(t)}) = -y^{(t)} \log \hat{y}^{(t)} - (1-y^{(t)}) \log (1-\hat{y}^{(t)})$$

loss for all time step. [for MultiClassification]

$$L(\hat{y}, y) = \sum_{t=1}^T L^{(t)}(\hat{y}^{(t)}, y^{(t)})$$

Total loss  
Page No. \_\_\_\_\_  
Date: / /



CE

Syntex

a Way

$$Z^0 = \text{Way } X^{(0)} + \text{Waa } a^{(0)} + ba,$$

$$a^{(1)} = \tanh(Z^0)$$

$$S^{(1)} = \text{Way} \cdot a^{(1)} + by$$

$$y_{\text{-cap}}^{(1)} = \text{Softmax}(S^{(1)})$$

$$L^{(1)} = \sum y_i^{(1)} \log(y_{\text{-cap}}^{(1)}).$$



$$Z^2 = \text{Way } X^{(2)} + \text{Waa } a^{(1)} + ba,$$

$$a^{(2)} = \tanh(Z^2)$$

$$S^{(2)} = \text{Way} \cdot a^{(2)} + by$$

$$y_{\text{-cap}}^{(2)} = \text{Softmax}(S^{(2)})$$

$$L^{(2)} = \sum y_i^{(2)} \log(y_{\text{-cap}}^{(2)}).$$

$$z^3 = W_{30} X^{(3)} + W_{31} a^{(2)} + b_3$$

$$a^{(3)} = \tanh(z^3)$$

$$s^{(3)} = \text{Way } a^{(3)} + b_4$$

$$y_{\text{-cap}}^{(3)} = \text{Softmax}(s^{(3)})$$

$$L^{(3)} = \sum y^{(3)} \log(y_{\text{-cap}}^{(3)})$$

$$z^4 = W_{40} X^{(4)} + W_{41} a^{(3)} + b_5$$

$$a^{(4)} = \tanh(z^4)$$

$$s^{(4)} = \text{Way } a^{(4)} + b_6$$

$$y_{\text{-cap}}^{(4)} = \text{Softmax}(s^{(4)})$$

$$L^{(4)} = \sum y^{(4)} \log(y_{\text{-cap}}^{(4)})$$

We have cost function

$$J(\theta) = -\frac{1}{m} \sum_{j=1}^m y_j \cdot \log(y_{\text{-cap}}_j)$$

$$J(\theta) = -\frac{1}{T} \sum_{t=1}^T \sum_{j=1}^m y_j^{(t)} \cdot \log(y_{\text{-cap}}_j^{(t)})$$

We have,

$$-\frac{1}{4} \left[ L^{(1)} + L^{(2)} + L^{(3)} + L^{(4)} \right] \leftarrow \text{ANL}$$

Average Negative log likelihood

We know that,

$$\frac{\partial L}{\partial A} = A - y, \quad \frac{\partial L}{\partial \theta} = (A - y) A'$$

from MN Notes -

Here,

$$(i) \frac{\partial L^{(4)}}{\partial w_{ay}} = [y_{cap}^{(4)} - y^{(4)}] \cdot a^{(4)}$$

from here we can clearly see that  $\frac{\partial L}{\partial w_{ay}}$  only depends on the values at the current time steps:  $y_{cap}^{(4)}, y^{(4)}, a^{(4)}$ .

If we have these calculate the gradient of  $w_{ay}$  a simple matrix multiplication

$$dL_{w_{ay}} = \frac{\partial L^{(4)}}{\partial w_{ay}} + \frac{\partial L^{(3)}}{\partial w_{ay}} + \frac{\partial L^{(2)}}{\partial w_{ay}} + \frac{\partial L^{(1)}}{\partial w_{ay}}.$$

This is change in  $L^{(4)}$  w.r.t  $w_{ay}$  at each time steps.

$$(ii) \frac{\partial L^{(4)}}{\partial w_{aa}} = \frac{\partial L^{(4)}}{\partial s^{(4)}} \cdot \frac{\partial s^{(4)}}{\partial a^{(4)}} \cdot \frac{\partial a^{(4)}}{\partial z^{(4)}} \cdot \frac{\partial z^{(4)}}{\partial w_{aa}}.$$

Now, note that we can't treat  $a^{(3)}$  as an constant because it depends on  $a^{(2)}$  which also depends on  $a^{(1)}, a^{(0)}$ . Then what we have.

$$\frac{\partial L^{(4)}}{\partial w_{aa}} = \sum_{k=1}^4 \frac{\partial L^{(4)}}{\partial a^{(k)}} \cdot \frac{\partial a^{(k)}}{\partial Q^{(k)}} \cdot \frac{\partial Q^{(k)}}{\partial w_{aa}}.$$

~~At  $t = 4$~~

$$= \left[ \frac{JL^{(4)}}{Ja^{(4)}} \cdot \frac{Ja^{(4)}}{Ja^{(1)}} \cdot \frac{Ja^{(1)}}{JWaa} + \frac{JL^{(4)}}{Ja^{(4)}} \cdot \frac{Ja^{(4)}}{Ja^{(2)}} \cdot \frac{Ja^{(2)}}{JWba} \right. \\ \left. + \frac{JL^{(4)}}{Ja^{(4)}} \cdot \frac{Ja^{(4)}}{Ja^{(3)}} \cdot \frac{Ja^{(3)}}{JWba} + \frac{JL^{(4)}}{Ja^{(4)}} \cdot \frac{Ja^{(4)}}{Ja^{(1)}} \cdot \frac{Ja^{(1)}}{Wba} \right]$$

$$K=2 = \left[ \frac{JL^{(4)}}{Ja^{(4)}} \cdot \frac{Ja^{(4)}}{Ja^{(3)}} \cdot \frac{Ja^{(3)}}{Ja^{(2)}} \cdot \frac{Ja^{(2)}}{JWba} + \right.$$

$$K=1 \quad \left[ \frac{JL^{(4)}}{Ja^{(4)}} \cdot \frac{Ja^{(4)}}{Ja^{(3)}} \cdot \frac{Ja^{(3)}}{Ja^{(2)}} \cdot \frac{Ja^{(2)}}{Ja^{(1)}} \cdot \frac{Ja^{(1)}}{JWba} + \right.$$

$$K=3 \quad \left[ \frac{JL^{(4)}}{Ja^{(4)}} \cdot \frac{Ja^{(4)}}{Ja^{(3)}} \cdot \frac{Ja^{(3)}}{Ja^{(2)}} + \right]$$

$$K=4 \quad \left. \frac{JL^{(4)}}{Ja^{(4)}} \cdot \frac{Ja^{(4)}}{JWba} \right]$$

We sum up the contributions of each time step to the gradient. because  $W_{aa}$  is used in every step up to the output we care about, we need to backpropagate gradients from  $t=4$  through the network all the way to  $t=0$ .

$$dW_{aa} = \frac{\partial L^{(4)}}{\partial W_{aa}} + \frac{\partial L^{(3)}}{\partial W_{aa}} + \frac{\partial L^{(2)}}{\partial W_{aa}} + \frac{\partial L^{(1)}}{\partial W_{aa}}$$

$$\star \star \Rightarrow \frac{\partial L^{(T)}}{\partial W_{aa}} = \sum_{k=1}^T \frac{\partial L^{(T)}}{\partial a^{(T)}} \left( \prod_{j=k+1}^T \frac{\partial a^{(j)}}{\partial a^{(j-1)}} \right) \cdot \frac{\partial a^{(k)}}{\partial W_{aa}}$$

This is the reason of exploding gradient problem.

$\prod_{j=k+1}^T \frac{\partial a^{(j)}}{\partial a^{(j-1)}}$  the derivative of a vector w.r.t vector, the result is a Jacobian Matrix.

This product can shrink to zero or explode to infinite

When it goes explode or Vanish?

Let,  $\lambda_1$  = largest singular value of  $W_{aa}$  ?

If you use tanh

If  $\lambda < 1$  : Vanishing gradient.

If  $\lambda > 1$  : Exploding gradient.

For sigmoid

If  $\lambda < 0.25$  : Vanishing gradient.

If  $\lambda > 0.25$  : Explode gradient.

One solution of this is Truncate backpropagate. Skip some timesteps.

1.  $\Rightarrow$  Using an L1 and L2 norm penalty on recurrent weights can help with exploding gradient. The L1 and L2 term can ensure that during training  $\lambda$  stay smaller than 1. and in this regime gradient can not explode. also, it help in vanishing problem. It ensure that output of sigmoid or any function lies in linear region i.e. its derivative also not too high or too low.

2. ~~\*~~ Scaling down the gradient [clipping]

It reduce Exploding problem.

$$\hat{g} = \frac{\partial L}{\partial w_{aa}}$$

If  $\|\hat{g}\| \geq \text{threshold}$  then

$$\hat{g} = \frac{\text{threshold}}{\|\hat{g}\|} \cdot \hat{g}$$

end if

What we do here, we clipping the norm of the exploded gradient.

At  $t = 2$ .

$$\frac{\partial L^{(2)}}{\partial w_{aa}} = \left[ \frac{\partial L^{(2)}}{\partial a^{(2)}} \cdot \frac{\partial a^{(2)}}{\partial a^{(1)}} \cdot \frac{\partial a^{(1)}}{\partial w_{aa}} + \right. \\ \left. \frac{\partial L^{(2)}}{\partial a^{(2)}} \cdot \frac{\partial a^{(2)}}{\partial a^{(1)}} \right].$$

$$(i) \quad \frac{\partial L^{(t)}}{\partial a^{(t)}} = \frac{\partial L^{(t)}}{\partial S^{(t)}} \cdot \frac{\partial S^{(t)}}{\partial a^{(t)}}$$

$$= [y_{-cap}^{(t)} - y^{(t)}] \cdot w_{ay}.$$

$$(ii) \quad \frac{\partial a^{(t)}}{\partial a^{(t-1)}} = \left[ 1 - \tanh(w_{xa}x^{(t)} + w_{aa}a^{(t-1)} + b_a)^2 \right] w_{aa}$$

$$(iii) \quad \frac{\partial a^{(t)}}{\partial w_{aa}} = \left[ 1 - \tanh(w_{xa}x^{(t)} + w_{aa}a^{(t-1)} + b_a)^2 \right] a^{(t-1)}.$$

$$(iii') \quad \frac{\partial L^{(4)}}{\partial w_{ax}} = \frac{\partial L^{(4)}}{\partial a^{(4)}} \cdot \frac{\partial a^{(4)}}{\partial w_{ax}}$$

$$= (y_{-cap}^{(4)} - y^{(4)}) \cdot w_{ay} \cdot x^{(4)}.$$

Here, also we cannot treat  $x^{(4)}$  as a constant because it varies at each time step.

$$\frac{JL^{(4)}}{Wa_x} = \sum_{K=1}^4 \frac{JL^{(4)}}{Ja^{(4)}} \cdot \frac{Ja^{(4)}}{Ja^K} \cdot \frac{Ja^K}{JWap}$$

At  $t = 2$

$$\frac{JL^{(2)}}{Wa_x} = \sum_{K=1}^2 \frac{JL^{(2)}}{Ja^{(2)}} \cdot \frac{Ja^{(2)}}{Ja^{(K)}} \cdot \frac{Ja^K}{JWap}$$

$$= \left[ \frac{JL^{(2)}}{Ja^{(2)}} \cdot \frac{Ja^{(2)}}{Ja^{(1)}} \cdot \frac{Ja^{(1)}}{JWap} + \right.$$

$$\left. \frac{JL^{(2)}}{Ja^{(2)}} \cdot \frac{Ja^{(2)}}{JWap} \right]$$

$$(ii) \frac{Ja^{(t)}}{JWap} = \left[ 1 - \tanh(Wap x^{(t)} + Wa_x a^{(t-1)} + ba) \right] x^{(t)}$$

$$dh_{Wap} \Rightarrow \frac{JL}{JWap} = \frac{JL^{(1)}}{JWap} + \frac{JL^{(2)}}{JWap} + \frac{JL^{(3)}}{JWap} + \frac{JL^{(4)}}{JWap}$$

$$\frac{JL^{(T)}}{JWap} = \sum_{K=1}^T \frac{JL^{(T)}}{Ja^{(T)}} \left( \prod_{j=K+1}^T \frac{Ja^{(j)}}{Ja^{(j-1)}} \right) \cdot \frac{Ja^{(K)}}{JWap}$$

$$(iv) \frac{JL^{(4)}}{Jba} = \frac{JL^{(4)}}{Ja^{(4)}} \cdot \frac{Ja^{(4)}}{Jba}$$

$$= [y_{-cap}^{(4)} - y^{(4)}] \cdot Way \cdot \frac{Ja^{(4)}}{Jba}$$

$$\frac{J_a^{(4)}}{J_{ba}} = \left[ i - \tanh(W_{ab}x^{(4)} + W_{aa}a^{(3)} + b_a) \right]^2.$$

$$d_{ba} = \frac{J_L^{(1)}}{J_{ba}} + \frac{J_L^{(2)}}{J_{ba}} + \frac{J_L^{(3)}}{J_{ba}} + \frac{J_L^{(4)}}{J_{ba}}.$$

$$(V) \quad \frac{J_L^{(4)}}{J_{by}} = \frac{J_L^{(4)}}{J_S^{(4)}} \cdot \frac{J_S^{(4)}}{J_{by}}$$

$$= [y_{cap}^{(4)} - y^{(4)}].$$

$$d_{by} = \sum_{i=0}^T \frac{J_L^{(i)}}{J_{ba}}.$$

$$(VI) \quad \frac{\partial L^{(4)}}{J_x^{(4)}} = \frac{J_L^{(4)}}{J_a^{(4)}} \cdot \frac{J_a^{(4)}}{J_x^{(4)}},$$

When we dealing with one to many structure then input is dependent on previous timestep output's.

$$d_{xt} \Rightarrow \frac{J_L^{(T)}}{J_x} = \frac{J_L^{(1)}}{J_x^{(1)}} + \frac{J_L^{(2)}}{J_x^{(2)}} + \frac{J_L^{(3)}}{J_x^{(3)}} + \frac{J_L^{(4)}}{J_x^{(4)}}$$

$$\frac{J_L^{(T)}}{J_x} = \sum_{K=1}^T \frac{J_L^{(T)}}{J_a^{(T)}} \left( \prod_{j=k+1}^T \frac{J_a^{(j)}}{J_a^{(j-1)}} \right) \cdot \frac{J_a^{(k)}}{J_x^{(k)}}$$

$$(Vip) \frac{\delta L^{(4)}}{\delta a_0} = \frac{\delta L^{(4)}}{\delta a^{(4)}} \cdot \frac{\delta a^{(4)}}{\delta a^{(3)}} \cdot \frac{\delta a^{(3)}}{\delta a^{(2)}} \cdot \frac{\delta a^{(2)}}{\delta a^{(1)}}$$

We also need to get optimal value of  $a_{prev}$  or  $a_0$  because we initialize it as random.

$$\delta a_0 = \frac{\delta L^{(1)}}{\delta a_0} + \frac{\delta L^{(2)}}{\delta a_0} + \frac{\delta L^{(3)}}{\delta a_0} + \frac{\delta L^{(4)}}{\delta a_0}$$

$$\frac{\delta L^{(1)}}{\delta a_0} = \frac{\delta L^{(1)}}{\delta a^{(1)}} \cdot \frac{\delta a^{(1)}}{\delta a^{(0)}}$$

$$\frac{\delta L^{(1)}}{\delta a^{(1)}} = [y - a_p^{(1)} - y^{(1)}] \cdot w_{11}$$

$$\frac{\delta a^{(1)}}{\delta a^{(0)}} = [1 - \tanh(w_{01}p^{(1)} + w_{02}a^{(0)} + b_1)]^2 w_{11}$$