

An Introduction to Statistical Relational Learning

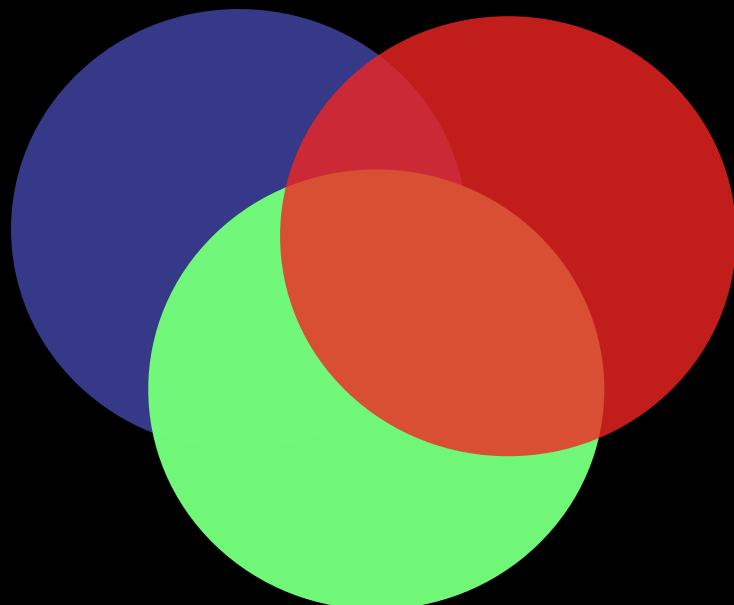
Luc De Raedt
Katholieke Universiteit Leuven

joint work with Kristian Kersting

Probabilistic Logic Learning*

One of the key open questions of artificial intelligence concerns

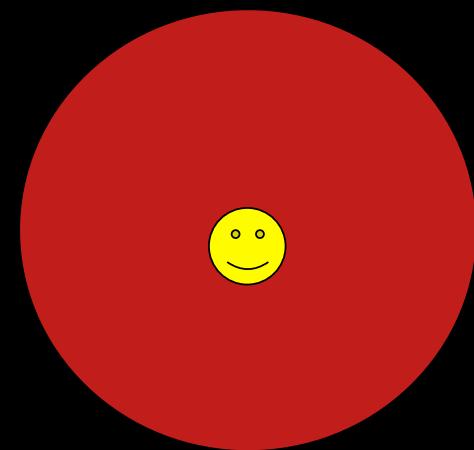
"probabilistic logic learning",
i.e. the integration of
probabilistic reasoning
with
**logical or relational
representations** and
machine learning.



*In the US, sometimes called **Statistical Relational Learning**

Probabilistic **Logic** Learning*

The knowledge representation
framework par excellence



Logic: Syntax

class(Book,Topic) :- Clause(s) *constant*
author(Book,Author), favorite-topic(Author,Topic).

the class of Book is Topic

the Book is written by Author

Topic is the favorite of the Author

favorite-topic(rowlings,fantasy) Fact(s)

the favorite topic of rowlings is fantasy

predicate/relation

Logic: three views

entailment/logical consequence

$H \models \text{class(harrypotter,fantasy)}$

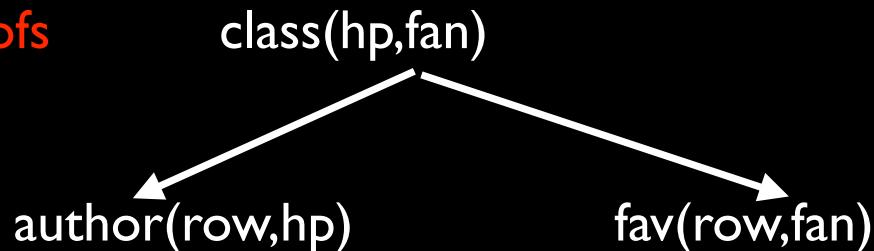
$\text{class(harrypotter,fantasy)}$ follows from H

interpretation $\{\text{class}(hp,fan), \text{author}(row,hp), \text{fav}(row,fan)\}$

is a model for H

is a possible world given H

proofs

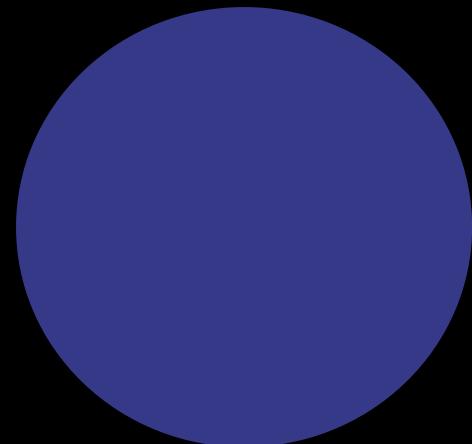


Probabilistic Logic Learning* 😊

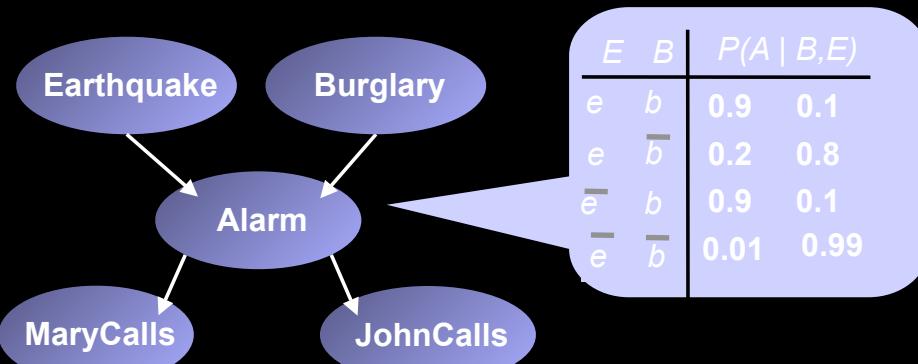
Reasoning about
uncertainty

Graphical models, such
as Bayesian and Markov
Networks

Probabilistic grammars
and databases.



Bayesian Networks



$$P(E, B, A, J, M) = P(E).P(B).P(A|E).P(A|B).P(J|A).P(M|A)$$

Structure &
Parameters

INTERPRETATION
STATE/DESCRIPTION
 $\{A, \neg E, \neg B, J, M\}$

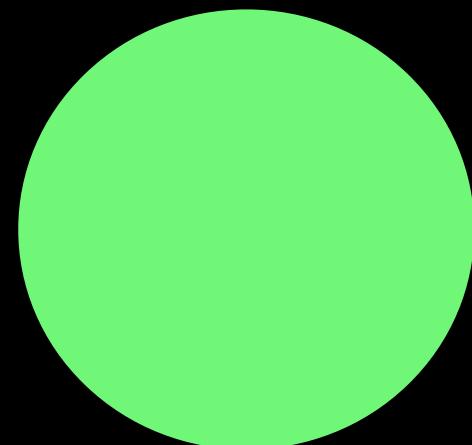
Probabilistic Logic Learning*

Learning the model from
data

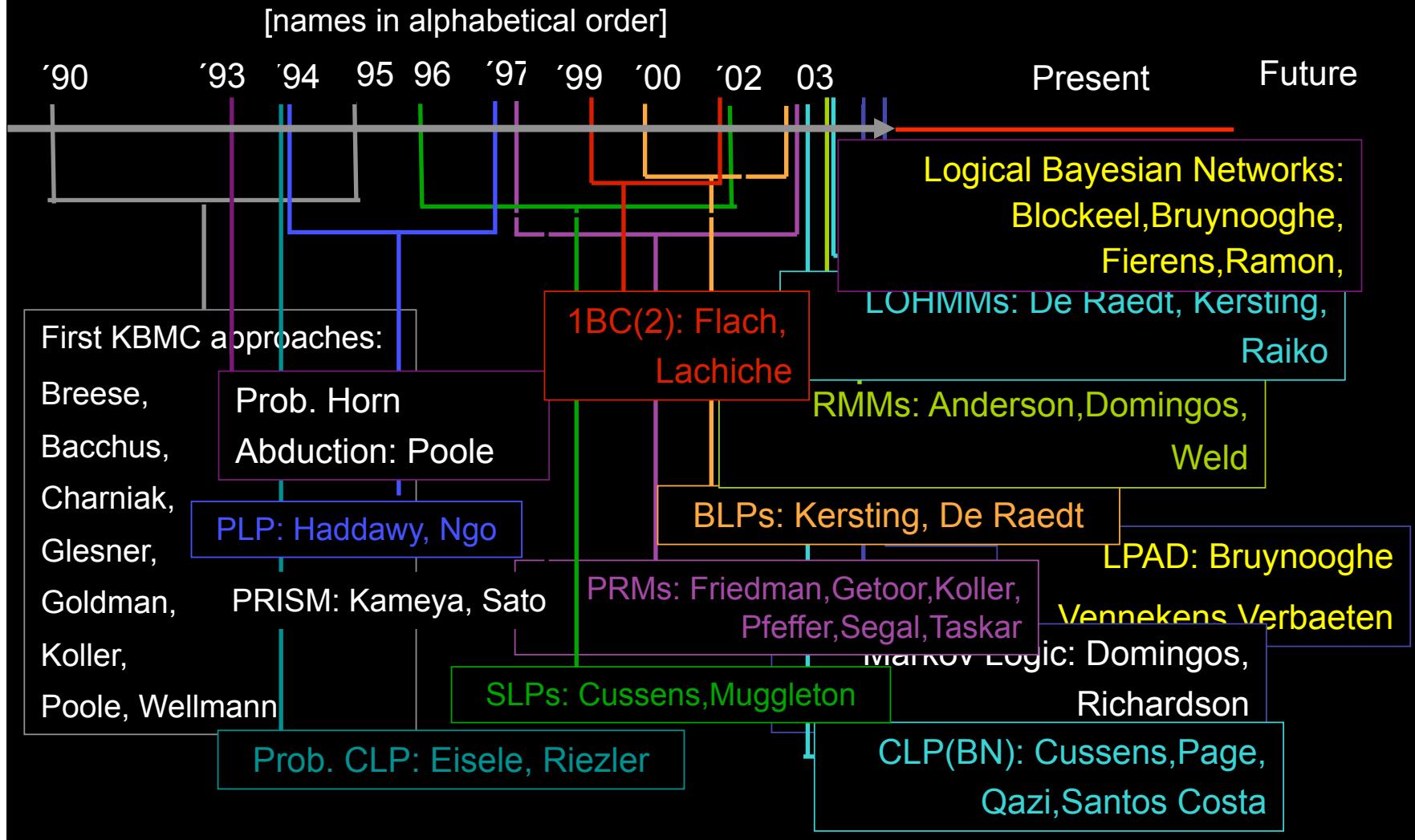


Structure &
Parameters

Logical &
Statistical



Some PLL formalisms



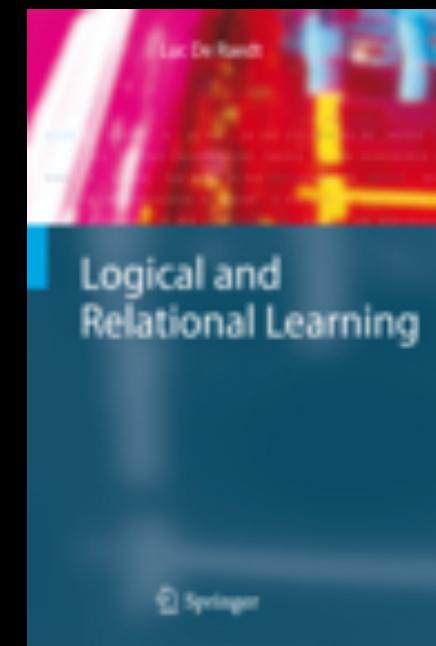
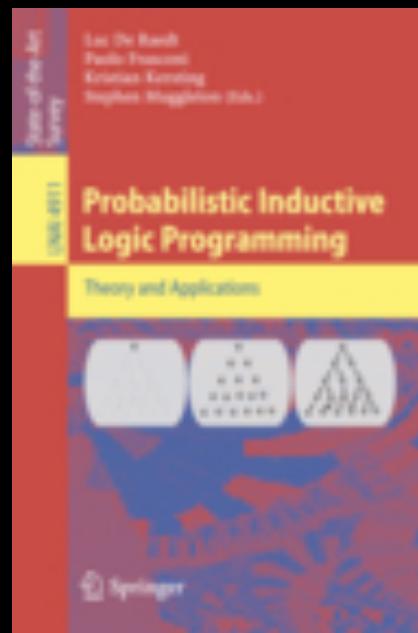
What I am not going to do

- Introduce yet another “lingua franca” or probabilistic logic
 - analogy with programming languages & knowledge representation
 - Talk a lot about logic

Rather

- I will start from logic + learning and extend it with probabilistic reasoning
- I will argue that traditional logic + learning settings apply to probabilistic logic learning as well
- I will argue that this is probably useful for applications and that learning is essential

Books



Also, survey paper De Raedt, Kersting, SIGKDD Explorations 03

Overview

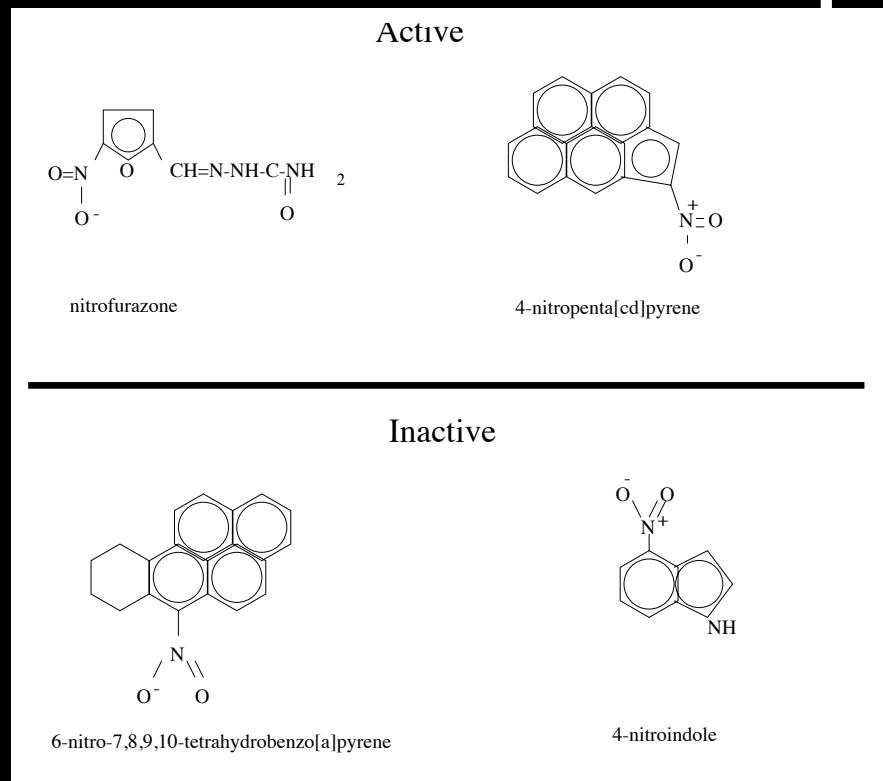
- A. Motivation for Probabilistic Logic Learning
- B. Logic + Learning
- C. Adding probability

A. Motivation

A. Motivation

A. I Some cases

Case I: Structure Activity Relationship Prediction



[Srinivasan et al. AlJ 96]

Structural alert:



General Purpose
Logic Learning System

Data = Set of Small Graphs

Uses and Produces
Knowledge

Using and Producing Knowledge

Logical and relational learning systems can use and produce knowledge (ILP)

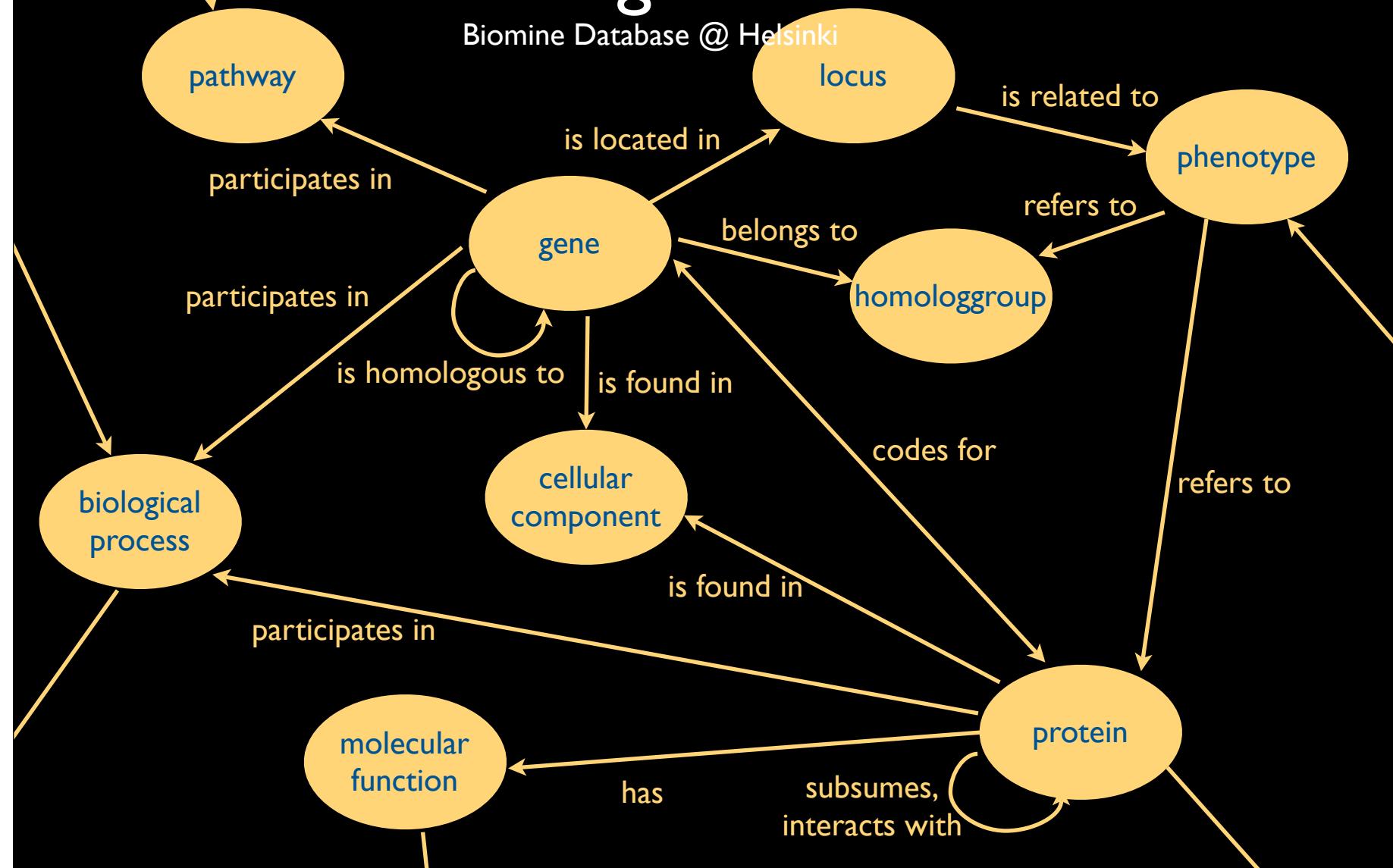
Result of learning task is understandable and interpretable

Logical and relational learning algorithms can use background knowledge, e.g. ring structures

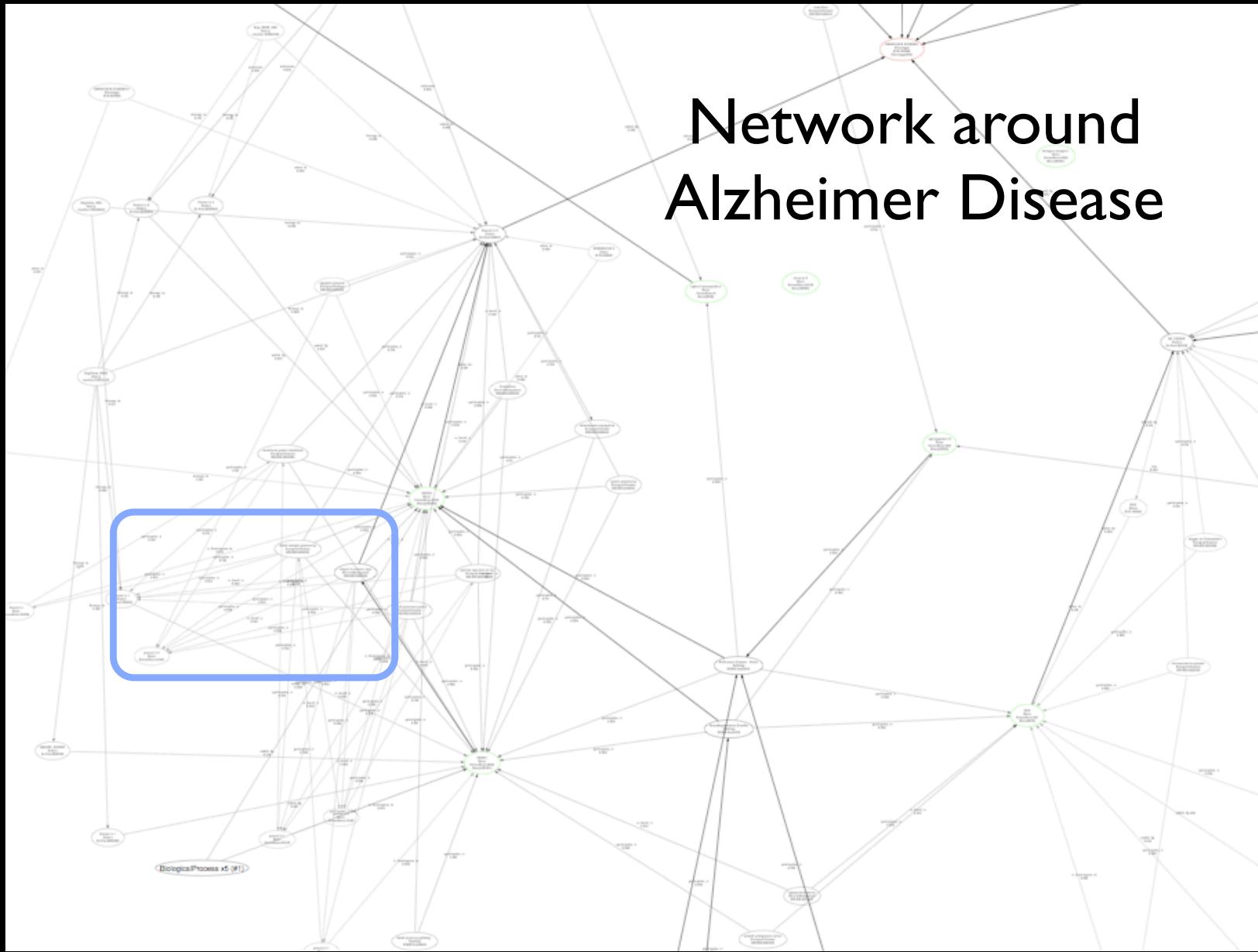
Data = Large (Probabilistic) Network

Case 2: Biological Networks

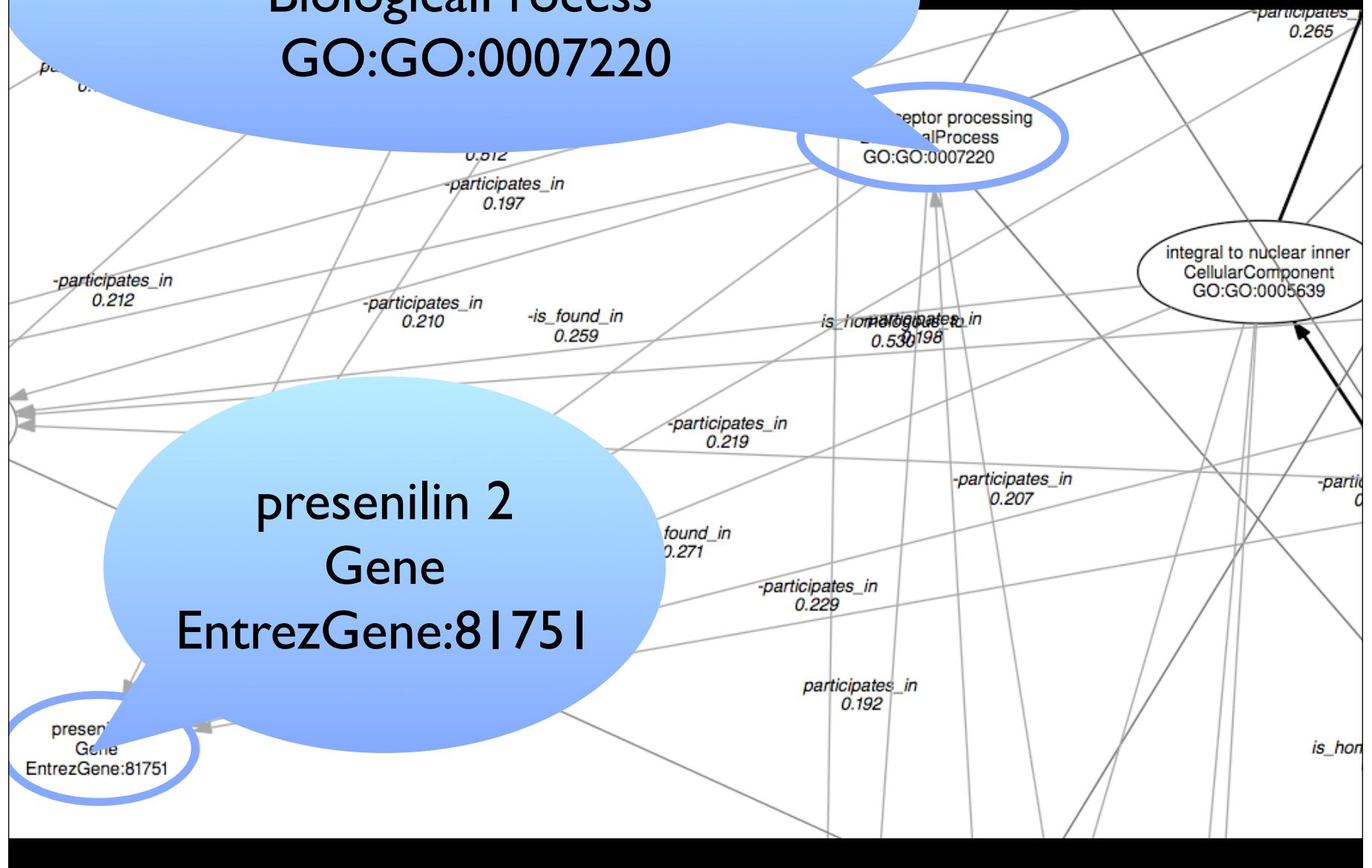
Biomine Database @ Helsinki

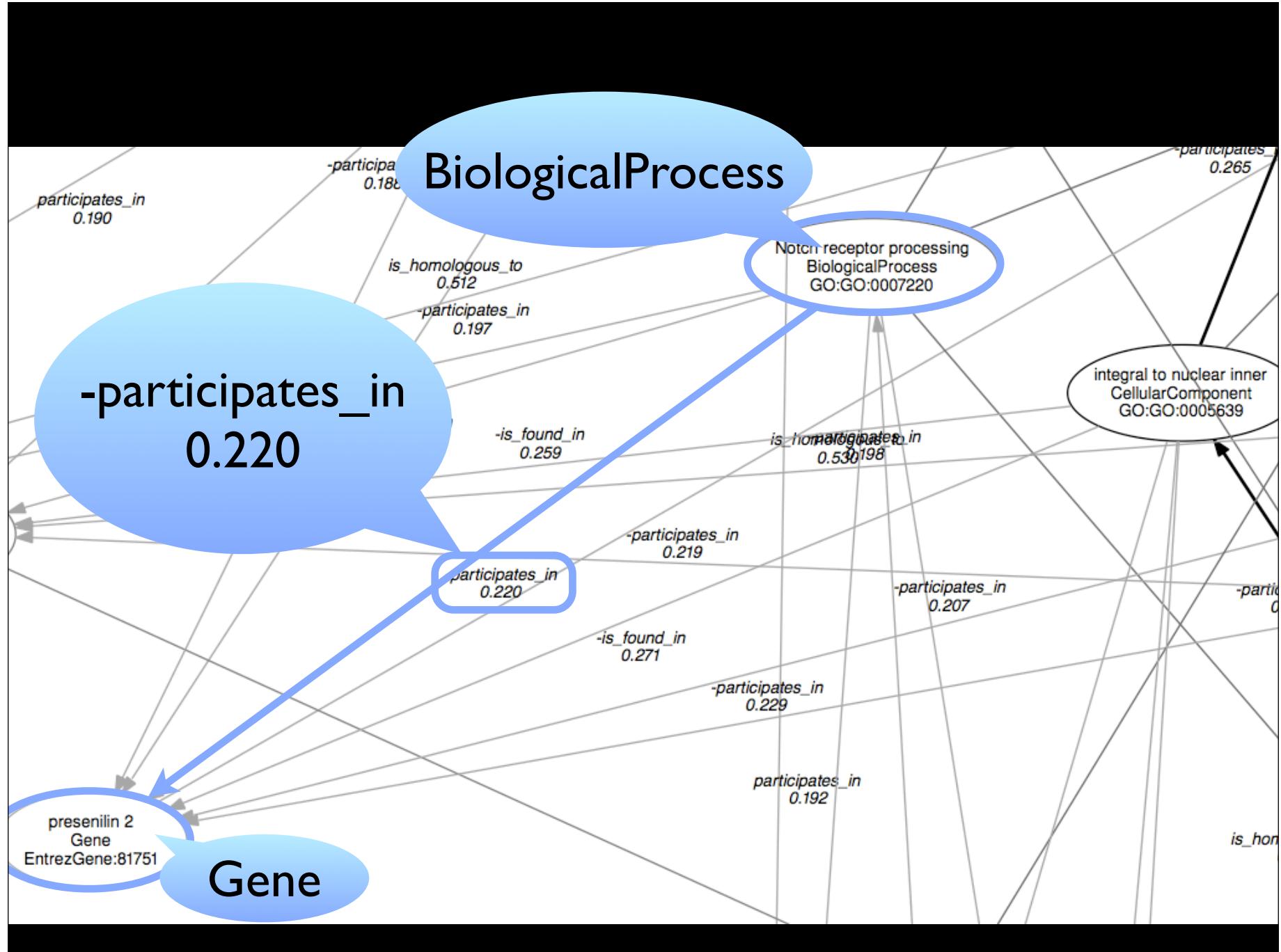


Network around Alzheimer Disease



Notch receptor processing
BiologicalProcess
GO:GO:0007220





Questions to ask

How to support the life scientist in using and discovering new knowledge in the network ?

- Is gene X involved in disease Y ?
- Should there be a link between gene X and disease Y ? If so, what type of link ?
- What is the probability that gene X is connected to disease Y ?
- Which genes are similar to X w.r.t. disease Y?
- Which part of the network provides the most information (network extraction) ?
- ...

See *Invited Talk on Tuesday*

Case 3: Evolving Networks



- *Travian*: A massively multiplayer real-time strategy game
 - Commercial game run by TravianGames GmbH
 - ~3.000.000 players spread over different “worlds”
 - ~25.000 players in one world

[Thon et al. ECML 08]



World Dynamics

Fragment of world with

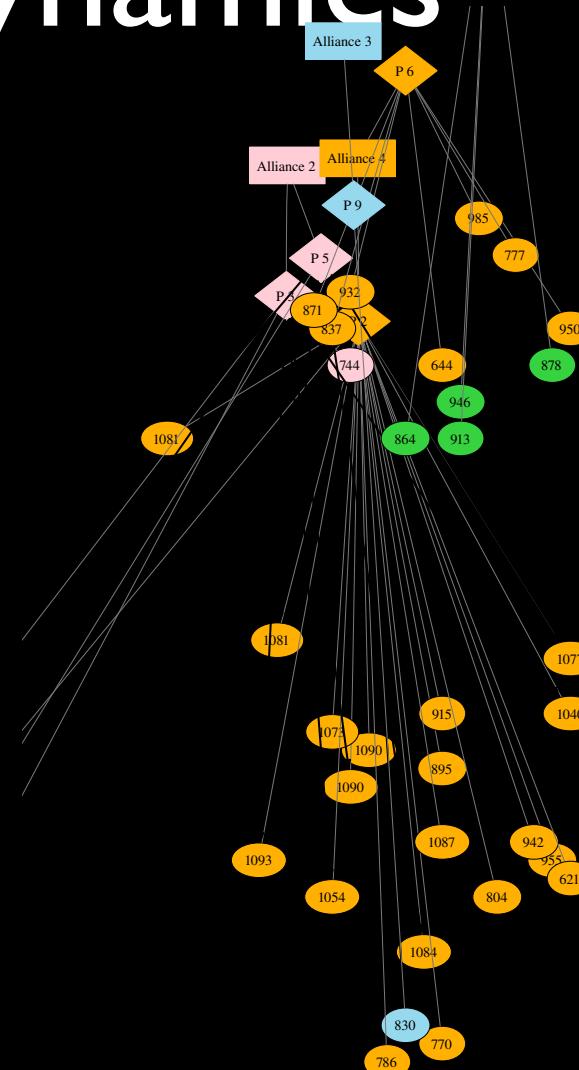
- ~10 alliances
- ~200 players
- ~600 cities

alliances color-coded

Can we build a model
of this world ?

Can we use it for playing
better ?

[Thon, Landwehr, De Raedt, ECML08]



World Dynamics

Fragment of world with

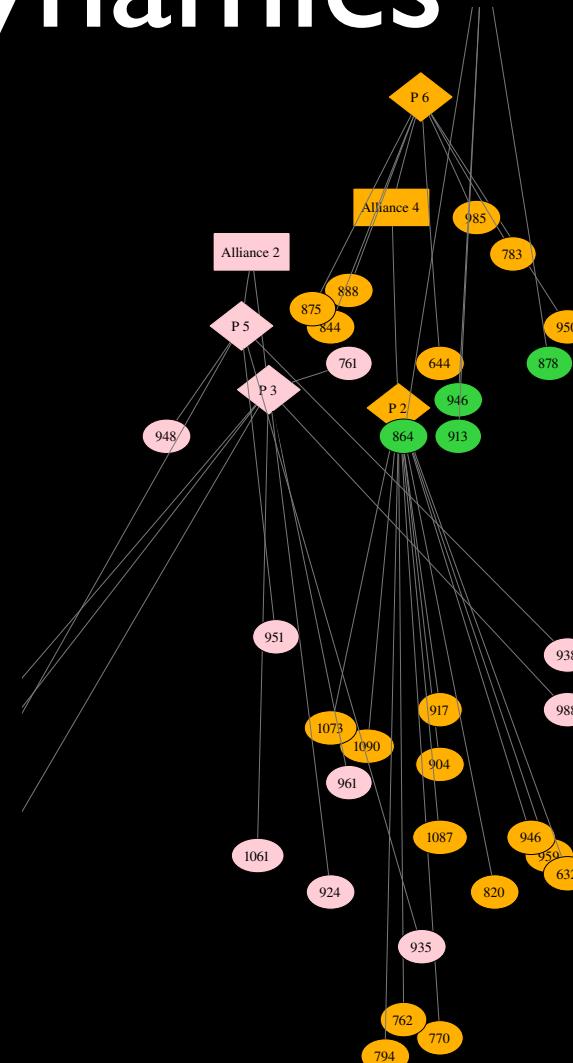
~10 alliances
~200 players
~600 cities

alliances color-coded

Can we build a model of this world ?

Can we use it for playing
better ?

[Thon, Landwehr, De Raedt, ECML08]



World Dynamics

Fragment of world with

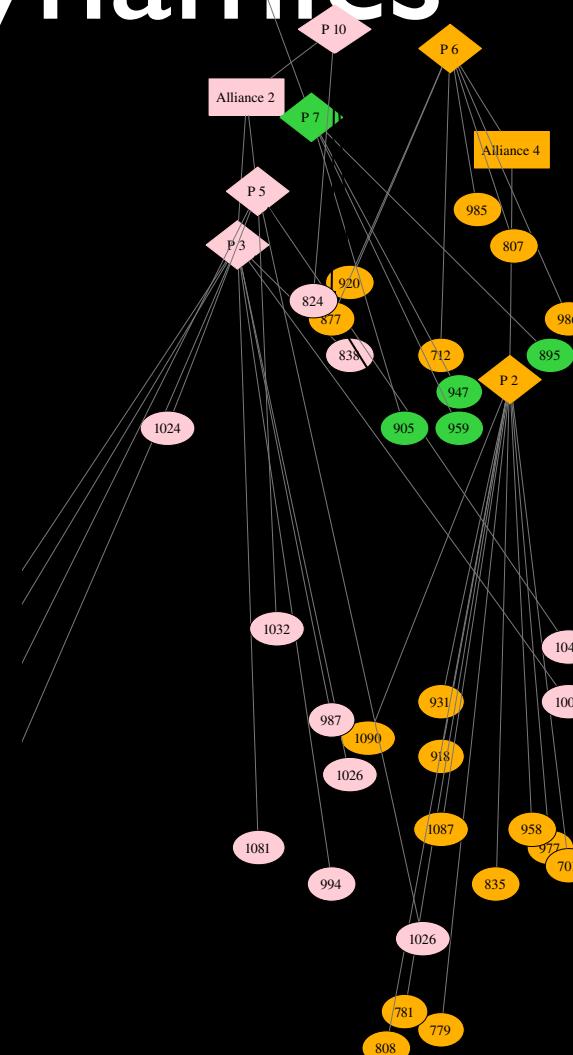
~10 alliances
~200 players
~600 cities

alliances color-coded

Can we build a model of this world ?

Can we use it for playing
better ?

[Thon, Landwehr, De Raedt, ECML08]



World Dynamics

Fragment of world with

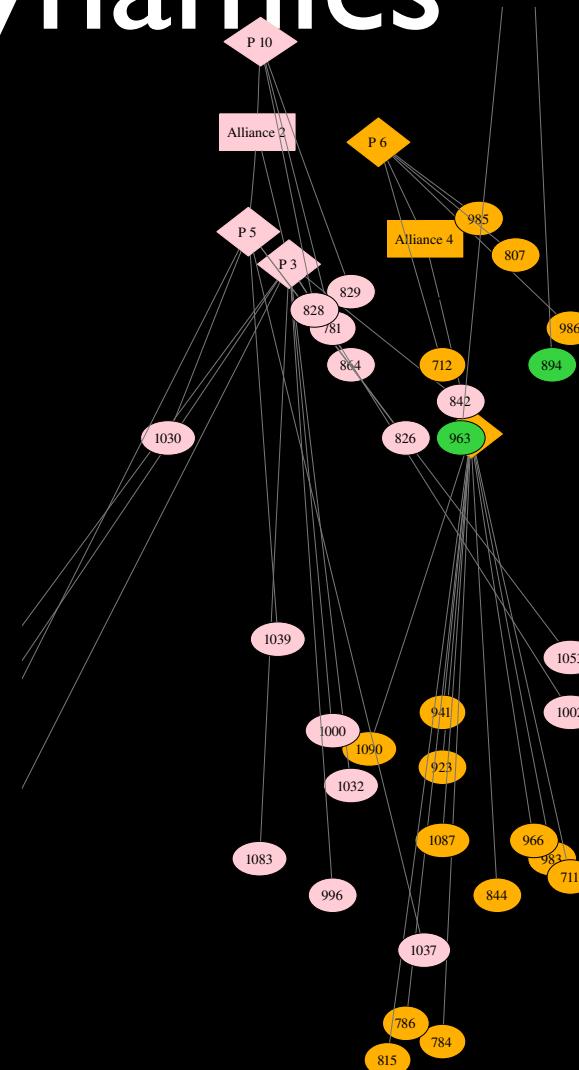
~10 alliances
~200 players
~600 cities

alliances color-coded

Can we build a model
of this world ?

Can we use it for playing
better ?

[Thon, Landwehr, De Raedt, ECML08]



World Dynamics

Fragment of world with

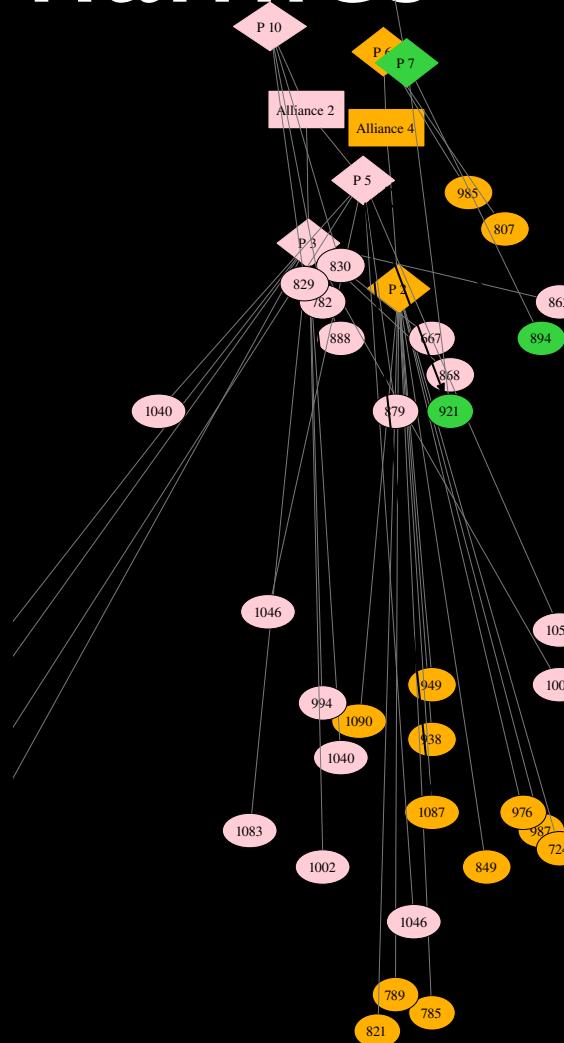
~10 alliances
~200 players
~600 cities

alliances color-coded

Can we build a model
of this world ?

Can we use it for playing
better ?

[Thon, Landwehr, De Raedt, ECML08]



World Dynamics

Fragment of world with

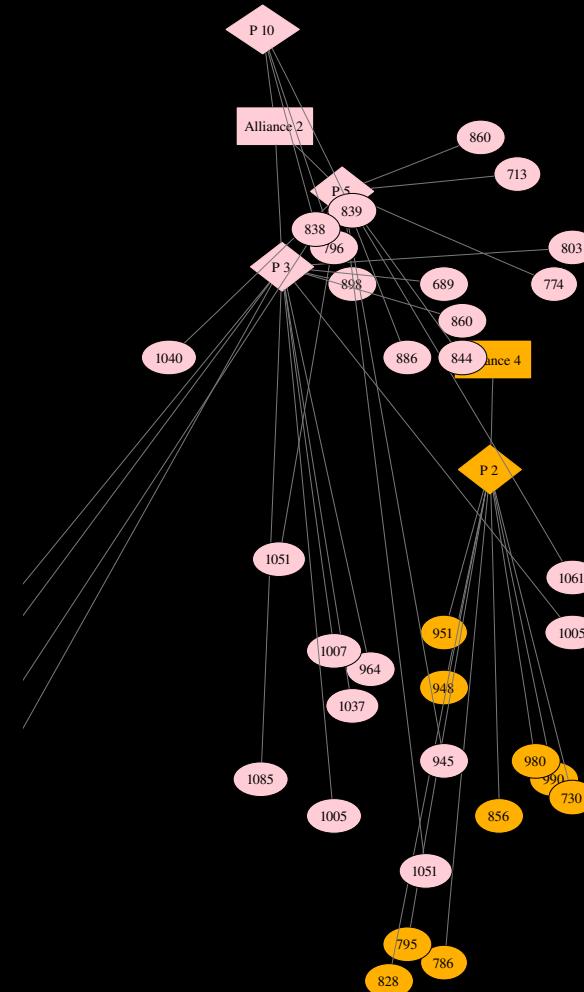
- ~10 alliances
- ~200 players
- ~600 cities

alliances color-coded

Can we build a model
of this world ?

Can we use it for playing
better ?

[Thon, Landwehr, De Raedt, ECML08]



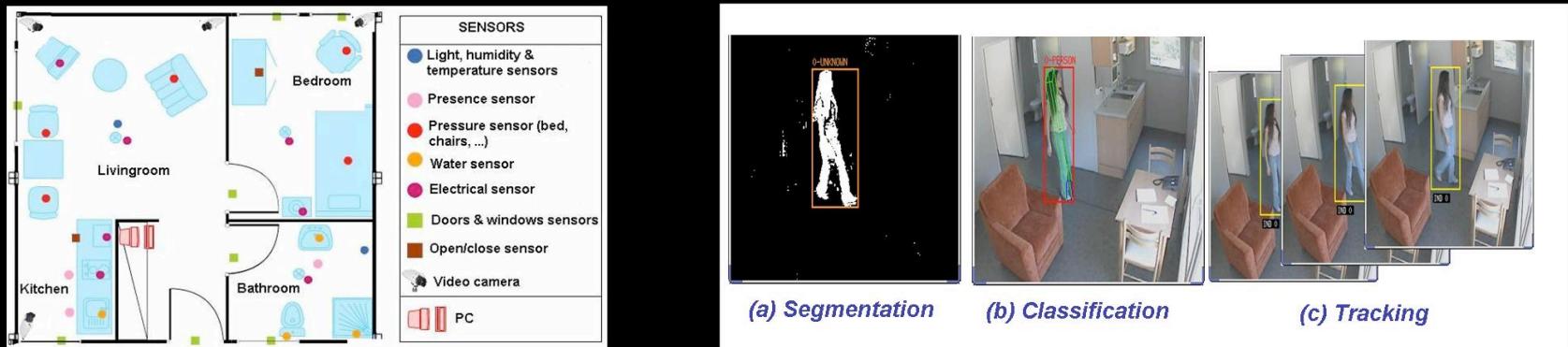
Emerging Data Sets

In many application areas :

- vision, surveillance, activity recognition, robotics, ...
- data in relational format are becoming available
- use of knowledge and reasoning is essential
- in Travian -- aka STRIPS representation

GerHome Example

Action and Activity Learning



(courtesy of Francois Bremond, INRIA-Sophia-Antipolis)

<http://www-sop.inria.fr/orion/personnel/Francois.Bremond/topicsText/gerhomeProject.html>

“Relational” Robotics

Robotics

- several lower level problems have been (more or less) solved
- time to bridge the gap between lower and higher level representations ?



[Kersting et al. Adv. Robotics 07]

Motivation for PLL

- Representational limitations of (standard)
 - probabilistic representations & logic
 - We: computational logic
- An example
 - graphs & networks
 - Essential : objects & relationships !! In many applications.

A. Motivation

A.2 *Representational Issues*

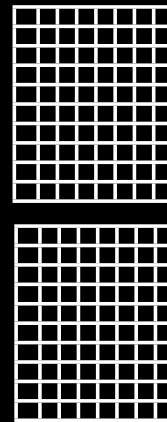
Represent the data Hierarchy

	<i>at</i>	<i>at</i>	<i>at</i>	<i>att</i>	<i>at</i>
<i>example</i>					

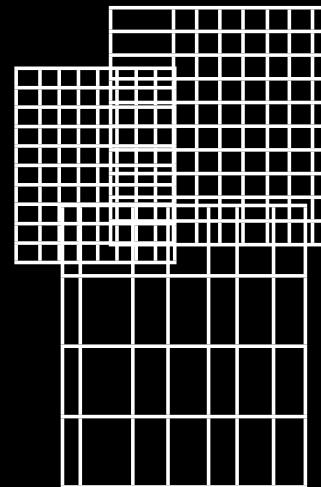
single-table
single-tuple
attribute-value

	<i>at</i>	<i>at</i>	<i>at</i>	<i>at</i>	<i>at</i>
<i>exampl</i>					

single-table
multiple-tuple
multi-instance



2 relations
edge / vertex
graphs & networks



multi-table
multiple-tuple
relational

Attribute-Value

	<i>at</i>	<i>at</i>	<i>at</i>	<i>att</i>	<i>at</i>
<i>example</i>					

Traditional Setting in Machine Learning
(cf. standard tools like Weka)

single-table
single-tuple
attribute-value

Multi-Instance

[Dietterich et al. AIJ 96]

	at	at	at	at	at
exampl					
exampl					
exampl					

*single-table
multiple-tuple
multi-instance*

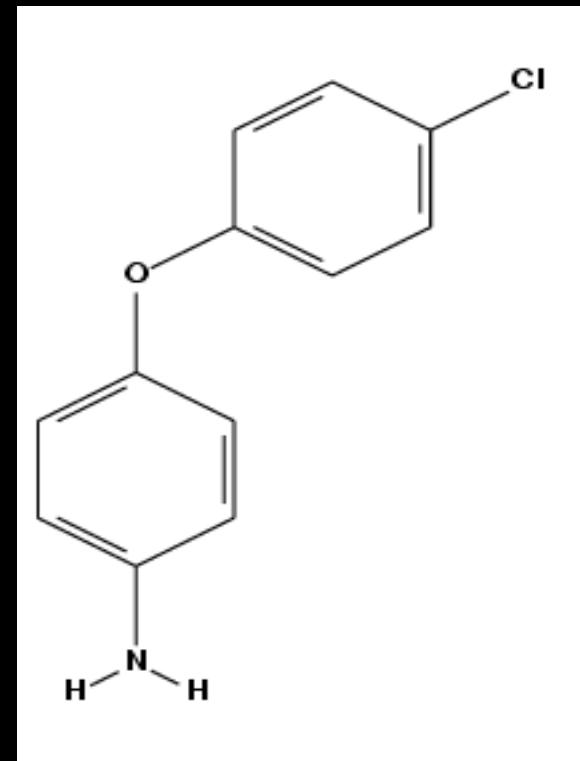
An example is positive if there exists a tuple in the example that satisfies particular properties

Boundary case between relational and propositional learning.

A lot of interest in past 10 years

Applications: vision, chemo-informatics, ...

Encoding Graphs



Encoding Graphs

atom(1,cl).

atom(2,c).

atom(3,c).

atom(4,c).

atom(5,c).

atom(6,c).

bond(3,4,s).

atom(7,c).

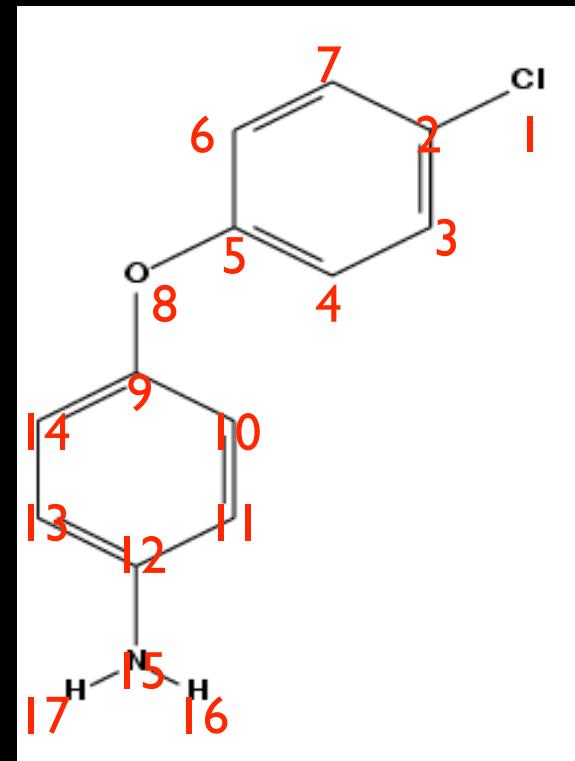
bond(1,2,s).

atom(8,o).

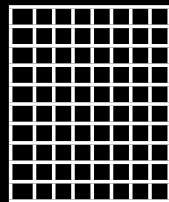
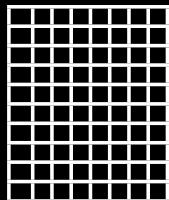
bond(2,3,d).

...

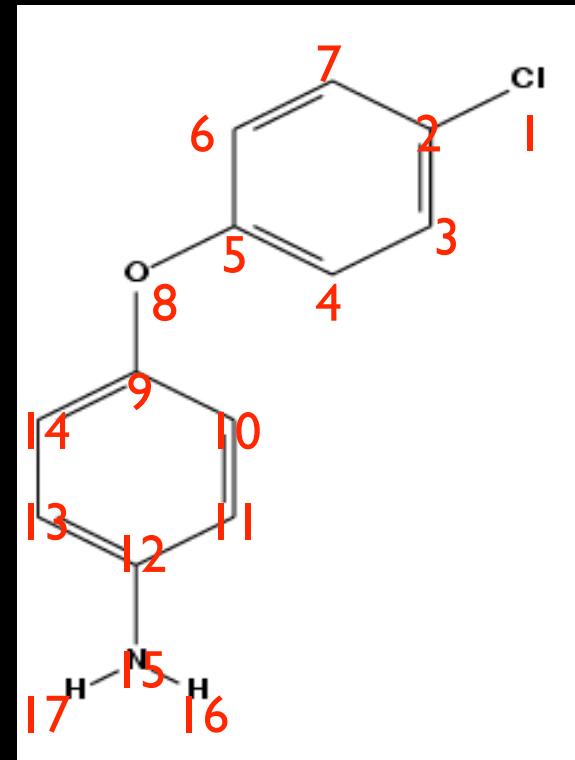
...



Encoding Graphs



*2 relations
edge / vertex
graphs & networks*



Encoding Graphs

atom(1,cl,21,0.297)

atom(2,c,21,0.187)

atom(3,c,21,-0.143)

atom(4,c,21,-0.143)

atom(5,c,21,-0.143)

atom(6,c,21,-0.143)

bond(3,4,s).

atom(7,c,21,-0.143)

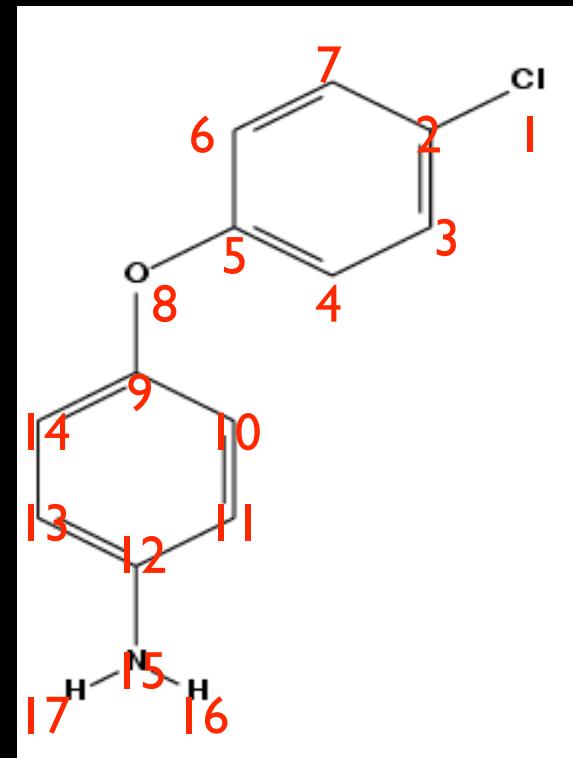
bond(1,2,s).

atom(8,o,52,0.98)

bond(2,3,d).

...

...



Note: add identifier for molecule

Encoding Knowledge

Use background knowledge in form of rules

- encode hierarchies

halogen(A):- atom(X,f)

halogen(A):- atom(X,cl)

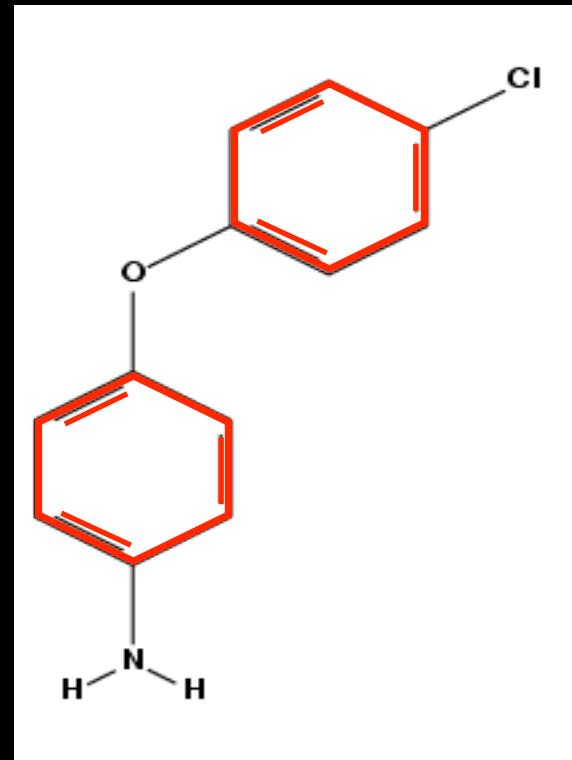
halogen(A):- atom(X,br)

halogen(A):- atom(X,i)

halogen(A):- atom(X,as)

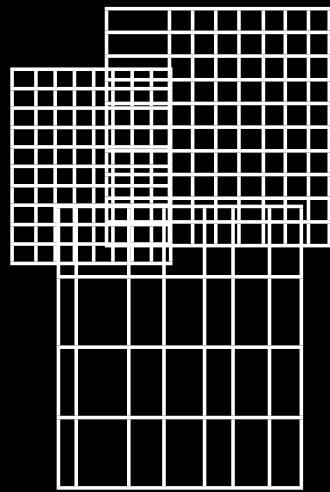
- encode functional group

benzene-ring :- ...

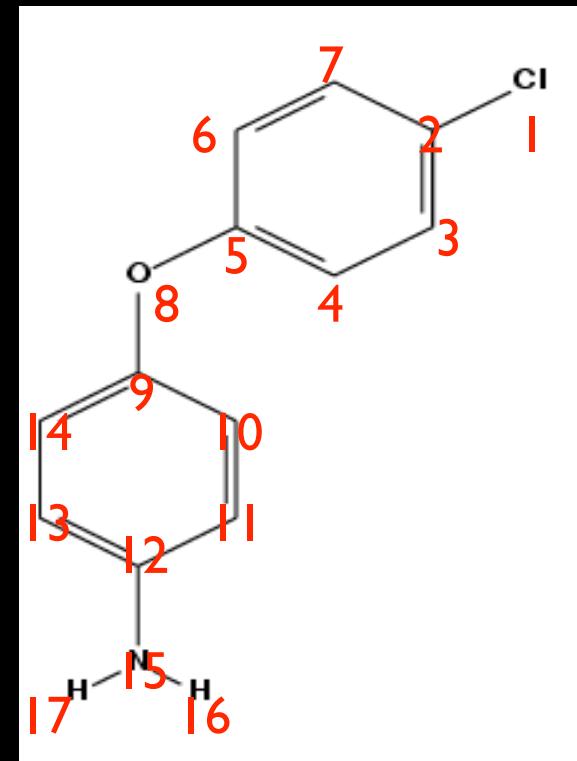


intentional versus extentional encodings

Relational Representation



*multi-table
multiple-tuple
relational*



Relational versus Graphs

Advantages Relational

- background knowledge in the form of rules, ontologies, features, ...
- relations of arity > 2 (but hypergraphs)
- graphs capture structure but annotations with many features/labels is non-trivial

Advantages Graphs

- efficiency and scalability
- full relational is more complex
- matrix operations

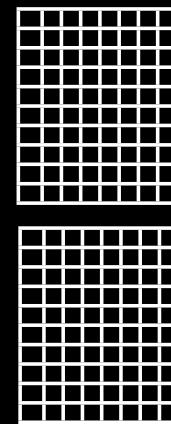
The Hierarchy

	<i>at</i>	<i>at</i>	<i>at</i>	<i>att</i>	<i>at</i>
<i>example</i>					

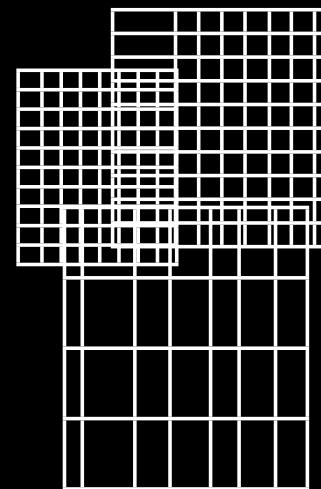
single-table
single-tuple
attribute-value

	<i>at</i>	<i>at</i>	<i>at</i>	<i>at</i>	<i>at</i>
<i>exampl</i>					

single-table
multiple-tuple
multi-instance



2 relations
edge / vertex
graphs & networks



multi-table
multiple-tuple
relational

Two questions

UPGRADING : Can we develop systems that work with richer representations (starting from systems for simpler representations)?

PROPOSITIONALISATION: Can we change the representation from richer representations to simpler ones ? (So we can use systems working with simpler representations)

B. Logic + Learning

B. Logic + Learning

B.I. Problem Settings

Traditional predictive ILP Problem

Given

- a set of positive and negative examples (**Pos, Neg**)
- a domain $L_0 \times L_1$
- a **target concept** $C : L_0 \times L_1 \rightarrow \{0, 1\}$
- a **language** L_h over $L_0 \times L_1$
- a **covers relation** on $L_0 \times L_1 \times L_h$

Concept-learning in a

logical/relational representation

Find

- A hypothesis h over L_h that covers all positive **Pos** and no negative **Neg** examples taking **B** into account

Three possible settings

- Learning from entailment
 - $\text{covers}(H,e)$ iff $H \models e$
- Learning from interpretations
 - $\text{covers}(H,e)$ iff e is a model for H
- Learning from proofs or traces.
 - $\text{covers}(H,e)$ iff e is proof given H

Three possible settings

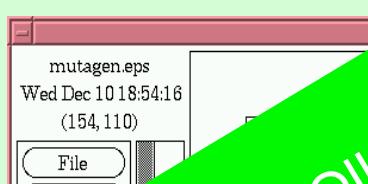
- Learning from entailment
 - $\text{covers}(H,e)$ iff $H \models e$
 - examples are facts / clauses

The Mutagenicity dataset

Background theory

```
molecule(225).  
logmutag(225,0.64).  
lumo(225,-1.785).  
logp(225,1.01).  
nitro(225,[f1_4,f1_8,f1_10,f1_9]).  
atom(225,f1_1,c,21,0.187).  
atom(''  
atom  
atom  
atom  
atom  
atom  
...  
ring  
hete
```

Applications



A screenshot of a Windows-style file explorer window. The title bar says 'File Explorer'. The main area shows a file named 'mutagen.eps' with the following details:
mutagen.eps
Wed Dec 10 18:54:16
(154, 110)
Below the details are two buttons: 'File' and 'Page'.

Applications

hypo

mut

Example

mutagen.eps
Wed Dec 10 18:54:16
(154, 110)

File
Page
M

e.g. FOIL, PROGOL, TIE

4-nitropenta[cd]pyrene

Inactive

Structural alert:

6-nitro-7,8,9,10-tetrahydrobenzo[a]pyrene

4-nitroindole

$\text{Y}=\text{Z}$

B \cup H

Three possible settings

- Learning from interpretations
 - $\text{covers}(H,e)$ iff e is a model for H
 - examples are Herbrand interpretations / state descriptions / possible worlds

Learning from interpretations

Examples

- **Positive:** { human(luc), human(lieve), male(luc), female(lieve) }
- **Negative:** { bat(dracula), male(dracula), vampire(dracula) }

Hypothesis

- $\text{human}(X) :- \text{male}(X)$

Learning from interpretations

Examples

- Positive: { human(luc), human(lieve), male(luc), male(lieve) }

Hypothesis (positives or

Systems:

e.g. LOGAN-H, Clau

dien, ICL, Warmr, ..

Applications

Enforcement of integrity constraints /
frequent patterns in relational databases

Three possible settings

- Learning from proofs or traces.
- $\text{covers}(H, e)$ iff e is proof given H

An example

```
sentence(A, B) :- noun_phrase(C, A, D), verb_phrase(C, D, B).  
noun_phrase(A, B, C) :- article(A, B, D), noun(C, D, C).  
verb_phrase(A, B, C) :- intransitive_verb(A, B, C).  
article(singular, A, B) :- terminal(A, B).  
article(singular, A, B) :- terminal(A, B).  
article(plural, A, B) :- terminal(A, B).  
noun(singular, A, B) :- terminal(A, B).  
noun(plural, A, B) :- terminal(A, B).  
intransitive_verb(simple, A, B).  
intransitive_verb(complex, A, B).  
terminal(simple, A, B).  
terminal(complex, A, B).
```

B) :- terminal(
B) :- termina
:- terminal,
:- termi

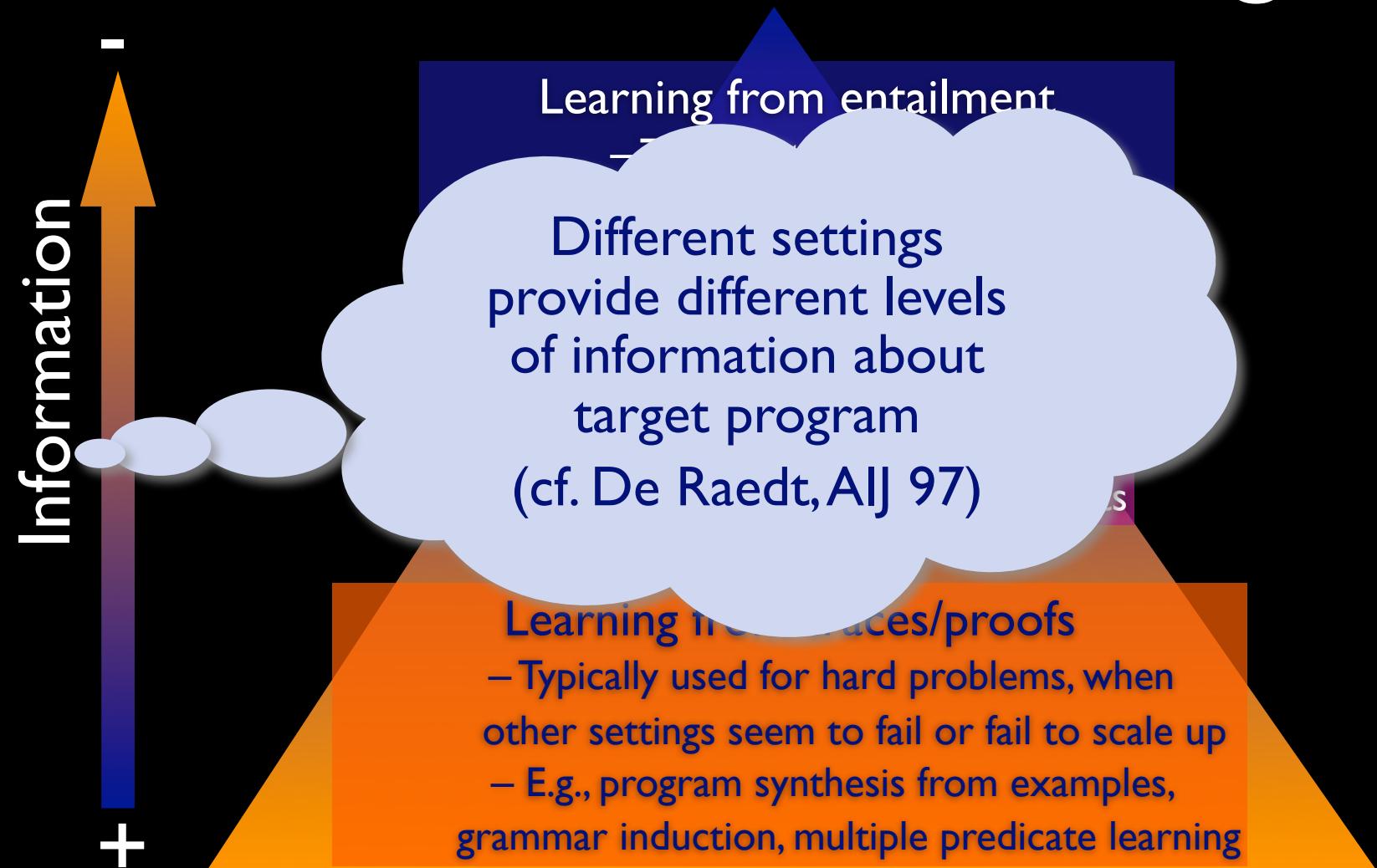
A Systems: Grammars, Tracy, MIS, ...
B, B)
cles, B

– Try grammar learning

Program synthesis

e.g. Spiro's MIS poses queries in order to reconstruct trace or proof

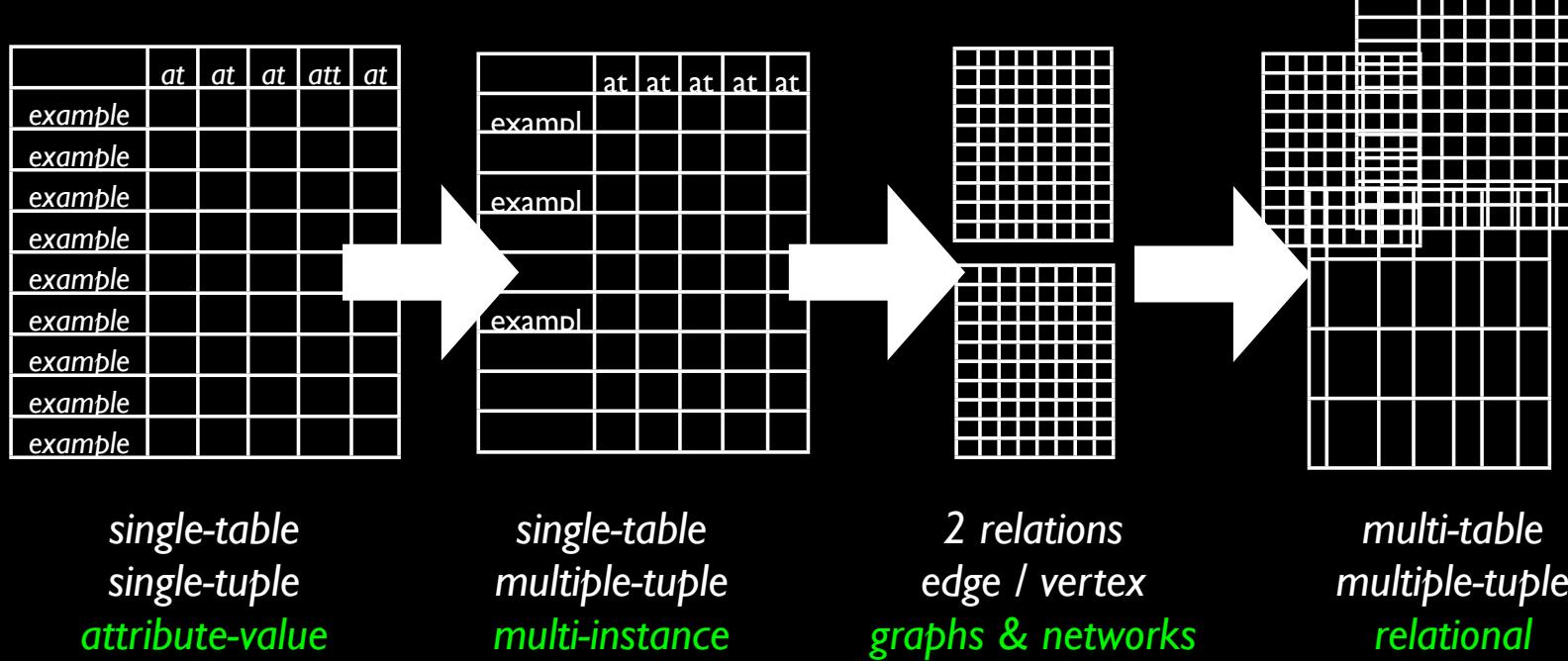
Use of different Settings



B. Logic + Learning

B.2. *Systems & Methodology*

Representational Hierarchy -- Systems



UPGRADING

Two messages

Logical and Relational Learning applies essentially to any machine learning and data mining task, not just concept-learning

- distance based learning, clustering, descriptive learning, reinforcement learning, bayesian approaches

there is a recipe that is being used to derive new LRL algorithms on the basis of propositional ones

- not the only way to LRL

Learning Tasks

- rule-learning & decision trees [Quinlan 90], [Blockeel 96]
- frequent and local pattern mining [Dehaspe 98]
- distance-based learning (clustering & instance-based learning) [Horvath, 01], [Ramon 00]
- probabilistic modeling (cf. statistical relational learning)
- reinforcement learning [Dzeroski et al. 01]
- kernel and support vector methods

Logical and relational representations can (and have been) used for all learning tasks and techniques

The RECIPE

Start from well-known propositional learning system

Modify representation and operators

- e.g. generalization/specialization operator, similarity measure, ...
- often use theta-subsumption as framework for generality

Build new system, retain as much as possible from propositional one

LRL Systems and techniques

FOIL ~ CN2 – Rule Learning (Quinlan MLJ 90)

Tilde ~ C4.5 – Decision Tree Learning (Blockeel & DR AIJ 98)

Warmr ~ Apriori – Association rule learning (Dehaspe 98)

Progol ~~ AQ – Rule learning (Muggleton NGC 95)

Graph miners ...

A case : FOIL

Learning from entailment -- the setting

$$B \cup H \models e$$

Given

```
molecule(225).  
logmutag(225, 0.64).  
lumo(225, -1.785).  
logp(225, 1.01).  
nitro(225, [f1_4, f1_8, f1_10, f1_9]).  
atom(225, f1_1, c, 21, 0.187).  
atom(225, f1_2, c, 21, -0.143).  
atom(225, f1_3, c, 21, -0.143).  
atom(225, f1_4, c, 21, -0.013).  
atom(225, f1_5, o, 52, -0.043).  
...  
ring_size_5(225, [f1_5, f1_1, f1_2, f1_3, f1_4]).  
hetero_aromatic_5_ring(225, [f1_5, f1_1, f1_2, f1_3, f1_4]).
```

background

```
bond(225, f1_1, f1_2, 7).  
bond(225, f1_2, f1_3, 7).  
bond(225, f1_3, f1_4, 7).  
bond(225, f1_4, f1_5, 7).  
bond(225, f1_5, f1_1, 7).  
bond(225, f1_8, f1_9, 2).  
bond(225, f1_8, f1_10, 2).  
bond(225, f1_1, f1_11, 1).  
bond(225, f1_11, f1_12, 2).  
bond(225, f1_11, f1_13, 1).
```

mutagenic(225), ...

examples

$$B \cup H \models e$$

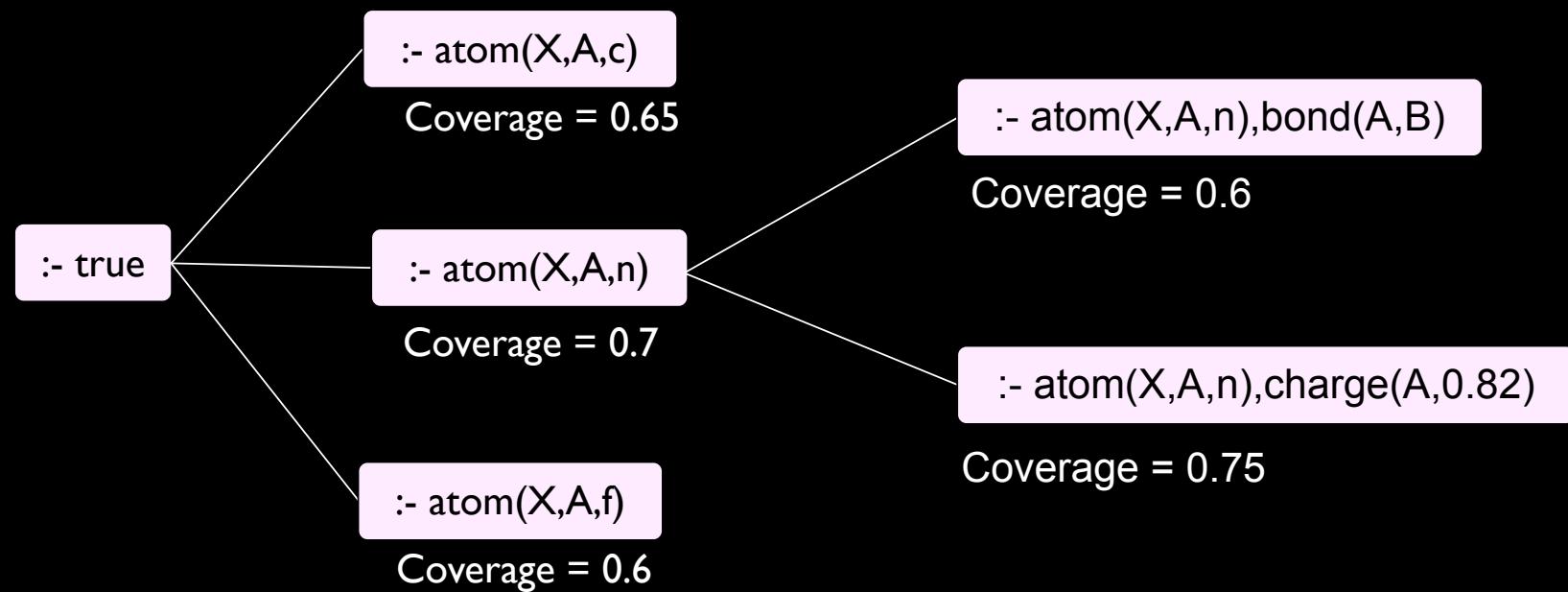
Find

```
mutagenic(M) :- nitro(M, R1), logp(M, C), C > 1 .
```

rules

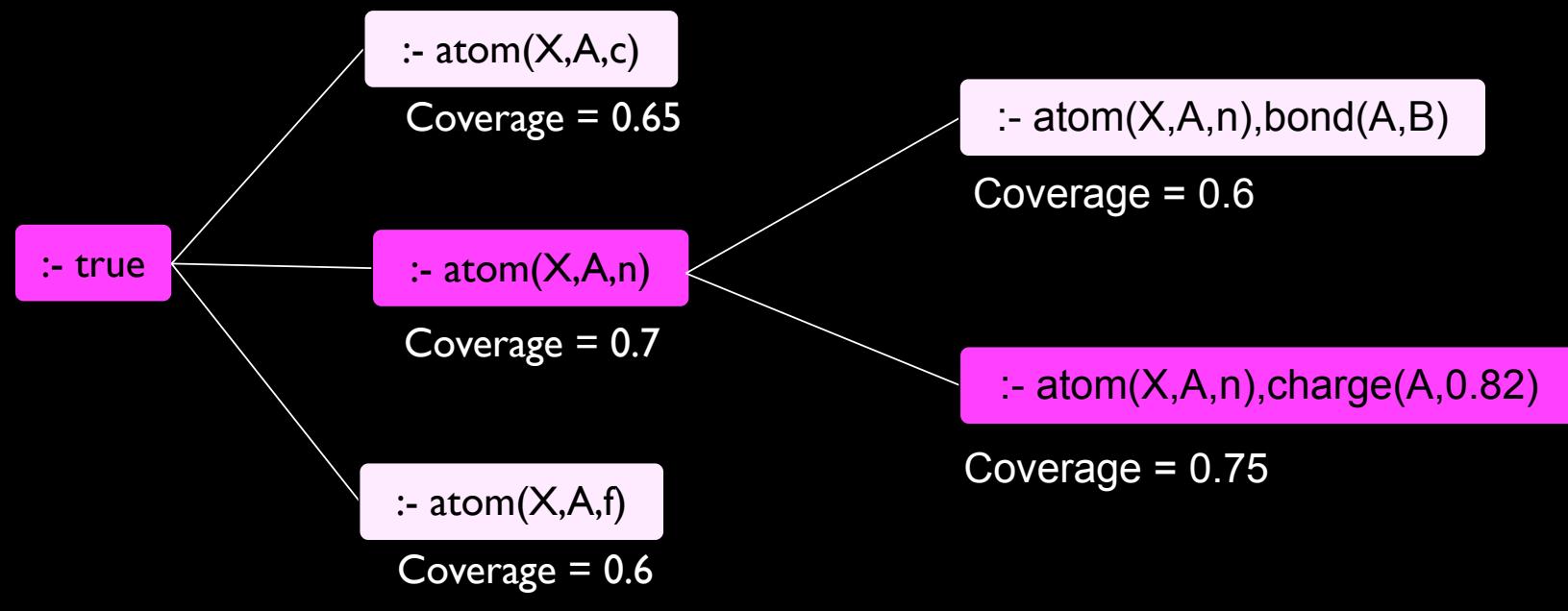
Searching for a rule

Greedy separate-and-conquer for rule set
Greedy general-to-specific search for single rule

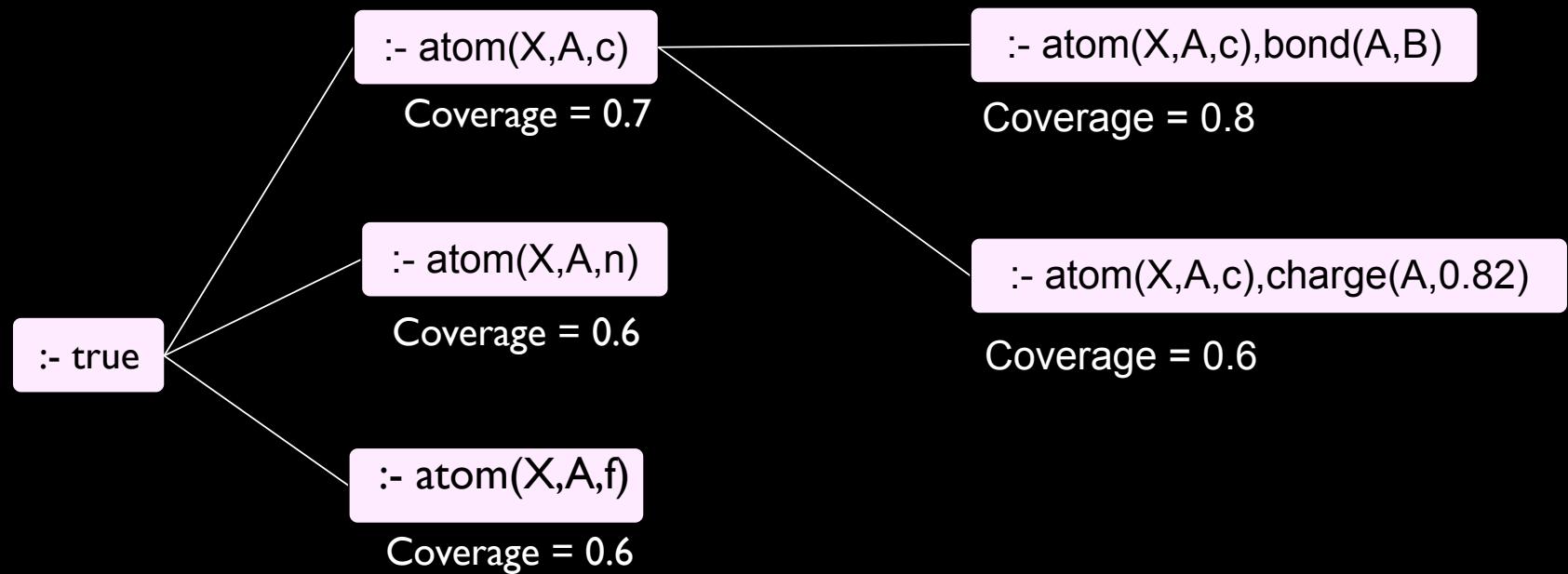


Searching for a rule

Greedy separate-and-conquer for rule set
Greedy general-to-specific search for single rule

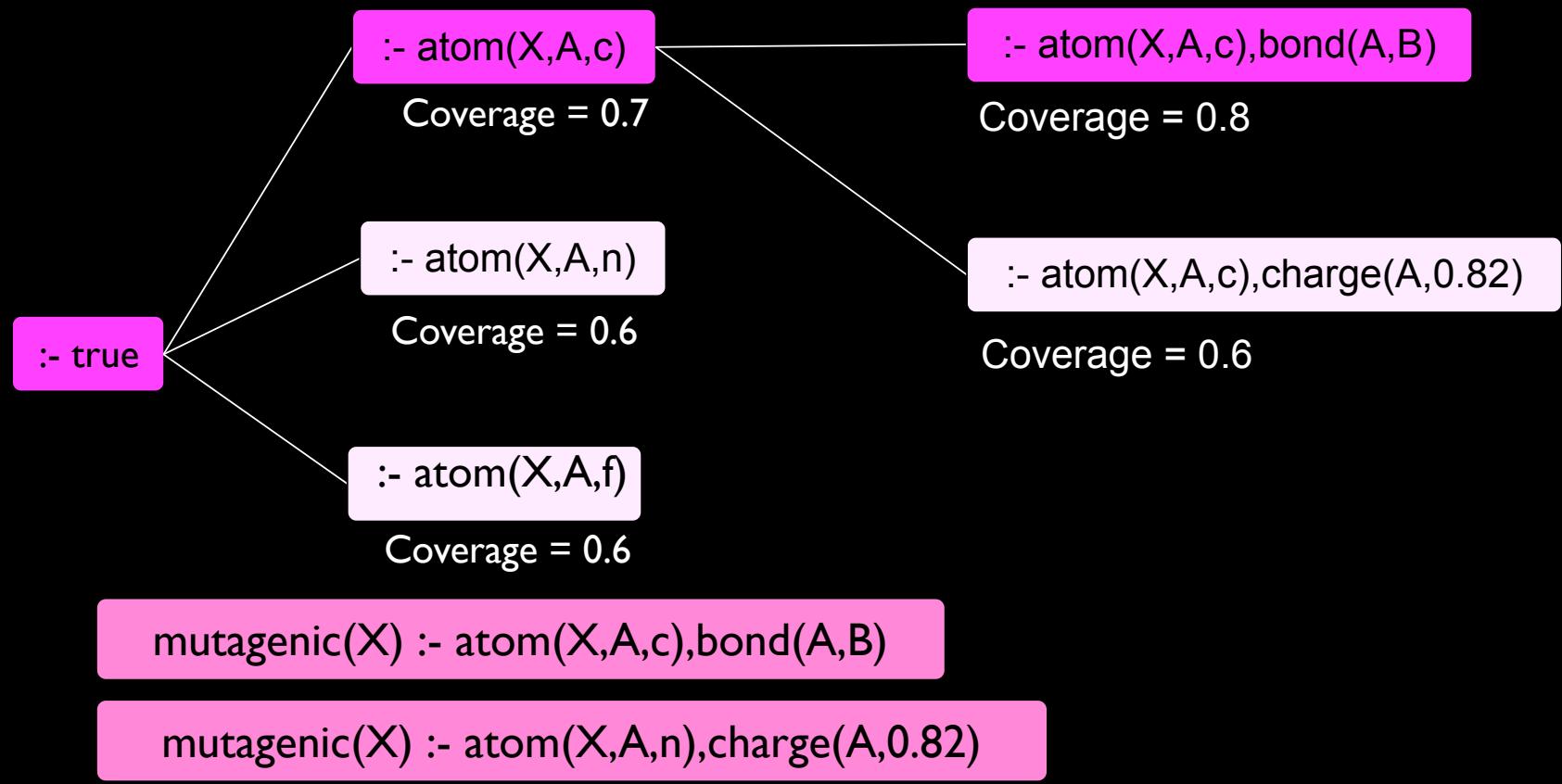


FOIL



mutagenic(X) :- atom(X,A,n),charge(A,0.82)

FOIL



mutagenic(X) :- atom(X,A,c),charge(A,0.45)

mutagenic(X) :- atom(X,A,c),bond(A,B)

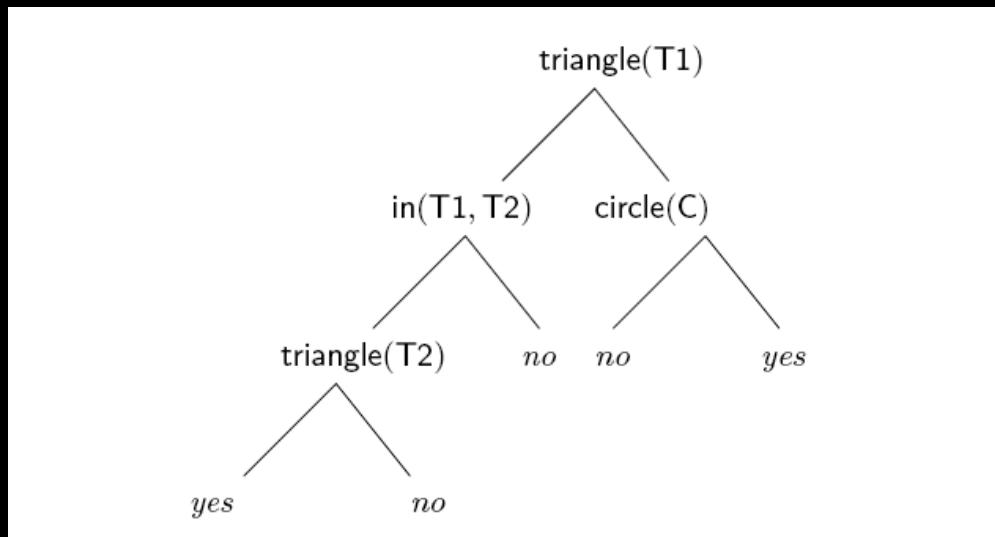
mutagenic(X) :- atom(X,A,n),charge(A,0.82)

Tilde

Logical Decision Trees (Blockeel & De Raedt Alj 98)

Negative Examples					
○	○	□ □	▽ ○ ▽	○ ▽ ▽	▽ △
□ ○ △	△	○ □	○ ○ △	○ □ ▽	○ ○ □
○ ▽	○ ○	▽	○ ○ ○	○ △	△ ▽ ○
Positive Examples					
△ △ □	○ △ △	○ ▽	▽	△ △ □	○ △ △
△ ○ ▽	△ ▽ △	△ ▽	▽ ○ ▽	△ □ ▽	△ ▽ □

A logical decision tree

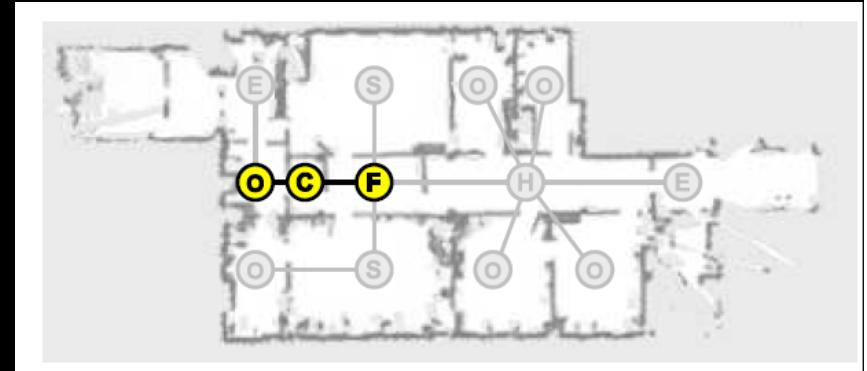
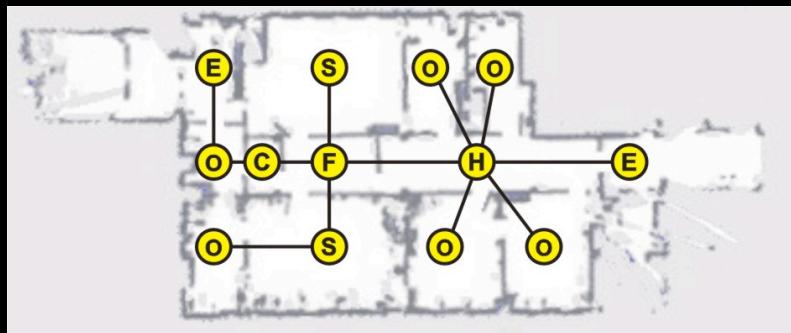


```
IF triangle(T1), in(T1, T2), triangle(T2) THEN Class = yes  
ELSIF triangle(T1), in(T1, T2) THEN Class = no  
ELSIF triangle(T1) THEN Class = no  
ELSIF circle(C) THEN Class = no  
ELSE Class = yes
```

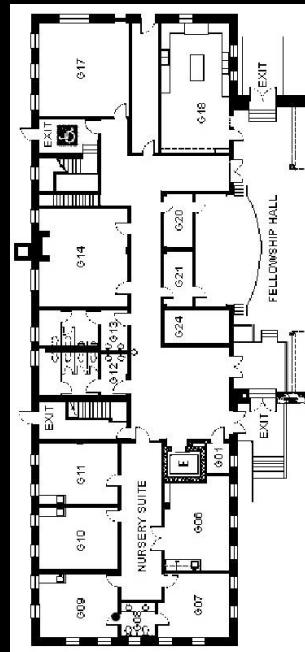
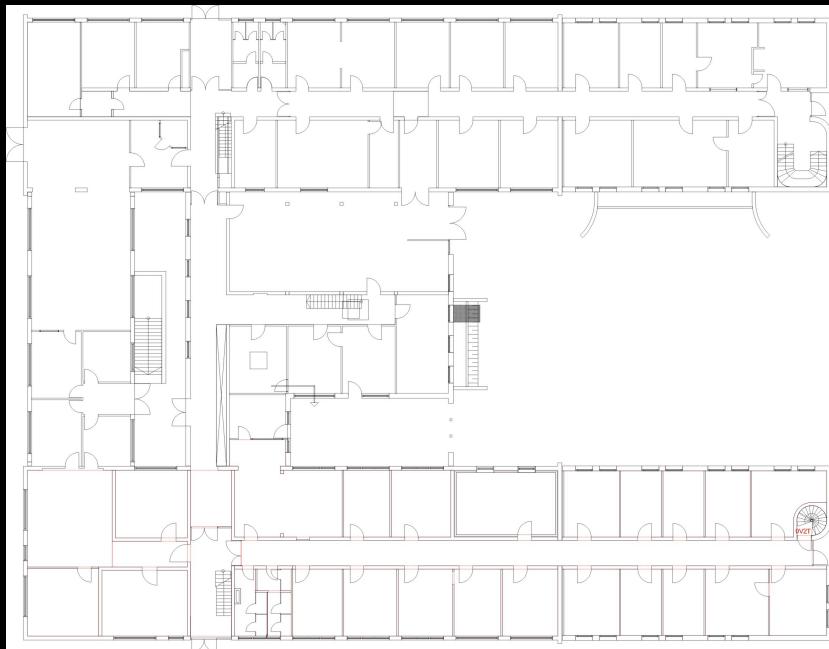
Search in unknown environment

Popular Algorithms

- Uninformed LRTA* [Koenig]
- Depth-first



Buildings



Relational MDPs

Abstract states:

$\text{in}(X), \text{conn}(X, Y),$
 $\text{room}(X), \text{hPassage}(Y)$

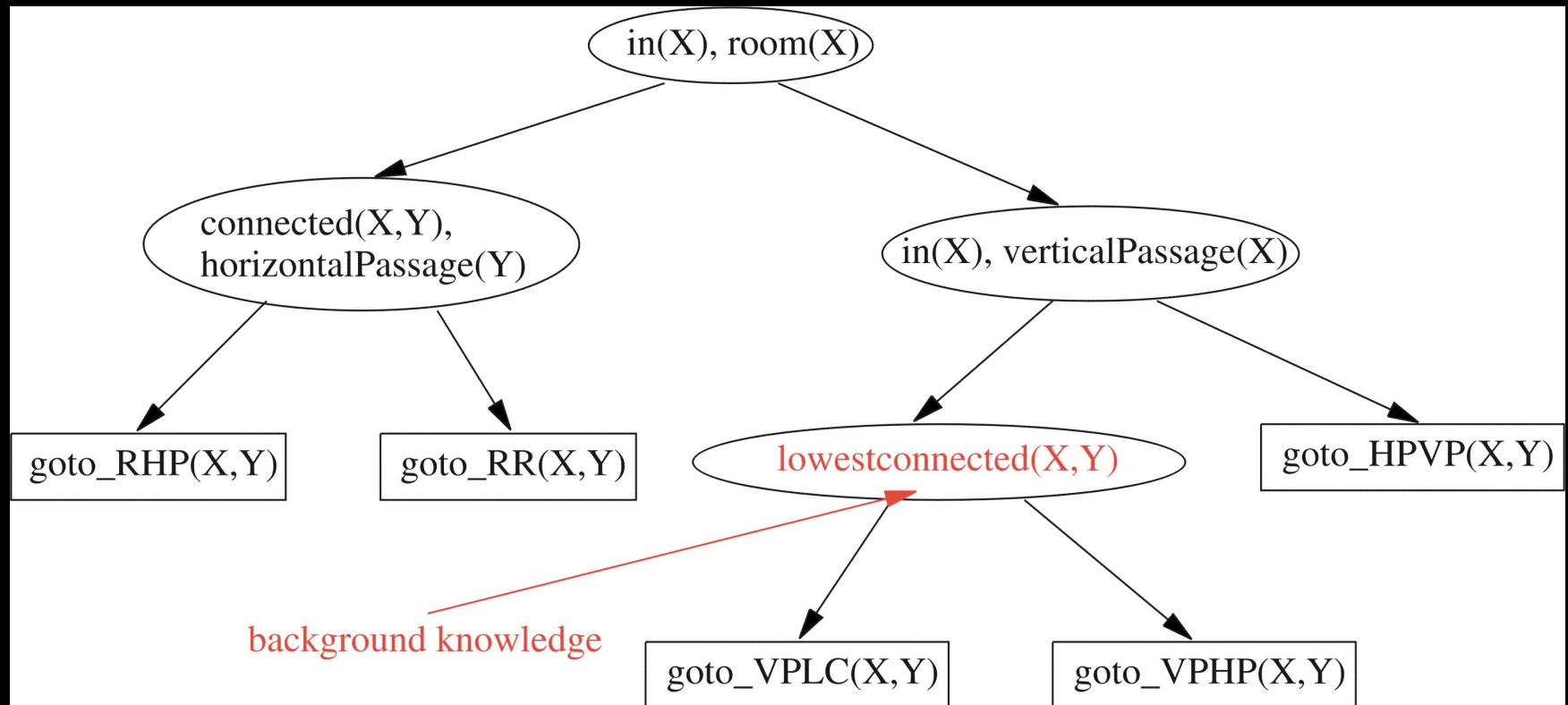
Abstract actions:

$\text{in}(Y), \text{conn}(X, Y),$
 $\text{hPassage}(Y), \text{room}(X),$ $\xleftarrow{\text{p: goto_RHP}(X, Y)}$ $\text{in}(X), \text{conn}(X, Y),$
 $X \neq Y$ $X \neq Y$

Abstract reward model:

10 \neg $\text{in}(l_1)$
0 \neg true

Policy



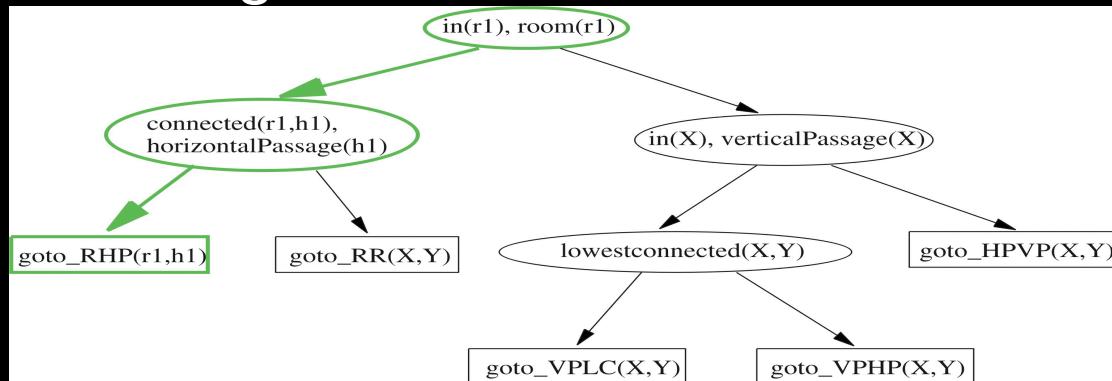
logical decision trees [Blockeel & De Raedt AIJ 98]

Using a Relational Policy

1. Start in State

$\underbrace{\text{in}(r_1), \text{conn}(r_1, h_1), \text{conn}(h_1, s_1)}_{\text{building description}}$

2. Evaluate logical decision tree



3. Execute Action

4. Observe New State

5. Repeat if goal not yet reached.

WARMR

PARTICIPANT Table

NAME	JOB	COMPANY	PARTY	R_NUMBER
adams	researcher	scuf	no	23
blake	president	jvt	yes	5
king	manager	ucro	no	78
miller	manager	jvt	yes	14
scott	researcher	scuf	yes	94
turner	researcher	ucro	no	81

COMPANY Table

COMPANY	TYPE
jvt	commercial
scuf	university
ucro	university

COURSE Table

COURSE	LENGTH	TYPE
cso	2	introductory
erm	3	introductory
so2	4	introductory
srw	3	advanced

SUBSCRIPTION Table

NAME	COURSE
adams	erm
adams	so2
adams	srw
blake	cso
blake	erm
king	cso
king	erm
king	so2
king	srw
miller	so2
scott	erm
scott	srw
turner	so2
turner	srw

What to count ? Keys.

PARTICIPANT Table

NAME	JOB	COMPANY	PARTY	R_NUMBER
adams	researcher	scuf	no	23
blake	president	jvt	yes	5
king	manager	ucro	no	78
miller	manager	jvt	yes	14
scott	researcher	scuf	yes	94
turner	researcher	ucro	no	81

SUBSCRIPTION Table

NAME	COURSE
adams	erm
adams	so2
adams	srw
blake	cso
blake	erm
king	cso
king	erm
king	so2
king	srw
miller	so2
scott	erm
scott	srw
turner	so2
turner	srw

COMPANY Table

COMPANY	TYPE
jvt	commercial
scuf	university
ucro	university

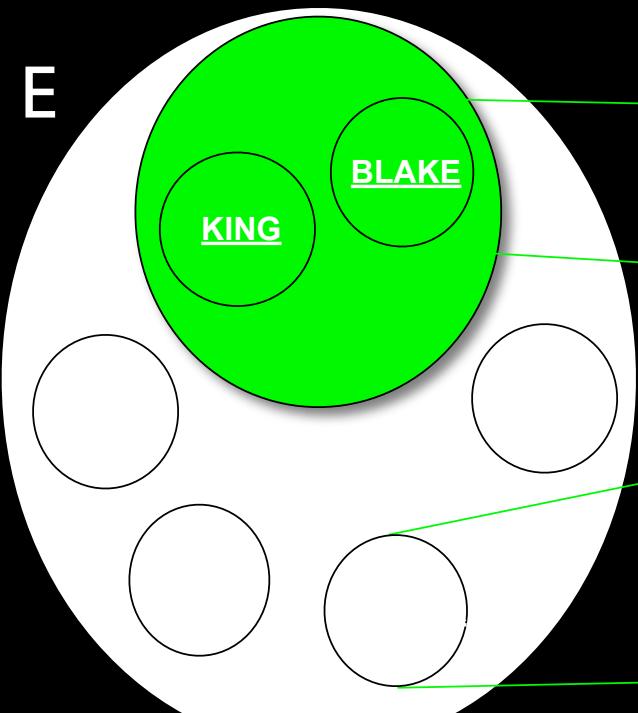
COURSE Table

COURSE	LENGTH	TYPE
cso	2	introductory
erm	3	introductory
so2	4	introductory
srw	3	advanced

Coverage

H: “*there is a subscription for a course of length less than 3*”

```
?- subscription(K,C), course(K,C,L,_), (L < 3).
```



```
participant(blake,president,jvt,yes,5).  
subscription(blake,cso).  
subscription(blake,erm).
```

```
course(cso,2,introductory).  
course(erm,3,introductory).  
course(so2,4,introductory).  
course(srw,3,advanced).
```

```
participant(turner,researcher,scuf,no,23).  
subscription(turner,erm).  
subscription(turner,so2).  
subscription(turner,srw).
```

K=blake
C = cso
L = 2
Yes

K=turner
No

Descriptive Data Mining

multi-relational database association rules over multiple relations

PARTICIPANT Table				
NAME	JOB	COMPANY	PARTY	R_NUMBER
adams	researcher	scuf	no	23
blake	president	jvt	yes	5
king	manager	ucro	no	78
miller	manager	jvt	yes	14
scott	researcher	scuf	yes	94
turner	researcher	ucro	no	81

COMPANY Table		COURSE Table		
COMPANY	TYPE	COURSE	LENGTH	TYPE
jvt	commercial	cso	2	introductory
scuf	university	erm	3	introductory
ucro	university	so2	4	introductory
		srw	3	advanced

NAME	COURSE
adams	erm
adams	so2
adams	srw
blake	cso
blake	erm
king	cso
king	erm
king	so2
king	srw
miller	so2
scott	erm
scott	srw
turner	so2
turner	srw

Selected,
Preprocessed,
and Transformed
Data

**"IF ?- participant(P,C,PA,X),
course(P,Y,advanced)
THEN ?-PA = no."**

**support: 20 %
confidence: 75 %**

**"IF participant follows an advanced course
THEN she skips the welcome party"**

**support: 20 %
confidence: 75 %**

Warmr

First order association rule :

- IF Query1 THEN Query2
- Shorthand for
- IF Query1 THEN Query1 and Query2
- to obtain variable bindings

IF ?- participant(P,C,PA,X), course(P,Y,advanced) THEN PA=no

IF ?- participant(P,C,PA,X), course(P,Y,advanced) succeeds for P
THEN ?- participant(P,C,PA,X), course(P,Y,advanced), PA=no

Warmr ~ Apriori

Works as Apriori :

- keeping track of frequent and infrequent queries
- order queries by theta-subsumption
- using special mechanism (bias) to declare type of queries searched for

Generalizes many of the specific variants of Apriori : item-sets, episodes, hierarchies, intervals, ...

Upgrading

Key ideas / contributions FOIL + Tilde + Warmr

- determine the representation of examples and hypotheses
- select the right type of coverage and generality (subsumption)
- keep existing algorithm (CN2 or C4.5) but replace operators (often theta-subsumption)
- keep search strategy
- fast implementation

Distance Based Learning

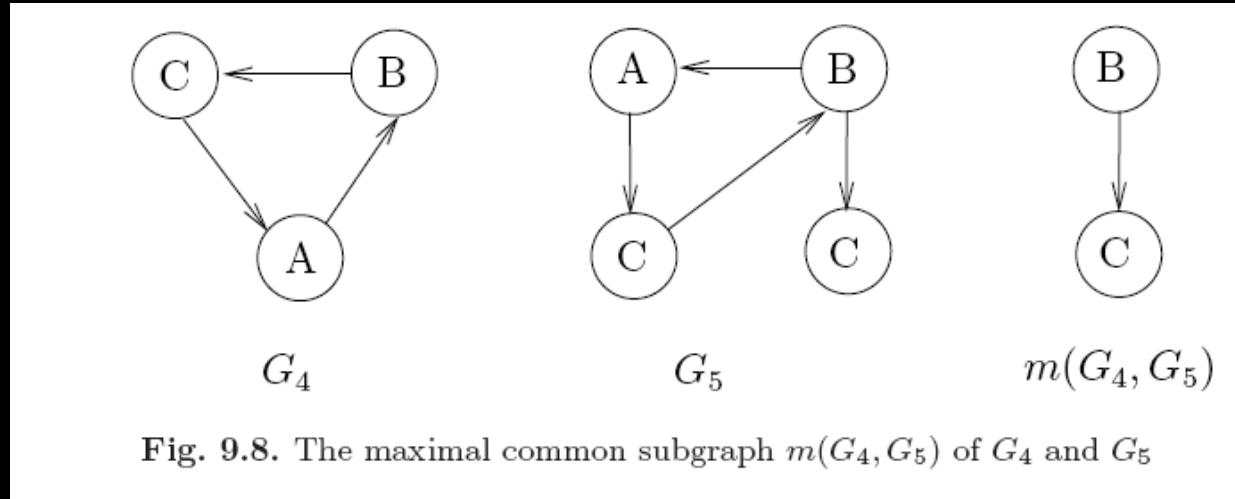


Fig. 9.8. The maximal common subgraph $m(G_4, G_5)$ of G_4 and G_5

$$d_{graph}(G_1, G_2) = |G_1| + |G_2| - 2|mcs(G_1, G_2)|$$

Distances

$|h|$ = the size of a hypothesis h

$|g| \leq |s|$ if g is more general than s

$\text{mgg}(t,u)$ = set of minimally general generalizations

$|\text{mgg}(t,u)|$ = max size of element in $\text{mgg}(t,u)$

under mild conditions the following is a metric

$$d_{\preceq}(x, y) = |x| + |y| - 2|\text{mgg}(x, y)|$$

The RECIPE

Relevant for ALL levels of the hierarchy

Still being applied across data mining,

- mining from graphs, trees, and sequences

Works in both directions

- upgrading and downgrading !!!

Mining from graphs or trees as downgraded Relational Learning

Many of the same problems / solutions apply to graphs as to relational representations

From Upgrading to Downgrading

Work at the right level of representation

- trade-off between **expressivity & efficiency**

The **old** challenge: **upgrade** learning techniques for simpler representations to richer ones.

The **new** challenge: **downgrade** more expressive ones to simpler ones for efficiency and scalability; e.g. graph miners.

Note: systems using rich representations form a **baseline**, and can be used to test out ideas.

Relevant also for ALL machine learning and data mining tasks

Learning Tasks

Logical and relational representations can (and have been) used for all learning tasks and techniques

- rule-learning & decision trees
- frequent and local pattern mining
- distance-based learning (clustering & instance-based learning)
- probabilistic modeling (cf. statistical relational learning)
- reinforcement learning
- kernel and support vector methods

C. Probabilistic Logic Learning

PLL: What Changes ?

- Clauses annotated with probability labels
 - E.g. in Sato's Prism, Muggleton's SLPs, Kersting and De Raedt's BLPs, ...
- Prob. covers relation $\text{covers}(e, H \cup B) = P(e | H, B)$
 - Likelihood of example given background and hypothesis
 - Probability distribution P over the different values e can take; so far only (true, false)
 - impossible examples have probability 0
- Knowledge representation issue
 - Define probability distribution on examples / individuals
 - What are these examples / individuals ?

Probabilistic generative ILP Problem

Given

- a set of examples E
- a background theory B
- a language Le to represent examples
- a language Lh to represent hypotheses
- a probabilistic covers P relation on $Le \times Lh$

Find

- hypothesis h^* maximizing some score based on the probabilistic covers relation; often some kind of maximum likelihood

PLL: Three Settings

- Probabilistic learning from interpretations
 - Bayesian logic programs, Koller's PRMs, Domingos' MLNs, Vennekens' LPADs
- Probabilistic learning from entailment
 - Muggleton's Stochastic Logic Programs, Sato's Prism, Poole's ICL, De Raedt et al.'s ProbLog
- Probabilistic learning from proofs
 - Learning the structure of SLPs; a tree-bank grammar based approach, Anderson et al.'s RMMs, Kersting et al.

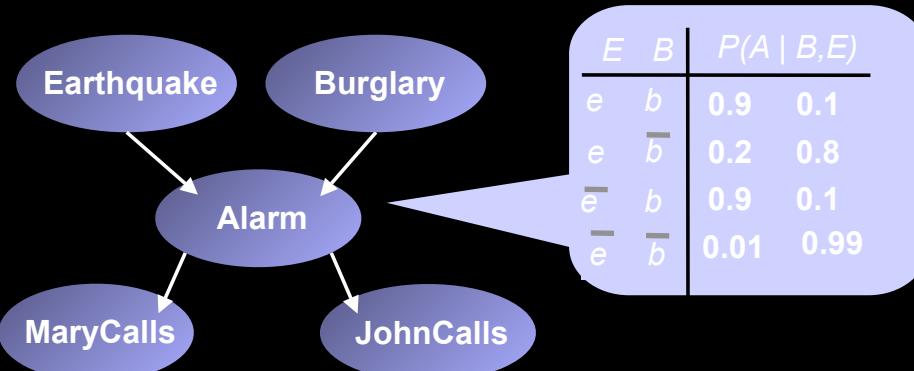
C. Probabilistic Logic Learning

C. I. Learning from Interpretations

Graphical Model

- Defines a joint probability distribution in a compact way; attractive from a computational perspective
- The structure encodes some conditional independency assumptions.
- Two types :
 - Directed -- Bayesian Network
 - Undirected -- Markov Network

Bayesian Networks



$$P(E, B, A, J, M) = P(E).P(B).P(A|E, B).P(J|A).P(M|A)$$

Assumption 1 (cf. Russell and Norvig [2004]) Each node X_i in the graph is conditionally independent of any subset \mathbf{A} of nodes that are non-descendants of X_i given a joint state of $\text{Pa}(X_i)$, that is, $P(X_i | \mathbf{A}, \text{Pa}(X_i)) = P(X_i | \text{Pa}(X_i))$.

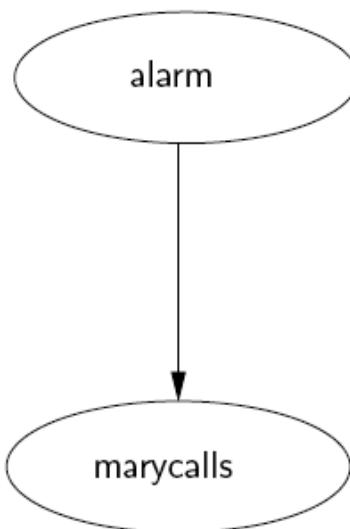
Using the Joint Prob. Distribution,
any query $P(Q | \text{evidence})$ can be answered

Learning

Two aspects:

- parameters, e.g. max likelihood
- structure learning

Toy Example



$$\frac{P(\text{alarm})}{(\lambda_0, 1 - \lambda_0)}$$

$$\begin{array}{ll} \text{alarm } P(\text{marycalls}) \\ \hline \text{true } (\lambda_1, 1 - \lambda_1) \\ \text{false } (\lambda_2, 1 - \lambda_2) \end{array}$$

Fig. 8.6. A simple Bayesian network. $P(a) = \lambda_0$; $P(m|a) = \lambda_1$; $P(m|\neg a) = \lambda_2$

Parameter Estimation

						complete data set	simply counting
A1	A2	A3	A4	A5	A6	X1	
true	true	false	true	false	false	X2	
false	true	true	true	false	false	⋮	
...	XM	
true	false	false	false	true	true		

Parameter Estimation

						incomplete data set
A1	A2	A3	A4	A5	A6	Real-world data: states of some random variables are missing E.g. medical diagnose: not all patient are subjects to all test Parameter reduction, e.g. clustering, ...
true	true	?	true	false	false	
?	true	?	?	false	false	
...	
true	false	?	false	true	?	

? missing value

On the Alarm Net

	alarm	marycalls
e_1	true	true
e_2	true	true
e_3	false	false
e_4	false	true
...		

Table 8.3. A completely observed data set.

	alarm	marycalls
e_1	true	true
e_2	true	?
e_3	false	false
e_4	?	true

Table 8.4. A data set with missing values.

maximize likelihood:
 $\max P(E | M)$

Derivation (fully obs)

$$P(a, \neg m) = \lambda_0(1 - \lambda_2)$$

Generalize :

$$P(E|M) = \lambda_0^{|a|} \cdot (1 - \lambda_0)^{|\neg a|} \cdot \lambda_1^{|a \wedge m|} \cdot (1 - \lambda_1)^{|a \wedge \neg m|} \cdot \lambda_2^{| \neg a \wedge m |} \cdot (1 - \lambda_2)^{|\neg a \wedge \neg m|}$$

Use log likelihood:

$$\begin{aligned} L = \log P(E|M) &= |a| \cdot \log \lambda_0 + |\neg a| \cdot \log(1 - \lambda_0) + \\ &\quad |a \wedge m| \cdot \log \lambda_1 + |a \wedge \neg m| \cdot \log(1 - \lambda_1) + \\ &\quad |\neg a \wedge m| \cdot \log \lambda_2 + |\neg a \wedge \neg m| \cdot \log(1 - \lambda_2) \end{aligned}$$

Derivation (fully obs) (2)

$$\begin{aligned}\frac{\partial L}{\partial \lambda_0} &= \frac{|a|}{\lambda_0} - \frac{|\neg a|}{1 - \lambda_0} \\ \frac{\partial L}{\partial \lambda_1} &= \frac{|a \wedge m|}{\lambda_1} - \frac{|a \wedge \neg m|}{1 - \lambda_1} \\ \frac{\partial L}{\partial \lambda_2} &= \frac{|\neg a \wedge m|}{\lambda_2} - \frac{|\neg a \wedge \neg m|}{1 - \lambda_2}\end{aligned}$$

Yielding as solutions

$$\lambda_0 = \frac{|a|}{|a| + |\neg a|} \tag{1}$$

$$\lambda_1 = \frac{|a \wedge m|}{|a \wedge m| + |a \wedge \neg m|} \tag{2}$$

$$\lambda_2 = \frac{|\neg a \wedge m|}{|\neg a \wedge m| + |\neg a \wedge \neg m|} \tag{3}$$

So, simply counting !

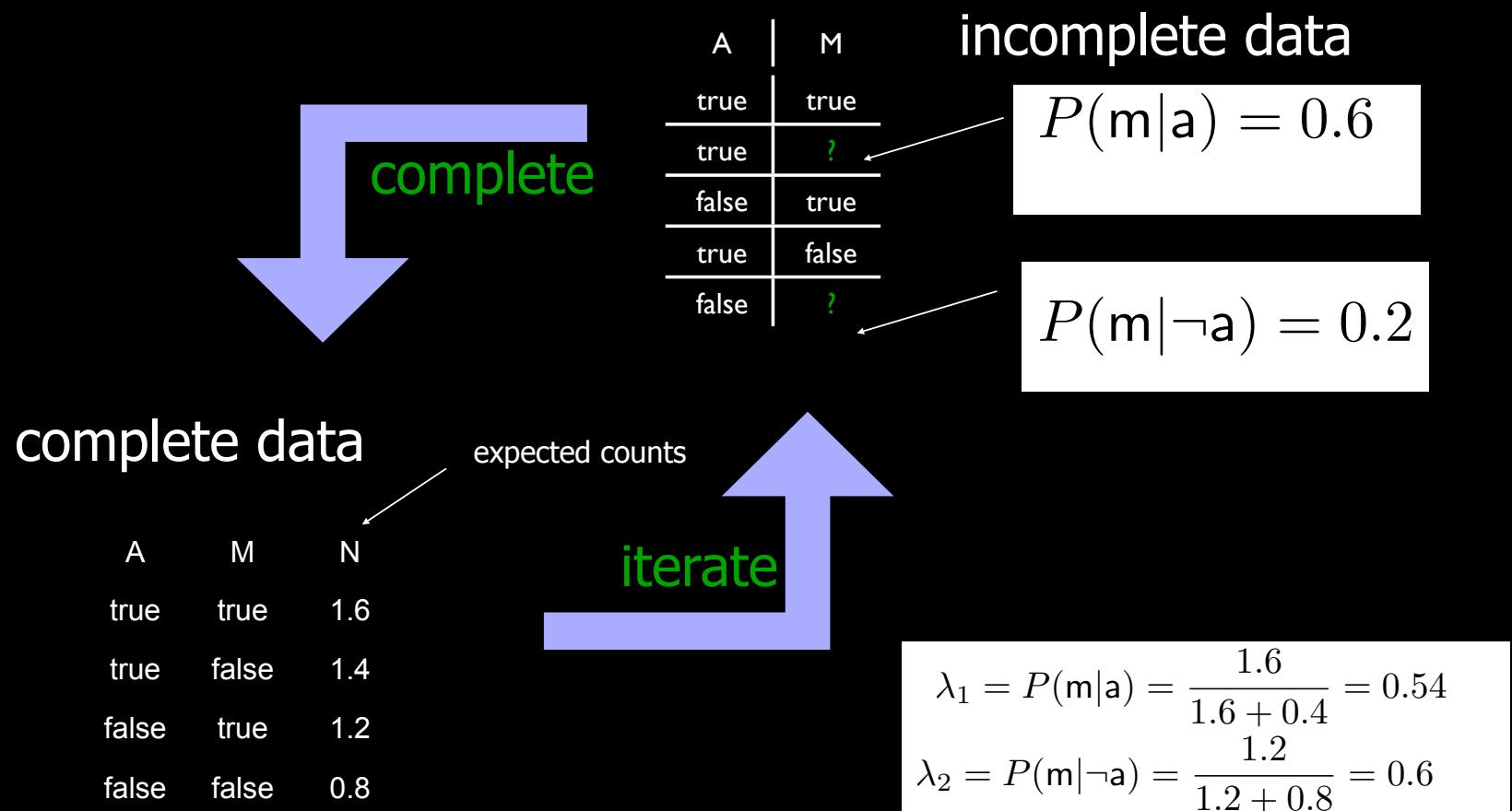
Parameter Estimation (part. obs)

Idea Expectation Maximization

Incomplete data ?

- Complete data (Imputation)
 - most probable?, average?, ... value
- Count
- Iterate

EM Idea: Complete the data



Parameter Tying

$$\begin{aligned} P(E|M) = & \lambda_0^{|a|} \cdot (1 - \lambda_0)^{|\neg a|} \cdot \lambda_1^{|a \wedge m|} \cdot (1 - \lambda_1)^{|a \wedge \neg m|} \cdot \lambda_2^{\neg a \wedge m} \cdot (1 - \lambda_2)^{\neg a \wedge \neg m} \\ & \lambda_3^{|a \wedge j|} \cdot (1 - \lambda_3)^{|a \wedge \neg j|} \cdot \lambda_4^{\neg a \wedge j} \cdot (1 - \lambda_4)^{\neg a \wedge \neg j} \end{aligned}$$

Assume $\lambda_1 = \lambda_3$ and $\lambda_2 = \lambda_4$. The likelihood simplifies to

$$\begin{aligned} P(E|M) = & \lambda_0^{|a|} \cdot (1 - \lambda_0)^{|\neg a|} \cdot \lambda_1^{|a \wedge m| + |a \wedge j|} \cdot (1 - \lambda_1)^{|a \wedge \neg m| + |a \wedge \neg j|} \cdot \\ & \lambda_2^{|\neg a \wedge m| + |\neg a \wedge j|} \cdot (1 - \lambda_2)^{|\neg a \wedge \neg m| + |\neg a \wedge \neg j|} \end{aligned}$$

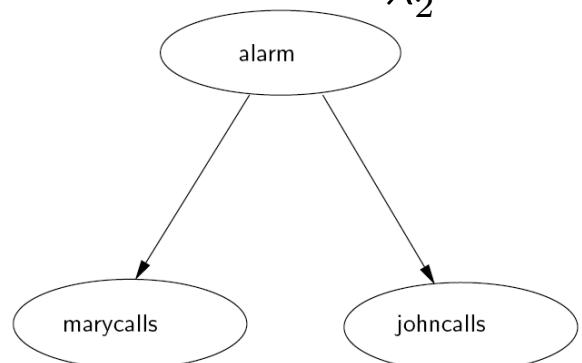


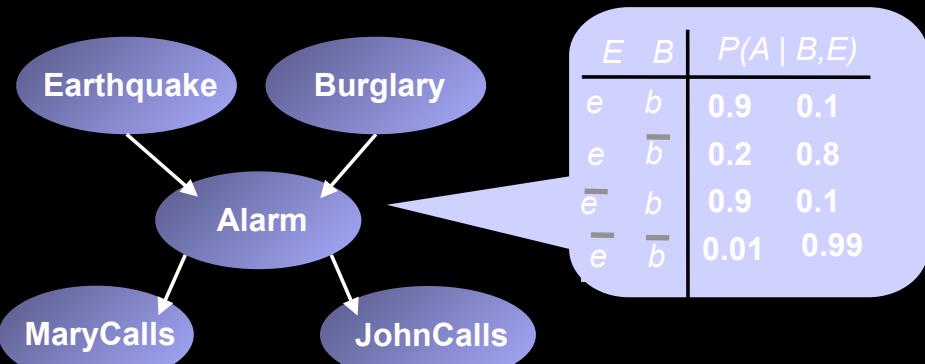
Fig. 8.7. A Bayesian network. $P(a) = \lambda_0$; $P(m|a) = \lambda_1$; $P(m|\neg a) = \lambda_2$

Structure Learning

Of Bayesian networks

- one approach greedily searches through the space of possible networks:
- iterate
 - perform operation on graph (add, delete, reverse arcs)
 - estimate parameters
 - evaluate quality

Bayesian Networks

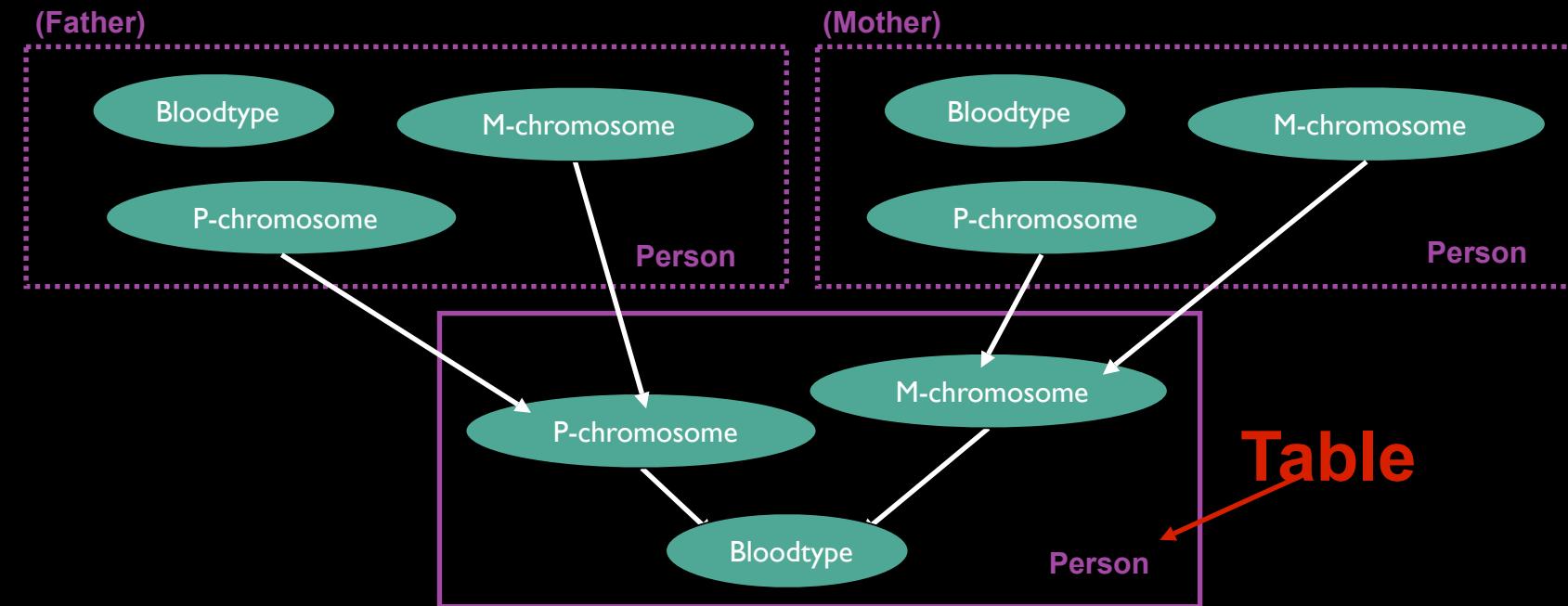


$$P(E, B, A, J, M) = P(E).P(B).P(A|E).P(A|B).P(J|A).P(M|A)$$

```
earthquake.  
burglary.  
alarm :- earthquake, burglary.  
marycalls :- alarm.  
johncalls:- alarm.
```

Probabilistic Relational Models (PRMs)

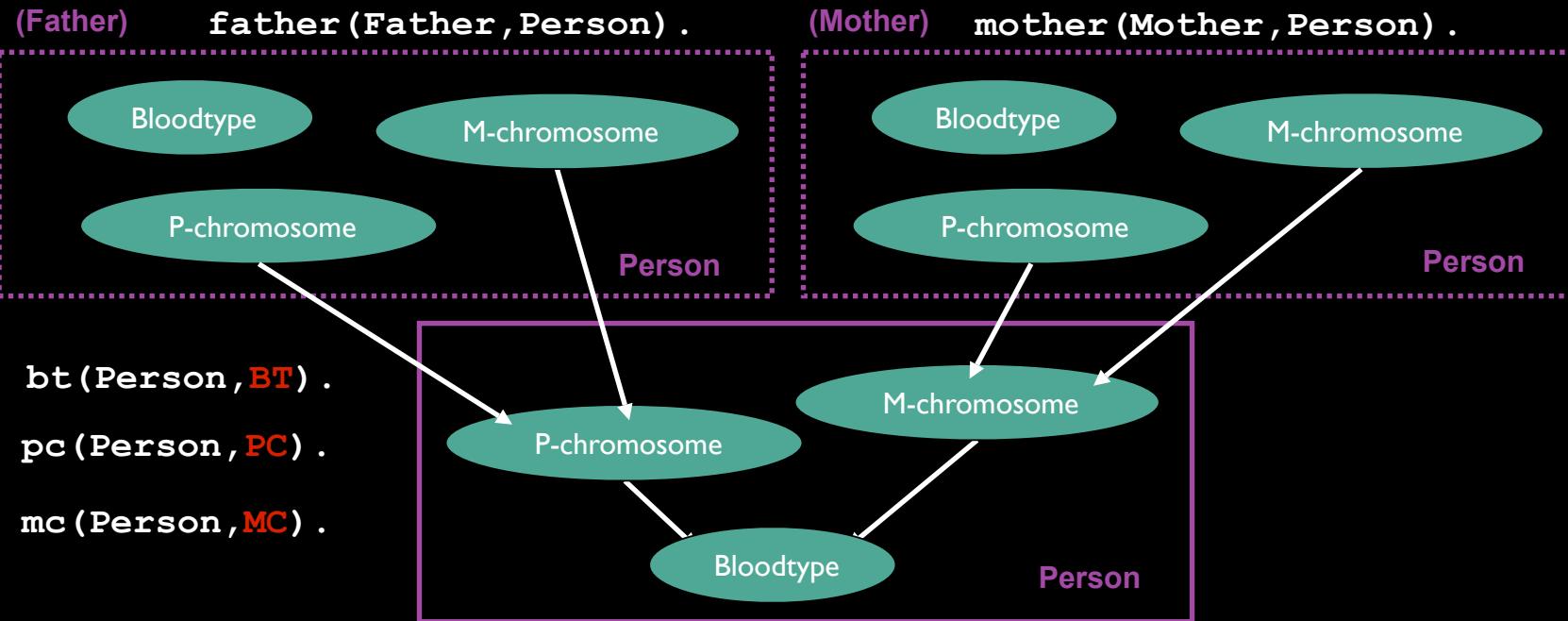
[Getoor,Koller, Pfeffer]



[Getoor,Koller, Pfeffer]

Probabilistic Relational Models (PRMs)

[Getoor, Koller, Pfeffer]



Dependencies (CPDs associated with):

bt(Person, BT) :- pc(Person, PC), mc(Person, MC).

pc(Person, PC) :- pc_father(Father, PCf), mc_father(Father, MCf).

View :

pc_father(Person, PCf) | father(Father, Person), pc(Father, PC).

...

Probabilistic Relational Models (PRMs) Bayesian Logic Programs (BLPs)

```

Extension
father(rex,fred).    mother(ann,fred).
father(brian,doro).   mother(utta, doro).
father(fred,henry).   mother(doro,henry).

```

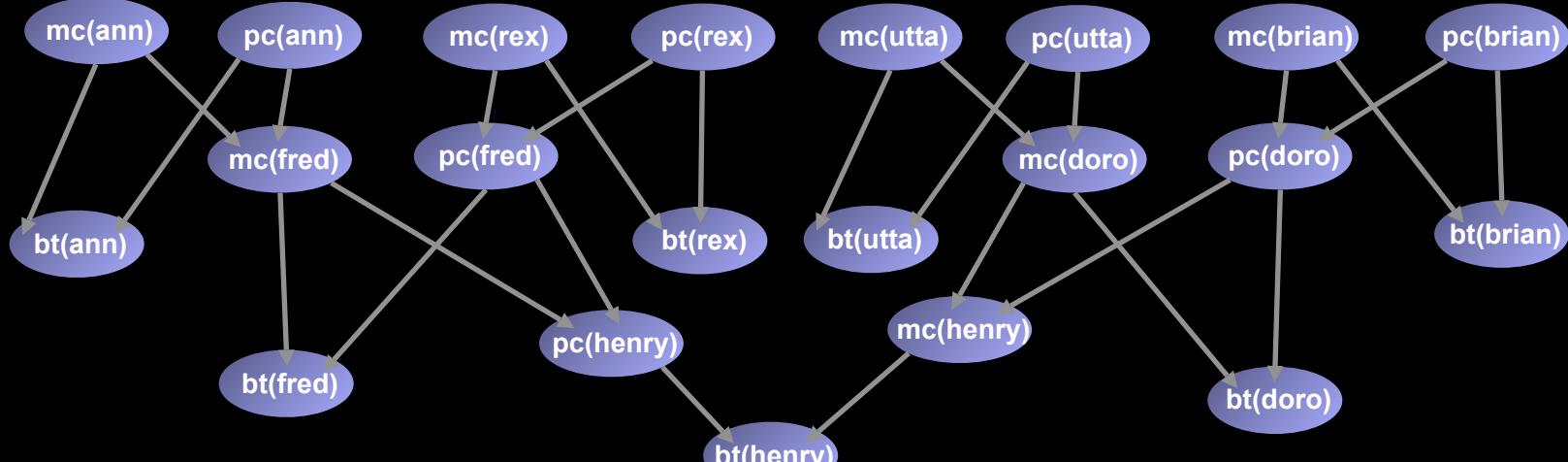
```

Intension
pc_father(Person,PCf) | father(Father,Person),pc(Father,PC).
...
mc(Person,MC) | pc_mother(Person,PCm), pc_mother(Person,MCM).
pc(Person,PC) | pc_father(Person,PCf), mc_father(Person,MCF).
bt(Person,BT) | pc(Person,PC), mc(Person,MC).

```

RV

State



Knowledge Based Model Construction

Extension + Intension => Probabilistic Model

Advantages

same intension used for multiple extensions

parameters are being shared / tied together

unification is essential

learning becomes feasible

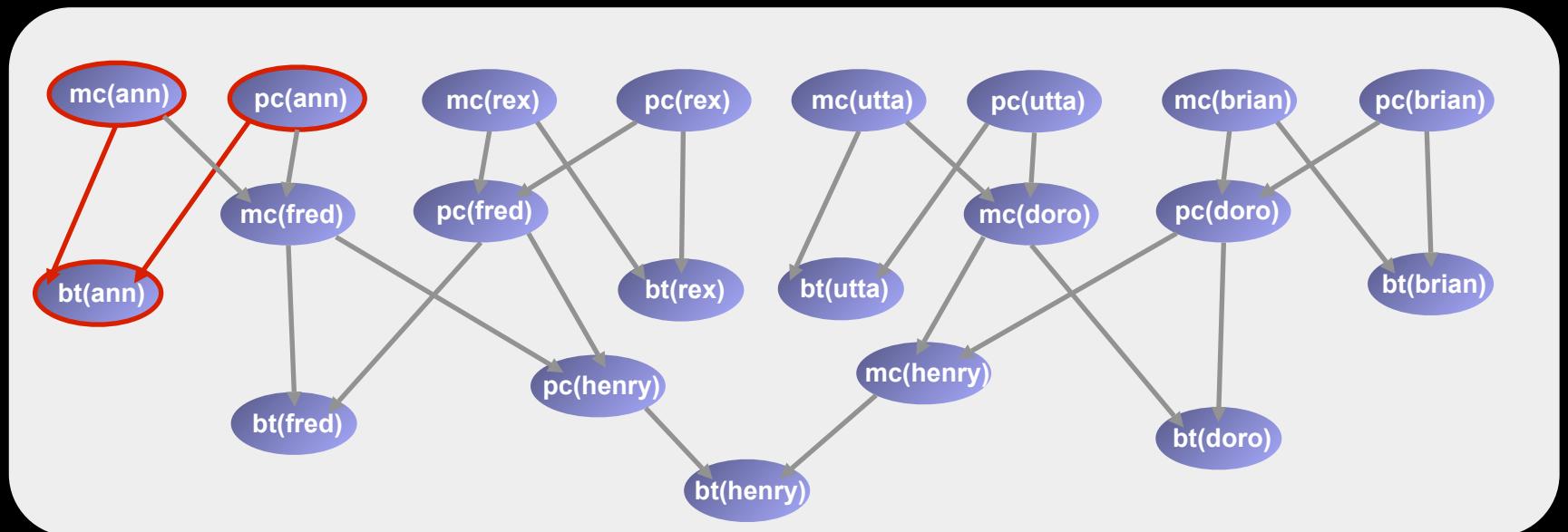
Typical use includes

prob. inference $P(Q | E)$, $P(bt(mary) | bt(john) = o -)$

max. likelihood parameter estimation & structure learning

Answering Queries

$P(bt(ann)) ?$



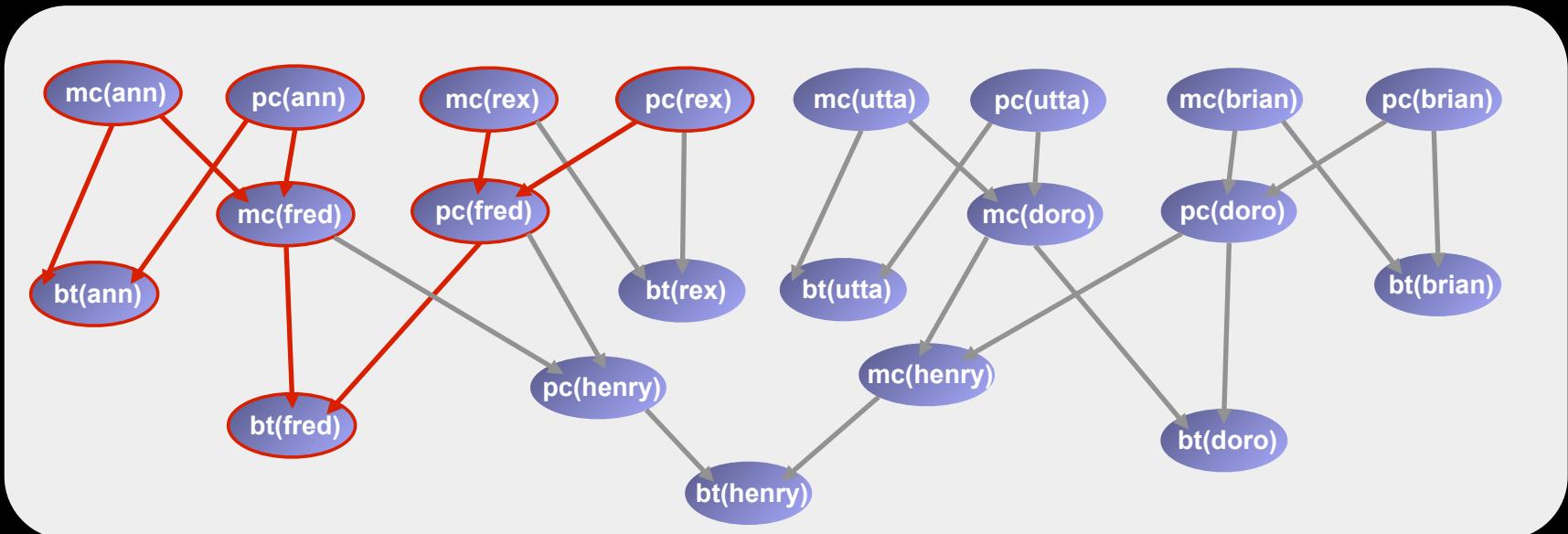
Answering Queries

Bayes' rule

$P(bt(ann), bt(fred)) ?$

$$P(bt(ann) | bt(fred)) =$$

$$\frac{P(bt(ann), bt(fred))}{P(bt(fred))}$$



Combining Partial Knowledge

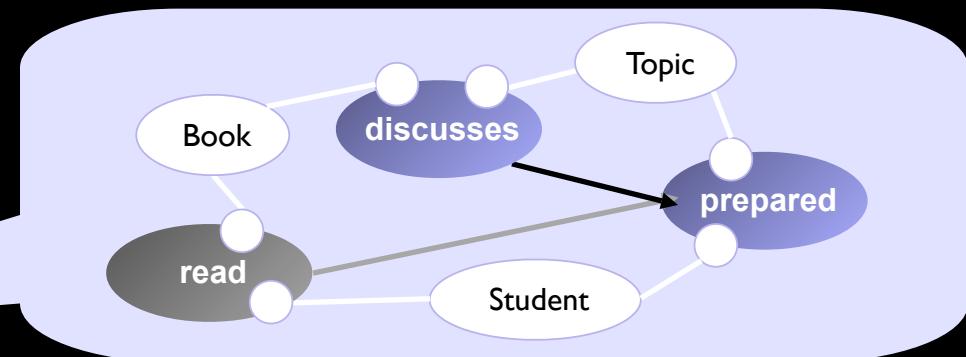
:

discusses/2

read/1

prepared/2

passes/1



prepared(Student, Topic) | read(Student, Book),
discusses(Book, Topic).

bn

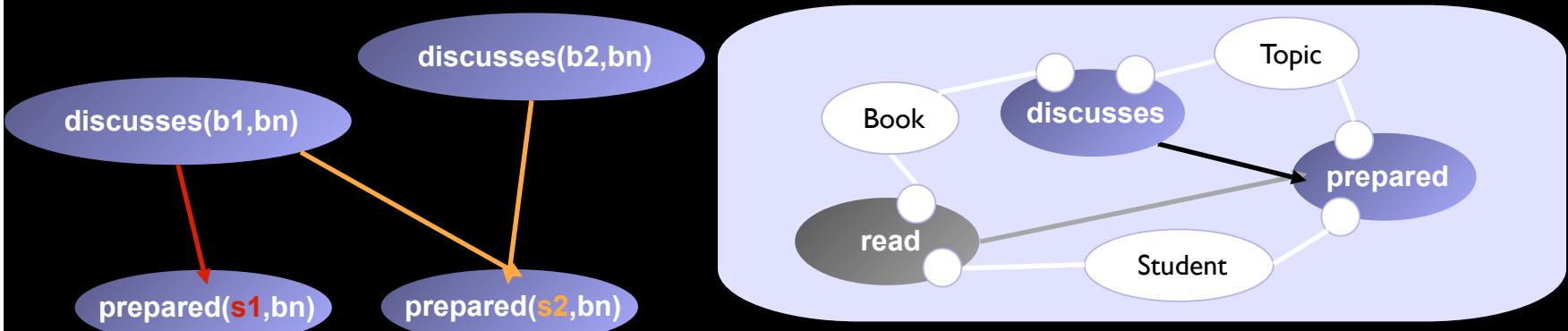
prepared

logic

passes

passes(Student) | prepared(Student, bn),
prepared(Student, logic).

Combining Partial Knowledge



```
prepared(Student, Topic) | read(Student, Book),  
discusses(Book, Topic).
```

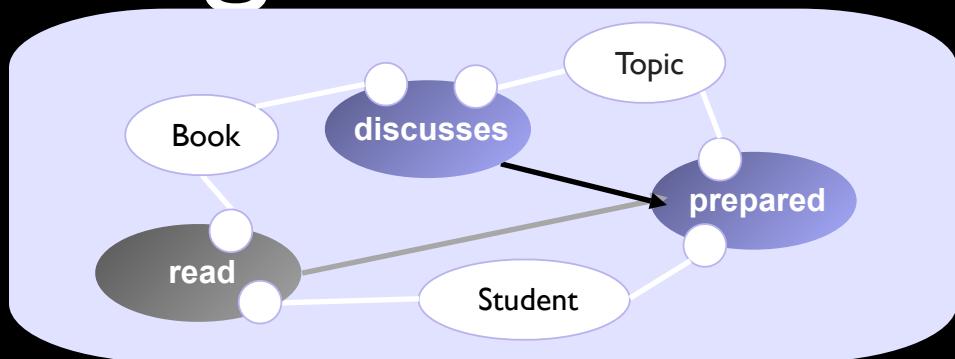
- variable # of parents for `prepared/2` due to `read/2`
 - whether a student prepared a topic depends on the books she read
- CPD only for one book-topic pair

Combining Rules

$P(A|B)$ and $P(A|C)$

CR

$P(A|B,C)$

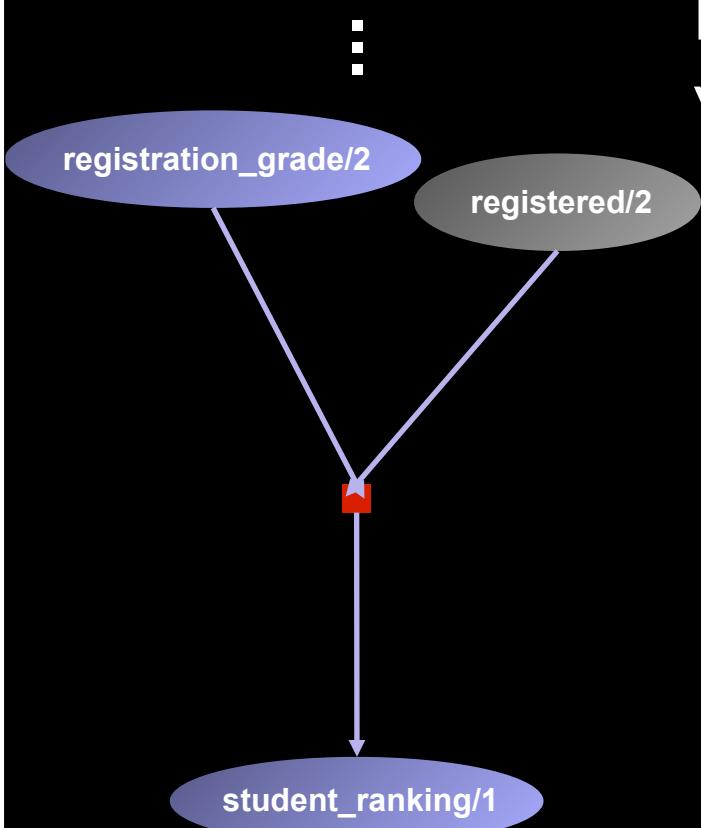


`prepared(Student,Topic) | read(Student,Book),
discusses(Book,Topic).`

Any algorithm which
has an empty output if and only if the input is empty
combines a set of CPDs into a single (combined) CPD
E.g. noisy-or, regression, ...

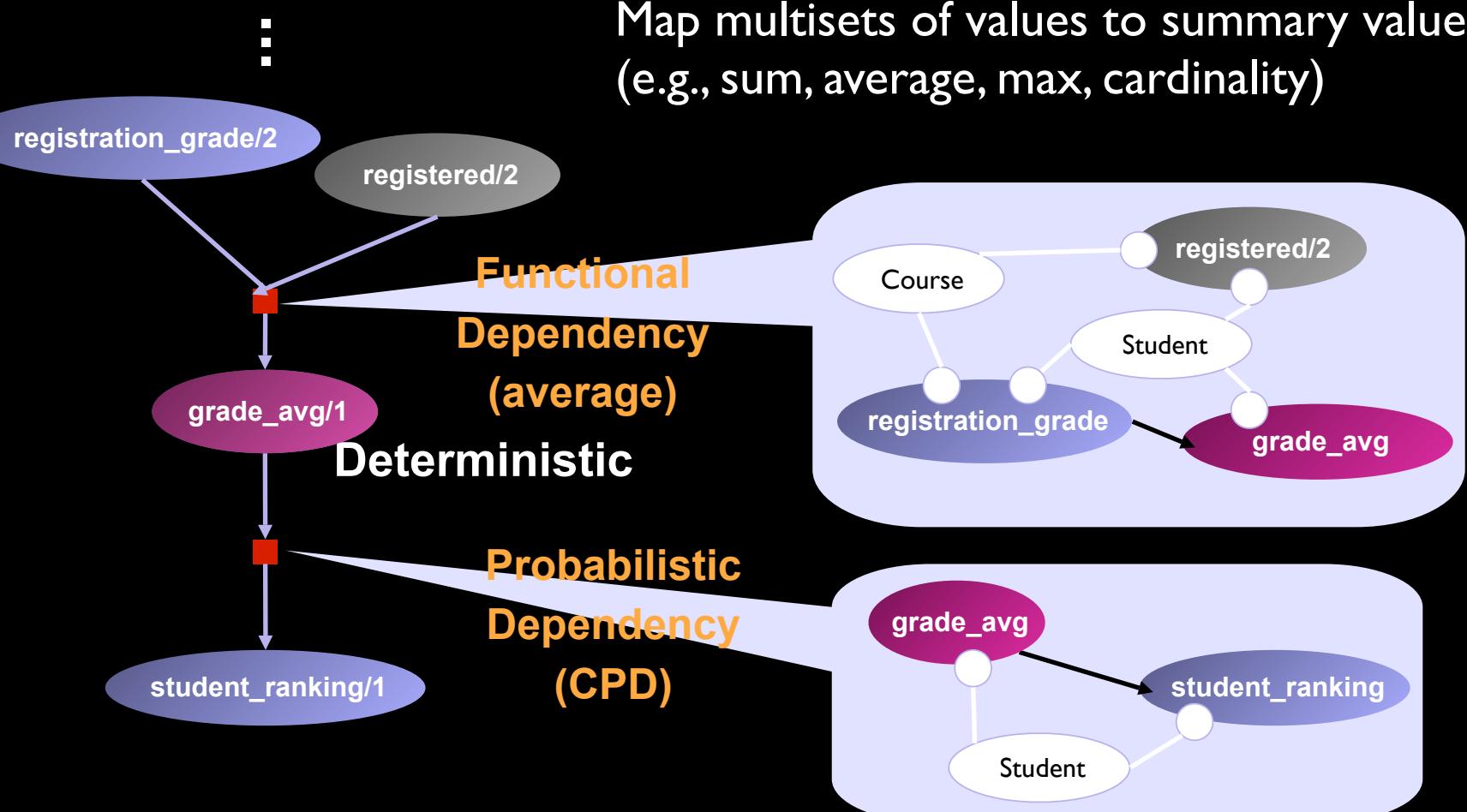
Aggregates

Map multisets of values to summary values (e.g., sum, average, max, cardinality)

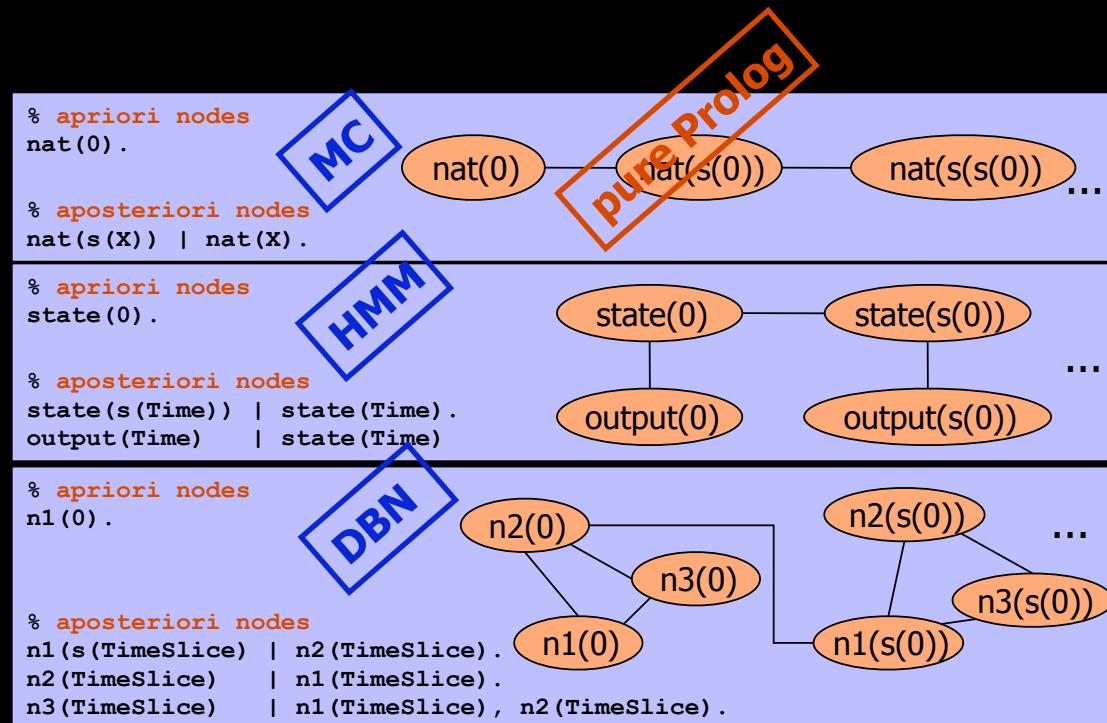


Aggregates

Map multisets of values to summary values
(e.g., sum, average, max, cardinality)



Bayesian Logic Programs



Prolog and Bayesian Nets as Special Case

Learning: the data

RVs + States = (partial) Herbrand interpretation
Probabilistic learning from interpretations

Family(1)

pc(brian)=b,
bt(ann)=a,
bt(brian)=?,
bt(dorothy)=a

Background

m(ann,dorothy),
f(brian,dorothy),
m(cecily,fred),
f(henry,fred),
f(fred,bob),
m(kim,bob),
...

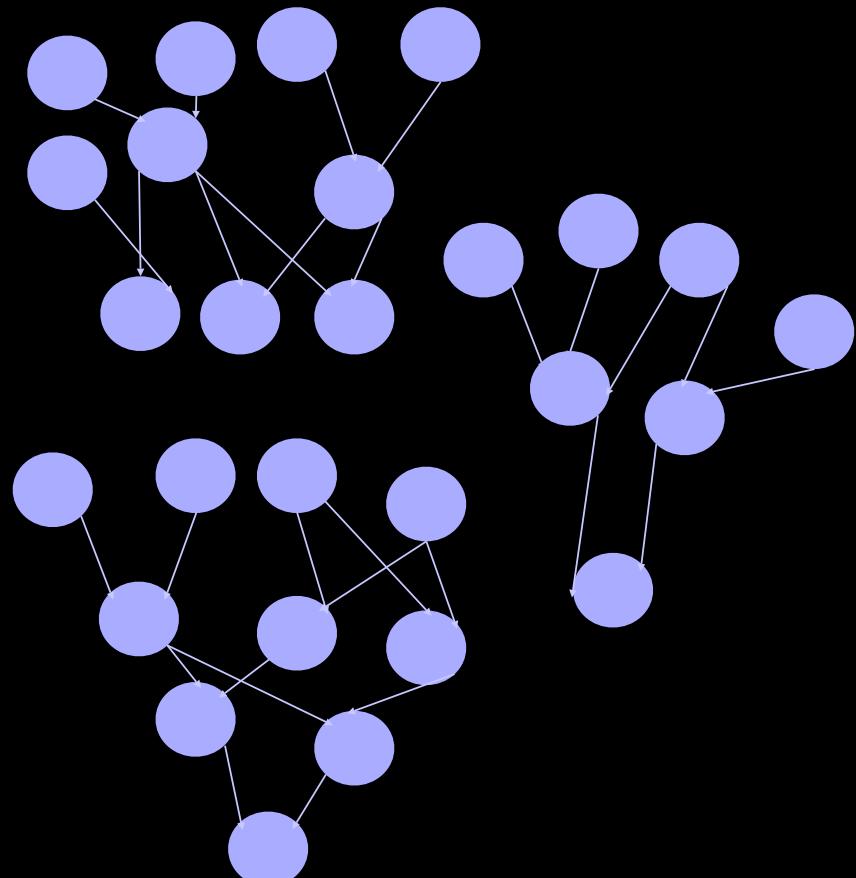
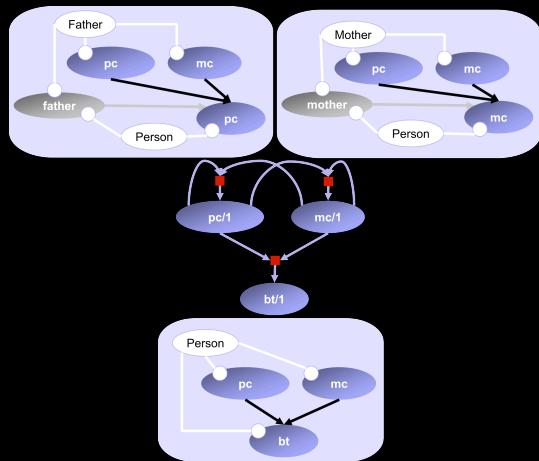
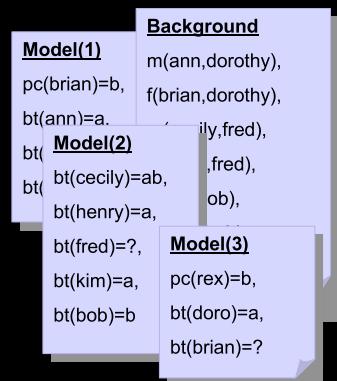
Family(2)

bt(cecily)=ab,
pc(henry)=a,
mc(fred)=?,
bt(kim)=a,
pc(bob)=b

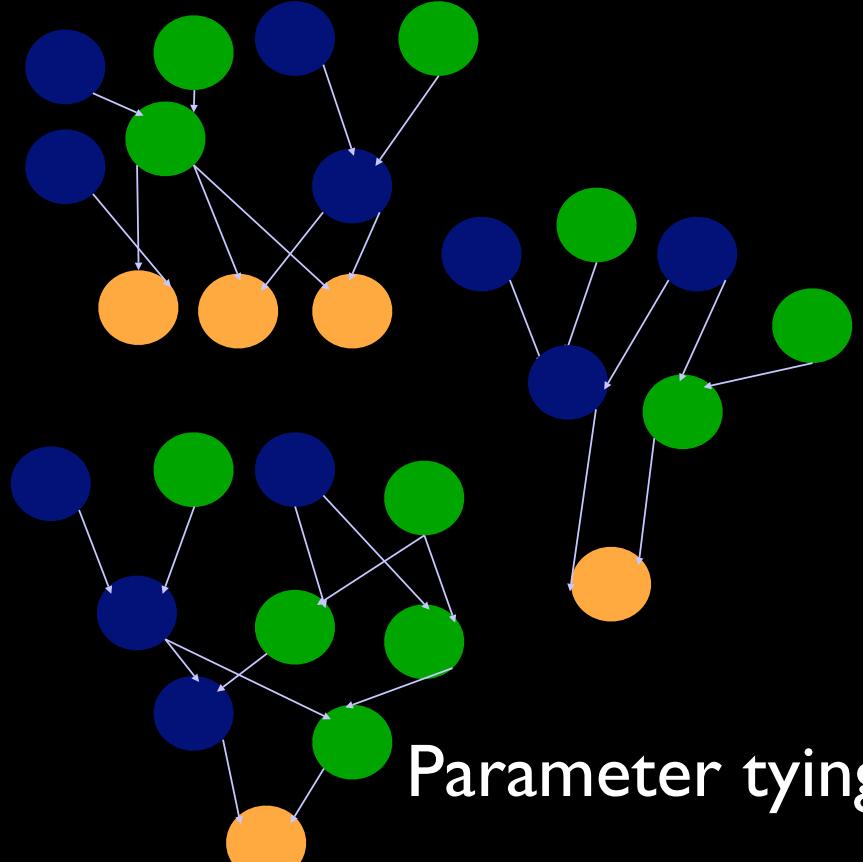
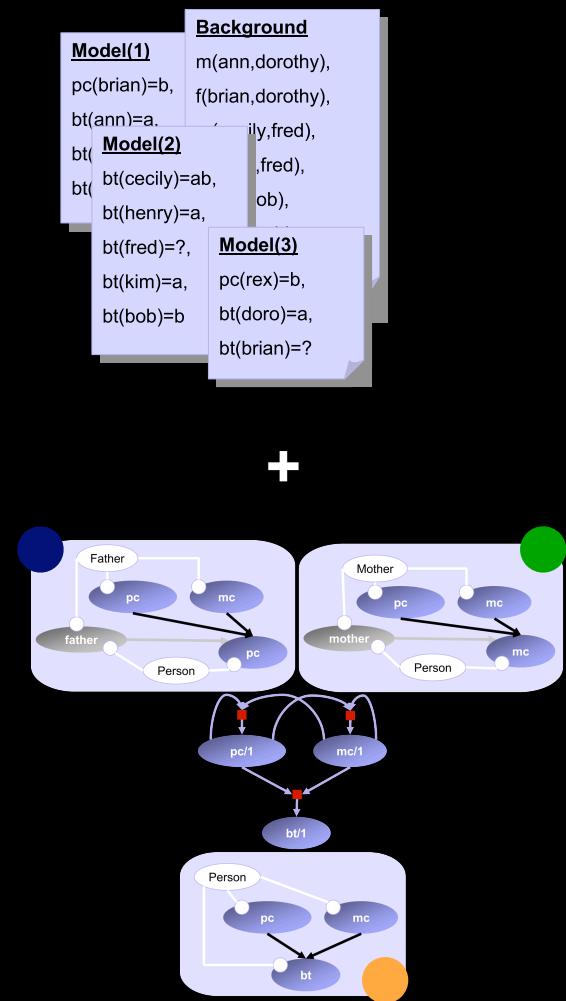
Family(3)

pc(rex)=b,
bt(doro)=a,
bt(brian)=?

Parameter Estimation



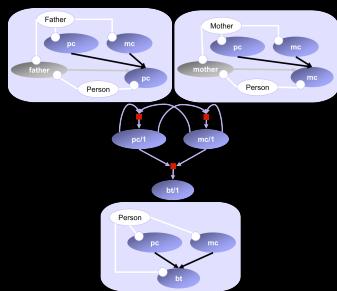
Parameter Estimation



Expectation Maximization

EM-algorithm:
iterate until convergence

Logic Program L



Initial Parameters q_0

**Current Model
(M, q_k)**

$$\sum_{\text{Ground Instance GI}} \sum_{\text{DataCase DC}} P(\text{head(GI)}, \text{body(GI)} | \text{DC})$$

$$\sum_{\text{Ground Instance GI}} \sum_{\text{DataCase DC}} P(\text{body(GI)} | \text{DC})$$

Expected counts of a clause

$$\sum_{\text{Ground Instance GI}} \sum_{\text{DataCase DC}} P(\text{head(GI)}, \text{body(GI)} | \text{DC})$$

Maximization
Update parameters
(ML, MAP)

Background	
Model(1)	m(ann,dorothy),
pc(brian)=b,	f(brian,dorothy),
bt(ann)=a.	lily,fred),
Model(2)	,fred),
bt(cecily)=ab,	ob),
bt(henry)=a,	
bt(fred)=?,	
bt(kim)=a,	
bt(bob)=b	
Model(3)	
pc(rex)=b,	
bt(doro)=a,	
bt(brian)=?	

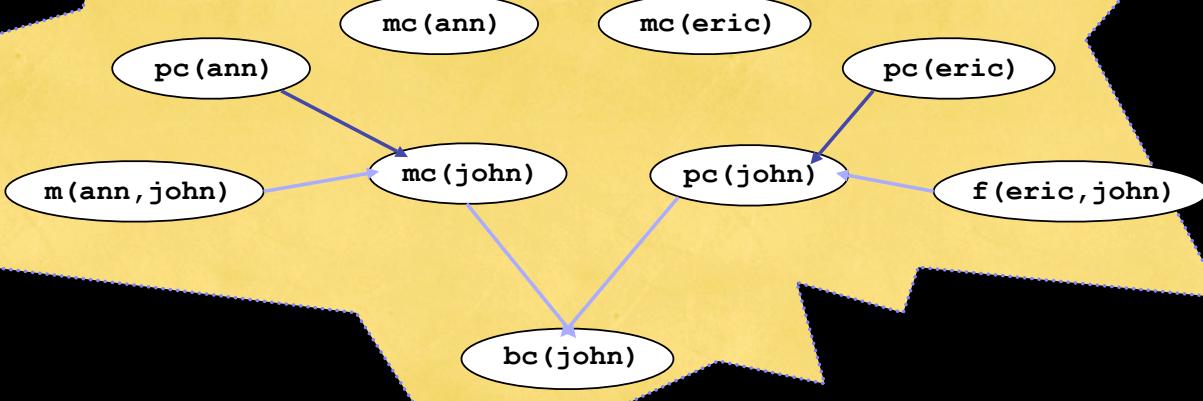
Structure Learning

Original program

```
mc(X) | m(M,X), mc(M), pc(M).  
pc(X) | f(F,X), mc(F), pc(F).  
bt(X) | mc(X), pc(X).
```

Initial hypothesis

```
mc(X) | m(M,X).  
pc(X) | f(F,X).  
bt(X) | mc(X).
```



Refinement

```
mc(X) | m(M,X).  
pc(X) | f(F,X).  
bt(X) | mc(X), pc(X).
```

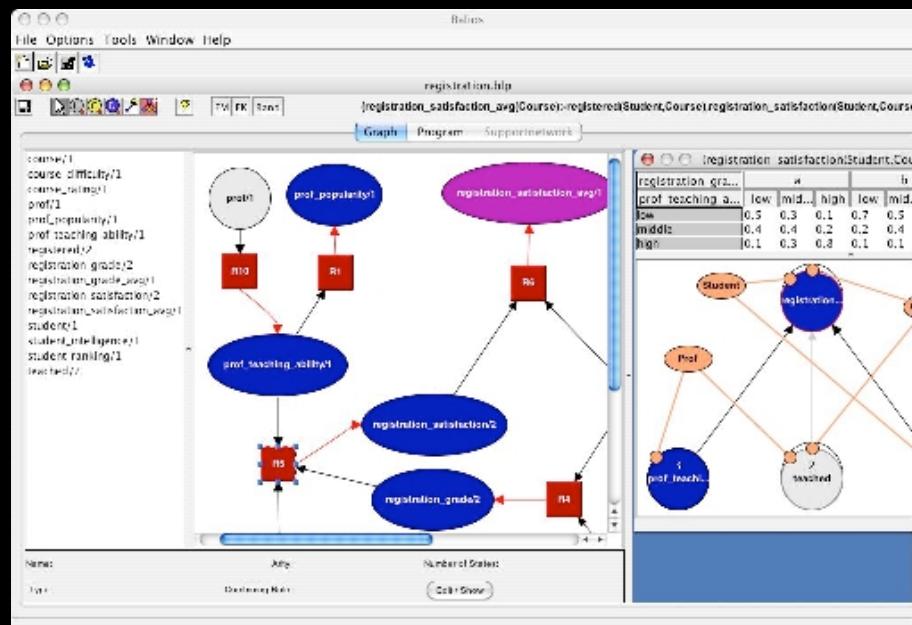
Refinement

```
mc(X) | m(M,X), pc(X).  
pc(X) | f(F,X).  
bt(X) | mc(X), pc(X).
```

Bayesian Logic Programs (BLPs)

- Unique probability distribution over Herbrand interpretations
 - Finite branching factor, finite proofs, no self-dependency
- Highlight
 - Separation of qualitative and quantitative parts
 - Functors
- Graphical Representation
- Discrete and continuous RV
- BNs, DBNs, HMMs, SCFGs, Prolog ...
- Turing-complete programming language
- Learning

Balios Tool



Undirected Probabilistic Relational

So far, **directed** graphical models only
Impose **acyclicity constraint**

Undirected graphical models do not impose
the acyclicity constraint

Undirected Probabilistic Relational

Two approaches

- Relational Markov Networks (**RMNs**) (Taskar et al.)
- Markov Logic Networks (**MLNs**) (Anderson et al.)

Idea

- Semantics defined in terms of Markov Networks (undirected graphical models)
- More natural for certain applications

RMNs ~ undirected PRM

MLNs ~ undirected BLP

A Markov Network/Random Field

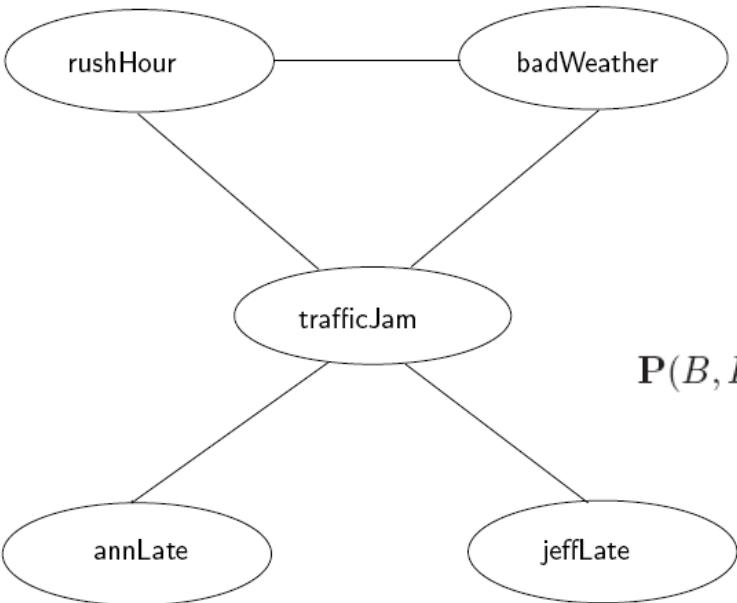


Fig. 8.3. A Markov network

$$P(X_1, \dots, X_n) = \frac{1}{Z} \prod_{c: \text{clique}} f_c(X_{c1}, \dots, X_{ck_c})$$

$$Z = \sum_{x_1, \dots, x_n} \prod_{c: \text{clique}} f_c(X_{c1} = x_{c1}, \dots, X_{ck_c} = x_{ck_c})$$

$$P(B, R, T, J, A) = \frac{1}{Z} f_{T,R,B}(T, R, B) f_{J,T}(J, T) f_{A,T}(A, T)$$

$$f_c(X_{c1}, \dots, X_{ck_c}) = e^{w_i F_i(X_{c1}, \dots, X_{ck_c})}$$

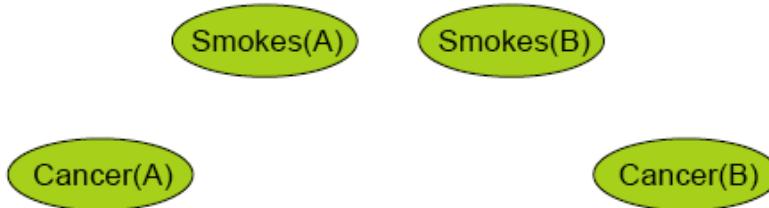
Assumption 2 Each node X_i in the graph is conditionally independent of any subset \mathbf{A} of nodes that are not neighbours of X_i given a joint state of $\text{Ne}(X_i)$, that is, $P(X_i | \mathbf{A}, \text{Ne}(X_i)) = P(X_i | \text{Ne}(X_i))$. The expression $\text{Ne}(X_i)$ refers to the set of all neighbours of X_i .

Markov Logic

1.5 $\forall x \text{ Smokes}(x) \Rightarrow \text{Cancer}(x)$

1.1 $\forall x, y \text{ Friends}(x, y) \Rightarrow (\text{Smokes}(x) \Leftrightarrow \text{Smokes}(y))$

Suppose we have two constants: Anna (A) and Bob (B)



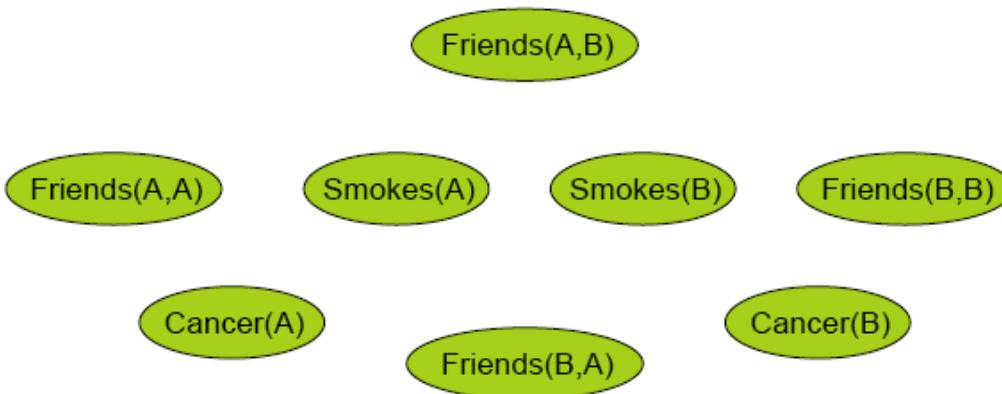
(slide and work by Domingos et al.)

Markov Logic

1.5 $\forall x \text{ Smokes}(x) \Rightarrow \text{Cancer}(x)$

1.1 $\forall x, y \text{ Friends}(x, y) \Rightarrow (\text{Smokes}(x) \Leftrightarrow \text{Smokes}(y))$

Suppose we have two constants: Anna (A) and Bob (B)



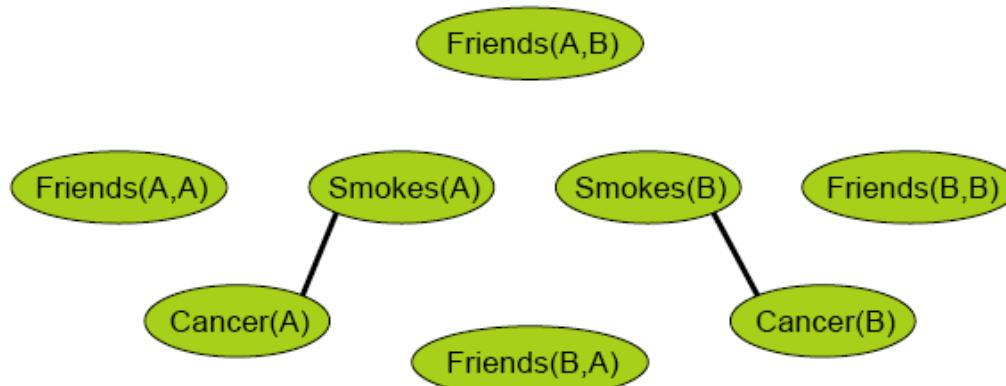
(slide and work by Domingos et al.)

Markov Logic

1.5 $\forall x \text{ Smokes}(x) \Rightarrow \text{Cancer}(x)$

1.1 $\forall x, y \text{ Friends}(x, y) \Rightarrow (\text{Smokes}(x) \Leftrightarrow \text{Smokes}(y))$

Suppose we have two constants: Anna (A) and Bob (B)



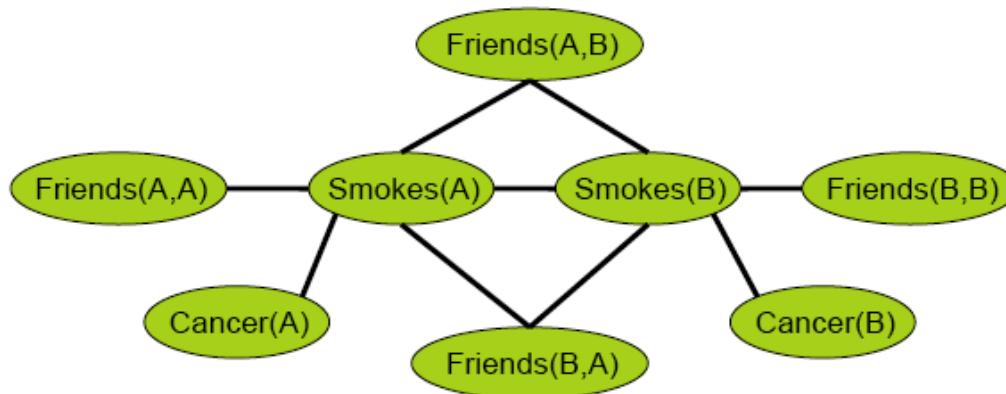
(slide and work by Domingos et al.)

Markov Logic

1.5 $\forall x \text{ Smokes}(x) \Rightarrow \text{Cancer}(x)$

1.1 $\forall x, y \text{ Friends}(x, y) \Rightarrow (\text{Smokes}(x) \Leftrightarrow \text{Smokes}(y))$

Suppose we have two constants: Anna (A) and Bob (B)



(slide and work by Domingos et al.)

Some interesting problems

- Inference is generally intractable
- Most Probable State Estimation
- Lazy Max Walk Sat algorithm

MLNs

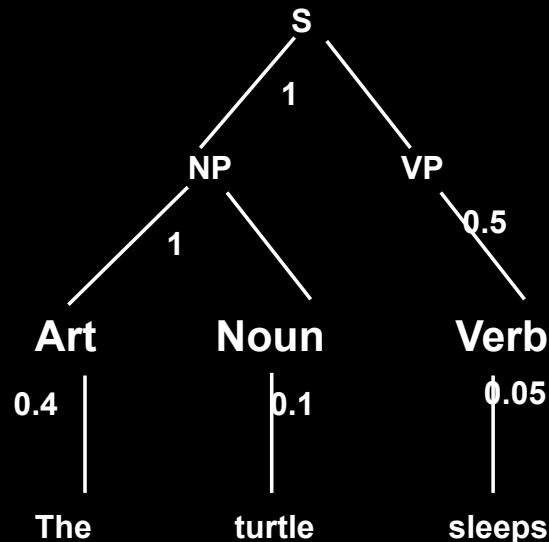
- Probably the best known and most elaborated approach to SRL
- Many different inference & learning algorithms known; discriminative, generative, structure learning, ...
- See website Pedro Domingos and Alchemy system

C. Probabilistic Logic Learning

C.2. Learning from proofs & entailment

Probabilistic Context Free Grammars

1.0 : S -> NP, VP
1.0 : NP -> Art, Noun
0.6 : Art -> a
0.4 : Art -> the
0.1 : Noun -> turtle
0.1 : Noun -> turtles
...
0.5 : VP -> Verb
0.5 : VP -> Verb, NP
0.05 : Verb -> sleep
0.05 : Verb -> sleeps
....



$$P(\text{parse tree}) = 1 \times 1 \times 0.5 \times 1 \times 0.4 \times 0.05$$

PCFGs

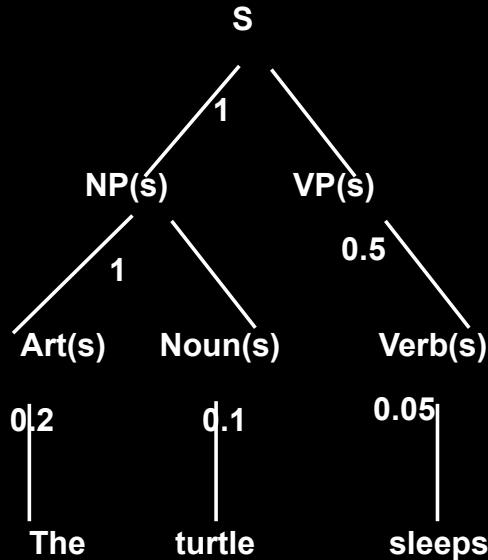
$P(\text{parse tree}) = \prod_i p_i c_i$
where p_i is the probability of rule i
and c_i the number of times
it is used in the parse tree

$$P(\text{sentence}) = \sum_{p:\text{parsetree}} P(p)$$

Observe that derivations always succeed, that is
 $S \rightarrow T, Q$ and $T \rightarrow R, U$
always yields
 $S \rightarrow R, U, Q$

Probabilistic DCG

1.0 $S \rightarrow NP(Num), VP(Num)$
1.0 $NP(Num) \rightarrow Art(Num), Noun(Num)$
0.6 $Art(sing) \rightarrow a$
0.2 $Art(sing) \rightarrow the$
0.2 $Art(plur) \rightarrow the$
0.1 $Noun(sing) \rightarrow turtle$
0.1 $Noun(plur) \rightarrow turtles$
...
0.5 $VP(Num) \rightarrow Verb(Num)$
0.5 $VP(Num) \rightarrow Verb(Num), NP(Num)$
0.05 $Verb(sing) \rightarrow sleep$
0.05 $Verb(plur) \rightarrow sleeps$
...



$$P(\text{derivation tree}) = 1 \times 1 \times .5 \times 1 \times .2 \times .05$$

1

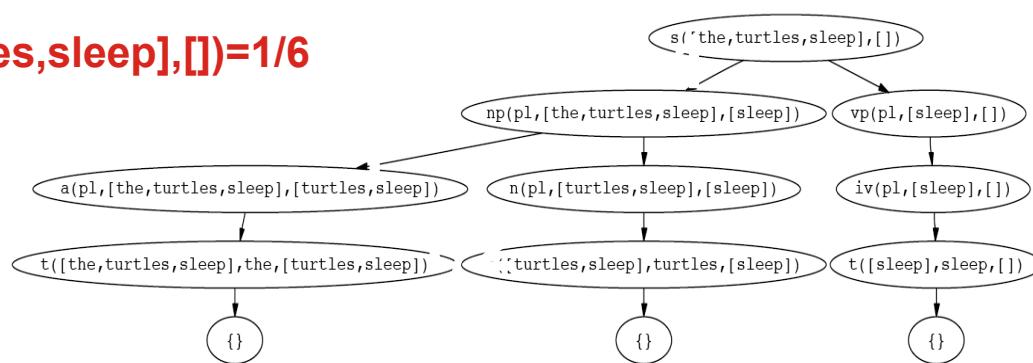
1/3

1/2

As a Stochastic Logic Program

```
sentence(A, B) :- noun_phrase(C, A, D), verb_phrase(C, D, B).  
noun_phrase(A, B, C) :- article(A, B, D), noun(A, D, C).  
verb_phrase(A, B, C) :- intransitive_verb(A, B, C).  
article(singular, A, B) :- terminal(A, a, B).  
article(singular, A, B) :- terminal(A, the, B).  
article(plural, A, B) :- terminal(A, the, B).  
noun(singular, A, B) :- terminal(A, turtle, B).  
noun(plural, A, B) :- terminal(A, turtles, B).  
intransitive_verb(singular, A, B) :- terminal(A, sleeps, B).  
intransitive_verb(plural, A, B) :- terminal(A, sleep, B).  
terminal([A|B], A, B).
```

$$P(s([\text{the}, \text{turtles}, \text{sleep}], [])) = 1/6$$



Probabilistic DCG

1.0 $S \rightarrow NP(Num), VP(Num)$

1.0 $NP(Num) \rightarrow Art(Num), Noun(Num)$

0.6 $Art(sing) \rightarrow a$

0.2 $Art(sing) \rightarrow the$

0.2 $Art(plur) \rightarrow the$

0.1 $Noun(sing) \rightarrow turtle$

0.1 $Noun(plur) \rightarrow turtles$

...

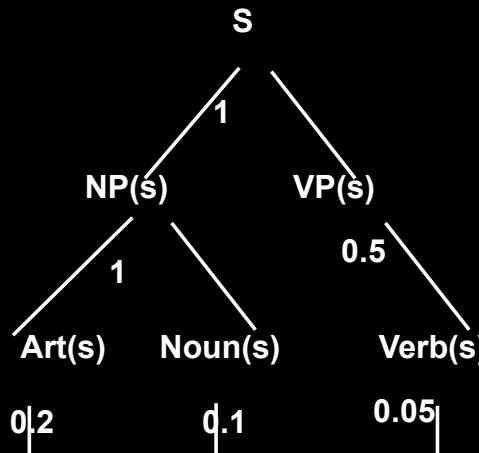
0.5 $VP \rightarrow VP$

0.5 $VP \rightarrow VP$

0.05 $Verb(sing) \rightarrow sleep$

0.05 $Verb(plur) \rightarrow sleeps$

...



What about “A turtles sleeps” ?

$$P(\text{derivation tree}) = 1 \times 1 \times .5 \times 1 \times .2 \times .05$$

SLPs

$$P_d(\text{derivation}) = \prod_i p_i^{c_i}$$

where p_i is the probability of rule i
and c_i the number of times
it is used in the parse tree

Observe that some derivations now fail due to unification, that
 $np(Num) \rightarrow art(Num)$, $noun(Num)$ and $art(sing) \rightarrow a$
 $noun(plural) \rightarrow turtles$

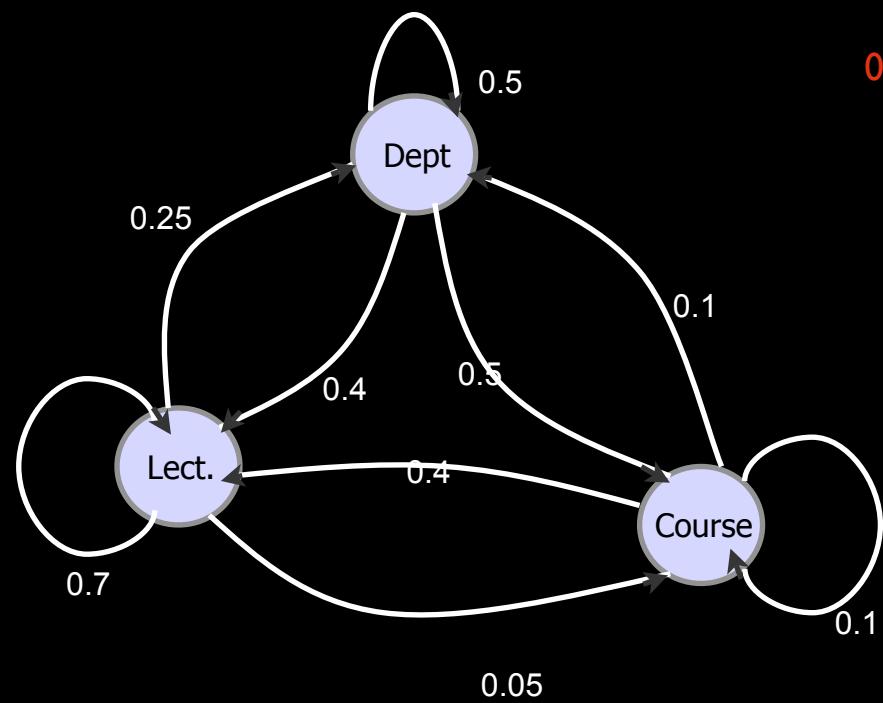
Normalization necessary

$$P_s(\text{proof}) = \frac{P_d(\text{proof})}{\sum_i P_d(\text{proof}_i)}$$

Example Application

- Consider traversing a university website
 - Pages are characterized by predicate
 - $\text{department}(\text{cs}, \text{nebel})$ denotes the page of cs following the link to nebel
 - Rules applied would be of the form
 - $\text{department}(\text{cs}, \text{nebel}) :- \text{prof}(\text{nebel}), \text{in}(\text{cs}), \text{co}(\text{ai}), \text{lecturer}(\text{nebel}, \text{ai}).$
 - $\text{pagetype1}(\text{t1}, \text{t2}) :- \text{type1}(\text{t1}), \text{type2}(\text{t2}), \text{type3}(\text{t3}), \text{pagetype2}(\text{t2}, \text{t3})$
 - SLP models probabilities over traces / proofs / web logs
 - $\text{department}(\text{cs}, \text{nebel}), \text{lecturer}(\text{nebel}, \text{ai007}), \text{course}(\text{ai007}, \text{burgard}), \dots$
- This is actually a Logical Markov Model

Logical Markov Model



0.4 department(D,L) :-
prof(L),
in(D),
co(C),
lecturer(L,C).

0.1 prof(nebel).
0.05 prof(burgard).
...

Learning SLPs from entailment

Entailment

Given

- a set of ground facts
 - $s(\text{[the, boy, rides]})$
- structure of the logic program

sentences

Find

- the parameters of the logic program that maximize the log likelihood

CFG

PCFG

Parameter estimation

- Failure adjusted maximization algorithm (Cussens)
 - Only facts observed
 - Proofs and Failures not observed !
- Inside-outside algorithm for Prism (Sato)
- EM based techniques

no failures

Learning SLPs from entailment

Entailment

Given

- a set of ground facts **sentences**
- ~~the structure of the logic program~~ **CFG**

Find

- the logic program and its parameters that maximizes the log likelihood

PCFG

Structure learning

- requires one to solve the underlying relational learning problem
- as well as the parameter estimation problem
- Very hard because the information contained in the examples is very limited.
- What about using richer examples ? proofs ?

Learning SLPs from entailment

Entailment

Given

- a set of proof trees **sentences**
- ~~the structure of the logic program~~ **CFG**

Find

- the logic program and its parameters that maximizes the log likelihood

PCFG

Structure learning

- analogy with learning from tree banks, cf. De Raedt et al. AAAI 05.

no failures

Learning from parse tree

$S \rightarrow NP(s), VP(s)$

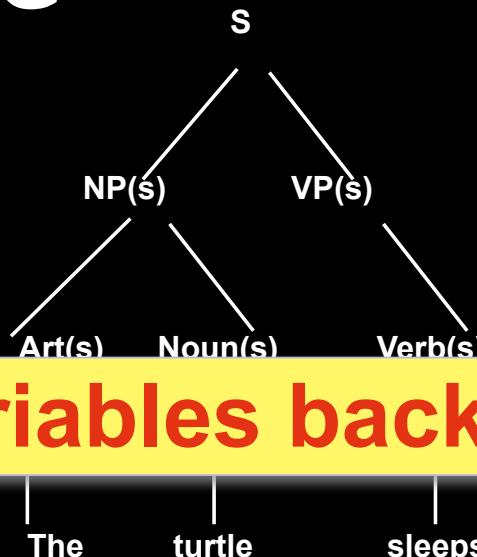
$NP(s) \rightarrow Art(s).$

$VP(s) \rightarrow Verb(s)$

$Art(s) \rightarrow the$

$Noun(s) \rightarrow turtle$

$Verb(s) \rightarrow sleeps$



How to get the variables back ?

Generalizing Rules in SLPs

Generalization in ILP

- Take two clauses for same predicate and replace them by the lgg under theta-subsumption (Plotkin), e.g

```
department(cs,nebel) :-
```

```
    prof(nebel), in(cs), course(ai), lect(nebel,ai).
```

```
department(cs,burgard) :-
```

```
    prof(burgard), in(cs), course(ai), lect(burgard,ai)
```

- Induce

```
department(cs,P) :-
```

```
    prof(P), in(cs), course(ai), lect(P,ai)
```

Strong logical constraints

Replacing the rules r_1 and r_2 by the l_{gg} should preserve the proofs !

So, two rules r_1 and r_2 should only be generalized when

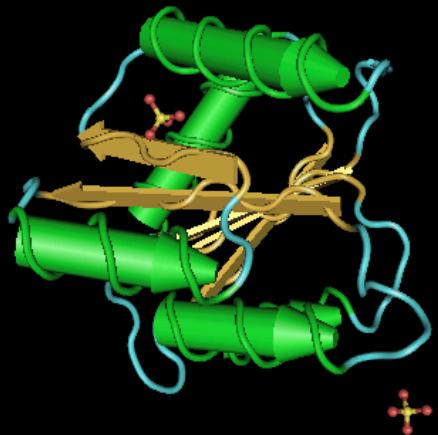
- There is a one to one mapping (with corresponding substitutions) between literals in r_1, r_2 and $l_{gg}(r_1, r_2)$

Exclude

- $\text{father}(j,a) :- m(j), f(a), \text{parent}(j,a)$
- $\text{father}(j,t) :- m(j), m(t), \text{parent}(j,t)$

[Kersting et al.; Kersting, Gaertner]

Protein Fold Recognition



```
[helix(h(right,3to10),5),  
 helix(h(right,alpha),13),  
 strand(null,7),  
 strand(minus,7),  
 strand(minus,5),  
 helix(h(right,3to10),5),...]
```

If one of the two similar proteins has a known structure, can build a rough model of the protein of unknown structure.

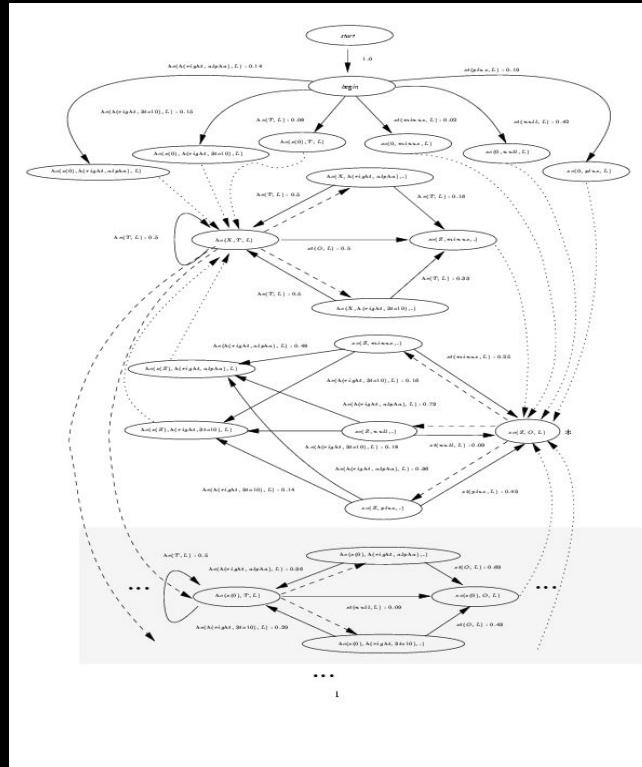
LoHMM

[Kersting et al., JAIR 2006]

~120 parameters

vs.

over 62000
parameters



Secondary structure of domains of proteins

(from PDB and SCOP)

fold1:TIM beta/alpha barrel fold, fold2: NAD(P)-binding Rossmann-fold fold23: Ribosomal protein L4, fold37: glucosamine 6-phosphate deaminase/isomerase old fold55: leucine aminopeptidas fold. 3187 logical sequences (> 30000 ground atoms)

Results

[Kersting et al.; Kersting, Gaertner; Landwehr et al.]

Accuracy: 74% vs. 82.7% (1622 vs. 1809 / 2187)

Majority vote: 43%

	fold1	fold2	fold23	fold37	fold55
precision	0.86 / 0.89	0.69 / 0.86	0.56 / 0.82	0.72 / 0.70	0.66 / 0.74
recall	0.78 / 0.87	0.67 / 0.81	0.71 / 0.85	0.66 / 0.72	0.96 / 0.86

Web Log Data

[Anderson et al., KDD 03]

Relational Markov Models

Log data of web sides

KDDCup 200

RMM

- Home()
- Boutique()
- Departments()
- Legcare_vendor()
- Lifestyles()
- Vendor()
- AssortmentDefault()
- Assortment(Assortment)
- ProductDetailLegcareDefault()
- ProductDetailLegcare(Product)
- ProductDetailLegwearDefault()
- ProductDetailLegwearProduct(Product)
- ProductDetailLegwearAssortment(Assortment)
- ProductDetailLegwearProdCollect(Product, Collection)
- ProductDetailLegwearProdAssort(Product, Assortment)
- ProductDetailLegwear(Product, Collection, Assortment)

Probabilities on Proofs

Two views

- stochastic logic programs define a prob. distribution over atoms for a given predicate.
 - The sum of the probabilities = 1.
 - Sampling. Like in probabilistic grammars.
- distribution semantics define a prob. distribution over possible worlds/interpretations. Degree of belief.

Probabilities on Facts

- Sato's PRISM and Poole's ICL and PHA use disjoint statements (or multi-valued switches)

disjoint($p_1:a_1, \dots, p_k:a_k$)

disjoint(0.5:head(C); 0.5:tail(C))

head(C) and tail(C) are mutually exclusive

false :- head(C), tail(C)

- for such statements: one of the alternatives is true, instead of true or false for a single fact.

PRISM (Sato) / ICL (Poole)

A logic program in which probability labels are attached to facts;

Clauses carry no probability label (or equiv. $P = 1$)

Probability distributions can be defined in a related/similar fashion on proofs / on explanations / on atoms - though some differences

Abductive reasoning with probabilities !

PRISM and ICL more expressive than SLPs

- SLPs can be easily transformed into ICL or PRISM

ICL / PRISM example

btype('A') :- (gtype(a,a); gtype(a,o); gtype(o,a)).

btype('B') :- (gtype(b,b); gtype(b,o); gtype(o,b)).

btype('O') :- gtype(o,o).

btype('AB') :- (gtype(a,b); gtype(b, a)).

gtype(X,Y) :- gene(father,X), gene(mother,Y).

gene(P,G) :- msw(gene,P,G). *probabilistic switch*

disjoint(p1: msw(gene,P,o),

p2: msw(gene,P,a),

p3: msw(gene,P,b))

outcomes + probabilities switch

Infer e.g. P(btype('A'))

Abduction

`btype('A') iff gtype(a,a) or gtype(a,o) or gtype(o,a)`

`gtype(a,a) iff msw(gene,father;a) and msw(gene,mother;a)`

`gtype(a,o) iff msw(gene,father;a) and msw(gene,mother;o)`

`gtype(o,a) iff msw(gene,father;o) and msw(gene,mother,a)`

So, `btype('A') = p1Xp1 + p1Xp2+p2Xp1`

Can be used for sampling as well.

ICL/PRISM Assumption

- The Exclusive Explanation Assumption :
 - avoid the NP-complete problem by assuming that all explanations (that is proofs, or sets $S \subseteq F$, are mutually exclusive)
 - under this condition (which can be enforced in different ways)
 - $P(A \vee B \vee C) = P(A) + P(B) + P(C)$

Learning from entailment

Given

- examples of the form $btype(x)$
- structure of program

Estimate parameters of switches

Nice results -- HMMs, PCFGs, some type of BNs as special case

Sato's Prism one of the most advanced PLL systems

Has been used by Cussens to learn SLPs

Learning from Entailment and Proofs

SLPs extend PCFGs as a representation

Proof-trees for SLPs correspond to parse-trees in PCFGs

Upgrading the learning from tree-banks setting for use in SLPs

Allows one also to elegantly model and study RMMs and LOHMMs

- Sequential relational / logical traces.

Defining probabilities on facts provides interesting perspective

- PRISM and ICL
- Probabilistic Abduction
- Also learning ...

C. Probabilistic Logic Learning

C.3. *Discriminative models & propositionalization*

Probabilistic *discriminative* ILP Problem

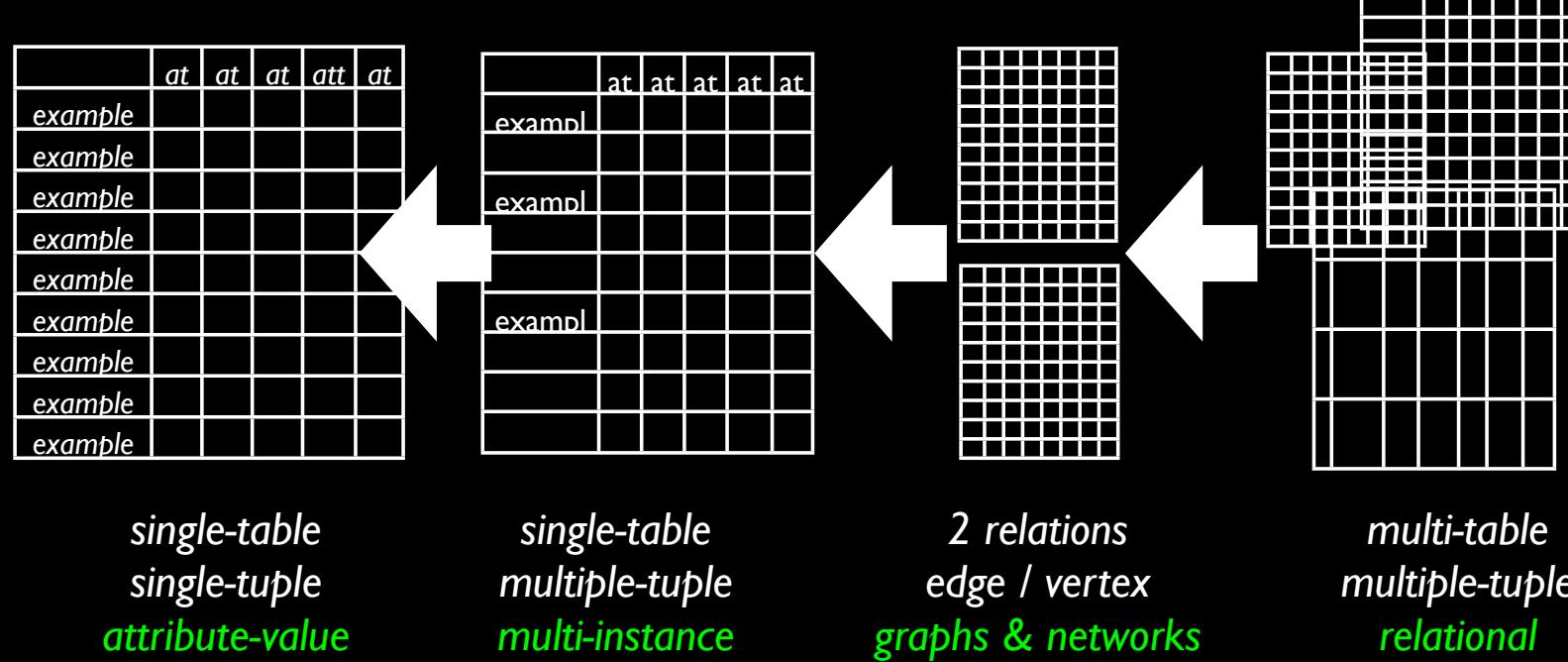
Given

- a set of examples **E, positive and negatives.**
- a background theory **B**
- a language **Le** to represent examples
- a language **Lh** to represent hypotheses
- a **probabilistic covers P** relation on **Le x Lh**

Find

- hypothesis **h*** maximizing some score based on the probabilistic covers relation; often some kind of conditional likelihood

Propositionalization



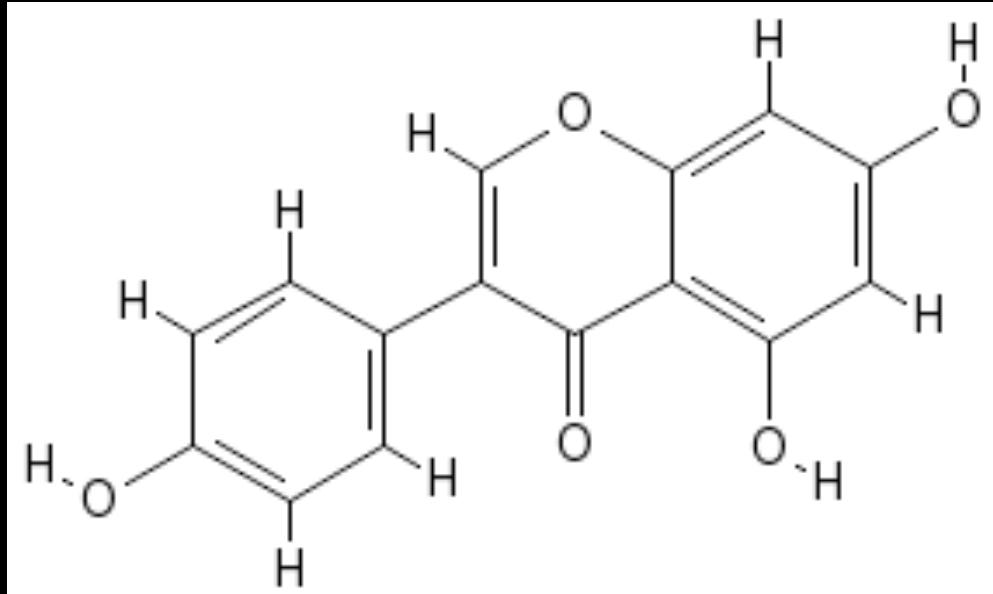
Downgrading the data ?

Two questions

UPGRADING : Can we develop systems that work with richer representations (starting from systems for simpler representations)?

PROPOSITIONALISATION: Can we change the representation from richer representations to simpler ones ? (So we can use systems working with simpler representations)

Relational Data and Clausal Theories



Examples E

$pos(m_1)$
 $neg(m_2)$
 $pos(m_3)$
 $pos(m_4)$
...

Background Knowledge B

$molecule(m_1)$

$atom(m_1, a_{11}, c)$

$atom(m_1, a_{12}, c)$

$bond(m_1, a_{12}, a_{11})$

$bond(m_1, a_{13}, a_{12})$

...

$molecule(m_2)$

$atom(m_2, a_{21}, o)$

$atom(m_2, a_{22}, n)$

$bond(m_2, a_{21}, a_{22})$

$charge(m_2, a_{21}, 0.82)$

...

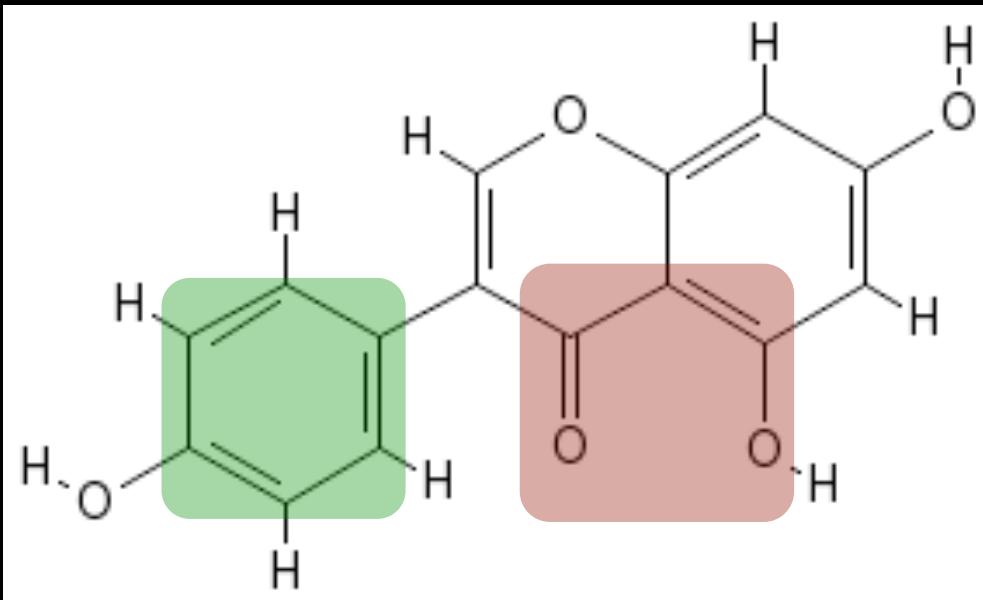
$aromaticRing(M) :-$

$bond(M, A, B, -), bond(M, B, C, =),$

$bond(M, C, D, -), bond(M, D, E, =),$

$bond(M, E, F, -), bond(M, F, A, =).$

Relational Data and Clausal Theories



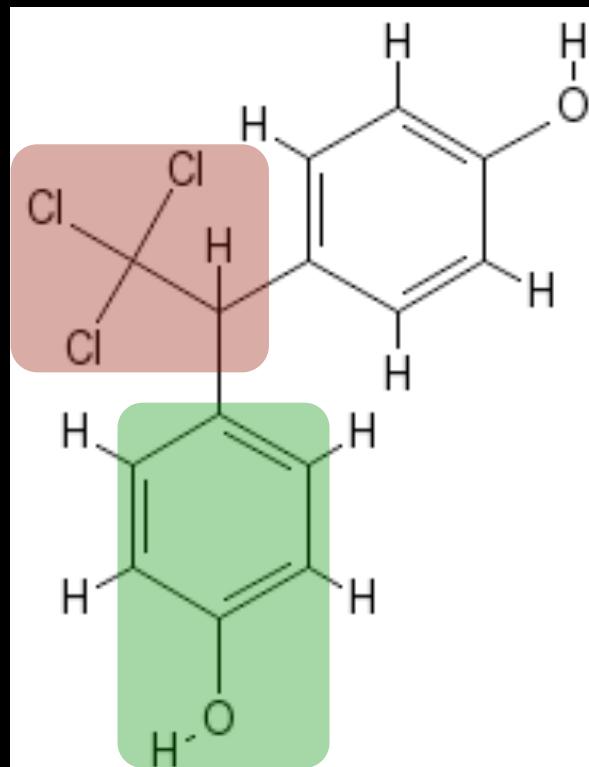
Clauses match on structural features of
e.g. molecules
Set of clauses is called theory or hypothesis

```
pos(X) ← bond(A,B,=), bond(B,C,-), bond(C,D,=)  
          bond(D,E,-), bond(E,F,=), bond(F,A,-)
```

```
pos(X) ← atom(A,o), bond(A,B,=), bond(B,C,-)  
          bond(C,D,=), bond(D,E,-), atom(E,o)
```

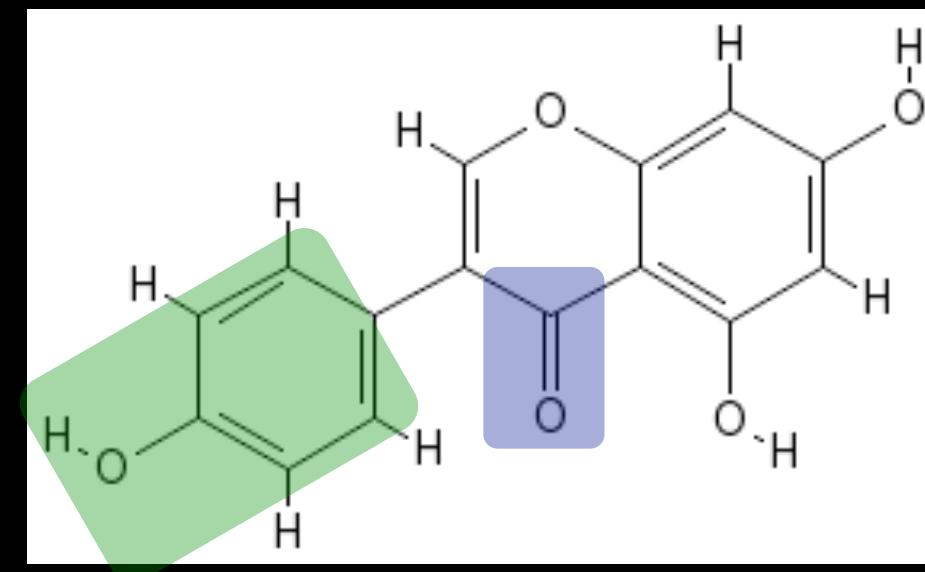
Propositionalization

Clauses c_1 , c_2 c_3



$$\varphi_H(e_1) = \begin{pmatrix} 1 \\ 1 \\ 0 \end{pmatrix}$$

$$\varphi_H(e_2) = \begin{pmatrix} 0 \\ 1 \\ 1 \end{pmatrix}$$



Logical Theories

Theory $H = \{c_1, \dots, c_k\}$ can be seen as a set of features

$$\varphi_H(e) = \begin{pmatrix} \pi_1 \\ \pi_2 \\ \dots \\ \pi_k \end{pmatrix} \quad \pi_i = \begin{cases} 1 & : c_i \text{ matches } e \\ 0 & : \text{otherwise} \end{cases}$$

Disjunction: example positive if one of the clauses matches
(hypothesis $c_1 \vee \dots \vee c_k$)

	$\varphi_H(e_1)$	$\varphi_H(e_2)$	$\varphi_H(e_3)$...
c_1	0	0	1	...
c_2	1	0	0	...
c_3	1	0	0	...
<i>class</i>	<i>pos</i>	<i>neg</i>	<i>pos</i>	...

How to learn H?

Traditional **static** propositionalization approaches

- a. ILP + SVM
 - 1. Run ILP system to generate small set of clauses
 - 2. Train SVM based on these features
- c. Pattern Mining + SVM
 - 1. Generate large set of clauses, e.g., all clauses that are frequent in the training data
 - 2. Train SVM based on these features

Problem: Two steps are independent

- we will not find the hypothesis maximizing our score
- Theory that yields optimal similarity measure different from best ILP theory!
- Information loss

nFOIL: Clauses + Probabilistic Model

nFOIL / tFOIL: combine clauses using a (tree augmented) naive Bayes model

$$\begin{aligned} P(p\theta \mid H, B) &= \frac{P(p\theta \mid c_1\theta, \dots, c_k\theta)}{P(c_1\theta, \dots, c_k\theta)} \\ &= \frac{P(c_1\theta, \dots, c_k\theta \mid p\theta)P(p\theta)}{P(c_1\theta, \dots, c_k\theta)} \\ &= \frac{\prod_{i=1}^k P(c_i\theta \mid p\theta)P(p\theta)}{P(c_1\theta, \dots, c_k\theta)} \end{aligned}$$

Instantiated class predicate Instantiated clause

Defines distribution over class variable for a given example
(substitution, instantiates to example)

kFOIL: Clauses + Kernel

kFOIL: combine clauses using a kernel machine

Idea: Theory H measures similarity

- e.g., standard dot product kernel $K_P(e_1, e_2, H) = \langle \varphi_H(e_1), \varphi_H(e_2) \rangle$

„number of features examples have in common“

- polynomial kernel $K_P(e_1, e_2, H) = (\langle \varphi_H(e_1), \varphi_H(e_2) \rangle + 1)^p$

„conjunctions of up to p features“

- Gaussian kernel $K_G(e_1, e_2, H) = \exp\left(-\frac{\|\varphi_H(e_1) - \varphi_H(e_2)\|^2}{2\sigma}\right)$

„symmetric difference: number of features which hold on only one example“

Classification can be obtained from similarity by kernel machine

How to learn H?

Looking for a hypothesis which is

- the best feature set in a naive Bayes model (nFOIL)
- defines the best similarity measure (kFOIL)

The optimization criterion is different from standard ILP

Use a score for clause search that takes into account the way the clauses are combined

- For nFOIL / tFOIL: Conditional likelihood (probabilistic coverage)

$$\prod_{e \in E} P(e | H, B) = \prod_{e \in E} \frac{\prod_{i=1}^k P(c_i \theta | e) P(e)}{P(c_1 \theta, \dots, c_k \theta)}$$

- For kFOIL: Performance of the kernel machine on the training data

The Score Function

Score of a theory H

- Theory defines kernel $K(e_i, e_j, H)$
- Parameters $\alpha_1, \dots, \alpha_n, b$ obtained by standard SVM training
- Given parameters, examples can be labeled according to

$$class(e) = sign\left(\sum_{i=1}^n \alpha_i y(e_i) K(e, e_i, H) + b\right) \quad (\text{classification})$$

$$f(e) = \sum_{i=1}^n (\alpha_i - \alpha_i^*) K(e, e_i, H) + b \quad (\text{regression})$$

- Most simple choice for $score(E, H, B)$ Accuracy/RMSE

How to learn H?

Dynamic propositionalization [Landwehr et. al. JMLR 06]

- Integrate the two steps
- Modify ILP algorithm: **search for clauses guided by performance of the combination method (NB/SVM)**

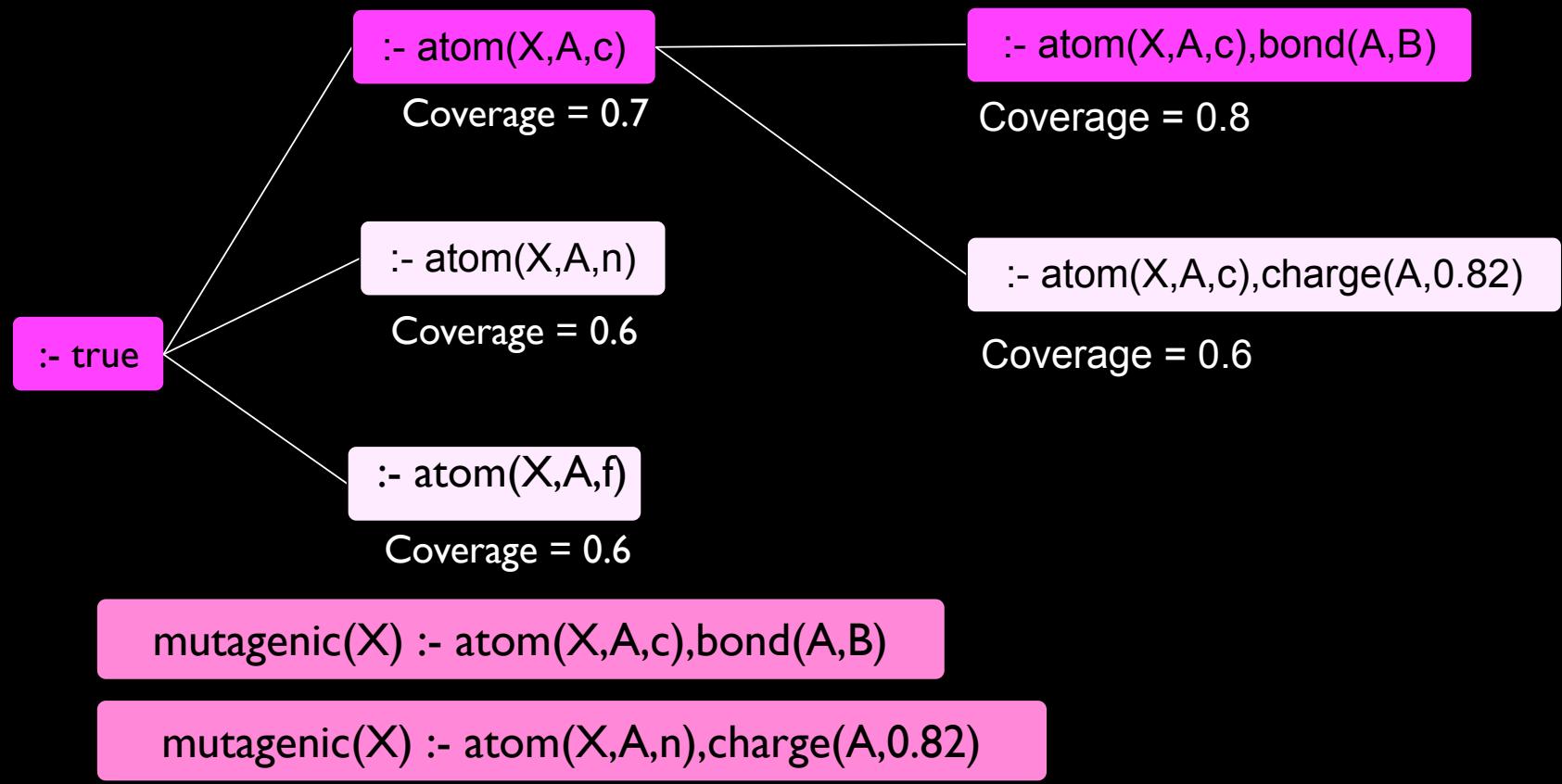
FOIL-style search

- Greedy general-to-specific search for clauses employing refinement operator

$$\rho(c) = \{c' \mid c' \text{ is minimal specialization of } c \text{ in } \mathcal{L}\}$$

- Greedy incremental search for clause set

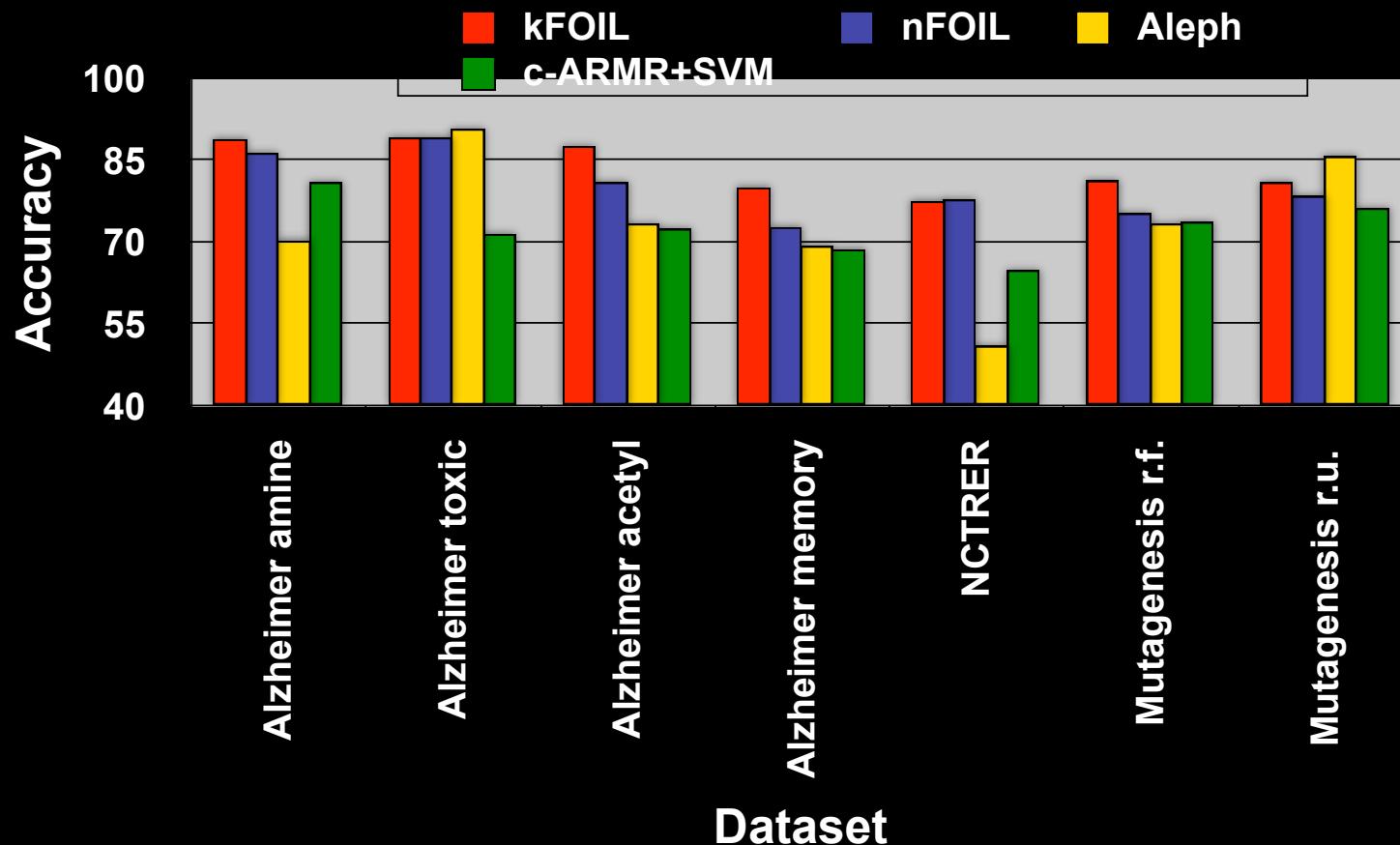
FOIL



Experimental Results: Classification

Comparison against Aleph (ILP) and c-ARMR+SVM (static propositionalization)

on several ILP benchmark datasets



Propositionalization

Dynamic propositionalization can yield more accurate models than ILP or static propositionalization

- key idea: guide clause search by score related to combination method
- typically yields clause sets which are better suited to the particular combination method

Added computational complexity if a costly clause combination method (such as SVM) is used

- However, complexity can be reduced by incrementally refining the combination model when the clause set changes
- This brings dynamic propositionalization computationally close to ILP (on the datasets we looked at)

Conclusions

A flavor of Probabilistic ILP or Statistical Relational Learning from a logical perspective

A definition of Probabilistic ILP

- based on a probabilistic coverage notion and annotated logic programs

Different Probabilistic ILP setting as for traditional ILP

- Learning from entailment (Parameter Est. SLPs/Prism)
- Learning from interpretations (BLPs, PRMs & Co)
- Learning from traces or proofs (SLPs, RMMs, LOHMMs)
- Discriminative ILP (FOIL + Naïve Bayes or Kernels)

Probabilistic Logic Learning

Excellent Area for Research & Ph.D. topics !

THANK YOU !