

# CS 6501 Natural Language Processing

## Recurrent Neural Networks

---

Yangfeng Ji

November 12, 2019

Department of Computer Science  
University of Virginia



ENGINEERING

1. Recurrent Neural Networks
2. RNN Language Modeling
3. Challenge of Training RNNs
4. Variants of RNNs

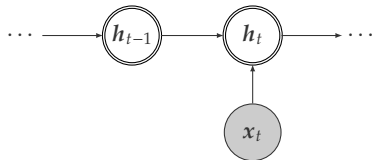
# Recurrent Neural Networks

---

A simple RNN is defined as

$$h_t = f(x_t, h_{t-1}) \quad (1)$$

where  $x_t$  and  $h_t$  is the input and hidden state at time  $t$ <sup>1</sup>



---

<sup>1</sup>Double circles indicate non-random nodes

# Transition Function

For the simplest case,  $f$  is an element-wise sigmoid function as

$$f(x_t, h_{t-1}) = \sigma(\mathbf{W}_h h_{t-1} + \mathbf{W}_i x_t + b) \quad (2)$$

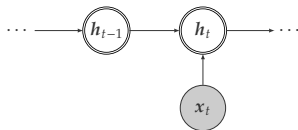
where

- ▶  $\mathbf{W}_h$ : parameter matrix for hidden states
- ▶  $\mathbf{W}_i$ : parameter matrix for inputs
- ▶  $b$ : bias term (also a parameter)

# Unfolding RNNs

Recursive:

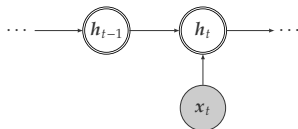
$$h_t = f(x_t, h_{t-1}) \quad (3)$$



# Unfolding RNNs

Recursive:

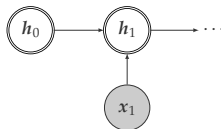
$$h_t = f(x_t, h_{t-1}) \quad (3)$$



Unfolded:

$$\begin{aligned} h_t &= f(x_t, f(x_{t-1}, h_{t-2})) \\ &= f(x_t, f(x_{t-1}, f(x_{t-2}, h_{t-3}))) \\ &= \dots \\ &= f(x_t, f(x_{t-1}, f(x_{t-2}, \dots f(x_1, h_0) \dots))) \end{aligned} \quad (4)$$

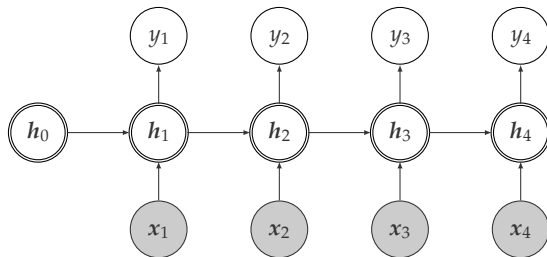
# Base Condition



$$h_t = f(x_t, f(x_{t-1}, f(x_{t-2}, \dots f(\mathbf{x}_1, \mathbf{h}_0) \dots))) \quad (5)$$

- ▶  $h_0$ : zero vector or parameter
- ▶  $x_1$ : input at time  $t = 1$





# For Sequential Modeling

Loss at single time step  $t$

$$L_t(y_t, \hat{y}_t) = \text{cross-entropy}(y_t, \hat{y}_t) \quad (6)$$

where  $y_t$  and  $\hat{y}_t = g(\mathbf{h}_t)$  are the ground truth and predicted output respectively.

The total loss is given as

$$\ell = \sum_{t=1}^T L_t(y_t, \hat{y}_t) \quad (7)$$

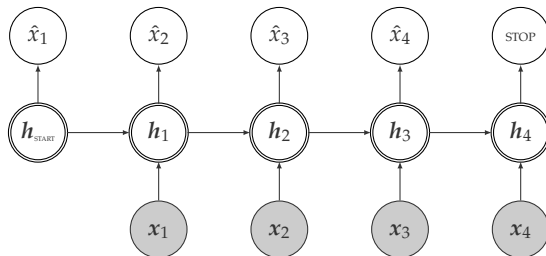
# RNN Language Modeling

---

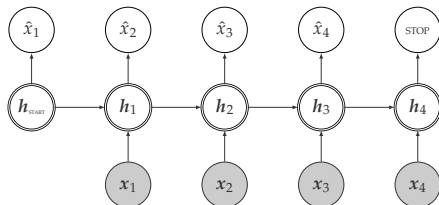
For a given sentence  $\{x_1, \dots, x_T\}$ , the input at time  $t$  is word embedding  $x_t$ . The probability distribution of next word  $X_{t+1}$

$$P(X_{t+1} = x) = \frac{\exp(w_{o,x}^\top h_t)}{\sum_{x'} \exp(w_{o,x'}^\top h_t)} \quad (8)$$

where  $w_{o,x}$  is the output weight vector related to word  $x$ .



# Special Cases



$\{\text{start}, x_1, \dots, x_T, \text{stop}\}$

- ▶ at time  $t = 1$

$$P(X_1 = x) \propto \exp(w_{0,x}^\top h_{\text{start}}) \quad (9)$$

- ▶ at time  $t = T$

$$P(X_T = \text{stop}) \propto \exp(w_{0,x}^\top h_T) \quad (10)$$

$$P(X_{t+1} = x) = \frac{\exp(\mathbf{w}_{o,x}^\top \mathbf{h}_t)}{\sum_{x'} \exp(\mathbf{w}_{o,x'}^\top \mathbf{h}_t)} \quad (11)$$

Options:

- ▶ Negative sampling (x)
- ▶ Class-factored softmax

# Class-factored Softmax: Definition

- ▶ Partition the vocab into  $K$  classes  $\{\mathcal{C}_1, \dots, \mathcal{C}_K\}$ , such that  $\mathcal{V} = \cup \mathcal{C}_k$  and  $\mathcal{C}_k \cap \mathcal{C}_{k'} = \emptyset$  for any  $k' \neq k$

[Baltescu and Blunsom, 2014]



# Class-factored Softmax: Definition

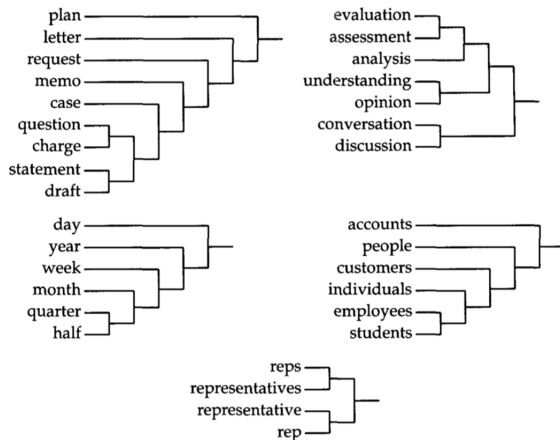
- ▶ Partition the vocab into  $K$  classes  $\{\mathcal{C}_1, \dots, \mathcal{C}_K\}$ , such that  $\mathcal{V} = \cup \mathcal{C}_k$  and  $\mathcal{C}_k \cap \mathcal{C}_{k'} = \emptyset$  for any  $k' \neq k$
- ▶ Define the probability distribution of word as

$$\begin{aligned} P(X_{t+1} = x; \mathbf{h}_t) &= P(X_{t+1} = x, C_{t+1} = c; \mathbf{h}_t) \\ &= P(X_{t+1} = x \mid C_{t+1} = c; \mathbf{h}_t) \\ &\quad \cdot P(C_{t+1} = c \mid \mathbf{h}_t) \end{aligned} \tag{12}$$

[Baltescu and Blunsom, 2014]

# Class-factored Softmax: Word clusters

## Brown clusters



[Brown et al., 1992]

# Computational Complexity

Given

- ▶  $|\mathcal{V}|$  is the vocab size
- ▶  $D$  is the dimension of hidden representations

Model	Training/Decoding
Standard	$\mathcal{O}( \mathcal{V}  \cdot D)$
Class-factored	$\mathcal{O}(\sqrt{ \mathcal{V} } \cdot D)$

Table: Computational complexities of different softmax functions.

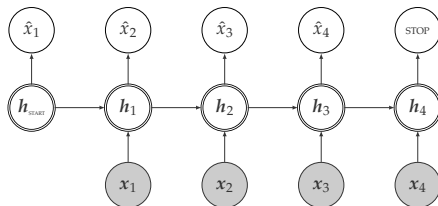
## Challenge of Training RNNs

---

# Backpropagation Through Time

Consider the gradient of  $\ell$  with respect to the network parameters  $\theta = \{\mathbf{W}_h, \mathbf{W}_i, \mathbf{b}\}$ ,

$$\frac{\partial \ell}{\partial \theta} = \sum_{t=1}^T \frac{\partial L_t}{\partial \theta} \quad (13)$$

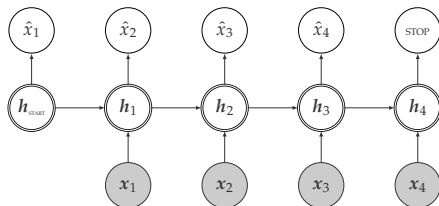


Backpropagation Through Time [Rumelhart et al., 1985, BPTT]

# Gradients

For each time step  $t$ , we have

$$\frac{\partial L_t}{\partial \theta} = \sum_{i=1}^t \left\{ \frac{\partial L_t}{\partial h_t} \cdot \left( \prod_{j=i}^{t-1} \frac{\partial h_{j+1}}{\partial h_j} \right) \cdot \frac{\partial h_i}{\partial \theta} \right\} \quad (14)$$



$$\frac{\partial L_t}{\partial \theta} = \sum_{i=1}^t \left\{ \frac{\partial L_t}{\partial \mathbf{h}_t} \cdot \left( \prod_{j=i}^{t-1} \frac{\partial \mathbf{h}_{j+1}}{\partial \mathbf{h}_j} \right) \cdot \frac{\partial \mathbf{h}_i}{\partial \theta} \right\} \quad (15)$$

- ▶ vanishing gradients
- ▶ exploding gradients

[Pascanu et al., 2013]

Solution: **norm clipping** [Pascanu et al., 2013].

Consider the gradient  $\mathbf{g} = \frac{\partial \ell}{\partial \theta}$ ,

$$\hat{\mathbf{g}} \leftarrow \tau \cdot \frac{\mathbf{g}}{\|\mathbf{g}\|} \quad (16)$$

when  $\|\mathbf{g}\| > \tau$ . Usually,  $\tau = 3$  or 5 in practice.



Solution:

- ▶ initialize parameters carefully
- ▶ replace hidden state transition function  $\sigma(\cdot)$  with other options

$$f(x_t, h_{t-1}) = \sigma(\mathbf{W}_h h_{t-1} + \mathbf{W}_i x_t + \mathbf{b}) \quad (17)$$

- ▶ LSTM [Hochreiter and Schmidhuber, 1997]
- ▶ GRU [Cho et al., 2014]

# Long Short-Term Memory

$$\mathbf{i}_t = \sigma(\mathbf{W}_{xi}\mathbf{x}_t + \mathbf{W}_{hi}\mathbf{h}_{t-1} + \mathbf{W}_{ci}\mathbf{c}_{t-1} + \mathbf{b}_i) \quad (18)$$

$$\mathbf{f}_t = \sigma(\mathbf{W}_{xf}\mathbf{x}_t + \mathbf{W}_{hf}\mathbf{h}_{t-1} + \mathbf{W}_{cf}\mathbf{c}_{t-1} + \mathbf{b}_f) \quad (19)$$

$$\mathbf{c}_t = \mathbf{f}_t \circ \mathbf{c}_{t-1} + \mathbf{i}_t \circ \tanh(\mathbf{W}_{xc}\mathbf{x}_t + \mathbf{W}_{hc}\mathbf{h}_{t-1} + \mathbf{b}_c) \quad (20)$$

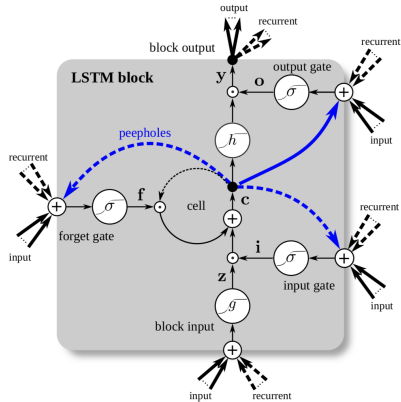
$$\mathbf{o}_t = \sigma(\mathbf{W}_{xo}\mathbf{x}_t + \mathbf{W}_{ho}\mathbf{h}_{t-1} + \mathbf{W}_{co}\mathbf{c}_t + \mathbf{b}_o) \quad (21)$$

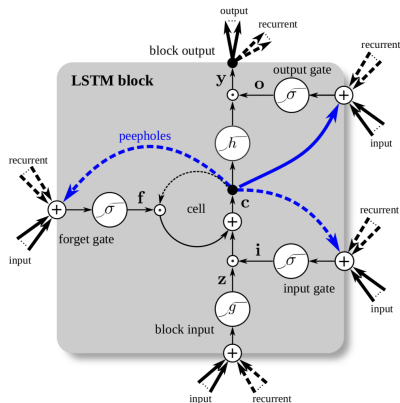
$$\mathbf{h}_t = \mathbf{o}_t \circ \tanh(\mathbf{c}_t) \quad (22)$$

$\circ$  is the element-wise multiplication

[Graves, 2013]

# LSTM





- Forget gate  $f_t$  — discounting on the memory cell
- Peephole connections (connections in blue color) [Gers and Schmidhuber, 2000]

# Gated Recurrent Units

A gated recurrent unit (GRU) was proposed in [Cho et al., 2014].

$$\mathbf{r}_t = \sigma(\mathbf{W}_{rx}\mathbf{x}_t + \mathbf{W}_{rh}\mathbf{h}_{t-1}) \quad (23)$$

$$\tilde{\mathbf{h}}_t = \tanh(\mathbf{W}_{hx}\mathbf{x}_t + \mathbf{W}_{hr}(\mathbf{r}_t \odot \mathbf{h}_{t-1})) \quad (24)$$

$$\mathbf{z}_t = \sigma(\mathbf{W}_{zx}\mathbf{x}_t + \mathbf{W}_{zh}\mathbf{h}_{t-1}) \quad (25)$$

$$\mathbf{h}_t = (1 - \mathbf{z}_t) \odot \mathbf{h}_{t-1} + \mathbf{z}_t \odot \tilde{\mathbf{h}}_t \quad (26)$$

$$(27)$$

Empirical results show GRU units are *comparable* to LSTM units [Chung et al., 2014].

## Variants of RNNs

---

- ▶ Bi-directional RNNs
- ▶ Stacked (*or* Multi-layer) LSTM
- ▶ Memory networks [Weston et al., 2014]

To construct a bi-directional RNN, we need another uni-directional RNN running from the end of the sequence to the beginning, as

$$\mathbf{u}_t = f(\mathbf{x}_t, \mathbf{u}_{t+1}). \quad (28)$$

where  $\mathbf{u}_t$  is the hidden state at time  $t$  in this new model.

[Schuster and Paliwal, 1997]



Use the hidden state  $\mathbf{h}_t^{(k)}$  from the current layer as input  $\mathbf{x}_t^{(k+1)}$  to the next layer [Sutskever et al., 2014],

$$\mathbf{x}_t^{(k+1)} = \mathbf{h}_t^{(k)} . \quad (29)$$

[Sutskever et al., 2014]

1. Recurrent Neural Networks
2. RNN Language Modeling
3. Challenge of Training RNNs
4. Variants of RNNs

# Reference



Baltescu, P. and Blunsom, P. (2014).  
Pragmatic neural language modelling in machine translation.  
*arXiv preprint arXiv:1412.7119*.



Brown, P. F., Desouza, P. V., Mercer, R. L., Pietra, V. J. D., and Lai, J. C. (1992).  
Class-based n-gram models of natural language.  
*Computational linguistics*, 18(4):467–479.



Cho, K., Van Merriënboer, B., Bahdanau, D., and Bengio, Y. (2014).  
On the properties of neural machine translation: Encoder-decoder approaches.  
*arXiv preprint arXiv:1409.1259*.



Chung, J., Gulcehre, C., Cho, K., and Bengio, Y. (2014).  
Empirical evaluation of gated recurrent neural networks on sequence modeling.  
*arXiv preprint arXiv:1412.3555*.



Gers, F. A. and Schmidhuber, J. (2000).  
Recurrent nets that time and count.  
In *Proceedings of the IEEE-INNS-ENNS International Joint Conference on Neural Networks. IJCNN 2000. Neural Computing: New Challenges and Perspectives for the New Millennium*, volume 3, pages 189–194. IEEE.



Graves, A. (2013).  
Generating sequences with recurrent neural networks.  
*arXiv preprint arXiv:1308.0850*.



Hochreiter, S. and Schmidhuber, J. (1997).  
Long short-term memory.  
*Neural computation*, 9(8):1735–1780.



Pascanu, R., Mikolov, T., and Bengio, Y. (2013).  
On the difficulty of training recurrent neural networks.  
In *International Conference on Machine Learning*, pages 1310–1318.



Rumelhart, D. E., Hinton, G. E., and Williams, R. J. (1985).