

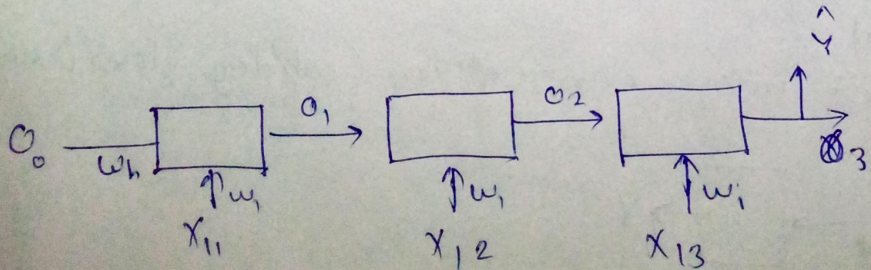
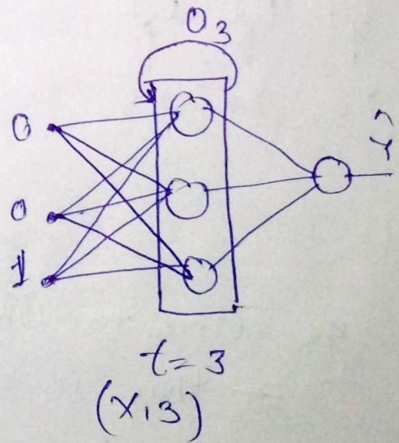
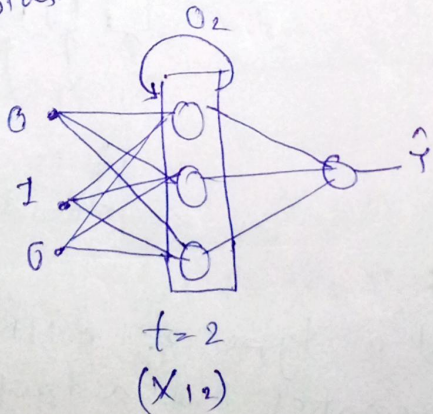
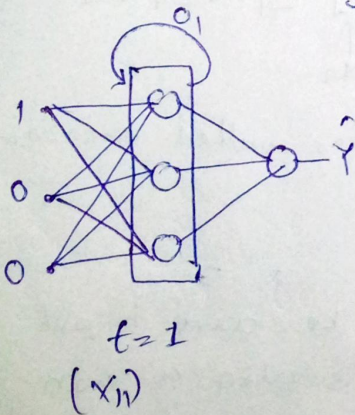
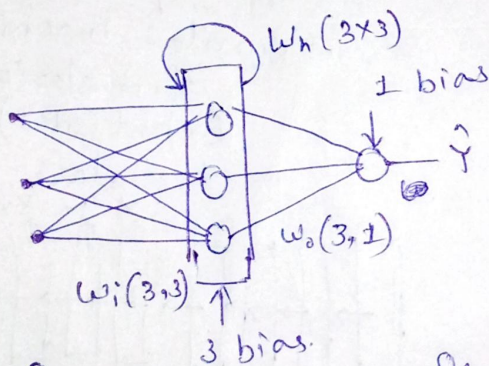
18 Oct 2024

BACK PROPAGATION IN RNN (BPTT)

- BPTT (Backpropagation Through Time)
- Here we take a Many to one RNN Example
- Example:

text	output
cat mat rat	1
rat rat mat	1
mat mat cat	0

	X	Y
x_1	[1 0 0] [0 1 0] [0 0 1]	1
x_2	[0 0 1] [0 0 1] [0 1 0]	1
x_3	[0 1 0] [0 1 0] [1 0 0]	0



$$o_1 = f(x_{11}w_i + o_0w_h)$$

$$o_2 = f(x_{12}w_i + o_1w_h)$$

$$o_3 = f(x_{13}w_i + o_2w_h)$$

$$\hat{y} = \text{sigmoid}(o_3 \cdot w_o)$$

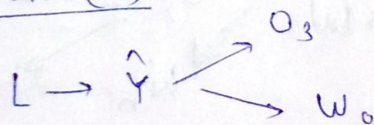
$$\text{loss}(L) = -y_i \log \hat{y}_i - (1 - y_i) \log (1 - \hat{y}_i)$$

- now minimize the loss using gradient descent where w_i, w_h, w_o values (optimal) where the $\text{loss}(L)$ is minimum by weights update

$$w_i = w_i - \eta \frac{\partial L}{\partial w_i} \quad \text{--- (i)} \quad w_h = w_h - \eta \frac{\partial L}{\partial w_h} \quad \text{--- (ii)}$$

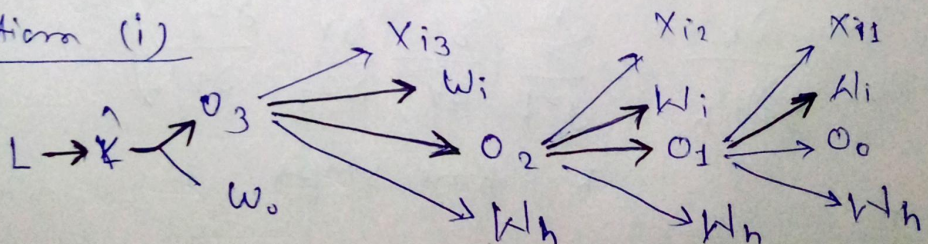
$$w_o = w_o - \eta \frac{\partial L}{\partial w_o} \quad \text{--- (iii)}$$

* in equation (iii)



$$\frac{\partial L}{\partial w_o} = \frac{\partial L}{\partial \hat{y}} \cdot \frac{\partial \hat{y}}{\partial w_o}$$

* in equation (i)



$$\begin{aligned} \frac{\partial L}{\partial w_i} &= \frac{\partial L}{\partial \hat{y}} \cdot \frac{\partial \hat{y}}{\partial o_3} \cdot \frac{\partial o_3}{\partial w_i} \\ \frac{\partial L}{\partial w_i} &= \frac{\partial L}{\partial \hat{y}} \cdot \frac{\partial \hat{y}}{\partial o_3} \cdot \frac{\partial o_3}{\partial o_2} \cdot \frac{\partial o_2}{\partial w_i} \end{aligned}$$

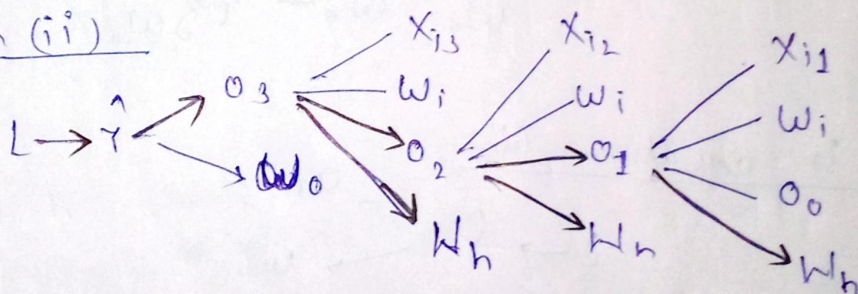
$$\frac{\partial L}{\partial w_i} = \frac{\partial L}{\partial \hat{y}} \cdot \frac{\partial \hat{y}}{\partial o_3} \cdot \frac{\partial o_3}{\partial w_i} + \frac{\partial L}{\partial \hat{y}} \cdot \frac{\partial \hat{y}}{\partial o_3} \cdot \frac{\partial o_3}{\partial o_2} \cdot \frac{\partial o_2}{\partial w_i} + \frac{\partial L}{\partial \hat{y}} \cdot \frac{\partial \hat{y}}{\partial o_3} \cdot \frac{\partial o_3}{\partial o_2} \cdot \frac{\partial o_2}{\partial o_1} \cdot \frac{\partial o_1}{\partial w_i}$$

OR

$$\frac{\partial L}{\partial w_i} = \sum_{j=1}^n \left(\frac{\partial L}{\partial \hat{y}} \cdot \frac{\partial \hat{y}}{\partial o_j} \cdot \frac{\partial o_j}{\partial w_i} \right)$$

$n = \text{time steps}$
~~in general case~~

* in equation (ii)



$$\frac{\partial L}{\partial w_h} = \frac{\partial L}{\partial \hat{y}} \cdot \frac{\partial \hat{y}}{\partial o_3} \cdot \frac{\partial o_3}{\partial w_h} + \frac{\partial L}{\partial \hat{y}} \cdot \frac{\partial \hat{y}}{\partial o_3} \cdot \frac{\partial o_3}{\partial o_2} \cdot \frac{\partial o_2}{\partial w_h} + \frac{\partial L}{\partial \hat{y}} \cdot \frac{\partial \hat{y}}{\partial o_3} \cdot \frac{\partial o_3}{\partial o_2} \cdot \frac{\partial o_2}{\partial o_1} \cdot \frac{\partial o_1}{\partial w_h}$$

OR

$$\frac{\partial L}{\partial w_h} = \sum_{j=1}^n \left(\frac{\partial L}{\partial \hat{y}} \cdot \frac{\partial \hat{y}}{\partial o_j} \cdot \frac{\partial o_j}{\partial w_h} \right)$$

$n = \text{time steps}$

PROBLEM WITH RNN

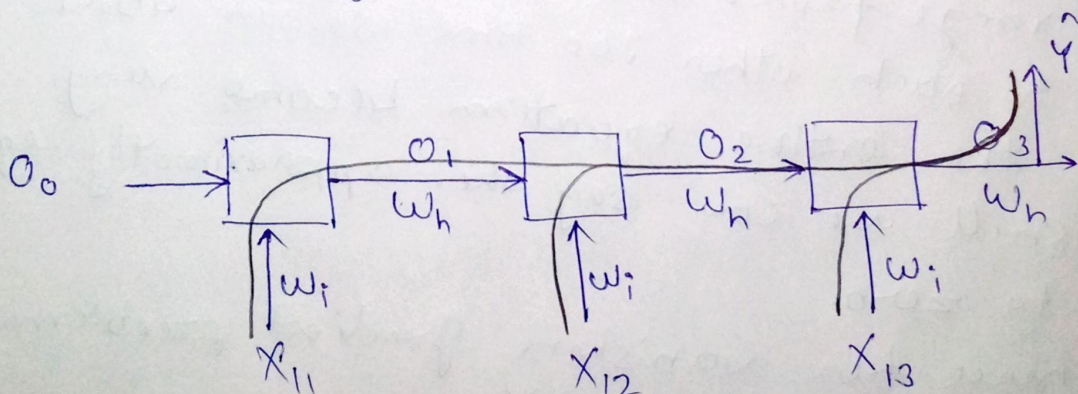
problem with RNN:

- problem of long term dependency.
- unstable Training/stagnated training.
- the above two problems are caused due to unstable gradients

① Problem of long term dependency:

- this problem is caused due to vanishing gradient
- $\frac{\partial L}{\partial w_i}$ depends on three terms.

$$\frac{\partial L}{\partial w_i} = \underbrace{\frac{\partial L}{\partial \hat{y}} \cdot \frac{\partial \hat{y}}{\partial o_3} \cdot \frac{\partial o_3}{\partial w_i}}_{\text{short term dependency}} + \underbrace{\frac{\partial L}{\partial \hat{y}} \cdot \frac{\partial \hat{y}}{\partial o_3} \cdot \frac{\partial o_3}{\partial o_2} \cdot \frac{\partial o_2}{\partial w_i}}_{\text{long term dependency}} + \underbrace{\frac{\partial L}{\partial \hat{y}} \cdot \frac{\partial \hat{y}}{\partial o_3} \cdot \frac{\partial o_3}{\partial o_2} \cdot \frac{\partial o_2}{\partial o_1} \cdot \frac{\partial o_1}{\partial w_i}}_{\text{long term dependency}}$$



- for more number of time steps the number of terms in $\frac{\partial L}{\partial w_i}$ is very long due to that. long term dependency accuracy (becomes small)

- hence the calculating the $\frac{\partial L}{\partial w_i}$ the major contribution is on short term dependency.

for very large time step.

$$\frac{\partial L}{\partial w_i} = \frac{\partial L}{\partial \hat{y}} \cdot \frac{\partial \hat{y}}{\partial o_{100}} \prod_{t=2}^{100} \left(\frac{\partial o_t}{\partial o_{t-1}} \right) \frac{\partial o_1}{\partial w_i}$$

where

$$o_t = \tanh(x_{it} w_i + o_{t-1} w_h)$$

$$\frac{\partial o_t}{\partial o_{t-1}} = \tanh'(x_{it} w_i + o_{t-1} w_h) w_h$$

↓
range 0 to 1

hence,

$$\frac{\partial L}{\partial w_i} = \frac{\partial L}{\partial \hat{y}} \cdot \frac{\partial \hat{y}}{\partial o_{100}} \prod_{t=2}^{100} \left(\tanh'(x_{it} w_i + o_{t-1} w_h) w_h \right) \frac{\partial o_1}{\partial w_i}$$

- $\tanh'()$ value range from 0 to 1 and w_h value range from 0 to 1 and they are multiply to each other 100 time. so. in above the entire equation become very small we can say the. approximately equal to zero.
- hence the vanishing gradient problem occur.

Solutions to resolve this problem:

- Use different activation functions.
- Better weight initialization, example.
we can use the identity matrix
- different type of RNN (skip RNN)
- use LSTM

② Unstable Training. (Exploding gradients):

- Here the long term dependency becomes very large. as compared to short term dependency

- This problem occurs may be due to very large learning rate
- use ReLU activation function with $W_h = 1$ initialisation.

Solutions to resolve this problem:

- use gradient clipping
- controlled learning rate
- use LSTM