

# RNNs as Directed Graphical Models

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

10.2.0 Overview

10.2.1 Teacher forcing for output-to-hidden RNNs

10.2.2 Computing the gradient in a RNN

**10.2.3 RNNs as Directed Graphical Models**

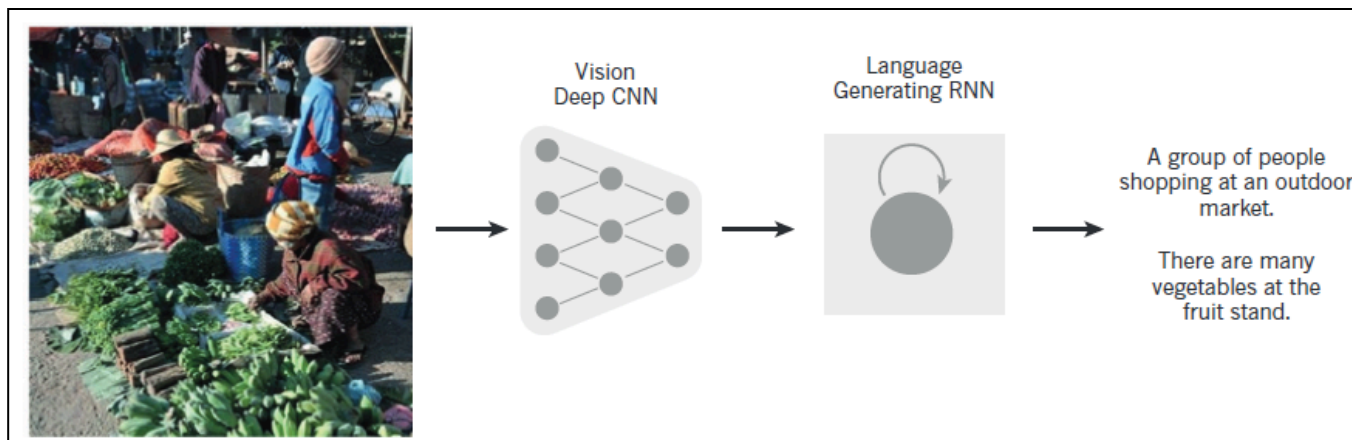
10.2.4 Modeling Sequences Conditioned on Context with RNNs

## 10.2.3 Topics in RNNs as Directed Graphical Models

- Motivation
- What are directed graphical models?
- RNNs as Directed Graphical Models
  - Conditioning Predictions on Past Values
  - Fully connected graphical model for a sequence
  - Hidden units as random variables
  - Efficiency of RNN parameterization
- How to draw samples from the model

# Motivation: From images to text

- Combining ConvNets and Recurrent Net Modules
- Caption generated by a recurrent neural network (RNN) taking as input:
  1. Representation generated by a deep CNN
  2. RNN trained to translate high-level representations of images into captions

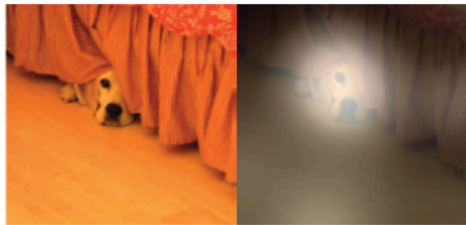


## Better translation of images to captions

- Different focus (lighter patches given more attention)



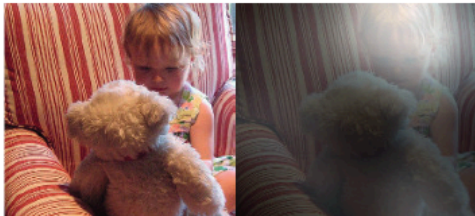
A woman is throwing a **frisbee** in a park.



A **dog** is standing on a hardwood floor.



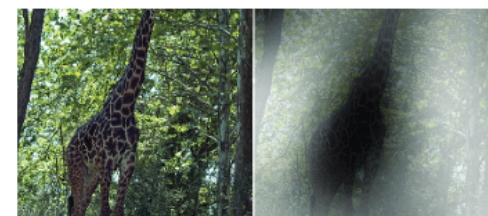
A **stop** sign is on a road with a mountain in the background



A little **girl** sitting on a bed with a teddy bear.



A group of **people** sitting on a boat in the water.



A giraffe standing in a forest with **trees** in the background.

- As it generates each word (**bold**), it exploits it to achieve better translation of images to captions

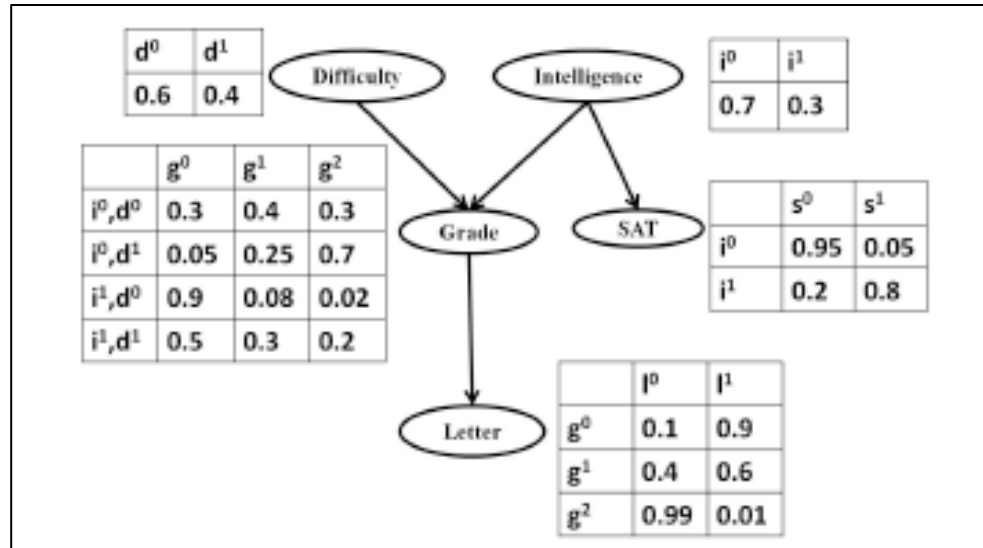
# What is a Directed graphical model?

- Specifies a joint distribution  $p(\mathbf{x}) = p(x_1, \dots, x_n)$  in terms of conditional distributions (called a Bayesian Network)

$$p(\mathbf{x}) = \prod_{i=1}^N p(x_i \mid pa(x_i))$$

- where  $p(x_i \mid pa(x_i))$  are conditional distributions
  - $pa$  represents parents in a directed graph
- Computational efficiency is obtained by omitting edges that do not correspond to strong interactions
- Model is generative:
  - We can easily sample from a BN (ancestral sampling)
- We will see here that RNNs have an interpretation as an efficient directed graphical model

# Example of a Directed Graphical Model (Bayesian Net)



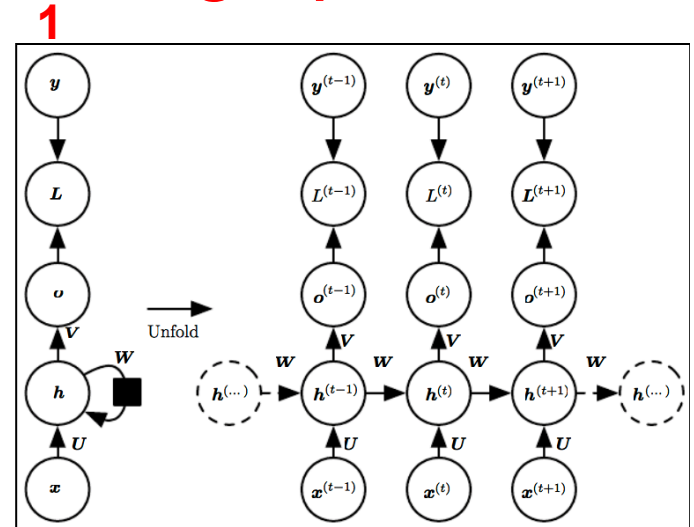
Instead of  $2 \times 2 \times 3 \times 2 \times 2 = 48$  parameters we need 18 parameters using the directed graphical model (which makes certain condition independence assumptions)

$$\begin{aligned}
 P(D, I, G, S, L) &= P(D)P(I)P(G|D, I)P(S|I)P(L|G) \\
 P(i^1, d^0, g^2, s^1, l^0) &= P(i^1)P(d^0)P(g^2 | i^1, d^0)P(s^1 | i^1)P(l^0 | g^2) \\
 &= 0.3 \cdot 0.6 \cdot 0.08 \cdot 0.8 \cdot 0.4 = 0.004608
 \end{aligned}$$

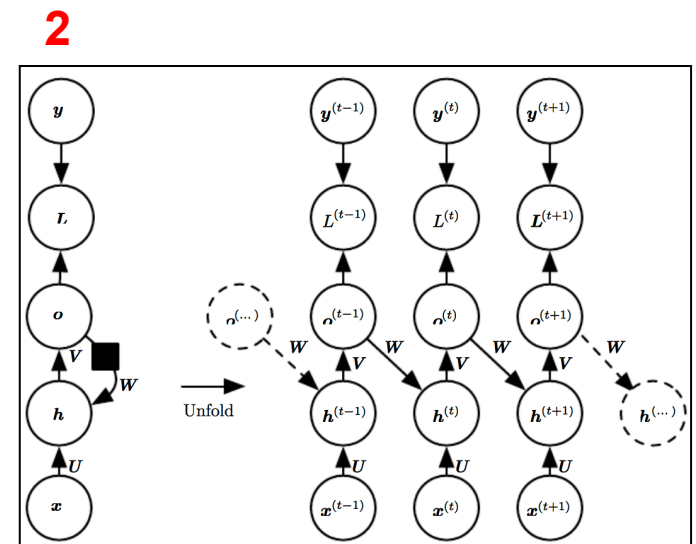


# CPDs of model depend on RNN Design pattern

1. Recurrent connections between hidden units



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



# Loss function for RNNs

- In feed-forward networks our goal is to minimize

$$-E_{\mathbf{x} \sim \hat{p}_{\text{data}}} \left[ \log p_{\text{model}}(\mathbf{x}) \right]$$

- Which is the cross-entropy between distribution of training set (data) and probability distribution defined by model
  - *Definition of cross entropy* between distributions  $p$  and  $q$  is
$$H(p, q) = E_p[-\log q] = H(p) + D_{\text{KL}}(p || q)$$
  - For discrete distributions  $H(p, q) = -\sum_{\mathbf{x}} p(\mathbf{x}) \log q(\mathbf{x})$
- As with a feedforward network, we wish to interpret the output of the RNN as a probability distribution
  - With RNNs the losses  $L^{(t)}$  are cross entropies between training targets  $\mathbf{y}^{(t)}$  and outputs  $\mathbf{o}^{(t)}$ 
    - Mean squared error is the cross-entropy loss associated with an output distribution that is a unit Gaussian

## Types of CPDs in an RNN

### 1. With a predictive log-likelihood training objective:

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

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

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

### 2. Alternatively 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. Do not feed past  $y$  values as inputs for next step prediction,  
directed model 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. 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  $\mathbf{x}$ .
  - Input at time step  $t$  is simply the output at time step  $t-1$
  - RNN defines a directed graphical model over  $\mathbf{y}$  variables
- Parameterize joint distribution of 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  $|$  is empty for  $t=1$

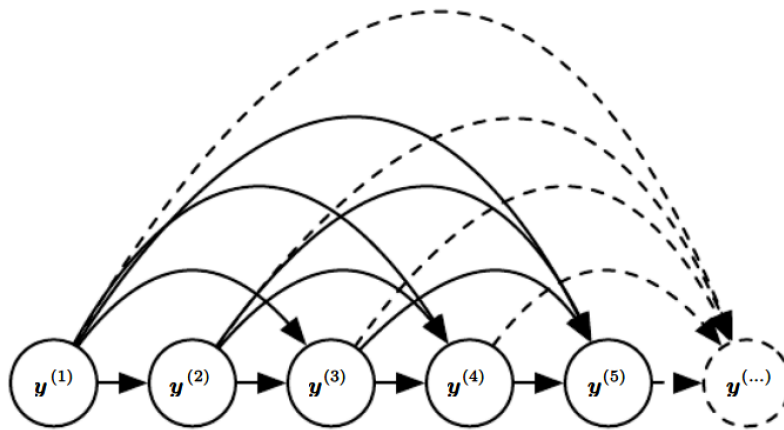
- Hence the negative log-likelihood of a set of values  $\{\mathbf{y}^{(1)}, \dots, \mathbf{y}^{(\tau)}\}$  according to such a model is

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

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

# Fully connected graphical model for a sequence

- Every past observation  $y^{(i)}$  may influence the conditional distribution of some  $y^{(t)}$  for  $t > 1$ 
  - Parameterizing this way is inefficient.



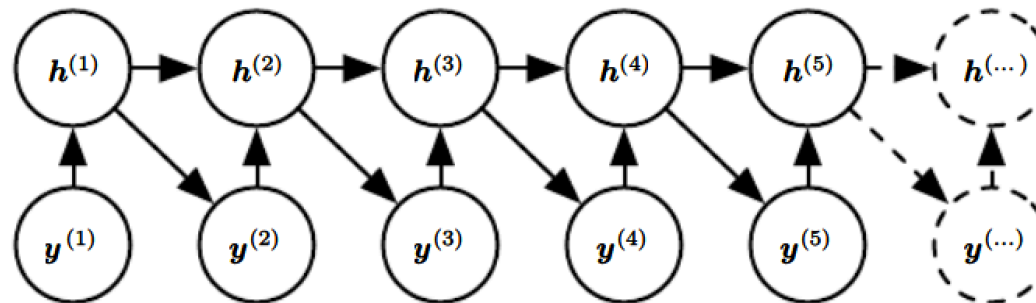
- RNNs obtain the same full connectivity but efficient parameterization as shown next

# Introducing the state variable in the PGM of an RNN

## Introduce 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

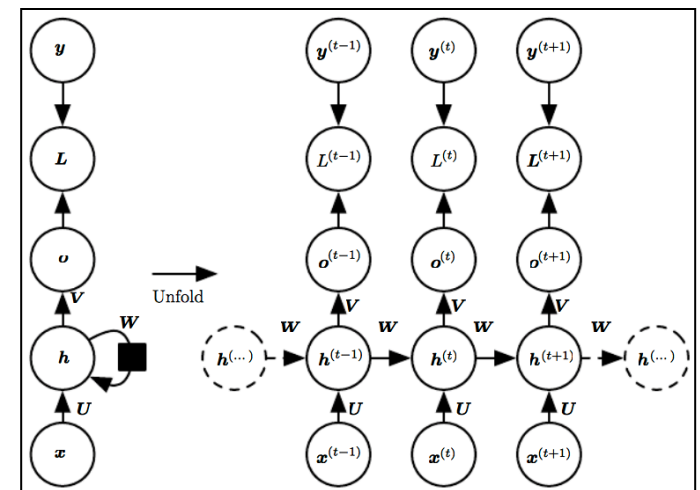
# Efficient parameterization in RNNs

- If  $y$  can take on  $k$  different values, the tabular representation would have  $O(k^\tau)$  parameters
  - By comparison because of parameter sharing the no of parameters in the RNN is  $O(1)$  as a function of sequence length
- Even with efficient parameterization of PGM, some operations are computationally challenging
  - Difficult to predict missing values in the middle of the sequence
- Price RNNs pay for reduced no of parameters is that optimizing parameters may be difficult



# Stationarity in RNNs

- Parameter sharing in RNNs relies on assuming same parameters can be used in different time steps
- Equivalently CPD over variables at time  $t+1$  given the variables at time  $t$  is stationary
- In principle it is possible to use  $t$  as an extra input at each time step
  - Learner discovers time dependence



# How to draw samples from the model

- Simply sample from the conditional distribution at each time step
- One additional complication
  - RNN must have some mechanism for determining the length of the sequence

# Methods for predicting end of sequence

- This can be achieved in various ways
  1. Adding a special symbol for end of sequence
    - When the symbol is generated, the sampling process stops
  2. Introduce an extra Bernoulli output
    - It represents the decision to either continue generation or halt generation at each time step
      - More general than previous method as it can be added to any RNN rather than only RNNs that output a sequence of symbols, e.g., RNN that emits sequence of real numbers
  3. Add an extra output to the model to predict  $\tau$ 
    - Model can sample value of  $\tau$  and then sample  $\tau$  steps worth of data