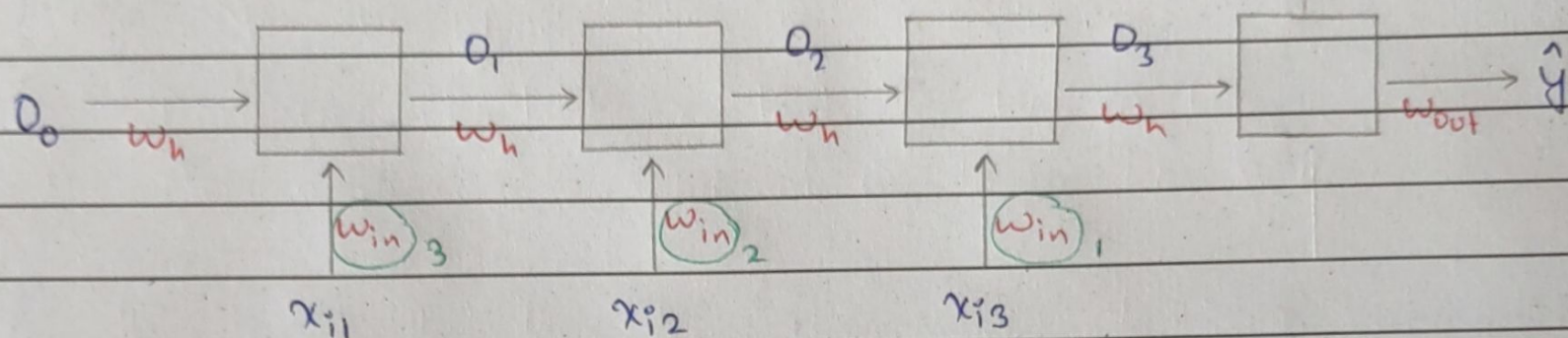




PROBLEMS WITH RNN



Problem-1: Long Term Dependency Problem

$$w_{in} = w_{in} - \eta \frac{\partial L}{\partial w_{in}} \quad \left. \frac{\partial L}{\partial w_{in}} = \frac{\partial L}{\partial \hat{y}} \times \frac{\partial \hat{y}}{\partial O_3} \times \frac{\partial O_3}{\partial w_{in1}} \right\} \text{Short term dependency}$$

$$w_{out} = w_{out} - \eta \frac{\partial L}{\partial w_{out}} \quad + \frac{\partial L}{\partial \hat{y}} \times \frac{\partial \hat{y}}{\partial O_3} \times \frac{\partial O_3}{\partial O_2} \times \frac{\partial O_2}{\partial w_{in2}}$$

$$w_h = w_h - \eta \frac{\partial L}{\partial w_h} \quad + \frac{\partial L}{\partial \hat{y}} \times \frac{\partial \hat{y}}{\partial O_3} \times \frac{\partial O_3}{\partial O_2} \times \frac{\partial O_2}{\partial O_1} \times \frac{\partial O_1}{\partial w_{in3}}$$

Long term dependency

This was for only 3 time steps. Now assume there are 100 time steps.

$$\Rightarrow \frac{\partial L}{\partial w_{in}} = \frac{\partial L}{\partial \hat{y}} \times \frac{\partial \hat{y}}{\partial O_{100}} \times \dots \times \frac{\partial O_1}{\partial w_{in}} \quad \left. \right\} \text{Long term dependency for 100 time steps}$$

$$\Rightarrow \frac{\partial L}{\partial w_{in}} = \frac{\partial L}{\partial \hat{y}} \times \frac{\partial \hat{y}}{\partial O_{100}} \times \prod_{t=100}^2 \left(\frac{\partial O_t}{\partial O_{t-1}} \right) \times \frac{\partial O_1}{\partial w_{in}}$$

$$O_t = \tanh(x_{it} w_{in} + O_{t-1} w_h)$$

$$\Rightarrow O_t = \tanh(x_{it} w_{in} + O_{t-1} w_h)$$

$$\Rightarrow \frac{\partial O_t}{\partial O_{t-1}} = \underbrace{\tanh'(x_{it} w_{in} + O_{t-1} w_h)}_{\text{Derivative of tanh is always b/w 0-1}} \cdot w_h$$

Derivative of tanh is
always b/w 0-1



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

Derivative of \tanh , i.e. \tanh' is b/w 0-1

Now also assume value of w_h to be b/w 0-1

In this case the long term dependency $\frac{\partial L}{\partial w_{in}}$ will become a vanishing gradient & will not contribute to a short term dependencies will have more responsibility in updating the parameters during backpropagation.

Problem-2: Unstable Training Problem

This usually occurs because of exploding gradient problem. Suppose you use ReLU instead of tanh in your RNN & initialized your weight recurrent weight w_h with 1, in that case your long term dependency will explode & dominate other dependencies during backpropagation due to which unstable gradient training may happen. Similar thing can happen if you have a large value learning rate η .