## Introduction to Probabilistic Graphical Models

Sargur Srihari srihari@cedar.buffalo.edu

## Topics

- 1. What are graphical models (or PGMs)
- 2. Directionality in graphs
- 3. Bayesian Networks
- 4. Markov Networks

## What are Graphical Models?

- They are *diagrammatic representations* of probability distributions
  - marriage between probability theory and graph theory
- Also called *probabilistic graphical models*
- They augment analysis instead of using pure algebra

# **Graph Directionality**

- Directed graphical models
  - directionality associated with arrows
- Bayesian networks (BNs)
  - Express causal relationships between random variables
- More popular in AI and statistics



- links without arrows
- Markov random fields (MRFs)
  - Better suited to express soft constraints between variables
- More popular in Vision
   and physics



PGMs

# **BN to Probability Distribution**

- Joint distribution defined by a graph is given by a product
- Product terms are conditional distributions of each node conditioned on variables corresponding to parents of that node in the graph

$$p(\mathbf{x}) = \prod_{k=1}^{K} p(x_k \mid pa_k)$$



- where  $x=(x_1,..,x_K)$  and  $pa_k$  denotes set of parents of  $a_k$
- By inspection, joint distribution given by  $p(x_1) \ p(x_2) \ p(x_3)p(x_4|x_1, x_{2}, x_3)p(x_5|x_{1}, x_3)p(x_6|x_4) \ p(x_7|x_4, x_5)$

## **BN Example: Student Network**



- CPDs:  $P(x_i \mid pa(x_i))$
- Joint Distribution:

$$P(\mathbf{x}) = P(x_1, x_2, x_3, x_4, x_5)$$
  

$$P(\mathbf{x}) = \prod_{i=1}^{5} P(x_i \mid pa(x_i))$$

 $P(D,I,G,S,L) = P(D)P(I)P(G \mid D,I)P(S \mid I)P(L \mid G)$ 

## Example of Probabilistic Query



• Posterior Marginal:  $P(I = i^1 | L = l^0, S = s^1) = ?$ 

$$P(I = i^{1}, L = l^{0}, S = s^{1}) = \sum_{D, I = i^{1}, G, L = l^{0}, S = s^{1}} P(D, I, G, S, L)$$

• Probability of Evidence:  $P(L = l^0, S = s^1) = ?$ 

$$P(L = l^{0}, S = s^{1}) = \sum_{D, I, G, L = l^{0}, S = s^{1}} P(D, I, G, S, L)$$

 Here we are asking for a specific probability rather than a full distribution

## **Generative Models**

- Many situations where we need to draw samples from a probability distribution
- Many methods of sampling exist
- Ancestral Sampling is relevant to graphical models
  - Given a graphical model we can specify how samples can be drawn from the joint or marginal distributions
    - Conditional can be drawn by setting the parent value to known value

## **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  $x_2^{\wedge},x_4^{\wedge}$  and discard remaining values

$$\{\hat{x}_{j\neq 2,4}\}$$

## Causal Model Example

- Models can express processes by which the observed data arose
- Object recognition task
  - Observed data point is an image of one of the objects
  - Latent variables are position and orientation of object
- To find posterior distribution over object
  - Integrate over all positions and orientations
- The graphical model captures causal process by which data is generated



#### **Generative Causal Model**

- Since data can be generated from causal model it is a generative model
- Ancestral sampling could be used to generate synthetic data
  - Mimics creation of observed data
  - Gives rise to "fantasy" data whose probability distribution is the same as model
  - Synthetic observation is informative in understanding form of distribution represented by model



Image pixel intensities are dependent on identity of object as well as its pose

## **Undirected Graphical Models**

- When no natural directionality exists
  - Offer a simpler perspective on directed graphs
    - Independence structure
    - Inference task
  - Also called Markov Networks or MRFs
- Can combine directed and undirected graphs
  - Unlike in previous discussion on BNs attention restricted to discrete state spaces

## A simple MRF: Misconception



## MRFs popular in computer vision

Labeling Pixels



$$P(Y \mid X) = \frac{1}{Z(X)} \tilde{P}(Y, X)$$
$$\tilde{P}(Y, X) = \prod_{i=1}^{m} \phi_i(D_i)$$
$$Z(X) = \sum_{Y} \tilde{P}(Y, X)$$

Conditional Random Field for pixel labeling Output distribution y given input x

$$P(\mathbf{y}|\mathbf{x}) = \frac{1}{Z(\mathbf{x})} \exp\{-E(\mathbf{y};\mathbf{x})\}\$$

Energy function over cliques c

$$E(\boldsymbol{y};\boldsymbol{x}) = \sum_{c} \psi_{c}(\boldsymbol{y}_{c};\boldsymbol{x}).$$

Unary and pairwise cliques

$$E(\boldsymbol{y};\boldsymbol{x}) = \sum_{i=1}^{n} \psi_{i}^{U}(\boldsymbol{y}_{i};\boldsymbol{x}) + \sum_{\boldsymbol{y} \in \mathcal{E}} \psi_{\boldsymbol{y}}^{P}(\boldsymbol{y}_{i},\boldsymbol{y}_{j};\boldsymbol{x}).$$

https://cacm.acm.org/magazines/2014/11/179821-scene<sub>14</sub>
 understanding-by-labeling-pixels/abstract

## Gibbs Distributions

- Generalize the idea of factor product to define an undirected parameterization of a distribution
- A distribution  $P_{\Phi}$  is a Gibbs distribution parameterized by a set of factors  $\Phi = \{\phi_1(D_1), \dots, \phi_k(D_k)\}$

$$\Phi = \left\{ \phi_1(D_1), \dots, \phi_K(D_K) \right\}$$

If it is defined as follows

$$P_{\Phi}(X_1,..X_n) = \frac{1}{Z}\tilde{P}(X_1,..X_n)$$

where

$$\tilde{P}(\mathbf{X}_1,\dots,\mathbf{X}_n) = \prod_{i=1}^m \phi_i(D_i)$$

is an unnomalized measure and

$$Z = \sum_{X_1,..X_n} \tilde{P}(X_1,..X_n)$$
 is a normalizing constant called the partition function

The factors do not necessarily represent the marginal distributions of the variables in their scopes

 $D_i$  are sets of random variables

## **Combining Local Models into Global**

- As in BNs parameterization of Markov network defines local interactions
- We combine local models by multiplying them

 $\phi_1(A,B) \cdot \phi_2(B,C) \cdot \phi_3(C,D) \cdot \phi_4(D,A)$ 

*Alice, Bob* and *Debbie* have misconception Charles does not

- $\phi_1(a^1, b^1) \cdot \phi_2(b^1, c^0) \cdot \phi_3(c^0, d^1) \cdot \phi_4(d^1, a^1) = 10.1.100.100 = 100,000$
- Convert to a legal distribution by performing a normalization

	ssig	nme	nt	Unnormalized
$a^0$	$b^0$	$c^0$	$d^0$	300000.
$a^0$	$b^0$	$c^0$	$d^1$	300000
$a^0$	$b^0$	$c^1$	$d^0$	300000
$a^0$	$b^0$	$c^1$	$d^1$	30
$a^0$	$b^1$	$c^0$	$d^0$	500
$a^0$	$b^1$	$c^0$	$d^1$	500
$a^0$	$b^1$	$c^1$	$d^0$	5000000
$a^0$	$b^1$	$c^1$	$d^1$	500
$a^1$	$b^0$	$c^0$	$d^0$	100
$a^1$	$b^0$	$c^0$	$d^1$	1000000
$a^1$	$b^0$	$c^1$	$d^0$	100
$a^1$	$b^0$	$c^1$	$d^1$	100
$a^1$	$b^1$	$c^0$	$d^0$	10
$a^1$	$b^1$	$c^0$	$d^1$	100000
$a^1$	$b^1$	$c^1$	$d^0$	100000
$a^1$	$b^1$	$c^1$	$d^1$	100000
	$\begin{array}{c} A \\ a^{0} \\ a^{0} \\ a^{0} \\ a^{0} \\ a^{0} \\ a^{0} \\ a^{1} \\ a^{$	$\begin{array}{c c} Assig\\ a^0 & b^0\\ a^0 & b^0\\ a^0 & b^0\\ a^0 & b^1\\ a^1 & b^0\\ a^1 & b^0\\ a^1 & b^0\\ a^1 & b^0\\ a^1 & b^1\\ a^1 & $	$\begin{array}{c cccc} Assignme\\ \hline a^0 & b^0 & c^0\\ a^0 & b^0 & c^1\\ a^0 & b^0 & c^1\\ a^0 & b^1 & c^0\\ a^0 & b^1 & c^0\\ a^0 & b^1 & c^1\\ a^0 & b^1 & c^1\\ a^1 & b^0 & c^0\\ a^1 & b^0 & c^0\\ a^1 & b^0 & c^1\\ a^1 & b^1 & c^0\\ a^1 & b^1 & c^0\\ a^1 & b^1 & c^1\\ a^1 & b^1 & c^1\\$	$\begin{array}{c ccccc} Assignment \\ \hline a^0 & b^0 & c^0 & d^0 \\ a^0 & b^0 & c^1 & d^0 \\ a^0 & b^0 & c^1 & d^1 \\ a^0 & b^1 & c^0 & d^0 \\ a^0 & b^1 & c^0 & d^1 \\ a^0 & b^1 & c^1 & d^0 \\ a^0 & b^1 & c^1 & d^1 \\ a^1 & b^0 & c^0 & d^0 \\ a^1 & b^0 & c^0 & d^1 \\ a^1 & b^0 & c^1 & d^1 \\ a^1 & b^1 & c^0 & d^1 \\ a^1 & b^1 & c^0 & d^1 \\ a^1 & b^1 & c^0 & d^1 \\ a^1 & b^1 & c^1 & d^0 \\ a^1 & b^1 & c^1 & d^1 \end{array}$

PGMs

## **Normalized Joint Distribution**

$$P(a,b,c,d) = \frac{1}{Z}\phi_1(a,b) \cdot \phi_2(b,c) \cdot \phi_3(c,d) \cdot \phi_4(d,a)$$

where

$$Z = \sum_{a,b,c,d} \phi_1(a,b) \cdot \phi_2(b,c) \cdot \phi_3(c,d) \cdot \phi_4(d,a)$$

*Z* is a normalizing constant called the Partition function

"Partition" originates from Markov Random Fields in statistical physics

Energy of a physical system of interacting atoms

"Function" is because Z is a function of the parameters

	4ssig	nme	nt	Unnormalized	Normalized
$a^0$	$b^0$	$c^0$	$d^0$	300000	0.04
$a^0$	$b^0$	$c^0$	$d^1$	300000	0.04
$a^0$	$b^0$	$c^1$	$d^0$	300000	0.04
$a^0$	$b^0$	$c^1$	$d^1$	30	$4.1 \cdot 10^{-6}$
$a^0$	$b^1$	$c^0$	$d^0$	500	$6.9 \cdot 10^{-5}$
$a^0$	$b^1$	$c^0$	$d^1$	500	$6.9 \cdot 10^{-5}$
$a^0$	$b^1$	$c^1$	$d^0$	5000000	0.69
$a^0$	$b^1$	$c^1$	$d^1$	500	$6.9 \cdot 10^{-5}$
$a^1$	$b^0$	$c^0$	$d^0$	100	$1.4 \cdot 10^{-5}$
$a^1$	$b^0$	$c^0$	$d^1$	1000000	0.14
$a^1$	$b^0$	$c^1$	$d^0$	100	$1.4 \cdot 10^{-5}$
$a^1$	$b^0$	$c^1$	$d^1$	100	$1.4 \cdot 10^{-5}$
$a^1$	$b^1$	$c^0$	$d^0$	10	$1.4 \cdot 10^{-6}$
$a^1$	$b^1$	$c^0$	$d^1$	100000	0.014
$a^1$	$b^1$	$c^1$	$d^0$	100000	0.014
$a^1$	$b^1$	$c^1$	$d^1$	100000	0.014

Z=7,201,840

## Log-linear Parameterization

- To convert factors in Gibbs parameterization to log-linear form:
  - Take negative natural logarithm of each potential
    - Requires potential to be positive (to take logarithm)

0

#### Features in a Markov Network

- If *D* is a subset of variables, feature *f*(*D*) is a function from *D* to *R* (a real value)
- Feature is a factor *without* a non-negativity requirement
- Given a set of k features  $\left[P(X_1,...,X_n) = \frac{1}{Z} \exp\left[-\sum_{i=1}^k w_i f_i(D_i)\right]\right] \left\{ f_1(D_1), \dots, f_k(D_k) \right\}$ 
  - where  $w_i f_i(D_i)$  is entry in energy function table, since

$$P(X_1,..,X_n) = \frac{1}{Z} \exp\left[-\sum_{i=1}^m \varepsilon_i(D_i)\right]$$

- Can have several functions over same scope, *k.ne.m* 
  - So can represent a standard set of table potentials

## **Example of Feature**

- Pairwise Markov Network  $\overleftarrow{A}$ —B— $\overleftarrow{C}$ 
  - Variables are binary
  - Three clusters:  $C_1 = \{A, B\}, C_2 = \{B, C\}, C_3 = \{C, A\}$
  - Log-linear model with features
    - $f_{00}(x,y)=1$  if x=0, y=0; 0 otherwise for x,y instance of  $C_i$
    - $f_{11}(x,y)=1$  if x=1, y=1 and 0 otherwise
  - Three data instances (A, B, C): (0, 0.0), (0, 1, 0), (1, 0, 0)
    - Unnormalized Feature counts are

 $E_{\tilde{P}}[f_{00}] = (3+1+1)/3 = 5/3$  $E_{\tilde{P}}[f_{11}] = (0+0+0)/3 = 0$ 

#### Definition of Log-linear model with features

- A distribution P is a log-linear model over H if
  - A set of k features  $F = \{f_1(D_1), .., f_k(D_k)\}$  where each  $D_i$  is a complete subgraph and a set of weights  $w_i$
- Such that

$$P(X_1,..,X_n) = \frac{1}{Z} \exp\left[-\sum_{i=1}^k w_i f_i(D_i)\right]$$

Note that k is the no of features, not no of subgraphs

## **Indicator Feature**

- A type of feature of particular interest
- Takes on value 1 for some values  $y \in Val(D)$ and 0 for others
- Example:
  - $D = \{A_1, A_2\}$  : each variable has *l* values  $a^1, ..., a^l$
  - Function  $\phi(A_1, A_2)$  is an indicator function for the event  $A_1 = A_2$ 
    - E.g., two super-pixels have the same greyscale







# Features in Text Analysis Compact for variables with large domains



Y = target variables

Srihari

X = known variables

- -X are words of text, Y are named entities
  - *B-PER*=Begin Person, *I-PER*=within person, *OTH*=Not entity
- Factors for word t:  $\Phi^{1}_{t}(Y_{t}, Y_{t+1}), \Phi^{2}_{t}(Y_{t}, X_{1}, ..., X_{T})$
- Features (hundreds of thousands):
  - Word itself (capitalized, in list of common names)
  - Aggregate features of sequence (>2 sports related)
  - $Y_t$  dependent on several words in a window of t
  - Skip chain CRF: connections between adjacent words & multiple occurrences of same word