

Modeling Sequences Conditioned on Context with RNNs

Sargur Srihari
srihari@buffalo.edu

10. Topics in Sequence Modeling

- Overview
- 1. Unfolding Computational Graphs
- 2. Recurrent Neural Networks
- 3. Bidirectional RNNs
- 4. Encoder-Decoder Sequence-to-Sequence Architectures
- 5. Deep Recurrent Networks
- 6. Recursive Neural Networks
- 7. The Challenge of Long-Term Dependencies
- 8. Echo-State Networks
- 9. Leaky Units and Other Strategies for Multiple Time Scales
- 10. LSTM and Other Gated RNNs
- 11. Optimization for Long-Term Dependencies
- 12. Explicit Memory

10.2 Topics in Recurrent Neural Networks

0. Overview

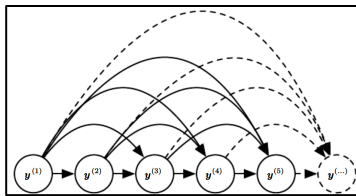
1. Teacher forcing for output-to-hidden RNNs
2. Computing the gradient in a RNN
3. RNNs as Directed Graphical Models
4. Modeling Sequences Conditioned on Context with RNNs

Graphical models of RNNs *without/with* inputs

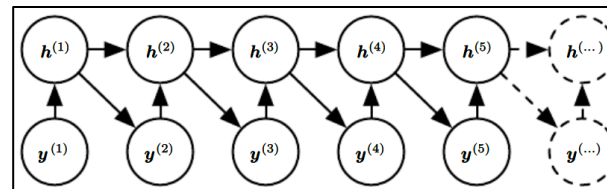
1. Directed graphical models of RNNs *without* inputs

- having a set of random variables $y^{(t)}$:

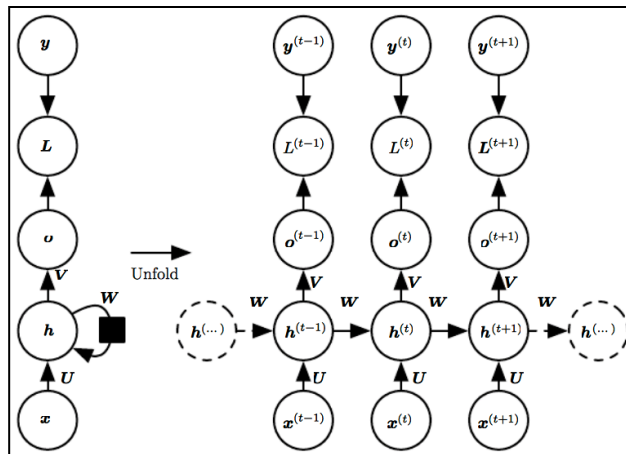
Fully connected graphical model



Efficient parameterization based on $h^{(t)} = f(h^{(t-1)}, x^{(t)}; \theta)$



2. RNNs do include a sequence of inputs $x^{(1)}, x^{(2)}, \dots, x^{(\tau)}$



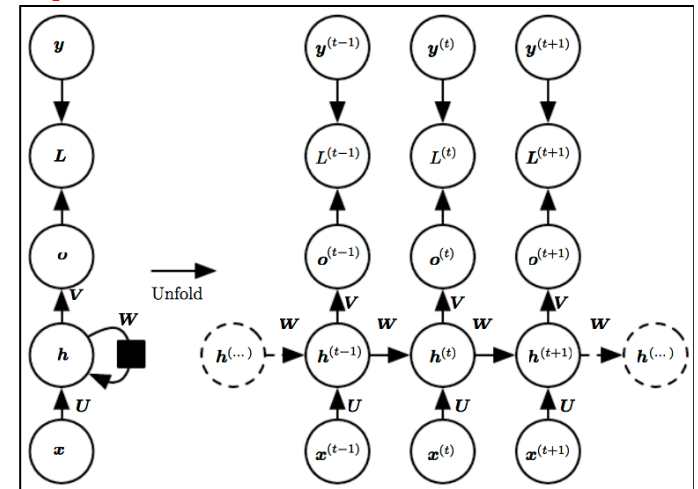
$$\begin{aligned} o^{(t)} &= c + Vh^{(t)} \\ h^{(t)} &= \tanh(a^{(t)}) \\ a^{(t)} &= b + Wh^{(t-1)} + Ux^{(t)} \end{aligned}$$

- RNNs allow the extension of the graphical model view to represent not only the joint distribution view over y variables but also a conditional distribution over y given x

CPDs of model depend on RNN Design pattern

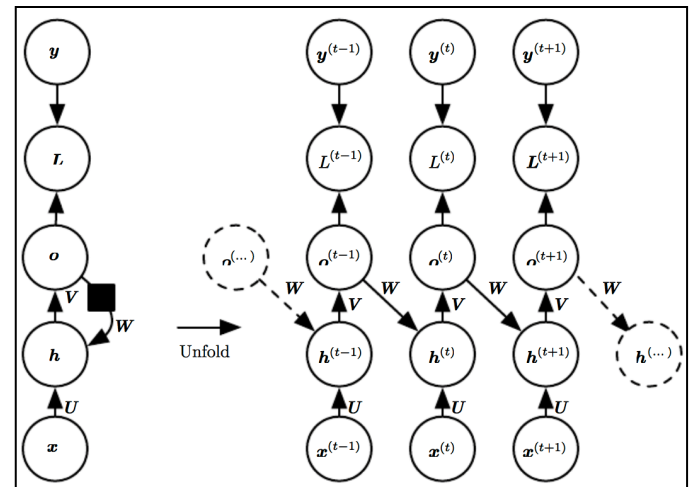
1

1. Recurrent connections between hidden units



2. Recurrent connections only from output at one time step to hidden units at next time step

2



Extending RNNs to represent conditional $P(\mathbf{y}|\mathbf{x})$

- A model representing a variable $P(\mathbf{y}; \boldsymbol{\theta})$ can be reinterpreted as a model representing a conditional distribution $P(\mathbf{y}|\boldsymbol{\omega})$ with $\boldsymbol{\omega}=\boldsymbol{\theta}$
- We can extend such a model to represent a distribution $P(\mathbf{y}|\mathbf{x})$ by using the same $P(\mathbf{y}|\boldsymbol{\omega})$ as before but making $\boldsymbol{\omega}$ a function of \mathbf{x}
- In the case of an RNN this can be achieved in several ways
 - Most common choices are described next

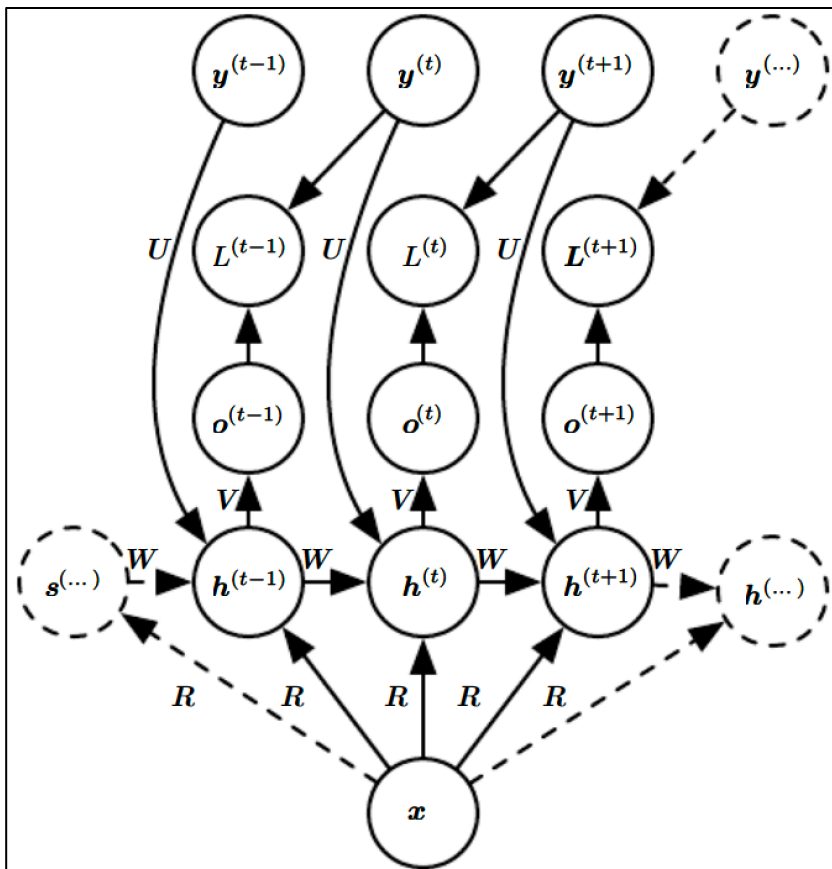
Taking a single vector \mathbf{x} as an extra input

- Instead of taking a sequence $\mathbf{x}^{(t)}$, $t = 1, \dots, \tau$ as input we can take a single vector \mathbf{x} as input
- When \mathbf{x} is a fixed-size vector we can simply make it an extra input of the RNN that generates the \mathbf{y} sequence
- Common ways of providing an extra input to RNN are
 1. An extra input at each time step, or
 2. As the initial state $\mathbf{h}^{(0)}$, or
 3. Both
- The first and common approach is illustrated next

Mapping vector x into distribution over sequences \mathbf{Y}

An extra input at each time step

Interaction between input x and hidden unit vector $h^{(t)}$ is parameterized by a newly introduced weight matrix R that was absent from the model with only y values



Appropriate for tasks such as image captioning

where a single image x is input which produces a sequence of words describing the image

Each element of the observed output $y^{(t)}$ of the observed output sequence serves both as input (for the current time step) and during training as target

RNN to receive a sequence of vectors $\mathbf{x}^{(t)}$ as input

- RNN described by $\mathbf{a}^{(t)} = \mathbf{b} + W\mathbf{h}^{(t-1)} + U\mathbf{x}^{(t)}$ corresponds to a conditional distribution $P(\mathbf{y}^{(1)}, \dots, \mathbf{y}^{(\tau)} | \mathbf{x}^{(1)}, \dots, \mathbf{x}^{(\tau)})$
- It makes a conditional independence assumption that this distribution factorizes as

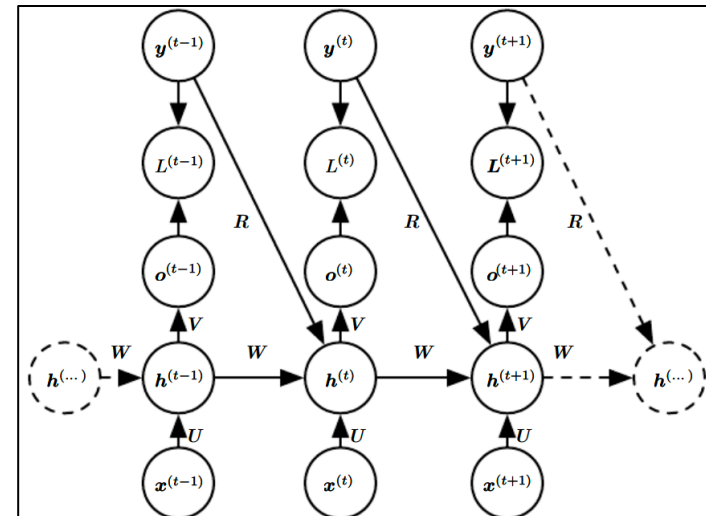
$$\prod_t P(\mathbf{y}^{(t)} | \mathbf{x}^{(1)}, \dots, \mathbf{x}^{(t)})$$

- To remove the conditional independence assumption, we can add connections from the output at time t to the hidden unit at time $t+1$ (see next slide)
 - The model can then represent arbitrary probability distributions over the \mathbf{y} sequence
- Limitation: both sequences must be of same length
 - Removing this restriction is discussed in Section 10.4

Removing conditional independence assumption

Connections from previous output to current state allow RNN to model arbitrary distribution over sequences of y

Conditional RNN mapping a variable-length sequence of x values into a distribution over sequences of y values of same length



Compare it to model that is only able to represent distributions in which the y values are conditionally independent from each other given x values

