# 1) RNN Language Model
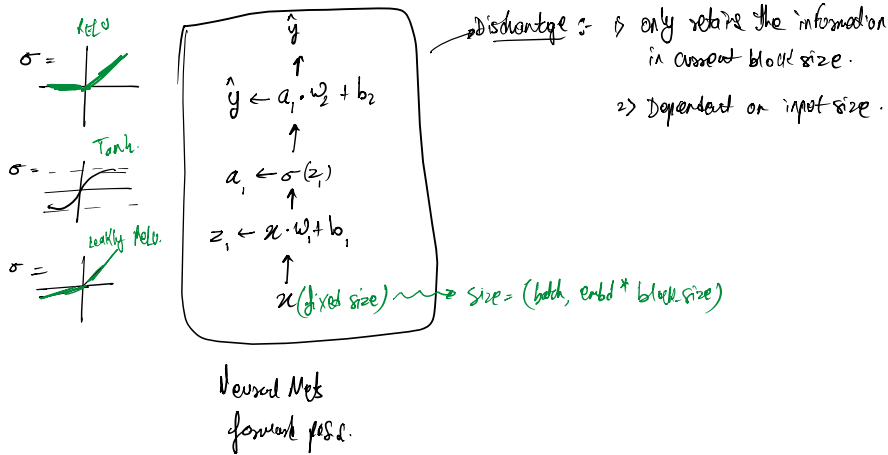
❑ **Architecture of a traditional RNN** — Recurrent neural networks, also known as RNNs, are a class of neural networks that allow previous outputs to be used as inputs while having hidden states. They are typically as follows:

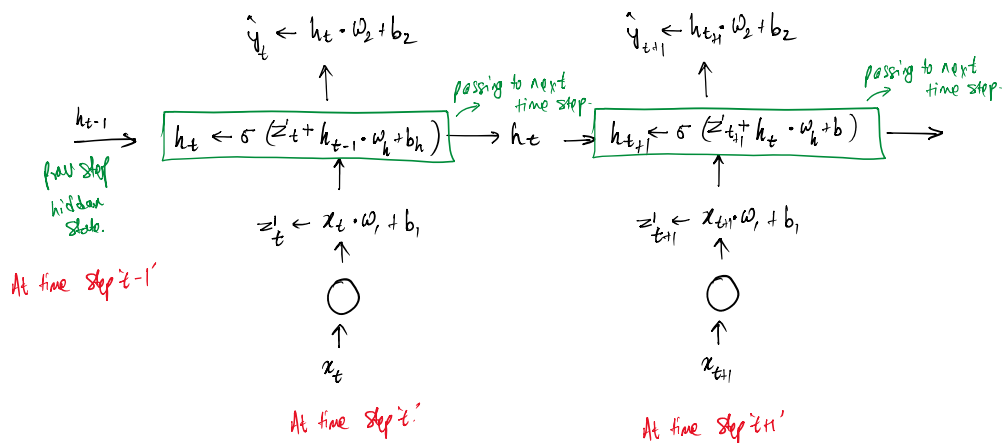| Advantages | Drawbacks |
|---|---|
| • Possibility of processing input of any length<br>• Model size not increasing with size of input<br>• Computation takes into account historical information<br>• Weights are shared across time | • Computation being slow<br>• Difficulty of accessing information from a long time ago<br>• Cannot consider any future input for the current state |

## MULTI-LAYER PERCEPTRON

RELU

$\sigma =$

Tanh

$\sigma =$

leakly ReLu

$\sigma =$

$\hat{y}$

$\hat{y} \leftarrow a_1 \cdot \omega_2 + b_2$

$a_1 \leftarrow \sigma(z_1)$

$z_1 \leftarrow x \cdot \omega_1 + b_1$

$x$ (fixed size) $\rightarrow$ size = (batch, embd * block_size)

Neural Nets
forward pass.

disadvantage :- 1) only retains the information in current block size.

2) Dependent on input size.

# Solution :-
Recurrent Neural Networks.
(RNN).

# RNN. → Transfer everything you learned to next hidden states time step

RNN - Architecture

$\hat{y}_t \leftarrow h_t \cdot \omega_2 + b_2$

$\hat{y}_{t+1} \leftarrow h_{t+1} \cdot \omega_2 + b_2$

passing to next time step.

$h_{t-1}$ → prev step hidden state.

$h_t \leftarrow \sigma(z'_t + h_{t-1} \cdot \omega_h + b_h)$ → $h_t$ → $h_{t+1} \leftarrow \sigma(z'_{t+1} + h_t \cdot \omega_h + b)$ →

passing to next time step.

$z'_t \leftarrow x_t \cdot \omega_1 + b_1$

$z'_{t+1} \leftarrow x_{t+1} \cdot \omega_1 + b_1$

$x_t$

$x_{t+1}$

At time step 't-1'

At time step 't'

At time step 't+1'

Basically we have a function, performing

$y_t \leftarrow f(\text{current-input}, \text{prev-hidden-state})$

$\Rightarrow \hat{y}_t \leftarrow f(x_t, h_{t-1})$

$\hat{y}_t \leftarrow (\sigma(x_t \cdot \omega_1 + h_{t-1} \cdot \omega_h + b)) \cdot \omega_2 + b_2$   # forward pass for RNN

$$f(x_t, h_{t-1}) = (\sigma(x_t \cdot \omega_1 + h_{t-1} \cdot \omega_h + b)) \cdot \omega_2 + b_2$$
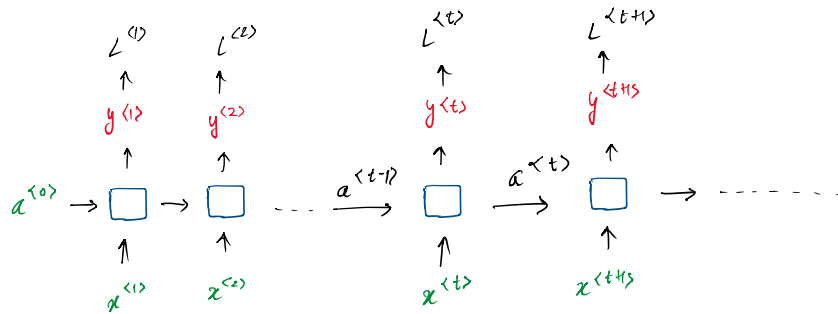
$$f(x_t, h_{t-1}) = \left(\sigma\left(x_t \cdot \omega_1 + h_{t-1} \cdot \omega_h + b\right)\right) \cdot \omega_2 + b_2$$

Activation function

information obtained from past hidden state

(past time step)

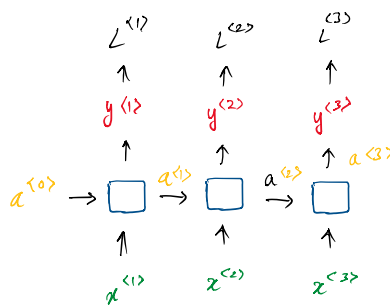Controls how much the past hidden state is important.

# Generally



$L^{<1>}$  $L^{<2>}$   $L^{<t>}$   $L^{<t+1>}$

$y^{<1>}$  $y^{<2>}$   $y^{<t>}$   $y^{<t+1>}$

$a^{<0>}$ → □ → □ --- $a^{<t-1>}$ □ $a^{<t>}$ □ → -----

$x^{<1>}$  $x^{<2>}$   $x^{<t>}$   $x^{<t+1>}$

where,

$$a^{<t>} = g_1\left(x^{<t>} \cdot \omega_x + a^{<t-1>} \cdot \omega_a + b\right)$$

$$y^{<t>} = g_2\left(a^{<t>} \cdot \omega_y + b_y\right)$$

$$L(\hat{y}, y) = \sum_{i=0}^{t} L^{<i>} \quad ; \quad L^{<t>} = L\left(y^{<t>}, \hat{y}\right)$$

# Training RNNs (Backpropagation)

Assume this Architecture.



$L^{<1>}$  $L^{<2>}$  $L^{<3>}$

$y^{<1>}$  $y^{<2>}$  $y^{<3>}$

$a^{<0>}$ → □ $a^{<1>}$ → □ $a^{<2>}$ → □ $a^{<3>}$

$x^{<1>}$  $x^{<2>}$  $x^{<3>}$

$$L = L^{<1>} + L^{<2>} + L^{<3>}$$

$$a^{<1>} = x^{<1>} \cdot \omega_x + a^{<0>} \cdot \omega_h + b$$

$$a^{<2>} = x^{<2>} \cdot \omega_x + a^{<1>} \cdot \omega_h + b$$
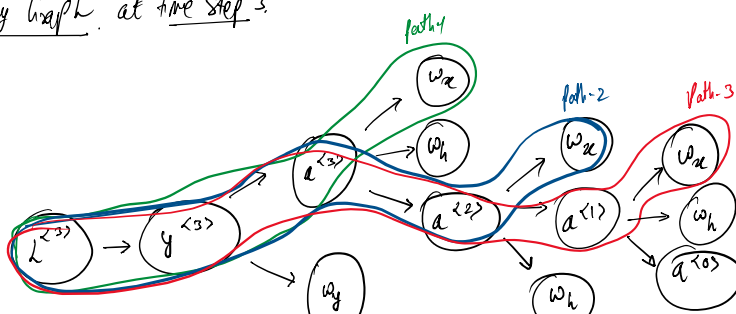
$$a^{<3>} = x^{<3>} \cdot \omega_x + a^{<2>} \cdot \omega_h + b$$

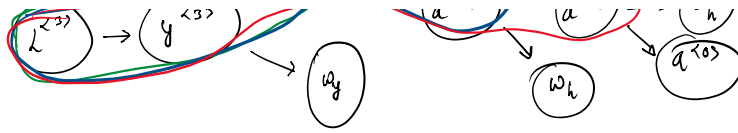$$y^{<1>} = a^{<1>} \cdot \omega_y + b_y$$

$$y^{<2>} = a^{<2>} \cdot \omega_y + b_y$$

$$y^{<3>} = a^{<3>} \cdot \omega_y + b_y$$

# Dependency Graph at time step 3



Path-1  Path-2  Path-3

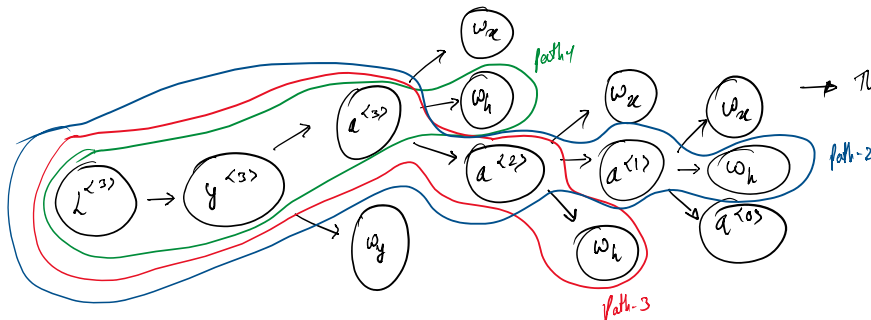→ This tells us how $a^{<3>} \to a^{<2>} \to a^{<1>}$ at time step = 3

Gradients :-

$$\frac{\partial L^{(3)}}{\partial \omega_{x}} = \left(\frac{\partial L^{(3)}}{\partial y^{(3)}} \times \frac{\partial y^{(3)}}{\partial a^{(3)}} \times \frac{\partial a^{(3)}}{\partial \omega_{x}}\right)_{\text{path-1}} + \left(\frac{\partial L^{(3)}}{\partial y^{(3)}} \times \frac{\partial y^{(3)}}{\partial a^{(3)}} \times \frac{\partial a^{(3)}}{\partial a^{(2)}} \times \frac{\partial a^{(2)}}{\partial \omega_{x}}\right)_{\text{path-2}} + \left(\frac{\partial L^{(3)}}{\partial y^{(3)}} \times \frac{\partial y^{(3)}}{\partial a^{(3)}} \times \frac{\partial a^{(3)}}{\partial a^{(2)}} \times \frac{\partial a^{(2)}}{\partial a^{(1)}} \times \frac{\partial a^{(1)}}{\partial \omega_{x}}\right)_{\text{path-3}}$$

$$\boxed{\frac{\partial L^{(T)}}{\partial \omega_{h}} = \sum_{t=1}^{T} \frac{\partial L^{(t)}}{\partial y^{(t)}} \times \frac{\partial y^{(t)}}{\partial a^{(t)}} \times \frac{\partial a^{(3)}}{\partial \omega_{R}}}$$   BackPropagation Through Time   (BPTT)

# Dependency Graph  at time step 3



→ This tells us how  $a^{(3)} \to a^{(2)} \to a^{(1)}$
at time step = 3

Gradients :-

$$\frac{\partial L^{(3)}}{\partial \omega} = \left(\frac{\partial L^{(3)}}{\partial y^{(3)}} \times \frac{\partial y^{(3)}}{\partial a^{(3)}} \times \frac{\partial a^{(3)}}{\partial \omega_{h}}\right)_{\text{path-1}} + \left(\frac{\partial L^{(3)}}{\partial y^{(3)}} \times \frac{\partial y^{(3)}}{\partial a^{(3)}} \times \frac{\partial a^{(3)}}{\partial a^{(2)}} \times \frac{\partial a^{(2)}}{\partial \omega_{h}}\right)_{\text{path-3}} + \left(\frac{\partial L^{(3)}}{\partial y^{(3)}} \times \frac{\partial y^{(3)}}{\partial a^{(3)}} \times \frac{\partial a^{(3)}}{\partial a^{(2)}} \times \frac{\partial a^{(2)}}{\partial a^{(1)}} \times \frac{\partial a^{(1)}}{\partial \omega_{h}}\right)_{\text{path-2}}$$

$$\boxed{\frac{\partial L^{(T)}}{\partial \omega_{h}} = \sum_{t=1}^{T} \frac{\partial L^{(t)}}{\partial y^{(t)}} \times \frac{\partial y^{(t)}}{\partial a^{(t)}} \times \frac{\partial a^{(3)}}{\partial \omega_{R}}}$$   BackPropagation Through Time   (BPTT)

→ Again Same General formula.

# Generalized BPTT

$$\boxed{\frac{\partial L^{(T)}}{\partial \omega_{h}} = \sum_{t=1}^{T} \frac{\partial L^{(t)}}{\partial y^{(t)}} \times \frac{\partial y^{(t)}}{\partial a^{(t)}} \times \frac{\partial a^{(3)}}{\partial \omega_{R}}}$$