

Lecture 13. RNNs and LSTMs

Honglak Lee

02/24/2025



Outline

- RNN Basics**

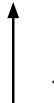
- Backpropagation through time (BPTT)
- The vanishing/exploding gradients problem
- Applications

RNN Basics

Sequence modeling

- Can we model sequences with standard NNs?

$$[\hat{y}_1, \hat{y}_2, \hat{y}_3, \hat{y}_4]$$



$$[x_1, x_2, x_3, x_4]$$

$$\text{where } x_t \in \mathbb{R}^d, h_t \in \mathbb{R}^m, y_t \in \mathbb{R}^n$$

5

RNN Basics

Sequence modeling

- Can we model sequences with standard NNs?
- Sure, but only if the length is fixed
- Input/output weight matrices grow with number of inputs/outputs!**

$$[\hat{y}_1, \hat{y}_2, \hat{y}_3, \hat{y}_4]$$

$$W^{out} \in \mathbb{R}^{4n \times m}$$

$$[x_1, x_2, x_3, x_4]$$

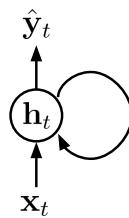
$$\text{where } x_t \in \mathbb{R}^d, h_t \in \mathbb{R}^m, y_t \in \mathbb{R}^n$$

7

RNN Basics

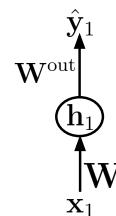
How can we model sequences in an efficient way?

- Use a model that reads each input one at a time
- Parameters can simply be reused (or shared) for each input in a recurrent computation



RNN Basics

Unrolling RNNs through time



$$h_1 = \sigma(Wx_1 + b)$$

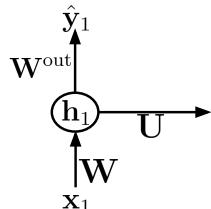
$$\hat{y}_1 = W^{out}h_1$$

10

11

RNN Basics

Unrolling RNNs through time

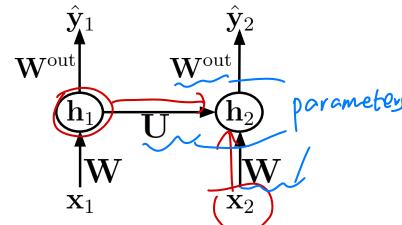


$$h_1 = \sigma(Wx_1 + b)$$

$$\hat{y}_1 = W^{out}h_1$$

RNN Basics

Unrolling RNNs through time



$$h_1 = \sigma(Wx_1 + b)$$

$$\hat{y}_1 = W^{out}h_1$$

$$h_2 = \sigma(Uh_1 + Wx_2 + b)$$

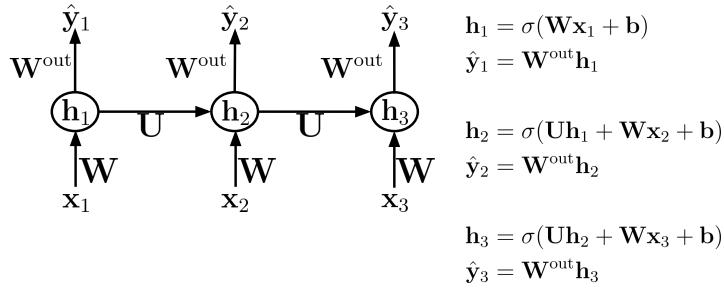
$$\hat{y}_2 = W^{out}h_2$$

12

13

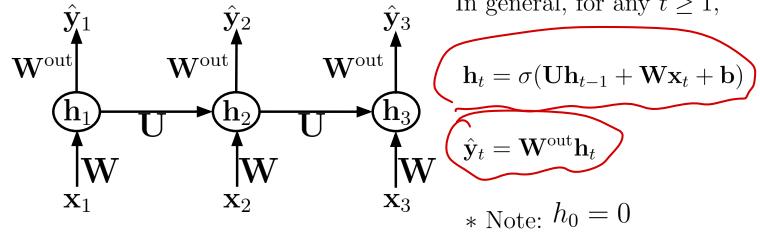
RNN Basics

Unrolling RNNs through time



RNN Basics

Unrolling RNNs through time

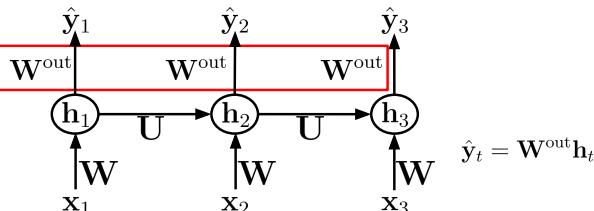


14

15

RNN Basics

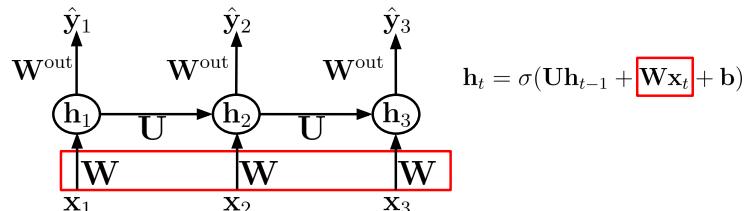
Unrolling RNNs through time



Weights are shared!

RNN Basics

Unrolling RNNs through time



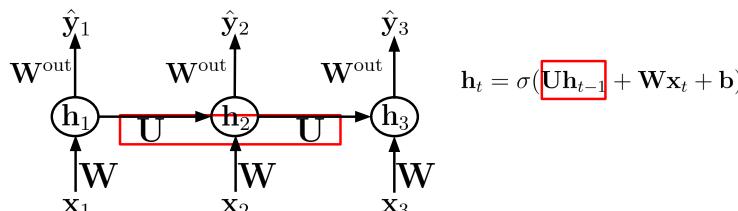
Weights are shared!

16

17

RNN Basics

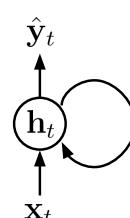
Unrolling RNNs through time



Weights are shared!

RNN Basics

Single step computation in RNNs



$$\begin{aligned} \mathbf{h}_t &= \sigma(\mathbf{U} \mathbf{h}_{t-1} + \mathbf{W} \mathbf{x}_t + \mathbf{b}) \\ \hat{y}_t &= \mathbf{W}^{\text{out}} \mathbf{h}_t \end{aligned}$$

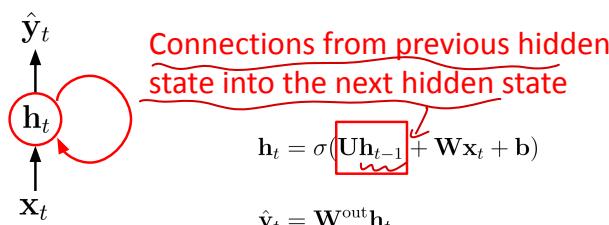
where $\mathbf{x}_t \in \mathbb{R}^d$, $\mathbf{h}_t \in \mathbb{R}^m$, $\mathbf{y}_t \in \mathbb{R}^n$,
 $\mathbf{U} \in \mathbb{R}^{m \times m}$, $\mathbf{W} \in \mathbb{R}^{m \times d}$, $\mathbf{W}^{\text{out}} \in \mathbb{R}^{n \times m}$, $\mathbf{b} \in \mathbb{R}^m$

18

19

RNN Basics

Single step computation in RNNs

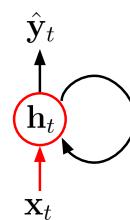


where $\mathbf{x}_t \in \mathbb{R}^d$, $\mathbf{h}_t \in \mathbb{R}^m$, $\mathbf{y}_t \in \mathbb{R}^n$,

$\mathbf{U} \in \mathbb{R}^{m \times m}$, $\mathbf{W} \in \mathbb{R}^{m \times d}$, $\mathbf{W}^{\text{out}} \in \mathbb{R}^{n \times m}$, $\mathbf{b} \in \mathbb{R}^m$

RNN Basics

Single step computation in RNNs



$$\begin{aligned} \mathbf{h}_t &= \sigma(\mathbf{U} \mathbf{h}_{t-1} + \mathbf{W} \mathbf{x}_t + \mathbf{b}) \\ \hat{y}_t &= \mathbf{W}^{\text{out}} \mathbf{h}_t \end{aligned}$$

where $\mathbf{x}_t \in \mathbb{R}^d$, $\mathbf{h}_t \in \mathbb{R}^m$, $\mathbf{y}_t \in \mathbb{R}^n$,

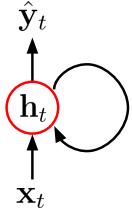
$\mathbf{U} \in \mathbb{R}^{m \times m}$, $\mathbf{W} \in \mathbb{R}^{m \times d}$, $\mathbf{W}^{\text{out}} \in \mathbb{R}^{n \times m}$, $\mathbf{b} \in \mathbb{R}^m$

20

21

RNN Basics

Single step computation in RNNs



Bias term

$$h_t = \sigma(Uh_{t-1} + Wx_t + b)$$

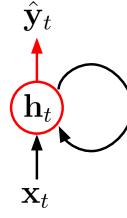
$$\hat{y}_t = W^{out}h_t$$

where $x_t \in \mathbb{R}^d$, $h_t \in \mathbb{R}^m$, $y_t \in \mathbb{R}^n$,
 $U \in \mathbb{R}^{m \times m}$, $W \in \mathbb{R}^{m \times d}$, $W^{out} \in \mathbb{R}^{n \times m}$, $b \in \mathbb{R}^m$

22

RNN Basics

Single step computation in RNNs



Connection from hidden unit to output

$$h_t = \sigma(Uh_{t-1} + Wx_t + b)$$

$$\hat{y}_t = W^{out}h_t$$

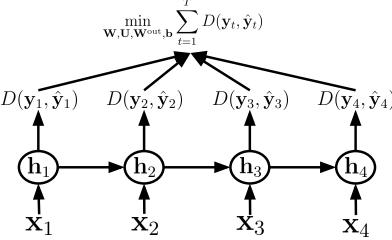
where $x_t \in \mathbb{R}^d$, $h_t \in \mathbb{R}^m$, $y_t \in \mathbb{R}^n$,
 $U \in \mathbb{R}^{m \times m}$, $W \in \mathbb{R}^{m \times d}$, $W^{out} \in \mathbb{R}^{n \times m}$, $b \in \mathbb{R}^m$

23

RNN Basics: Objective function

RNN general objective function

- Compute loss (if any) for each step in the sequence prediction and aggregate
- Gradient flows through all unrolled steps in the RNN



24

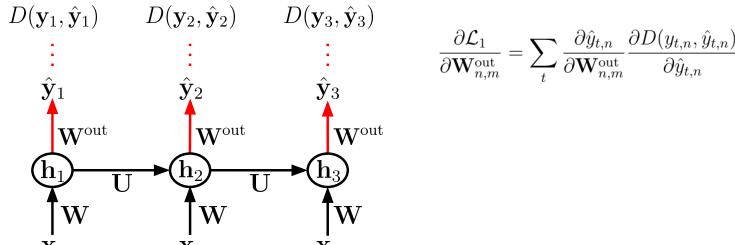
Outline

- RNN Basics
- Backpropagation through time (BPTT)**
- The vanishing/exploding gradients problem
- Applications

Backpropagation through time (BPTT)

- \mathcal{L}_t aggregates the loss for from time t to T:

$$\mathcal{L}_t = \sum_{\tau=t}^T D(y_\tau, \hat{y}_\tau)$$

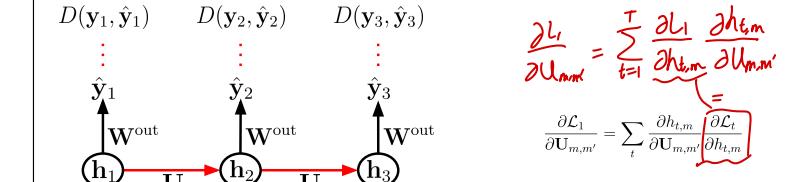


For the output parameters, we only consider the gradient of the current loss and add them up.

Backpropagation through time (BPTT)

- \mathcal{L}_t aggregates the loss for from time t to T:

$$\mathcal{L}_t = \sum_{\tau=t}^T D(y_\tau, \hat{y}_\tau)$$



For the other parameters, we use a chain rule involving the hidden activations at each time step and add them up.

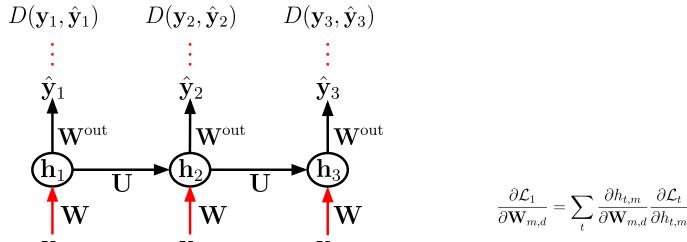
26

27

Backpropagation through time (BPTT)

- \mathcal{L}_t aggregates the loss for from time t to T:

$$\mathcal{L}_t = \sum_{\tau=t}^T D(y_\tau, \hat{y}_\tau)$$

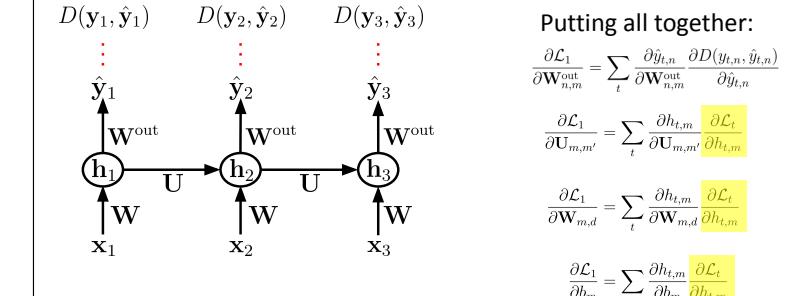


For the other parameters, we use a chain rule involving the hidden activations at each time step and add them up.

Backpropagation through time (BPTT)

- \mathcal{L}_t aggregates the loss for from time t to T:

$$\mathcal{L}_t = \sum_{\tau=t}^T D(y_\tau, \hat{y}_\tau)$$

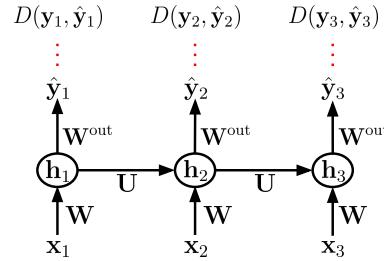


28

29

Backpropagation through time (BPTT)

- \mathcal{L}_t aggregates the loss for from time t to T :
$$\mathcal{L}_t = \sum_{\tau=t}^T D(\mathbf{y}_\tau, \hat{\mathbf{y}}_\tau)$$



Summary (simpler notation):

$$\text{if } \theta \in \{\mathbf{W}^{\text{out}}\}$$

$$\frac{\partial \mathcal{L}_1}{\partial \theta} = \sum_t \frac{\partial \hat{y}_{t,n}}{\partial \theta} \frac{\partial D(y_{t,n}, \hat{y}_{t,n})}{\partial \hat{y}_{t,n}}$$

$$\text{if } \theta \in \{\mathbf{W}, \mathbf{U}, \mathbf{b}\}$$

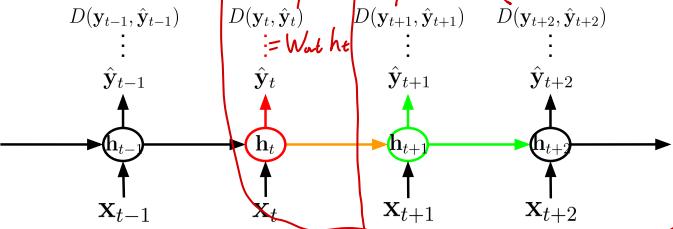
$$\frac{\partial \mathcal{L}_1}{\partial \theta} = \sum_t \sum_m \frac{\partial h_{t,m}}{\partial \theta} \frac{\partial \mathcal{L}_t}{\partial h_{t,m}}$$

Note: we are not done yet! We still need to calculate $\frac{\partial \mathcal{L}_t}{\partial h_{t,m}}$ (see the next slides)

31

Backprop through time for calculating $\frac{\partial \mathcal{L}_t}{\partial h_{t,m}}$

- The general recursion formula for the gradient w.r.t. h_t :



recursion (for general t):

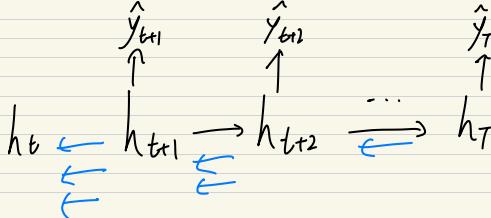
$$\frac{\partial \mathcal{L}_t}{\partial h_{t,m}} = \frac{\partial D(\mathbf{y}_t, \hat{\mathbf{y}}_t)}{\partial h_{t,m}} + \sum_{m'} \frac{\partial h_{t+1,m'}}{\partial h_{t,m}} \frac{\partial \mathcal{L}_{t+1}}{\partial h_{t+1,m'}}$$

$$\text{where } \mathcal{L}_t = \sum_{\tau=t}^T D(\mathbf{y}_\tau, \hat{\mathbf{y}}_\tau)$$

37

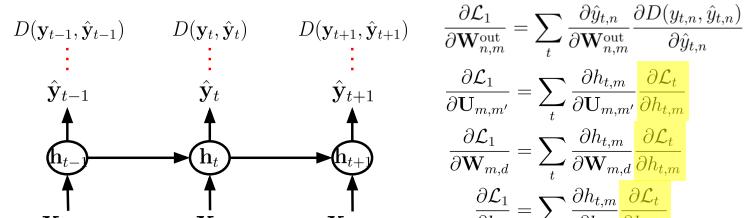
$$\frac{\partial \mathcal{L}_{t+1}}{\partial h_{t,m}} = \frac{\partial \mathcal{L}_{t+1}}{\partial h_{t+1}} \frac{\partial h_{t+1}}{\partial h_{t,m}} = \sum_{m'} \frac{\partial \mathcal{L}_{t+1}}{\partial h_{t+1,m'}} \frac{\partial h_{t+1,m'}}{\partial h_{t,m}}$$

$$\left(\text{since } \mathbf{y}_t = \mathbf{W}^{\text{out}} \mathbf{h}_t, \forall \tau \geq t+1 \Rightarrow \frac{\partial \mathbf{y}_t}{\partial h_{t,m}} = \frac{\partial \mathbf{y}_t}{\partial h_{t+1}} \frac{\partial h_{t+1}}{\partial h_{t,m}} \right)$$



Backpropagation through time (BPTT)

- Full summary: $\mathcal{L}_t = \sum_{\tau=t}^T D(\mathbf{y}_\tau, \hat{\mathbf{y}}_\tau)$



$$\text{Recursion} \rightarrow \frac{\partial \mathcal{L}_t}{\partial h_{t,m}} = \frac{\partial D(\mathbf{y}_t, \hat{\mathbf{y}}_t)}{\partial h_{t,m}} + \sum_{m'} \frac{\partial h_{t+1,m'}}{\partial h_{t,m}} \frac{\partial \mathcal{L}_{t+1}}{\partial h_{t+1,m'}}$$

39

Outline

- RNN Basics
- Backpropagation through time (BPTT)
- The vanishing/exploding gradients problem**
- Applications

The vanishing/exploding gradients problem

- The hidden-to-hidden connections in standard RNNs can cause gradient vanishing during backprop of the loss in the last step.

$$\begin{aligned} \frac{\partial \mathcal{L}_T}{\partial \mathbf{h}_t} &= \left[\frac{\partial \mathbf{h}_T}{\partial \mathbf{h}_t} \right]^\top \frac{\partial \mathcal{L}_T}{\partial \mathbf{h}_T} \\ &= \left(\prod_{i=t+1}^T \left[\frac{\partial \mathbf{h}_i}{\partial \mathbf{h}_{i-1}} \right]^\top \right) \frac{\partial \mathcal{L}_T}{\partial \mathbf{h}_T} = \left(\prod_{i=t+1}^T [\text{diag}(\mathbf{h}_i(1-\mathbf{h}_i)) \mathbf{U}]^\top \right) \frac{\partial \mathcal{L}_T}{\partial \mathbf{h}_T} \\ &= \left(\prod_{i=t+1}^T \mathbf{U}^\top \text{diag}(\mathbf{h}_i(1-\mathbf{h}_i)) \right) \frac{\partial \mathcal{L}_T}{\partial \mathbf{h}_T} \end{aligned}$$

Note: $\frac{\partial \mathcal{L}_T}{\partial \mathbf{h}_t}$: mx1 vector $\left[\frac{\partial \mathcal{L}_T}{\partial \mathbf{h}_t} \right]_i = \frac{\partial \mathcal{L}_T}{\partial \mathbf{h}_{t,i}}$ $\frac{\partial \mathbf{h}_T}{\partial \mathbf{h}_t}$: mxm matrix $\left[\frac{\partial \mathbf{h}_T}{\partial \mathbf{h}_t} \right]_{i,j} = \frac{\partial \mathbf{h}_T}{\partial \mathbf{h}_{t,j}}$

44

The vanishing/exploding gradients problem

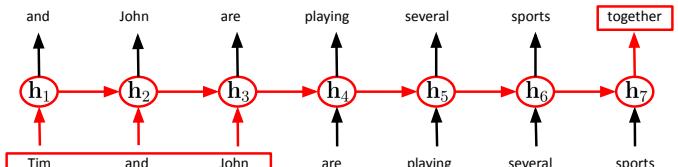
- The hidden-to-hidden connections in standard RNNs can cause gradient vanishing during backprop of the loss in the last step.
- This is caused by the gradients with respect to the activation function being multiplied through time!
- Gradient can explode if the norm of \mathbf{U} is large

$$\begin{aligned} \frac{\partial \mathcal{L}_T}{\partial \mathbf{h}_t} &= \left[\frac{\partial \mathbf{h}_T}{\partial \mathbf{h}_t} \right]^\top \frac{\partial \mathcal{L}_T}{\partial \mathbf{h}_T} \\ &= \left(\prod_{i=t+1}^T \mathbf{U}^\top \text{diag}(\mathbf{h}_i(1-\mathbf{h}_i)) \right) \frac{\partial \mathcal{L}_T}{\partial \mathbf{h}_T} \end{aligned}$$

(Gradient can also vanish if the norm of \mathbf{U} is small.)

The vanishing/exploding gradients problem

- The hidden-to-hidden connections in standard RNNs can cause gradient vanishing during backprop of the loss in the last step.
- This is caused by the gradients with respect to the activation function being multiplied through time!
- Gradient can also explode if the norm of weight (e.g., \mathbf{U}) is large.
- Long term dependencies in the sequence become affected!



46

47

The vanishing/exploding gradients problem

Gradient clipping

- A simple trick to prevent gradient explosion is to limit them to a max value.
- If the gradient is larger than η , then clip the norm to η :

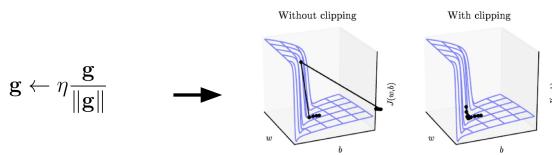


Diagram credit: Goodfellow et al., Deep Learning

48

The vanishing/exploding gradients problem

Additional tricks

- Identity initialization:** Initialize the RNN weights to be the identity function, initialize the bias to zero and use ReLU activation.
 - This helps with stability during training
- Weight Regularization:** Regularize the RNN weights to prevent large activations.
 - This prevents exploding gradients.

The vanishing/exploding gradients problem

Initialization tricks: Weight orthogonality*

- Orthogonal matrices have the property of preserving norms:
$$\|Uh\| = \|h\|$$
- By the definition of operator norms, given matrices **A** and **B**, and vector **v**, we have:
$$\|Av\| \leq \|A\|\|v\| \quad \text{and} \quad \|AB\| \leq \|A\|\|B\|$$

* Recall: operator norm for a matrix A

$$\begin{aligned} \|A\|_{op} &= \inf\{c \geq 0 : \|Av\| \leq c\|v\| \text{ for all } v \in R^d\} \\ &= \sup_{\{v \in R^d, v \neq 0\}} \frac{\|Av\|}{\|v\|} \end{aligned}$$

51

The vanishing/exploding gradients problem

Initialization tricks: Weight orthogonality*

- Now, given the gradient of the loss with respect to h_t :
$$\frac{\partial \mathcal{L}_T}{\partial h_T} = \left[\frac{\partial h_T}{\partial h_t} \right]^\top \frac{\partial \mathcal{L}_T}{\partial h_t} = \prod_{k=t}^{T-1} \left[\frac{\partial h_{k+1}}{\partial h_k} \right]^\top \frac{\partial \mathcal{L}_T}{\partial h_T} = \prod_{k=t}^{T-1} [U^\top D_{k+1}] \frac{\partial \mathcal{L}_T}{\partial h_T} = \prod_{k=t}^{T-1} [U^\top D_{k+1}] \frac{\partial \mathcal{L}_T}{\partial h_T}$$
- Where D_{k+1} is a diagonal matrix of activation function gradients. Therefore, given an orthogonal weight matrix **U**, we have:

$$\left\| \frac{\partial \mathcal{L}_T}{\partial h_t} \right\| = \left\| \left(\prod_{k=t}^{T-1} U^\top D_{k+1} \right) \frac{\partial \mathcal{L}_T}{\partial h_T} \right\| \leq \left(\prod_{k=t}^{T-1} \|U^\top D_{k+1}\| \right) \left\| \frac{\partial \mathcal{L}_T}{\partial h_T} \right\| = \left(\prod_{k=t}^{T-1} \|D_{k+1}\| \right) \left\| \frac{\partial \mathcal{L}_T}{\partial h_T} \right\|$$

\uparrow
 $\|Av\| \leq \|A\|\|v\|$

Arjovsky, M., Shah, A., & Bengio, Y. Unitary evolution recurrent neural networks. In ICML 2016.

57

The vanishing/exploding gradients problem

Initialization tricks: Weight orthogonality*

- Now, given the gradient of the loss with respect to h_t :
$$\frac{\partial \mathcal{L}_T}{\partial h_t} = \left[\frac{\partial h_T}{\partial h_t} \right]^\top \frac{\partial \mathcal{L}_T}{\partial h_T} = \prod_{k=t}^{T-1} \left[\frac{\partial h_{k+1}}{\partial h_k} \right]^\top \frac{\partial \mathcal{L}_T}{\partial h_T} = \prod_{k=t}^{T-1} [U^\top D_{k+1}] \frac{\partial \mathcal{L}_T}{\partial h_T}$$
- Where D_{k+1} is a diagonal matrix of activation function gradients. Therefore, given an orthogonal weight matrix **U**, we have:

$$\left\| \frac{\partial \mathcal{L}_T}{\partial h_t} \right\| = \left\| \left(\prod_{k=t}^{T-1} U^\top D_{k+1} \right) \frac{\partial \mathcal{L}_T}{\partial h_T} \right\| \leq \left(\prod_{k=t}^{T-1} \boxed{\|U^\top D_{k+1}\|} \right) \left\| \frac{\partial \mathcal{L}_T}{\partial h_T} \right\| = \left(\prod_{k=t}^{T-1} \boxed{\|D_{k+1}\|} \right) \left\| \frac{\partial \mathcal{L}_T}{\partial h_T} \right\|$$

\uparrow
 $\|Uh\| = \|h\|$

58

Arjovsky, M., Shah, A., & Bengio, Y. Unitary evolution recurrent neural networks. In ICML 2016.

60

The vanishing/exploding gradients problem

Long Short-Term Memory (LSTM)

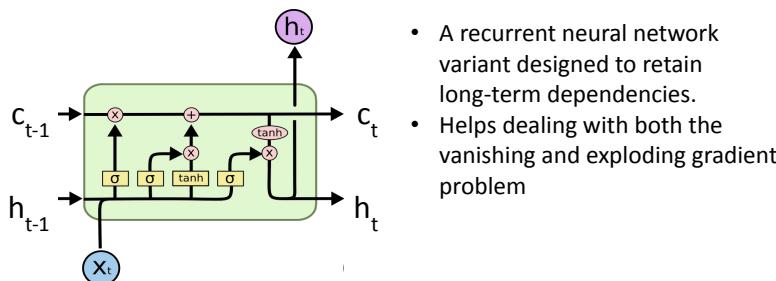


Diagram credit: <http://colah.github.io/posts/2015-08-Understanding-LSTMs>

The vanishing/exploding gradients problem

Long Short-Term Memory (LSTM)

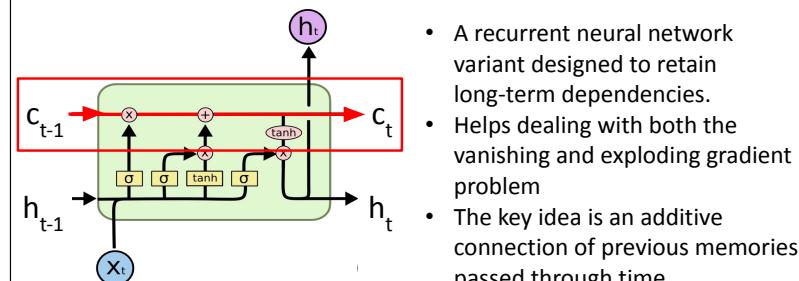


Diagram credit: <http://colah.github.io/posts/2015-08-Understanding-LSTMs>

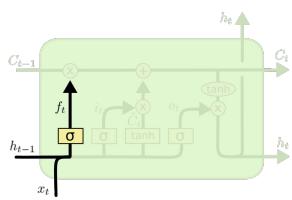
61

62

The vanishing/exploding gradients problem

Long Short-Term Memory (LSTM)

- The forget gate allows LSTM to choose to zero out part of previous memories and let others through.



$$f_t = \sigma(W_f \cdot [h_{t-1}, x_t] + b_f)$$

Notation: \cdot is a matrix/vector product
 $*$ is an element-wise product

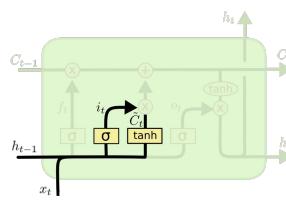
Diagram credit: <http://colah.github.io/posts/2015-08-Understanding-LSTMs>

63

The vanishing/exploding gradients problem

Long Short-Term Memory (LSTM)

- The input gate behaves similar to the forget gate with new inputs.



$$i_t = \sigma(W_i \cdot [h_{t-1}, x_t] + b_i)$$

$$\tilde{C}_t = \tanh(W_C \cdot [h_{t-1}, x_t] + b_C)$$

Notation: \cdot is a matrix/vector product
 $*$ is an element-wise product

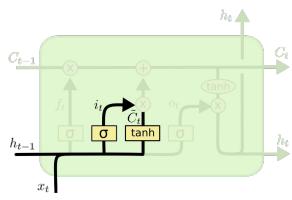
Diagram credit: <http://colah.github.io/posts/2015-08-Understanding-LSTMs>

64

The vanishing/exploding gradients problem

Long Short-Term Memory (LSTM)

- The input gate behaves similar to the forget gate with new inputs.
- New information is computed from the current input and previous hidden units.



$$i_t = \sigma(W_i \cdot [h_{t-1}, x_t] + b_i)$$

$$\tilde{C}_t = \tanh(W_C \cdot [h_{t-1}, x_t] + b_C)$$

Notation: \cdot is a matrix/vector product
 $*$ is an element-wise product

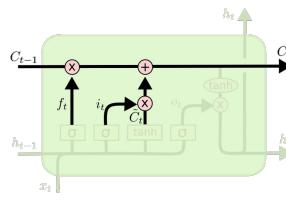
Diagram credit: <http://colah.github.io/posts/2015-08-Understanding-LSTMs>

65

The vanishing/exploding gradients problem

Long Short-Term Memory (LSTM)

- Memories to be passed to the next steps are computed using the forget gate on the previous memories and the input gate on the current information found in the sequence



$$f_t = \sigma(W_f \cdot [h_{t-1}, x_t] + b_f)$$

$$C_t = f_t * C_{t-1} + i_t * \tilde{C}_t$$

Notation: \cdot is a matrix/vector product
 $*$ is an element-wise product

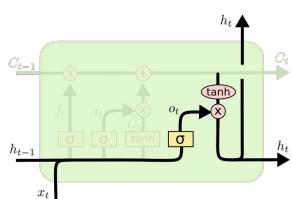
Diagram credit: <http://colah.github.io/posts/2015-08-Understanding-LSTMs>

66

The vanishing/exploding gradients problem

Long Short-Term Memory (LSTM)

- An output gate is computed to choose information from the current memories for the next hidden state in the LSTM.



$$o_t = \sigma(W_o \cdot [h_{t-1}, x_t] + b_o)$$

$$h_t = o_t * \tanh(C_t)$$

Notation: \cdot is a matrix/vector product
 $*$ is an element-wise product

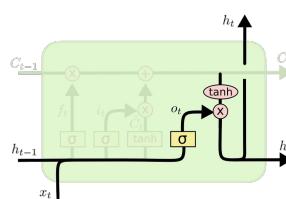
Diagram credit: <http://colah.github.io/posts/2015-08-Understanding-LSTMs>

67

The vanishing/exploding gradients problem

Long Short-Term Memory (LSTM)

- An output gate is computed to choose information from the current memories for the next hidden state in the LSTM.
- Next hidden state is computed from the current memories and gate.



$$o_t = \sigma(W_o \cdot [h_{t-1}, x_t] + b_o)$$

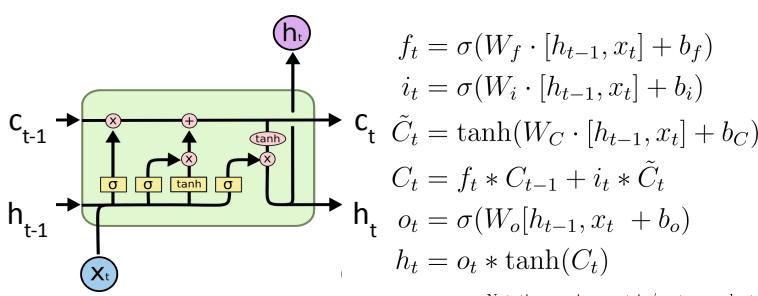
$$h_t = o_t * \tanh(C_t)$$

Notation: \cdot is a matrix/vector product
 $*$ is an element-wise product

Diagram credit: <http://colah.github.io/posts/2015-08-Understanding-LSTMs>

68

Long Short-Term Memory (LSTM)



$$f_t = \sigma(W_f \cdot [h_{t-1}, x_t] + b_f)$$

$$i_t = \sigma(W_i \cdot [h_{t-1}, x_t] + b_i)$$

$$\tilde{C}_t = \tanh(W_C \cdot [h_{t-1}, x_t] + b_C)$$

$$C_t = f_t * C_{t-1} + i_t * \tilde{C}_t$$

$$o_t = \sigma(W_o \cdot [h_{t-1}, x_t] + b_o)$$

$$h_t = o_t * \tanh(C_t)$$

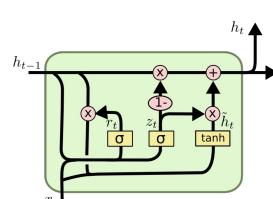
Notation: \cdot is a matrix/vector product
 $*$ is an element-wise product

Diagram credit: <http://colah.github.io/posts/2015-08-Understanding-LSTMs>

Gated Recurrent Unit (GRU)

- A simplified variation of LSTM

- The input / forget gates are simplified into a single gate z_t



$$z_t = \sigma(W_z \cdot [h_{t-1}, x_t])$$

$$r_t = \sigma(W_r \cdot [h_{t-1}, x_t])$$

$$\tilde{h}_t = \tanh(W \cdot [r_t * h_{t-1}, x_t])$$

$$h_t = (1 - z_t) * h_{t-1} + z_t * \tilde{h}_t$$

Notation: \cdot is a matrix/vector product
 $*$ is an element-wise product

Diagram credit: <http://colah.github.io/posts/2015-08-Understanding-LSTMs>

70

Gated Recurrent Unit (GRU)

- A simplified variation of LSTM
 - The input / forget gates are simplified into a single gate z_t
 - C_t removed and consolidated with h_t
 - r_t introduced for gating (similar role as o_t for LSTM)
- Diagram credit: <http://colah.github.io/posts/2015-08-Understanding-LSTMs>
- Notation: \cdot is a matrix/vector product
 \otimes is an element-wise product
- Equations:
- $$z_t = \sigma(W_z \cdot [h_{t-1}, x_t])$$
- $$r_t = \sigma(W_r \cdot [h_{t-1}, x_t])$$
- $$\tilde{h}_t = \tanh(W \cdot [r_t * h_{t-1}, x_t])$$
- $$h_t = (1 - z_t) * h_{t-1} + z_t * \tilde{h}_t$$

Multi-layer RNN

- To make RNNs more powerful, we can stack multiple layers of RNNs on top of each other

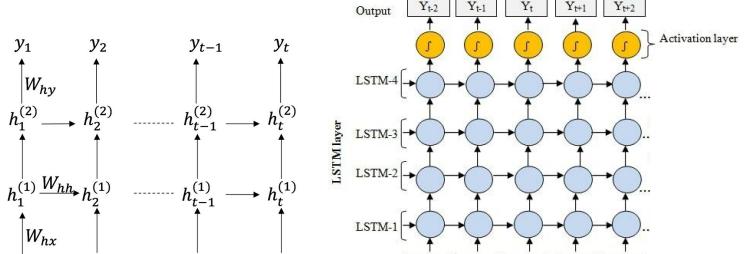


Diagram credit: <https://arxiv.org/pdf/1707.08262.pdf>

Diagram credit: https://link.springer.com/chapter/10.1007/978-3-030-86970-0_24

Outline

- RNN Basics
- Backpropagation through time (BPTT)
- The vanishing/exploding gradients problem
- Applications**

Applications

Action Recognition

- Determine the action occurring in frame sequences

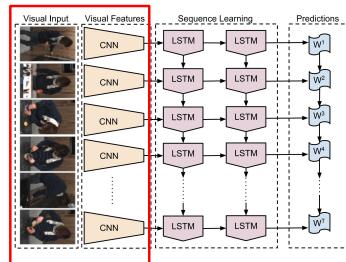


Figure credit: <http://www.thumos.info>

Applications

Action Recognition

- Encode images in compact representation



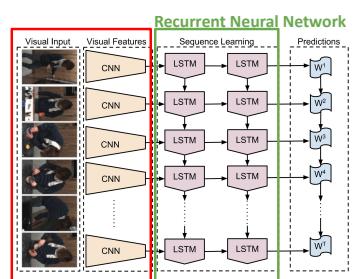
Convolutional Neural Network

Figure credit: Jeff Donahue et al, 2015

Applications

Action Recognition

- Encode images in compact representation
- Feed them to a multi-layer LSTM



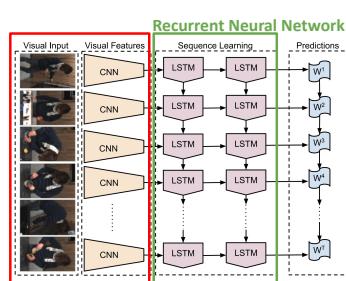
Convolutional Neural Network

Figure credit: Jeff Donahue et al, 2015

Applications

Action Recognition

- Encode images in compact representation
- Feed them to a multi-layer LSTM
- Average predictions

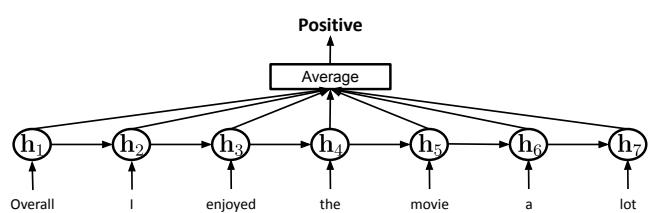


Convolutional Neural Network

Figure credit: Jeff Donahue et al, 2015

Applications: Sentence Classification

- You can train RNNs on text!
- Read a full sentence (or paragraph) and give it a rating
- Example: Movie reviews

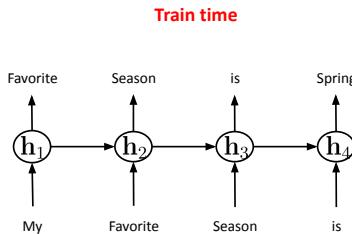


71

72

Applications: Language modeling

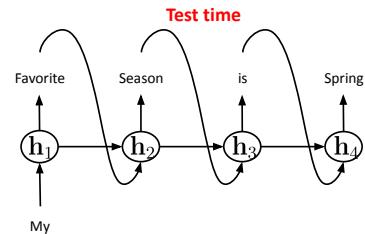
- You can train RNNs on text!
- Let the RNN read one word at a time and try to predict the next word
- This task requires long-term dependencies between words
- It can generate as many words as we let it generate



Slide credit: Abigail See 79

Applications: Language modeling

- You can train RNNs on text!
- Let the RNN read one word at a time and try to predict the next word
- This task requires long-term dependencies between words
- It can generate as many words as we let it generate



Slide credit: Abigail See 80

Applications: Language modeling

- You can train RNNs on any kind of text, and have it generate the same the style of text
- RNN trained on **Obama** speeches



The United States will step up to the cost of a new challenges of the American people that will share the fact that we created the problem. They were attacked and so that they have to say that all the task of the final days of war that I will not be able to get this done.

Slide credit: Abigail See 81

Applications: Language modeling

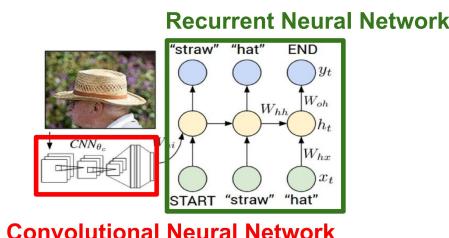
- You can train RNNs on any kind of text, and have it generate the same the style of text
- RNN trained on **Harry Potter** scripts



"No idea," said Nearly Headless Nick, casting low close by Cedric, carrying the last bit of treacle Charms, from Harry's shoulder, and to answer him the common room perched upon it, four arms held a shining knob from when the spider hadn't felt it seemed. He reached the teams too.

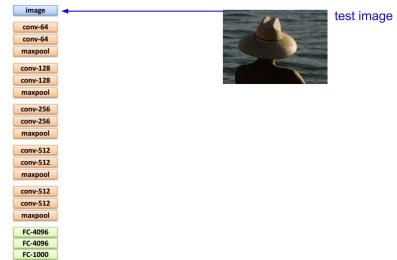
Slide credit: Abigail See 82

Applications: Image Captioning



Slide credit: Fei fei Li, Justin Johnson and Serena Yeung 83

Applications: Image Captioning



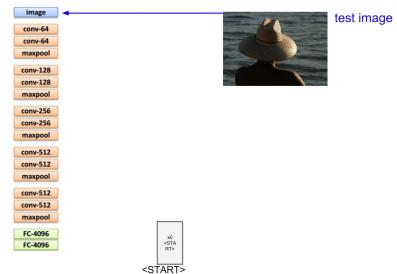
Slide credit: Fei fei Li, Justin Johnson and Serena Yeung 84

Applications: Image Captioning



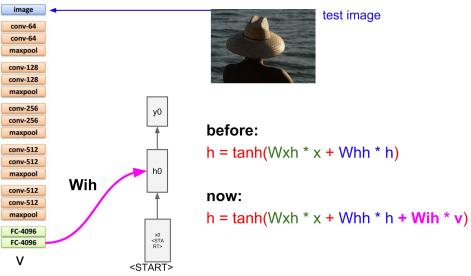
Slide credit: Fei fei Li, Justin Johnson and Serena Yeung 85

Applications: Image Captioning



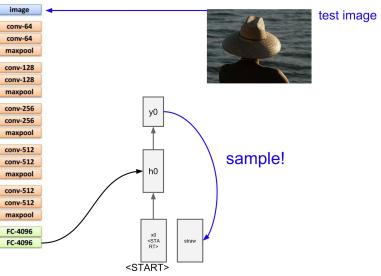
Slide credit: Fei fei Li, Justin Johnson and Serena Yeung 86

Applications: Image Captioning



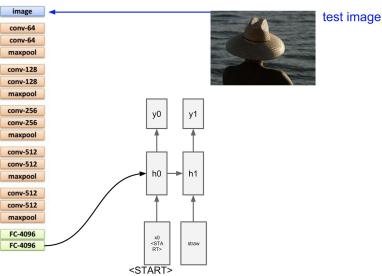
Slide credit: Fei fei Li, Justin Johnson and Serena Yeung 87

Applications: Image Captioning



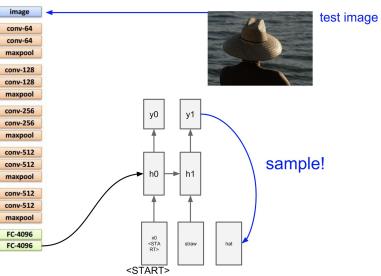
Slide credit: Fei fei Li, Justin Johnson and Serena Yeung 88

Applications: Image Captioning



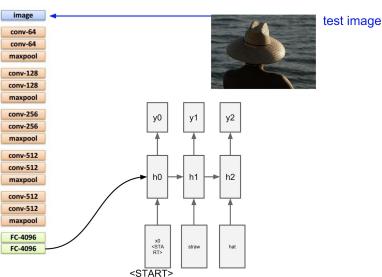
Slide credit: Fei fei Li, Justin Johnson and Serena Yeung 89

Applications: Image Captioning



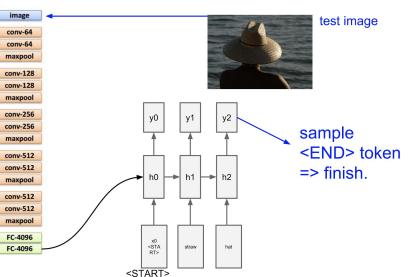
Slide credit: Fei fei Li, Justin Johnson and Serena Yeung 90

Applications: Image Captioning



Slide credit: Fei fei Li, Justin Johnson and Serena Yeung 91

Applications: Image Captioning



Slide credit: Fei fei Li, Justin Johnson and Serena Yeung 92

Applications: Image Captioning

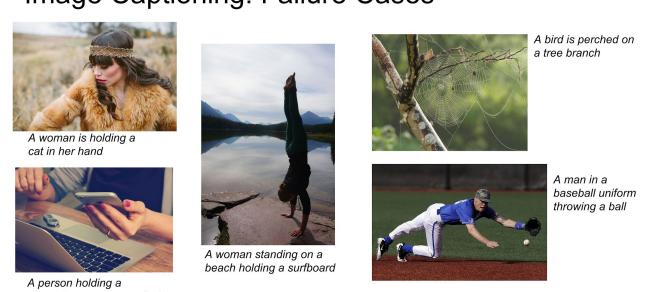
Image Captioning: Example Results



Slide credit: Fei fei Li, Justin Johnson and Serena Yeung 93

Applications: Image Captioning

Image Captioning: Failure Cases



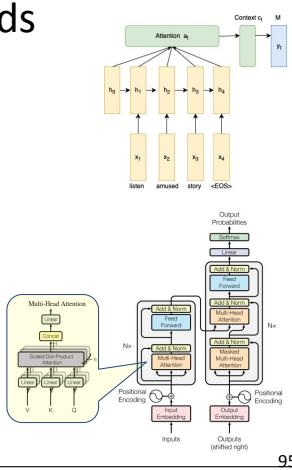
Slide credit: Fei fei Li, Justin Johnson and Serena Yeung 94

Remarks on SOTA methods

- Recent RNN/LSTM models incorporate attention mechanism, which allows for better modeling of long-term dependencies and more stable gradient flow.
- Many of the current SOTA methods are based on Transformer architectures, which use self-attention mechanism at the core.
(Next lecture!)

References:

- [Attention in RNNs](#)
- [Machine Translation With Sequence To Sequence Models](#)
- [Attention Is All You Need](#), 2017. (The first paper which proposed the transformer architecture.)
- [The Illustrated Transformer](#)



95

Summary

- RNN Basics
- Backpropagation through time (BPTT)
 - Backprop with recursive computation due to hidden-to-hidden connections over time
- The vanishing/exploding gradients problem
 - Tricks: Gradient clipping, Weight Initialization, Orthogonal weight matrices, LSTM/GRU, Transformers, etc.
- Applications
 - Video classification, text classification, image/video captioning, machine translation, etc.

96