

# RNN Recurrent Neural Networks

- Used for 'sequential' data  
(where order/time matters)

Given a sequence,  
predict next item in sequence

- Like CNNs look at neighbors (Spatial models),  
RNNs look at past & future (temporal models)  
(sometimes)

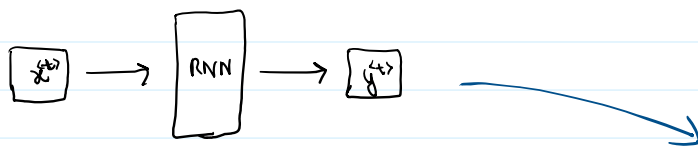
e.g  $x = [1, 2, 3, 4, 5, 6, 7, 8, 9, 10]$   
 $t_1 \quad t_2 \quad t_3 \quad \dots \quad t_{10}$

given  $x^{(t)}$ , predict  $x^{(t+1)}$

let  $x^{(t+1)}$  i.e. our target be  $y$  & since we're using  $x^t$  at time  $t$ , let's call it  $y^t$   
i.e.  $y^{(t)} = x^{(t+1)}$

$\Rightarrow$  given  $x^{(t)}$ , predict  $y^{(t)}$

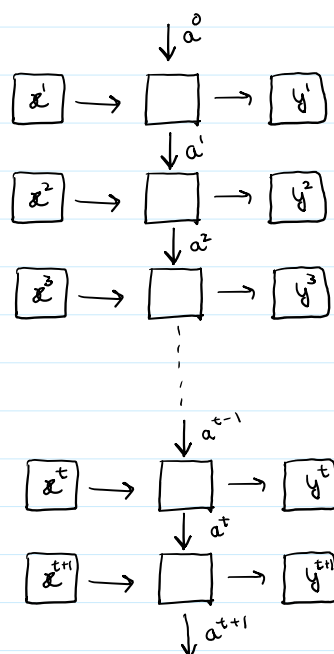
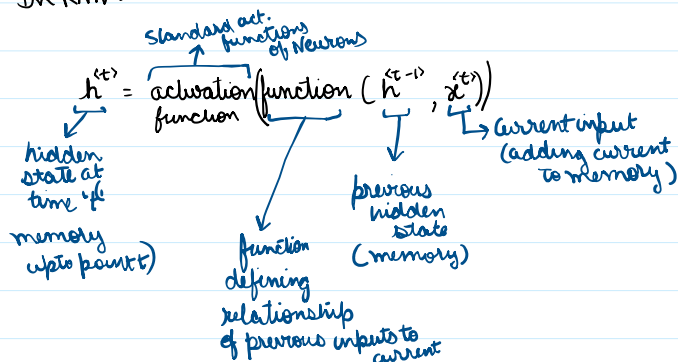
e.g at  $t=1$ ,  $x^{(1)} = 1$ ,  $y^{(1)} = x^{(2)} = 2$   
 $t=2$ ,  $x^{(2)} = 2$ ,  $y^{(2)} = x^{(3)} = 3$



We want our RNN to sort of keep a memory of every  $x^i$  that passes through  
& we do it through the hidden state that model learns  
this hidden state stores memory upto a point in sequence ( $t$ )

In a standard neural network a hidden layer is usually defined as:  
 $h = \text{activation}(Wx + b)$

In RNN:



Basically a FeedForward NN  
'unrolled' across time

Now for function describing  $(h_{t-1}, x_t)$ , we use linear function

$$f(h_{t-1}, x_t) = \underbrace{W_{hh}}_{\text{learned coeff. (weights)}} h^{(t-1)} + \underbrace{W_{hx}}_{\text{learned coeff. (weights)}} x^{(t)} + \underbrace{b_h}_{\text{bias}}$$

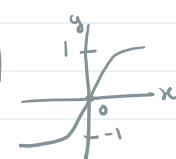
∴ the hidden layer becomes:

$$\hat{h}^{(t)} = \underset{\substack{\uparrow \\ \text{tanh}}}{a} (W_{hh} \hat{h}^{(t-1)} + W_{xh} \hat{x}^{(t)} + b_h)$$

& ∴ the output becomes:

$$\hat{y}^{(t)} = \underset{\substack{\uparrow \\ \text{None}}}{a} (W_{hy} \hat{h}^{(t)} + b_y)$$

★  $\tanh(x) \rightarrow [-1, 1]$   
 centered around 0  
 as  $x \rightarrow \infty$   
 $(\tanh(x))' \rightarrow 0$   
 i.e. gradient vanishes  
 Default for PyTorch's RNN



Let's walk through our example

with a simple RNN which we define to have an,

input size = 1

hidden size = 1

output size = 1

tanh activation

$$\hat{h}^{(t)} = \tanh(W_{hh} \hat{h}^{(t-1)} + W_{xh} \hat{x}^{(t)} + b_h)$$

$$\hat{y}^{(t)} = W_{hy} \hat{h}^{(t)} + b_y$$

Step 1

Initialize Weights

$$\begin{matrix} W_{hh} & , & W_{xh} & , & b_h & , & W_{hy} & , & b_y & , & h^0 \\ 0.5 & & 0.8 & & 0 & & 1.2 & & 0 & & 0 \end{matrix}$$

Step 2

Forward pass

Compute  $\hat{h}^{(t)}$  &  $\hat{y}^{(t)}$  for every  $t$  where  $t=1 \dots t-1$  i.e.  $(1, 10)$  in our case

$$\begin{aligned} \text{for } t=1, \quad \hat{h}^{(1)} &= \tanh(0.5(0) + 0.8(1) + 0) = 0.664 \\ \hat{y}^{(1)} &= 1.2(0.664 + 0) = 0.797 \end{aligned}$$

$$\begin{aligned} t=2, \quad \hat{h}^{(2)} &= \tanh(0.5(0.664) + 0.8(2) + 0) = 0.958 \\ \hat{y}^{(2)} &= 1.2(0.958) = 1.150 \\ &\dots \text{ and so on.} \end{aligned}$$

Step 3

Calculate Loss function

Forgot to mention before, but let's take MSE

At time  $t$

$$L^{(t)} = \frac{1}{2} (\underbrace{\hat{y}^{(t)}}_{\text{Predicted}} - \underbrace{y^{(t)}}_{\text{Original}})^2$$

For full sequence,

$$L_{\text{total}} = \sum_t L^{(t)} = \sum_t \frac{1}{2} (\hat{y}^{(t)} - y^{(t)})^2$$

objective function to minimize

For our example, this'll be,

$$L^{(1)} = \frac{1}{2} (0.797 - 2)^2 = 0.719$$

$$L^{(2)} = \frac{1}{2} (1.150 - 3)^2 = 1.711$$

.....  $L^{(9)}$

Step 4

Backpropagation Through Time

BPTT

compute gradients

$$\Rightarrow \frac{\partial L}{\partial w_{hh}}, \frac{\partial L}{\partial w_{hn}}, \frac{\partial L}{\partial w_{hy}}, \frac{\partial L}{\partial b_h}, \frac{\partial L}{\partial b_y}$$

this is why we call it 'unrolling'

$$\frac{\partial L^{(t)}}{\partial \hat{y}} = \hat{y}^{(t)} - y^{(t)}$$

where,

$$\frac{\partial L^{(t)}}{\partial w_{hh}} = \frac{\partial L^{(t)}}{\partial \hat{y}} \cdot \frac{\partial \hat{y}^{(t)}}{\partial h^{(t)}} \cdot \frac{\partial h^{(t)}}{\partial w_{hh}} = (\hat{y}^{(t)} - y^{(t)}) w_{hy} (1 - h^{(t)})^2 h^{(t)-1}$$

$$\frac{\partial L^{(t)}}{\partial w_{zh}} = \frac{\partial L^{(t)}}{\partial \hat{y}} \cdot \frac{\partial \hat{y}^{(t)}}{\partial h^{(t)}} \cdot \frac{\partial h^{(t)}}{\partial w_{zh}} = (\hat{y}^{(t)} - y^{(t)}) w_{hy} (1 - h^{(t)})^2 x^{(t)}$$

$$\frac{\partial L^{(t)}}{\partial w_{hy}} = \frac{\partial L^{(t)}}{\partial \hat{y}} \cdot \frac{\partial \hat{y}^{(t)}}{\partial w_{hy}} = (\hat{y}^{(t)} - y^{(t)}) (h^{(t)})$$

$$\frac{\partial L^{(t)}}{\partial b_h} = \frac{\partial L^{(t)}}{\partial \hat{y}} \cdot \frac{\partial \hat{y}^{(t)}}{\partial h^{(t)}} \cdot \frac{\partial h^{(t)}}{\partial b_h} = (\hat{y}^{(t)} - y^{(t)}) w_{hy} (1 - h^{(t)})^2$$

$$\frac{\partial L^{(t)}}{\partial b_y} = \frac{\partial L^{(t)}}{\partial \hat{y}} \cdot \frac{\partial \hat{y}}{\partial b_y} = (\hat{y}^{(t)} - y^{(t)}) (1)$$

Apply gradient descent

$$\Theta_{\text{new}} = \Theta_{\text{existing}} - \alpha \frac{\partial L}{\partial \Theta}$$

New param      Existing param      Learning Rate      Gradient computed w.r.t param above

$$\Theta = \{w_{hy}, w_{zh}, w_{hh}, b_h, b_y\}$$

