

# Undirected Graphical Models

Sargur Srihari

[srihari@cedar.buffalo.edu](mailto:srihari@cedar.buffalo.edu)

# Topics

- Directed versus Undirected graphical models
- Components of a Markov Network
- Independence Properties
- Parameterization
- Gibbs Distributions and Markov Networks
- Reduced Markov Networks

# Directed vs Undirected

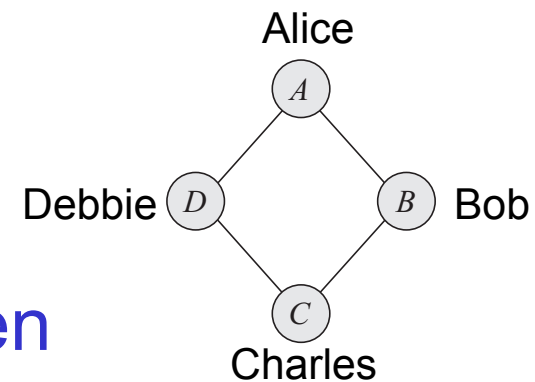
- Bayesian Networks= Directed Graphical Models
  - Useful in many real-world domains
- Undirected Graphical Models
  - When no natural directionality exists betw. variables
  - Offer a simpler perspective on directed graphs
    - Independence structure
    - Inference task
  - Also called Markov Networks or MRFs
- Can combine directed and undirected graphs
- Attention restricted to discrete state spaces

# Example to Motivate Undirected Graphs

## 1. Four students study in pairs to work on homework

Alice and Bob are friends  
Bob and Charles study together  
Charles and Debbie argue with each other  
Debbie and Alice study together  
Alice and Charles can't stand each other  
Bob and Debbie had relationship ended badly

A Social Network



## 2. Professor may have mis-spoken

e.g., on a machine learning topic

## 3. Students may have figured out the problem

e.g., by thinking about issue  
or by studying textbook

## 4. Students transmit this understanding to his/her study partner

# Modeling the Misconception Problem

Four binary random variables representing whether or not student has misconception

$$X \in \{A, B, C, D\}$$

$x^1$  : student has the misconception

$x^0$  : Student does not have a misconception

## Prior Probabilities

$a^0$	$a^1$	$b^0$	$b^1$
0.3	0.7	0.2	0.8

$a^0$  = has misconception

$a^1$  = no misconception

# Modeling Influences Using a BN

Probability of misconception of one person depends on whether their study partner has a misconception

Alice and Charles never speak *directly*

Thus  $A$  and  $C$  should be conditionally independent given  $B$  and  $D$

We need to model  $(A \perp C | \{B, D\})$

Consider this Proposed Bayesian Network

It does model  $(A \perp C | \{B, D\})$

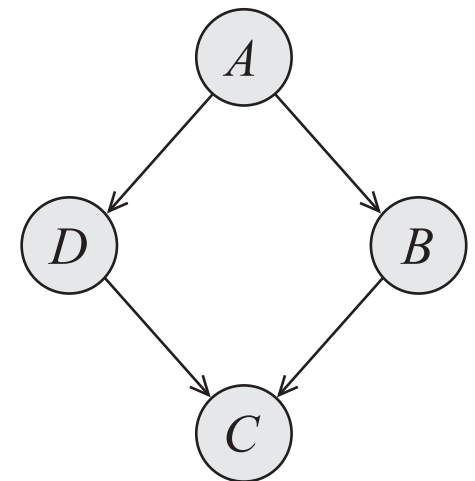
since the path between  $A$  and  $C$  is blocked when  $B, D$  are known

But also means  $B$  and  $D$  are independent given only  $A$

since V-structure through  $C$  implies blockage when  $C$  is not known

But dependent given both  $A$  and  $C$

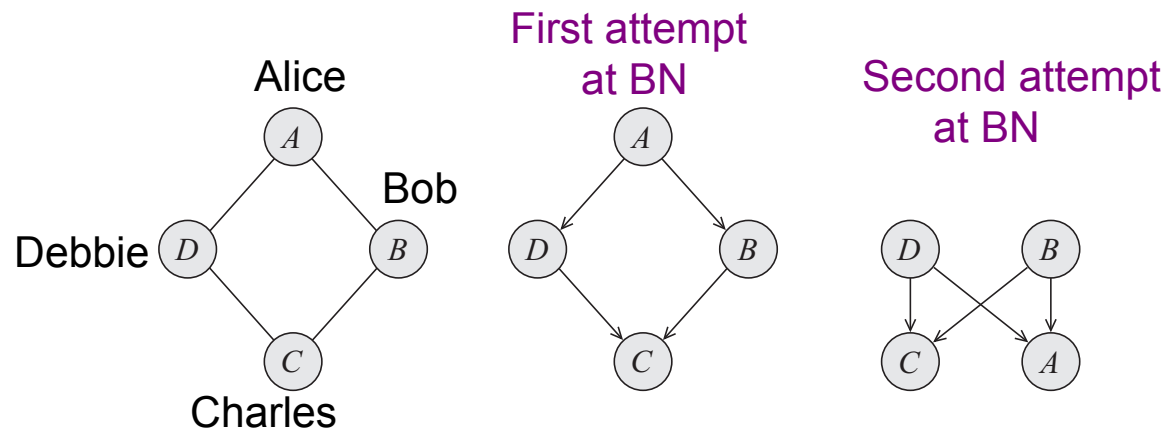
since V-structure through  $C$  implies no blockage when  $C$  is known



# Lack of Perfect Map in BN

## Misconception Example

We need to model  $(A \perp C | B, D), (B \perp D | A, C)$



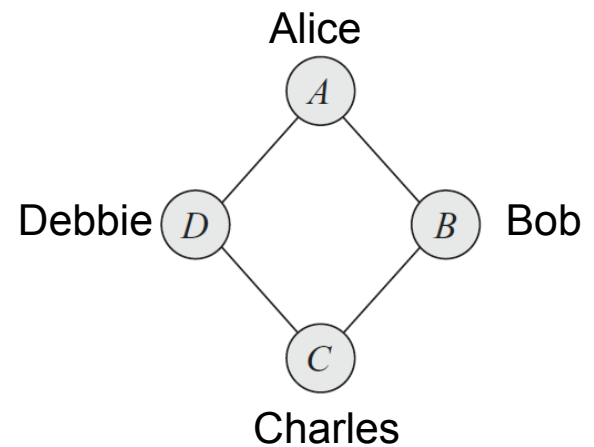
In both (b) and (c)  $(A \perp C | B, D)$  holds  
but  $(B \perp D | A, C)$  due to v-structure  $D \rightarrow C \leftarrow B$  in d-separation

No perfect map since independences imposed by BN are inappropriate for the distribution

In a perfect map the graph precisely captures the independencies in the given distribution

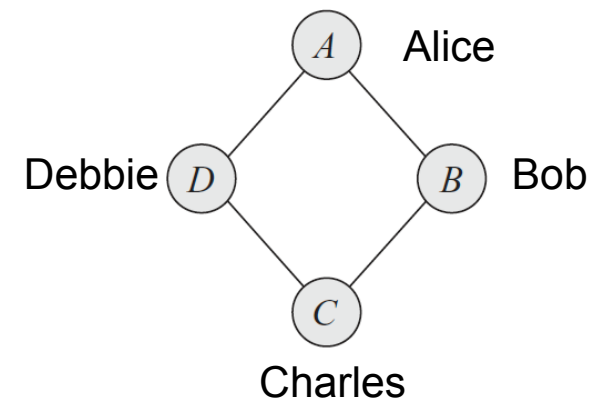
# Drawbacks of BN in Example

- Independences imposed by BN are inappropriate for the distribution
- Interaction between the variables are symmetric



# Components of a Markov Network

- Nodes correspond to variables
- Edges correspond to direct probabilistic interactions
- How to parameterize it?
  - Interaction is not directed
    - No need for CPDs
  - Need a more symmetric parameterization
    - Capture Affinities between variables
      - Bob and Alice more often agree than disagree
    - Associate with A, B a general purpose function
      - Called a “Factor”

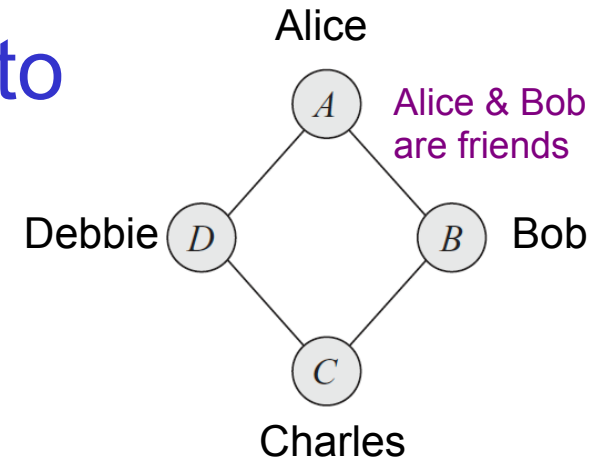


# Factor definition

- Function from  $Val(D)$  to  $R$ 
  - Where  $D$  is a set of random variables
- Attention restricted to non-negative factors
  - $Val(A,B)$  to  $R^+$
  - Value associated with  $(a,b)$  indicates affinity between them
  - Higher value indicates higher compatibility
- Factor captures compatibility between variables in its scope

# Example of a Factor

- Factor is not normalized to sum to one
- Values need not be in  $[0, 1]$
- $\Phi_1(A, B)$  asserts that  $A$  and  $B$ 
  - are more likely to agree than disagree
  - More weight to when they agree than when they disagree
    - Translates to higher probability for disagreeing values



$$\phi_1[A, B]$$

$a^0$	$b^0$	30
$a^0$	$b^1$	5
$a^1$	$b^0$	1
$a^1$	$b^1$	10

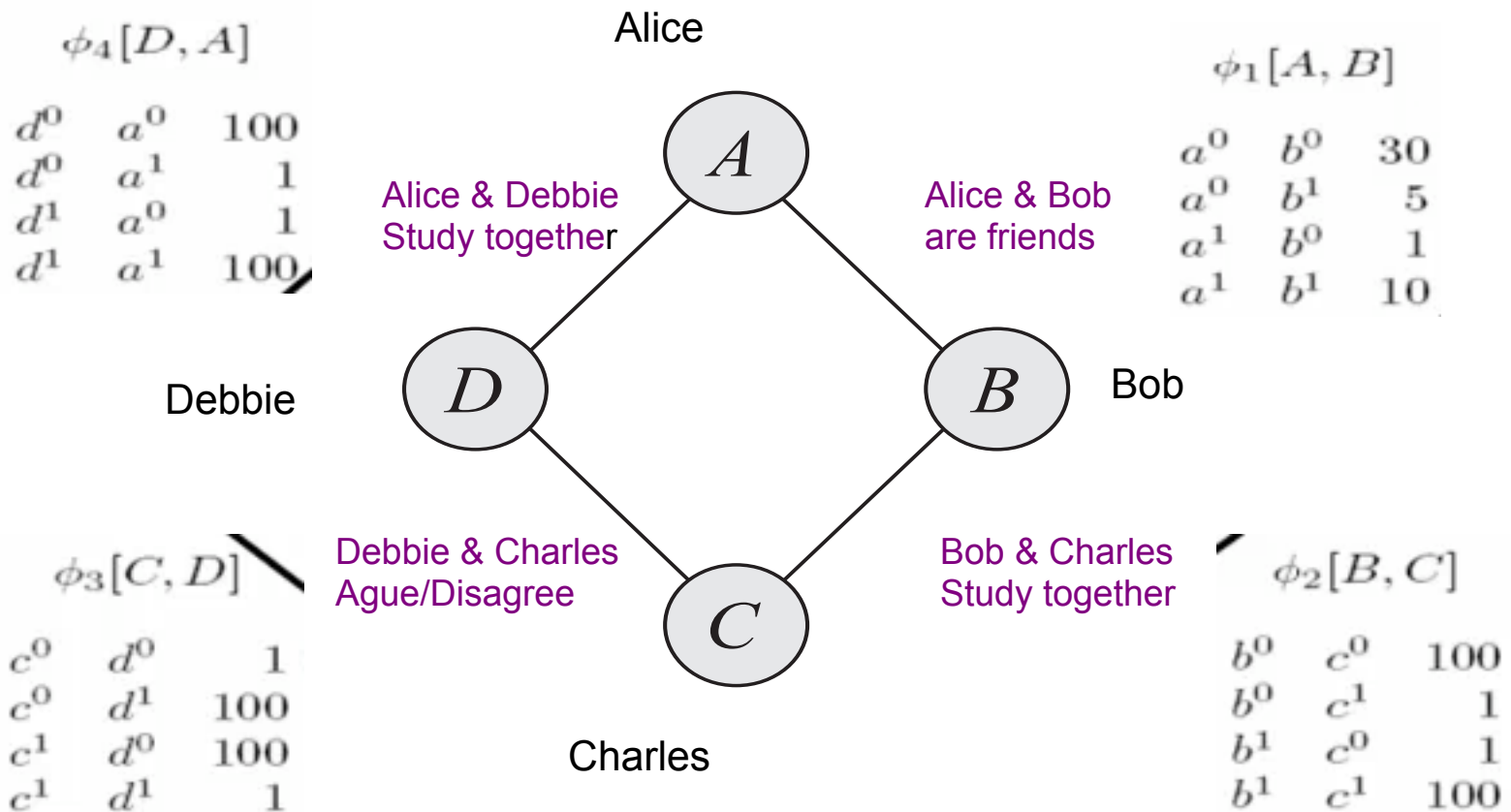
$a^0$  = has misconception

$a^1$  = no misconception

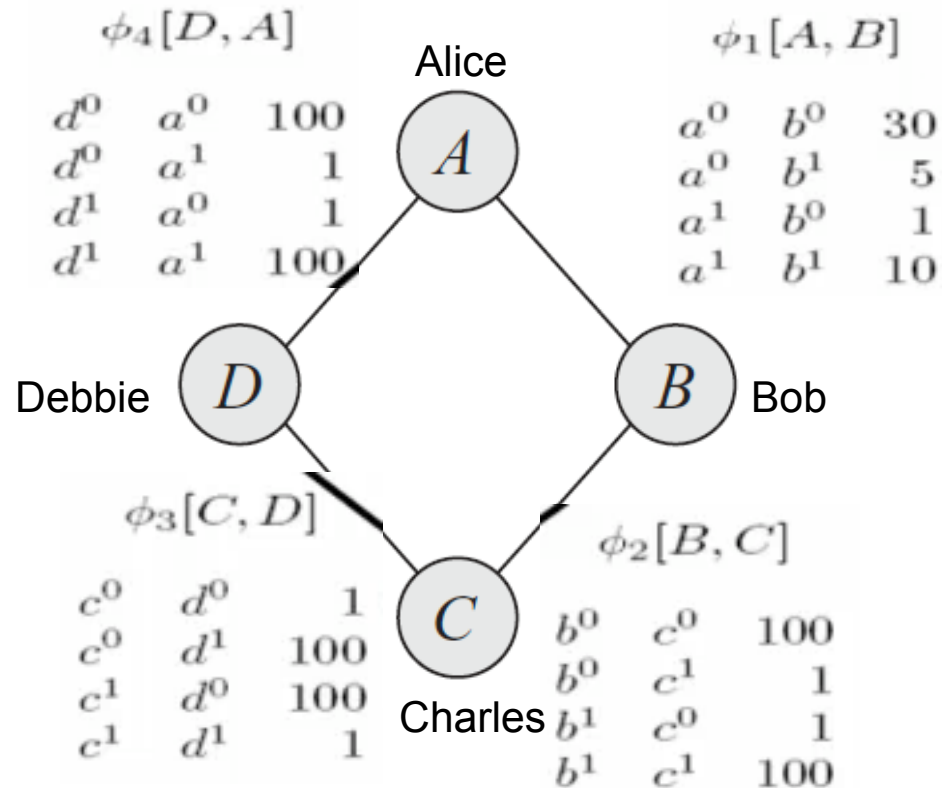
$b^0$  = has misconception

$b^1$  = no misconception

# Factors for Misconception Example



# Pairwise Compatibility Factors



$a^0$  = has misconception

$a^1$  = no misconception

$b^0$  = has misconception

$b^1$  = no misconception

Charles and Debbie  
are incompatible

Most likely instantiations are  
when they disagree

Factor is similar to CPD  
For each combination there is a value

# 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

$$\phi_1(a^1, b^1) \cdot \phi_2(b^1, c^0) \cdot \phi_3(c^0, d^1) \cdot \phi_4(d^1, a^1) = 10.1.100.100 = 100,000$$

Assignment				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

- Convert to a legal distribution by performing a normalization

# 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

Assignment				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

$$Z = 7,201,840$$

We can obtain any desired probability from the joint distribution

$P(b^0) = 0.268$ : Bob is 26% likely to have a misconception

$P(b^1 | c^0) = 0.06$ : if Charles does not have the misconception, Bob is 6% likely.

Most probable joint probability:  $P(a^0, b^1, c^1, d^0) = 0.69$

Alice, Debby have no misconception, Bob and Charles have misconception

# Independence Properties

- Tight connection between factorization of the distribution and independencies
- $P$  supports  $(X \perp Y|Z)$  iff we can write distribution as  $P(X) = \Phi_1(X,Z)\Phi_2(Y,Z)$

$$P(A,B,C,D) = \underbrace{\left[ \frac{1}{Z} \phi_1(A,B) \cdot \phi_2(B,C) \right]}_{\text{Factor with } \{B, \{A,C\}\}} \cdot \underbrace{\left[ \phi_3(C,D) \cdot \phi_4(D,A) \right]}_{\text{Factor with } \{D, \{A,C\}\}}$$

Factor with  $\{B, \{A,C\}\}$

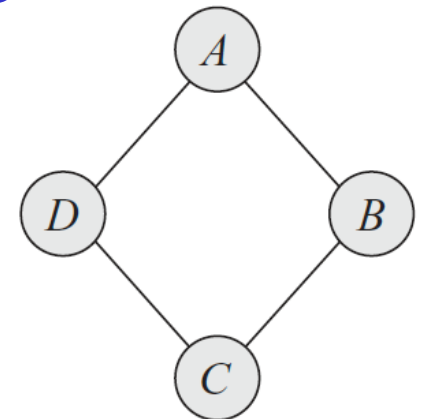
Factor with  $\{D, \{A,C\}\}$

– Therefore  $(B \perp D|A,C)$

– By grouping factors with  $\{A, \{B,D\}\}$  and  $\{C, \{B,D\}\}$

- We get  $(A \perp C|B,D)$

– Precisely the independencies we could not get with a Bayesian network



# Parameterization

- Not as intuitive as CPDs in Bayesian networks
- Significantly harder to estimate from data
- Factors represent both joint distributions and CPDs
  - A joint distribution over  $D$  is a factor over  $D$ 
    - Note that a product of factors is a factor
    - Specifies real no for every value of  $D$
  - A conditional distribution  $P(X|U)$  is a factor over  $\{X\}$   
 $UU$
- Should satisfy normalization constraints
  - No restriction on parameters in a factor

$\phi_1[A, B]$		
$a^0$	$b^0$	30
$a^0$	$b^1$	5
$a^1$	$b^0$	1
$a^1$	$b^1$	10

# Factors

- Factor describes *compatibilities* between different values of variables in its scope
- A graph can be parameterized by associating a set of factors with it
- Only associating a factor with each edge is insufficient for a full distribution

# Pairwise Factors

- Associating parameters with each edge is insufficient to specify arbitrary joint distribution
- Consider a fully connected graph
  - There are no independence assumptions
- If all variables are binary
  - Each factor over an edge has 4 parameters
  - Total number of parameters is  $4 \binom{n}{2}$  (a quadratic in  $n$ )
  - An arbitrary distribution needs  $2^n - 1$  parameters
- Pairwise factors have insufficient parameters
  - Parameters cannot be just associated with edges

# Definition of Factor Product

Let  $X$ ,  $Y$  and  $Z$  be three disjoint sets of variables

$\Phi_1(X, Y)$  and  $\Phi_2(Y, Z)$  are two factors

Factor product is defined to match up the common part of  $Y$

$$\Psi(X, Y, Z) = \Phi_1(X, Y) \cdot \Phi_2(Y, Z)$$

$\Phi_1(A, B)$

$a^1$	$b^1$	0.5
$a^1$	$b^2$	0.8
$a^2$	$b^1$	0.1
$a^2$	$b^2$	0
$a^3$	$b^1$	0.3
$a^3$	$b^2$	0.9

$\Phi_2(B, C)$

$b^1$	$c^1$	0.5
$b^1$	$c^2$	0.7
$b^2$	$c^1$	0.1
$b^2$	$c^2$	0.2



$\psi(A, B, C)$

$a^1$	$b^1$	$c^1$	$0.5 \cdot 0.5 = 0.25$
$a^1$	$b^1$	$c^2$	$0.5 \cdot 0.7 = 0.35$
$a^1$	$b^2$	$c^1$	$0.8 \cdot 0.1 = 0.08$
$a^1$	$b^2$	$c^2$	$0.8 \cdot 0.2 = 0.16$
$a^2$	$b^1$	$c^1$	$0.1 \cdot 0.5 = 0.05$
$a^2$	$b^1$	$c^2$	$0.1 \cdot 0.7 = 0.07$
$a^2$	$b^2$	$c^1$	$0 \cdot 0.1 = 0$
$a^2$	$b^2$	$c^2$	$0 \cdot 0.2 = 0$
$a^3$	$b^1$	$c^1$	$0.3 \cdot 0.5 = 0.15$
$a^3$	$b^1$	$c^2$	$0.3 \cdot 0.7 = 0.21$
$a^3$	$b^2$	$c^1$	$0.9 \cdot 0.1 = 0.09$
$a^3$	$b^2$	$c^2$	$0.9 \cdot 0.2 = 0.18$

Two factors are multiplied in a way that matches-up the common part  $Y$

# Factor Product and Bayesian network

- Both joint distributions and CPDs are factors
- Chain rule of Bayesian networks defines joint distribution as product of factors

$$P(A,B)=P(A)P(B|A)$$

– Letting  $\phi_{X_i}(X_i, Pa_{X_i})$  represent  $P(X_i|PaX_i)$  we have

$$P(X_1, \dots, X_n) = \prod_i \phi_{X_i}$$

# Gibbs Distributions

- Generalizes the idea of factor product
- A distribution  $P_{\Phi}$  is a Gibbs distribution parameterized by a set of factors

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

$D_i$  are sets of random variables

– If it is defined as follows

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

where

$$\tilde{P}(X_1, \dots, X_n) = \prod_{i=1}^m \phi_i(D_i)$$

is an unnormalized measure and

$$Z = \sum_{X_1, \dots, X_n} \tilde{P}(X_1, \dots, X_n)$$

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

# Influence of Factors

- Factors do not represent marginal probabilities of the variables within their scope (as we will see)
- A factor is only one contribution to the overall joint distribution
- Distribution as a whole has to take into consideration contributions from all factors involved
- Example given next

# Factors .ne. Marginal Probabilities

## Joint Distribution

Assignment				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

## Marginal distribution over Alice, Bob

$a^0$	$b^0$	0.13
$a^0$	$b^1$	0.69
$a^1$	$b^0$	0.14
$a^1$	$b^1$	0.04

## Factor over Alice, Bob

$\phi_1[A, B]$

$a^0$	$b^0$	30
$a^0$	$b^1$	5
$a^1$	$b^0$	1
$a^1$	$b^1$	10

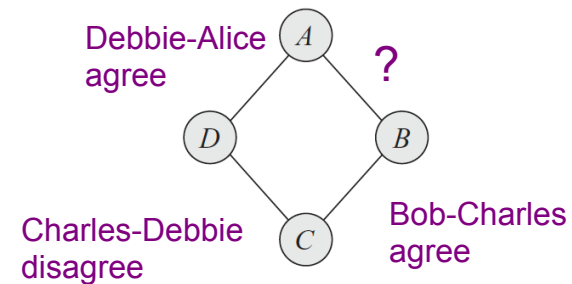
Most likely:  $a^0, b^1$  (disagree)       $a^0, b^0$  (agree)

Because probability takes into account influence of other factors

$\phi_2[B, C]$			$\phi_3[C, D]$		
$b^0$	$c^0$	100	$c^0$	$d^0$	1
$b^0$	$c^1$	1	$c^0$	$d^1$	100
$b^1$	$c^0$	1	$c^1$	$d^0$	100
$b^1$	$c^1$	100	$c^1$	$d^1$	1

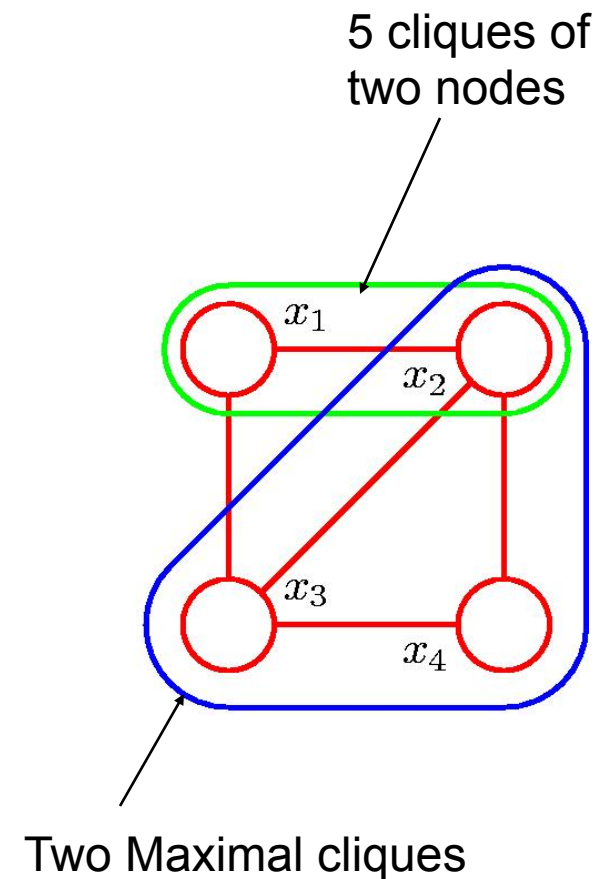
$\phi_4[D, A]$

$d^0$	$a^0$	100
$d^0$	$a^1$	1
$d^1$	$a^0$	1
$d^1$	$a^1$	100



# Clique in a graph

- Subset of nodes in graph such that there exists a link between all pairs of nodes in subset
  - Set of nodes in clique are fully connected
- Maximal Clique
  - Not possible to include any other nodes in the graph in the set without ceasing to be a clique



# Gibbs distribution and Markov Network

- A distribution  $P_\Phi$  with  $\Phi = \{\phi_1(D_1), \dots, \phi_K(D_K)\}$ 
  - factorizes over a Markov network  $H$  if each  $D_k$  is a complete subgraph (clique) of  $H$
- Factors that parameterize network are called clique potentials
- Number of factors are reduced, e.g.,

Four Maximal Cliques:

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

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

with possible parameter settings:

$$\phi_1[A, B]$$

$a^0$	$b^0$	30
$a^0$	$b^1$	5
$a^1$	$b^0$	1
$a^1$	$b^1$	10

$$\phi_2[B, C]$$

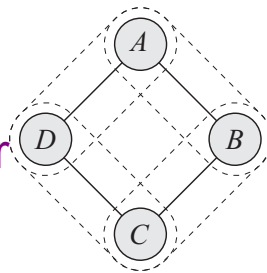
$b^0$	$c^0$	100
$b^0$	$c^1$	1
$b^1$	$c^0$	1
$b^1$	$c^1$	100

$$\phi_3[C, D]$$

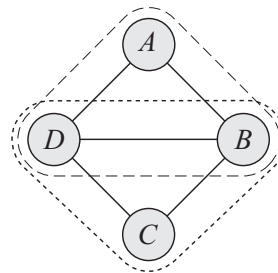
$c^0$	$d^0$	1
$c^0$	$d^1$	100
$c^1$	$d^0$	100
$c^1$	$d^1$	1

$$\phi_4[D, A]$$

$d^0$	$a^0$	100
$d^0$	$a^1$	1
$d^1$	$a^0$	1
$d^1$	$a^1$	100



(a)



(b)

Two Maximal Cliques:

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

$\{B, C, D\}$

# Maximal Cliques obscure structure present

- No of parameters reduced by allowing factors only for maximal cliques
  - But obscures structure present
- Example: All  $n$  variables are binary
  - Pairwise cliques
    - Each factor over an edge has 4 parameters
    - Total number of parameters is  $4 {}^n C_2$ 
      - Quadratic number of parameters
  - An arbitrary distribution (Single clique)
    - needs  $2^n - 1$  parameter
      - Exponential number of parameters

# Factors as Cliques

- Functions of maximal cliques
- Set of variables in clique  $C$  is denoted  $x_C$
- Joint distribution is written as a product of *potential functions*

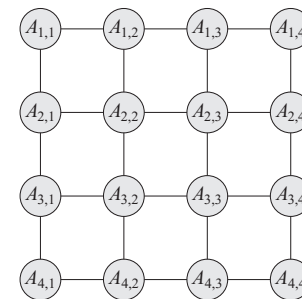
$$p(\mathbf{x}) = \frac{1}{Z} \prod_C \psi_C(\mathbf{x}_C)$$

- Where  $Z$ , the *partition function*, is a normalization constant

$$Z = \sum_{\mathbf{x}} \prod_C \psi_C(\mathbf{x}_C)$$

# Pairwise Markov Network

- Subclass of Markov networks commonly encountered, e.g.,
  1. Ising model and Boltzmann machines
  2. Computer Vision
- All factors are over single variables Or pairs of variables
- A pairwise Markov network over graph  $H$ 
  - has a set of node potentials  $\Phi(X_i)$  and edge potentials  $\Phi(X_i, X_j)$
- Although simple they pose a challenge for inference algorithms
- Special Case when arranged as a grid



# Reduced Markov Networks

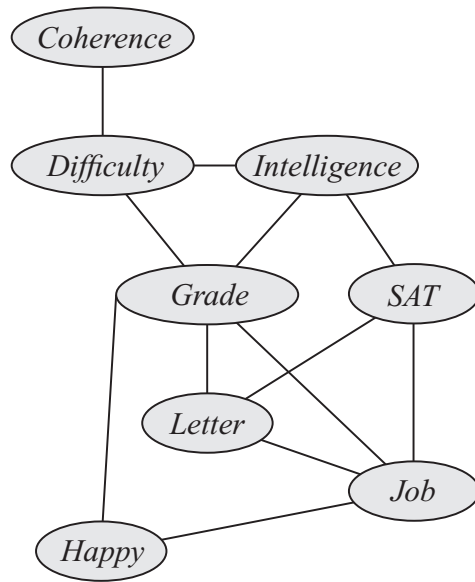
$a^1$	$b^1$	$c^1$	$0.5 \cdot 0.5 = 0.25$
$a^1$	$b^1$	$c^2$	$0.5 \cdot 0.7 = 0.35$
$a^1$	$b^2$	$c^1$	$0.8 \cdot 0.1 = 0.08$
$a^1$	$b^2$	$c^2$	$0.8 \cdot 0.2 = 0.16$
$a^2$	$b^1$	$c^1$	$0.1 \cdot 0.5 = 0.05$
$a^2$	$b^1$	$c^2$	$0.1 \cdot 0.7 = 0.07$
$a^2$	$b^2$	$c^1$	$0 \cdot 0.1 = 0$
$a^2$	$b^2$	$c^2$	$0 \cdot 0.2 = 0$
$a^3$	$b^1$	$c^1$	$0.3 \cdot 0.5 = 0.15$
$a^3$	$b^1$	$c^2$	$0.3 \cdot 0.7 = 0.21$
$a^3$	$b^2$	$c^1$	$0.9 \cdot 0.1 = 0.09$
$a^3$	$b^2$	$c^2$	$0.9 \cdot 0.2 = 0.18$

Factor is reduced to the context  $C = c^1$

$a^1$	$b^1$	$c^1$	0.25
$a^1$	$b^2$	$c^1$	0.08
$a^2$	$b^1$	$c^1$	0.05
$a^2$	$b^2$	$c^1$	0
$a^3$	$b^1$	$c^1$	0.15
$a^3$	$b^2$	$c^1$	0.09

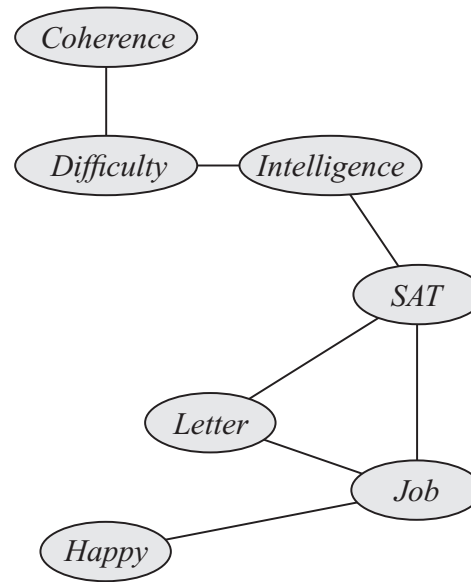
Conditioning distribution on some assignment  $\mathbf{u}$  to set of variables  $\mathbf{U}$   
 Conditioning a distribution is to eliminate all entries in the joint distribution that are inconsistent with a subset of variables  $\mathbf{U}=\mathbf{u}$  and renormalizing remaining entries to sum to 1

# Reduced Markov Network Example



(a)

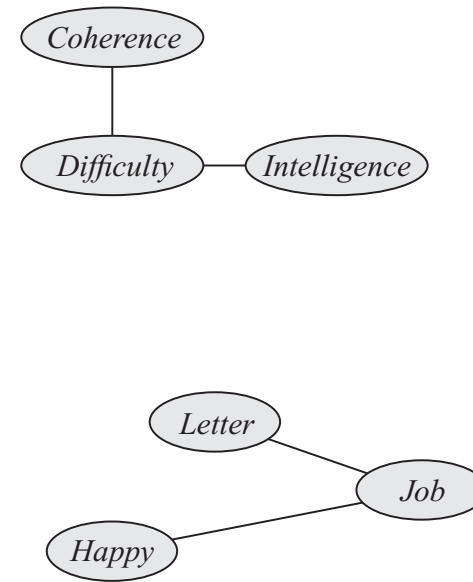
Initial set of factors



(b)

Reduced to context

$G=g$



(c)

Reduced to context

$G=g, S=s$