

Recurrent Neural Networks

CS114B Lab 11

Kenneth Lai

April 23, 2021

Sequence Labeling

- ▶ Suppose we observe a list of words X . What are the respective parts of speech Y ?

Sequence Labeling

- ▶ Suppose we observe a list of words X . What are the respective parts of speech Y ?
 - ▶ What is $P(Y|X)$?

Sequence Labeling

- ▶ Suppose we observe a list of words X . What are the respective parts of speech Y ?
 - ▶ What is $P(Y|X)$?
- ▶ Generative approach: Hidden Markov Models

Sequence Labeling

- ▶ Suppose we observe a list of words X . What are the respective parts of speech Y ?
 - ▶ What is $P(Y|X)$?
- ▶ Generative approach: Hidden Markov Models
- ▶ Discriminative approaches:
 - ▶ Conditional random fields
 - ▶ Structured perceptrons

Sequence Labeling

- ▶ Suppose we observe a list of words X . What are the respective parts of speech Y ?
 - ▶ What is $P(Y|X)$?
- ▶ Generative approach: Hidden Markov Models
- ▶ Discriminative approaches:
 - ▶ Conditional random fields
 - ▶ Structured perceptrons
 - ▶ Neural networks

Sequence Labeling

- ▶ Discriminative approaches:

Sequence Labeling

- ▶ Discriminative approaches:
 - ▶ At each time step, use local features to compute local scores, and use the Viterbi algorithm to make predictions for the whole sentence

Sequence Labeling

- ▶ Discriminative approaches:
 - ▶ At each time step, use local features to compute local scores, and use the Viterbi algorithm to make predictions for the whole sentence
 - ▶ Conditional random fields
 - ▶ Structured perceptrons

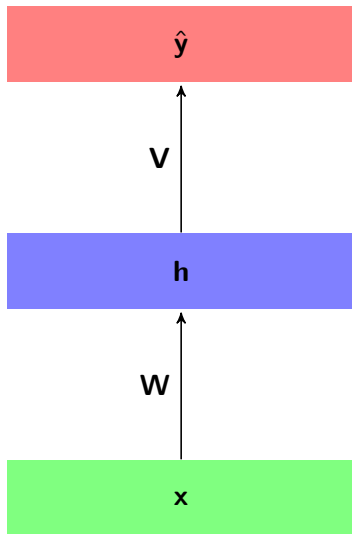
Sequence Labeling

- ▶ Discriminative approaches:
 - ▶ At each time step, use local features to compute local scores, and use the Viterbi algorithm to make predictions for the whole sentence
 - ▶ Conditional random fields
 - ▶ Structured perceptrons
 - ▶ Use features from other time steps, but make independent predictions at each time step

Sequence Labeling

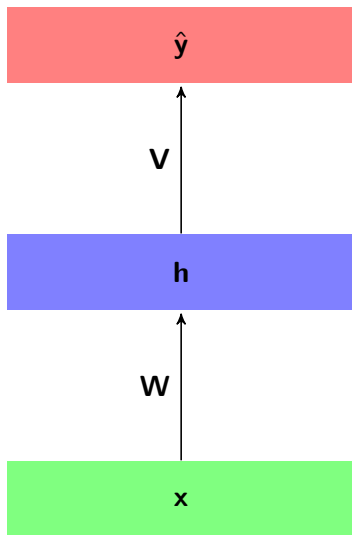
- ▶ Discriminative approaches:
 - ▶ At each time step, use local features to compute local scores, and use the Viterbi algorithm to make predictions for the whole sentence
 - ▶ Conditional random fields
 - ▶ Structured perceptrons
 - ▶ Use features from other time steps, but make independent predictions at each time step
 - ▶ Neural networks

Feedforward Neural Networks



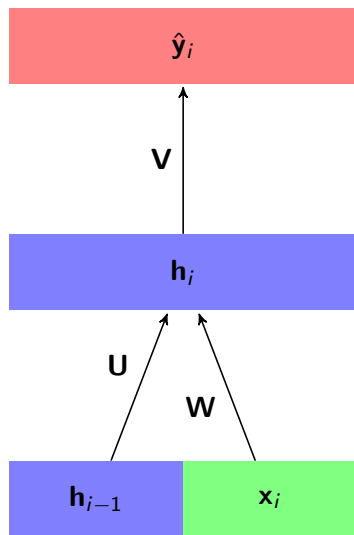
- ▶ Output layer
- ▶ Hidden layer(s)
- ▶ Input layer
- ▶ $h = g(x \cdot W)$
- ▶ $\hat{y} = g(h \cdot V)$

Feedforward Neural Networks

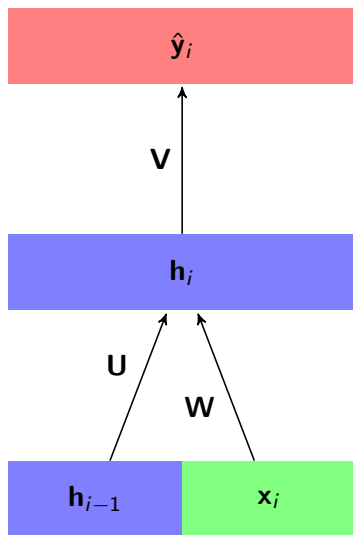


- ▶ Output layer
- ▶ Hidden layer(s)
- ▶ Input layer
- ▶ $h = g(x \cdot W)$
- ▶ $\hat{y} = g(h \cdot V)$
 - ▶ We will assume that the dummy feature 1 is part of x and h , and that the bias term is part of W and V , etc.

Neural Networks for Sequence Labeling

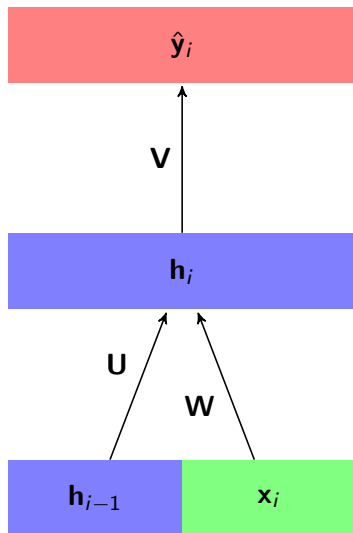


Neural Networks for Sequence Labeling



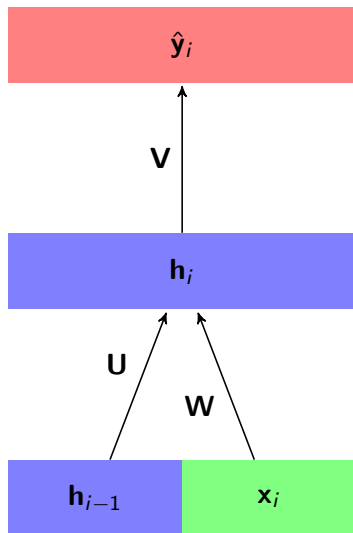
- At each time i , the input to the neural network consists of:

Neural Networks for Sequence Labeling



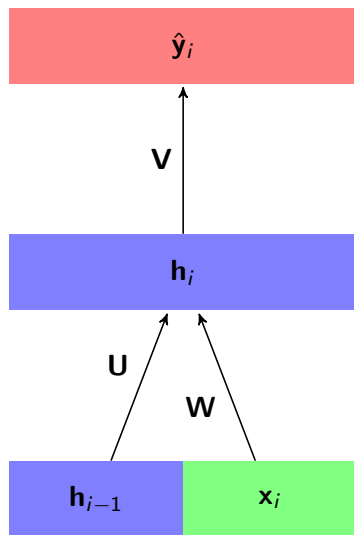
- ▶ At each time i , the input to the neural network consists of:
 - ▶ Current word (or other input) vector x_i

Neural Networks for Sequence Labeling



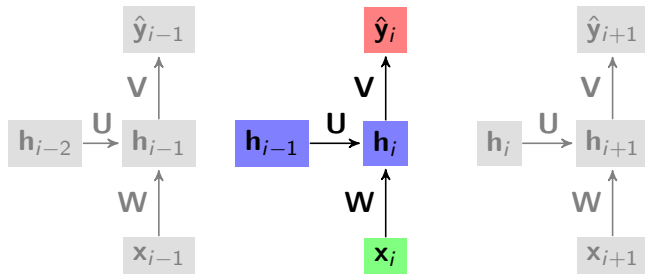
- ▶ At each time i , the input to the neural network consists of:
 - ▶ Current word (or other input) vector \mathbf{x}_i
 - ▶ History/(past) context vector \mathbf{h}_{i-1}

Neural Networks for Sequence Labeling

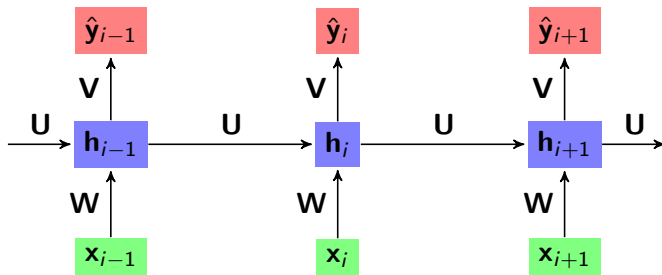


- ▶ At each time i , the input to the neural network consists of:
 - ▶ Current word (or other input) vector \mathbf{x}_i
 - ▶ History/(past) context vector \mathbf{h}_{i-1}
- ▶ $\mathbf{h}_i = g(\mathbf{x}_i \cdot \mathbf{W} + \mathbf{h}_{i-1} \cdot \mathbf{U})$

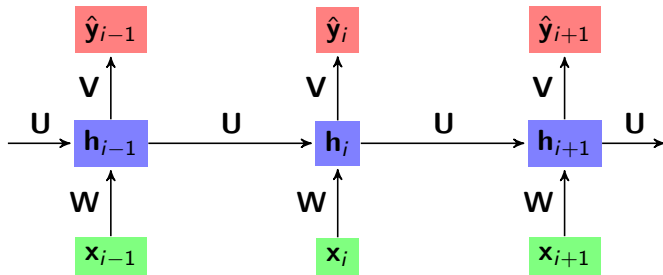
Neural Networks for Sequence Labeling



Neural Networks for Sequence Labeling

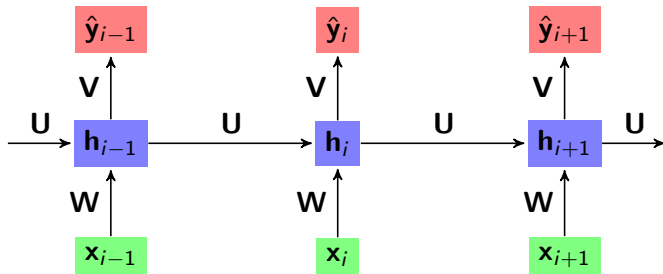


Neural Networks for Sequence Labeling



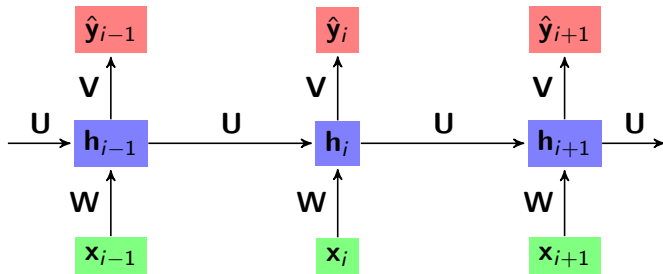
- The output of the hidden state at one time step is the history/past context input for the next time step!

Neural Networks for Sequence Labeling



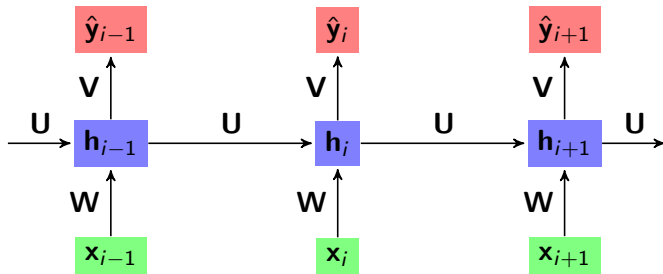
- ▶ The output of the hidden state at one time step is the history/past context input for the next time step!
- ▶ What context information is embedded in \mathbf{h}_{i-1} ?

Neural Networks for Sequence Labeling



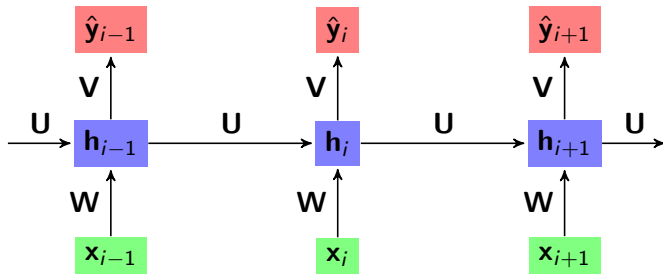
- ▶ The output of the hidden state at one time step is the history/past context input for the next time step!
- ▶ What context information is embedded in \mathbf{h}_{i-1} ?
 - ▶ Previous word \mathbf{x}_{i-1}
 - ▶ Previous context \mathbf{h}_{i-2}

Neural Networks for Sequence Labeling



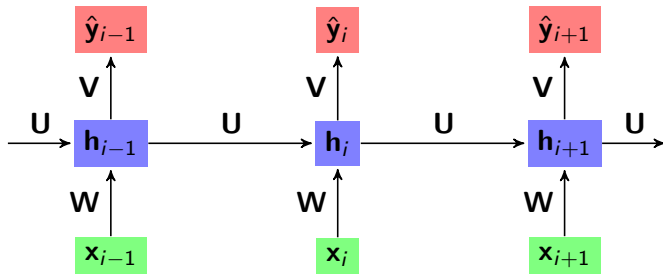
- ▶ The output of the hidden state at one time step is the history/past context input for the next time step!
- ▶ What context information is embedded in h_{i-1} ?
 - ▶ Previous word x_{i-1}
 - ▶ Previous context h_{i-2}
 - ▶ Previous previous word x_{i-2}
 - ▶ Previous previous context h_{i-3}

Neural Networks for Sequence Labeling



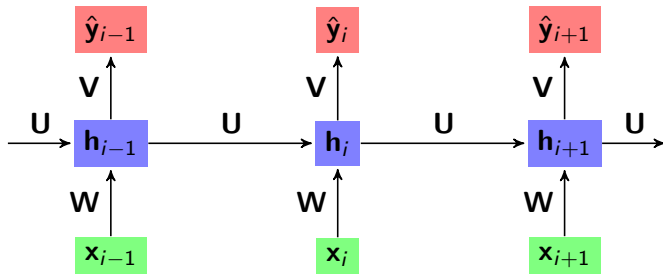
- ▶ The output of the hidden state at one time step is the history/past context input for the next time step!
- ▶ What context information is embedded in \mathbf{h}_{i-1} ?
 - ▶ All previous words

Neural Networks for Sequence Labeling



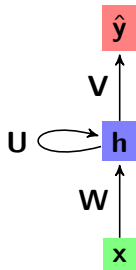
- ▶ The output of the hidden state at one time step is the history/past context input for the next time step!
- ▶ What context information is embedded in \mathbf{h}_{i-1} ?
 - ▶ All previous words
 - ▶ What about previous parts of speech (as in HMMs, CRFs, structured perceptrons)?

Neural Networks for Sequence Labeling

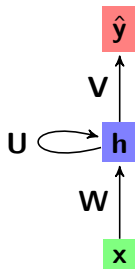


- ▶ The output of the hidden state at one time step is the history/past context input for the next time step!
- ▶ What context information is embedded in \mathbf{h}_{i-1} ?
 - ▶ All previous words
 - ▶ What about previous parts of speech (as in HMMs, CRFs, structured perceptrons)?
 - ▶ To learn more, take StatNLP in the fall!

Recurrent Neural Networks

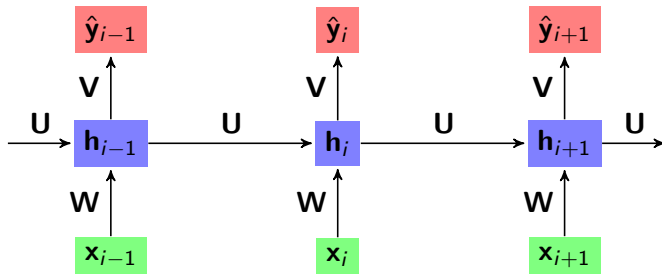


Recurrent Neural Networks

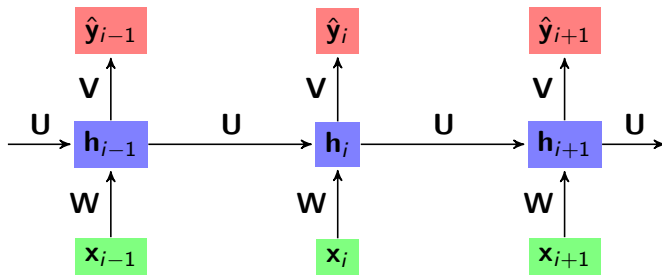


- Neural networks in which the output of a layer in one time step is input to a layer in the next time step

RNN Language Models

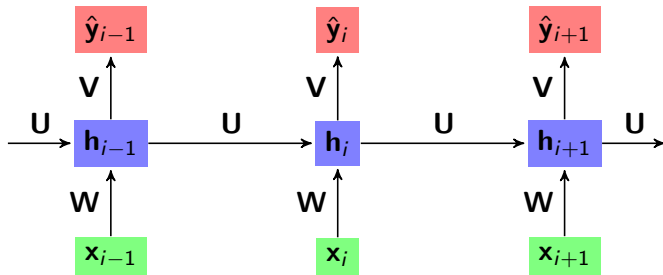


RNN Language Models



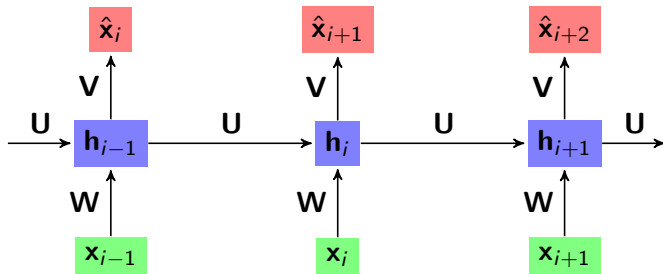
- Sequence labeling: predict current tag given current word, history

RNN Language Models



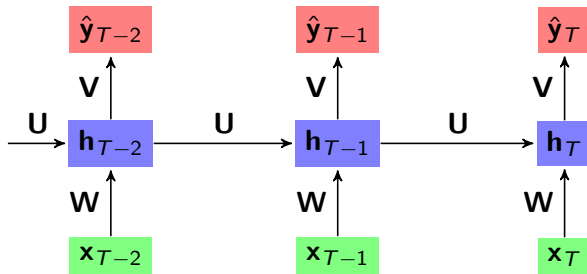
- ▶ Sequence labeling: predict current tag given current word, history
- ▶ Language modeling: predict next word given current word, history

RNN Language Models

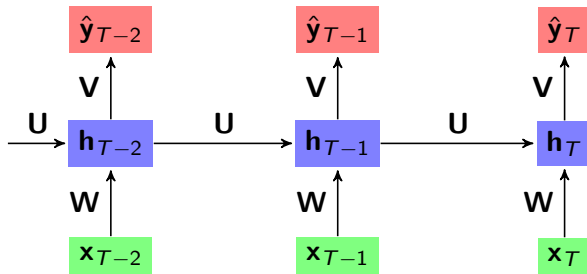


- ▶ Sequence labeling: predict current tag given current word, context
- ▶ Language modeling: predict next word given current word, context

RNNs for Text Classification

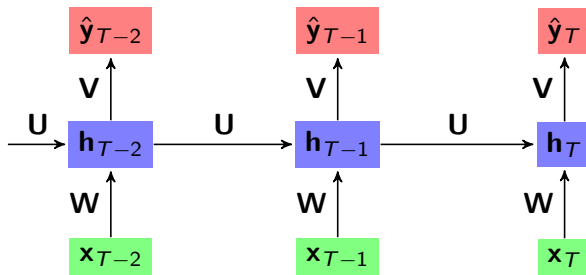


RNNs for Text Classification



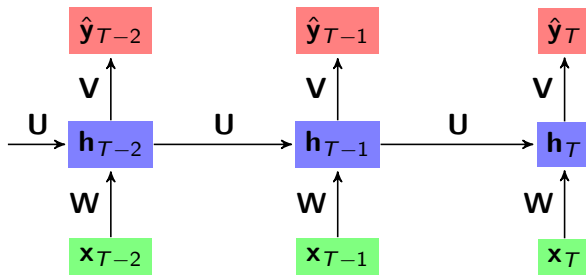
- What context information is embedded in \mathbf{h}_T ?

RNNs for Text Classification



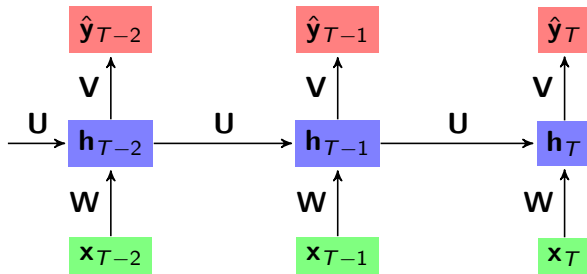
- ▶ What context information is embedded in \mathbf{h}_T ?
 - ▶ Current word \mathbf{x}_T
 - ▶ Context \mathbf{h}_{T-1}

RNNs for Text Classification



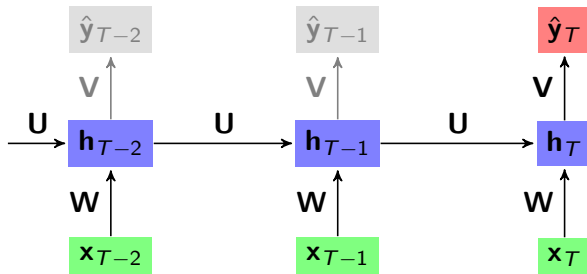
- What context information is embedded in \mathbf{h}_T ?
 - All words (i.e. the whole text)

RNNs for Text Classification



- ▶ What context information is embedded in \mathbf{h}_T ?
 - ▶ All words (i.e. the whole text)
- ▶ Use \mathbf{h}_T to predict class $\hat{\mathbf{y}}_T$ of entire document

RNNs for Text Classification

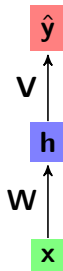


- ▶ What context information is embedded in \mathbf{h}_T ?
 - ▶ All words (i.e. the whole text)
- ▶ Use \mathbf{h}_T to predict class $\hat{\mathbf{y}}_T$ of entire document
 - ▶ Ignore other outputs

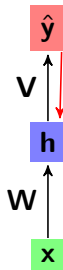
Backpropagation

- ▶ For each matrix of weights \mathbf{W} , starting from the output and working backwards:
 - ▶ Compute gradient $(\nabla L)^{[\mathbf{W}]}$
- ▶ For each matrix of weights \mathbf{W} :
 - ▶ Move in direction of negative gradient

Backpropagation

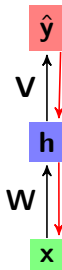


Backpropagation



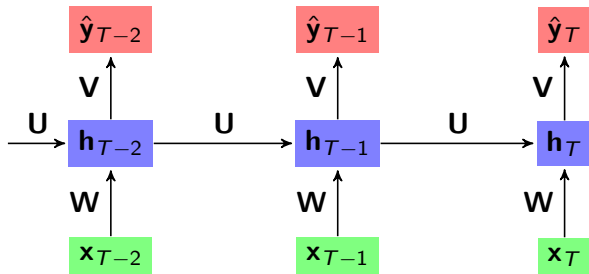
- Compute gradient $(\nabla L)^{[V]}$

Backpropagation

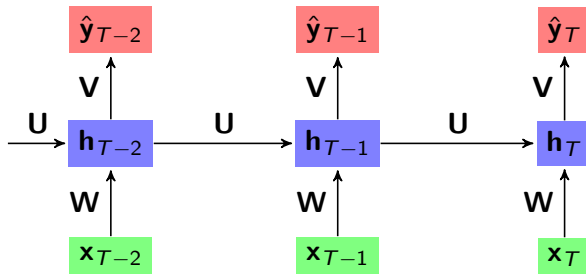


- ▶ Compute gradient $(\nabla L)^{[V]}$
- ▶ Use $(\nabla L)^{[V]}$ to compute gradient $(\nabla L)^{[W]}$

Backpropagation Through Time

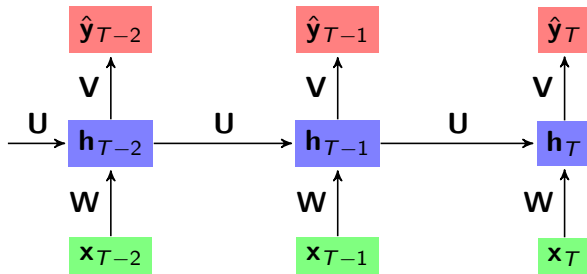


Backpropagation Through Time



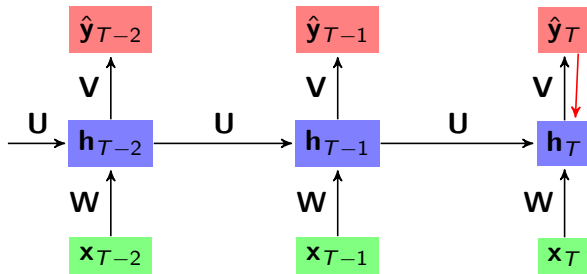
- Start at the end of the text and work backwards

Backpropagation Through Time



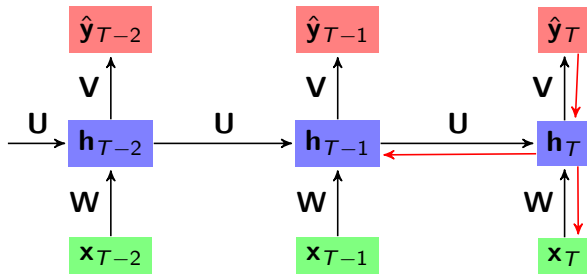
- ▶ Start at the end of the text and work backwards
 - ▶ Let $(\nabla L)_{i,j}^{[W]}$ denote the part of the gradient for weight matrix \mathbf{W} at time i that comes from the output at time j

Backpropagation Through Time



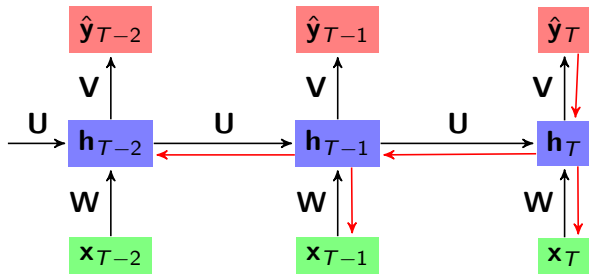
- ▶ Start at the end of the text and work backwards
 - ▶ Compute gradient $(\nabla L)_{T,T}^{[V]}$

Backpropagation Through Time



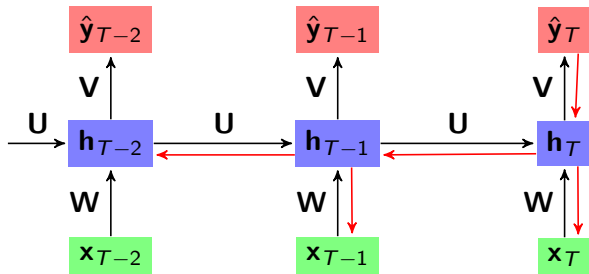
- ▶ Start at the end of the text and work backwards
 - ▶ Compute gradient $(\nabla L)_{T,T}^{[V]}$
 - ▶ Use $(\nabla L)_{T,T}^{[V]}$ to compute gradients $(\nabla L)_{T,T}^{[W]}$ and $(\nabla L)_{T,T}^{[U]}$

Backpropagation Through Time



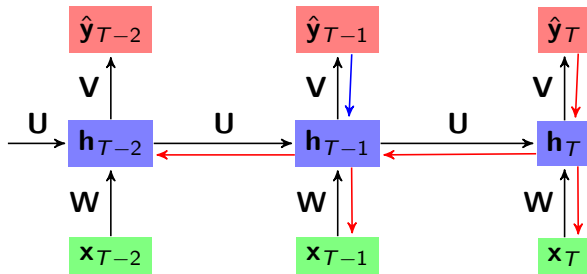
- ▶ Start at the end of the text and work backwards
 - ▶ Compute gradient $(\nabla L)_{T,T}^{[V]}$
 - ▶ Use $(\nabla L)_{T,T}^{[V]}$ to compute gradients $(\nabla L)_{T,T}^{[W]}$ and $(\nabla L)_{T,T}^{[U]}$
 - ▶ Use $(\nabla L)_{T,T}^{[U]}$ to compute gradients $(\nabla L)_{T-1,T}^{[W]}$ and $(\nabla L)_{T-1,T}^{[U]}$

Backpropagation Through Time



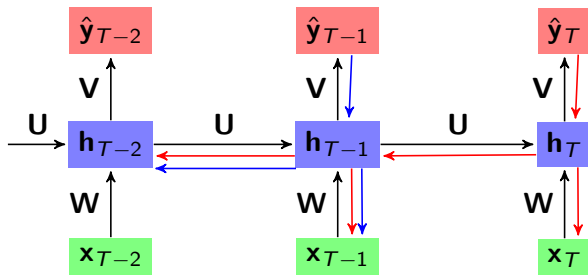
- ▶ Start at the end of the text and work backwards
 - ▶ Compute gradient $(\nabla L)_{T,T}^{[\mathbf{V}]}$
 - ▶ Use $(\nabla L)_{T,T}^{[\mathbf{V}]}$ to compute gradients $(\nabla L)_{T,T}^{[\mathbf{W}]}$ and $(\nabla L)_{T,T}^{[\mathbf{U}]}$
 - ▶ Use $(\nabla L)_{T,T}^{[\mathbf{U}]}$ to compute gradients $(\nabla L)_{T-1,T}^{[\mathbf{W}]}$ and $(\nabla L)_{T-1,T}^{[\mathbf{U}]}$
 - ▶ etc.

Backpropagation Through Time



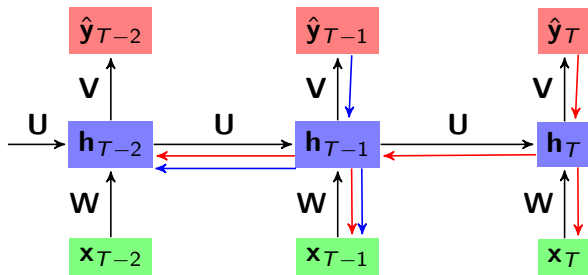
- ▶ Start at the end of the text and work backwards
 - ▶ Compute gradient $(\nabla L)_{T-1, T-1}^{[V]}$

Backpropagation Through Time



- ▶ Start at the end of the text and work backwards
 - ▶ Compute gradient $(\nabla L)_{T-1, T-1}^{[V]}$
 - ▶ Use $(\nabla L)_{T-1, T-1}^{[V]}$ to compute gradients $(\nabla L)_{T-1, T-1}^{[W]}$ and $(\nabla L)_{T-1, T-1}^{[U]}$

Backpropagation Through Time



- ▶ Start at the end of the text and work backwards
 - ▶ Compute gradient $(\nabla L)_{T-1, T-1}^{[V]}$
 - ▶ Use $(\nabla L)_{T-1, T-1}^{[V]}$ to compute gradients $(\nabla L)_{T-1, T-1}^{[W]}$ and $(\nabla L)_{T-1, T-1}^{[U]}$
 - ▶ etc.

Backpropagation Through Time

- ▶ The overall gradient for a weight matrix \mathbf{W} is the sum of the gradients at each time i from each output $\hat{\mathbf{y}}_j$

Backpropagation Through Time

- ▶ The overall gradient for a weight matrix \mathbf{W} is the sum of the gradients at each time i from each output $\hat{\mathbf{y}}_j$

- ▶
$$(\nabla L)^{[\mathbf{W}]} = \sum_{j=1}^T \sum_{i=1}^j (\nabla L)_{i,j}^{[\mathbf{W}]}$$

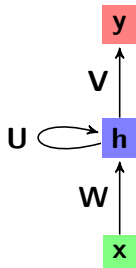
Backpropagation Through Time

- ▶ The overall gradient for a weight matrix \mathbf{W} is the sum of the gradients at each time i from each output $\hat{\mathbf{y}}_j$

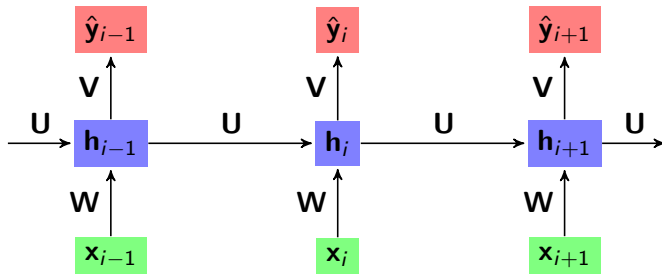
- ▶ $(\nabla L)^{[\mathbf{W}]} = \sum_{j=1}^T \sum_{i=1}^j (\nabla L)_{i,j}^{[\mathbf{W}]}$

- ▶ Then move in direction of negative gradient

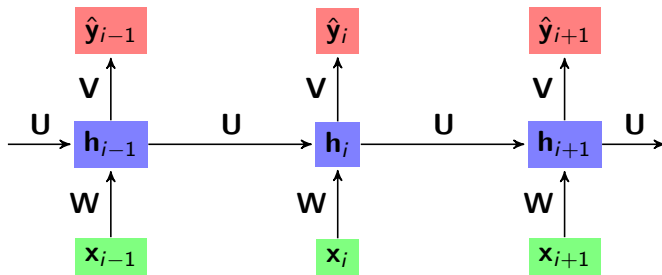
Recurrent Neural Networks



Recurrent Neural Networks

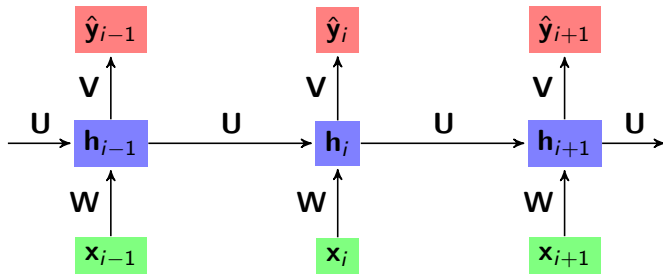


Recurrent Neural Networks



- Output \hat{y}_i depends on hidden state h_i (i.e. current word x_i and history/(past) context h_{i-1})

Recurrent Neural Networks



- ▶ Output \hat{y}_i depends on hidden state h_i (i.e. current word x_i and history/(past) context h_{i-1})
- ▶ What about future context?

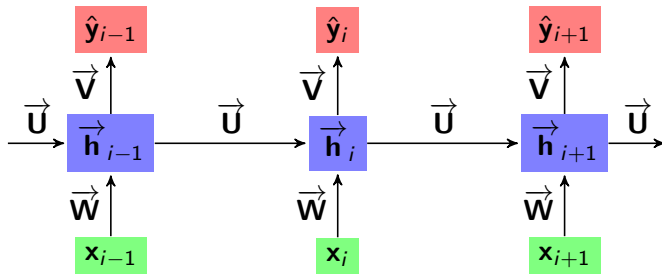
Bidirectional RNNs

- ▶ Idea: Train two RNNs: passing the input into one forward and one backward

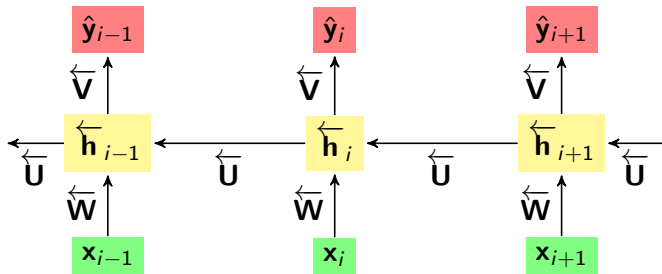
Bidirectional RNNs

- ▶ Idea: Train two RNNs: passing the input into one forward and one **backward**
- ▶ Output $\hat{\mathbf{y}}_i$ depends on forward hidden state $\overrightarrow{\mathbf{h}}_i$ and backward hidden state $\overleftarrow{\mathbf{h}}_i$

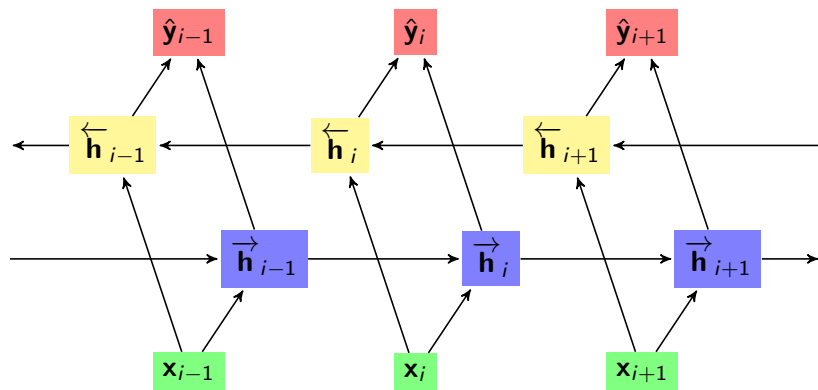
Forward RNN



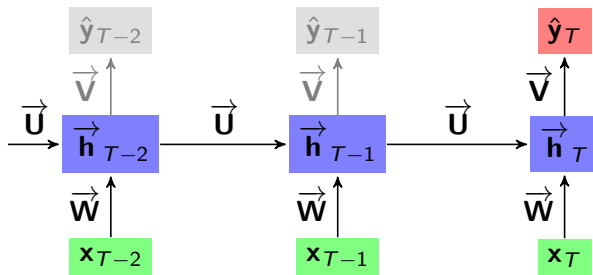
Backward RNN



Bidirectional RNN

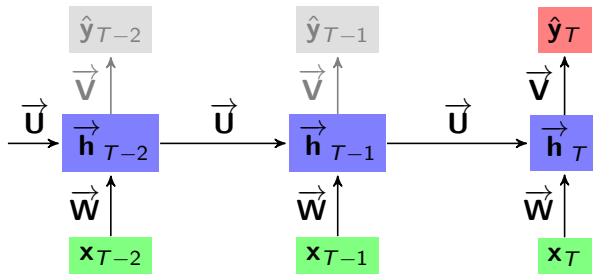


Bidirectional RNNs for Text Classification



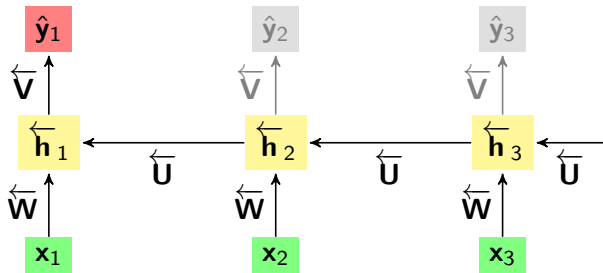
- \vec{h}_T encodes the whole text

Bidirectional RNNs for Text Classification



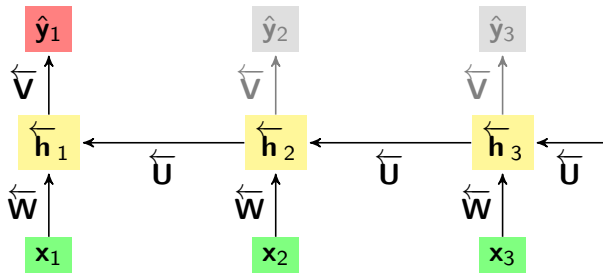
- ▶ \vec{h}_T encodes the whole text
 - ▶ Use \vec{h}_T to predict class \hat{y}_T of entire document

Bidirectional RNNs for Text Classification



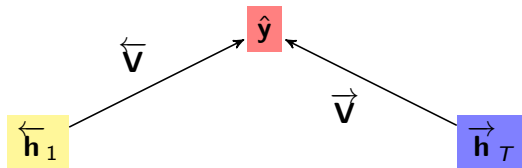
- ▶ \overrightarrow{h}_T encodes the whole text
 - ▶ Use \overrightarrow{h}_T to predict class \hat{y}_T of entire document
- ▶ \overleftarrow{h}_1 also encodes the whole text

Bidirectional RNNs for Text Classification

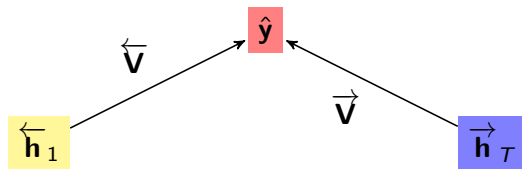


- ▶ \overrightarrow{h}_T encodes the whole text
 - ▶ Use \overrightarrow{h}_T to predict class \hat{y}_T of entire document
- ▶ \overleftarrow{h}_1 also encodes the whole text
 - ▶ Use \overleftarrow{h}_1 to predict class \hat{y}_1 of entire document

Bidirectional RNNs for Text Classification



Bidirectional RNNs for Text Classification



- Use $\overrightarrow{\mathbf{h}}_T$ and $\overleftarrow{\mathbf{h}}_1$ to predict class $\hat{\mathbf{y}}$ of entire document

Context and Long-Distance Dependencies

- ▶ \mathbf{h}_{i-1} encodes the (past, in a forward RNN) context $\mathbf{x}_1, \dots, \mathbf{x}_{i-1}$

Context and Long-Distance Dependencies

- ▶ \mathbf{h}_{i-1} encodes the (past, in a forward RNN) context $\mathbf{x}_1, \dots, \mathbf{x}_{i-1}$
 - ▶ But mostly \mathbf{x}_{i-1} , less \mathbf{x}_{i-2} , even less \mathbf{x}_{i-3}, \dots , very little \mathbf{x}_1

Context and Long-Distance Dependencies

- ▶ \mathbf{h}_{i-1} encodes the (past, in a forward RNN) context $\mathbf{x}_1, \dots, \mathbf{x}_{i-1}$
 - ▶ But mostly \mathbf{x}_{i-1} , less \mathbf{x}_{i-2} , even less \mathbf{x}_{i-3}, \dots , very little \mathbf{x}_1
- ▶ Context is **local**

Context and Long-Distance Dependencies

- ▶ Example: subject-verb agreement

Context and Long-Distance Dependencies

- ▶ Example: subject-verb agreement
- ▶ The flights the airline was cancelling were full.

Context and Long-Distance Dependencies

- ▶ Example: subject-verb agreement
- ▶ The flights the **airline was** cancelling were full.

Context and Long-Distance Dependencies

- ▶ Example: subject-verb agreement
- ▶ The flights the **airline was** cancelling were full.
 - ▶ The context for “**was**” is mostly “**airline**”

Context and Long-Distance Dependencies

- ▶ Example: subject-verb agreement
- ▶ The **flights** the **airline** **was** cancelling **were** full.
 - ▶ The context for “**was**” is mostly “**airline**”

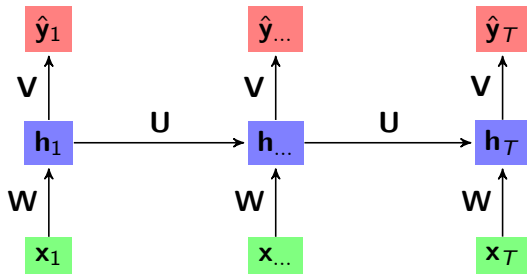
Context and Long-Distance Dependencies

- ▶ Example: subject-verb agreement
- ▶ The **flights** the **airline** **was** cancelling **were** full.
 - ▶ The context for “**was**” is mostly “**airline**”
 - ▶ The context for “**were**” is mostly “cancelling”, “**was**”, “**airline**”

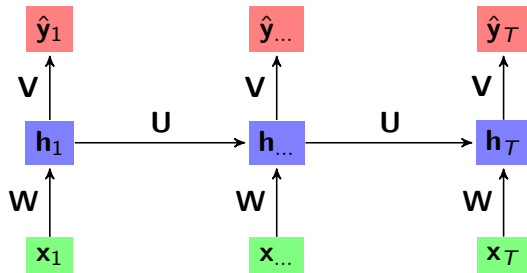
Context and Long-Distance Dependencies

- ▶ Example: subject-verb agreement
- ▶ The **flights** the **airline** **was** cancelling **were** full.
 - ▶ The context for “**was**” is mostly “**airline**”
 - ▶ The context for “**were**” is mostly “cancelling”, “**was**”, “**airline**”
 - ▶ Very little “**flights**”

Vanishing and Exploding Gradients

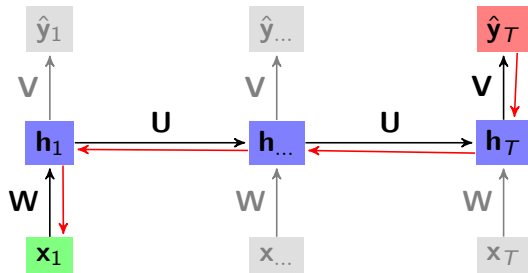


Vanishing and Exploding Gradients



- What is $(\nabla L)_{1,T}^{[W]}$?

Vanishing and Exploding Gradients



- What is $(\nabla L)_{1,T}^{[W]}$?

Vanishing and Exploding Gradients

- ▶ For all layers i :
 - ▶ $(\nabla L)^{[i]} = (\mathbf{a}^{[i-1]})^T \cdot \delta^{[i]}$

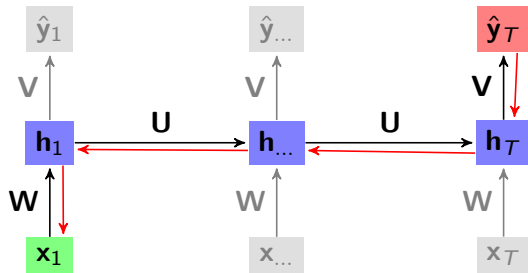
Vanishing and Exploding Gradients

- ▶ For all layers i :
 - ▶ $(\nabla L)^{[i]} = (\mathbf{a}^{[i-1]})^T \cdot \delta^{[i]}$
 - ▶ For simplicity, we will assume that the minibatch size $m = 1$

Vanishing and Exploding Gradients

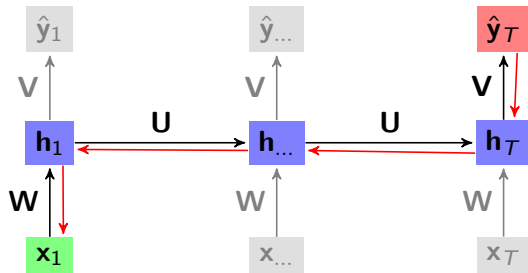
- ▶ For all layers i :
 - ▶ $(\nabla L)^{[i]} = (\mathbf{a}^{[i-1]})^T \cdot \delta^{[i]}$
 - ▶ For simplicity, we will assume that the minibatch size $m = 1$
- ▶ For an output layer \mathcal{L} :
 - ▶ $\delta^{[\mathcal{L}]} = \hat{\mathbf{y}} - \mathbf{y}$
- ▶ For a hidden layer i :
 - ▶ $\delta^{[i]} = (\delta^{[i+1]} \cdot (\mathbf{W}^{[i+1]})^T) \odot g'(\mathbf{z}^{[i]})$

Vanishing and Exploding Gradients



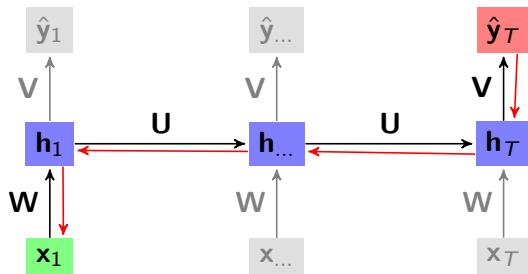
$$(\nabla L)_{1,T}^{[W]} = \mathbf{x}^T \cdot \delta^{[h_1]}$$

Vanishing and Exploding Gradients



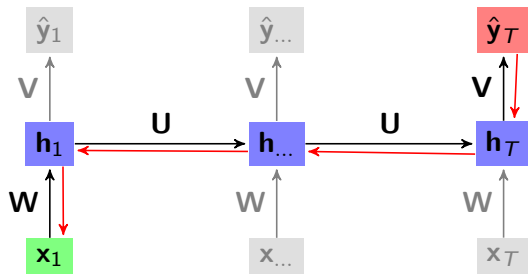
$$\begin{aligned}(\nabla L)_{1,T}^{[W]} &= \mathbf{x}^T \cdot \delta^{[h_1]} \\ &= \mathbf{x}^T \cdot ((\delta^{[h_2]} \cdot \mathbf{U}^T) \odot g'(\mathbf{z}^{[h_1]}))\end{aligned}$$

Vanishing and Exploding Gradients



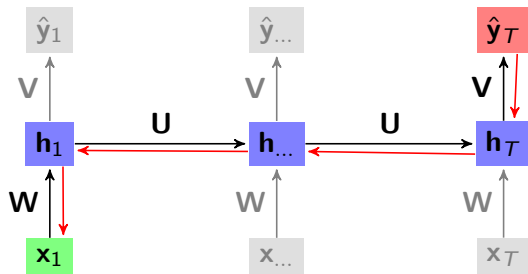
$$\begin{aligned}(\nabla L)_{1,T}^{[W]} &= \mathbf{x}^T \cdot \delta^{[h_1]} \\&= \mathbf{x}^T \cdot ((\delta^{[h_2]} \cdot \mathbf{U}^T) \odot g'(\mathbf{z}^{[h_1]})) \\&= \mathbf{x}^T \cdot (((\delta^{[h_3]} \cdot \mathbf{U}^T) \odot g'(\mathbf{z}^{[h_2]})) \cdot \mathbf{U}^T) \odot g'(\mathbf{z}^{[h_1]}))\end{aligned}$$

Vanishing and Exploding Gradients



$$\begin{aligned}(\nabla L)_{1,T}^{[W]} &= \mathbf{x}^T \cdot \delta^{[h_1]} \\&= \mathbf{x}^T \cdot ((\delta^{[h_2]} \cdot \mathbf{U}^T) \odot g'(\mathbf{z}^{[h_1]})) \\&= \mathbf{x}^T \cdot (((\delta^{[h_3]} \cdot \mathbf{U}^T) \odot g'(\mathbf{z}^{[h_2]})) \cdot \mathbf{U}^T) \odot g'(\mathbf{z}^{[h_1]})) \\&= \dots\end{aligned}$$

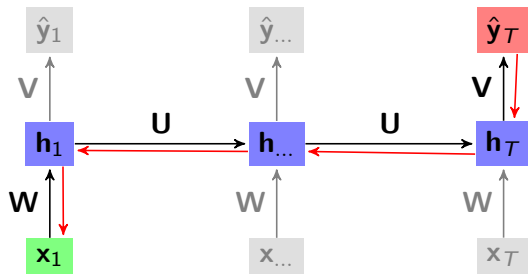
Vanishing and Exploding Gradients



$$\begin{aligned}(\nabla L)_{1,T}^{[W]} &= \mathbf{x}^T \cdot \delta^{[h_1]} \\&= \mathbf{x}^T \cdot ((\delta^{[h_2]} \cdot \mathbf{U}^T) \odot g'(\mathbf{z}^{[h_1]})) \\&= \mathbf{x}^T \cdot (((\delta^{[h_3]} \cdot \mathbf{U}^T) \odot g'(\mathbf{z}^{[h_2]}) \cdot \mathbf{U}^T) \odot g'(\mathbf{z}^{[h_1]})) \\&= \dots\end{aligned}$$

- If weights/derivatives are small, vanishing gradient

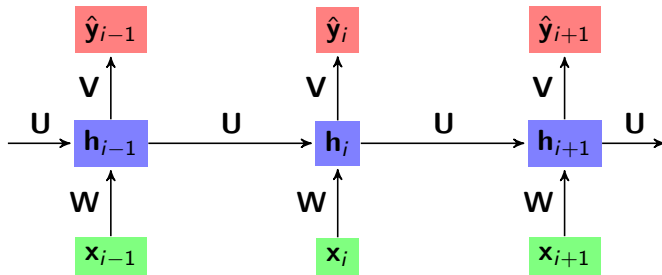
Vanishing and Exploding Gradients



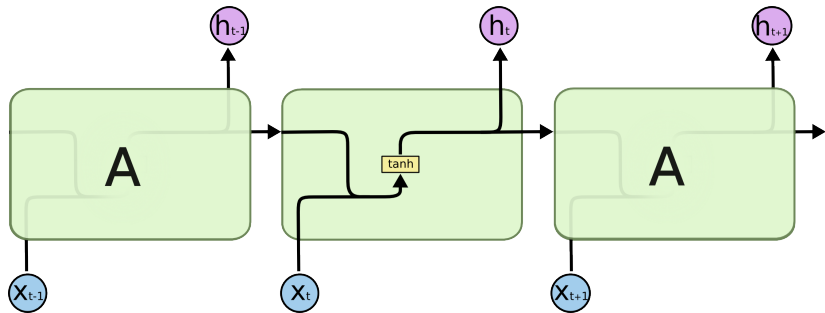
$$\begin{aligned}(\nabla L)_{1,T}^{[W]} &= \mathbf{x}^T \cdot \delta^{[h_1]} \\&= \mathbf{x}^T \cdot ((\delta^{[h_2]} \cdot \mathbf{U}^T) \odot g'(\mathbf{z}^{[h_1]})) \\&= \mathbf{x}^T \cdot (((\delta^{[h_3]} \cdot \mathbf{U}^T) \odot g'(\mathbf{z}^{[h_2]})) \cdot \mathbf{U}^T) \odot g'(\mathbf{z}^{[h_1]})) \\&= \dots\end{aligned}$$

- ▶ If weights/derivatives are small, vanishing gradient
- ▶ If weights/derivatives are large, exploding gradient

Simple RNN

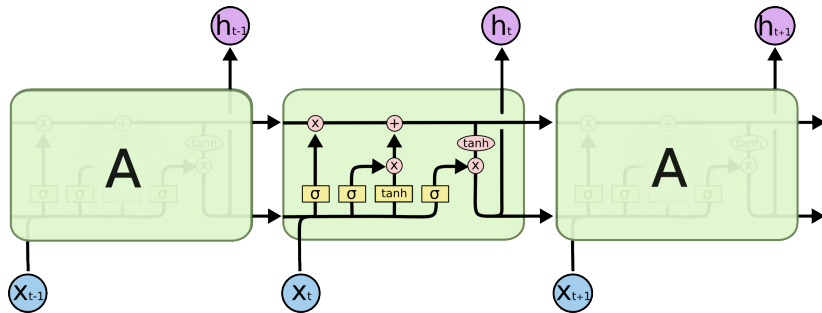


Simple RNN



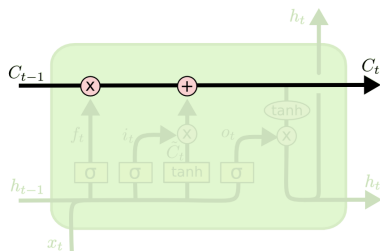
Source

Long Short-Term Memory



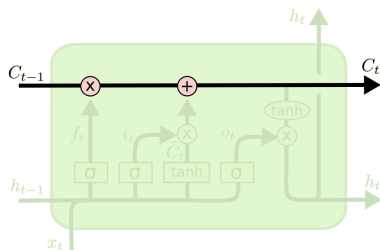
Source

Long Short-Term Memory



Source

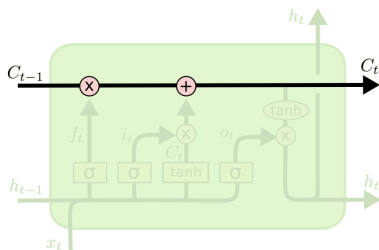
Long Short-Term Memory



Source

- Separate memory (cell) state

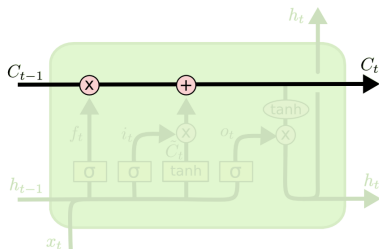
Long Short-Term Memory



Source

- ▶ Separate memory (cell) state
 - ▶ Reading from and writing to memory controlled by **gates**

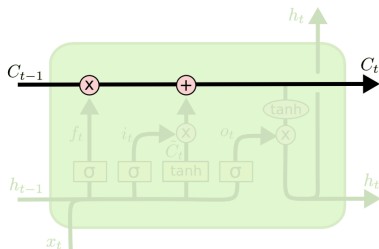
Long Short-Term Memory



Source

- ▶ Separate memory (cell) state
 - ▶ Reading from and writing to memory controlled by **gates**
 - ▶ Each gate contains one or two neural network layers

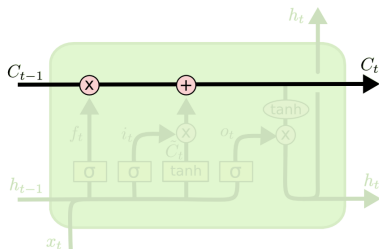
Long Short-Term Memory



Source

- ▶ Separate memory (cell) state
 - ▶ Reading from and writing to memory controlled by **gates**
 - ▶ Each gate contains one or two neural network layers
 - ▶ State **persists** across time

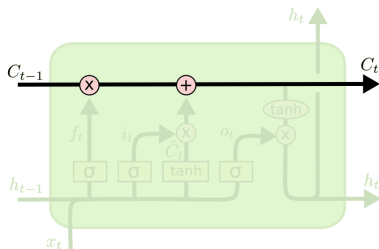
Long Short-Term Memory



Source

- ▶ Separate memory (cell) state
 - ▶ Reading from and writing to memory controlled by **gates**
 - ▶ Each gate contains one or two neural network layers
 - ▶ State **persists** across time
 - ▶ May remember information from long ago

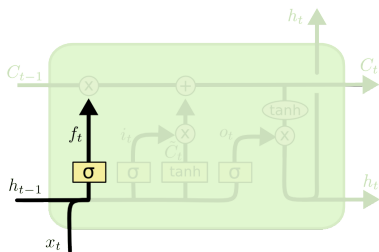
Long Short-Term Memory



Source

- ▶ Separate memory (cell) state
 - ▶ Reading from and writing to memory controlled by **gates**
 - ▶ Each gate contains one or two neural network layers
 - ▶ State **persists** across time
 - ▶ May remember information from long ago
 - ▶ Gradients for memory don't decay with time

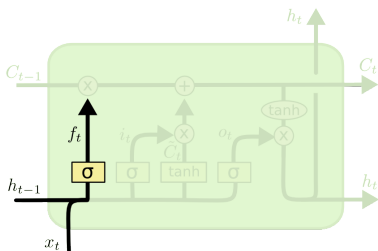
Forget Gate



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

Source

Forget Gate

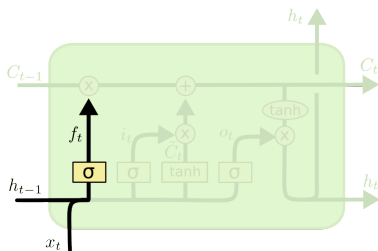


$$f_t = \sigma(W_f \cdot [h_{t-1}, x_t] + b_f)$$

Source

- Neural network layer with logistic activation function

Forget Gate

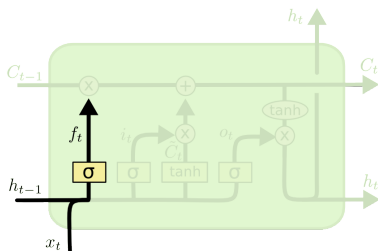


$$f_t = \sigma(W_f \cdot [h_{t-1}, x_t] + b_f)$$

Source

- ▶ Neural network layer with logistic activation function
- ▶ Element-wise multiplication of forget gate output with memory state

Forget Gate

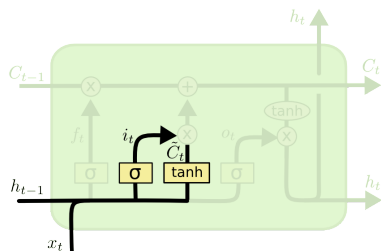


$$f_t = \sigma(W_f \cdot [h_{t-1}, x_t] + b_f)$$

Source

- ▶ Neural network layer with logistic activation function
- ▶ Element-wise multiplication of forget gate output with memory state
 - ▶ **Mask**: What parts of memory to forget/remember?

Input Gate

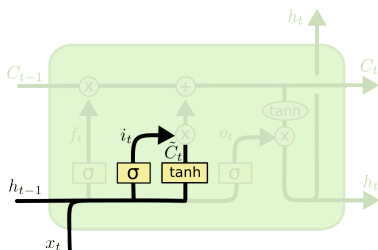


$$i_t = \sigma(W_i \cdot [h_{t-1}, x_t] + b_i)$$

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

Source

Input Gate



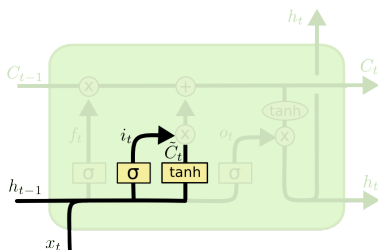
$$i_t = \sigma(W_i \cdot [h_{t-1}, x_t] + b_i)$$

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

Source

- Two parts

Input Gate



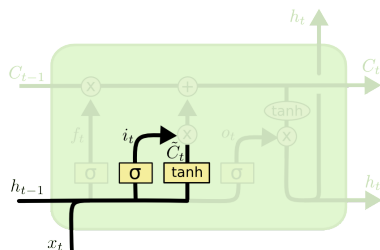
$$i_t = \sigma(W_i \cdot [h_{t-1}, x_t] + b_i)$$

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

Source

- ▶ Two parts
 1. Candidate choice

Input Gate



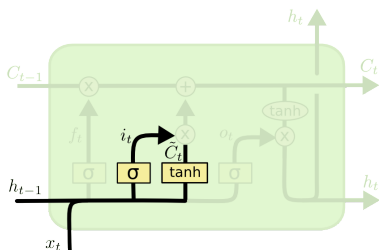
$$i_t = \sigma(W_i \cdot [h_{t-1}, x_t] + b_i)$$

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

Source

- ▶ Two parts
 1. Candidate choice
 - ▶ Logistic activation function

Input Gate



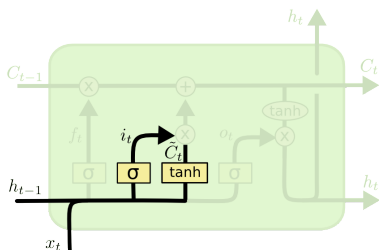
$$i_t = \sigma(W_i \cdot [h_{t-1}, x_t] + b_i)$$

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

Source

- ▶ Two parts
 1. Candidate choice
 - ▶ Logistic activation function
 - ▶ What parts of memory to update?

Input Gate



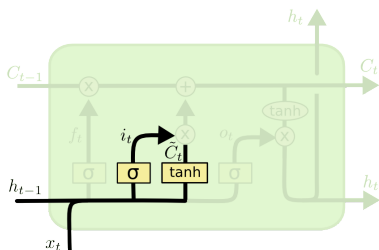
$$i_t = \sigma(W_i \cdot [h_{t-1}, x_t] + b_i)$$

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

Source

- ▶ Two parts
 1. Candidate choice
 - ▶ Logistic activation function
 - ▶ What parts of memory to update?
 2. Candidate values

Input Gate



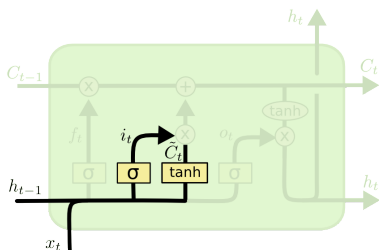
$$i_t = \sigma(W_i \cdot [h_{t-1}, x_t] + b_i)$$

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

Source

- ▶ Two parts
 1. Candidate choice
 - ▶ Logistic activation function
 - ▶ What parts of memory to update?
 2. Candidate values
 - ▶ Tanh activation function

Input Gate

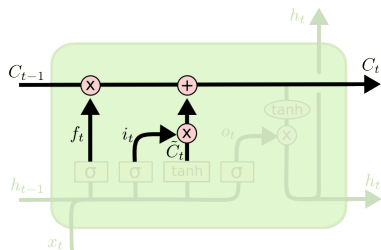


$$i_t = \sigma(W_i \cdot [h_{t-1}, x_t] + b_i)$$
$$\tilde{C}_t = \tanh(W_C \cdot [h_{t-1}, x_t] + b_C)$$

Source

- ▶ Two parts
 1. Candidate choice
 - ▶ Logistic activation function
 - ▶ What parts of memory to update?
 2. Candidate values
 - ▶ Tanh activation function
 - ▶ How much to update them by?

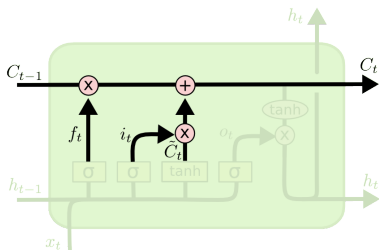
Input Gate



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

Source

Input Gate

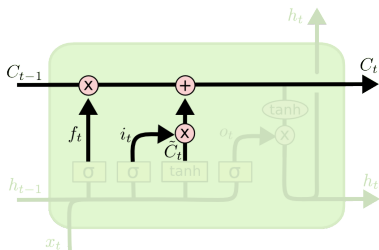


$$C_t = f_t * C_{t-1} + i_t * \tilde{C}_t$$

Source

- Element-wise multiplication of two outputs

Input Gate

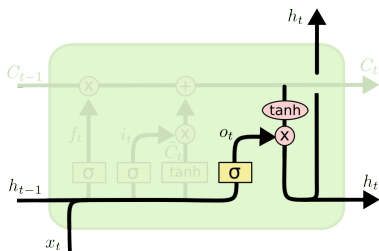


$$C_t = f_t * C_{t-1} + i_t * \tilde{C}_t$$

Source

- ▶ Element-wise multiplication of two outputs
- ▶ Then element-wise addition with memory state

Output Gate

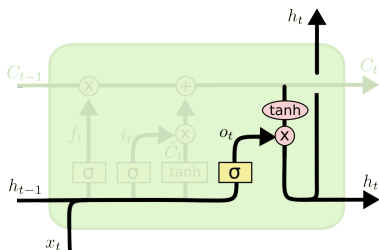


$$o_t = \sigma(W_o [h_{t-1}, x_t] + b_o)$$

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

Source

Output Gate



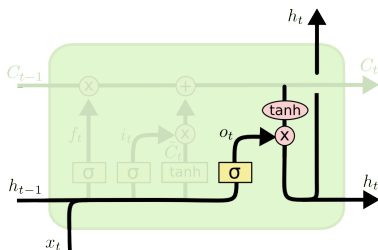
Source

- Logistic activation function

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

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

Output Gate



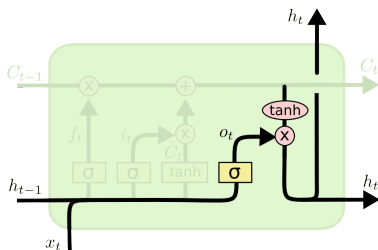
Source

- ▶ Logistic activation function
 - ▶ What parts of memory to output?

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

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

Output Gate



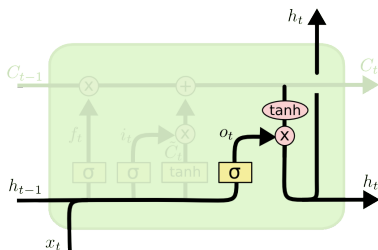
$$o_t = \sigma(W_o [h_{t-1}, x_t] + b_o)$$

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

Source

- ▶ Logistic activation function
 - ▶ What parts of memory to output?
- ▶ Element-wise multiplication with \tanh of memory state

Output Gate



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

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

Source

- ▶ Logistic activation function
 - ▶ What parts of memory to output?
- ▶ Element-wise multiplication with \tanh of memory state
 - ▶ This is the “hidden layer output” that gets passed on to the output layer/next time step