

CS 6501 Natural Language Processing

Recurrent Neural Networks

Yangfeng Ji

October 29, 2018

Department of Computer Science
University of Virginia



ENGINEERING

Overview

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

Recurrent Neural Networks

A simple RNN is defined as

$$\mathbf{h}_t = f(\mathbf{x}_t, \mathbf{h}_{t-1}) \quad (1)$$

where \mathbf{x}_t and \mathbf{h}_t is the input and hidden state at time t , and \mathbf{b} is the bias. \mathbf{W}_h and \mathbf{W}_i are the weight matrices for hidden states and inputs respectively.

Transition Function

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

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

Unfolding RNNs

Recursive:

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

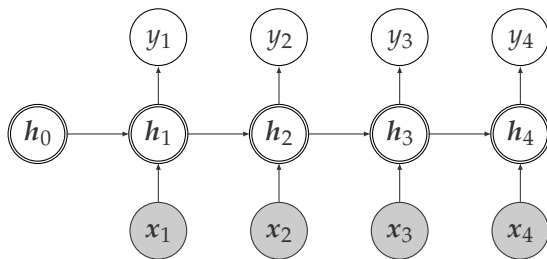
Unfolded:

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

Base Condition

$$h_t = f(x_t, f(x_{t-1}, f(x_{t-2}, \cdots f(x_1, h_0) \cdots))) \quad (5)$$

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



Loss at single time step t

$$L_t(y_t, \hat{y}_t) = \|y_t - \hat{y}_t\|_2^2 \quad (6)$$

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

$$\ell = \sum_{t=1}^T L_t \quad (7)$$

RNN Language Modeling

RNN Language Models

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

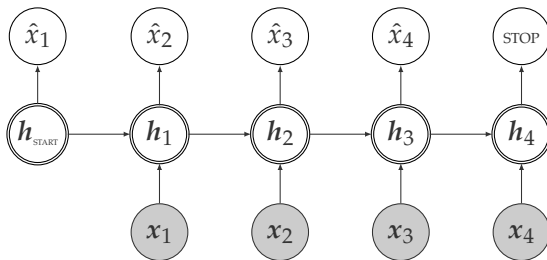
$$\{\text{start}, x_1, \dots, x_T, \text{stop}\}$$

At time $t = 1$

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

At time $t = T$

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



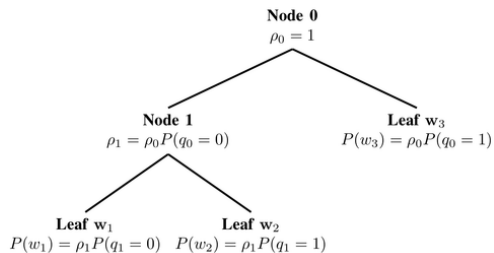
Normalization Term

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

Options:

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

Hierarchical Softmax



Class-factored Softmax: Definition

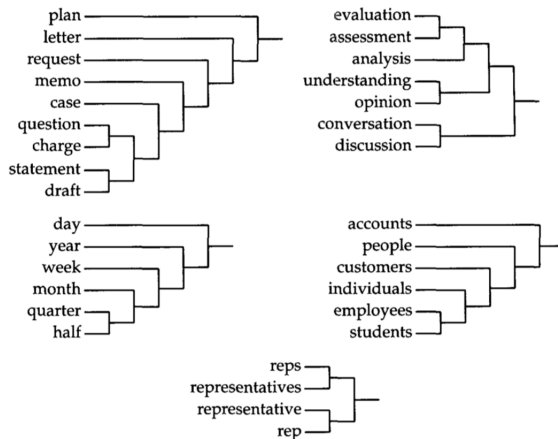
- ▶ Partition the vocab into K classes $\{\mathcal{C}_1, \dots, \mathcal{C}_K\}$, such that $\mathcal{V} = \cup \mathcal{C}_k$ and $\mathcal{C}_k \cap \mathcal{C}_{k'} = \emptyset$ for any $k' \neq k$
- ▶ Define the probability distribution of word as

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

[Baltescu and Blunsom, 2014]

Class-factored Softmax: Word clusters

Brown clusters



[Brown et al., 1992]

Plot

Computational Complexity

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

Table: Computational complexities of different softmax functions.

Challenge of Training RNNs

Backpropagation Through Time

The algorithm used to train RNNs is called *Backpropagation Through Time* [Rumelhart et al., 1985, BPTT].

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

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

For each time step t , we have

$$\frac{\partial L_t}{\partial \theta} = \sum_{i=1}^t \frac{\partial L_t}{\partial h_t} \frac{\partial h_t}{\partial h_i} \frac{\partial h_i}{\partial \theta} \quad (14)$$

where both $\frac{\partial L_t}{\partial h_t}$ and $\frac{\partial h_i}{\partial \theta}$ are the intermediate gradients

Gradients (Cont.)

Computation of $\frac{\partial h_t}{\partial h_i}$ requires the chain rule in calculus,

$$\frac{\partial h_t}{\partial h_i} = \prod_{j=i+1}^t \frac{\partial h_j}{\partial h_{j-1}}. \quad (15)$$

which can be justified by the unfolded version of h_t

$$\begin{aligned}
h_t &= f(x_t, f(x_{t-1}, h_{t-2})) \\
&= f(x_t, f(x_{t-1}, f(x_{t-2}, h_{t-3}))) \\
&= \dots \\
&= f(x_t, f(x_{t-1}, f(x_{t-2}, \dots f(x_1, h_0) \dots))) \quad (16)
\end{aligned}$$

Challenges

$$\frac{\partial \mathbf{h}_t}{\partial \mathbf{h}_i} = \prod_{j=i+1}^t \frac{\partial \mathbf{h}_j}{\partial \mathbf{h}_{j-1}}. \quad (17)$$

- ▶ vanishing gradients
- ▶ exploding gradients

[Pascanu et al., 2013]

Exploding Gradients

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

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

$$\hat{\mathbf{g}} \leftarrow \tau \cdot \frac{\mathbf{g}}{\|\mathbf{g}\|} \quad (18)$$

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

Vanishing Gradients

Solution:

- ▶ initialize parameters carefully
- ▶ replace hidden state function $\sigma()$ with other options
 - ▶ LSTM [Hochreiter and Schmidhuber, 1997]
 - ▶ GRU [Cho et al., 2014]

Long Short-Term Memory

$$\mathbf{i}_t = \sigma(\mathbf{W}_{xi}\mathbf{x}_t + \mathbf{W}_{hi}\mathbf{h}_{t-1} + \mathbf{W}_{ci}\mathbf{c}_{t-1} + \mathbf{b}_i) \quad (19)$$

$$\mathbf{f}_t = \sigma(\mathbf{W}_{xf}\mathbf{x}_t + \mathbf{W}_{hf}\mathbf{h}_{t-1} + \mathbf{W}_{cf}\mathbf{c}_{t-1} + \mathbf{b}_f) \quad (20)$$

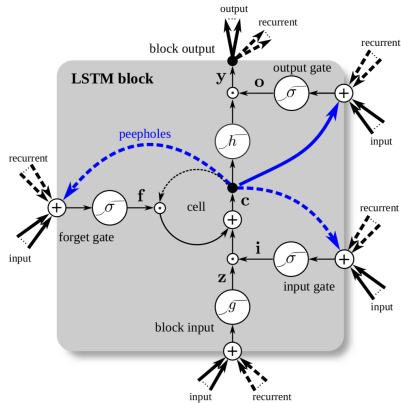
$$\mathbf{c}_t = \mathbf{f}_t \circ \mathbf{c}_{t-1} + \mathbf{i}_t \circ \tanh(\mathbf{W}_{xc}\mathbf{x}_t + \mathbf{W}_{hc}\mathbf{h}_{t-1} + \mathbf{b}_c) \quad (21)$$

$$\mathbf{o}_t = \sigma(\mathbf{W}_{xo}\mathbf{x}_t + \mathbf{W}_{ho}\mathbf{h}_{t-1} + \mathbf{W}_{co}\mathbf{c}_t + \mathbf{b}_o) \quad (22)$$

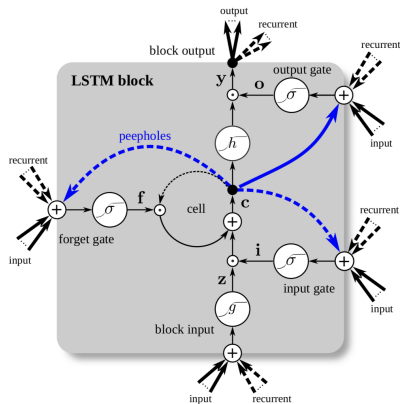
$$\mathbf{h}_t = \mathbf{o}_t \circ \tanh(\mathbf{c}_t) \quad (23)$$

\circ is the element-wise multiplication [Graves, 2013]

LSTM



LSTM



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

A Simple LSTM

A LSTM without forget gate and peephole connections

$$\mathbf{i}_t = \sigma(\mathbf{W}_{xi}\mathbf{x}_t + \mathbf{W}_{hi}\mathbf{h}_{t-1} + \mathbf{b}_i) \quad (24)$$

$$\mathbf{o}_t = \sigma(\mathbf{W}_{xo}\mathbf{x}_t + \mathbf{W}_{ho}\mathbf{h}_{t-1} + \mathbf{b}_o) \quad (25)$$

$$\mathbf{c}_t = \mathbf{c}_{t-1} + \mathbf{i}_t \circ \tanh(\mathbf{W}_{xc}\mathbf{x}_t + \mathbf{W}_{hc}\mathbf{h}_{t-1} + \mathbf{b}_c) \quad (26)$$

$$\mathbf{h}_t = \mathbf{o}_t \circ \tanh(\mathbf{c}_t) \quad (27)$$

[Greff et al., 2017]

Gated Recurrent Units

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

$$\mathbf{r}_t = \sigma(\mathbf{W}_{rx}\mathbf{x}_t + \mathbf{W}_{rh}\mathbf{h}_{t-1}) \quad (28)$$

$$\tilde{\mathbf{h}}_t = \tanh(\mathbf{W}_{hx}\mathbf{x}_t + \mathbf{W}_{hr}(\mathbf{r}_t \odot \mathbf{h}_{t-1})) \quad (29)$$

$$\mathbf{z}_t = \sigma(\mathbf{W}_{zx}\mathbf{x}_t + \mathbf{W}_{zh}\mathbf{h}_{t-1}) \quad (30)$$

$$\mathbf{h}_t = (1 - \mathbf{z}_t) \odot \mathbf{h}_{t-1} + \mathbf{z}_t \odot \tilde{\mathbf{h}}_t \quad (31)$$

$$(32)$$

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

Variants of RNNs

Overview

- ▶ Bi-directional RNNs
- ▶ Stacked (*or* Multi-layer) LSTM
- ▶ Memory networks [Weston et al., 2014]
- ▶ Recurrent neural network grammars [Dyer et al., 2016]

Bi-directional RNNs

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

$$\mathbf{u}_t = f(\mathbf{x}_t, \mathbf{u}_{t+1}). \quad (33)$$

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

[Schuster and Paliwal, 1997]

Stacked LSTM

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

$$x_t^{(k+1)} = h_t^{(k)}. \quad (34)$$

[Sutskever et al., 2014]

Applications

Applications

- ▶ Language modeling
- ▶ POS tagging
- ▶ Named entity recognition
- ▶ Code switch
- ▶ ...

Example

Atlantis touched down at Kennedy Space Center

Example

[Atlantis]_{MSIC} touched down at [Kennedy Space
Center]_{LOC}

Example

[Atlantis]_{MSIC} touched down at [Kennedy Space
Center]_{LOC}

Tag set

- ▶ B: beginning
- ▶ I: inside
- ▶ O: outside

Category

- ▶ Person
- ▶ Location
- ▶ Organization
- ▶ Msic

Example

[Atlantis]_{MSIC} touched down at [Kennedy Space
Center]_{LOC}

Tag set

- ▶ B: beginning
- ▶ I: inside
- ▶ O: outside

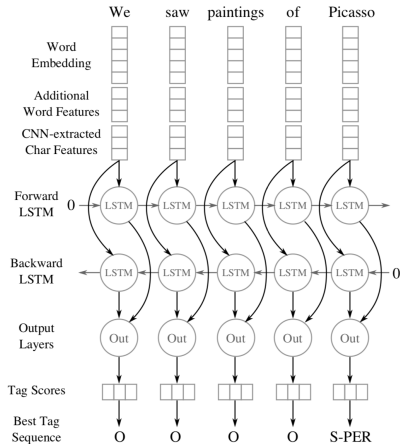
Category

- ▶ Person
- ▶ Location
- ▶ Organization
- ▶ Msic

Atlantis	touched	down	at	Kennedy	Space	Center	.
B _{MSIC}	O	O	O	B _{LOC}	I _{LOC}	I _{LOC}	O

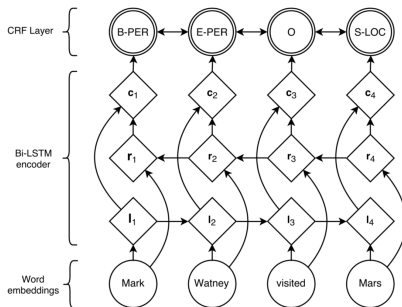
NER (Cont.)

As classification



NER (Cont.)

As sequential modeling (CRFs)



Summary

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

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.



Dyer, C., Kuncoro, A., Ballesteros, M., and Smith, N. A. (2016).
Recurrent neural network grammars.
arXiv preprint arXiv:1602.07776.



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.



Greff, K., Srivastava, R. K., Koutník, J., Steunebrink, B. R., and Schmidhuber, J. (2017).
Lstm: A search space odyssey.