### Scalable Uncertainty Management 06 – Markov Logic

Rainer Gemulla

July 13, 2012

### Overview

In this lecture

- Statistical relational learning (SRL)
- Introduction to probabilistic graphical models (PGM)
- Basics of undirected models (called Markov networks)
- Markov logic as a template for undirected models
- Basics of inference in Markov logic networks

Not in this lecture

- Directed models (called Bayesian networks)
- Other SRL approaches (such as probabilistic relational models)
- High coverage and in-depth discussion of inference
- Learning Markov logic networks

# Outline



- Probabilistic Graphical Models
  - Introduction
  - Preliminaries
- 3 Markov Networks
- 4 Markov Logic Networks
  - Grounding Markov logic networks
  - Log-Linear Models

### 5 Inference in MLNs

- Basics
- Exact Inference
- Approximate Inference

### 5 Summary

# Correlations in probabilistic databases

- Simple probabilistic models
  - Tuple-independent databases
  - Block-disjoint independent databases
  - Key/foreign key constraints, ...
- $\bullet$  Correlations (mainly) through  $\mathcal{RA}$  queries/views
  - Any discrete probability distribution can be modeled
  - Queries describe *precisely* how result is derived

E	Example (Nell)										
	NellExtraction NellSource										
	Subject	Pattern	Object	Source	$\mathbb{P}$		Source	$\mathbb{P}$			
	Sony	produces	Walkman	1	0.96	1	1	0.99			
	IBM	produces	PC	1	0.96	ĺ	2	0.1			
	IBM	produces	PC	2	1				_		
	Microsoft	produces	MacOS	2	0.9						
	AlbertEinstein bornIn Ulm 1 0.9 Produces										
Subject Object $\mathbb{P}$							$\mathbb{P}$				
	$Produces(x, y) \leftarrow NellExtraction(x, 'produces', y, s),$						Sony	W	/alkman	0.9504	1

IRM

Microsoft

PC

MacOS

0.95536

0.09

NellSource(s)

# Statistical relational learning (I)

#### **Smoking and Quitting in Groups**

Researchers studying a network of 12,067 people found that smokers and nonsmokers tended to cluster in groups of close friends and family members. As more people quit over the decades, remaining groups of smokers were increasingly pushed to the periphery of the social network.



Does John smoke?

Learn correlations from structured data, then apply to new data.

# Statistical relational learning (II)

- Goal: Declarative modelling of correlations in structured data
- Idea: Use (subsets of) first-order logic
  - Very expressive formalism; lots of knowledge bases use it
  - Symmetry:  $\forall x. \forall y. Friends(x, y) \iff Friends(y, x)$
  - Everybody has a friend:  $\forall x. \exists y. Friends(x, y)$
  - ▶ Transitivity:  $\forall x.\forall y.\forall z.Friends(x, y) \land Friends(y, z) \implies Friends(x, z)$
  - Smoking causes cancer: ∀x.Smokes(x) ⇒ Cancer(x)
  - Friends have similar smoking habits:  $\forall x. \forall y. Friends(x, y) \implies (Smokes(x) \iff Smokes(y))$
- Problem: Real-world knowledge is incomplete, contradictory, complex
  - $\rightarrow$  Above rules do not generally hold, but they are "likely" to hold!
- Approach: Combine first-order logic with probability theory
  - Expressiveness of first-order logic
  - Principled treatment of uncertainty using probability theory

There are many approaches of this kind. Our focus is on *Markov logic*, a recent and very successful language.

### Markov logic networks

#### Definition

A Markov logic network is a set of pairs  $(F_i, w_i)$ , where  $F_i$  is a formula in first-order logic and the weight  $w_i$  is a real number.

#### Example

Smoking causes cancer 155

$$5 \downarrow \forall x. Smokes(x) \implies Cancer(x)$$

Friends have similar smoking habits  $\forall x. \forall y. Friends(x, y) \implies (Smokes(x) \iff Smokes(y))$ 1.1 {

- Formulas may or may not hold
- Weights express confidence
  - High positive weight  $\rightarrow$  confident that formula holds
  - High negative weight  $\rightarrow$  confident that formula does not hold
  - But careful: weights actually express confidence of certain "groundings" of a formula and *not* the formula as a whole (more later)
- Formulas may introduce complex correlations

# Simple MLN for entity resolution

#### Which citations refer to the same publication?

author	Richardson, Matt and Domingos, Pedro	M. Richardson and P. Domingos	Domingos, Pedro and Richardson, Matthew
title	Markov Logic Networks	Markov logic networks	Markov Logic: A Unifying Framework for Statistical Relational Learning
year	2006	2006	2007

#### // predicates

HasToken(token, field, citation) SameCitation(citation, citation)

// e.g., HasToken('Logic', 'title ', C1) SameField(field, citation, citation) // Semantic equality of values in a field // Semantic equality of citations

#### // formulas

 $HasToken(+t, +f, c1) \cap HasToken(+t, +f, c2) => SameField(+f, c1, c2)$ SameField(+f, c1, c2) => SameCitation(c1, c2)SameCitation(c1, c2)  $\hat{}$  SameCitation(c2, c3) => SameCitation(c1, c3)

Rule weights are usually learned from data. The same rule may have different weights for different constants (indicated by "+").

### Alchemy

- Alchemy is well-known software package for Markov logic
- Developed at University of Washington
- Supports a wide range of tasks
  - Structure learning
  - Weight learning
  - Probabilistic inference
- Has been used for wide range of applications
  - Information extraction
  - Social network modeling
  - Entity resolution
  - Collective classification
  - Link prediction
- Check out http://alchemy.cs.washington.edu/
  - Code
  - Real-world datasets
  - Real-world Markov logic networks
  - Literature

# From Markov logic to graphical models (example)

Friends						
Name1	Name2	Value				
Anna	Bob	Yes				
Bob	Anna	Yes				
Anna	Anna	Yes				
Bob	Bob	Yes				





### Probabilistic databases and graphical models

	Probabilistic databases	Graphical models	
Probabilistic model	Simple (disjoint-independent tuples)	Complex (independencies given by graph)	
Query	Complex (e.g., $\exists x. \exists y. R(x, y) \land S(x)$ )	Simple (e.g., $\mathbb{P}(X_1, X_2   Z_1, Z_2, Z_3)$ )	
Network	Dynamic (database + query)	Static (Bayesian or Markov network)	
Complexity measured in size of	Database	Network	
Complexity parameter	Query	Treewidth	
System	Extension to RDBMS	Stand-alone	
Hybrid approaches have many potential applications and are under active research.			

# Outline



- Probabilistic Graphical Models
  - Introduction
  - Preliminaries

#### 3 Markov Networks

- 4 Markov Logic Networks
  - Grounding Markov logic networks
  - Log-Linear Models

#### 5 Inference in MLNs

- Basics
- Exact Inference
- Approximate Inference

### 5 Summary

# Outline



Introduction to Markov Logic Networks

- 2 Probabilistic Graphical Models
  - Introduction
  - Preliminaries

### 3 Markov Networks

#### 4 Markov Logic Networks

- Grounding Markov logic networks
- Log-Linear Models

### Inference in MLNs

- Basics
- Exact Inference
- Approximate Inference



### Reasoning with uncertainty

- Goal: Automated reasoning system
  - Take all available information
    - (e.g., patient information: symptoms, test results, personal data)
  - Reach conclusions

(e.g., which diseases the patient has, which medication to give)

- Desiderata
  - Separation of knowledge and reasoning
    - \* Declarative, model-based representation of knowledge
    - ★ General suite of reasoning algorithms, applicable to many domains
  - Principled treatment of uncertainty
    - ★ Partially observed data
    - ★ Noisy observations
    - ★ Non-deterministic relationships
- Lots of applications
  - medical diagnosis, fault diagnosis, analysis of genetic and genomic data, communication and coding, analysis of marketing data, speech recognition, *natural language understanding*, segmenting and denoising images, social network analysis, ...

### Probabilistic models

- Multiple interrelated aspects may relate to the reasoning task
  - Possible diseases
  - Hundreds of symptoms and diagnostic tests
  - Personal characteristics
- Characterize data by a set of random variables
  - Flu (yes / no)
  - Hayfever (yes / no)
  - Season (Spring / Sommer / Autumn / Winter)
  - Congestion (yes / no)
  - MusclePain (yes / no)
  - $\rightarrow$  Variables and their domain are important design decision
- Ø Model dependencies by a joint distribution
  - Diseases, season, and symptoms are correlated
  - ▶ Probabilistic models construct joint probability space → 2 · 2 · 4 · 2 · 2 outcomes (64 values, 63 non-redundant)
  - Given joint probability space, interesting questions can be answered

 $\mathbb{P}(\mathsf{Flu} | \mathsf{Season} = \mathsf{Spring}, \mathsf{Congestion}, \neg\mathsf{MusclePain})$ 

#### Specifying a joint distribution is infeasible in general!

# Probabilistic graphical models

- A graph-based representation of *direct* probabilistic interactions
- A break-down of high-dimensional distributions into smaller *factors* (here: 63 vs. 17 non-redundant parameters)
- A compact representation of a set of (conditional) independencies



### Main components

- Representation
  - Tractability
    - \* Variables tend to interact *directly* only with very few others
    - Natural and compact encoding as graphical model
  - Transparency
    - $\star\,$  Models can be understood/evaluated by human experts
- Inference
  - Answer queries using the distribution as model of the world
  - Work on graph structure
    - $\rightarrow$  orders of magnitude faster than working on joint probability
- Learning
  - Learn a model from data that captures past experience to a good approximation
  - Human experts may provide rough guidance
  - ▶ Details filled in by fitting the model to the data → Often better reflection of domain than hand-constructed models, sometimes surprising insights

Graphical models exploit locality structure that appears in many distributions that arise in practice.

# Outline



- Probabilistic Graphical Models
  - Introduction
  - Preliminaries

### 3 Markov Networks

#### 4 Markov Logic Networks

- Grounding Markov logic networks
- Log-Linear Models

### Inference in MLNs

- Basics
- Exact Inference
- Approximate Inference



### Notation

Let **X** and **Y** be sets of random variables with domain Dom(X) and Dom(Y). Let  $x \in Dom(X)$  and  $y \in Dom(Y)$ .

Expression	Shortcut notation
$\mathbb{P}\left( \left. \mathbf{X}=\mathbf{x} \right.  ight)$	$\mathbb{P}(\mathbf{x})$
$\mathbb{P}\left( \left. X=x \mid Y=y  ight.  ight)$	$\mathbb{P}(\mathbf{x} \mid \mathbf{y})$
$orall \mathbf{x}. \ \mathbb{P}\left( \ \mathbf{X} = \mathbf{x} \  ight) = f(\mathbf{x})$	$\mathbb{P}\left( \left. \mathbf{X} \right.  ight) = f(\mathbf{X})$
$orall \mathbf{x}.orall \mathbf{y}. \ \mathbb{P}\left( \mathbf{X} = \mathbf{x} \mid \mathbf{Y} = \mathbf{y}  ight) = f(\mathbf{x}, \mathbf{y})$	$\mathbb{P}\left(\left. \mathbf{X} \mid \mathbf{Y} \right.  ight) = f(\mathbf{X},\mathbf{Y})$

- $\mathbb{P}(X)$  and  $\mathbb{P}(X \mid Y)$  are entire probability distributions
- Can be thought of as functions from  $\mathsf{Dom}(X) \to [0,1]$  or  $(\mathsf{Dom}(X),\mathsf{Dom}(Y)) \to [0,1]$ , respectively
- f<sub>y</sub>(X) = P(X | y) is often referred to as conditional probability distribution (CPD)
- For discrete variables, may be represented as a table (CPT)

# Conditional independence

Definition

Let X, Y and Z be sets of random variables. X and Y are said to be *conditionally independent* given Z if and only if

 $\mathbb{P}\left(\left. \mathsf{X},\mathsf{Y} \mid \mathsf{Z}\right.\right) = \mathbb{P}\left(\left. \mathsf{X} \mid \mathsf{Z}\right.\right) \mathbb{P}\left(\left. \mathsf{Y} \mid \mathsf{Z}\right.\right).$ 

We write  $(\mathbf{X} \perp \mathbf{Y} \mid \mathbf{Z})$  for this conditional independence statement. If  $\mathbf{Z} = \emptyset$ , we write  $(\mathbf{X} \perp \mathbf{Y})$  for marginal independence.

### Example



 $(F \perp H \mid S), (C \perp S \mid F, H)$  $(M \perp H, C, S \mid F)$ 

$$\mathbb{P}(S, F, H, M, C)$$
  
=  $\mathbb{P}(S) \cdot \mathbb{P}(F | S) \cdot \mathbb{P}(H | S)$   
 $\cdot \mathbb{P}(C | F, H) \cdot \mathbb{P}(M | F)$ 

# Properties of conditional independence

Theorem

In general,  $(X \perp Y)$  does not imply nor is implied by  $(X \perp Y \mid Z)$ 

The following relationships hold:

$$\begin{array}{ccc} (\mathsf{X} \perp \mathsf{Y} \mid \mathsf{Z}) \iff (\mathsf{Y} \perp \mathsf{X} \mid \mathsf{Z}) & (symmetry) \\ (\mathsf{X} \perp \mathsf{Y}, \mathsf{W} \mid \mathsf{Z}) \implies (\mathsf{X} \perp \mathsf{Y} \mid \mathsf{Z}) & (decomposition) \\ (\mathsf{X} \perp \mathsf{Y}, \mathsf{W} \mid \mathsf{Z}) \implies (\mathsf{X} \perp \mathsf{Y} \mid \mathsf{Z}, \mathsf{W}) & (weak \ union) \\ \mathsf{X} \perp \mathsf{W} \mid \mathsf{Z}, \mathsf{Y}) \land (\mathsf{X} \perp \mathsf{Y} \mid \mathsf{Z}) \implies (\mathsf{X} \perp \mathsf{Y}, \mathsf{W} \mid \mathsf{Z}) & (contraction) \end{array}$$

For positive distributions and mutally disjoint sets X, Y, Z, W:

 $(\textbf{X} \perp \textbf{Y} \mid \textbf{Z}, \textbf{W}) \land (\textbf{X} \perp \textbf{W} \mid \textbf{Z}, \textbf{Y}) \implies (\textbf{X} \perp \textbf{Y}, \textbf{W} \mid \textbf{Z}) \quad (\textit{intersection})$ 

#### Proof.

Discussed in exercise group.

# Querying a distribution (1)

Consider a joint distribution on a set of variables  $\ensuremath{\mathcal{X}}$ 

- $\bullet~$  Let  $\textbf{E}\subseteq \mathcal{X}$  be a set of *evidence variables* that takes values e
- Let  $\mathbf{W} = \mathcal{X} \mathbf{E}$  be the set of *latent variables*
- Let  $\mathbf{Y} \subseteq \mathbf{W}$  be a set of *query variables*
- Let  $\mathbf{Z} = \mathbf{W} \mathbf{Y}$  be the set of *non-query variables*

### Example

- $\bullet \ \mathcal{X} = \{ \, \mathsf{Season}, \mathsf{Congestion}, \mathsf{MusclePain}, \mathsf{Flu}, \mathsf{Hayfever} \, \}$
- $E = \{ Season, Congestion, MusclePain \}$
- $\mathbf{e} = \{ \text{ Spring}, \text{Yes}, \text{No} \}$
- $\mathbf{W} = \{ \mathsf{Flu}, \mathsf{Hayfever} \}$
- $\mathbf{Y} = \{ \mathsf{Flu} \}$
- $\mathbf{Z} = \{ Hay fever \}$

# Querying a distribution (2)

- Conditional probability query
  - Compute the *posterior distribution* of the query variables  $\mathbb{P}(\mathbf{Y} \mid \mathbf{e})$
- MAP query
  - Compute the most likely value of the latent variables MAP(W | e) = argmax<sub>w</sub> P(w | e) = argmax<sub>w</sub> P(w, e)
- Marginal MAP query
  - Compute the most likely value of the query variables  $MAP(\mathbf{Y} \mid \mathbf{e}) = argmax_{\mathbf{y}} \mathbb{P}(\mathbf{y} \mid \mathbf{e}) = argmax_{\mathbf{y}} \sum_{\mathbf{z}} \mathbb{P}(\mathbf{y}, \mathbf{z}, \mathbf{e})$

### Example

$\mathbb{P}(\mathbf{W} \mid \mathbf{e})$	Flu	¬Flu
Hayfever	5%	35%
$\neg$ Hayfever	40%	20%

- $\label{eq:prince} { \ensuremath{\mathbb Q}} \ { \ensuremath{\mathbb P}} \ ( \ { \ensuremath{\mathsf{Flu}}} \ | \ { \ensuremath{\mathsf{Spring}}}, { \ensuremath{\mathsf{Congestion}}}, \neg { \ensuremath{\mathsf{MusclePain}}} \ ) \ { \ensuremath{\mathsf{Yes}}} \ ( \ensuremath{\mathsf{45\%}} \ ), \ { \ensuremath{\mathsf{No}}} \ ( \ensuremath{\mathsf{55\%}} \ ) \ ) \$
- MAP(Flu | Spring, Congestion,  $\neg$ MusclePain)  $\rightarrow$  No flu (!)

# Querying graphical models

- Graphical models induce conditional independences
- Queries reason about dependencies between variables

Can we evaluate queries more efficiently given a graphical model and its associated independences?

#### Example

Independence properties help inference!



Table known to satisfy (F $\perp$ H   <b>E</b>							
$\mathbb{P}(\mathbf{W} \mid \mathbf{e})$	Flu ¬Flu						
Hayfever	24%	16%	40%				
$\neg Hayfever$	36%	24%	60%				
	60%	40%					

Thus, for example, monotonicity is now known to hold for MAP:  $MAP(Flu, Hayfever | \mathbf{E}) = (MAP(Flu | \mathbf{E}), MAP(Hayfever | \mathbf{E}))$ 

# Outline

Introduction to Markov Logic Networks

- Probabilistic Graphical Models
  - Introduction
  - Preliminaries

### 3 Markov Networks

- 4 Markov Logic Networks
  - Grounding Markov logic networks
  - Log-Linear Models

### 5 Inference in MLNs

- Basics
- Exact Inference
- Approximate Inference

### 5 Summary

### Misconception example

#### Example

• Alice, Bob, Charles, and Debbie study in pairs for the SUM exam



- Lecturer misspoke in class, giving rise to a possible misconception
- Some students figured out the problem, others did not

Which of the students has the misconception?

- If A does not have the misconception, he may help B and D  $\rightarrow$  Students influence each other
- If A has the misconception, he may be helped by B and D  $\rightarrow$  Influence has no natural "direction"
- A does not study with  $C \rightarrow No$  direct influence between A and C

### Markov network

### Definition

A *Markov network* is an undirected graph  $\mathcal{H} = (\mathcal{X}, \mathcal{E})$ , where  $\mathcal{X}$  is a set of random variables and  $\mathcal{E} \subseteq \mathcal{X} \times \mathcal{X}$  is the set of edges.



We will see that Markov networks encode a set of conditional independence assumptions between its variables.

### Local models

#### Definition

Let **D** be a set of random variables. A *factor*  $\phi$  is a function from  $\text{Dom}(\mathbf{D}) \to \mathbb{R}$ . A factor is *nonnegative* if has range  $\mathbb{R}^+$ . The set **D** is called the *scope* of the factor and is denoted  $\text{Scope}[\phi]$ .

We restrict attention to nonnegative factors.

#### Example



• Factors describe "compatibility" between values (not normalized)

- $\phi_1$ : More "weight" when A and B agree than when they disagree
- $\phi_1$ : More weight when A and B are both right than when both are wrong
- $\phi_1$ : If they disagree, more weight when A is right than when B is right

В

С

D

## Combining local models

#### Definition

Let X, Y, Z be three disjoint sets of random variables and let  $\phi_1(X, Y)$ and  $\phi_2(Y, Z)$  be two factors. The *factor product*  $\psi = \phi_1 \times \phi_2$  is given by the factor  $\psi : \text{Dom}(X, Y, Z) \to \mathbb{R}$  with

$$\psi(\mathsf{X},\mathsf{Y},\mathsf{Z})=\phi_1(\mathsf{X},\mathsf{Y})\cdot\phi_2(\mathsf{Y},\mathsf{Z}).$$

#### Example

$\begin{array}{c ccccccccccccccccccccccccccccccccccc$	$\begin{array}{c ccccccccccccccccccccccccccccccccccc$
---	---

# Factor products and the product rule of probability

Recall the product rule of probability

$$\mathbb{P}(\mathbf{X}, \mathbf{Y}) = \mathbb{P}(\mathbf{Y}) \mathbb{P}(\mathbf{X} \mid \mathbf{Y}).$$



MusclePain M P Yes 0.1 No 0.9 Flu | MusclePain

the second se						
M	F	$\mathbb{P}$				
'es	Yes	0.8				
'es	No	0.2				
lo	Yes	0.1				
lo	No	0.9				

 Flu,
 MusclePain

 M
 F
 P

 Yes
 Yes
 0.08

 Yes
 No
 0.02

 No
 Yes
 0.09

 No
 No
 0.81

- Set  $\phi_1(\mathsf{MusclePain}) = \mathbb{P}(\mathsf{MusclePain})$
- Set  $\phi_2(MusclePain, Flu) = \mathbb{P}(Flu | MusclePain)$
- Set  $\psi(\mathsf{MusclePain},\mathsf{Flu}) = \mathbb{P}(\mathsf{MusclePain},\mathsf{Flu})$
- Then  $\psi = \phi_1 \times \phi_2$

Factor products generalize the product rule of probability.

### Gibbs distribution

#### Definition

A distribution  $\mathbb{P}_{\Phi}$  is a *Gibbs distribution* parameterized by a set of factors  $\Phi = \{ \phi_1(\mathbf{D}_1), \dots, \phi_m(\mathbf{D}_m) \}$  if it is defined by

$$\mathbb{P}_{\Phi}(X_{1},\ldots,X_{n}) = \frac{1}{Z} \,\tilde{\mathbb{P}}_{\Phi}(X_{1},\ldots,X_{n})$$
$$\tilde{\mathbb{P}}_{\Phi}(X_{1},\ldots,X_{n}) = \phi_{1}(\mathbf{D}_{1}) \times \phi_{2}(\mathbf{D}_{2}) \times \cdots \times \phi_{m}(\mathbf{D}_{m})$$
$$Z = \sum_{X_{1},\ldots,X_{n}} \tilde{\mathbb{P}}_{\Phi}(X_{1},\ldots,X_{n})$$

Here,  $\tilde{\mathbb{P}}_{\Phi}(X_1, \ldots, X_n)$  is an *unnormalized measure* and Z a normalizing constant called the *partitioning function*.

- Factors contribute to the overall joint distribution
- Overall dist. takes into consideration the contribution from all factors

A set of factors defines a Gibbs distribution, i.e., a joint probability distribution over all variables.

### Gibbs distribution for Misconception example

$\frown$	A	В	С	D	P	$\mathbb{P}$
(A)		$b_0$	<i>c</i> <sub>0</sub>	$d_0$	300,000	0.04
$\overrightarrow{D}$ $\overrightarrow{B}$	<i>a</i> 0	$b_0$	<i>c</i> <sub>0</sub>	$d_1$	300,000	0.04
	<i>a</i> 0	$b_0$	<i>c</i> <sub>1</sub>	$d_0$	300,000	0.04
(C)	<i>a</i> 0	$b_0$	<i>c</i> 1	$d_1$	30	$4.1 \cdot 10^{-6}$
	$a_0$	$b_1$	<i>c</i> <sub>0</sub>	$d_0$	500	$6.9\cdot10^{-5}$
$\begin{array}{c cccc} A & B & \phi_1 \\ \hline 0 & (0 & c_2) \\ \hline \end{array} \qquad \begin{array}{c ccccc} B & C & \phi_2 \\ \hline 0 & (0 & c_2) \\ \hline \end{array}$	$a_0$	$b_1$	<i>c</i> <sub>0</sub>	$d_1$	500	$6.9\cdot10^{-5}$
$a^{0} b^{0} 30 = b^{0} c^{0} 100$	<i>a</i> 0	$b_1$	<i>c</i> 1	$d_0$	5,000,000	0.69
$\begin{vmatrix} a^0 & b^1 & 5 \\ 1 & 0 & c^1 & 1 \end{vmatrix}$	<i>a</i> 0	$b_1$	<i>c</i> <sub>1</sub>	$d_1$	500	$6.9\cdot10^{-5}$
$\begin{vmatrix} a^{1} & b^{0} & 1 \\ 1 & c^{1} & c^{0} & 1 \end{vmatrix}$	$a_1$	$b_0$	<i>c</i> <sub>0</sub>	$d_0$	100	$1.4 \cdot 10^{-5}$
$a^{1} b^{1}   10 $ $b^{1}   c^{1}   100$	$a_1$	$b_0$	<i>c</i> 0	$d_1$	1,000,000	0.14
$C D \phi_{0} D A \phi_{1}$	$a_1$	$b_0$	<i>c</i> <sub>1</sub>	$d_0$	100	$1.4 \cdot 10^{-5}$
$c^{0} d^{0} 1$ $d^{0} a^{0} 100$	$a_1$	$b_0$	<i>c</i> 1	$d_1$	100	$1.4 \cdot 10^{-5}$
$c^{0} d^{1} 100 d^{0} a^{1} 1$	$a_1$	$b_1$	<i>c</i> <sub>0</sub>	$d_0$	10	$1.4 \cdot 10^{-6}$
$c^{1} d^{0} 100 d^{1} a^{0} 1$	$a_1$	$b_1$	<i>c</i> <sub>0</sub>	$d_1$	100,000	0.014
$c^{1} d^{1} 1$ $d^{1} a^{1} 100$	$a_1$	$b_1$	$c_1$	$d_0$	100,000	0.014
	$a_1$	$b_1$	<i>c</i> <sub>1</sub>	$d_1$	100,000	0.014
				7 =	= 7 201 840	3

32 / 78

# Factorization and factor graphs

### Definition

A distribution  $\mathbb{P}_{\Phi}$  with  $\Phi = \{ \phi_1(\mathbf{D}_1), \dots, \phi_m(\mathbf{D}_m) \}$  factorizes over a Markov network  $\mathcal{H}$  if each  $\mathbf{D}_i$  is a complete subgraph of  $\mathcal{H}$ . The factors  $\phi_i$  are often called *clique potentials*.



### Active paths

### Definition

Let  $X_1 - \ldots - X_k$  be a path in  $\mathcal{H} = (\mathcal{X}, \mathcal{E})$ . Let  $\mathbf{Z} \subseteq \mathcal{X}$  be a set of observed variables. The path  $X_1 - \ldots - X_k$  is *active* given  $\mathbf{Z}$  if  $X_i \notin \mathbf{Z}$  for  $1 \leq i \leq k$ .

#### Example



All active paths given A:

- *D*-*C*
- *C*-*B*
- *D*-*C*-*B*

Some inactive paths given A:

- *D*–*A*–*B*
- *C*-*D*-*A*-*B*

# Separation and independencies for Markov networks

### Definition

We say that a set of nodes Z separates X and Y in  $\mathcal{H}$ , denoted  $\sup_{\mathcal{H}} (X; Y \mid Z)$ , if there is no active path between any node in X and any node in Y given Z. We associate with  $\mathcal{H}$  the following set of independencies:

$$\mathcal{I}(\mathcal{H}) = \{ (\mathbf{X} \perp \mathbf{Y} \mid \mathbf{Z}) : \mathsf{sep}_{\mathcal{H}}(\mathbf{X}; \mathbf{Y} \mid \mathbf{Z}) \}$$



- $\emptyset$  does not separate any nodes
- { A } does not separate any nodes
- $\{A, C\}$  separates  $\{B\}$  and  $\{D\}$
- $\{A, B, C\}$  does not separate any nodes

$$\mathcal{I}(\mathcal{H}) = \{ (B \perp D \mid A, C), (D \perp B \mid A, C) \\ (A \perp C \mid B, D), (C \perp A \mid B, D) \}$$

# Relationship Gibbs distributions and Markov networks

#### Definition

- Let P be a probability distribution over X. Define I(P) to be the set of independence assertions of the form (X ⊥ Y | Z) that hold in P.
- A Markov network  $\mathcal{H}$  is an *I-map* for  $\mathbb{P}$  if  $\mathcal{I}(\mathcal{H}) \subseteq \mathcal{I}(\mathbb{P})$ .

#### Theorem

Soundness  $(\rightarrow)$ 

Let  $\mathbb{P}$  be a distribution and  $\mathcal{H}$  be a Markov network over  $\mathcal{X}$ . If  $\mathbb{P}$  is a Gibbs distribution that factorizes over  $\mathcal{H}$ , then  $\mathcal{H}$  is an I-map for  $\mathbb{P}$ .

### Theorem (Hammersley-Clifford theorem)

Soundness ( $\leftarrow$ )

Let  $\mathbb{P}$  be a positive distribution and  $\mathcal{H}$  be a Markov network over  $\mathcal{X}$ . If  $\mathcal{H}$  is an I-map for  $\mathbb{P}$ , then  $\mathbb{P}$  is a Gibbs distribution that factorizes over  $\mathcal{H}$ .

#### Theorem

Completeness

If X and Y are not separated given **Z** in  $\mathcal{H}$ , then X and Y are dependent for some distribution  $\mathbb{P}$  that factorizes over  $\mathcal{H}$ .
## Application: Image denoising



## Application: Stanford Named Entity Recognizer

Named Entity Recognition (NER) labels sequences of words in a text which are the names of things, such as person and company names, or gene and protein names.



- Local evidence often strong clue for label
- Long-range evidence (label consistency) helps when local evidence is insufficient

## Outline



- Probabilistic Graphical Models
  - Introduction
  - Preliminaries

#### 3 Markov Networks

- 4 Markov Logic Networks
  - Grounding Markov logic networks
  - Log-Linear Models

#### 5 Inference in MLNs

- Basics
- Exact Inference
- Approximate Inference

#### 5 Summary

## Outline



Introduction to Markov Logic Networks

- 2 Probabilistic Graphical Models
  - Introduction
  - Preliminaries

#### 3 Markov Networks

4 Markov Logic Networks Grounding Markov logic networks Log-Linear Models

#### Inference in MLNs

- Basics
- Exact Inference
- Approximate Inference



## Semantics of Markov logic networks

#### Definition

A Markov logic network  $L = \{ (F_i, w_i) \}$  is a template for constructing Markov networks. Given a set of constants C, a ground Markov logic  $M_{L,C}$  specifies a distribution over the possible worlds as follows

$$\mathbb{P}\left(\mathbf{X}=\mathbf{x}\right)\propto\exp\left[\sum_{i}w_{i}n_{i}(\mathbf{x})
ight],$$

where  $n_i(\mathbf{x})$  is the number of "true groundings" of formula  $F_i$  in the possible world  $\mathbf{x}$ .

- A possible world **x** is likely if
  - It satisfies many groundings with positive weight
  - If satisfies few groundings with negative weight
  - It satisfies groundings with high positive weight
  - It does not satisfy groundings with high negative weight

How many true groundings does a formula have?

•  $F_1 = M(A)$ 



•  $F_2 = M(A) \lor M(B)$ (M(A)) (M(D)) How many true groundings does a formula have? (2)



Conjunctions in FOL are sensitive to noise: If just one of the conjuncts is unsatisfied, the formula is also unsatisfied. MLNs count how many of the conjuncts are true and thus are less sensitive to noise. How many true groundings does a formula have? (3)



Disjunctions in FOL are insensitive to noise, so we are fine.

## Grounding a formula in Markov logic

Let *F* be a formula and  $C = \{c_1, \ldots, c_d\}$  be a set of constants. Conceptually, we obtain the set G(F) of ground formulas as follows:

- Whenever a subformulas of form ∃x.F'(x) occurs, replace by (F'(c<sub>1</sub>) ∨ · · · ∨ F'(c<sub>d</sub>))
- ② Convert the formula to form ∀x.F'(x), where F' is in conjunctive normal form and is quantifier-free, optionally simplify, denote result by cnf(F)

3 For all 
$$\mathbf{c} \in C^{|\mathbf{x}|}$$
, set  $G(F, \mathbf{c}) = \{ G : G \text{ is a clause in } F'(\mathbf{c}) \}$ 

• Set 
$$G(F) = \left\{ G(F, \mathbf{c}) : \mathbf{c} \in C^{|\mathbf{x}|} \right\}$$

#### Example

• 
$$C = \{A, B\}$$

• 
$$F_1 = \forall x. \text{Smokes}(x) \implies \text{Cancer}(x)$$

 ${\small 1} {\small 0} {\small No existential quantifiers} \rightarrow {\small nothing to do}$ 

$$cnf(F_1) = \forall x. \neg S(x) \lor C(x)$$

**3** 
$$G(F_1, A) = \{ \neg S(A) \lor C(A) \}$$
  
 $C(F_1, A) = \{ \neg S(A) \lor C(A) \}$ 

 $G(F_1, B) = \{ \neg S(B) \lor C(B) \}$  $G(F_1) = \{ \{ \neg S(A) \lor C(A) \}, \{ \neg S(B) \lor C(B) \} \}$ 

# Grounding a formula (example)

#### Example

- $C = \{A, B\}$
- $F_2 = \forall x. \forall y. Friends(x, y) \implies (Smokes(x) \iff Smokes(y))$ • No existential quantifiers  $\rightarrow$  nothing to do
  - $cnf(F_2) = \forall x.\forall y.[\neg F(x,y) \lor S(x) \lor \neg S(y)] \land [\neg F(x,y) \lor \neg S(x) \lor S(y)]$

$$G(F_2, (A, A)) = \{ \neg F(A, A) \lor S(A) \lor \neg S(A), \neg F(A, A) \lor \neg S(A) \lor S(A) \} G(F_2, (A, B)) = \{ \neg F(A, B) \lor S(A) \lor \neg S(B), \neg F(A, B) \lor \neg S(A) \lor S(B) \} G(F_2, (B, A)) = \{ \neg F(B, A) \lor S(B) \lor \neg S(A), \neg F(B, A) \lor \neg S(B) \lor S(A) \} G(F_2, (B, B)) = \{ \neg F(B, B) \lor S(A) \lor \neg S(B), \neg F(B, B) \lor \neg S(A) \lor S(B) \}$$

$$G(F_2) = \{\{ \neg F(A, A) \lor S(A) \lor \neg S(A), \neg F(A, A) \lor \neg S(A) \lor S(A) \}, \\ \{ \neg F(A, B) \lor S(A) \lor \neg S(B), \neg F(A, B) \lor \neg S(A) \lor S(B) \}, \\ \{ \neg F(B, A) \lor S(B) \lor \neg S(A), \neg F(B, A) \lor \neg S(B) \lor S(A) \}, \\ \{ \neg F(B, B) \lor S(A) \lor \neg S(B), \neg F(B, B) \lor \neg S(A) \lor S(B) \} \}$$

## Grounding a Markov logic network

Given an MLN  $\{(F_i, w_i)\}$  and a set of constants C.

- Create a Boolean variable  $R(\mathbf{c})$  for each predicate that occurs in one of the formulas and each  $\mathbf{c} \in C^m$ , where *m* is the arity of the relation
- **2** For each formula  $F_i$ 
  - Ground  $F_i$  to obtain  $G(F_i)$
  - **2** For each ground set of clauses  $G(F_i, \mathbf{c}) \in G(F_i)$ 
    - Split weight evenly among clauses:  $w'_i = w_i / |G(F_i, \mathbf{c})|$
    - **2** For each clause  $F_{ij}$  in  $G(F_i, \mathbf{c})$ , create a factor

$$\phi(\mathbf{D}_{ij}) = w'_i f_{ij}(\mathbf{D}_{ij}),$$

where  $\mathbf{D}_{ij}$  is the set of variables that occur in  $F_{ij}$ , and

 $f_{ij}(\mathbf{D}_{ij}) = \begin{cases} 1 & \text{if } j\text{-th clause in in } G(F_i, \mathbf{c}) \text{ is satisfied for assignment } \mathbf{D}_{ij} \\ 0 & \text{otherwise} \end{cases}$ 

is an "indicator feature" with weight  $w'_i$ .

The weight of a ground CNF formula is split evenly among its clauses.

## Grounding a Markov logic network (example)

F <sub>1</sub> : 1.5 {	Smoking causes cancer $\forall x. \text{Smokes}(x) \implies \text{Cancer}(x)$
$F_2: 1.1 \{$	Friends have similar smoking habits $\forall x. \forall y. Friends(x, y) \implies (Smokes(x) \iff Smokes(y))$



## Outline



Introduction to Markov Logic Networks

- 2 Probabilistic Graphical Models
  - Introduction
  - Preliminaries

#### 3 Markov Networks

#### 4 Markov Logic Networks

- Grounding Markov logic networks
- Log-Linear Models

#### Inference in MLNs

- Basics
- Exact Inference
- Approximate Inference



## Log-linear model

#### Definition

A positive distribution  $\mathbb P$  is a log-linear model over a Markov network  $\mathcal H$  if it is associated with

- a set of *features*  $\mathcal{F} = \{ f_1(\mathbf{D}_1), \dots, f_m(\mathbf{D}_m) \}$ , where each  $\mathbf{D}_i$  is a complete subgraph in  $\mathcal{H}$
- a set of weights  $w_1, \ldots, w_m$

such that

$$\mathbb{P}(X_1,\ldots,X_n)\propto \exp\left[\sum_{i=1}^m w_i f_i(\mathbf{D}_i)\right].$$

The terms  $\epsilon_i(\mathbf{D}_i) = -w_i f_i(\mathbf{D}_i)$  are called *energy functions*.

 $\log \mathbb{P}(X_1, \ldots, X_n)$  is a linear combination of the the features. The linearity allows us to detect and *eliminate* redundancy in the features (using standard linear algebra techniques).

## From factors to features

#### Definition

Let **D** be a subset of variables. An *indicator feature* is a function  $f(\mathbf{D}) : \mathbf{D} \to \{0, 1\}.$ 

#### Theorem

Every factor of a graphical model on discrete variables can be expressed in terms of a linear combination of weighted indicator features.

#### Proof (Boolean case).

Consider a factor  $\phi(X_1, \ldots, X_k)$  on k Boolean variables. Let  $\Theta$  be the set of all assignments of values to  $X_1, \ldots, X_k$ . Set

$$w_{\theta} = \ln \phi(X_{1}[\theta], \dots, X_{k}[\theta])$$
(constants)  
$$f_{\theta}(X_{1}, \dots, X_{k}) = \begin{cases} 1 & \text{if } X_{1} = X_{1}[\theta], \dots, X_{k} = X_{k}[\theta] \\ 0 & \text{otherwise} \end{cases}$$
(indicator features)  
$$\ln \phi(X_{1}, \dots, X_{k}) = \sum_{\theta \in \Theta} w_{\theta} f_{\theta}(X_{1}, \dots, X_{k})$$
(decomposition)

# From factors to features (example)

#### Example

Consider three friends with similiar interests and let A, B, C be Boolean variables that indicate whether each of the friends likes football.



#### We have

$$\ln \phi(A, B, C) = \sum_{\theta} w_{\theta} f_{\theta}(A, B, C) = 2.3 \cdot f_{FFF}(A, B, C) + 2.3 \cdot f_{TTT}(A, B, C).$$

Even more compact:  $\ln \phi(A, B, C) = 2.3 \cdot I_{ABC \vee \neg A \neg B \neg C}$ 

## From Gibbs distribution to log-linear models

#### Theorem

Every positive Gibbs distribution  $\mathbb{P}$  over  $\mathcal{H}$  on Boolean variables  $X_1, \ldots, X_n$  has a log-linear model over  $\mathcal{H}$  with only indicator features and vice versa.

# Proof. $\mathbb{P}(X_1, \dots, X_n) = \frac{1}{Z} \prod_{i=1}^m \phi_i(\mathbf{D}_i)$ $= \frac{1}{Z} \exp\left[\sum_{i=1}^m \ln \phi_i(\mathbf{D}_i)\right]$ $= \frac{1}{Z} \exp\left[\sum_{i=1}^m \sum_{\theta \in \Theta_{\mathbf{D}_i}} w_\theta f_\theta(\mathbf{D}_i)\right].$

Markov logic networks are "templates" for constructing loglinear models. Any positive Gibbs distribution with finitedomain variables can be modeled.

# Outline



- 2 Probabilistic Graphical Models
  - Introduction
  - Preliminaries
- 3 Markov Networks
- 4 Markov Logic Networks
  - Grounding Markov logic networks
  - Log-Linear Models

#### 5 Inference in MLNs

- Basics
- Exact Inference
- Approximate Inference

#### 5 Summary

## Outline



Introduction to Markov Logic Networks

- 2 Probabilistic Graphical Models
  - Introduction
  - Preliminaries

## 3 Markov Networks

#### 4 Markov Logic Networks

- Grounding Markov logic networks
- Log-Linear Models

## Inference in MLNs

- Basics
- Exact Inference
- Approximate Inference



## Inference in probabilistic graphical models

- Recall the queries of interest
  - Conditional probability query
  - 2 MAP query
  - Marginal MAP query

#### Definition

Let  $\mathbb{P}_{\Phi}$  be a Gibbs distribution over variables  $\{X, X_1, \ldots, X_n\}$ .

• The  $\mathbb{P}_{\Phi}$ -decision problem asks whether  $\mathbb{P}_{\Phi}(X = x) > 0$ ,

**2** The  $\mathbb{P}_{\Phi}$ -probability computation problem asks for  $\mathbb{P}_{\Phi}(X = x)$ .

## Complexity of inference in probabilistic graphical models

#### Theorem

The  $\mathbb{P}_{\Phi}$ -decision problem is NP-complete,  $\mathbb{P}_{\Phi}$ -probability computation is #P-hard.

#### Proof (by reduction from 3-SAT and #3-SAT).

Take a 3-SAT formula  $\Psi = C_1 \wedge C_2 \wedge \ldots \wedge C_m$  over variables  $\mathcal{X} = \{X_1, X_2, \ldots, X_n\}$ . Consider the following Gibbs distribution  $\mathbb{P}_{\Phi}$  over Boolean variables:



Here,  $\forall_i (C_i, \mathbf{X}_i) = 1$  if for assignment  $\mathbf{X}_i$  the truth value of clause  $C_i$  equals variable  $C_i$ , else  $\forall_i (C_i, \mathbf{X}_i) = 0$ ; similarly for  $\wedge$ -factors.  $\mathbb{P}_{\Phi}$  can be computed in polynomial time in the size of  $\Psi$ . Assertion 1 follows since  $\mathbb{P}_{\Phi} (X = \text{TRUE}) > 0$  if and only if  $\Psi$  is satisfiable.  $\mathbb{P}_{\Phi} (X = \text{TRUE}) = \mathbb{P} (\Psi)$  where  $\mathbb{P} (X_i = \text{TRUE}) = 1/2$  and the  $\{X_i\}$  are i.i.d. Assertion 2 follows since  $\#\Psi = 2^n \mathbb{P} (\Psi) = 2^n \mathbb{P}_{\Phi} (X = \text{TRUE})$ .

## Queries in Markov logic

- Standard PGM queries, e.g.,  $\mathbb{P}(\mathsf{Smokes}(\mathsf{B}),\mathsf{Cancer}(\mathsf{B}) \mid \mathsf{Smokes}(\mathsf{A}) \land \mathsf{Friends}(\mathsf{A},\mathsf{B}) \land \dots)$ 
  - $\rightarrow$  #P-hard
- More general queries of form "What is the probability that formula F<sub>1</sub> holds given that formula F<sub>2</sub> holds?", e.g.,

   P(∃x.Cancer(x) | ∀x.Smokes(x))
- Let L be an MLN and C be a set of constants

$$\mathbb{P}(F_1 | F_2, L, C) = \mathbb{P}(F_1 | F_2, M_{L,C})$$

$$= \frac{\mathbb{P}(F_1 \wedge F_2 | M_{L,C})}{\mathbb{P}(F_2 | M_{L,C})}$$

$$= \frac{\sum_{\mathbf{x} \in \mathcal{X}_{F_1} \cap \mathcal{X}_{F_2}} \mathbb{P}(\mathcal{X} = \mathbf{x} | M_{L,C})}{\sum_{\mathbf{x} \in \mathcal{X}_{F_2}} \mathbb{P}(\mathcal{X} = \mathbf{x} | M_{L,C})},$$

where  $\mathcal{X}_F$  is the set of worlds in which F holds

We focus on standard PGM queries.

## Outline



Introduction to Markov Logic Networks

- 2 Probabilistic Graphical Models
  - Introduction
  - Preliminaries

## 3 Markov Networks

#### 4 Markov Logic Networks

- Grounding Markov logic networks
- Log-Linear Models

### Inference in MLNs

- Basics
- Exact Inference
- Approximate Inference



## Naive approach

В

'*c*<sup>1</sup>

#### Exponential in number of variables!

	A	В	С	D	$\tilde{\mathbb{P}}$	$\mathbb{P}$
	<i>a</i> 0	$b_0$	<i>c</i> 0	$d_0$	300,000	0.04
	<i>a</i> 0	$b_0$	<i>c</i> 0	$d_1$	300,000	0.04
)	<i>a</i> 0	$b_0$	$c_1$	$d_0$	300,000	0.04
	<i>a</i> 0	$b_0$	<i>c</i> <sub>1</sub>	$d_1$	30	$4.1 \cdot 10^{-6}$
<u> </u>	<i>a</i> 0	$b_1$	<i>c</i> 0	$d_0$	500	$6.9\cdot10^{-5}$
$C \phi_2$	<i>a</i> 0	$b_1$	<i>c</i> <sub>0</sub>	$d_1$	500	$6.9\cdot10^{-5}$
$c^{0}$ 100	<i>a</i> 0	$b_1$	<i>c</i> 1	$d_0$	5,000,000	0.69
$c^{1}$ 1	<i>a</i> 0	$b_1$	<i>c</i> <sub>1</sub>	$d_1$	500	$6.9 \cdot 10^{-5}$
$c^0$ 1	<i>a</i> 1	$b_0$	<i>c</i> <sub>0</sub>	$d_0$	100	$1.4 \cdot 10^{-5}$
<i>c</i> <sup>1</sup> 100	a <sub>1</sub>	$b_0$	<i>c</i> 0	$d_1$	1,000,000	0.14
A d.	a <sub>1</sub>	$b_0$	<i>c</i> 1	$d_0$	100	$1.4 \cdot 10^{-5}$
$\rho_{4}^{0} = \frac{\phi_{4}}{100}$	$a_1$	$b_0$	<i>c</i> 1	$d_1$	100	$1.4 \cdot 10^{-5}$
a 100	$a_1$	$b_1$	<i>c</i> 0	$d_0$	10	$1.4 \cdot 10^{-6}$
a 1 $a^0 1$	$a_1$	$b_1$	<i>c</i> 0	$d_1$	100,000	0.014
$a \downarrow 1$	a <sub>1</sub>	$b_1$	$c_1$	$d_0$	100,000	0.014
a 100	$a_1$	$b_1$	$c_1$	$d_1$	100,000	0.014
				7 =	= 7.201.840	60 / 1



1

## Grounding with evidence (1)

Denote by M the weighted ground clauses in a ground Markov logic network  $M_{L,C}$ . Given evidence **E**, we can partition M into:

- Clauses  $M_1$  that involve only observed variables
- **2** Clauses  $M_2$  that involve both observed and latent variables
- Solution  $M_3$  that involve only latent variables

$$\log \mathbb{P}(\mathbf{W} \mid \mathbf{E}) = -\log Z + \sum_{\substack{\phi = (f, w) \in M}} wf(\mathbf{W}_f, \mathbf{E}_f)$$
$$= -\log Z + \underbrace{\sum_{\substack{(f, w) \in M_1 \\ \text{Constant}}} wf(\mathbf{E}_f)}_{\text{Constant}} + \underbrace{\sum_{\substack{(f, w) \in M_2 \\ \text{Constant}}} wf(\mathbf{W}_f, \mathbf{E}_f) + \sum_{\substack{(f, w) \in M_3 \\ \text{W} \in M_3}} wf(\mathbf{W}_f)$$



## Grounding with evidence (2)

Denote by M the weighted ground clauses in a ground Markov logic network  $M_{L,C}$ . Given evidence **E**, we can partition M into:

- Clauses  $M_1$  that involve only observed variables
- ② Clauses  $M_2$  that involve both observed and latent variables
- **③** Clauses  $M_3$  that involve only latent variables

$$\log \mathbb{P}(\mathbf{W} \mid \mathbf{E}) = -\log Z + \sum_{\substack{\phi = (f, w) \in M}} wf(\mathbf{W}_f, \mathbf{E}_f)$$
$$= -\log Z' + \sum_{\substack{(f, w) \in M_2 \\ \text{Replace observed variables by their values}} wf(\mathbf{W}_f, \mathbf{E}_f) + \sum_{\substack{(f, w) \in M_3 \\ \text{Replace observed variables by their values}} wf(\mathbf{W}_f)$$



## Grounding with evidence (3)

Denote by M the weighted ground clauses in a ground Markov logic network  $M_{L,C}$ . Given evidence **E**, we can partition M into:

- Clauses  $M_1$  that involve only observed variables
- ② Clauses  $M_2$  that involve both observed and latent variables
- **③** Clauses  $M_3$  that involve only latent variables

$$\begin{split} \operatorname{og} \mathbb{P} \left( \mathbf{W} \mid \mathbf{E} \right) &= -\log Z + \sum_{\phi = (f, w) \in M} \operatorname{wf} (\mathbf{W}_f, \mathbf{E}_f) \\ &= -\log Z' + \sum_{(f, w) \in M'_2} \operatorname{wf} (\mathbf{W}_f) + \sum_{(f, w) \in M_3} \operatorname{wf} (\mathbf{W}_f) \\ &= -\log Z' + \sum_{(f, w) \in M'} \operatorname{wf} (\mathbf{W}_f) \end{split}$$

No observed variables are left. Gives rise to efficient grounding methods.

#### Example



$$M'_{2} = \{ \phi'_{23}, \phi'_{24}, \phi'_{25}, \phi'_{26}, \phi'_{27}, \phi'_{28} \}$$

$$M_{3} = \{ \phi_{12} \}$$

$$M' = M'_{2} \cup M_{3}$$

$$\phi_{24} = \neg F(A, B) \lor \neg S(A) \lor S(B)$$

$$\phi'_{24} = \text{FALSE} \lor \text{FALSE} \lor S(B)$$

$$= S(B)$$

$$63.778$$

# MAP inference for MLNs (1)

#### Example

What is the most likely world for a given Markov logic network?



Corresponds to weighted CNF formula:  $\Psi = f_{11} \wedge f_{12} \wedge f_{23} \wedge f_{24} \wedge f_{25} \wedge f_{26} \wedge f_{27} \wedge f_{28}$ 

# MAP inference for MLNs (2)

#### Definition

Consider a CNF formula F over variables  $\mathcal{X}$ , in which each of the clauses  $f_1, \ldots, f_m$  is associated with a corresponding weight  $w_1, \ldots, w_m$ . The Weighted MAX-SAT problem is to find an assignment  $\mathbf{x}^* \in \mathcal{X}_F$  that maximizes the sum of the weights of satisfied clauses, i.e.,  $\mathbf{x}^* = \operatorname{argmax}_{\mathbf{x}} \sum_i w_i f_i$ .

Consider the following transformation:

$$\arg\max_{\mathbf{x}} \mathbb{P}(\mathbf{x}) = \arg\max_{\mathbf{x}} \left[ \frac{1}{Z} \exp \sum_{\substack{(f,w) \in M_{L,C} \\ \mathbf{x}}} wf(\mathbf{x}) \right]$$
$$= \arg\max_{\mathbf{x}} \sum_{\substack{(f,w) \in M_{L,C} \\ F}} w_{i} \frac{f(\mathbf{x})}{f_{i}} = \mathbf{x}^{*}$$

There are many algorithms and solvers for Weighted MAX-SAT, both exact and approximate. Specialized algorithms for MLNs do exist; they try to reduce grounding by computing  $M_{L,C}$  only partially.

# MAP inference for MLNs (3)



MAP world characterizes distribution well MAP world not distinguished from other terize only a part of words the distribution

MAP estimates provide the "most consistent" world, i.e., the world that satisfies most of the rules. This world *may or may not* characterize the entire distribution well.

# Variable elimination (idea)

Goal: Eliminate non-query variables from the graph.



 $c^{1}$ 

1000

 $a^1 b^1$ 

graph represents  $\mathbb{P}(A, C, D)$ .

01 / 10

## Variable elimination (why it works)

Recall that

$$\mathbb{P}(A, B, C, D) = \frac{1}{Z}\phi_1(A, B) \times \phi_2(B, C) \times \phi_3(C, D) \times \phi_4(D, A)$$

and thus

$$\begin{split} \mathbb{P}(A, C, D) &= \mathbb{P}(A, b^{0}, C, D) + \mathbb{P}(A, b^{1}, C, D) \\ &= \frac{1}{Z} [\phi_{1}(A, b^{0}) \times \phi_{2}(b^{0}, C) \times \phi_{3}(C, D) \times \phi_{4}(D, A) \\ &+ \phi_{1}(A, b^{1}) \times \phi_{2}(b^{1}, C) \times \phi_{3}(C, D) \times \phi_{4}(D, A)] \\ &= \frac{1}{Z} \left[ \left\{ \sum_{b \in \{b^{0}, b^{1}\}} \phi_{1}(A, b) \times \phi_{2}(b, C) \right\} \times \phi_{3}(C, D) \times \phi_{4}(D, A) \right] \\ &= \frac{1}{Z} \left[ \phi_{12}(A, C) \times \phi_{3}(C, D) \times \phi_{4}(D, A) \right] \end{split}$$

# Variable elimination (remarks)

- Also called sum-product variable elimination
- Whenever we eliminate a variable B
  - We remove all factors connected to B
  - We introduce a single factor that is connected to the neighbors of B
  - If B has k neighbors, the new factor has  $2^k$  rows
    - $\rightarrow$  Potentially exponential blow-up
- Computational cost
  - Dominated by sizes of intermediate factors
  - Depends strongly on elimination ordering
  - NP-hard to find optimal ordering
  - Lots of useful heuristics exist
  - "Conditioning" can be used to avoid large factors for increased processing time
- Similar observations give rise to other important algorithms, e.g., "message passing" in "clique trees"

## Outline



Introduction to Markov Logic Networks

- 2 Probabilistic Graphical Models
  - Introduction
  - Preliminaries

#### 3 Markov Networks

#### 4 Markov Logic Networks

- Grounding Markov logic networks
- Log-Linear Models

### Inference in MLNs

- Basics
- Exact Inference
- Approximate Inference

## Summary

## Sampling methods

- Also called particle-based approximate inference
- Idea: Obtain samples from the distribution underlying the graphical model
- If samples were independent, we could count how often each variables is true/false and apply the sampling theorem
- $\bullet\,$  Sampling is much more difficult in Markov networks  $\to\,$  samples are generally dependent
  - Goal is to minimize the dependencies
  - More samples needed than "implied" by the sampling theorem
  - $\blacktriangleright$  If dependencies vanish between far-apart samples  $\rightarrow$  correctness and convergence
- Many techniques
  - Forward sampling (for directed models)
  - Likelihood weighting
  - Importance sampling
  - Gibbs sampling
  - Other Markov Chain Monte Carlo (MCMC) methods
  - Collapsed particles

# Gibbs sampling (idea)

Gibbs sampling is a simple algorithm to sample from  $\mathbb{P}(X, Y)$ . It is used when it is hard to sample from  $\mathbb{P}(X, Y)$ , but easy to sample from  $\mathbb{P}(X | Y)$  and  $\mathbb{P}(Y | X)$ .

- Pick an initial point  $(x_0, y_0)$
- **2** For n = 1, 2, ...

• Generate 
$$x_n \sim \mathbb{P}(X \mid Y = y_{n-1})$$


# Gibbs sampling (idea)

Gibbs sampling is a simple algorithm to sample from  $\mathbb{P}(X, Y)$ . It is used when it is hard to sample from  $\mathbb{P}(X, Y)$ , but easy to sample from  $\mathbb{P}(X | Y)$  and  $\mathbb{P}(Y | X)$ .

- Pick an initial point  $(x_0, y_0)$
- **2** For n = 1, 2, ...
  - Generate  $x_n \sim \mathbb{P}(X \mid Y = y_{n-1})$
  - **2** Generate  $y_n \sim \mathbb{P}(Y \mid X = x_n)$



73 / 78

### Gibbs sampling for Markov networks

Recall that

$$\mathbb{P}(A, B, C, D) = \frac{1}{Z}\phi_1(A, B) \times \phi_2(B, C) \times \phi_3(C, D) \times \phi_4(D, A).$$

Sampling from  $\mathbb{P}(A, B, C, D)$  is hard but sampling from

$$\mathbb{P}(A \mid B, C, D) = \frac{\mathbb{P}(A, B, C, D)}{\mathbb{P}(B, C, D)}$$
$$= \frac{\frac{1}{Z}[\phi_1(A, B) \times \phi_2(B, C) \times \phi_3(C, D) \times \phi_4(D, A)]}{\frac{1}{Z} \sum_{a \in \{a^0, a^1\}} [\phi_1(a, B) \times \phi_2(B, C) \times \phi_3(C, D) \times \phi_4(D, a)]}$$
$$= \frac{\phi_1(A, B) \times \phi_4(D, A)}{\sum_{a \in \{a^0, a^1\}} \phi_1(a, B) \times \phi_4(D, a)}$$

is easy. Only the factors connected to A remain.

When resampling a variable A, we only have to look at the factors connected to A, and thus only the subset of variables connected to A. These variables are called the *Markov blanket* of A.

## Gibbs sampling for Markov networks (remarks)

- Variables are picked according to a *schedule* 
  - $\rightarrow$  sequential, random, . . .
- An instance of the more general class of MCMC methods
  - Markov chains describe how the sampling process moves through the set of worlds
  - Irreducible if all worlds can be reached from all other worlds
  - Convergence speed depends on how fast the sampling process moves (*mixing time*)



Gibbs sampling works Gibbs sampling works Gibbs sampling does not well (fast mixing) reasonable (slow mixing) work (not irreducible)

• MCMC methods can perform "bigger" steps than Gibbs sampling; they change multiple variables simultaneously

# Outline



- Probabilistic Graphical Models
  - Introduction
  - Preliminaries
- 3 Markov Networks
- 4 Markov Logic Networks
  - Grounding Markov logic networks
  - Log-Linear Models

#### 5 Inference in MLNs

- Basics
- Exact Inference
- Approximate Inference



#### Lessons learned

- Probabilistic databases and graphical models focus on different aspects of probabilistic reasoning
- Probabilistic graphical models
  - Describe and reason about probability distributions and independencies
  - Exploit locality structure (conditional independence)
  - Main components: representation, inference, learning
- Markov logic
  - Combines first-order logic and probability theory
  - Set of formulas with weights
  - Template for generating undirected graphical models
- Inference
  - #P-hard in general
  - MAP inference on MLNs corresponds to Weighted MAX-SAT
  - Exact methods for probability computation (e.g., variable elimination) may work when graph has no dense regions
  - Approximate methods often based on MCMC sampling
  - Gibbs sampling is the simplest MCMC method; it changes one variable at a time

## Suggested reading

- Daphne Koller, Nir Friedman *Probabilistic Graphical Models: Principles and Techniques* The MIT Press, 2009
- Matthew Richardson and Pedro Domingos Markov Logic Networks Machine Learning, 62(1-2), pp. 107–136, 2006
- Michael Mitzenmacher, Eli Upfal *Probability and Computing: Randomized Algorithms and Probabilistic Analysis* Cambridge University Press, 2005
- http://alchemy.cs.washington.edu/