

# RNNs as PGMs

Sargur Srihari  
srihari@buffalo.edu

# 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

# 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

## Topics

- RNNs as Directed Graphical Models
- Conditioning Predictions on Past Values
- Fully connected model versus RNN

## Recurrent Networks as Directed PGMs

- With a predictive log-likelihood training objective, such as

$$L\left(\left\{\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(t)}\right\}, \left\{\mathbf{y}^{(1)}, \dots, \mathbf{y}^{(t)}\right\}\right) = \sum_t L^{(t)} = -\sum_t \log p_{\text{model}}\left(\mathbf{y}^{(t)} \mid \left\{\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(t)}\right\}\right)$$

we train the RNN to estimate the conditional distribution of next sequence element  $\mathbf{y}^{(t)}$  given past inputs. This may mean that we maximize log-likelihood

$$\log p(\mathbf{y}^{(t)} \mid \mathbf{x}^{(1)}, \dots, \mathbf{x}^{(t)})$$

or if the model includes connections from output at one time step to the next

$$\log p(\mathbf{y}^{(t)} \mid \mathbf{x}^{(1)}, \dots, \mathbf{x}^{(t)}, \mathbf{y}^{(1)}, \dots, \mathbf{y}^{(t-1)})$$

## Conditioning predictions on past $y$ values

- Decomposing joint probability over sequence of  $y$  values as a series of one-step probabilistic predictions  
is one way to capture the full joint distribution across the whole sequence.
- 1. When we do not feed past  $y$  values as inputs that condition the next step prediction,  
directed PGM contains no edges from any  $y^{(i)}$  in the past to current  $y^{(t)}$ .  
In this case outputs  $y^{(i)}$  are conditionally independent given the  $x$  values
- 2. When we do feed the actual  $y$  values (not their prediction, but the actual observed or generated values) back into the network  
the directed PGM contains edges from all  $y^{(i)}$  values in the past to the current  $y^{(t)}$  value.

## A simple example

- RNN models a sequence of scalar random variables  $Y = \{\mathbf{y}^{(1)}, \dots, \mathbf{y}^{(\tau)}\}$  with no additional inputs  $\mathbf{x}$ .
  - The input at time step  $t$  is simply the output at time step  $t-1$
  - The RNN then defines a directed PGM over the  $\mathbf{y}$  variables
- We parameterize the joint distribution of these observations using the chain rule for conditional probabilities:

$$P(Y) = P(\mathbf{y}^{(1)}, \dots, \mathbf{y}^{(\tau)}) = \prod_{t=1}^{\tau} P(\mathbf{y}^{(t)} \mid \mathbf{y}^{(t-1)}, \mathbf{y}^{(t-2)}, \dots, \mathbf{y}^{(1)})$$

where the right hand side of the bar  $\mid$  is empty for  $t=1$

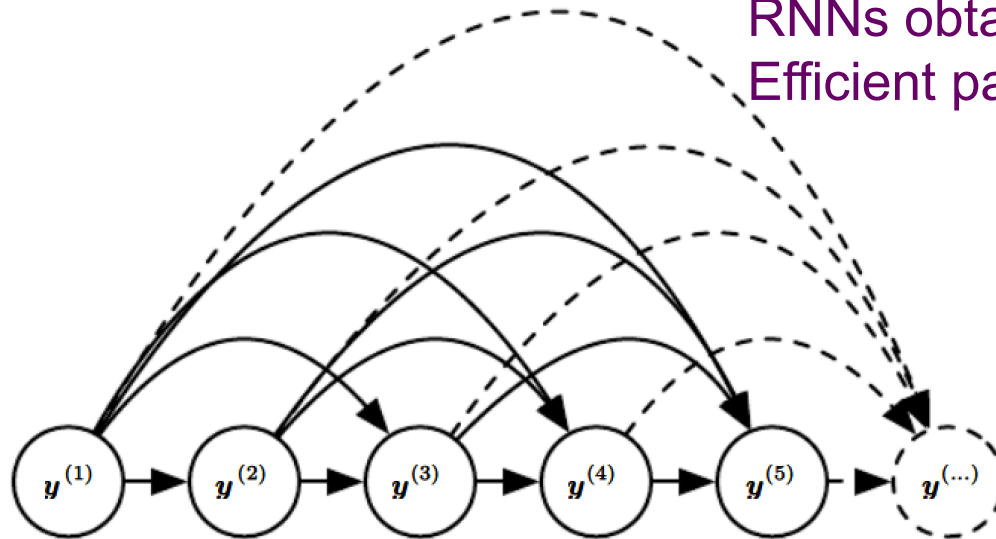
- Hence the negative log-likelihood of a set of values  $\{\mathbf{y}^{(1)}, \dots, \mathbf{y}^{(\tau)}\}$  according to such a model is

$$L = \sum_t L^{(t)}$$

$$\text{where } L^{(t)} = -\sum_t \log P(\mathbf{y}^{(t)} = \mathbf{y}^{(t)} \mid \mathbf{y}^{(t-1)}, \mathbf{y}^{(t-2)}, \dots, \mathbf{y}^{(1)})$$

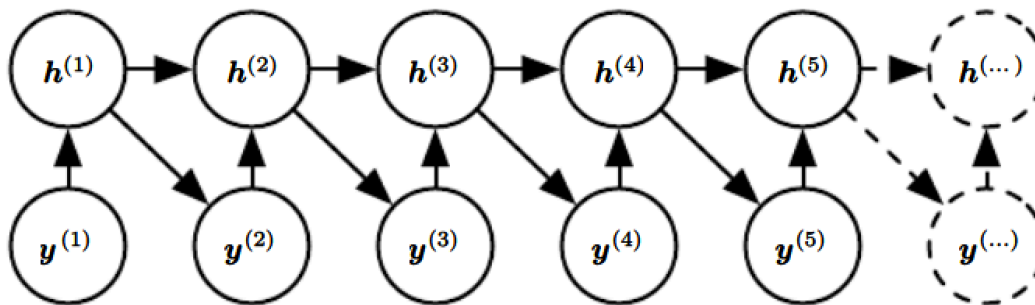
# Fully connected graph model vs RNN

Every one of the past observations  $y^{(i)}$  may influence the conditional distribution of some  $y^{(t)}$ . Parameterizing this way is inefficient. RNNs obtain the same connectivity but Efficient parameterization as shown below



Introducing the state variable in the PGM of RNN, even though it is a deterministic function of its inputs, helps see we get efficient parameterization

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



Every stage in the sequence for  $h^{(t)}$  and  $y^{(t)}$  involves the same structure and can share the parameters with other stages