ECE 1508: Applied Deep Learning Chapter 6: Recurrent NNs

Ali Bereyhi

ali.bereyhi@utoronto.ca

Department of Electrical and Computer Engineering University of Toronto

Winter 2025

We now want to train this basic RNN



Let's consider training for only one sample

We assume that we have a sequence of labels v[1],...,v[T]
We know some v[t] could be empty, e.g., in many-to-one scenario
So could be some of x[t]'s, e.g., in one-to-many scenario

We however have no problem with that!

We now want to train this basic RNN



Let's consider training for only one sample

The loss in general can be written as $\hat{R} = \mathcal{L} (\mathbf{y}[1:T], \mathbf{v}[1:T])$ where we use shorten notation $\mathbf{y}[1:T] = \mathbf{y}[1], \dots, \mathbf{y}[T]$

We can think of such a diagram



- + It seems to be hard to get back from \hat{R} to each $\mathbf{y}[t]$
- Well! That's true, but we have some remedy for it

We can think of such a diagram



For the moment, we assume that we can

compute the $\nabla_{\mathbf{y}[t]} \hat{R}$ for all $t \equiv$ move backward from \hat{R} to any $\mathbf{y}[t]$

Backward Pass Through Time

Starting from \hat{R} , say we want to find $abla_{\mathbf{W}_{\mathrm{m}}}\hat{R}$

• Since \hat{R} is a function of $\mathbf{y}[1:T]$, we should write a vectorized chain rule

$$\nabla_{\mathbf{W}_{\mathrm{m}}} \hat{R} = \nabla_{\mathbf{y}[1]} \hat{R} \circ \nabla_{\mathbf{W}_{\mathrm{m}}} \mathbf{y}[1] + \ldots + \nabla_{\mathbf{y}[T]} \hat{R} \circ \nabla_{\mathbf{W}_{\mathrm{m}}} \mathbf{y}[T]$$
$$= \sum_{t=1}^{T} \nabla_{\mathbf{y}[t]} \hat{R} \circ \nabla_{\mathbf{W}_{\mathrm{m}}} \mathbf{y}[t]$$

• Since we assumed that we have $\nabla_{\mathbf{y}[t]} \hat{R}$, our main task is to find

 $abla_{\mathbf{W}_{\mathrm{m}}}\mathbf{y}[t]$

for all t: we could hence compute it for a general t

 $\, \, \downarrow \, \,$ This is a tensor-like gradient, i.e., $[
abla_{\mathbf{W}_{\mathrm{m}}} y_1[t], \ldots,
abla_{\mathbf{W}_{\mathrm{m}}} y_M[t]]$

• Apparently, we should apply chain rule for several times!

,

Backward Pass Through Time

- + But, this is going to be exhausting?
- Well, we could again follow Albert Einstein advice!



"Everything should be made as simple as possible, but not simpler!"

Let's consider a dummy RNN with all variables being scalar

- **1** We have $y[t] = f(w_2h[t])$
- 2 We have $h[t] = f(w_1x[t] + w_mh[t-1])$
- **3** We start with hidden state h[0]

So the diagram gets simplified as below



Starting from \hat{R} , say we want to find $\nabla_{w_{\rm m}}\hat{R}$: our main task is to find

Let's start the computation: say we have passed forward through the RNN

 $\frac{\partial y[t]}{\partial w_m}$

• we now have y[t], h[t] and x[t] for all t

To go backward, we note that $y[t] = f(w_2h[t])$

• y[t] is a function of h[t]: so we write the chain rule as

$$rac{\partial y[t]}{\partial w_{\mathrm{m}}} = rac{\partial y[t]}{\partial h[t]} rac{\partial h[t]}{\partial w_{\mathrm{m}}}$$

- \downarrow we can compute the first term as $\partial y[t]/\partial h[t] = w_2 \dot{f}(w_2 h[t])$

$$\frac{\partial y[t]}{\partial w_{\rm m}} = w_2 \dot{f}(w_2 h[t]) \frac{\partial h[t]}{\partial w_{\rm m}}$$

We keep going backward, by noting that $h[t] = f(w_1x[t] + w_mh[t-1])$

- h[t] is a function of $w_{\rm m}$ and h[t-1]:¹
 - $\, { \, { \scriptstyle { \scriptstyle \leftarrow } } } \, \, \textit{let's write} \, \, h[t] = g \, (w_{\rm m}, h[t-1])$
- So we write the chain rule as

$$\frac{\partial h[t]}{\partial w_{\rm m}} = \frac{\partial g}{\partial w_{\rm m}} \underbrace{\frac{\partial w_{\rm m}}{\partial w_{\rm m}}}_{1} + \frac{\partial g}{\partial h[t-1]} \frac{\partial h[t-1]}{\partial w_{\rm m}}$$
$$= \frac{\partial g}{\partial w_{\rm m}} + \frac{\partial g}{\partial h[t-1]} \frac{\partial h[t-1]}{\partial w_{\rm m}}$$

¹We ignore w_1 and x[t], as they are obviously not functions of $w_{
m m}$

$$\frac{\partial \boldsymbol{y[t]}}{\partial \boldsymbol{w_{\mathrm{m}}}} = w_2 \dot{f}(w_2 h[t]) \frac{\partial h[t]}{\partial \boldsymbol{w_{\mathrm{m}}}}$$

We keep going backward, by noting that h[t] = f(z[t])

• Let's define $z[t] = w_1 x[t] + w_m h[t-1]$

•
$$h[t] = g(w_{\rm m}, h[t-1])$$

- \downarrow we can compute $\partial g / \partial w_{\rm m} = h[t-1]\dot{f}(z[t])$
- \downarrow we can compute $\partial g/\partial h[t-1] = w_{\rm m}\dot{f}(z[t])$
- So we can simplify the chain rule as

$$\frac{\partial h[t]}{\partial w_{\rm m}} = h[t-1]\dot{f}(\boldsymbol{z}[t]) + w_{\rm m}\dot{f}(\boldsymbol{z}[t])\frac{\partial h[t-1]}{\partial w_{\rm m}}$$
$$= \dot{f}(\boldsymbol{z}[t])\left(h[t-1] + w_{\rm m}\frac{\partial h[t-1]}{\partial w_{\rm m}}\right)$$

$$\frac{\partial \boldsymbol{y[t]}}{\partial w_{\mathrm{m}}} = w_2 \dot{f}(w_2 h[t]) \dot{f}(\boldsymbol{z[t]}) \left(h[t-1] + w_{\mathrm{m}} \frac{\partial h[t-1]}{\partial w_{\mathrm{m}}} \right)$$

We keep going backward: h[t-1] = f(z[t-1])

- $\qquad \qquad \textbf{ Recall that } \boldsymbol{z[t-1]} = \boldsymbol{w_1}\boldsymbol{x[t-1]} + \boldsymbol{w_m}\boldsymbol{h[t-2]}$
- We just need to replace t with t 1 in the last derivation

$$\frac{\partial h[t-1]}{\partial w_{\rm m}} = \dot{f} \left(z[t-1] \right) \left(h[t-2] + w_{\rm m} \frac{\partial h[t-2]}{\partial w_{\rm m}} \right)$$

Backpropagation at time t is hence described by a pair of recursive equations

• At time t, we compute

$$\frac{\partial y[t]}{\partial w_{\rm m}} = w_2 \dot{f}(w_2 h[t]) \frac{\partial h[t]}{\partial w_{\rm m}}$$

• For each $i = t, t - 1, \dots, 1$, we use

$$\frac{\partial h[i]}{\partial w_{\rm m}} = \dot{f}\left(z[i]\right) \left(h[i-1] + w_{\rm m} \frac{\partial h[i-1]}{\partial w_{\rm m}}\right)$$

• We stop at i = 1, where we get

$$\frac{\partial h[1]}{\partial w_{\mathrm{m}}} = \dot{f}\left(z[1]\right) \left(h[0] + w_{\mathrm{m}} \underbrace{\frac{\partial h[0]}{\partial w_{\mathrm{m}}}}_{0}\right) = \dot{f}\left(z[1]\right) h[0]$$



Key Point

Propagating back in time is described via a recursive equation

Recursive equation has an interesting feature

• Say we have backpropagated from time t

$$\frac{\partial \boldsymbol{y[t]}}{\partial w_{\mathrm{m}}} = w_2 \dot{f}(w_2 h[t]) \frac{\partial h[t]}{\partial w_{\mathrm{m}}}$$

 $\downarrow \text{ We hence have } \partial h[t]/\partial w_{\mathrm{m}}, \partial h[t-1]/\partial w_{\mathrm{m}}, \dots, \partial h[1]/\partial w_{\mathrm{m}}$

- Now, if we want to backpropagate from t-1
 - $\, \, \downarrow \,$ We already have $\partial h[t-1]/\partial w_{\mathrm{m}}, \partial h[t-2]/\partial w_{\mathrm{m}}, \dots, \partial h[1]/\partial w_{\mathrm{m}}$

Moral of Story

Pass forward till end of sequence and backpropagate to the beginning just once

Backpropagation Through Time (BPTT)

BPTT():

- **1** Start at t with $\nabla_{\mathbf{W}_{m}} \mathbf{y}[t] = \mathbf{W}_{2} \circ \dot{f}(\mathbf{W}_{2}\mathbf{h}[t]) \circ \nabla_{\mathbf{W}_{m}}\mathbf{h}[t]$
- **2** Go back in time as $\nabla_{\mathbf{W}_{m}} \mathbf{h}[i] = f(\mathbf{z}[i]) \circ (\mathbf{h}[i-1] + \mathbf{W}_{m} \circ \nabla_{\mathbf{W}_{m}} \mathbf{h}[i-1])$
- **3** Stop at i = 1 with $\nabla_{\mathbf{W}_{m}} \mathbf{h}[1] = \dot{f}(\mathbf{z}[1]) \circ \mathbf{h}[0]$
- + It looks very similar to backpropagation in deep NNs!
- Exactly! Even simple RNN is very deep through time
- + Don't we experience vanishing or exploding gradient then!
- Yes! Let's check it out

BPTT: Vanishing Gradient

Let's expand the gradient in our example

$$\begin{aligned} \frac{\partial y[t]}{\partial w_{\rm m}} &= w_2 \dot{f}(w_2 h[t]) \dot{f}(z[t]) h[t-1] \\ &+ w_2 \dot{f}(w_2 h[t]) w_{\rm m} \dot{f}(z[t]) \dot{f}(z[t-1]) h[t-2] \\ &+ w_2 \dot{f}(w_2 h[t]) w_{\rm m}^2 \dot{f}(z[t]) \dot{f}(z[t-1]) \dot{f}(z[t-2]) h[t-3] \\ &+ \cdots \\ &+ w_2 \dot{f}(w_2 h[t]) w_{\rm m}^{i-1} \left(\prod_{j=0}^{i-1} \dot{f}(z[t-j])\right) h[t-i] \\ &+ \cdots \\ &+ w_2 \dot{f}(w_2 h[t]) w_{\rm m}^{t-1} \left(\prod_{j=0}^{i-1} \dot{f}(z[t-j])\right) h[0] \end{aligned}$$

BPTT: Vanishing Gradient

This says that gradient of hidden state in *i* steps back in time is multiplied by

$$w_2 w_{\mathrm{m}}^{i-1} \dot{f}(w_2 h[t]) \left(\prod_{j=0}^{i-1} \dot{f}\left(z[t-j] \right) \right)$$

Recall deep NNs: we could see two cases

- If $\dot{f}(\cdot) > 1$ most of the time
 - → very old hidden states can explode the gradient
 - $\, \downarrow \,$ this usually does not happen, because most activations are not like that
- If $\dot{f}(\cdot) < 1$ most of the time
 - └→ very old hidden states have pretty much no impact on the gradient
 - → this means that the RNN can train only up to a finite memory

Handling Vanishing Gradient Through Time

- + Cant we use the same remedy as in deep NNs?
- Yes and No!

It is important to note that vanishing gradient through time is a bit different: we are still updating weights with large enough gradients, but these gradients have no information of long-term memory

Solutions to vanishing gradient through time are mainly

- Use $tanh(\cdot)$ activation
- Use truncated BPTT
 - → Repeat short-term BPTTs every couple of time intervals
- Invoke gating approach
 - L→ This is the most sophisticated approach
 - → It was known for a long time, but received attention much later!