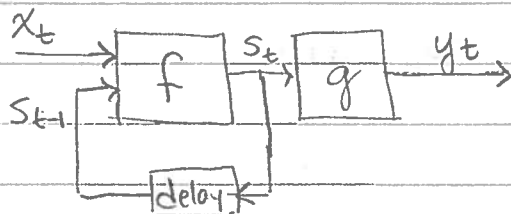


L 10-1

# Recurrent neural networks

Learning to be a state machine!

Recall:



Assume  $x_t$  is  $l \times 1$

$s_t$  is  $m \times 1$

$y_t$  is  $V \times 1$

$f_1$  and  $f_2$  are activation fns

Change in notation to match latex: "v"

Parameters  $W^{sx}$  is  $m \times l$

$W^{ss}$  is  $m \times m$

$W^o$  is  $V \times m$

$W_\phi$  is  $m \times 1$

$W_\phi^o$  is  $V \times 1$

always good to check dimensions!

$$s_t = f_1(W^{sx} x_t + W^{ss} s_{t-1} + W_\phi)$$

$$y_t = f_2(W^o s_t + W_\phi^o)$$

Training seq to seq (c.f. Dancing cheek to cheek)

Training set:  $[(x^{(1)}, y^{(1)}), \dots, (x^{(q)}, y^{(q)})]$

- $x^{(i)}$  and  $y^{(i)}$  are length  $n^{(i)}$  sequences
- the sequences in the same pair are the same length
- sequences in different pairs may have different lengths

Loss function: let  $\theta = (W^{sx}, W^{ss}, W^o, W_\phi, W_\phi^o)$

$$L(\theta) = \sum_{i=1}^q \text{Loss}_{\text{seq}}(y^{(i)}, \text{RNN}(x^{(i)}; \theta))$$

need loss function on sequences.

Usually:

$$\text{Loss}_{\text{seq}}(y^{(i)}, p^{(i)}) = \sum_{j=1}^{n^{(i)}} \text{Loss}(y_j^{(i)}, p_j^{(i)})$$

sorry :)

Change in notation  
to match later: "v"

$S_t$  is  $m \times 1$

$y_t$  is  $V \times 1$

activation fcn's

Parameters  $W^{sx}$  is  $m \times l$

$W^{ss}$  is  $m \times m$

$W^0$  is  $V \times m$

$W_\phi$  is  $m \times 1$

$W_\phi^0$  is  $V \times 1$

always  
good to  
check  
dimensions!

$$S_t = f_1(W^{sx} x_t + W^{ss} S_{t-1} + W_\phi)$$

$$y_t = f_2(W^0 S_t + W_\phi^0)$$

Training seq to seq (c.f. Dancing cheek to cheek)

Training set:  $[(x^{(1)}, y^{(1)}), \dots, (x^{(q)}, y^{(q)})]$

- $x^{(i)}$  and  $y^{(i)}$  are length  $n^{(q)}$  sequences
- the sequences in the same pair are the same length
- sequences in different pairs may have different lengths

Loss function: let  $\Theta = (W^{sx}, W^{ss}, W^0, W_\phi, W_\phi^0)$

$$L(\Theta) = \sum_{i=1}^q \text{Loss}_{\text{seq}}(y^{(i)}, \text{RNN}(x^{(i)}; \Theta))$$

result of running transduce on  $x$

need loss function on sequences.

Usually:

$$\text{Loss}_{\text{seq}}(y^{(i)}, p^{(i)}) = \sum_{j=1}^{n^{(q)}} \text{Loss}(y_j^{(i)}, p_j^{(i)})$$

sorry :)

sum of losses along sequence

note: this is consistent  
with latex'd notes but inconsistent  
with some of my prev notes which  
puts arguments in other order

L10-2

Our choice of per-element loss and output activation function  $f_2$  go together, and depend on the problem. Might be

$y_i$	$f_2$	Loss
real	linear	squared
one-hot	softmax	NLL
	etc.	

$f_1$  is typically tanh

How to find  $\theta$  to minimize loss on our training data?  
gradient descent!

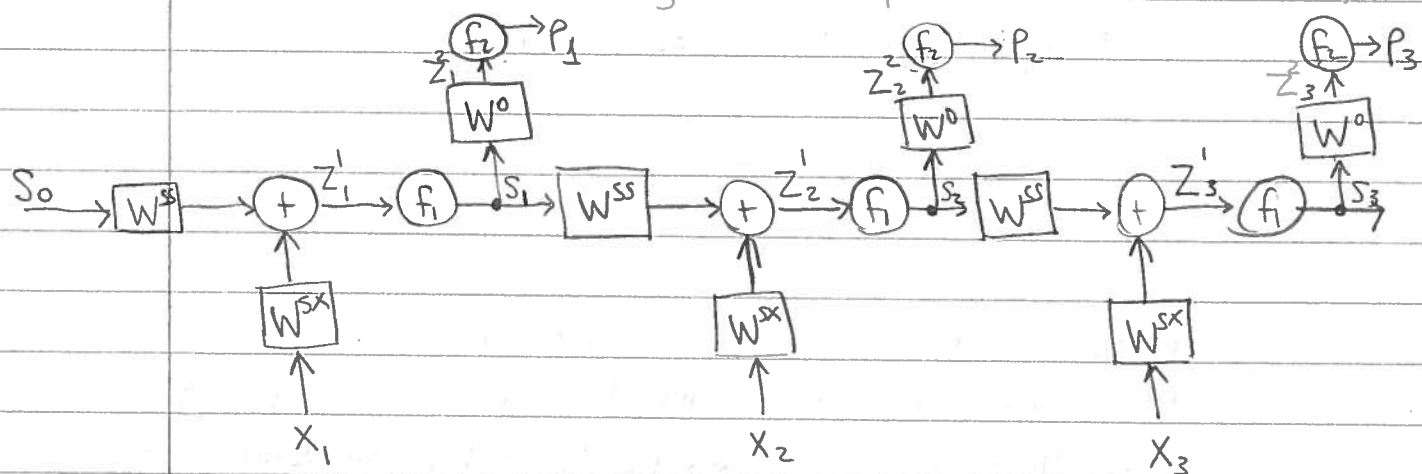
(BPTT)

Easy way to think about it: backpropagation through time

Do basic SGD:

① Sample training pair  $(x, y)$ ; let its length be  $n$

② "Unroll" the RNN to length  $T$  (picture for  $T=3$  below)



③ Do forward pass computing:

$$z'_t = W^{sx} x_t + W^{ss} s_{t-1} + W_\phi$$

$$s_t = f_1(z'_t)$$

$$z^2_t = W^0 s_t + W_\phi$$

$$P_t = f_2(z^2_t)$$

'predicted' output =  $P_1 \dots P_T$

L10-3

④ Do backward pass to compute the gradients.

We need

for all our  $W^{ss}$ ,  $W^{sx}$

$$\frac{\partial L}{\partial W} = \sum_{u=1}^n \frac{\partial L_u}{\partial W}$$

$$= \sum_{u=1}^n \sum_{t=1}^n \frac{\partial L_u}{\partial S_t} \cdot \frac{\partial S_t}{\partial W}$$

Total derivative  
Sum over all  
paths through which  
 $W$  affects  $L_t$

$$= \sum_{t=1}^n \frac{\partial S_t}{\partial W} \cdot \sum_{u=1}^n \frac{\partial L_u}{\partial S_t}$$

re-organizing

$$= \sum_{t=1}^n \frac{\partial S_t}{\partial W} \cdot \sum_{u=t}^n \frac{\partial L_u}{\partial S_t}$$

$S_t$  only affects  
 $L_u, L_{u+1}, \dots, L_n$

$$= \sum_{t=1}^n \frac{\partial S_t}{\partial W} \left[ \frac{\partial L_t}{\partial S_t} + \underbrace{\sum_{u=t+1}^n \frac{\partial L_u}{\partial S_t}}_{\text{define } \delta_t^{S_t}} \right]$$

$\delta_t^{S_t}$  is the dependence of the loss on steps  
after  $t$  on the state at time  $t$ .

We can compute this backwards, with  $t$  going from  
 $n$  down to 1.

L 10-4

$$F_t = \sum_{u=t+1}^n \text{Loss}(y_u, p_u)$$

future loss

Trickiest part is figuring out how early states contribute to late losses.

At the last stage,  $F_n = 0$  so  $\delta^{s_n} = 0$

Working backwards:

$$\delta^{s_{t-1}} = \frac{\partial}{\partial s_{t-1}} \sum_{j=t}^n \text{Loss}(y_j, p_j)$$

$$= \frac{\partial s_t}{\partial s_{t-1}} \cdot \frac{\partial}{\partial s_t} \sum_{j=t}^n \text{Loss}(y_j, p_j)$$

$$= \frac{\partial s_t}{\partial s_{t-1}} \cdot \frac{\partial}{\partial s_t} \left[ \text{Loss}(y_t, p_t) + \sum_{j=t+1}^n \text{Loss}(y_j, p_j) \right]$$

$$= \frac{\partial s_t}{\partial s_{t-1}} \cdot \left[ \frac{\partial \text{Loss}(y_t, p_t)}{\partial s_t} + \delta^{s_t} \right]$$

Now, we need

$$\frac{\partial \text{Loss}(y_t, p_t)}{\partial s_t} = \frac{\partial z_t^2}{\partial s_t} \cdot \frac{\partial \text{Loss}(y_t, p_t)}{\partial z_t^2}$$

(m x 1)

(m x v)

(v x 1)

$W_0^T$

$$\frac{\partial s_t}{\partial s_{t-1}} = \frac{\partial z_t'}{\partial s_{t-1}} \cdot \frac{\partial s_t}{\partial z_t'} = W_{ss}^T * f'(z_t')$$

(m x m)

(m x m)

(m x m)

↑ not dot!

$W_{ss}^T$

$\begin{bmatrix} f'(z_{t1}) \\ \vdots \\ 0 \end{bmatrix}$

$f'(z_{tm})$

$$\text{So: } \delta^{s_{t-1}} = W_{ss}^T * f'(z_t') \cdot \left( W_0^T \frac{\partial L_t}{\partial z_t^2} + \delta^{s_t} \right)$$

At the last stage,  $F_n = 0$  so  $\delta^{S_n} = 0$

Working backwards:

$$\delta^{S_{t-1}} = \frac{\partial S_{t-1}}{\partial S_t} \sum_{j=t}^n \text{Loss}(y_j, p_j)$$

$$= \frac{\partial S_t}{\partial S_{t-1}} \cdot \frac{\partial}{\partial S_t} \sum_{j=t}^n \text{Loss}(y_j, p_j)$$

$$= \frac{\partial S_t}{\partial S_{t-1}} \cdot \frac{\partial}{\partial S_t} \left[ \text{Loss}(y_t, p_t) + \sum_{j=t+1}^n \text{Loss}(y_j, p_j) \right]$$

$$= \frac{\partial S_t}{\partial S_{t-1}} \cdot \left[ \frac{\partial \text{Loss}(y_t, p_t)}{\partial S_t} + \delta^{S_t} \right]$$

Now, we need

$$\frac{\partial \text{Loss}(y_t, p_t)}{\partial S_t} = \frac{\partial Z_t^2}{\partial S_t} \cdot \frac{\partial \text{Loss}(y_t, p_t)}{\partial Z_t^2}$$

$(m \times 1) \quad (m \times v) \quad (v \times 1)$

$$W_0^T$$

$$\frac{\partial S_t}{\partial S_{t-1}} = \frac{\partial Z_t'}{\partial S_{t-1}} \cdot \frac{\partial S_t}{\partial Z_t'} = W_{ss}^T * f'(Z_t')$$

$(m \times m) \quad (m \times m) \quad (m \times m)$

$$W_{ss}^T \begin{bmatrix} f'(z_{t1}) \\ \vdots \\ 0 \end{bmatrix}$$

↑ not dot!

$$\text{So: } \delta^{S_{t-1}} = \underbrace{W_{ss}^T * f'(Z_t')}_{\frac{\partial S_t}{\partial S_{t-1}}} \cdot \underbrace{\left( W_0^T \frac{\partial L_t}{\partial Z_t^2} + \delta^{S_t} \right)}_{\frac{\partial F_{t-1}}{\partial S_t}}$$

$$\frac{\partial S_t}{\partial S_{t-1}}$$

$$\frac{\partial F_{t-1}}{\partial S_t}$$

(10-5)

Weight updates

As we iterate backwards, sum up

$$\frac{\partial L}{\partial W^{ss}} += \frac{\partial S_t}{\partial W^{ss}} \cdot \frac{\partial F_{t-1}}{\partial S_t} = \frac{\partial Z'_t}{\partial W^{ss}} \cdot \frac{\partial S_t}{\partial Z'_t} \cdot \frac{\partial F_{t-1}}{\partial S_t}$$

$S_{t-1}$       $f'(Z'_t)$      prev result

$$\frac{\partial L}{\partial W^{sx}} += \frac{\partial S_t}{\partial W^{sx}} \cdot \frac{\partial F_{t-1}}{\partial S_t} = \frac{\partial Z'_t}{\partial W^{sx}} \cdot \frac{\partial S_t}{\partial Z'_t} \cdot \frac{\partial F_{t-1}}{\partial S_t}$$

$X_t$       $f'(Z'_t)$      prev result

Handle  $W^0$  separately (it's easy):

$$\frac{\partial L}{\partial W^0} = \sum_{t=1}^n \frac{\partial L_t}{\partial W^0} = \sum_{t=1}^n \frac{\partial L_t}{\partial Z_t^2} \cdot \frac{\partial Z_t^2}{\partial W^0}$$

Assume we have  $\partial L_t / \partial Z_t^2 = (p_t - y_t)$   
 (true for quadratic, softmax-NLL, etc)

On each iteration

$$\frac{\partial L}{\partial W_0} += (p_t - y_t) \cdot S_t^T$$

$v \times m$       $(v \times 1)$       $(1 \times m)$

whew!

Can train a language model to predict next

token:  $c_t = \text{RNN}(c_1, \dots, c_{t-1})$

$c$  is character  
or word

by making training pairs

$X = (\langle \text{start} \rangle, c_1, c_2, c_3, \dots, c_T)$

$y = (c_1, c_2, c_3, \dots, c_T)$

$S_{t-1}$   $f'(z'_t)$  prev result

$$\frac{\partial L}{\partial W^{sx}} += \frac{\partial S_t}{\partial W^{sx}} \cdot \frac{\partial F_{t-1}}{\partial S_t} = \frac{\partial z'_t}{\partial W^{sx}} \cdot \frac{\partial S_t}{\partial z'_t} \cdot \frac{\partial F_{t-1}}{\partial S_t}$$

$x_t$   $f'(z'_t)$  prev result

Handle  $W^0$  separately (it's easy):

$$\frac{\partial L}{\partial W^0} = \sum_{t=1}^n \frac{\partial L_t}{\partial W^0} = \sum_{t=1}^n \frac{\partial L_t}{\partial z_t^2} \cdot \frac{\partial z_t^2}{\partial W^0}$$

Assume we have  $\partial L_t / \partial z_t^2 = (p_t - y_t)$   
(true for quadratic, softmax-NLL, etc)

On each iteration

$$\frac{\partial L}{\partial W_0} += (p_t - y_t) \cdot S_t^T$$

$v \times m$   $(v \times 1)$   $(1 \times m)$

whew!

Can train a language model to predict next

token:  $c_t = \text{RNN}(c_1, \dots, c_{t-1})$

$c$  is character  
or word

by making training pairs

$x = (\langle \text{start} \rangle, c_1, c_2, c_3, \dots, c_T)$

$y = (c_1, c_2, c_3, c_4, \dots, \langle \text{end} \rangle)$