

# Lecture 18: Long/Short-Term Memory

Mark Hasegawa-Johnson

University of Illinois

ECE 417: Multimedia Signal Processing



- 1 Review: Recurrent Neural Networks
- 2 Vanishing/Exploding Gradient
- 3 Running Example: a Pocket Calculator
- 4 Regular RNN
- 5 Forget Gate
- 6 Long Short-Term Memory (LSTM)
- 7 Backprop for an LSTM
- 8 Weight Gradient for an LSTM
- 9 Conclusion

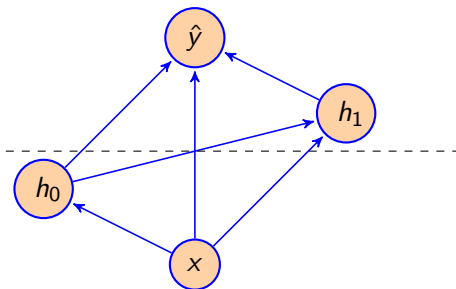
# Outline

- 1 Review: Recurrent Neural Networks
- 2 Vanishing/Exploding Gradient
- 3 Running Example: a Pocket Calculator
- 4 Regular RNN
- 5 Forget Gate
- 6 Long Short-Term Memory (LSTM)
- 7 Backprop for an LSTM
- 8 Weight Gradient for an LSTM
- 9 Conclusion

# Review: Partial and Total Derivatives

- The **total derivative** symbol,  $\frac{d\mathcal{L}}{dh_k}$ , **always means the same thing**: derivative including the contributions of all paths from  $h_k$  to  $\mathcal{L}$ .
- The **partial derivative** symbol,  $\frac{\partial\mathcal{L}}{\partial h_k}$ , can mean **different things in different equations** (because different equations might hold constant a different set of other variables).
- There is a notation we can use to specify **which** other variables are being held constant:  $\frac{\partial\mathcal{L}}{\partial h_k}(\hat{y}_1, \hat{y}_6, \hat{y}_{10}, h_1, \dots, h_N)$  means “hold  $\hat{y}_1, \hat{y}_6, \hat{y}_{10}$ , and  $h_1, \dots, h_{k-1}, h_{k+1}, \dots, h_N$  constant.”

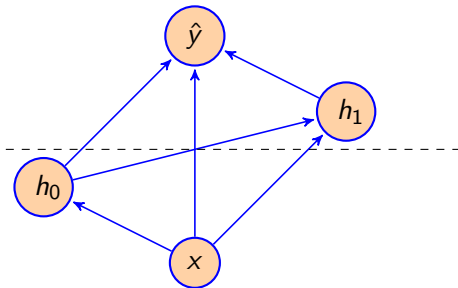
# Back-Propagation in terms of Flow Graphs



In order to find the derivative of an output w.r.t. any intermediate variables, one strategy that works is:

- 1 Draw a dashed line across the graph just downstream of the desired intermediate variables.
- 2 Apply the chain rule, with a summation across all edges that cross the dashed line.

# Back-Propagation in terms of Flow Graphs

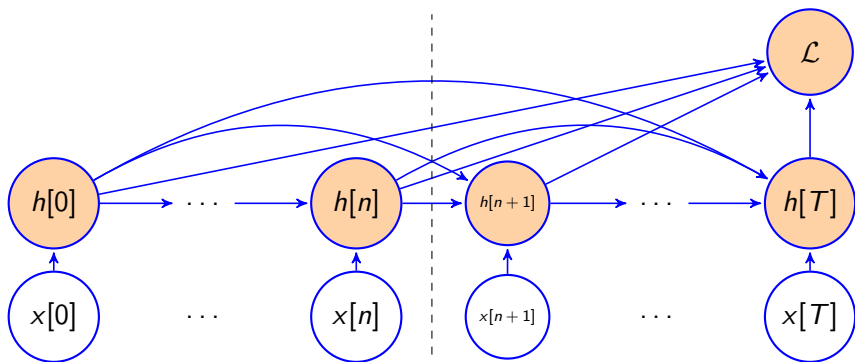


$$\frac{\partial \hat{y}}{\partial h_0}(x, h_0) = \frac{d\hat{y}}{dh_1} \frac{\partial h_1}{\partial h_0}(x, h_0, h_1) + \frac{d\hat{y}}{d\hat{y}} \frac{\partial \hat{y}}{\partial h_0}(x, h_0, h_1)$$

$$\frac{\partial \hat{y}}{\partial x}(x, h_0) = \frac{d\hat{y}}{dh_1} \frac{\partial h_1}{\partial x}(x, h_0, h_1) + \frac{d\hat{y}}{d\hat{y}} \frac{\partial \hat{y}}{\partial x}(x, h_0, h_1)$$

Notice:  $\frac{\partial \hat{y}}{\partial x}(x, h_0)$  does **not** include  $\frac{d\hat{y}}{dh_0} \frac{\partial h_0}{\partial x}$ .

# Back-Propagation Through Time



$$\frac{d\mathcal{L}}{dh[n]} = \frac{\partial\mathcal{L}}{\partial h[n]} + \sum_{m=1}^{T-n} \frac{d\mathcal{L}}{dh[n+m]} \frac{\partial h[n+m]}{\partial h[n]}$$

# Outline

- 1 Review: Recurrent Neural Networks
- 2 Vanishing/Exploding Gradient**
- 3 Running Example: a Pocket Calculator
- 4 Regular RNN
- 5 Forget Gate
- 6 Long Short-Term Memory (LSTM)
- 7 Backprop for an LSTM
- 8 Weight Gradient for an LSTM
- 9 Conclusion



# Vanishing/Exploding Gradient

- The “vanishing gradient” problem refers to the tendency of  $\frac{dh[n+m]}{dh[n]}$  to disappear, exponentially, when  $m$  is large.
- The “exploding gradient” problem refers to the tendency of  $\frac{dh[n+m]}{dh[n]}$  to explode toward infinity, exponentially, when  $m$  is large.
- If the largest feedback coefficient is  $|w[m]| > 1$ , then you get exploding gradient. If  $|w[m]| < 1$ , you get vanishing gradient.

## Example: A Memorizer Network

Suppose that we have a very simple RNN:

$$h[n] = wx[n] + uh[n - 1]$$

Suppose that  $x[n]$  is only nonzero at time 0:

$$x[n] = \begin{cases} x_0 & n = 0 \\ 0 & n \neq 0 \end{cases}$$

Suppose that, instead of measuring  $x[0]$  directly, we are only allowed to measure the output of the RNN  $m$  time-steps later. Our goal is to learn  $w$  and  $u$  so that  $h[m]$  remembers  $x_0$ , thus:

$$\mathcal{L} = \frac{1}{2} (h[m] - x_0)^2$$

# Example: A Memorizer Network

Now, how do we perform gradient update of the weights? If

$$h[n] = wx[n] + uh[n - 1]$$

then

$$\begin{aligned} \frac{d\mathcal{L}}{dw} &= \sum_n \left( \frac{d\mathcal{L}}{dh[n]} \right) \frac{\partial h[n]}{\partial w} \\ &= \sum_n \left( \frac{d\mathcal{L}}{dh[n]} \right) x[n] = \left( \frac{d\mathcal{L}}{dh[0]} \right) x_0 \end{aligned}$$

But the error is defined as

$$\mathcal{L} = \frac{1}{2} (h[m] - x_0)^2$$

so

$$\frac{d\mathcal{L}}{dh[0]} = u \frac{d\mathcal{L}}{dh[1]} = u^2 \frac{d\mathcal{L}}{dh[2]} = \dots = u^m (h[m] - x_0)$$

## Example: Vanishing Gradient

So we find out that the gradient, w.r.t. the coefficient  $w$ , is either exponentially small, or exponentially large, depending on whether  $|u| < 1$  or  $|u| > 1$ :

$$\frac{d\mathcal{L}}{dw} = x_0 (h[m] - x_0) u^m$$

In other words, if our application requires the neural net to wait  $m$  time steps before generating its output, then the gradient is exponentially smaller, and therefore training the neural net is exponentially harder.

## Exponential Decay

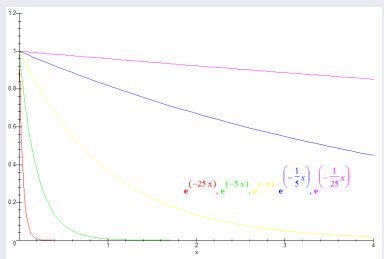


Image CC-SA-3.0, PeterQ, Wikipedia

# Outline

- 1 Review: Recurrent Neural Networks
- 2 Vanishing/Exploding Gradient
- 3 Running Example: a Pocket Calculator**
- 4 Regular RNN
- 5 Forget Gate
- 6 Long Short-Term Memory (LSTM)
- 7 Backprop for an LSTM
- 8 Weight Gradient for an LSTM
- 9 Conclusion

# Notation

Today's lecture will try to use notation similar to the Wikipedia page for LSTM.

- $x[t]$  = input at time  $t$
- $y[t]$  = target/desired output
- $c[t]$  = LSTM memory cell
- $h[t]$  = LSTM output
- $u$  = feedback coefficient
- $w$  = feedforward coefficient
- $b$  = bias

# Running Example: a Pocket Calculator

The rest of this lecture will refer to a toy application called “pocket calculator.”

## Pocket Calculator

- When  $x[t] > 0$ , add it to the current tally:  
 $c[t] = c[t - 1] + x[t]$ .
- When  $x[t] = 0$ ,
  - 1 Print out the current tally,  $h[t] = c[t - 1]$ , and then
  - 2 Reset the tally to zero,  $c[t] = 0$ .

## Example Signals

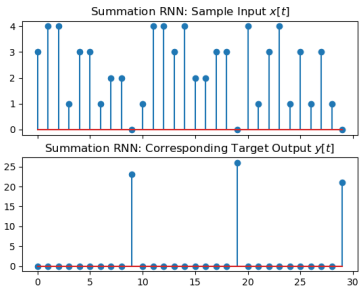
**Input:**  $x[t] = 1, 2, 1, 0, 1, 1, 1, 0$

**Target Output:**  $y[t] = 0, 0, 0, 4, 0, 0, 0, 3$

## Pocket Calculator

- When  $x[t] > 0$ , add it to the current tally:  
 $c[t] = c[t - 1] + x[t]$ .
- When  $x[t] = 0$ ,
  - 1 Print out the current tally,  $h[t] = c[t - 1]$ , and then
  - 2 Reset the tally to zero,  $c[t] = 0$ .

## Pocket Calculator





# Outline

- 1 Review: Recurrent Neural Networks
- 2 Vanishing/Exploding Gradient
- 3 Running Example: a Pocket Calculator
- 4 Regular RNN**
- 5 Forget Gate
- 6 Long Short-Term Memory (LSTM)
- 7 Backprop for an LSTM
- 8 Weight Gradient for an LSTM
- 9 Conclusion

# One-Node One-Tap Linear RNN

Suppose that we have a very simple RNN:

$$\text{Excitation: } c[t] = x[t] + uh[t - 1]$$

$$\text{Activation: } h[t] = \sigma_h(c[t])$$

where  $\sigma_h()$  is some feedback nonlinearity. In this simple example, let's just use  $\sigma_h(c[t]) = c[t]$ , i.e., no nonlinearity.

**GOAL:** Find  $u$  so that  $h[t] \approx y[t]$ . In order to make the problem easier, we will only score an "error" when  $y[t] \neq 0$ :

$$\mathcal{L} = \frac{1}{2} \sum_{t:y[t]>0} (h[t] - y[t])^2$$

## RNN: $u = 1$ ?

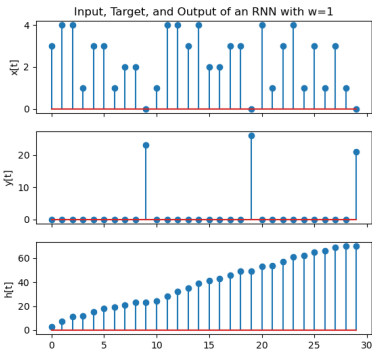
Obviously, if we want to just add numbers, we should just set  $u = 1$ . Then the RNN is computing

$$\text{Excitation: } c[t] = x[t] + h[t - 1]$$

$$\text{Activation: } h[t] = \sigma_h(c[t])$$

That works until the first zero-valued input. But then it just keeps on adding.

## RNN with $u = 1$

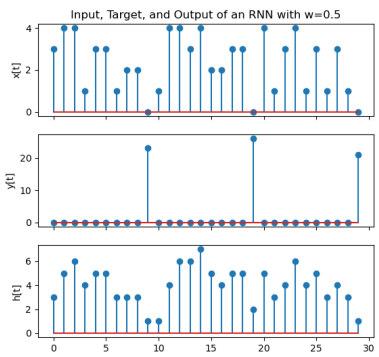


## RNN: $u = 0.5$

Can we get decent results using  $u = 0.5$ ?

- Advantage: by the time we reach  $x[t] = 0$ , the sum has kind of leaked away from us ( $c[t] \approx 0$ ), so a hard-reset is not necessary.
- Disadvantage: by the time we reach  $x[t] = 0$ , the sum has kind of leaked away from us ( $h[t] \approx 0$ ).

## RNN with $u = 0.5$



# Gradient Descent

$$c[t] = x[t] + uh[t - 1]$$
$$h[t] = \sigma_h(c[t])$$

Let's try initializing  $u = 0.5$ , and then performing gradient descent to improve it. Gradient descent has five steps:

- 1 **Forward Propagation:**  $c[t] = x[t] + uh[t - 1]$ ,  $h[t] = \sigma_h(c[t])$ .
- 2 **Synchronous Backprop:**  $\epsilon[t] = \partial\mathcal{L}/\partial c[t]$ .
- 3 **Back-Prop Through Time:**  $\delta[t] = d\mathcal{L}/dc[t]$ .
- 4 **Weight Gradient:**  $d\mathcal{L}/du = \sum_t \delta[t]h[t - 1]$
- 5 **Gradient Descent:**  $u \leftarrow u - \eta d\mathcal{L}/du$

# Gradient Descent

Excitation:  $c[t] = x[t] + uh[t - 1]$

Activation:  $h[t] = \sigma_h(c[t])$

$$\text{Error: } \mathcal{L} = \frac{1}{2} \sum_{t:y[t]>0} (h[t] - y[t])^2$$

So the back-prop stages are:

**Synchronous Backprop:**  $\epsilon[t] = \frac{\partial \mathcal{L}}{\partial c[t]} = \begin{cases} (h[t] - y[t]) & y[t] > 0 \\ 0 & \text{otherwise} \end{cases}$

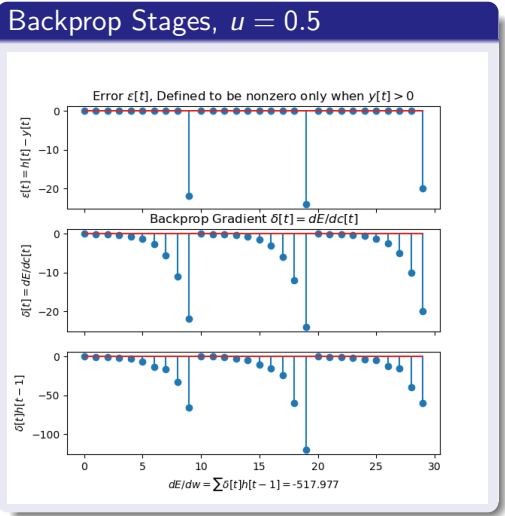
**BPTT:**  $\delta[t] = \frac{d\mathcal{L}}{dc[t]} = \epsilon[t] + u\delta[t + 1]$

**Weight Gradient:**  $\frac{d\mathcal{L}}{du} = \sum_t \delta[t]h[t - 1]$

### Backprop Stages

$$\epsilon[t] = \begin{cases} (h[t] - y[t]) & y[t] > 0 \\ 0 & \text{otherwise} \end{cases}$$

$$\delta[t] = \epsilon[t] + u\delta[t + 1]$$

$$\frac{d\mathcal{L}}{du} = \sum_t \delta[t]h[t - 1]$$


# Vanishing Gradient and Exploding Gradient

- Notice that, with  $|u| < 1$ ,  $\delta[t]$  tends to vanish exponentially fast as we go backward in time. This is called the **vanishing gradient** problem. It is a big problem for RNNs with long time-dependency, and for deep neural nets with many layers.
- If we set  $|u| > 1$ , we get an even worse problem, sometimes called the **exploding gradient** problem.

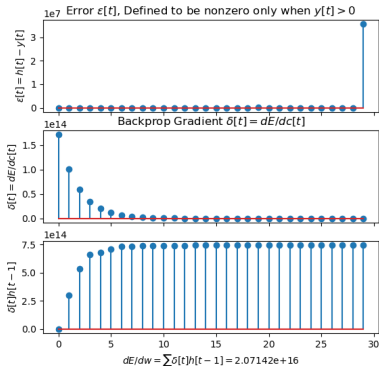
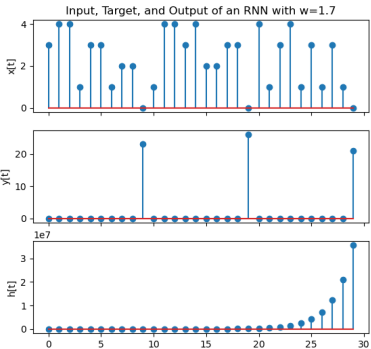


## RNN, $u = 1.7$

$$c[t] = x[t] + uh[t - 1]$$

## RNN, $u = 1.7$

$$\delta[t] = \epsilon[t] + u\delta[t + 1]$$



# Outline

- 1 Review: Recurrent Neural Networks
- 2 Vanishing/Exploding Gradient
- 3 Running Example: a Pocket Calculator
- 4 Regular RNN
- 5 Forget Gate**
- 6 Long Short-Term Memory (LSTM)
- 7 Backprop for an LSTM
- 8 Weight Gradient for an LSTM
- 9 Conclusion

# Hochreiter and Schmidhuber's Solution: The Forget Gate

Instead of multiplying by the same weight,  $u$ , at each time step, Hochreiter and Schmidhuber proposed: let's make the feedback coefficient a function of the input!

$$\text{Excitation: } c[t] = x[t] + f[t]h[t - 1]$$

$$\text{Activation: } h[t] = \sigma_h(c[t])$$

$$\text{Forget Gate: } f[t] = \sigma_g(w_f x[t] + u_f h[t - 1] + b_f),$$

where  $\sigma_h()$  and  $\sigma_g()$  might be different nonlinearities. In particular, it's OK for  $\sigma_h()$  to be linear ( $\sigma_h(c) = c$ ), but  $\sigma_g()$  should be clipped so that  $0 \leq f[t] \leq 1$ , in order to avoid gradient explosion.

# The Forget-Gate Nonlinearity

The forget gate is

$$f[t] = \sigma_g (w_f x[t] + u_f h[t - 1] + b_f)$$

where  $\sigma_g()$  is some nonlinearity such that  $0 \leq \sigma_g() \leq 1$ . Two such nonlinearities are worth knowing about.

# Forget-Gate Nonlinearity #1: CReLU

The first useful nonlinearity is the CReLU (clipped rectified linear unit), defined as

$$\sigma_g(w_f x + u_f h + b_f) = \min(1, \max(0, w_f x + u_f h + b_f))$$

- The CReLU is particularly useful for **knowledge-based design**. That's because  $\sigma(1) = 1$  and  $\sigma(0) = 0$ , so it is relatively easy to design the weights  $w_f$ ,  $u_f$ , and  $b_f$  to get the results you want.
- The CReLU is not very useful, though, if you want to choose your weights using **gradient descent**. What usually happens is that  $w_f$  grows larger and larger for the first 2-3 epochs of training, and then suddenly  $w_f$  is so large that  $\dot{\sigma}(w_f x + u_f h + b_f) = 0$  for all training tokens. At that point, the gradient is  $d\mathcal{L}/dw = 0$ , so further gradient-descent training is useless.

# Forget-Gate Nonlinearity #1: Logistic Sigmoid

The second useful nonlinearity is the logistic sigmoid, defined as:

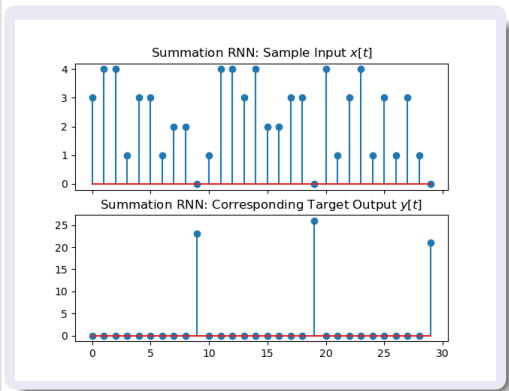
$$\sigma_g(w_f x + u_f h + b_f) = \frac{1}{1 + e^{-(w_f x + u_f h + b_f)}}$$

- The logistic sigmoid is particularly useful for **gradient descent**. That's because its gradient is defined for all values of  $w_f$ . In fact, it has a really simple form, that can be written in terms of the output:  $\dot{\sigma} = \sigma(1 - \sigma)$ .
- The logistic sigmoid is not as useful for **knowledge-based design**. That's because  $0 < \sigma < 1$ : as  $x \rightarrow -\infty$ ,  $\sigma(x) \rightarrow 0$ , but it never quite reaches it. Likewise as  $x \rightarrow \infty$ ,  $\sigma(x) \rightarrow 1$ , but it never quite reaches it.

## Pocket Calculator

- When  $x[t] > 0$ , accumulate the input, and print out nothing.
- When  $x[t] = 0$ , print out the accumulator, then reset.

...but the “print out nothing” part is not scored, only the accumulation. Furthermore, nonzero input is always  $x[t] \geq 1$ .

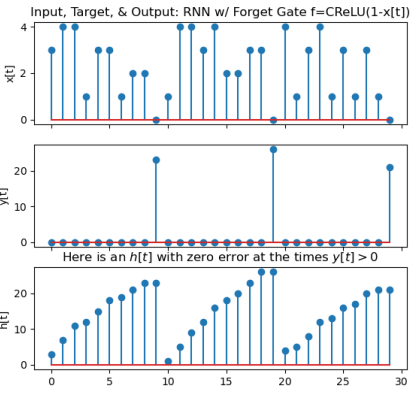


## Pocket Calculator

With zero error, we can approximate the pocket calculator as

- When  $x[t] \geq 1$ , accumulate the input.
- When  $x[t] = 0$ , print out the accumulator, then reset.

$$\mathcal{L} = \frac{1}{2} \sum_{t: y[t] > 0} (h[t] - y[t])^2 = 0$$





# Forget-Gate Implementation of the Pocket Calculator

It seems like we can approximate the pocket calculator as:

- When  $x[t] \geq 1$ , accumulate the input:  $c[t] = x[t] + h[t - 1]$ .
- When  $x[t] = 0$ , print out the accumulator, then reset:  $c[t] = x[t]$ .

So it seems that we just want the forget gate set to

$$f[t] = \begin{cases} 1 & x[t] \geq 1 \\ 0 & x[t] = 0 \end{cases}$$

This can be accomplished as

$$f[t] = \text{CReLU}(x[t]) = \max(0, \min(1, x[t]))$$

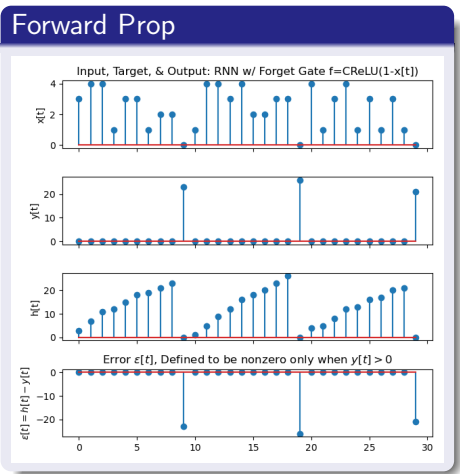
## Forget Gate Implementation of the Pocket Calculator

$$c[t] = x[t] + f[t]h[t - 1]$$

$$h[t] = c[t]$$

$$f[t] = \text{CReLU}(x[t])$$

... but the error is large!



# What Went Wrong?

- The forget gate correctly turned itself on (remember the past) when  $x[t] > 0$ , and turned itself off (forget the past) when  $x[t] = 0$ .
- Unfortunately, we don't want to forget the past when  $x[t] = 0$ . We want to forget the past on the **next time step after**  $x[t] = 0$ .
- Coincidentally, we also don't want any output when  $x[t] > 0$ . The error criterion doesn't score those samples, but maybe it should.

# Outline

- 1 Review: Recurrent Neural Networks
- 2 Vanishing/Exploding Gradient
- 3 Running Example: a Pocket Calculator
- 4 Regular RNN
- 5 Forget Gate
- 6 Long Short-Term Memory (LSTM)**
- 7 Backprop for an LSTM
- 8 Weight Gradient for an LSTM
- 9 Conclusion

## Long short-term memory

[S Hochreiter, J Schmidhuber](#) - Neural computation, 1997 - [ieeexplore.ieee.org](http://ieeexplore.ieee.org)

Learning to store information over extended time intervals by recurrent backpropagation takes a very long time, mostly because of insufficient, decaying error backflow. We briefly review Hochreiter's (1991) analysis of this problem, then address it by introducing a novel ...

☆ [Cite](#) [Cited by 56649](#) [Related articles](#) [All 50 versions](#) [↔](#)

# Long Short-Term Memory (LSTM)

The LSTM solves those problems by defining two types of memory, and three types of gates. The two types of memory are

- 1 The “cell,”  $c[t]$ , corresponds to the excitation in an RNN.
- 2 The “output” or “prediction,”  $h[t]$ , corresponds to the activation in an RNN.

The three gates are:

- 1 The cell remembers the past only when the forget gate is on,  $f[t] = 1$ .
- 2 The cell accepts input only when the input gate is on,  $i[t] = 1$ .
- 3 The cell is output only when the output gate is on,  $o[t] = 1$ .

# Long Short-Term Memory (LSTM)

The three gates are:

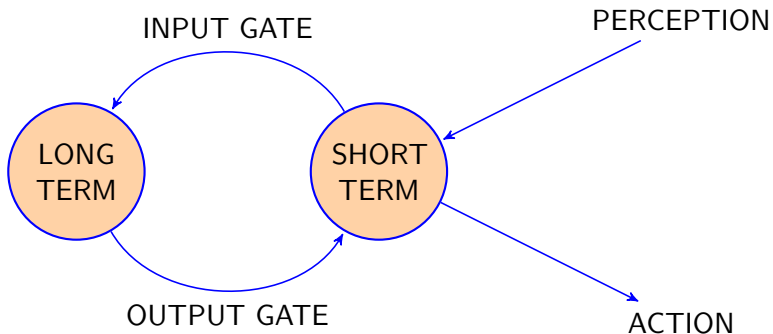
- 1 The cell remembers the past only when the forget gate is on,  $f[t] = 1$ .
- 2 The cell accepts input only when the input gate is on,  $i[t] = 1$ .

$$c[t] = f[t]c[t - 1] + i[t]\sigma_h(w_c x[t] + u_c h[t - 1] + b_c)$$

- 3 The cell is output only when the output gate is on,  $o[t] = 1$ .

$$h[t] = o[t]\sigma_h(c[t])$$

# Characterizing Human Memory



$$\Pr \{ \text{remember} \} = p_{LTM} e^{-t/T_{LTM}} + (1 - p_{LTM}) e^{-t/T_{STM}}$$



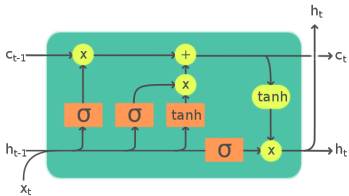
# When Should You Remember?

$$c[t] = f[t]c[t - 1] + i[t]\sigma_h(w_c x[t] + u_c h[t - 1] + b_c)$$

$$h[t] = o[t]\sigma_h(c[t])$$

- ① The forget gate is a function of current input and past output,  $f[t] = \sigma_g(w_f x[t] + u_f h[t - 1] + b_f)$
- ② The input gate is a function of current input and past output,  $i[t] = \sigma_g(w_i x[t] + u_i h[t - 1] + b_i)$
- ③ The output gate is a function of current input and past output,  $o[t] = \sigma_g(w_o x[t] + u_o h[t - 1] + b_o)$

# Neural Network Model: LSTM



**Legend:**

Layer	ComponentwiseCopy	Concatenate

$$i[t] = \text{input gate} = \sigma(w_i x[t] + u_i h[t - 1] + b_i)$$

$$o[t] = \text{output gate} = \sigma(w_o x[t] + u_o h[t - 1] + b_o)$$

$$f[t] = \text{forget gate} = \sigma(w_f x[t] + u_f h[t - 1] + b_f)$$

$$c[t] = \text{memory cell} = f[t]c[t - 1] + i[t]\tanh(w_c x[t] + u_c h[t - 1] + b_c)$$

$$h[t] = \text{output} = o[t]\tanh(c[t])$$

### Example: Pocket Calculator

$$c[t] = f[t]c[t - 1] + i[t]x[t]$$

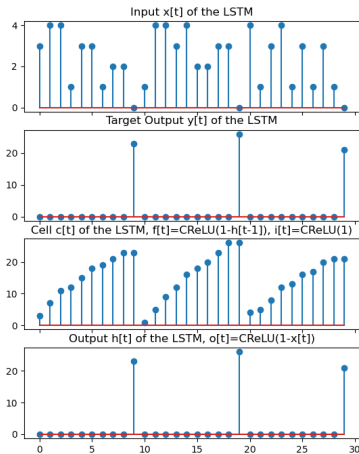
$$i[t] = 1 \text{ always}$$

$$h[t] = o[t]c[t]$$

$$o[t] = \begin{cases} 1 & x[t] = 0 \\ 0 & x[t] \geq 1 \end{cases}$$

$$f[t] = \begin{cases} 1 & h[t - 1] = 0 \\ 0 & h[t - 1] \geq 0 \end{cases}$$

### Forward Prop



### Example: Pocket Calculator

$$c[t] = f[t]c[t - 1] + i[t]x[t]$$

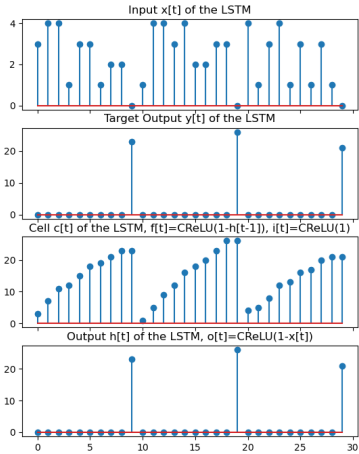
$$i[t] = \text{CReLU}(1)$$

$$h[t] = o[t]c[t]$$

$$o[t] = \text{CReLU}(1 - x[t])$$

$$f[t] = \text{CReLU}(1 - h[t - 1])$$

### Forward Prop



# Outline

- 1 Review: Recurrent Neural Networks
- 2 Vanishing/Exploding Gradient
- 3 Running Example: a Pocket Calculator
- 4 Regular RNN
- 5 Forget Gate
- 6 Long Short-Term Memory (LSTM)
- 7 Backprop for an LSTM**
- 8 Weight Gradient for an LSTM
- 9 Conclusion

# Backprop for a normal RNN

In a normal RNN, each epoch of gradient descent has five steps:

- ① **Forward-prop:** find the node excitation and activation, moving forward through time.
- ② **Synchronous backprop:** find the partial derivative of error w.r.t. node excitation at each time, assuming all other time steps are constant.
- ③ **Back-prop through time:** find the total derivative of error w.r.t. node excitation at each time.
- ④ **Weight gradient:** find the total derivative of error w.r.t. each weight and each bias.
- ⑤ **Gradient descent:** adjust each weight and bias in the direction of the negative gradient



# Backprop for an LSTM

In an LSTM, we'll implement each epoch of gradient descent with five steps:

- 1 **Forward-prop:** find all five of the variables at each time step, moving forward through time.
- 2 **Synchronous backprop:** find the partial derivative of error w.r.t.  $h[t]$ .
- 3 **Back-prop through time:** find the total derivative of error w.r.t. each of the five variables at each time, starting with  $h[t]$ .
- 4 **Weight gradient:** find the total derivative of error w.r.t. each weight and each bias.
- 5 **Gradient descent:** adjust each weight and bias in the direction of the negative gradient



# Synchronous Back-Prop: the Output

Suppose the error term is

$$\mathcal{L} = \frac{1}{2} \sum_{t=-\infty}^{\infty} (h[t] - y[t])^2$$

Then the first step, in back-propagation, is to calculate the partial derivative w.r.t. the prediction term  $h[t]$ :

$$\epsilon_h[t] = \frac{\partial \mathcal{L}}{\partial h[t]} = h[t] - y[t]$$

# Synchronous Back-Prop: the other variables

Remember that the error is defined only in terms of the output,  $h[t]$ . So, actually, partial derivatives with respect to the other variables are all zero!

$$\epsilon_i[t] = \frac{\partial \mathcal{L}}{\partial i[t]} = 0$$

$$\epsilon_o[t] = \frac{\partial \mathcal{L}}{\partial o[t]} = 0$$

$$\epsilon_f[t] = \frac{\partial \mathcal{L}}{\partial f[t]} = 0$$

$$\epsilon_c[t] = \frac{\partial \mathcal{L}}{\partial c[t]} = 0$$

# Back-Prop Through Time

Back-prop through time is really tricky in an LSTM, because four of the five variables depend on the previous time step, either on  $h[t - 1]$  and/or  $c[t - 1]$ :

$$i[t] = \sigma_g(w_i x[t] + u_i h[t - 1] + b_i)$$

$$o[t] = \sigma_g(w_o x[t] + u_o h[t - 1] + b_o)$$

$$f[t] = \sigma_g(w_f x[t] + u_f h[t - 1] + b_f)$$

$$c[t] = f[t]c[t - 1] + i[t]\sigma_h(w_c x[t] + u_c h[t - 1] + b_c)$$

$$h[t] = o[t]\sigma_h(c[t])$$

# Back-Prop Through Time

Taking the partial derivative of each variable at time  $t$  w.r.t. the variables at time  $t - 1$ , we get

$$\frac{\partial i[t]}{\partial h[t-1]} = \dot{\sigma}_g(w_i x[t] + u_i h[t-1] + b_i) u_i$$

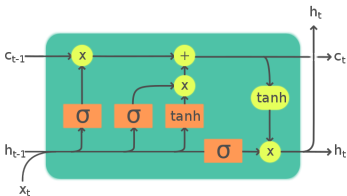
$$\frac{\partial o[t]}{\partial h[t-1]} = \dot{\sigma}_g(w_o x[t] + u_o h[t-1] + b_o) u_o$$

$$\frac{\partial f[t]}{\partial h[t-1]} = \dot{\sigma}_g(w_f x[t] + u_f h[t-1] + b_f) u_f$$

$$\frac{\partial c[t]}{\partial h[t-1]} = i[t] \dot{\sigma}_h(w_c x[t] + u_c h[t-1] + b_c) u_c$$

$$\frac{\partial c[t]}{\partial c[t-1]} = f[t]$$

# Back-Prop Through Time



Legend:

	Layer		ComponentwiseCopy	↗	↘
--	-------	--	-------------------	---	---

We can then combine all of those together to get:

$$\frac{d\mathcal{L}}{dh[t]} = \frac{\partial \mathcal{L}}{\partial h[t]} + \sum_{\xi \in \{i, o, f, c\}} \frac{d\mathcal{L}}{d\xi[t+1]} \frac{\partial \xi[t+1]}{\partial h[t]}$$

# Back-Prop Through Time

Back-propagation for all of the other variables is easier, since only  $c[t]$  has any direct connection from the current time to the next time:

$$\frac{d\mathcal{L}}{dc[t]} = \frac{d\mathcal{L}}{dh[t]} \frac{\partial h[t]}{\partial c[t]} + \frac{d\mathcal{L}}{dc[t+1]} \frac{\partial c[t+1]}{\partial c[t]}$$

$$\frac{d\mathcal{L}}{do[t]} = \frac{d\mathcal{L}}{dh[t]} \frac{\partial h[t]}{\partial o[t]}$$

$$\frac{d\mathcal{L}}{di[t]} = \frac{d\mathcal{L}}{dc[t]} \frac{\partial c[t]}{\partial i[t]}$$

$$\frac{d\mathcal{L}}{df[t]} = \frac{d\mathcal{L}}{dc[t]} \frac{\partial c[t]}{\partial f[t]}$$

# Outline

- 1 Review: Recurrent Neural Networks
- 2 Vanishing/Exploding Gradient
- 3 Running Example: a Pocket Calculator
- 4 Regular RNN
- 5 Forget Gate
- 6 Long Short-Term Memory (LSTM)
- 7 Backprop for an LSTM
- 8 Weight Gradient for an LSTM**
- 9 Conclusion

# Weight Gradient for an LSTM

Unlike BPTT, the weight gradient for an LSTM is really easy!  
There are a lot of different weights, but each of them has an influence on only one variable:

$$i[t] = \sigma_g(w_i x[t] + u_i h[t-1] + b_i)$$

$$o[t] = \sigma_g(w_o x[t] + u_o h[t-1] + b_o)$$

$$f[t] = \sigma_g(w_f x[t] + u_f h[t-1] + b_f)$$

$$c[t] = f[t]c[t-1] + i[t]\sigma_h(w_c x[t] + u_c h[t-1] + b_c)$$

$$h[t] = o[t]\sigma_h(c[t])$$



# Weight Gradient for an LSTM

Differentiating, we get:

$$\frac{\partial \mathcal{L}}{\partial w_i} = \sum_t \frac{d\mathcal{L}}{di[t]} \dot{\sigma}_g(\cdot) x[t], \quad \frac{\partial \mathcal{L}}{\partial u_i} = \sum_t \frac{d\mathcal{L}}{di[t]} \dot{\sigma}_g(\cdot) h[t-1],$$

$$\frac{\partial \mathcal{L}}{\partial w_o} = \sum_t \frac{d\mathcal{L}}{do[t]} \dot{\sigma}_g(\cdot) x[t], \quad \frac{\partial \mathcal{L}}{\partial u_o} = \sum_t \frac{d\mathcal{L}}{do[t]} \dot{\sigma}_g(\cdot) h[t-1],$$

$$\frac{\partial \mathcal{L}}{\partial w_f} = \sum_t \frac{d\mathcal{L}}{df[t]} \dot{\sigma}_g(\cdot) x[t], \quad \frac{\partial \mathcal{L}}{\partial u_f} = \sum_t \frac{d\mathcal{L}}{df[t]} \dot{\sigma}_g(\cdot) h[t-1],$$

$$\frac{\partial \mathcal{L}}{\partial w_c} = \sum_t \frac{d\mathcal{L}}{dc[t]} i[t] \dot{\sigma}_h(\cdot) x[t], \quad \frac{\partial \mathcal{L}}{\partial u_c} = \sum_t \frac{d\mathcal{L}}{dc[t]} i[t] \dot{\sigma}_h(\cdot) h[t-1]$$

# Outline

- 1 Review: Recurrent Neural Networks
- 2 Vanishing/Exploding Gradient
- 3 Running Example: a Pocket Calculator
- 4 Regular RNN
- 5 Forget Gate
- 6 Long Short-Term Memory (LSTM)
- 7 Backprop for an LSTM
- 8 Weight Gradient for an LSTM
- 9 Conclusion**

- RNNs suffer from either exponentially decreasing memory (if  $|w| < 1$ ) or exponentially increasing memory (if  $|w| > 1$ ). This is one version of a more general problem sometimes called the **gradient vanishing** problem.
- The forget gate solves that problem by making the feedback coefficient a function of the input.
- LSTM defines two types of memory (cell=excitation=“long-term memory,” and output=activation=“short-term memory”), and three types of gates (input, output, forget).

Each epoch of LSTM training has the same steps as in a regular RNN:

- 1 Forward propagation: find  $h[t]$ .
- 2 Synchronous backprop: find the time-synchronous partial derivatives  $\frac{\partial \mathcal{L}}{\partial h[t]}$ .
- 3 BPTT: find the total derivatives  $\frac{d\mathcal{L}}{dh[t]}$ ,  $\frac{d\mathcal{L}}{dc[t]}$ ,  $\frac{d\mathcal{L}}{do[t]}$ ,  $\frac{d\mathcal{L}}{df[t]}$ , and  $\frac{d\mathcal{L}}{di[t]}$ .
- 4 Weight gradient: find the weight gradients  $\frac{\partial \mathcal{L}}{\partial u_f}$  and so on.
- 5 Gradient descent: update the weights, e.g.,

$$u_f \leftarrow u_f - \eta \frac{\partial \mathcal{L}}{\partial u_f}$$