



## Undirected Graphical Models I

Lecture 4 – Apr 6, 2011  
CSE 515, Statistical Methods, Spring 2011

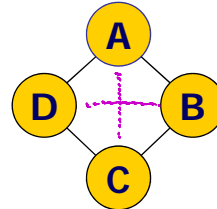
Instructor: Su-In Lee  
University of Washington, Seattle

## Bayesian Network Representation

- Directed acyclic graph structure
  - Conditional parameterization
  - Independencies in graphs
  - From distribution to BN graphs
- Conditional probability distributions (CPDs)
  - Table
  - Deterministic
  - Context-specific (Tree, Rule CPDs)
  - Independence of causal influence (Noisy OR, GLMs)
  - Continuous variables
  - Hybrid models

## The *Misconception* Example

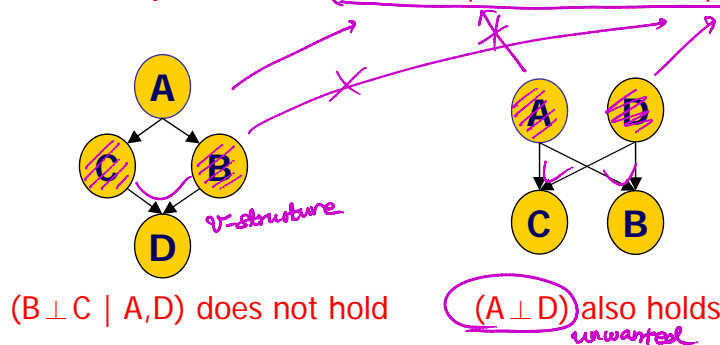
- Four students get together in pairs to work on HWs: **Alice, Bob, Charles, Debbie**
- Only the following pairs meet: (A&B), (B&C), (C&D), (D&A)
- Let's say that the prof accidentally misspoke in class
  - Each student may subsequently have figured out the problem.
  - In subsequent study pairs, they may transmit this newfound understanding to their partners.
- Consider 4 binary random variables
  - A, B, C, D: whether the student has the misconception or not.
- Independence assumptions?
  - (A ⊥ C | B, D), (B ⊥ D | A, C)
- Can we find the P-map for these?



3

## Reminder: Perfect Maps

- G is a **perfect map (P-map)** for P if  $I(P) = I(G)$
- Does every distribution have a P-map?
  - No: some structures cannot be represented in a BN
    - Independencies in P: (A ⊥ D | B, C) and (B ⊥ C | A, D)



## Representing Dependencies

- $(A \perp D \mid B, C)$  and  $(B \perp C \mid A, D)$ 
  - **Cannot** be modeled with a Bayesian network.
  - **Can** be modeled with an **undirected graphical models** (Markov networks).

## Undirected Graphical Models (Informal)

- **Nodes** correspond to random variables
- **Edges** correspond to direct probabilistic interaction
  - An interaction not mediated by any other variables in the network.

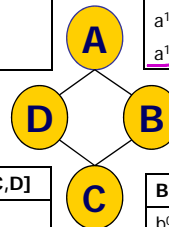
- How to **parameterize**?

- **Local factor models** are attached to sets of nodes

- Factor elements are positive
- Do not have to sum to 1
- Represent *affinities*, *not a prob.*  
compatibilities

A	D	$\pi_1[A, C]$
$a^0$	$d^0$	100
$a^0$	$d^1$	1
$a^1$	$d^0$	1
$a^1$	$d^1$	100

A	B	$\pi_2[A, B]$
$a^0$	$b^0$	30
$a^0$	$b^1$	5
$a^1$	$b^0$	1
$a^1$	$b^1$	10



C	D	$\pi_3[C, D]$
$c^0$	$d^0$	1
$c^0$	$d^1$	100
$c^1$	$d^0$	100
$c^1$	$d^1$	1

B	C	$\pi_4[B, C]$
$b^0$	$c^0$	100
$b^0$	$c^1$	1
$b^1$	$c^0$	1
$b^1$	$c^1$	1000

## Undirected Graphical Models (Informal)

- Represents joint distribution

- Unnormalized factor

$$F(a, b, c, d) = \pi_1[a, b] \pi_2[a, c] \pi_3[b, d] \pi_4[c, d]$$

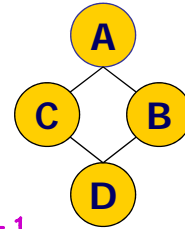
- Probability =  $\pi_1[A=a, B=b] \pi_2[A=a, C=c] \pi_3[B=b, D=d] \pi_4[C=c, D=d]$  — factor multiplication

$$P(a, b, c, d) = \frac{1}{Z} \pi_1[a, b] \pi_2[a, c] \pi_3[b, d] \pi_4[c, d]$$

$Z \rightarrow$  normalization  $\sum_{a,b,c,d} p(a,b,c,d) = 1$

- Partition function

$$Z = \sum_{a,b,c,d} \pi_1[a, b] \pi_2[a, c] \pi_3[b, d] \pi_4[c, d]$$



- As undirected graphical models represent joint distributions, they can be used for answering queries.

## Undirected Graphical Models Blurb

- Useful when edge directionality cannot be assigned

- Simpler interpretation of structure

- Simpler inference
  - Simpler independency structure

- Harder to learn parameters/structures *why?*

$$\text{eg. } Z = \sum_{x_1} \dots \sum_{x_n} F(x_1, \dots, x_n)$$

- We will also see models with combined directed and undirected edges *(eg. conditional random fields)*

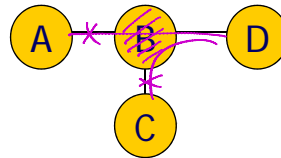
- Markov networks

# Markov Network Structure

- Undirected graph H
  - Nodes  $X_1, \dots, X_n$  represent random variables
- H encodes independence assumptions
  - A path  $X_1-X_2-\dots-X_k$  is **active** if none of the  $X_i$  variables along the path are observed
  - **X** and **Y** are separated in H given **Z** if there is no active path between any node  $x \in X$  and any node  $y \in Y$  given **Z**
    - Denoted  $\text{sep}_H(X;Y|Z)$



$$D \perp \{A, C\} \mid B$$



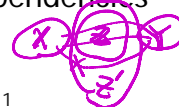
**Global** independencies associated with H:

$$I(H) = \{(X \perp Y | Z) : \text{sep}_H(X;Y|Z)\}$$

9

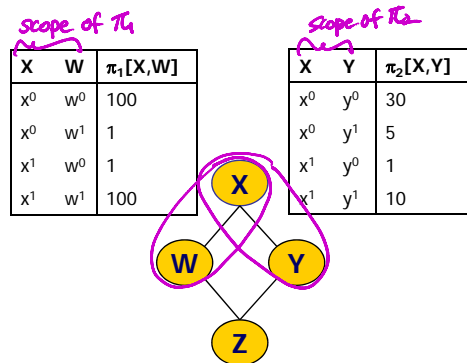
# Relationship with Bayesian Network

- Bayesian network
  - **Local** independencies → Independence by **d-separation** (**global**)
- Markov network  $I_H(G)$ 
  - **Global** independencies → **Local** independencies
- Can all independencies encoded by Markov networks be encoded by Bayesian networks?
  - No, counter example –  $(A \perp B \mid C, D)$  and  $(C \perp D \mid A, B)$
- Can all independencies encoded by Bayesian networks be encoded by Markov networks?
  - No, immoral v-structures (explaining away)
- Markov networks encode monotonic independencies
  - If  $\text{sep}_H(X;Y|Z)$  and  $Z \subseteq Z'$  then  $\text{sep}_H(X;Y|Z')$



## Markov Network Factors

- A **factor** (or "potential") is a function from value assignments of a set of random variables  $\mathbf{D}$  to real positive numbers  $\mathbb{R}^+$ 
  - The set of variables  $\mathbf{D}$  is the **scope** of the factor
- Factors generalize the notion of CPDs
  - Every CPD is a factor (with additional constraints)



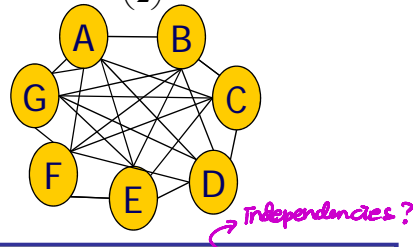
11

## Factors and Joint Distribution

- Can we represent any joint distribution by using only factors that are defined on edges?
  - **No!** Compare # of parameters
  - Example:  $n$  binary RVs
    - Joint distribution has  $2^n - 1$  independent parameters
    - Markov network with edge factors has  $4 \binom{n}{2}$  parameters

Needed:  $2^7 - 1 = 127!$

Edge parameters:  $4 \cdot \binom{7}{2} = 84$



- Factors introduce constraints on joint distribution

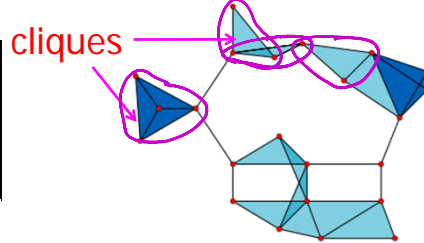
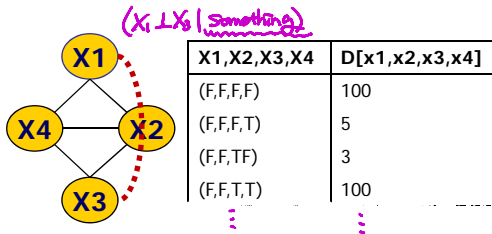
## Factors and Graph Structure

- Are there constraints imposed on the network structure  $H$  by a factor whose scope is  $\mathbf{D}$ ?
  - Hint 1: think of the independencies that must be satisfied
  - Hint 2: generalize from the basic case of  $|\mathbf{D}|=2$

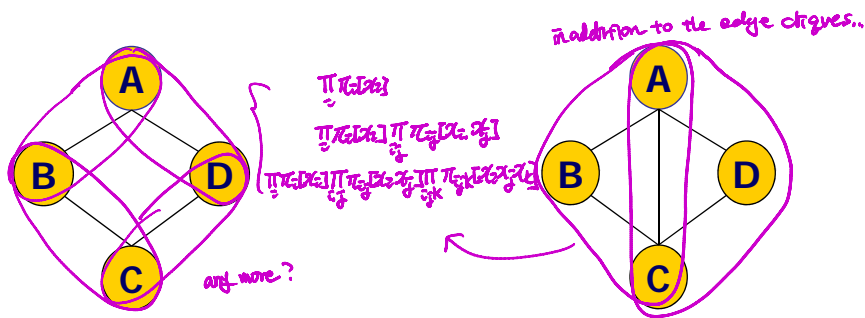


The induced subgraph over  $\mathbf{D}$  must be a clique (fully connected)

**Why?** otherwise two unconnected variables may be independent by blocking the active path between them, contradicting the direct dependency between them in the factor over  $\mathbf{D}$



## Markov Network Factors: Examples



Maximal cliques

- $\{A, B\}$
- $\{B, C\}$
- $\{C, D\}$
- $\{A, D\}$

Maximal cliques

- $\{A, B, C\}$
- $\{A, C, D\}$

## Markov Network Distribution

- A distribution  $P$  factorizes over  $H$  if it has:
  - A set of subsets  $\mathbf{D}_1, \dots, \mathbf{D}_m$  where each  $\mathbf{D}_i$  is a complete (fully connected) subgraph in  $H$  *clique*
  - Factors  $\pi_1[\mathbf{D}_1], \dots, \pi_m[\mathbf{D}_m]$  such that

$$P(X_1, \dots, X_n) = \frac{1}{Z} f(X_1, \dots, X_n) = \frac{1}{Z} \prod \pi_i[\mathbf{D}_i]$$

where un-normalized factor:  $f(X_1, \dots, X_n) = \prod \pi_i[\mathbf{D}_i]$

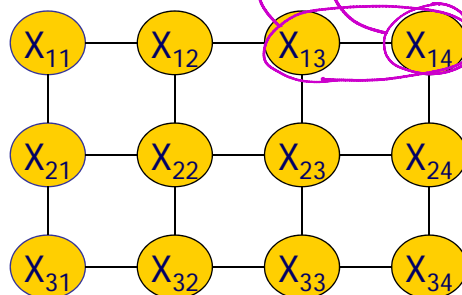
$$Z = \sum_{X_1, \dots, X_n} f(X_1, \dots, X_n) = \sum_{X_1, \dots, X_n} \prod \pi_i[\mathbf{D}_i]$$

- $Z$  is called the **partition function**
- $P$  is also called a **Gibbs distribution** over  $H$

## Pairwise Markov Networks

- A **pairwise Markov network** over a graph  $H$  has:
  - A set of **node potentials**  $\{\pi[X_i] : i=1, \dots, n\}$
  - A set of **edge potentials**  $\{\pi[X_i, X_j] : X_i, X_j \in H\}$

- Example:





## Logarithmic Representation

- We represent energy potentials by **applying a log transformation to the original potentials**
  - $\pi[\mathbf{D}] = \exp(-\varepsilon[\mathbf{D}])$  where  $\varepsilon[\mathbf{D}] = -\ln \pi[\mathbf{D}]$
  - Any Markov network parameterized with factors can be converted to a logarithmic representation
  - The log-transformed potentials can take on any real value
  - The joint distribution decomposes as

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

Log P(X) is a linear function of  $\varepsilon_i[\mathbf{D}_i]$

## I-Maps and Factorization

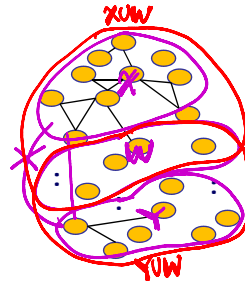
- Independence mappings (I-map)
  - I(P) – set of independencies ( $\mathbf{X} \perp \mathbf{Y} \mid \mathbf{Z}$ ) in P
  - I-map – independencies by a graph is a subset of I(P)
- Bayesian Networks
  - Factorization and reverse factorization theorems
    - G is an I-map of P iff P factorizes as  $P(X_1, \dots, X_n) = \prod_{i=1}^n P(X_i \mid Pa(X_i))$
- Markov Networks
  - Factorization and reverse factorization theorems
    - H is an I-map of P iff P factorizes as  $P(X_1, \dots, X_n) = \frac{1}{Z} \prod \pi_i[\mathbf{D}_i]$  *proof?*

## Reverse Factorization

■  $P(X_1, \dots, X_n) = \frac{1}{Z} \prod \pi_i[D_i] \rightarrow H$  is an I-map of  $P$

■ Proof:

- Let  $X, Y, W$  be any three **disjoint** sets of variables such that  $W$  separates  $X$  and  $Y$  in  $H$
- We need to show  $(X \perp Y | W) \in I(P)$



■ **Case 1:**  $X \cup Y \cup W = U$  (all variables)

- As  $W$  separates  $X$  and  $Y$  there are no direct edges between  $X$  and  $Y$
- any clique in  $H$  is fully contained in  $X \cup W$  or  $Y \cup W$
- Let  $I_X$  be cliques in  $X \cup W$  and  $I_Y$  be cliques in  $Y \cup W$  (not in  $I_X$ )

$$\rightarrow P(X_1, \dots, X_n) = \frac{1}{Z} \prod_{i \in I_X} \pi_i[D_i] \prod_{i \in I_Y} \pi_i[D_i] = \frac{1}{Z} f(X, W) g(Y, W)$$

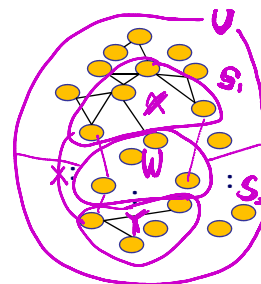
→  $(X \perp Y | W) \in I(P)$

## Reverse Factorization

■  $P(X_1, \dots, X_n) = \frac{1}{Z} \prod \pi_i[D_i] \rightarrow H$  is an I-map of  $P$

■ Proof:

- Let  $X, Y, W$  be any three **disjoint** sets of variables such that  $W$  separates  $X$  and  $Y$  in  $H$
- We need to show  $(X \perp Y | W) \in I(P)$



■ **Case 2:**  $X \cup Y \cup W \subset U$  (all variables)

- Let  $S = U - (X \cup Y \cup W)$
- $S$  can be partitioned into two disjoint sets  $S_1$  and  $S_2$  such that  $W$  separates  $X \cup S_1$  and  $Y \cup S_2$  in  $H$
- From case 1, we can derive  $(X, S_1 \perp Y, S_2 | W) \in I(P)$
- From decomposition of independencies

→  $(X \perp Y | W) \in I(P)$

## Factorization

- If H is an I-map of P then  $P(X_1, \dots, X_n) = \frac{1}{Z} \prod \pi_i[D_i]$
- Holds only for positive distributions P
  - Hammerly-Clifford theorem
- Defer proof

## Relationship with Bayesian Network

- Bayesian Networks *local → global*
  - **Semantics** defined via **local** independencies  $I_L(G)$ .
  - **Global** independencies induced by d-separation
  - Local and global independencies equivalent since one implies the other
- Markov Networks *global → local*
  - **Semantics** defined via **global** separation property  $I(H)$
  - Can we define the induced **local** independencies?
    - We show two definitions (call them “**Local Markov assumptions**”)
    - All three definitions (global and two local) are equivalent only for positive distributions P

## Pairwise Independencies

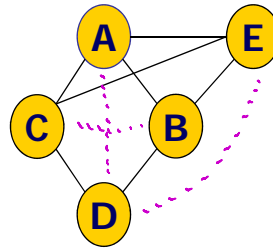
- Every pair of disconnected nodes are separated given all other nodes in the network
- Formally:  $I_p(H) = \{ (X \perp Y \mid U - \{X, Y\}) : X - Y \notin H \}$

Example:

$(A \perp D \mid B, C, E)$

$(B \perp C \mid A, D, E)$

$(D \perp E \mid A, B, C)$



## Local Independencies

- Every node is independent of all other nodes given its immediate neighboring nodes in the network  
Markov blank of  $X$ ,  $MB_H(X)$
- Formally:  $I_L(H) = \{ (X \perp U - \{X\} - MB_H(X) \mid MB_H(X)) : X \in H \}$

Example:

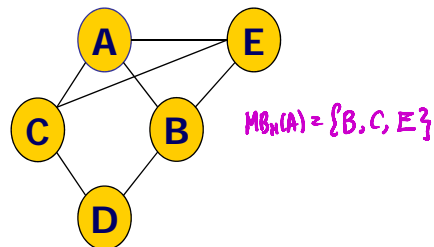
$(A \perp D \mid B, C, E)$

$(B \perp C \mid A, D, E)$

$(C \perp B \mid A, D, E)$

$(D \perp E, A \mid B, C)$

$(E \perp D \mid A, B, C)$



## Relationship Between Properties

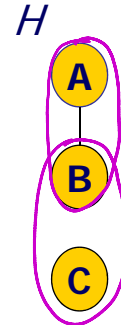
- Let  $I(H)$  be the **global separation** independencies
- Let  $I_L(H)$  be the **local (Markov blanket)** independencies
- Let  $I_p(H)$  be the **pairwise** independencies
  
- For **any** distribution  $P$ :
  - $I(H) \rightarrow I_L(H)$ 
    - The assertion in  $I_L(H)$ , that a node is independent of all other nodes given its neighbors, is part of the separation independencies since there is no active path between a node and its non-neighbors given its neighbors
  
  - $I_L(H) \rightarrow I_p(H)$ 
    - Follows from the monotonicity of independencies in Markov networks (if  $(\mathbf{X} \perp \mathbf{Y} | \mathbf{Z})$  and  $\mathbf{Z} \subseteq \mathbf{Z}'$  then  $(\mathbf{X} \perp \mathbf{Y} | \mathbf{Z}')$ )

## Relationship Between Properties

- Let  $I(H)$  be the **global separation** independencies
- Let  $I_L(H)$  be the **local (Markov blanket)** independencies
- Let  $I_p(H)$  be the **pairwise** independencies
  
- For any **positive** distribution  $P$ :
  - $I_p(H) \rightarrow I(H)$ 
    - Proof relies on intersection property for probabilities  $(\mathbf{X} \perp \mathbf{Y} | \mathbf{Z}, \mathbf{W})$  and  $(\mathbf{X} \perp \mathbf{W} | \mathbf{Z}, \mathbf{Y}) \rightarrow (\mathbf{X} \perp \mathbf{Y}, \mathbf{W} | \mathbf{Z})$  which holds in general only for positive distributions
    - **Details on the textbook**
  
  - Thus, for positive distributions
    - $I(H) \leftrightarrow I_L(H) \leftrightarrow I_p(H)$
  
  - How about a non-positive distribution?

## The Need for Positive Distribution

- Let  $P$  satisfy
  - $A$  is uniformly distributed
  - $A=B=C$
- $P$  satisfies  $I_P(H)$ 
  - $(B \perp C|A), (A \perp C|B)$   
(since each variable determines all others)
- $P$  does not satisfy  $I_L(H)$ 
  - $(C \perp A,B)$  needs to hold according to  $I_L(H)$  but does not hold in the distribution



## Constructing Markov Network for $P$

- **Goal:** Given a distribution, we want to construct a Markov network which is an I-map of  $P$
- Complete (fully connected) graphs will satisfy but are not interesting
- Minimal I-maps: A graph  $G$  is a minimal I-Map for  $P$  if:
  - $G$  is an I-map for  $P$
  - Removing any edge from  $G$  renders it not an I-map
- **Goal:** construct a graph which is a minimal I-map of  $P$

## Constructing Markov Network for P

- If P is a positive distribution, then  $I(H) \leftrightarrow I_L(H) \leftrightarrow I_P(H)$ 
  - Thus, sufficient to construct a network that satisfies  $I_P(H)$
- Construction algorithm
  - For every  $(X,Y)$  add edge if  $(X \perp Y | \mathcal{U} - \{X,Y\})$  does not hold in P
- **Theorem: network is minimal and unique I-map**
  - Proof:
    - **I-map** follows since  $I_P(H)$  by construction and  $I(H)$  by equivalence
    - **Minimality** follows since deleting an edge implies  $(X \perp Y | \mathcal{U} - \{X,Y\})$   
But, we know by construction that this does not hold in P since we added the edge in the construction process
    - **Uniqueness** follows since any other I-map has at least these edges and to be minimal cannot have additional edges

## Summary: Markov Network Representation

- Independencies in graph H
  - **Global** independencies  $I(H) = \{(X \perp Y | Z) : \text{sep}_H(X; Y | Z)\}$
  - **Local** independencies  $I_L(H) = \{(X \perp \mathcal{U} - \{X\} - \text{MB}_H(X) | \text{MB}_H(X)) : X \in H\}$
  - **Pairwise** independencies  $I_P(H) = \{(X \perp Y | \mathcal{U} - \{X,Y\}) : X - Y \notin H\}$
  - For any positive distribution P, they are equivalent.
- **(Reverse) factorization theorem: I-map  $\leftrightarrow$  factorization**
- Markov network **factors**
  - Has to encompass cliques
  - Maximal cliques, edge potentials
- **Pairwise Markov network**
  - Node/ edge potentials
  - Application in vision (image segmentation)
- **What next?**
- **Log-linear model**
  - Log-transformation of potentials
  - Features instead of factors
- Constructing Markov networks from Bayesian networks
- “Partially” directed graph (e.g. Conditional Random Fields)

## Acknowledgement

- These lecture notes were generated based on the slides from Prof Eran Segal.