#### Sampling from Graphical Models in Deep Learning

#### Sargur N. Srihari srihari@cedar.buffalo.edu

#### Topics in Structured PGMs for Deep Learning

- 0. Overview
- 1. Challenge of Unstructured Modeling
- 2. Using graphs to describe model structure
- 3.Sampling from graphical models
- 4. Advantages of structured modeling
- 5.Learning about Dependencies
- 6.Inference and Approximate Inference

7.The deep learning approach to structured probabilistic models

1. Ex: The Restricted Boltzmann machine

Srihari

#### Topics in Sampling from graphical models

- Sampling from directed models (BNs)
- Ancestral Sampling
- Sampling from undirected models (MNs)
- Gibbs Sampling
  - Gibbs sampling from a Gaussian
  - Gibbs sampling from an RBM
  - Sampling Termination

### Sampling from BNs

- Graphical models facilitate drawing samples
   from a model
- One advantage of using a directed graphical model is that a procedure called *ancestral sampling* can produce samples from the joint distribution represented by the model





4

#### **Ancestral Sampling**

- Sort the variables x<sub>i</sub> in the graph into a topological ordering, so that for all i and j, j is greater than i if x<sub>i</sub> is a parent of x<sub>j</sub>
  - i.e., first sample  $x_1 \sim P(x_1)$ , then sample  $P(x_2 | Pa_G(x_2))$ , etc until finally we sample  $P(x_n | Pa_G(x_n))$



To generate one sample:

- 1. Sample x<sub>1</sub>\* from Pr(x<sub>1</sub>)
- 2. Sample  $x_2^*$  from  $Pr(x_2 | x_1^*)$
- 3. Sample  $x_4^*$  from  $Pr(x_4 | x_1^*, x_2^*)$
- 4. Sample  $x_3^*$  from  $Pr(x_3 | x_2^*, x_4^*)$
- 5. Sample  $x_5^*$  from  $Pr(x_5 | x_3^*)$

• If each conditional  $p(x_i | Pa_G(x_i))$  is easy to sample, then whole model is easy to sample

# Role of topological sorting

• Topological sorting guarantees that we can read the conditional distributions in equation

 $p(\mathbf{x}) = \prod_i p(\mathbf{x}_i \mid Pa_{\mathcal{G}}(\mathbf{x}_i)).$ 

and sample from them in order

 Without topological sorting, we might attempt to sample a variable before its parents are available

# Algorithm: Ancestral Sampling

- Start with lowest numbered node
- Draw a sample from the distribution  $p(x_1)$  which we call  $\hat{x}_1$
- Work through each of the nodes in order
  - For node *n* we draw a sample from conditional distribution  $p(x_n | pa_n)$
  - Where parent variables are set to their sampled values
- Once final variable  $x_K$  is sampled
  - Achieved objective of obtaining a single sample from joint distribution
- To sample from marginal distribution
  - Sample from full distribution and discard unnecessary values
  - E.g., to draw from distribution  $p(x_2, x_4)$  simply sample from full distribution, retain values  $\hat{x}_2$ ,  $\hat{x}_4$ , and discard remaining values  $\{\hat{x}_{i\neq 2,4}\}$

# Sampling from Undirected graphs

- Ancestral sampling applies only to BNs
- We can sample from MNs by converting them to BNs
  - But involves solving intractable inference problems
    - To determine marginal for root nodes of directed graph
  - Or introducing so many edges that the resulting directed model becomes intractable
- So drawing samples from an MN is an expensive multi-pass process

# **Gibbs Sampling**

- Simplest approach for sampling from an MN
  - Assume MN over an *n*-dimensional vector of random variables x
  - Visit each variable  $x_i$  and draw a sample conditioned on all other variables, i.e., from  $p(x_i|x_{-i})$
  - Due to separation properties of MNs, we can equivalently condition on only the neighbors of  $x_i$



— . . . . .

— . . . . .

#### Gibbs Sampling with M variables

- Initialize first sample:  $\{z_i, i = 1, ..., M\}$
- For t=1,..,T, T = no of samples
  - Sample  $z_1^{(\tau+1)} \sim p(z_1 | z_2^{(\tau)}, z_3^{(\tau)}, ..., z_M^{(\tau)})$
  - Sample  $z_2^{(\tau+1)} \sim p(z_2|z_1^{(\tau+1)}, z_3^{(\tau)}, ..., z_M^{(\tau)})$
  - Sample  $z_j^{(\tau+1)} \sim p(z_j | z_1^{(\tau+1)}, \dots, z_{j-1}^{(\tau+1)}, z_{j+1}^{(\tau)}, \dots, z_M^{(\tau)})$

- Sample  $z_M^{(\tau+1)} \sim p(z_M | z_1^{(\tau+1)}, z_2^{(\tau+1)}, \dots, z_{M-1}^{(\tau+1)})$ 

•  $p(z_j|z_j)$  is called a *full conditional* for variable *j* 

#### Gibbs sampling from a Gaussian



Deep Learning

# Gibbs sampling from an RBM



Probability of  $(v_i, h_i)$  is

$$p(\mathbf{v}_i, \mathbf{h}_i; \mathbf{W}) = \frac{\exp\left(\mathbf{v}_i \mathbf{W} \mathbf{h}_i\right)}{\sum_{\forall \mathbf{v}, \mathbf{h}} \exp\left(\mathbf{v} \mathbf{W} \mathbf{h}\right)}$$

#### Gibbs sampling



•  $\mathbf{v} \in \{0, 1\}^n$  as the visible states

•  $\mathbf{h} \in \{0,1\}^m$  as the hidden states and with joint distribution

$$p(\mathbf{v}, \mathbf{h}; \mathbf{W}) = \frac{1}{Z} \exp\left(-\mathbf{v}^T \mathbf{W} \mathbf{h}\right)$$

where

•  $\mathbf{W} \in \mathbb{R}^{n \times m}$  are the model parameters

• 
$$Z = \sum_{\mathbf{v},\mathbf{h}\in\{0,1\}^{nm}} \exp\left(-\mathbf{v}^T \mathbf{W} \mathbf{h}\right)$$
 is the partition function,

• and  $E(\mathbf{v}, \mathbf{h}) = \mathbf{v}^T \mathbf{W} \mathbf{h}$  is the energy function.

#### The conditional probability

 $p(\mathbf{v}_i = 1 | \mathbf{h}; \mathbf{W}) = \sigma(-(\mathbf{W}\mathbf{h})_i - \mathbf{W}_{i,0}),$ where  $\sigma(\cdot)$  is the sigmoid function and  $W_{i,0}$  is the bias of visible unit *i*; and the conditional probability  $p(\mathbf{h}_j = 1 | \mathbf{v}; \mathbf{W}) = \sigma(-(\mathbf{v}^T \mathbf{W})_j - \mathbf{W}_{0,j}),$ 

where  $W_{i,0}$  is the bias of hidden unit *j*.

We can use Gibbs sampling to alternate between sampling v and h, assymptotically sampling from P(v, h; W).

Using Gibbs sampling we can sample from  $p(\mathbf{v}, \mathbf{h}; \mathbf{W}^{\star}).$  **Q3:** How can we use Gibbs sampling to sample new images from  $p(\mathbf{v}; \mathbf{W}^{\star}) == \sum_{\forall \mathbf{h}} p(\mathbf{v}, \mathbf{h}; \mathbf{W}^{\star})?$  **A3:** The answer is straightforward. Gibbs sampling gives us a sample from  $(\mathbf{v}_{\infty}, \mathbf{h}_{\infty}) \sim p(\mathbf{v}, \mathbf{h}; \mathbf{W}),$ 

thus the example  $\mathbf{v}_{\infty}$  is sampled from  $p(\mathbf{v};\mathbf{W}).$ 

#### Final states after 1000 Gibbs sampling steps



Srihari

# Gibbs sampling from a BN

 Gibbs Sampling can also be used to sample from any PGM based on the Markov Blanket of the random variables

 $X_{\text{Alice}} \in \{0, 1\}, X_{\text{Bob}} \in \{0, 1\}, X_{\text{Carol}} \in \{0, 1\},$ 



1. 
$$t \leftarrow 0$$
.  
2.  $x_{Alice}^{(0)} = 1, x_{Bob}^{(0)} = 1, x_{Carol}^{(0)} = 1$   
3.  $x_{Bob}^{(t+1)} = p(x_{Bob} | x_{Alice}^{(t)}, x_{Carol}^{(t)})$   
4.  $x_{Alice}^{(t+1)} = p(x_{Alice} | x_{Bob}^{(t+1)})$   
5.  $x_{Carol}^{(t+1)} = p(x_{Carol} | x_{Bob}^{(t+1)})$   
6.  $t \leftarrow t + 1$ ; GOTO 3:

 $(x_{\text{Alice}}^{(\infty)}, x_{\text{Bob}}^{(\infty)}, x_{\text{Carol}}^{(\infty)}) \sim p(X_{\text{Alice}}, X_{\text{Bob}}, X_{\text{Carol}})$ 

#### Gibbs Sampling Termination

- After one pass through the PGM and a sample all *n* variables, we still do not have a fair sample from *p*(**x**)
- Instead we must repeat the process and resample all *n* variables using the updated values of the neighbors
- Asymptotically after many repetitions, process converges to sampling from correct distribution
- Difficult to determine when the samples have reached a sufficiently accurate approximation <sup>14</sup>

#### Another example of sampling convergence

#### An energy-based generative model MCMC sampling algorithm: Metropolis-adjusted Langevin (MALA)

Samples at the beginning of the chain (i.e. simply from the ordinary generator,  $z \sim N(0,I)$ ).



Generated samples after 100 iterations. Chain is smoothly walking on the image manifold and changing semantically meaningful and coherent aspects of images

