

# TTIC 31170: Planning, Learning, and Estimation for Robotics and Artificial Intelligence

Spring 2021

Matthew Walter  
TTI-Chicago

## Lecture 2: Graphical Models (Continued)

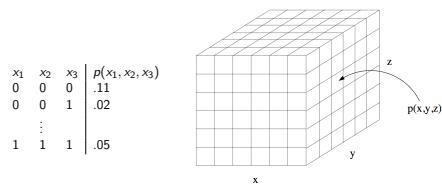
# Key challenges

- **Represent** the world as a collection of random variables  $X_1, \dots, X_n$  with joint distribution  $P(X_1, \dots, X_n)$ 
  - How do you compactly represent the joint distribution?
  - Directed and undirected graphical models
- **Infer** distribution  $P(X_i | X_1 = x_1, \dots, X_m = x_m)$  over random variable(s) based upon evidence
  - How do you make inference tractable?
  - Exact vs. approximate
- **Learn** the distribution from available data

2

## What is a good representation?

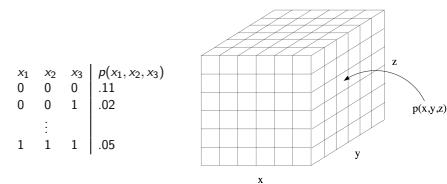
- Properties of a good representation:
  - Explicit
  - Modular (generalizable)
  - Permits efficient computation (i.e., learning & inference)



3

## What is a good representation?

- Properties of a good representation:
  - Explicit
  - Modular (generalizable)
  - Permits efficient computation (i.e., learning & inference)

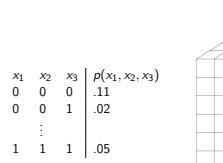


A naive representation requires a number of parameters that is exponential in the number of variables

3

## What is a good representation?

- Properties of a good representation:
  - Explicit
  - Modular (generalizable)
  - Permits efficient computation (i.e., learning & inference)

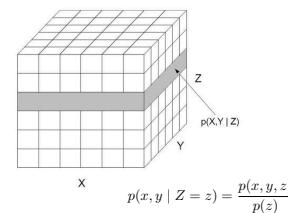


3

The amount of data required to learn these parameters is exponential in the number of variables

## What is a good representation?

- Properties of a good representation:
  - Explicit
  - Modular (generalizable)
  - Permits efficient computation (i.e., learning & inference)

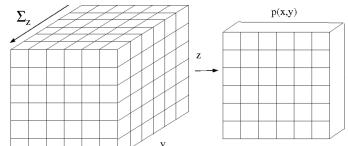


4

## What is a good representation?

- Properties of a good representation:

- Explicit
- Modular (generalizable)
- Permits efficient computation (i.e., learning & inference)



$$p(x, y) = \sum_{z \in \mathcal{Z}} p(x, y, z)$$

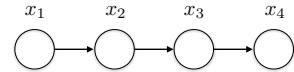
$$p(x) = \sum_{y \in \mathcal{Y}} p(x, y)$$

5

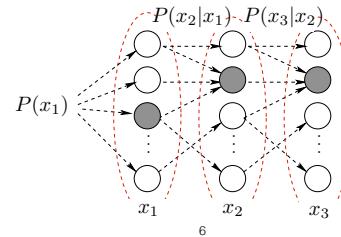
Cost of inference is exponential in the number of random variables

## Representation: Explicit

- Graphical model: Representation in terms of variables & dependencies:



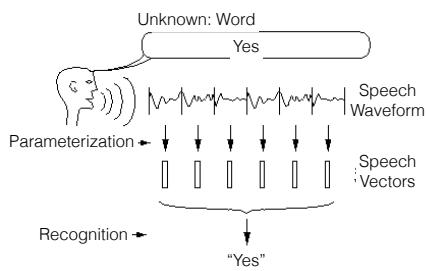
- Transition diagram: Representation in terms of state transitions



6

## Graphical models: Examples

- Hidden Markov model as a Bayesian network
  - Example: Speech recognition

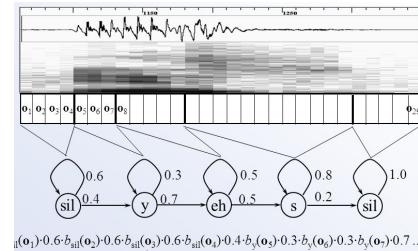


7

## Graphical models: Examples

- Hidden Markov model as a Bayesian network
  - Example: Speech recognition

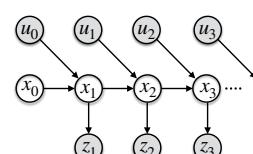
$$O = \{o_1, o_2, \dots, o_T\} \quad P(w_i|O) = \frac{P(O|w_i)P(w_i)}{P(O)} = \frac{P(O|w_i)P(w_i)}{P(O)}$$



8

## Graphical models: Examples

- Dynamic Bayesian networks
  - Example: Robot localization (and mapping)



$x_t$ : robot's pose

$u_t$ : control input (action)

$z_t$ : measurement

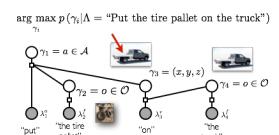
9

## Graphical models: Examples

- Factor graph
  - Example: Natural language symbol grounding



$$\arg \max_{\gamma_i} p(\gamma_i | \Lambda) = \text{"Put the tire pallet on the truck"}$$

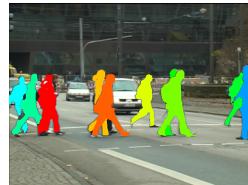
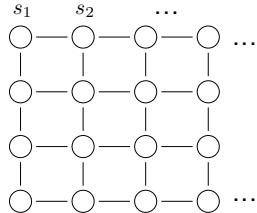


$\gamma_i$ : Groundings (objects, actions)  
 $\lambda_i$ : Phrases in instruction

10

## Graphical models: Examples

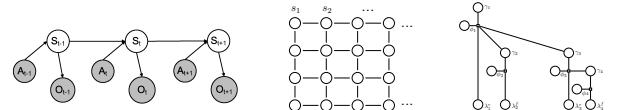
- Lattice models (e.g., Ising model) as a Markov random field
  - Example: Semantic image segmentation



11

## Graphical models

- Graph semantics:
  - graph → separation properties → independence
- Expression of probability distributions:
  - independence → family of distributions
- Inference and estimation:
  - graph structure → efficient computation



12

## Topics

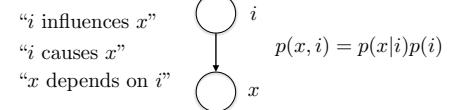
- Representations and graphical models
- Directed graphical models (Bayesian networks)
  - graphs and independence
- Undirected graphical models (Markov random fields)
  - graphs, independence
  - Bayesian networks as undirected models
- Inference

13

## Bayesian networks

- Directed acyclic graphs (DAGs)
  - Nodes represent variables
  - Directed edges express dependencies

A mixture model as  
a Bayesian network



14

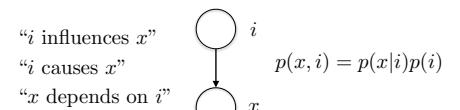
## Bayesian networks

- Directed acyclic graphs (DAGs)
    - Nodes represent variables
    - Directed edges express dependencies
- $i$  influences  $x$
- $i$  causes  $x$
- $x$  depends on  $i$
- A mixture model as  
a Bayesian network
- $p(x, i) = p(x|i)p(i)$

## Bayesian networks

- Directed acyclic graphs (DAGs)
  - Nodes represent variables
  - Directed edges express dependencies

A mixture model as  
a Bayesian network

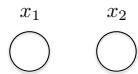


$$p(x_v) = \prod_{v \in \mathcal{V}} p(x_v | \pi_v) \quad \text{where } \pi_v \text{ are the parents of } v$$

15

16

## Example: Coin tosses



17

## Example: Coin tosses

$$P(x_1) = \begin{matrix} 0.5 \\ 0.5 \end{matrix} \quad \begin{matrix} x_1 \\ \text{circle} \end{matrix} \quad \begin{matrix} x_2 \\ \text{circle} \end{matrix} \quad P(x_2) = \begin{matrix} 0.5 \\ 0.5 \end{matrix}$$

18

## Example: Coin tosses

$$P(x_1) = \begin{matrix} 0.5 \\ 0.5 \end{matrix} \quad \begin{matrix} x_1 \\ \text{circle} \end{matrix} \quad \begin{matrix} x_2 \\ \text{circle} \end{matrix} \quad P(x_2) = \begin{matrix} 0.5 \\ 0.5 \end{matrix}$$

same?  $x_3 \in \{\text{yes, no}\}$

$$P(x_3|x_1, x_2) = \begin{array}{c|ccccc} & \text{hh} & \text{ht} & \text{th} & \text{tt} \\ \hline \text{y} & 1.0 & 0.0 & 0.0 & 1.0 \\ \text{n} & 0.0 & 1.0 & 1.0 & 0.0 \end{array}$$

19

## Example: Coin tosses

$$P(x_1) = \begin{matrix} 0.5 \\ 0.5 \end{matrix} \quad \begin{matrix} x_1 \\ \text{circle} \end{matrix} \quad \begin{matrix} x_2 \\ \text{circle} \end{matrix} \quad P(x_2) = \begin{matrix} 0.5 \\ 0.5 \end{matrix}$$

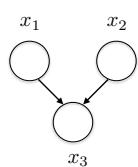
same?  $x_3 \in \{\text{yes, no}\}$

$$P(x_3|x_1, x_2) = \begin{array}{c|ccccc} & \text{hh} & \text{ht} & \text{th} & \text{tt} \\ \hline \text{y} & 1.0 & 0.0 & 0.0 & 1.0 \\ \text{n} & 0.0 & 1.0 & 1.0 & 0.0 \end{array}$$

- Two levels of description
  - graph structure (dependencies and independencies)
  - associated probability distribution

20

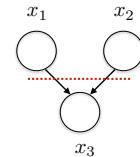
## What does the graph tell us?



21

## What does the graph tell us?

- $x_1$  and  $x_2$  are *marginally independent*

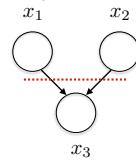


Knowing  $x_1$  doesn't tell us anything about  $x_2$

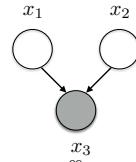
22

## What does the graph tell us?

- $x_1$  and  $x_2$  are *marginally independent*



- $x_1$  and  $x_2$  become *conditionally dependent* if we know  $x_3$



## Driving example

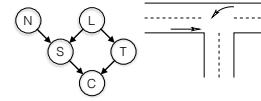
N: X obeys rules of road

L: Traffic light

S: X decides to stop?

T: The other car turns left?

C: Crash?



24

## Driving example

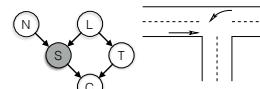
N: X obeys rules of road

L: Traffic light

S: X decides to stop?

T: The other car turns left?

C: Crash?

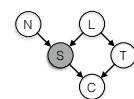


- If we know that X decided to stop (S), can knowing whether X obeys the law (N) tell us anything about the other car turning (T)?

25

## Graph, independence, d-separation

- Are N and T independent given S?

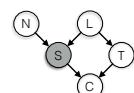


N: X obeys rules of road  
L: Traffic light  
S: X decides to stop?  
T: The other car turns left?  
C: Crash?

26

## Graph, independence, d-separation

- Are N and T independent given S?



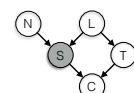
Definition: Variables N and T are *d-separated* given S if S separates them in the moralized ancestral graph

(d in d-separation stands for "dependence")

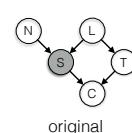
27

## Graph, independence, d-separation

- Are N and T independent given S?



Definition: Variables N and T are *d-separated* given S if S separates them in the moralized ancestral graph

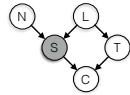


original

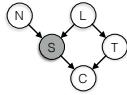
28

## Graph, independence, d-separation

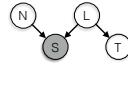
- Are N and T independent given S?



**Definition:** Variables N and T are **d-separated** given S if S separates them in the moralized ancestral graph



original

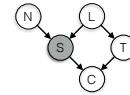


ancestral

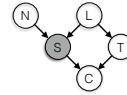
29

## Graph, independence, d-separation

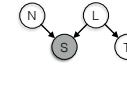
- Are N and T independent given S?



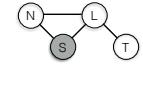
**Definition:** Variables N and T are **d-separated** given S if S separates them in the moralized ancestral graph



original



ancestral

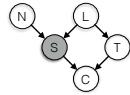


moralized ancestral

30

## Graph, independence, d-separation

- Are N and T independent given S?



**Definition:**  
S separates N and T  
if S is a parent of both N and T

N and T are **not**  
d-separated  
given S

N and T are **not**  
independent  
given S

original

ancestral

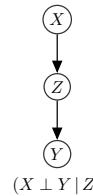
moralized ancestral

30

## Graph, independence, d-separation

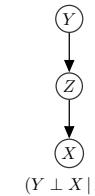
- Bayesian network structure implies conditional independencies

**Cascade**  
(Markov chain; causal/evidential trail)



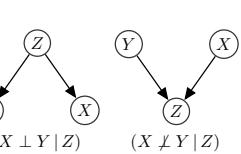
$$(X \perp\!\!\!\perp Y | Z)$$

**Common Cause**



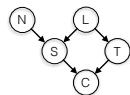
$$(Y \perp\!\!\!\perp X | Z)$$

**Common Effect**  
(v-structure)



31

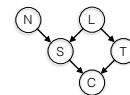
## Graphs and distributions



- A graph is a compact representation of a large collection of independence properties

32

## Graphs and distributions



- A graph is a compact representation of a large collection of independence properties

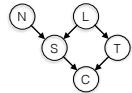
**Theorem:** Any probability distribution that is consistent with a directed graph  $G$  has to factor according to “node given parents”:

$$P(X|G) = \prod_{i=1}^d P(x_i|x_{\pi_i})$$

where  $\pi_i$  are the parents of  $x_i$  and  $d$  is the number of nodes in the graph

33

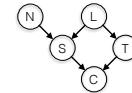
## Graphs and distributions



$$P(X|G) = \prod_{i=1}^d P(x_i|x_{\pi_i})$$

34

## Graphs and distributions

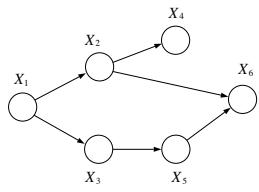


$$P(X|G) = \prod_{i=1}^d P(x_i|x_{\pi_i})$$

$$P(CSTNL) = P(C|S,T)P(T|L)P(S|N,L)P(N)P(L)$$

34

## Graphs and distributions

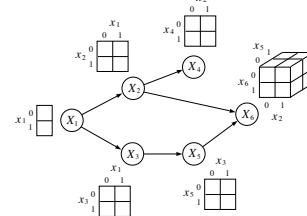


Representational (storage, learning, & inference) complexity

- **Joint distribution:** Exponential in the number of variables
- **Bayesian Network:** Exponential in the number of parents of each node, linear in the number of nodes

35

## Graphs and distributions



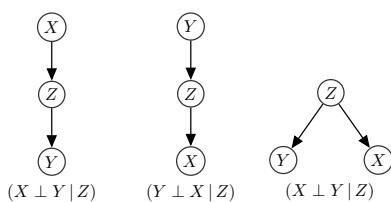
Representational (storage, learning, & inference) complexity

- **Joint distribution:** Exponential in the number of variables
- **Bayesian Network:** Exponential in the number of parents of each node, linear in the number of nodes

36

## Graphs and distributions

- Different Bayesian networks can encode the same conditional independencies (and the same distributions)



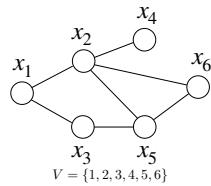
37

## Topics

- Representations and graphical models
- Directed graphical models (Bayesian networks)
  - graphs and independence
- Undirected graphical models (Markov random fields)
  - graphs, independence
  - Bayesian networks as undirected models
- Inference

38

## Undirected graphical models

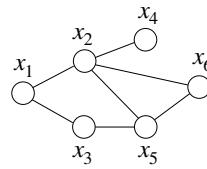


$$G = (V, E)$$

$$V = \{1, 2, 3, 4, 5, 6\}$$

39

## Undirected graphical models



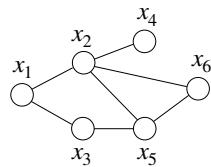
$$G = (V, E)$$

Definition: A *clique*  $C \subseteq V$  is a subset of fully-connected nodes

$$C = \{\{1,2\}, \{1,3\}, \{2,4\}, \{3,5\}, \{2,5,6\}\}$$

40

## Undirected graphical models

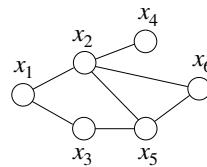


$$G = (V, E)$$

Definition: A non-negative *potential function*  $\psi_C(x_C)$  is associated with each clique  $C$

41

## Undirected graphical models



$$G = (V, E)$$

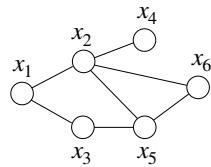
Theorem: The probability distribution for an undirected graph can be expressed as a product of clique potentials

$$p(x_V) = \frac{1}{Z} \prod_{C \in C} \psi_C(x_C)$$

where  $Z$  is a normalization constant

42

## Undirected graphical models



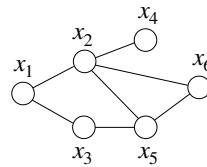
$$G = (V, E)$$

$$p(x_V) = \frac{1}{Z} \prod_{C \in C} \psi_C(x_C)$$

$$p(x_V) = \frac{1}{Z} \psi(x_1, x_2) \psi(x_1, x_3) \psi(x_2, x_4) \psi(x_3, x_5) \psi(x_2, x_5, x_6)$$

43

## Undirected graphical models



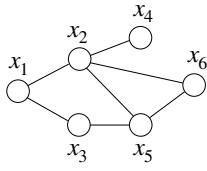
$$G = (V, E)$$

Definition: A non-negative *potential function*  $\psi_C(x_C)$  is associated with each clique  $C$

Example: Tabular

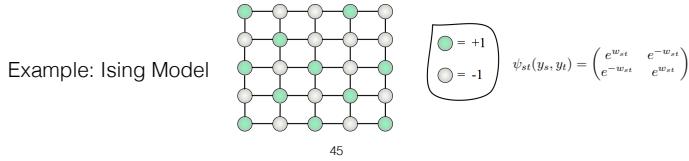
44

## Undirected graphical models

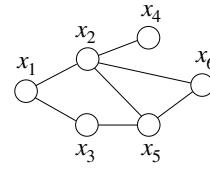


$$G = (V, E)$$

**Definition:** A non-negative *potential function*  $\psi_C(x_C)$  is associated with each clique  $C$



## Undirected graphical models



$$G = (V, E)$$

**Definition:** A non-negative *potential function*  $\psi_C(x_C)$  is associated with each clique  $C$

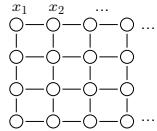
Example: Gaussian Markov random field

$$\prod_i \exp\left(-\frac{1}{2}(\Lambda_{ii}x_i^2 + 2\eta_i x_i)\right) \cdot \prod_{i,j:i \neq j} \exp\left(-\Lambda_{ij}x_i x_j\right)$$

46

## Undirected graphical models

- For example, a lattice model with binary variables  $x_i \in \{-1, 1\}$  (e.g., foreground vs. background) and pairwise edges

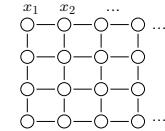


$$\begin{aligned} p(x_1, \dots, x_n) &= \frac{1}{Z} \prod_{C \in \mathcal{C}} \psi_C(x_C) \\ &= \frac{1}{Z} \prod_{(i,j) \in E} \psi_{i,j}(x_i, x_j) \end{aligned}$$

47

## Undirected graphical models

- For example, a lattice model with binary variables  $x_i \in \{-1, 1\}$  (e.g., foreground vs. background) and pairwise edges



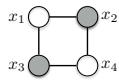
$$\begin{aligned} p(x_1, \dots, x_n) &= \frac{1}{Z} \prod_{(i,j) \in E} \psi_{i,j}(x_i, x_j) \\ &= \frac{1}{Z} \prod_{(i,j) \in E} \exp(J_{ij}x_i x_j) \end{aligned}$$

where  $J_{ij}$  specifies the “interaction strength” between  $x_i$  and  $x_j$

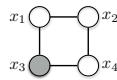
48

## Graph semantics

- Graph semantics for undirected graph come from separation



$x_1$  and  $x_4$  are independent given  $x_2$  and  $x_3$



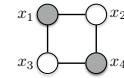
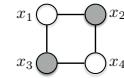
$x_1$  and  $x_4$  are not independent given  $x_3$

49

## Graph semantics: Comparison

- Directed and undirected graphs are complementary

The following two independence properties can not be captured *simultaneously* with a Bayesian network

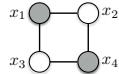
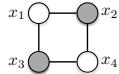


50

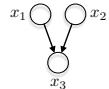
## Graph semantics: Comparison

- Directed and undirected graphs are complementary

The following two independence properties can not be captured *simultaneously* with a Bayesian network



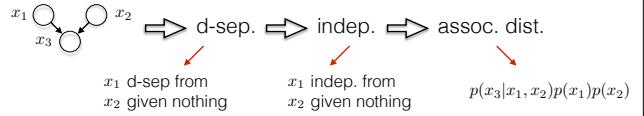
An undirected graph can not express marginal independence



51

## Summary: Graphs and Probabilities

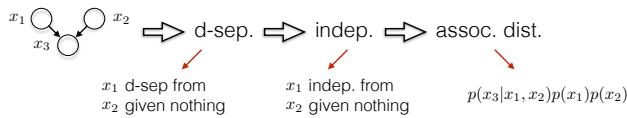
- Directed graphical models (Bayesian networks)



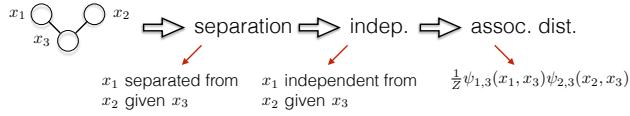
52

## Summary: Graphs and Probabilities

- Directed graphical models (Bayesian networks)



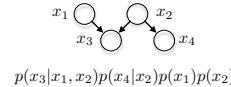
- Undirected graphical models (Markov random fields)



53

## Graph transformations

- Directed graphical models can be converted into undirected graphical models via **moralization** (aka "marry the parents"):

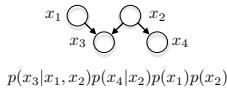


$p(x_3|x_1, x_2)p(x_4|x_2)p(x_1)p(x_2)$

54

## Graph transformations

- Directed graphical models can be converted into undirected graphical models via **moralization** (aka "marry the parents"):

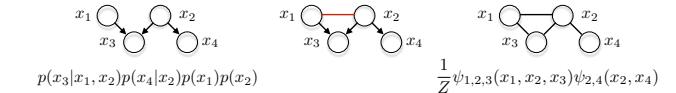


$p(x_3|x_1, x_2)p(x_4|x_2)p(x_1)p(x_2)$

55

## Graph transformations

- Directed graphical models can be converted into undirected graphical models via **moralization** (aka "marry the parents"):



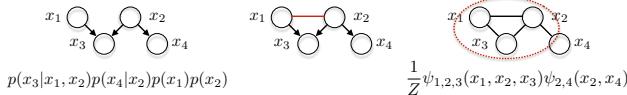
$p(x_3|x_1, x_2)p(x_4|x_2)p(x_1)p(x_2)$

$\frac{1}{Z} \psi_{1,2,3}(x_1, x_2, x_3) \psi_{2,4}(x_2, x_4)$

56

## Graph transformations

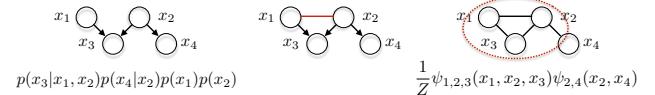
- Directed graphical models can be converted into undirected graphical models via **moralization** (aka “marry the parents”):



56

## Graph transformations

- Directed graphical models can be converted into undirected graphical models via **moralization** (aka “marry the parents”):



- The distribution doesn't change, only the graph
- The undirected graph is consistent with the distribution associated with the original graph

57

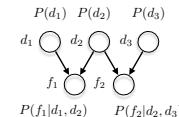
## Topics

- Representations and graphical models
- Directed graphical models (Bayesian networks)
  - graphs and independence
- Undirected graphical models (Markov Random Fields)
  - graphs, independence
  - Bayesian networks as undirected models
- Inference

58

## Inference

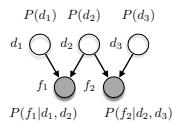
- Event prediction example: Binary event variables  $d$  and possible observations  $f$



59

## Inference

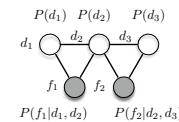
- Event prediction example: Binary event variables  $d$  and possible observations  $f$



60

## Inference

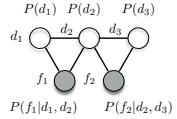
- Moralize the directed graphical model



61

## Inference

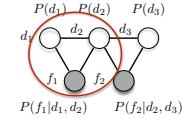
- Moralize the directed graphical model



62

## Inference

- Moralize the directed graphical model

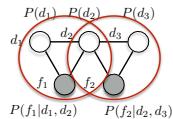


$$\psi_{1,2}(d_1, d_2) = P(d_1)P(d_2)P(f_1^*|d_1, d_2)$$

62

## Inference

- Moralize the directed graphical model



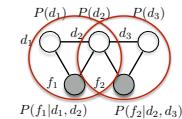
$$\psi_{1,2}(d_1, d_2) = P(d_1)P(d_2)P(f_1^*|d_1, d_2)$$

$$\psi_{2,3}(d_2, d_3) = P(d_3)P(f_2^*|d_2, d_3)$$

62

## Inference

- Moralize the directed graphical model



$$\psi_{1,2}(d_1, d_2) = P(d_1)P(d_2)P(f_1^*|d_1, d_2)$$

$$\psi_{2,3}(d_2, d_3) = P(d_3)P(f_2^*|d_2, d_3)$$

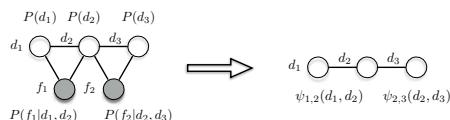
- Join distribution as a product of "interaction potentials"

$$P(d_1, d_2, d_3, \text{data}) = \psi_{1,2}(d_1, d_2) \cdot \psi_{2,3}(d_2, d_3)$$

62

## Inference

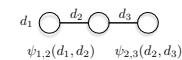
- Moralize the directed graphical model



$$P(d_1, d_2, d_3, \text{data}) = \psi_{1,2}(d_1, d_2) \cdot \psi_{2,3}(d_2, d_3)$$

63

## Marginalization and messages



- It suffices to calculate the following marginals:

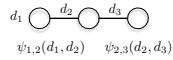
$$\begin{aligned} P(d_1, \text{data}) &= \sum_{d_2, d_3} P(d_1, d_2, d_3, \text{data}) \\ P(d_2, \text{data}) &= \sum_{d_1, d_3} P(d_1, d_2, d_3, \text{data}) \\ P(d_3, \text{data}) &= \sum_{d_1, d_2} P(d_1, d_2, d_3, \text{data}) \end{aligned}$$

We can use these marginals to calculate desired probabilities:

$$P(d_1|\text{data}) = \frac{P(d_1, \text{data})}{P(\text{data})} = \frac{P(d_1, \text{data})}{\sum_{d'_1} P(d'_1, \text{data})}$$

64

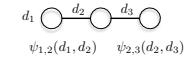
## Marginalization and messages



$$P(d_2, d_3, \text{data}) = \sum_{d_1} P(d_1, d_2, d_3, \text{data})$$

65

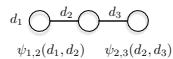
## Marginalization and messages



$$\begin{aligned} P(d_2, d_3, \text{data}) &= \sum_{d_1} P(d_1, d_2, d_3, \text{data}) \\ &= \sum_{d_1} \psi_{1,2}(d_1, d_2) \cdot \psi_{2,3}(d_2, d_3) \end{aligned}$$

66

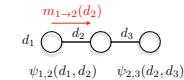
## Marginalization and messages



$$\begin{aligned} P(d_2, d_3, \text{data}) &= \sum_{d_1} P(d_1, d_2, d_3, \text{data}) \\ &= \sum_{d_1} \psi_{1,2}(d_1, d_2) \cdot \psi_{2,3}(d_2, d_3) \\ &= \left[ \sum_{d_1} \psi_{1,2}(d_1, d_2) \right] \cdot \psi_{2,3}(d_2, d_3) \end{aligned}$$

67

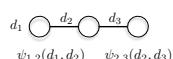
## Marginalization and messages



$$\begin{aligned} P(d_2, d_3, \text{data}) &= \sum_{d_1} P(d_1, d_2, d_3, \text{data}) \\ &= \sum_{d_1} \psi_{1,2}(d_1, d_2) \cdot \psi_{2,3}(d_2, d_3) \\ &= \left[ \sum_{d_1} \psi_{1,2}(d_1, d_2) \right] \cdot \psi_{2,3}(d_2, d_3) \\ &= m_{1 \rightarrow 2}(d_2) \cdot \psi_{2,3}(d_2, d_3) \end{aligned}$$

68

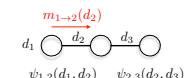
## Marginalization and messages



$$P(d_3, \text{data}) = \sum_{d_2} P(d_2, d_3, \text{data})$$

69

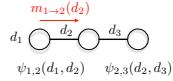
## Marginalization and messages



$$\begin{aligned} P(d_3, \text{data}) &= \sum_{d_2} P(d_2, d_3, \text{data}) \\ &= \sum_{d_2} m_{1 \rightarrow 2}(d_2) \cdot \psi_{2,3}(d_2, d_3) \end{aligned}$$

70

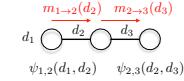
## Marginalization and messages



$$\begin{aligned} P(d_3, \text{data}) &= \sum_{d_2} P(d_2, d_3, \text{data}) \\ &= \sum_{d_2} m_{1 \rightarrow 2}(d_2) \cdot \psi_{23}(d_2, d_3) \\ &= \left[ \sum_{d_2} m_{1 \rightarrow 2}(d_2) \cdot \psi_{23}(d_2, d_3) \right] \end{aligned}$$

71

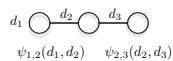
## Marginalization and messages



$$\begin{aligned} P(d_3, \text{data}) &= \sum_{d_2} P(d_2, d_3, \text{data}) \\ &= \sum_{d_2} m_{1 \rightarrow 2}(d_2) \cdot \psi_{23}(d_2, d_3) \\ &= \left[ \sum_{d_2} m_{1 \rightarrow 2}(d_2) \cdot \psi_{23}(d_2, d_3) \right] \\ &= m_{2 \rightarrow 3}(d_3) \end{aligned}$$

72

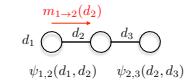
## Marginalization and messages



$$P(d_2, \text{data}) = \sum_{d_3} P(d_2, d_3, \text{data})$$

73

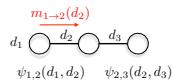
## Marginalization and messages



$$\begin{aligned} P(d_2, \text{data}) &= \sum_{d_3} P(d_2, d_3, \text{data}) \\ &= \sum_{d_3} m_{1 \rightarrow 2}(d_2) \cdot \psi_{23}(d_2, d_3) \end{aligned}$$

74

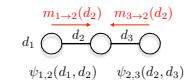
## Marginalization and messages



$$\begin{aligned} P(d_2, \text{data}) &= \sum_{d_3} P(d_2, d_3, \text{data}) \\ &= \sum_{d_3} m_{1 \rightarrow 2}(d_2) \cdot \psi_{23}(d_2, d_3) \\ &= m_{1 \rightarrow 2}(d_2) \cdot \left[ \sum_{d_3} \psi_{23}(d_2, d_3) \right] \end{aligned}$$

75

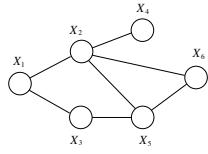
## Marginalization and messages



$$\begin{aligned} P(d_2, \text{data}) &= \sum_{d_3} P(d_2, d_3, \text{data}) \\ &= \sum_{d_3} m_{1 \rightarrow 2}(d_2) \cdot \psi_{23}(d_2, d_3) \\ &= m_{1 \rightarrow 2}(d_2) \cdot \left[ \sum_{d_3} \psi_{23}(d_2, d_3) \right] \\ &= m_{1 \rightarrow 2}(d_2) \cdot m_{3 \rightarrow 2}(d_2) \end{aligned}$$

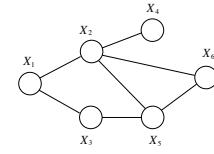
76

## Marginalization



77

## Marginalization

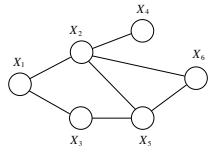


$$p(x_1) = \sum_{x_2} \sum_{x_3} \sum_{x_5} \sum_{x_5} \sum_{x_6} p(x_1, x_2, x_3, x_4, x_5, x_6)$$

$$\propto \sum_{x_2} \sum_{x_3} \sum_{x_5} \sum_{x_6} \psi(x_1, x_2) \psi(x_1, x_3) \psi(x_2, x_4) \psi(x_3, x_5) \psi(x_2, x_5, x_6)$$

78

## Marginalization



$$p(x_1) = \sum_{x_2} \sum_{x_3} \sum_{x_5} \sum_{x_6} p(x_1, x_2, x_3, x_4, x_5, x_6)$$

$$\propto \sum_{x_2} \sum_{x_3} \sum_{x_5} \sum_{x_6} \psi(x_1, x_2) \psi(x_1, x_3) \psi(x_2, x_4) \psi(x_3, x_5) \psi(x_2, x_5, x_6)$$

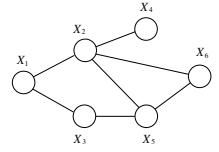
Nominally, each sum is applied to 6 variables  $\Rightarrow \mathcal{O}(r^6)$

79

## Marginalization

- We can reduce the complexity by exploiting the distribution's structure to be smart about the marginalization order

$$p(x_1) = \frac{1}{Z} \sum_{x_2} \psi(x_1, x_2) \sum_{x_3} \psi(x_1, x_3) \sum_{x_4} \psi(x_2, x_4) \sum_{x_5} \psi(x_3, x_5) \sum_{x_6} \psi(x_2, x_5, x_6)$$



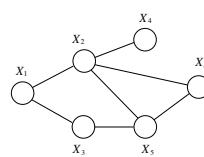
80

## Marginalization

- We can reduce the complexity by exploiting the distribution's structure to be smart about the marginalization order

$$p(x_1) = \frac{1}{Z} \sum_{x_2} \psi(x_1, x_2) \sum_{x_3} \psi(x_1, x_3) \sum_{x_4} \psi(x_2, x_4) \sum_{x_5} \psi(x_3, x_5) \sum_{x_6} \psi(x_2, x_5, x_6)$$

$$= \frac{1}{Z} \sum_{x_2} \psi(x_1, x_2) \sum_{x_3} \psi(x_1, x_3) \sum_{x_4} \psi(x_2, x_4) \sum_{x_5} \psi(x_3, x_5) m_6(x_2, x_5)$$



81

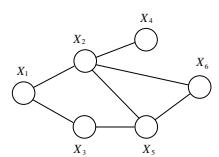
## Marginalization

- We can reduce the complexity by exploiting the distribution's structure to be smart about the marginalization order

$$p(x_1) = \frac{1}{Z} \sum_{x_2} \psi(x_1, x_2) \sum_{x_3} \psi(x_1, x_3) \sum_{x_4} \psi(x_2, x_4) \sum_{x_5} \psi(x_3, x_5) \sum_{x_6} \psi(x_2, x_5, x_6)$$

$$= \frac{1}{Z} \sum_{x_2} \psi(x_1, x_2) \sum_{x_3} \psi(x_1, x_3) \sum_{x_4} \psi(x_2, x_4) \sum_{x_5} \psi(x_3, x_5) m_6(x_2, x_5)$$

$$= \frac{1}{Z} \sum_{x_2} \psi(x_1, x_2) \sum_{x_3} \psi(x_1, x_3) m_5(x_2, x_3) \sum_{x_4} \psi(x_2, x_4)$$

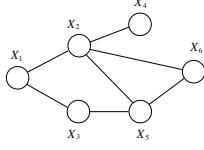


82

## Marginalization

- We can reduce the complexity by exploiting the distribution's structure to be smart about the marginalization order

$$\begin{aligned}
 p(x_1) &= \frac{1}{Z} \sum_{x_2} \psi(x_1, x_2) \sum_{x_3} \psi(x_1, x_3) \sum_{x_4} \psi(x_2, x_4) \sum_{x_5} \psi(x_3, x_5) \sum_{x_6} \psi(x_2, x_5, x_6) \\
 &= \frac{1}{Z} \sum_{x_2} \psi(x_1, x_2) \sum_{x_3} \psi(x_1, x_3) \sum_{x_4} \psi(x_2, x_4) \sum_{x_5} \psi(x_3, x_5) m_6(x_2, x_5) \\
 &= \frac{1}{Z} \sum_{x_2} \psi(x_1, x_2) \sum_{x_3} \psi(x_1, x_3) m_5(x_2, x_3) \sum_{x_4} \psi(x_2, x_4) \\
 &= \frac{1}{Z} \sum_{x_2} \psi(x_1, x_2) m_4(x_2) \sum_{x_3} \psi(x_1, x_3) m_5(x_2, x_3)
 \end{aligned}$$

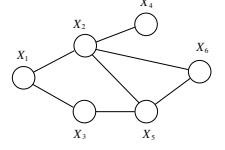


83

## Marginalization

- We can reduce the complexity by exploiting the distribution's structure to be smart about the marginalization order

$$\begin{aligned}
 p(x_1) &= \frac{1}{Z} \sum_{x_2} \psi(x_1, x_2) \sum_{x_3} \psi(x_1, x_3) \sum_{x_4} \psi(x_2, x_4) \sum_{x_5} \psi(x_3, x_5) \sum_{x_6} \psi(x_2, x_5, x_6) \\
 &= \frac{1}{Z} \sum_{x_2} \psi(x_1, x_2) \sum_{x_3} \psi(x_1, x_3) \sum_{x_4} \psi(x_2, x_4) \sum_{x_5} \psi(x_3, x_5) m_6(x_2, x_5) \\
 &= \frac{1}{Z} \sum_{x_2} \psi(x_1, x_2) \sum_{x_3} \psi(x_1, x_3) m_5(x_2, x_3) \sum_{x_4} \psi(x_2, x_4) \\
 &= \frac{1}{Z} \sum_{x_2} \psi(x_1, x_2) m_4(x_2) \sum_{x_3} \psi(x_1, x_3) m_5(x_2, x_3) \\
 &= \frac{1}{Z} \sum_{x_2} \psi(x_1, x_2) m_4(x_2) m_3(x_1, x_2)
 \end{aligned}$$



84

## Marginalization

- We can reduce the complexity by exploiting the distribution's structure to be smart about the marginalization order

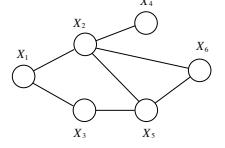
$$\begin{aligned}
 p(x_1) &= \frac{1}{Z} \sum_{x_2} \psi(x_1, x_2) \sum_{x_3} \psi(x_1, x_3) \sum_{x_4} \psi(x_2, x_4) \sum_{x_5} \psi(x_3, x_5) \boxed{\sum_{x_6} \psi(x_2, x_5, x_6)} \quad r \times r \times r \text{ matrix} \\
 &= \frac{1}{Z} \sum_{x_2} \psi(x_1, x_2) \sum_{x_3} \psi(x_1, x_3) \sum_{x_4} \psi(x_2, x_4) \sum_{x_5} \psi(x_3, x_5) m_6(x_2, x_5) \\
 &= \frac{1}{Z} \sum_{x_2} \psi(x_1, x_2) \sum_{x_3} \psi(x_1, x_3) m_5(x_2, x_3) \sum_{x_4} \psi(x_2, x_4) \\
 &= \frac{1}{Z} \sum_{x_2} \psi(x_1, x_2) m_4(x_2) \sum_{x_3} \psi(x_1, x_3) m_5(x_2, x_3) \\
 &= \frac{1}{Z} \sum_{x_2} \psi(x_1, x_2) m_4(x_2) m_3(x_1, x_2) \\
 &= \frac{1}{Z} m_2(x_1),
 \end{aligned}$$

85

## Marginalization

- We can reduce the complexity by exploiting the distribution's structure to be smart about the marginalization order

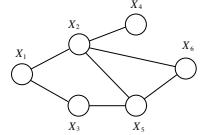
$$\begin{aligned}
 p(x_1) &= \frac{1}{Z} \sum_{x_2} \psi(x_1, x_2) \sum_{x_3} \psi(x_1, x_3) \sum_{x_4} \psi(x_2, x_4) \sum_{x_5} \psi(x_3, x_5) \boxed{\sum_{x_6} \psi(x_2, x_5, x_6)} \quad r \times r \times r \text{ matrix} \\
 &= \frac{1}{Z} \sum_{x_2} \psi(x_1, x_2) \sum_{x_3} \psi(x_1, x_3) \sum_{x_4} \psi(x_2, x_4) \sum_{x_5} \psi(x_3, x_5) m_6(x_2, x_5) \\
 &= \frac{1}{Z} \sum_{x_2} \psi(x_1, x_2) \sum_{x_3} \psi(x_1, x_3) m_5(x_2, x_3) \sum_{x_4} \psi(x_2, x_4) \\
 \mathcal{O}(r^3) &= \frac{1}{Z} \sum_{x_2} \psi(x_1, x_2) m_4(x_2) \sum_{x_3} \psi(x_1, x_3) m_5(x_2, x_3) \\
 &= \frac{1}{Z} \sum_{x_2} \psi(x_1, x_2) m_4(x_2) m_3(x_1, x_2) \\
 &= \frac{1}{Z} m_2(x_1),
 \end{aligned}$$



85

## Elimination order

- Summations (eliminations) introduce intermediate edges

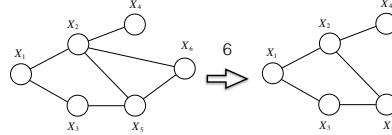


$$p(x_1) = \frac{1}{Z} \sum_{x_2} \psi(x_1, x_2) \sum_{x_3} \psi(x_1, x_3) \sum_{x_4} \psi(x_2, x_4) \sum_{x_5} \psi(x_3, x_5) \sum_{x_6} \psi(x_2, x_5, x_6)$$

86

## Elimination order

- Summations (eliminations) introduce intermediate edges

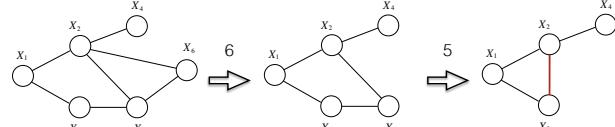


$$\begin{aligned}
 p(x_1) &= \frac{1}{Z} \sum_{x_2} \psi(x_1, x_2) \sum_{x_3} \psi(x_1, x_3) \sum_{x_4} \psi(x_2, x_4) \sum_{x_5} \psi(x_3, x_5) \sum_{x_6} \psi(x_2, x_5, x_6) \\
 &= \frac{1}{Z} \sum_{x_2} \psi(x_1, x_2) \sum_{x_3} \psi(x_1, x_3) \sum_{x_4} \psi(x_2, x_4) \sum_{x_5} \psi(x_3, x_5) m_6(x_2, x_5)
 \end{aligned}$$

87

## Elimination order

- Summations (eliminations) introduce intermediate edges

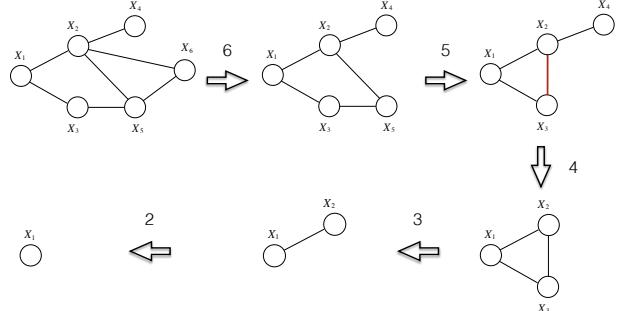


$$\begin{aligned}
 p(x_1) &= \frac{1}{Z} \sum_{x_2} \psi(x_1, x_2) \sum_{x_3} \psi(x_1, x_3) \sum_{x_4} \psi(x_2, x_4) \sum_{x_5} \psi(x_3, x_5) \sum_{x_6} \psi(x_2, x_5, x_6) \\
 &= \frac{1}{Z} \sum_{x_2} \psi(x_1, x_2) \sum_{x_3} \psi(x_1, x_3) \sum_{x_4} \psi(x_2, x_4) \sum_{x_5} \psi(x_3, x_5) m_6(x_2, x_5) \\
 &= \frac{1}{Z} \sum_{x_2} \psi(x_1, x_2) \boxed{\sum_{x_3} \psi(x_1, x_3) m_5(x_2, x_3)} \sum_{x_4} \psi(x_2, x_4)
 \end{aligned}$$

88

## Elimination order

- Summations (eliminations) introduce intermediate edges



89

## Elimination order

- Worst-case summation complexity is a function of the largest intermediate clique (e.g., 3 element clique)
- Goal: Choose summation order (*elimination order*) with smallest intermediate clique
  - e.g., we went from  $\mathcal{O}(r^6)$  to  $\mathcal{O}(r^3)$

**Definition:** The *treewidth* is the minimum over elimination order of the maximal clique size

- Note parallels with solving set of linear equations

90

## Elimination order

- Worst-case summation complexity is a function of the largest intermediate clique (e.g., 3 element clique)

- Goal: Find elimination order that achieves treewidth  $\leq k$
  - e.g., we went from  $\mathcal{O}(r^6)$  to  $\mathcal{O}(r^3)$
- Finding elimination order that achieves treewidth is NP-hard**

**Definition:** The *treewidth* is the minimum over elimination order of the maximal clique size

- Note parallels with solving set of linear equations

90