

ECML PKDD 2005

9th international conference on machine learning and knowledge discovery in databases



Porto, Portugal, 3 to 7 October



# Probabilistic Inductive Logic Programming

\* L. De Raedt, K. Kersting. "Probabilistic inductive Logic Programming". In S. Ben-David, J. Case and A. Maruoka, editors, Proceedings of the 15th International Conference on Algorithmic Learning Theory (ALT-2004), pages 19-36. Padova, Italy, October 2-5, 2004.

\* L. De Raedt, K. Kersting. "Probabilistic Logic Learning". In ACM-SIGKDD Explorations, special issue on Multi-Relational Data Mining, S. Dzeroski and L. De Raedt, editors, Vol. 5(1), pp. 31-48, July 2003.

Luc De Raedt, Kristian Kersting

Machine Learning Lab, Institute for Computer Science  
University of Freiburg, Germany



16th ECML, 9th PKDD, Porto, Portugal, Oct. 7th, 2005

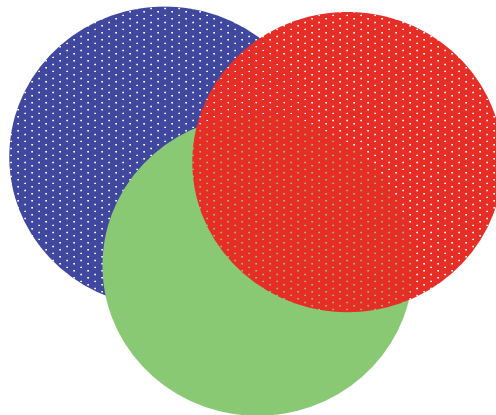


## Probabilistic Logic Learning\*

One of the key open questions of artificial intelligence concerns

"probabilistic logic learning",

i.e. the integration of probabilistic reasoning with first order logic representations and machine learning.



\*In the US, sometimes called Statistical Relational Learning 

16th ECML, 9th PKDD, Porto, Portugal, Oct. 7th, 2005



## Outline

1. Motivation / Introduction
2. Inductive Logic Programming (ILP)
  - Logic
  - Learning setting, cover relation
  - Learning from entailment, interpretations, and traces/proofs
3. Probabilistic ILP
  - Learning setting, probabilistic cover relation
4. Probabilistic Learning from
  - Interpretations, entailment, and traces/proofs
5. Discriminative ILP
6. Conclusions

16th ECML, 9th PKDD, Porto, Portugal, Oct. 7th, 2005



## Web Mining / Linked Bibliographic Data / Recommendation Systems / ...

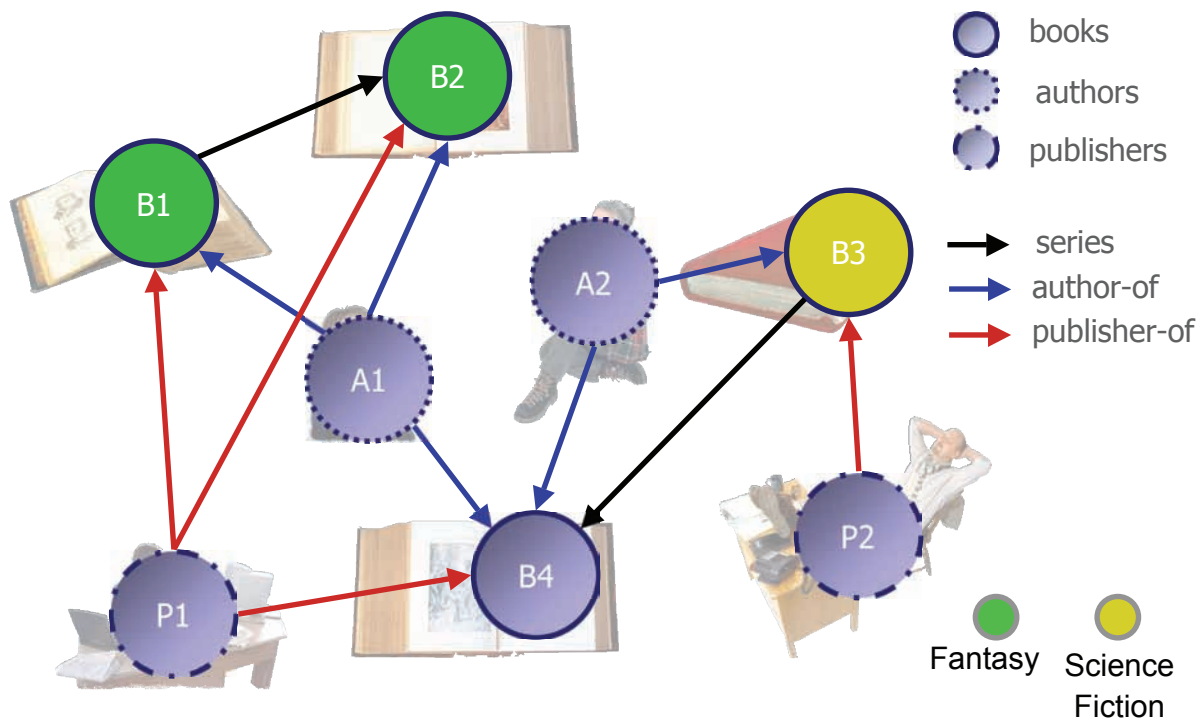


[illustration inspired by Lise Getoor]

16th ECML, 9th PKDD, Porto, Portugal, Oct. 7th, 2005



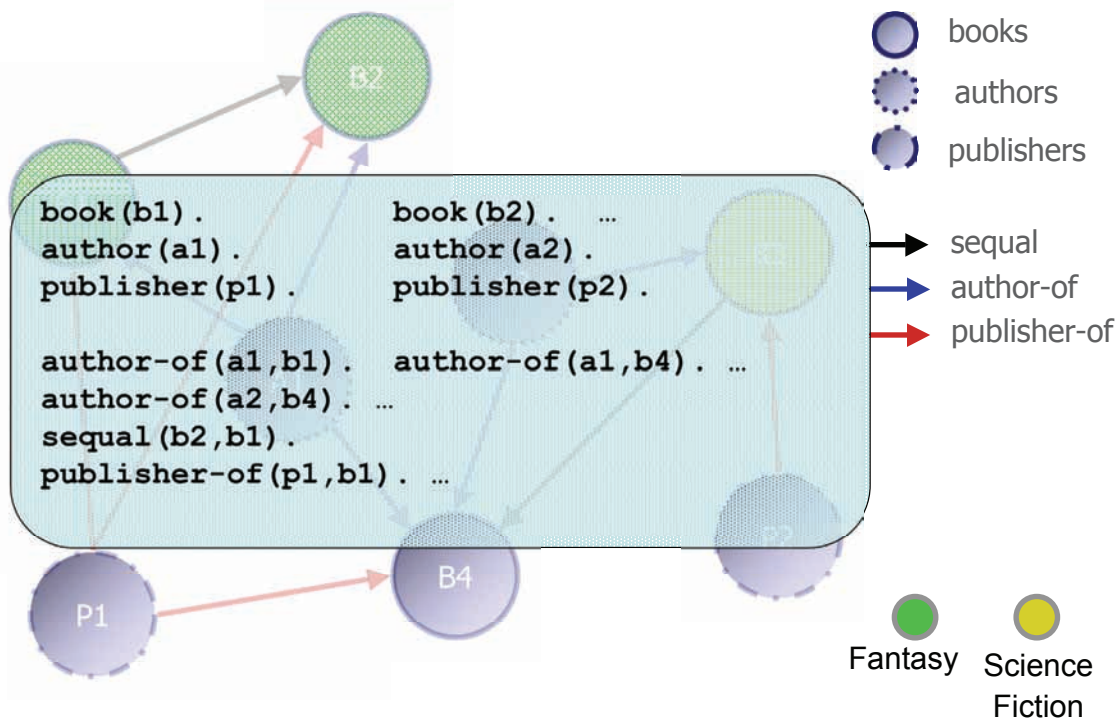
# Web Mining / Linked Bibliographic Data / Recommendation Systems / ...



16th ECML, 9th PKDD, Porto, Portugal, Oct. 7th, 2005



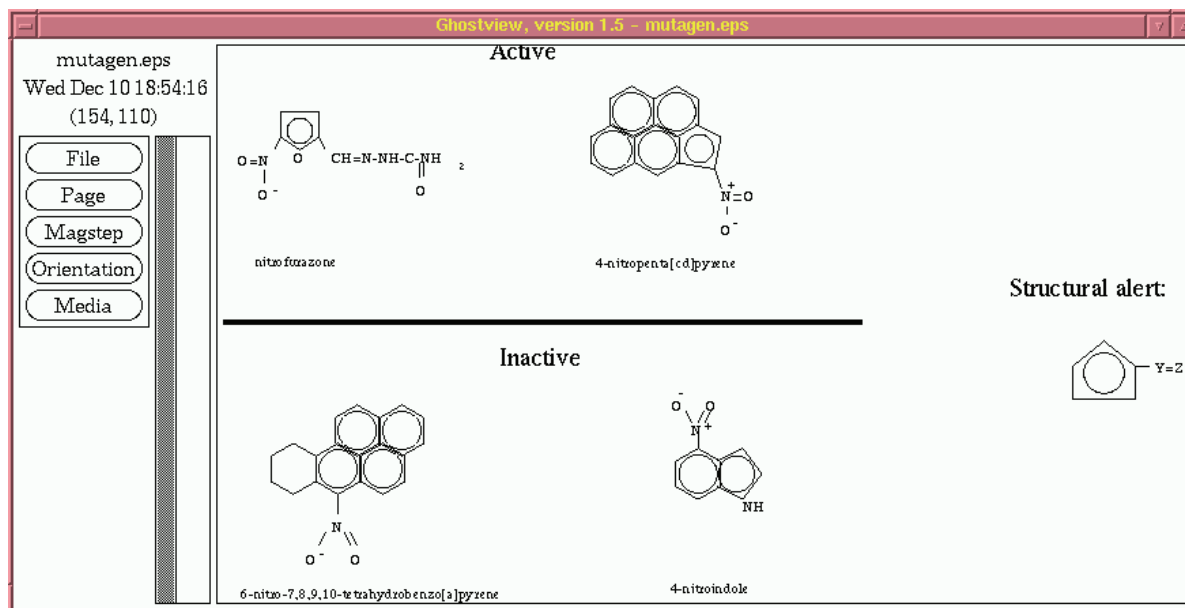
# Web Mining / Linked Bibliographic Data / Recommendation Systems / ...



16th ECML, 9th PKDD, Porto, Portugal, Oct. 7th, 2005



## Structure Activity Relationship Prediction



16th ECML, 9th PKDD, Porto, Portugal, Oct. 7th, 2005



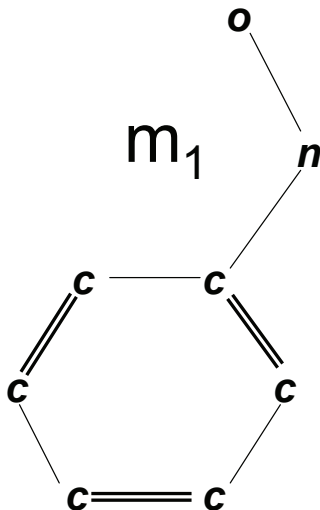
## Scientific Applications

- Discovering
  - New knowledge (readily interpretable)
  - With general purpose relational learning or inductive logic programming systems
  - Published in journals of the scientific application domain
  - Use of domain knowledge

16th ECML, 9th PKDD, Porto, Portugal, Oct. 7th, 2005



# Molecules

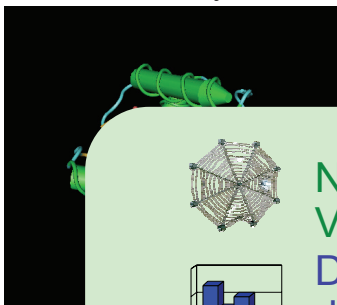


molecule ( $m_1$ )  
atom ( $m_1, a_{11}, c$ )  
atom ( $m_1, a_{12}, n$ )  
bond ( $m_1, a_{11}, a_{12}$ )  
charge ( $m_1, a_{11}, 0.82$ )  
...



# ... other Real World Applications

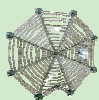
Protein Secondary Structure



Social Networks



Data Cleaning



Not flat but structured domains  
Variable #objects and relations

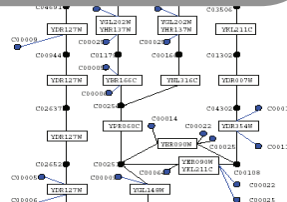
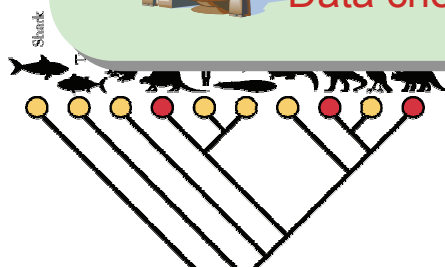


Dealing with noisy data, missing data and hidden variables



Knowledge Acquisition Bottleneck,  
Data cheap

Phylog



Interpretation





# Why do we need Probabilistic ILP\* ?

\*sometimes called statistical relational learning (SRL)

Real World Applications

## Probabilistic Logics

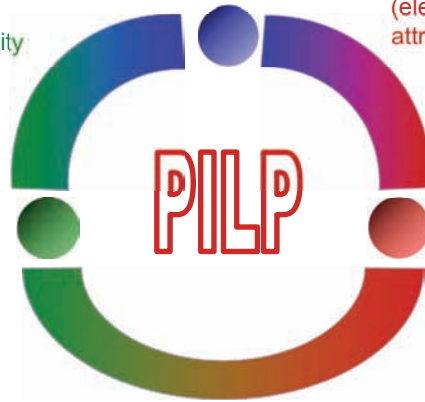
- no learning: too expensive to handcraft models
- + soft reasoning, expressivity



## Uncertainty

## Statistical Learning (SL)

- attribute-value representations: some learning problems cannot (elegantly) be described using attribute value representations
- + soft reasoning, learning



## Structured Domains



## Machine Learning



Inductive Logic Programming (ILP)  
Multi-Relational Data Mining (MRDM)

- crisp reasoning: some learning problems cannot (elegantly) be described without explicit handling of uncertainty
- + expressivity, learning

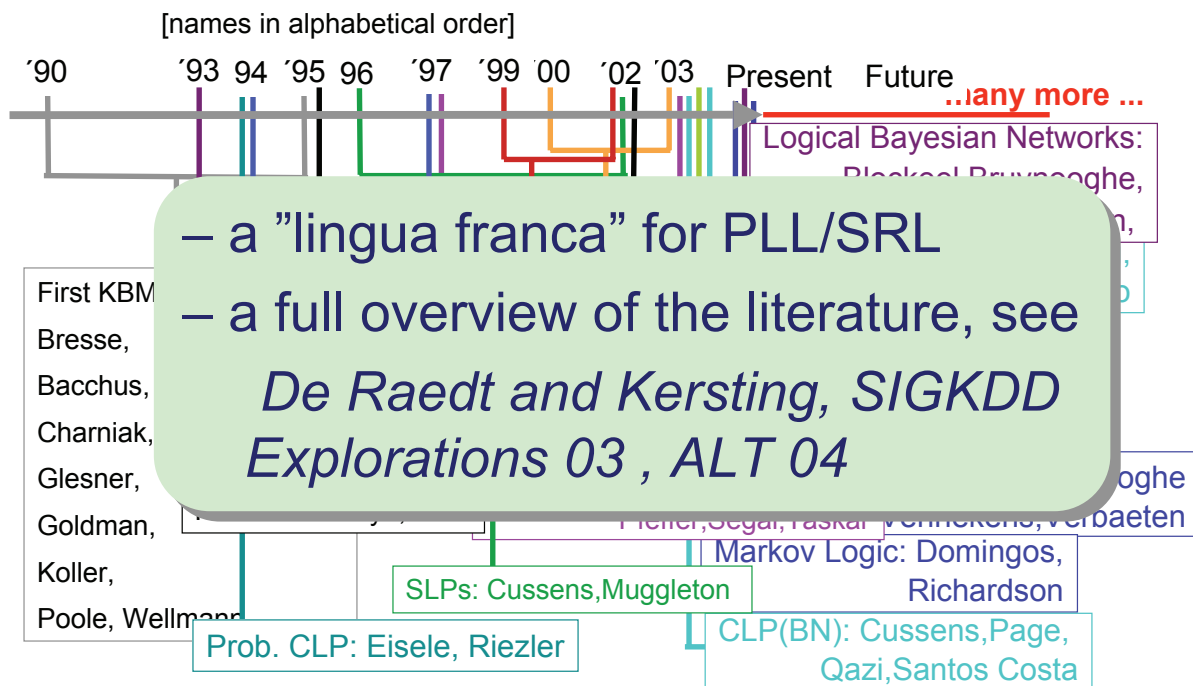


# The Tutorial's Approach is ...

- Start from **ILP settings** + **extend** them with **probabilistic methods**
  - Learning from entailment
  - Learning from interpretations
  - Learning from traces or proofs
- Hence, **probabilistic ILP**
- Provide insight in some **logical** issues
- Focus on learning ... but also relevant to KR
  - Probabilities on facts, interpretations, proofs



... but NOT ...



– a "lingua franca" for PLL/SRL

– a full overview of the literature, see *De Raedt and Kersting, SIGKDD Explorations 03, ALT 04*

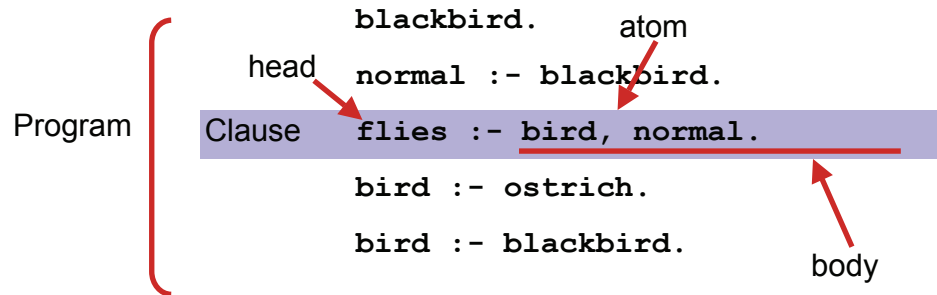


## Outline

1. Motivation / Introduction
2. Inductive Logic Programming (ILP)
  - Logic
  - Learning setting, cover relation
  - Learning from entailment, interpretations, and traces/proofs
3. Probabilistic ILP
  - Learning setting, probabilistic cover relation
4. Probabilistic Learning from
  - Interpretations, entailment, and traces/proofs
5. Discriminative ILP
6. Conclusions



## Propositional Logic



**Clauses:** IF `bird` and `normal` are true THEN `flies` is true

**Herbrand Base (HB) = all atoms in the program**

`blackbird`, `normal`, `flies`, `ostrich`, `bird`



## Model Theoretic Semantics - Restrictions on Possible Worlds -

- Herbrand Interpretation
  - Truth assignments to all elements of HB
- An interpretation is a **model** of a clause C
  - ⇔ If the body of C holds then the head holds, too.

```
blackbird.  
normal :- blackbird.  
flies :- bird, normal.  
bird :- ostrich.  
bird :- blackbird.  
  
{blackbird,normal,bird,flies}
```



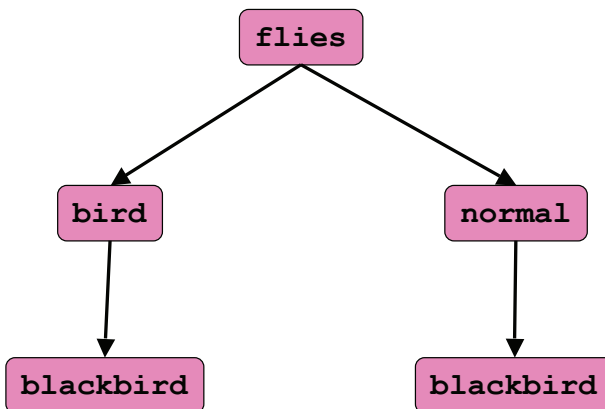


# Proof Theoretic Semantics

```

blackbird.
normal :- blackbird.
flies :- bird, normal.
bird :- ostrich.
bird :- blackbird.

```



# Upgrading - continued

## Full Clausal Logic

Functors aggregate objects

**Substitution:** Maps variables to terms:  $\{M / \text{ann}\}$ :

$\text{mc}(P, a) :- \text{mother}(\text{ann}, P), \text{pc}(\text{ann}, a), \text{mc}(\text{ann}, a).$

## Relational Clausal Logic

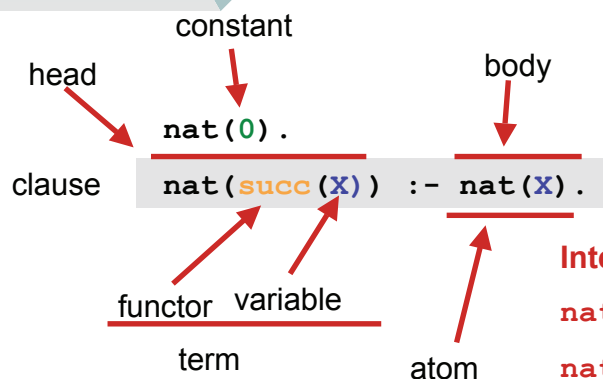
Constants and variables refer to objects

**Herbrand base:** set of ground atoms (no variables):

$\{\text{mc}(\text{fred}, \text{fred}), \text{mc}(\text{rex}, \text{fred}), \dots\}$

## Propositional Clausal Logic

Expressions can be true or false



**Interpretations can be infinite!**

$\text{nat}(0), \text{nat}(\text{succ}(0)),$   
 $\text{nat}(\text{succ}(\text{succ}(0))), \dots$



## Motivation

- We shall start from **logic + learning**
  - Inductive logic programming and relational learning
- **Methodological** practice in ILP
  - Many ILP systems obtained by **upgrading** propositional or **attribute-value learning systems**
  - Is the methodology also applicable to Probabilistic ILP ?
- Inductive logic programming has studied **several settings for learning**
  - Do they also apply to Probabilistic ILP ?



## Traditional ILP Problem

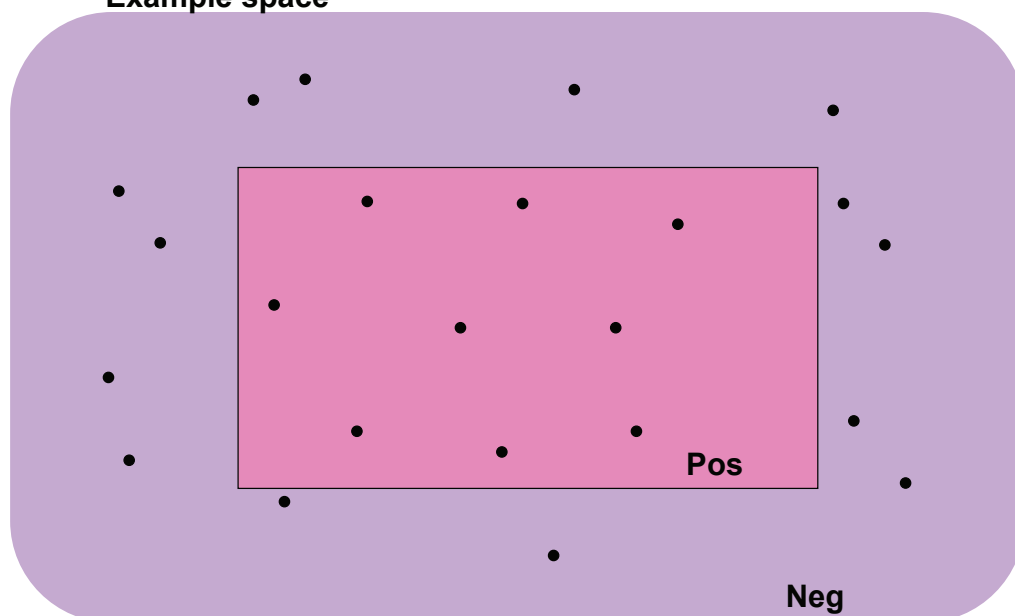
- **Given**
  - a set of positive and negative examples (**Pos, Neg**)
  - a background theory **B**
  - ?
  - ?
  - ?
- **Find**
  - A hypothesis **h** over **Lh** that covers all positive **Pos** and no negative **Neg** examples taking **B** into account

Concept-learning in a  
logical/relational representation



## Traditional ILP Problem

Example space



16th ECML, 9th PKDD, Porto, Portugal, Oct. 7th, 2005



## Three possible choices

- **Entailment**
  - Covers(H,e) iff  $H \models e$
- **Interpretations**
  - Covers(H,e) iff e is a model for H
- **Proofs**
  - Covers(H,e) iff e is a proof for H

16th ECML, 9th PKDD, Porto, Portugal, Oct. 7th, 2005

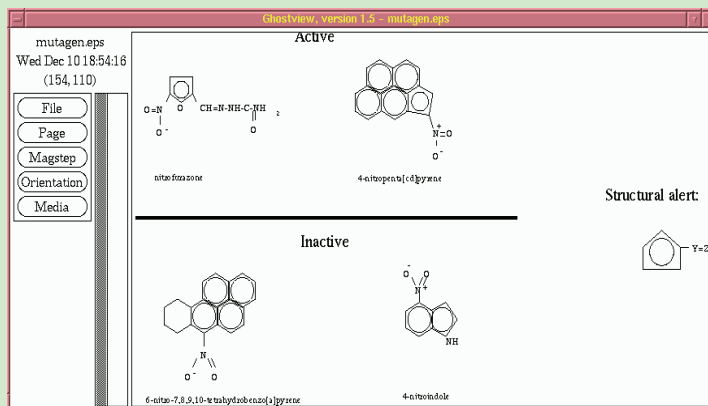


# Learning from entailment

- Examples are **facts** (or clauses)
- An example **e** is **covered** by a hypothesis **h** if and only if  $B U h \models e$

## Applications

- We
- F
- G



vasan),



# 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(225,1,1).
atom(225,2,1).
atom(225,3,1).
atom(225,4,1).
atom(225,5,1).
atom(225,6,1).
atom(225,7,1).
atom(225,8,1).
atom(225,9,1).
atom(225,10,1).
atom(225,11,1).
atom(225,12,1).
atom(225,13,1).
atom(225,14,1).
atom(225,15,1).
atom(225,16,1).
atom(225,17,1).
atom(225,18,1).
atom(225,19,1).
atom(225,20,1).
atom(225,21,1).
atom(225,22,1).
atom(225,23,1).
atom(225,24,1).
atom(225,25,1).

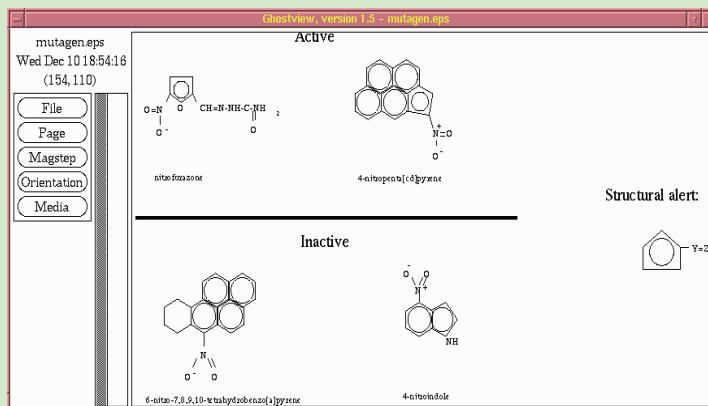
```

```

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_1,f1_9,2).
bond(225,f1_9,f1_10,2).
bond(225,f1_10,f1_11,1).
bond(225,f1_11,f1_12,2).
bond(225,f1_12,f1_13,1).

```

## Applications



hypo

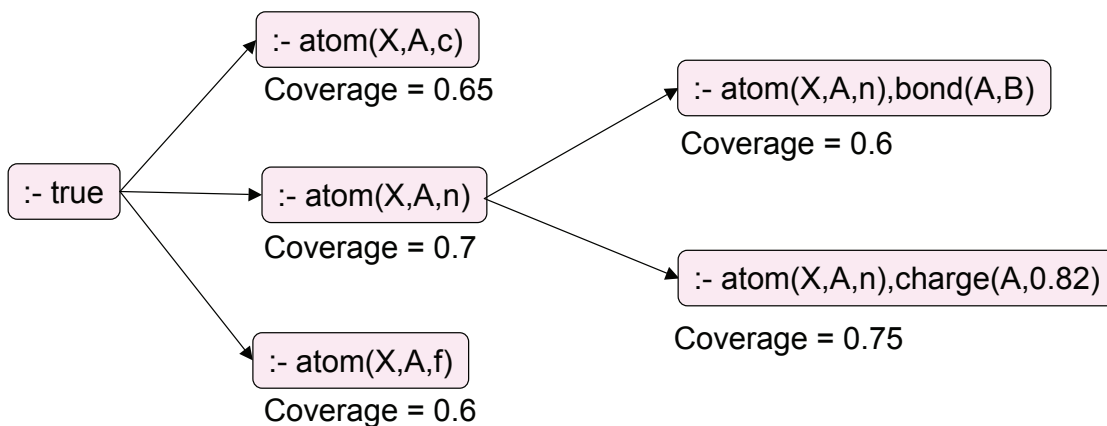
mut

Examp



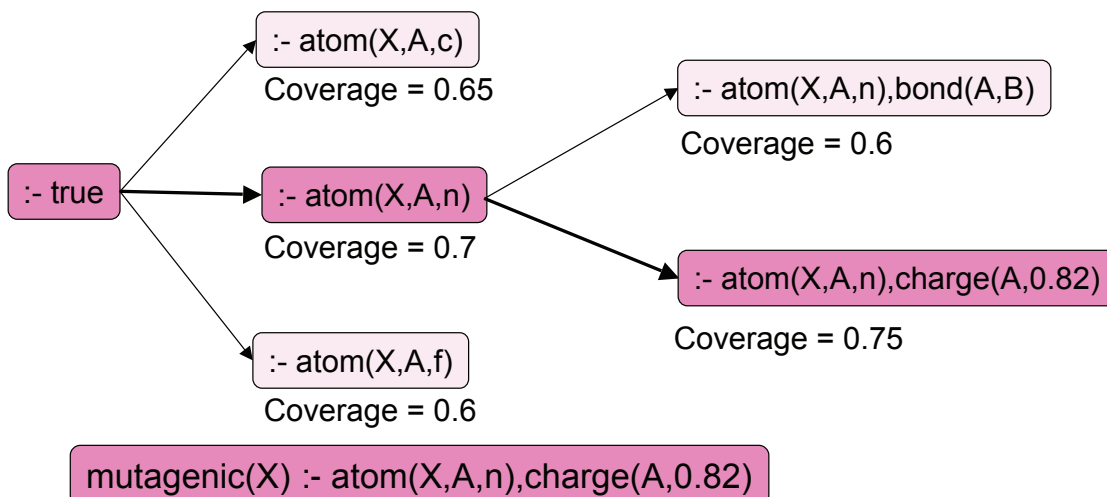
## Example ILP Algorithm: FOIL (Quinlan 1990)

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



## Example ILP Algorithm: FOIL (Quinlan 1990)

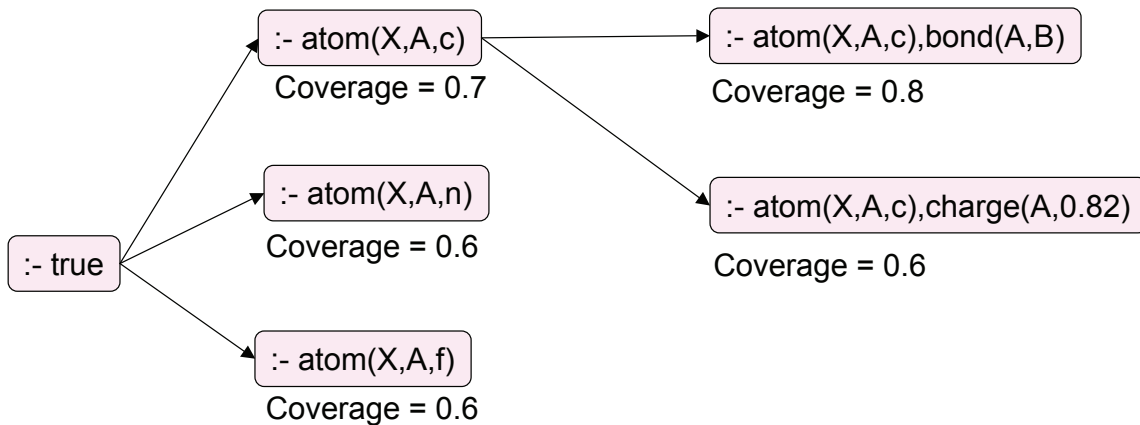
- Greedy separate-and-conquer search for clause set
- Greedy general-to-specific search for single clause





## Example ILP Algorithm: FOIL (Quinlan 1990)

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

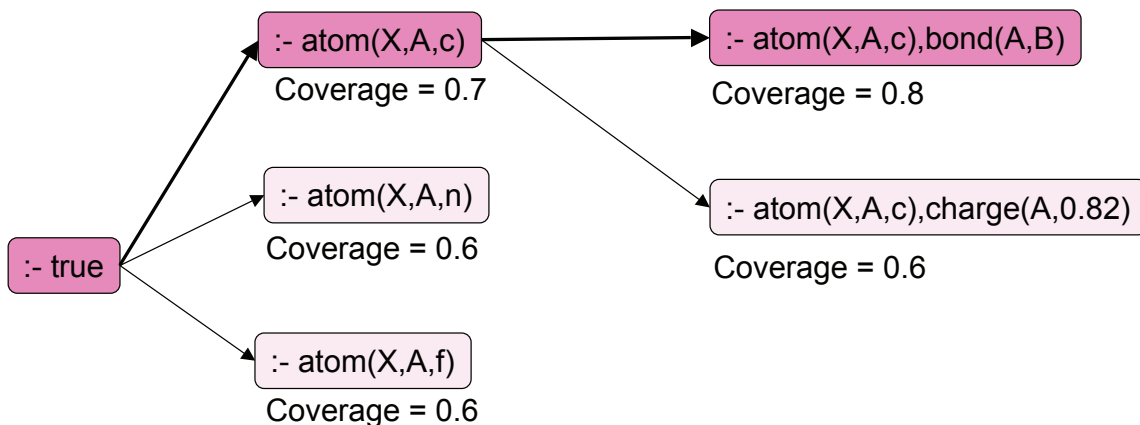


`mutagenic(X) :- atom(X,A,n),charge(A,0.82)`



## Example ILP Algorithm: FOIL (Quinlan 1990)

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



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

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



## Example ILP Algorithm: FOIL (Quinlan 1990)

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

⋮

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)



## ILP Principles

- Searching the space of hypotheses
  - Using ordering such as theta-subsumption
  - Using refinement operators
- Limitations of traditional machine learning
  - Dealing with structured data instead of feature vector, attribute value, boolean, propositional etc. representations
  - Employing background knowledge
  - Interpretability of results
- Application areas
  - Chemo- and bio-informatics, e.g. predictive toxicology
  - Language learning and information retrieval
  - Ecological applications
  - ....



## Learning from Interpretations

- Examples are (Herbrand) **interpretations**, i.e., sets of ground facts
- An example  $e$  is **covered** by a hypothesis  $h$  if and only if the example is a **model** for the hypothesis  $h$
- Well known examples
  - Logan-H (Khardon), Claudien (De Raedt and Dehaspe), Warmr (Dehaspe), ...



## An example

- **Examples**
  - **Positive**: { human(luc), human(lieve), male(luc), female(lieve) }
  - **Negative**: { bat(dracula), male(dracula), vampire(dracula) }
  - ...
- **Discriminative Hypothesis**
  - Learning from positives and negatives (possibles / impossibles)
  - human(X) :- male(X)
- Interpretation  $I$  is a **model** for clause  $h :- b_1, \dots, b_n$  iff for all  $\theta$  such that  $\{b_1\theta, \dots, b_n\theta\} \subseteq I$ , we have that  $h\theta$  in  $I$
- Consider  $\{X=dracula\}$  for negative





## An example

- **Examples**
  - **Positive:** { human(luc), human(lieve), male(luc), female(lieve)}
- **Characteristic Hypothesis**
  - Learning from positives (possibles) only
  - human(X) :- female(X)
  - human(X) :- male(X)

### Applications

- Finding integrity constraints / frequent patterns in relational databases



## Learning from Traces/Proofs

- Examples are **proof trees**
- An example e is **covered** by a hypothesis h if and only if e is a legal **proof tree** in h

### Applications

- W
  - Tree bank grammar learning
  - Program synthesis
    - Shapiro's MIS poses queries in order to reconstruct trace or proof



## 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(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).
```

## Applications

- Tree bank grammar learning
- Program synthesis
  - Shapiro's MIS poses queries in order to reconstruct trace or proof

a(pl,  
t([the,



## Use of different Settings

Information ↑  
+

Learning from entailment

Different settings provide different levels of information about target program (cf. De Raedt, AIJ 97)

Learning from traces

- Typically used for hard problems, when other settings seem to fail or fail to scale up
- E.g., program synthesis from examples, grammar induction, multiple predicate learning



## Outline

1. Motivation / Introduction
2. Inductive Logic Programming (ILP)
  - Logic
  - Learning setting, cover relation
  - Learning from entailment, interpretations, and traces/proofs
3. Probabilistic ILP
  - Learning setting, probabilistic cover relation
4. Probabilistic Learning from
  - Interpretations, entailment, and traces/proofs
5. Discriminative ILP
6. Conclusions



## Probabilistic ILP: What Changes?

- Clauses annotated with **probability labels**
  - E.g. in Sato's **Prism**, Eisele and Muggleton's **SLPs**, Kersting and De Raedt's **BLPs**, ...
- Prob. covers relation  $\text{covers}(e, H \cup B) = P(e \mid H, B)$ 
  - Probability 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 ILP Problem

### Given

- a set of examples **E**
- a background theory **B**
- a language **L<sub>e</sub>** to represent examples
- a language **L<sub>h</sub>** to represent hypotheses
- a **probabilistic covers P** relation on **L<sub>e</sub> x L<sub>h</sub>**

### Find

- hypothesis **h\*** maximizing some score based on the probabilistic covers relation



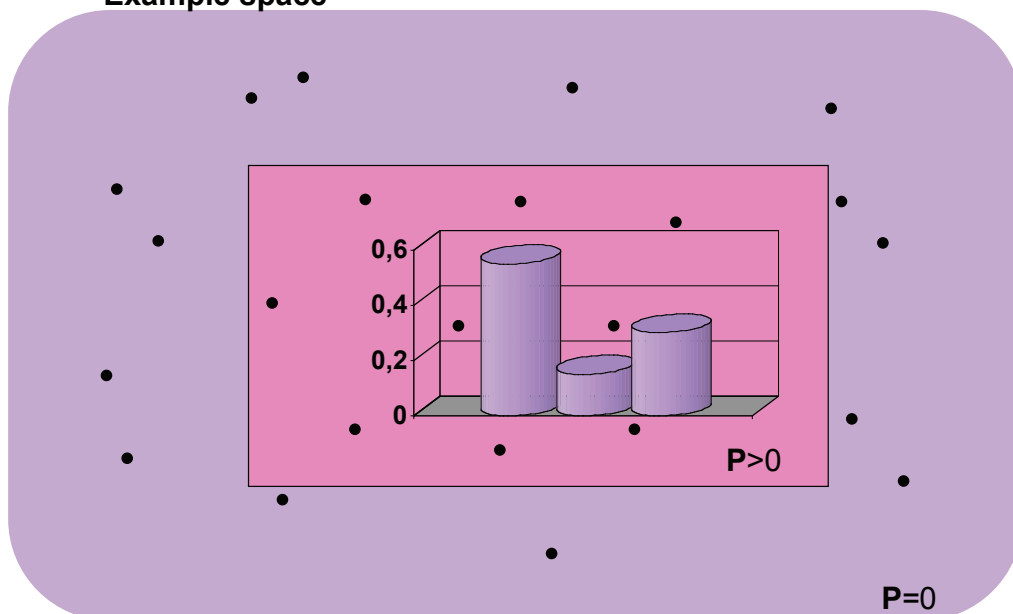
## Probabilistic ILP: Three Issues

- **Defining L<sub>h</sub> and P**
  - Clauses + Probability Labels
- **Learning Methods**
  - **Parameter Estimation**
    - Learning probability labels for fixed clauses
  - **Structure learning**
    - Learning both components



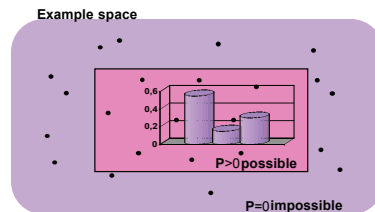
## Probabilistic ILP Problem

### Example space



## Probabilistic ILP: Two Objectives

- **Generative Learning**
  - Estimate joint probability distribution
  - E.g., likelihood + iid
$$h^* = \arg \max_h \mathbf{P}(\mathbf{e}|h, B)$$
$$= \arg \max_h \prod_i \mathbf{P}(e_i|h, B)$$

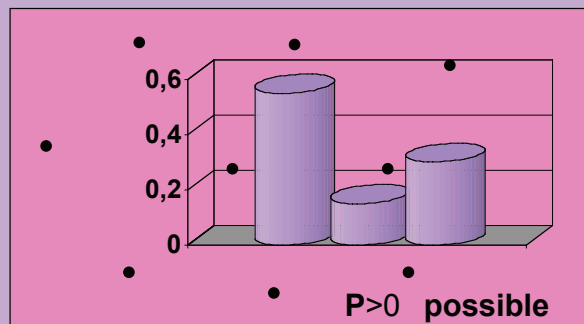


- **Discriminative Learning**
  - Estimate conditional prob. distribution over some predicates given evidence for the others



# Probabilistic ILP Problem

Example space

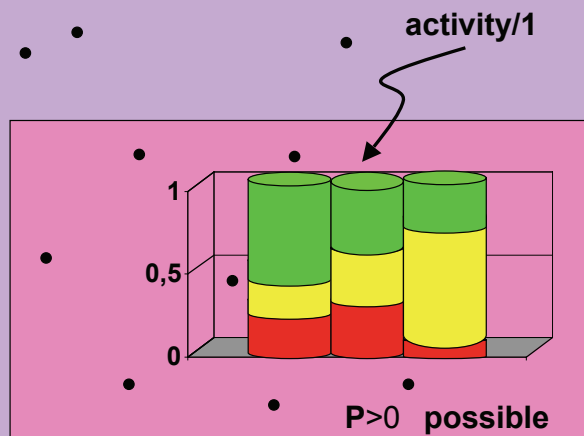


P=0 impossible



# Probabilistic ILP Problem

Example space



P=0 impossible

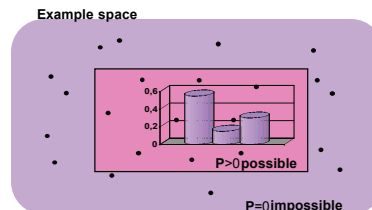


## Probabilistic ILP: Two Objectives

- **Generative Learning**

- Estimate joint probability distribution
- E.g., likelihood

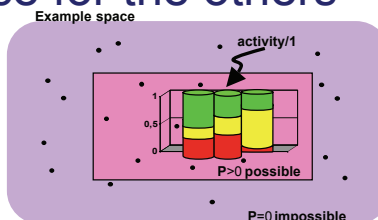
$$h^* = \arg \max_h \mathbf{P}(e|h, B)$$



- **Discriminative Learning**

- Estimate conditional prob. distribution over some predicates given evidence for the others
- E.g., conditional likelihood

$$c^* = \arg \max_c \mathbf{P}(e|c, h \setminus c, B)$$
$$= \arg \max_c \prod_i \mathbf{P}(e_i|c, h \setminus c, B)$$



## Probabilistic ILP: Three Settings

- **Probabilistic learning from entailment**

- Eichele and Muggleton's Stochastic Logic Programs, Sato's Prism, Poole's ICL

- **Probabilistic learning from proofs**

- Learning the structure of SLPs; a tree-bank grammar based approach, Logical HMMs, Anderson's RMMs

- **Probabilistic learning from interpretations**

- Bayesian logic programs, Koller's PRMs, Domingos' MLNs, Vennekens' LPADs, Taskar's RMNs



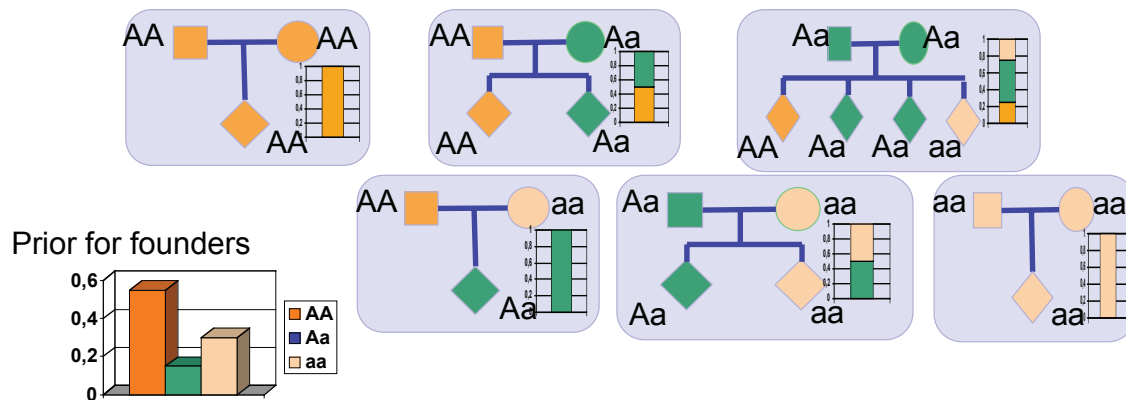
## Outline

1. Motivation / Introduction
2. Inductive Logic Programming (ILP)
  - Logic
  - Learning setting, cover relation
  - Learning from entailment, interpretations, and traces/proofs
3. Probabilistic ILP
  - Learning setting, probabilistic cover relation
4. Probabilistic Learning from
  - Interpretations, entailment, and traces/proofs
5. Discriminative ILP
6. Conclusions



## Blood Type / Genetics/ Breeding

- 2 Alleles: A and a
- Probability of Genotypes AA, Aa, aa ?







# Bayesian Networks [Pearl]

[illustration inspired by Kevin Murphy]

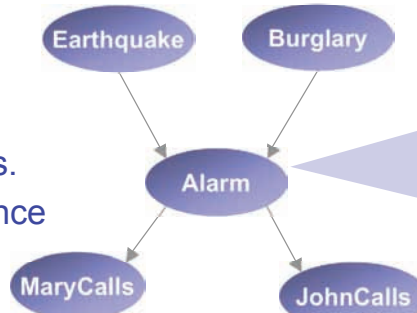
Compact representation of joint probability distributions

$P(E,B,A,M,J)$

## Qualitative part:

Directed acyclic graph

- Nodes - random vars.
- Edges - direct influence



E	B	P(A   B,E)	
e	b	0.9	0.1
e	$\bar{b}$	0.2	0.8
$\bar{e}$	b	0.9	0.1
$\bar{e}$	$\bar{b}$	0.01	0.99

## Together:

Define a unique distribution in a compact, factored form

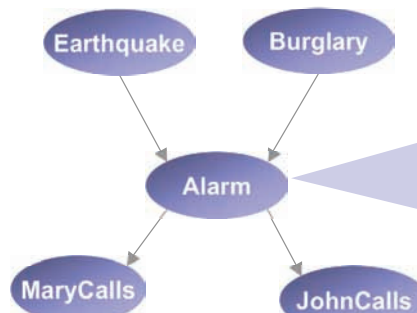
$$P(E,B,A,M,J) = P(E) * P(B) * P(A|E,B) * P(M|A) * P(J|A)$$

## Quantitative part:

Set of conditional probability distributions



# Bayesian Networks [Pearl]



E	B	P(A   B,E)	
e	b	0.9	0.1
e	$\bar{b}$	0.2	0.8
$\bar{e}$	b	0.9	0.1
$\bar{e}$	$\bar{b}$	0.01	0.99

$$\begin{aligned}
 P(j) = & P(j|a) * P(m|a) * P(a|e,b) * P(e) * P(b) \\
 & + P(j|a) * P(m|a) * P(a|e,\bar{b}) * P(e) * P(\bar{b}) \\
 & \dots \\
 & + P(j|\bar{a}) * P(\bar{m}|\bar{a}) * P(\bar{a}|\bar{e},\bar{b}) * P(\bar{e}) * P(\bar{b})
 \end{aligned}$$



# Bayesian nets: Parameter Estimation

complete data set  
simply counting

A1	A2	A3	A4	A5	A6	
true	true	false	true	false	false	X1
false	true	true	true	false	false	X2
...	...	...	...	...	...	⋮
true	false	false	false	true	true	XM



# Bayesian nets: Parameter Estimation

incomplete data set

Real-world data: states of some random variables are missing

A1	A2	hidden/ latent A3	A4	A5	A6
true	true	?	true	false	false
?	true	?	?	false	false
...	...	...	...	...	...
true	false	?	false	true	?

- E.g. medical diagnose: not all patient are subjects to all test
- Parameter reduction, e.g. clustering, ...

missing value



## Expectation-Maximization (EM): Idea

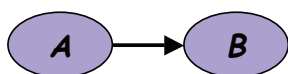
- In the case of complete data, ML parameter estimation is easy:
  - simply counting (1 iteration)

### Incomplete data ?

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



## EM Idea: Complete the data



incomplete data

A	B
true	true
true	?
false	true
true	false
false	?

$P(B = \text{true} | A = \text{true}) = 0.6$   
 $P(B = \text{true} | A = \text{false}) = 0.2$

complete data

A	B	N
true	true	1.6
true	false	1.4
false	true	1.2
false	false	0.8

expected counts

count + iterate

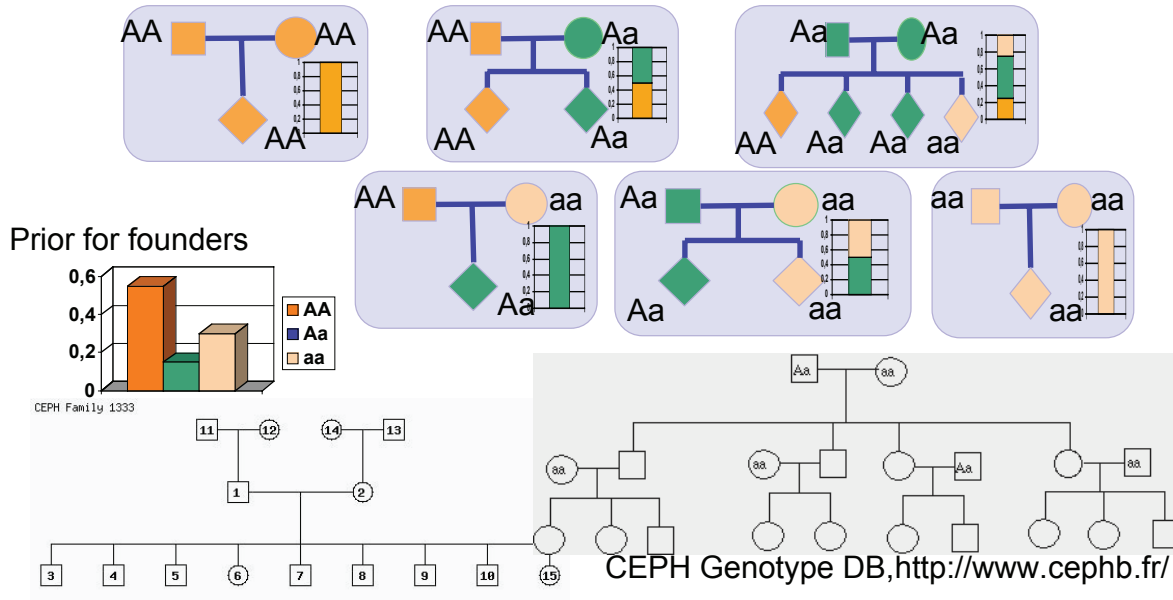
$$\theta_{B=\text{true}|A=\text{true}} = \frac{1.6}{1.6 + 1.4} = 0.54$$

$$\theta_{B=\text{true}|A=\text{false}} = \frac{1.2}{1.2 + 0.8} = 0.6$$



# Blood Type / Genetics/ Breeding

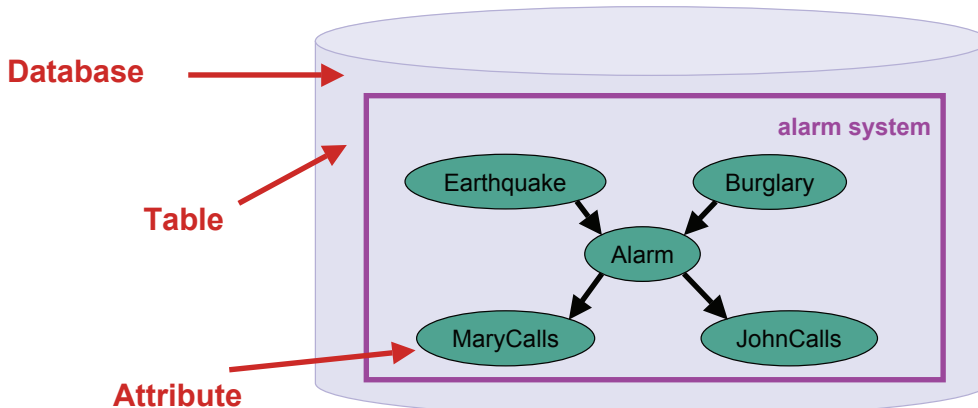
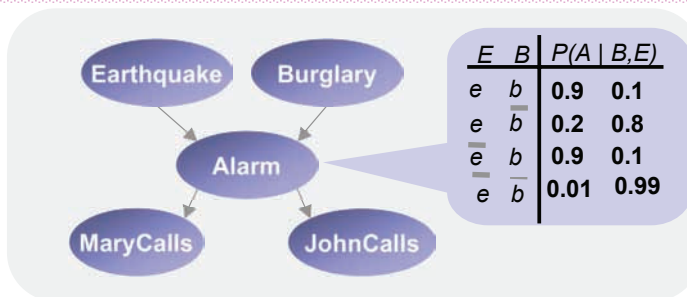
- 2 Alleles: A and a
- Probability of Genotypes AA, Aa, aa ?



# Probabilistic Relational Models (PRMs)

[Getoor, Koller, Pfeffer]

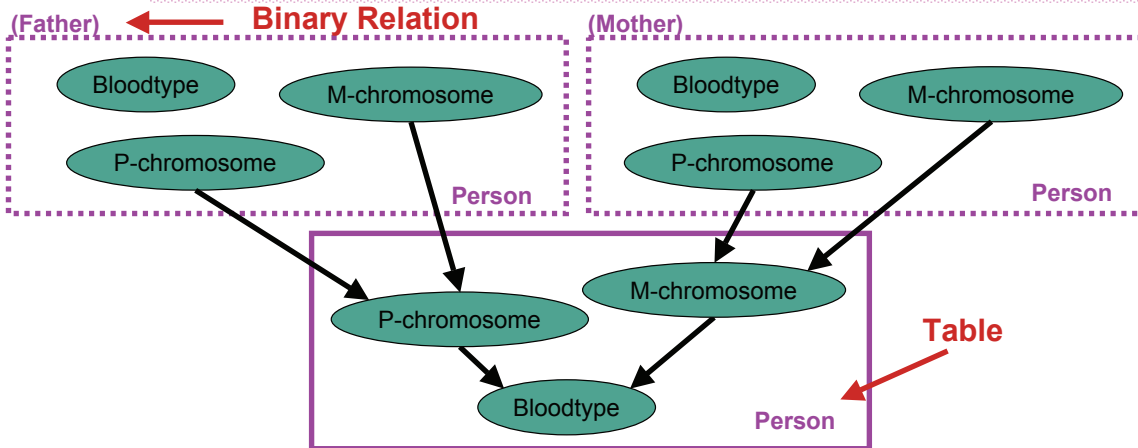
- Database theory
- Entity-Relationship Models
  - Attributes = RVs





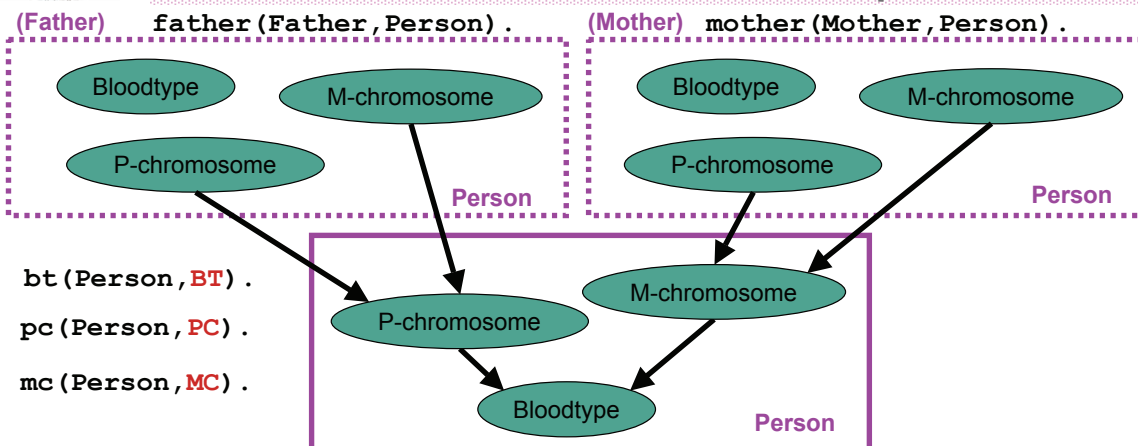
# Probabilistic Relational Models (PRMs)

[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)

[Getoor, Koller, Pfeffer]

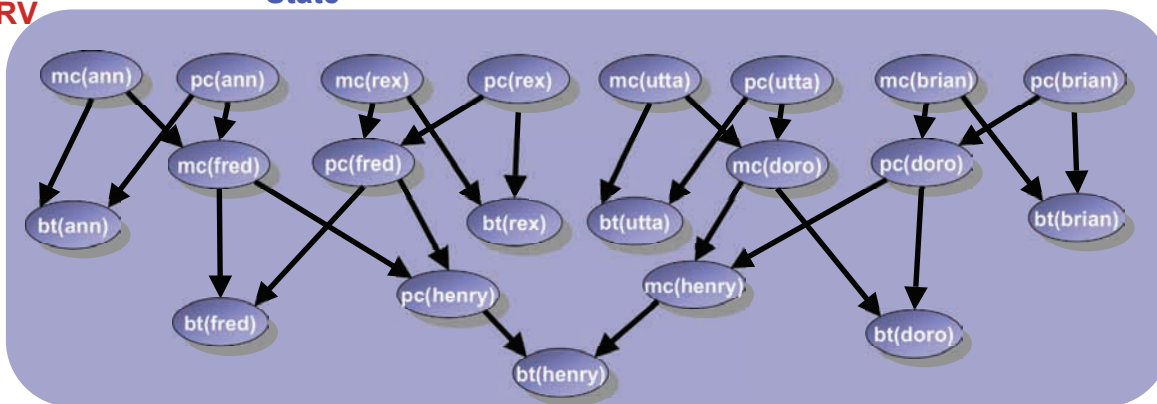
```

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

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 ← (red arrow)  
 ← (blue arrow) State



16th ECML, 9th PKDD, Porto, Portugal, Oct. 7th, 2005



[Segal et al.]

# Gene Regulation

- System Biology
- Gene expression: two-phase process
  1. Gene is transcribed into mRNA Measured by gene expression microarrays
  2. mRNA is translated Protein
- Genes that are similar expressed are often coregulated and involved in the same cellular processes
- Clustering: identification of clusters of genes and/or experiments that share similar expression patterns

16th ECML, 9th PKDD, Porto, Portugal, Oct. 7th, 2005



## PRM Application: Gene Regulation

[Segal et al.]

- System Biology: heterogenous data
- Limitations of Clustering:
  - Similarities over all measurements
  - Difficult to incorporate readily background knowledge such as clinical data or experimental details

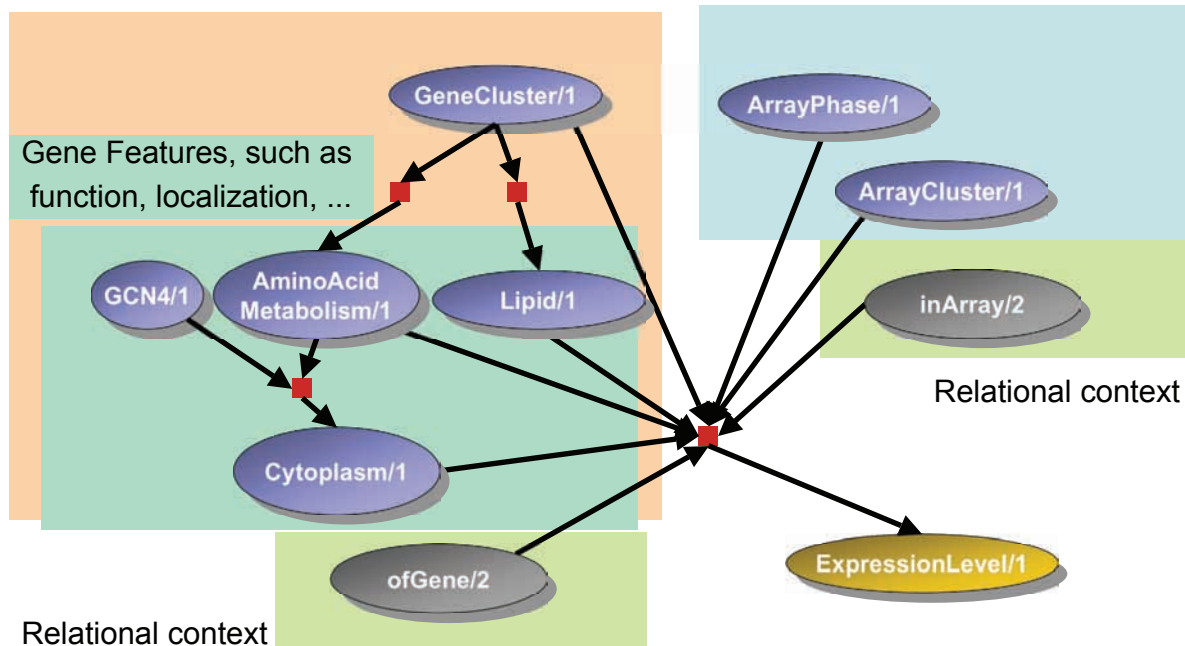


## PRM Application: Gene Regulation

[Segal et al., simplified representation]

Gene Cluster

Array Cluster





## PRM Application: Gene Regulation

[Segal et al.]

- Synthetic data: 1000 genes, 90 arrays (= 90.000 measurements), each gene 15 functions and 30 transcription factors.

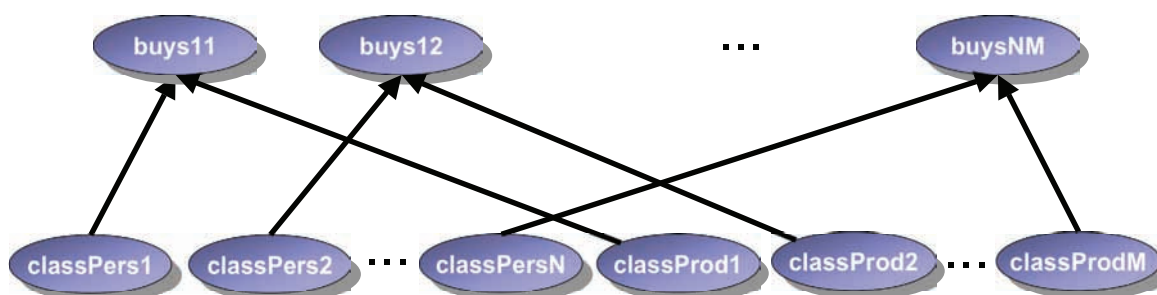
	Cluster recovery	
	Naive Bayes	PRMs
Simulated data	90.8±0.42	98.4±1.07
Noisy simulated data	76.7±1.42	88.1±1.52



[Getoor, Sahami]

## PRM Application: Collaborative Filtering

- User preference relationships for products / information.
- Traditionally: single dyadic relationship between the objects.



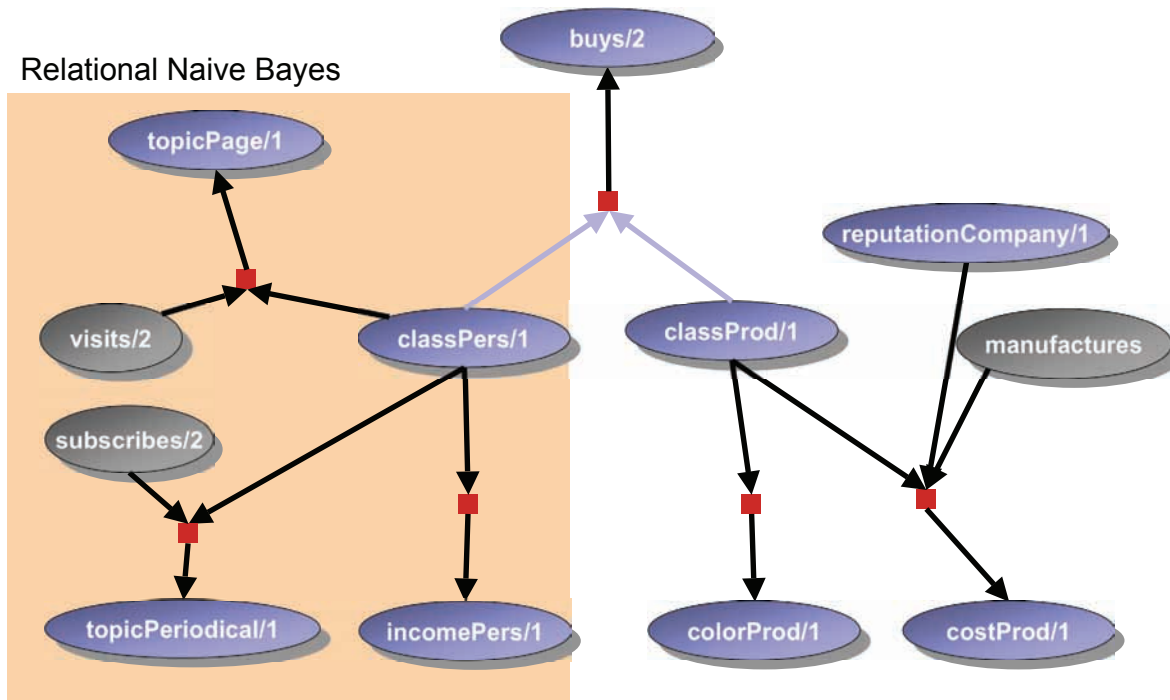




## PRM Application: Collaborative Filtering

[Getoor, Sahami; simplified representation]

Relational Naive Bayes



16th ECML, 9th PKDD, Porto, Portugal, Oct. 7th, 2005



## Probabilistic Relational Models (PRMs)

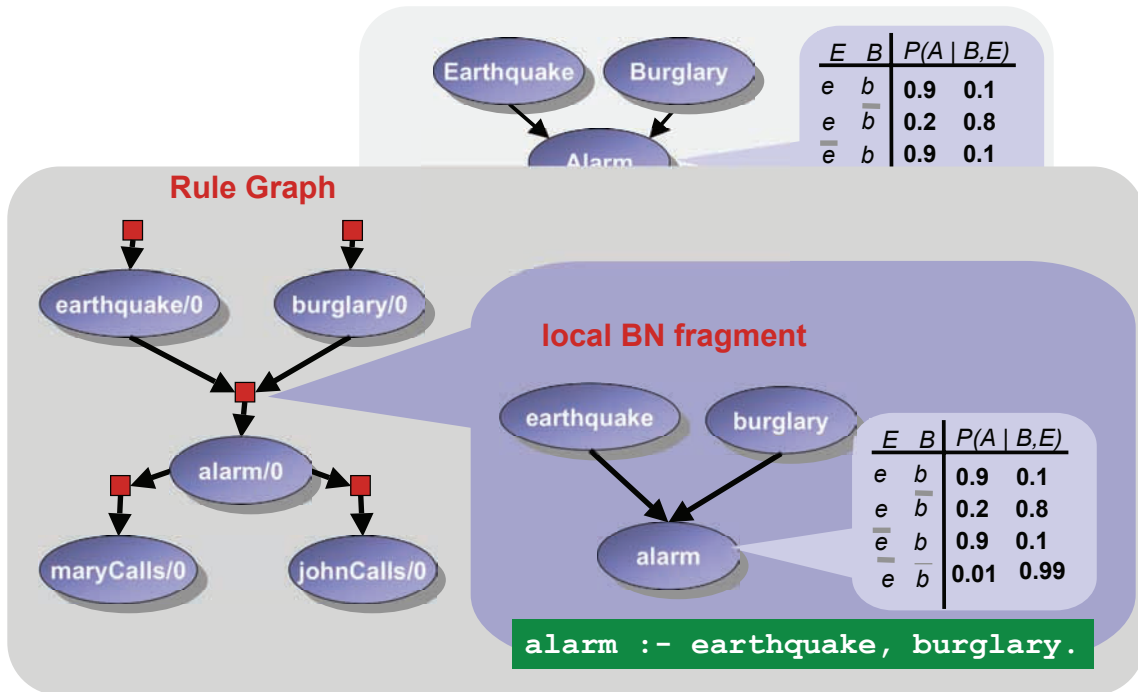
[Koller, Pfeffer, Getoor]

- Database View
- Unique Probability Distribution over finite Herbrand interpretations
  - No self-dependency
- Discrete and continuous RV
- BN used to do inference
- Graphical Representation
- BNs + extensions: hierarchical PRMs, dynamic PRMs, relational uncertainty types
- Learning

16th ECML, 9th PKDD, Porto, Portugal, Oct. 7th, 2005



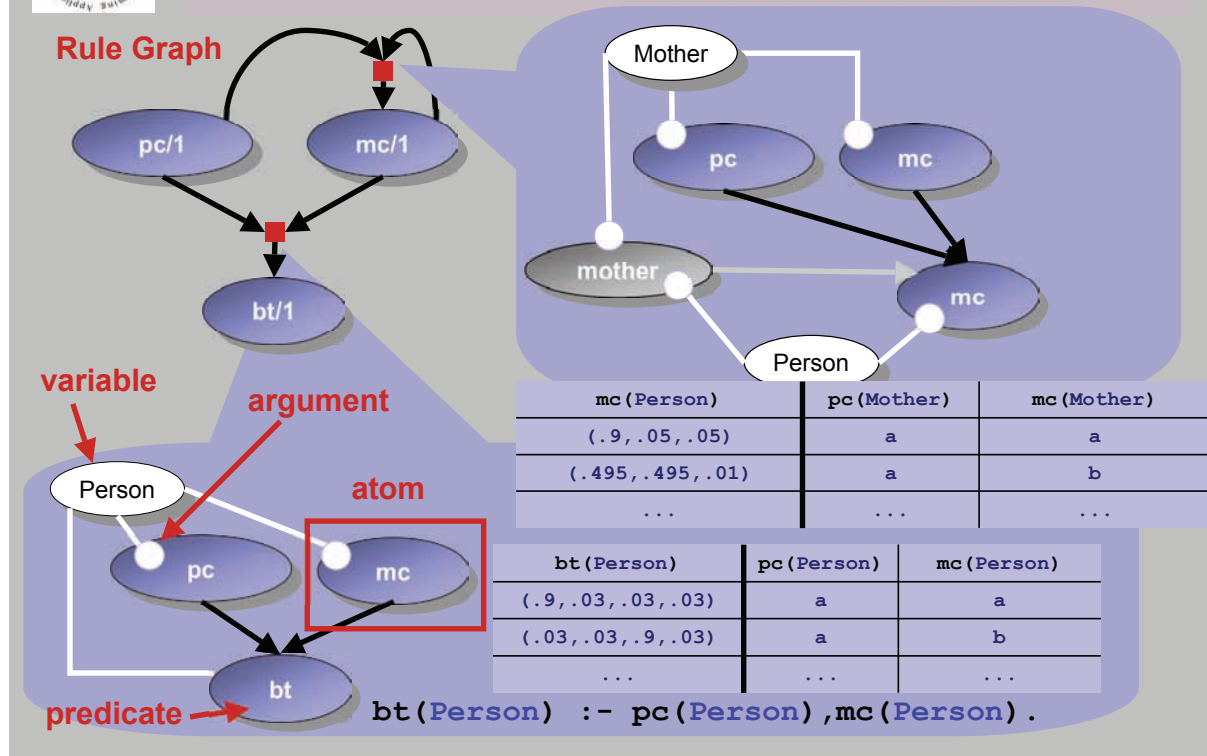
# Bayesian Logic Programs (BLPs)



16th ECML, 9th PKDD, Porto, Portugal, Oct. 7th, 2005



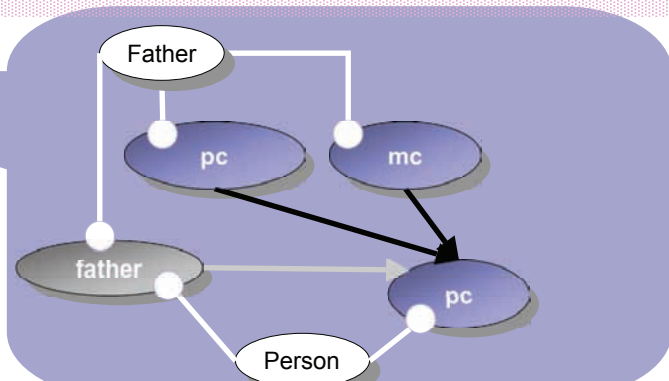
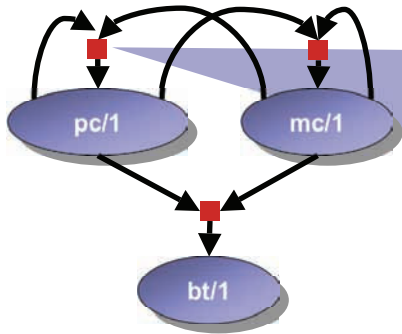
# Bayesian Logic Programs (BLPs)



16th ECML, 9th PKDD, Porto, Portugal, Oct. 7th, 2005



# Bayesian Logic Programs (BLPs)



pc(Person)	pc(Father)	mc(Father)
(.9, .05, .05)	a	a
(.495, .495, .01)	a	b
...	...	...

$mc(Person) \mid mother(Mother, Person), pc(Mother), mc(Mother).$   
 $pc(Person) \mid father(Father, Person), pc(Father), mc(Father).$   
 $bt(Person) \mid pc(Person), mc(Person).$

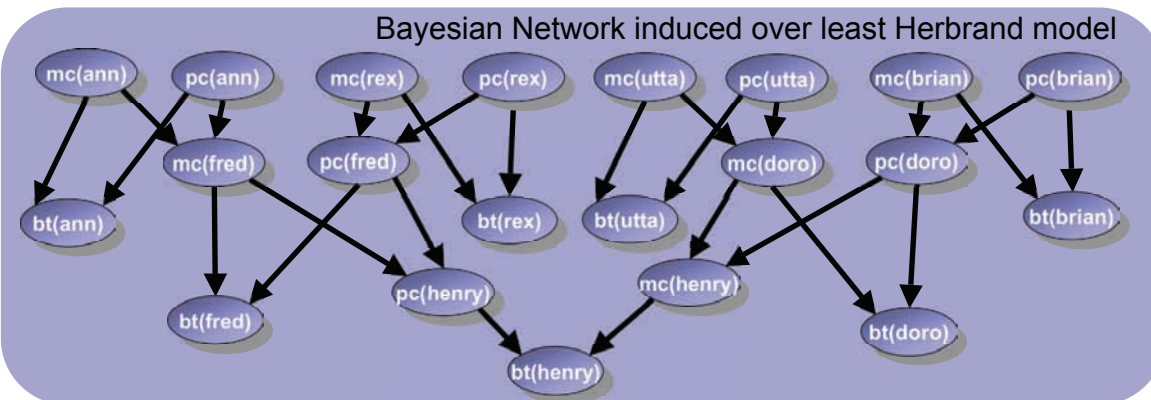


# Bayesian Logic Programs (BLPs)

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

$mc(Person) \mid mother(Mother, Person), pc(Mother), mc(Mother).$   
 $pc(Person) \mid father(Father, Person), pc(Father), mc(Father).$   
 $bt(Person) \mid pc(Person), mc(Person).$

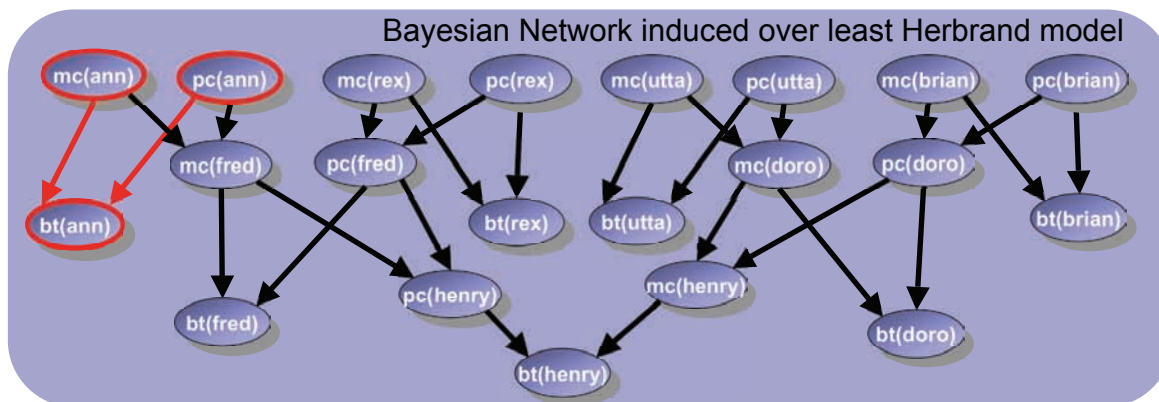
Bayesian Network induced over least Herbrand model





# Answering Queries

$P(bt(ann) ) ?$



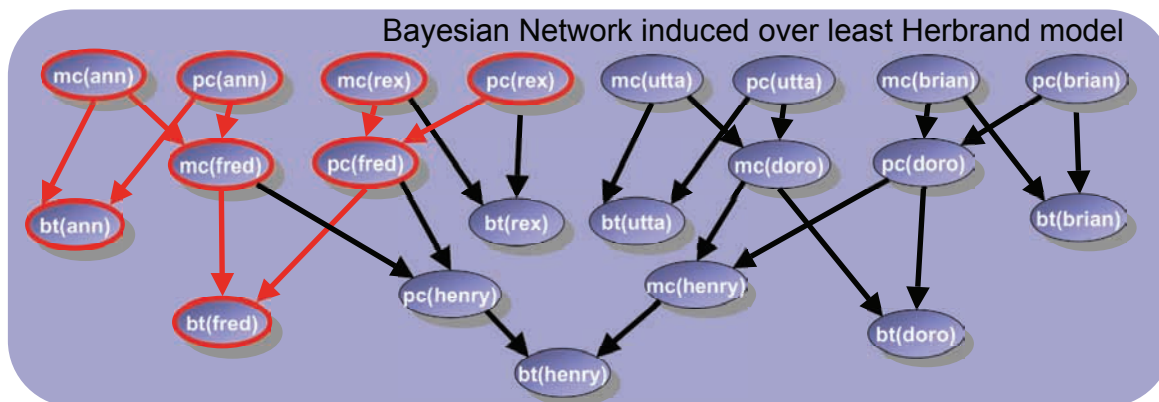
# Answering Queries

Bayes' rule

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

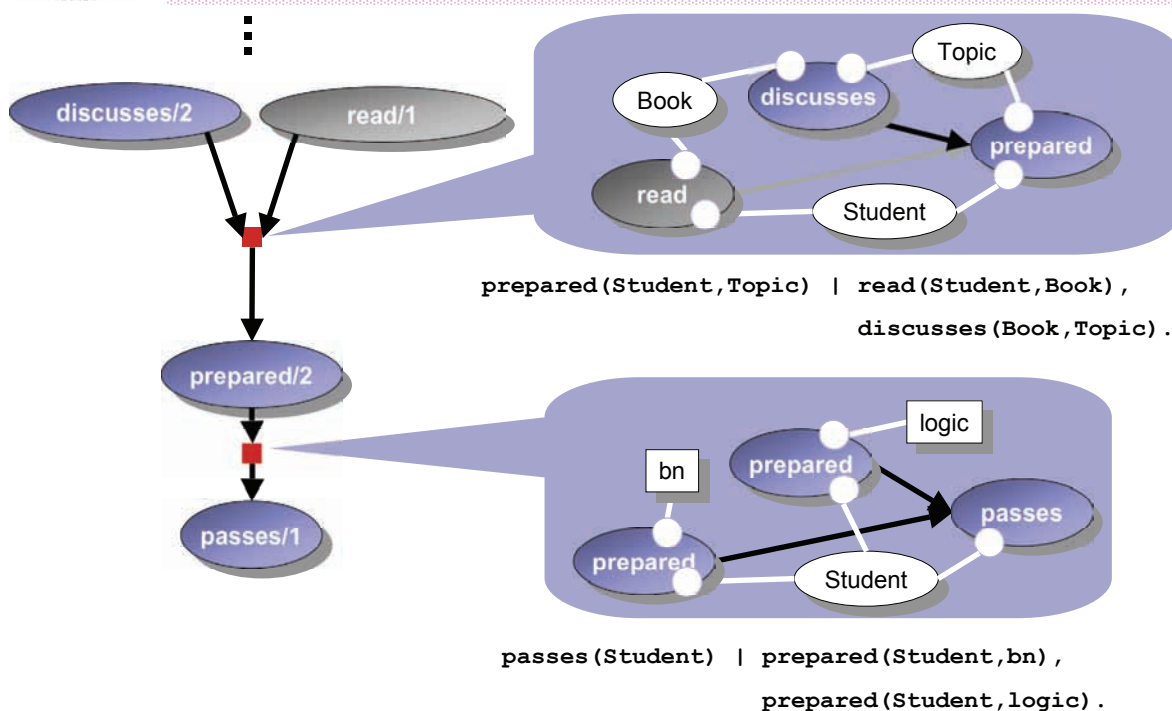
$$\frac{P(bt(ann), bt(fred))}{P(bt(fred))}$$

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

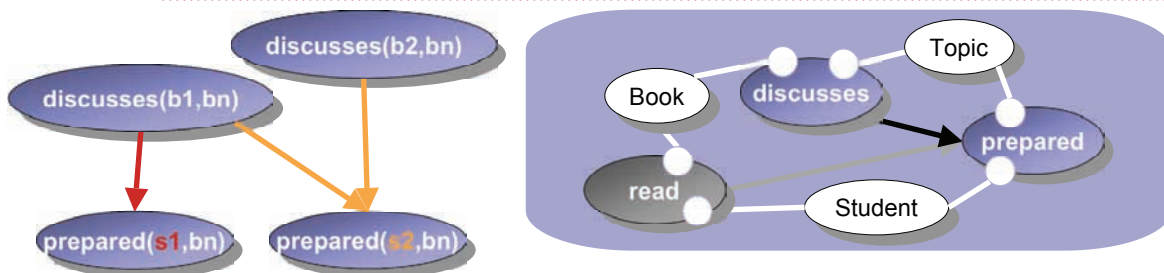




## Combining Partial Knowledge



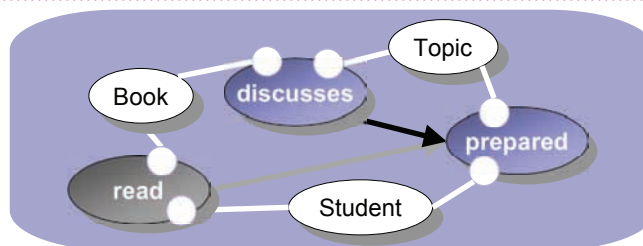
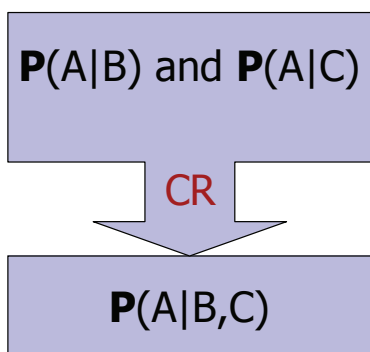
## Combining Partial Knowledge



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



## Combining Rules

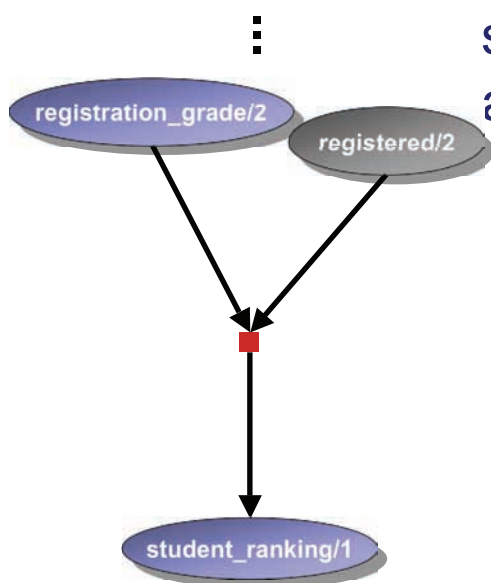


- Any algorithm which
  - has an empty output if and only if the input is empty
  - combines a set of CPDs into a single (combined) CPD
- E.g. noisy-or, regression, ...



## Aggregates

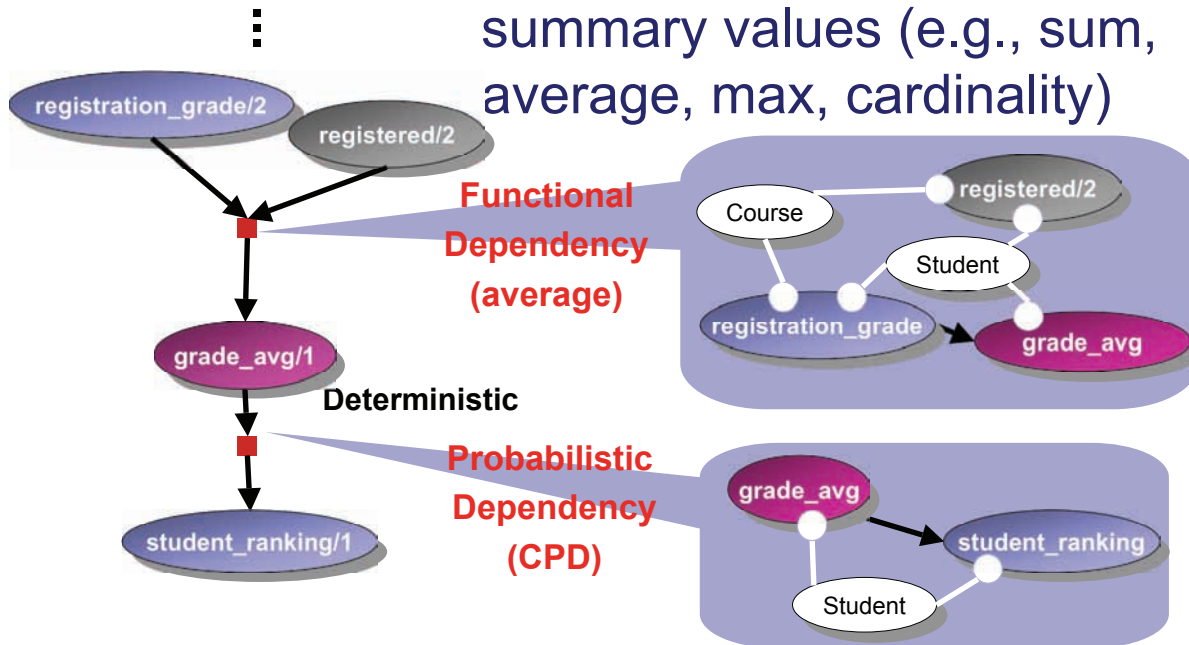
Map multisets of values to summary values (e.g., sum, average, max, cardinality)





## Aggregates

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

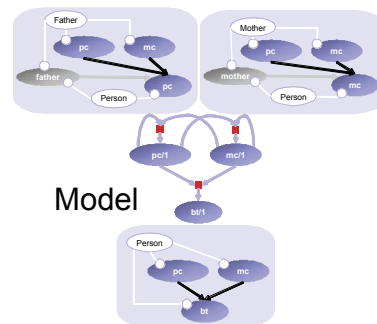
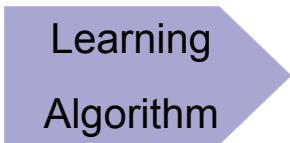
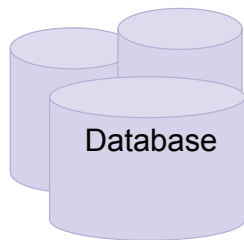


## Bayesian Logic Programs (BLPs)

- 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 ...
- Subsume PRMs
- Learning



# Learning Tasks



- Parameter Estimation
  - Numerical Optimization Problem
- Model Selection
  - Combinatorial Search



# What is the data about?

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(3)

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

### Family(2)

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



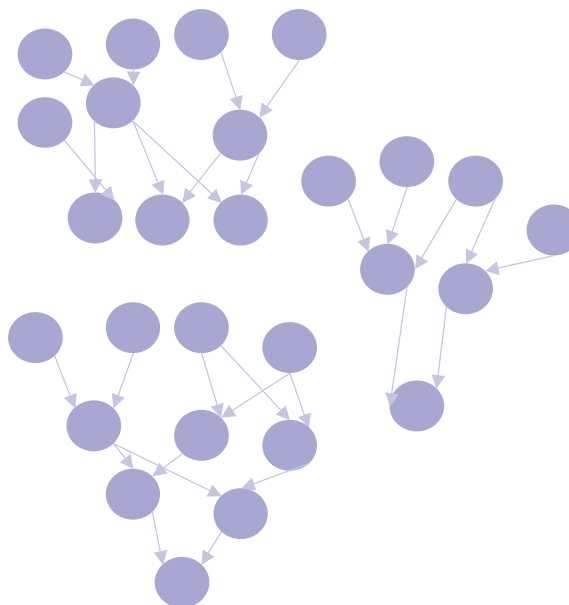
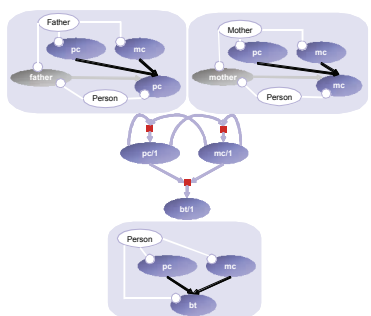


# Parameter Estimation

**Model(1)**  
 m(ann,dorothy),  
 pc(brian)=b,  
 bt(ann)=a,  
 bt(cecily)=ab,  
 bt(henry)=a,  
 bt(fred)=?,  
 bt(kim)=a,  
 bt(bob)=b

**Model(2)**  
 f(brian,dorothy),  
 f(fred),  
 f(fred),  
 f(job),  
 pc(rex)=b,  
 bt(doro)=a,  
 bt(brian)=?

**Background**

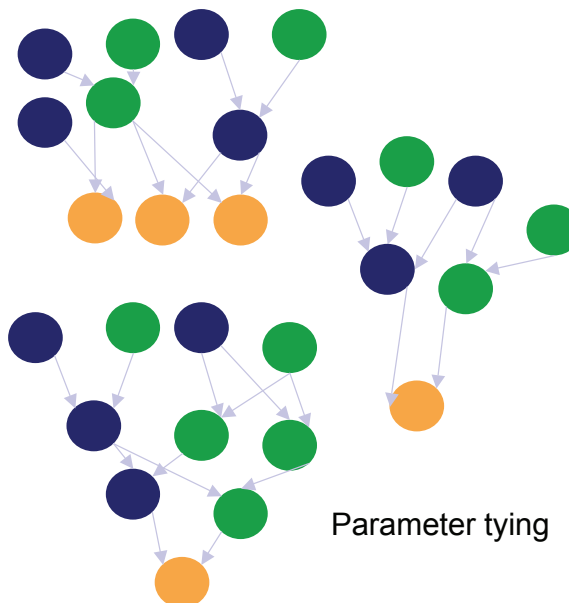
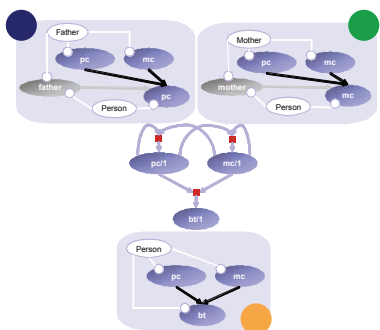


# Parameter Estimation

**Model(1)**  
 m(ann,dorothy),  
 pc(brian)=b,  
 bt(ann)=a,  
 bt(cecily)=ab,  
 bt(henry)=a,  
 bt(fred)=?,  
 bt(kim)=a,  
 bt(bob)=b

**Model(2)**  
 f(brian,dorothy),  
 f(fred),  
 f(fred),  
 f(job),  
 pc(rex)=b,  
 bt(doro)=a,  
 bt(brian)=?

**Background**

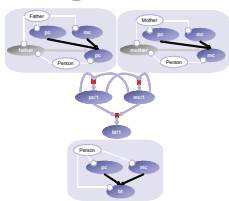




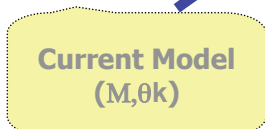
# Expectation Maximization

EM-algorithm:  
iterate until convergence

Logic Program L



Initial Parameters  $\theta_0$



Expectation

Inference

Expected counts of a clause

$$\sum_{GI} \sum_{DC} P(\text{head}(GI), \text{body}(GI) \mid DC)$$

Maximization  
Update parameters  
(ML, MAP)

$$\sum_{GI} \sum_{DC} P(\text{head}(GI), \text{body}(GI) \mid DC)$$

$$\sum_{GI} \sum_{DC} P(\text{body}(GI) \mid DC)$$

Model(1)	Background
pc(brian)=b.	m(ann,dorothy),
bt(ann)=a.	f(brian,dorothy),
bt(cecily)=ab.	ily,fred),
bt(henry)=a.	,fred),
bt(fred)=?	bob),
bt(kim)=a.	pc(alex)=b.
bt(bob)=b	bt(doro)=a.
	bt(brian)=?



# Model Selection

- Combination of ILP and BN learning
- Combinatorial search for hypo  $M^*$  s.t.
  - $M^*$  logically covers the data  $D$
  - $M^*$  is optimal w.r.t. some scoring function score, i.e.,  $M^* = \text{argmax}_M \text{score}(M,D)$ .
- Again, the prob. ILP approach *highlights*
  - Refinement operators
  - Background knowledge
  - Language bias
  - Search bias



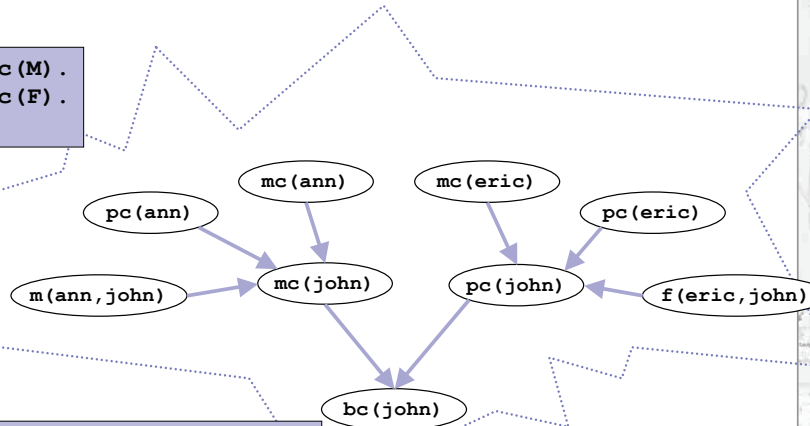
# Example

## 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).

```



## Data cases

```

{m(ann, john)=true, pc(ann)=a, mc(ann)=?,
 f(eric, john)=true, pc(eric)=b, mc(eric)=a,
 mc(john)=ab, pc(john)=a, bt(john) = ? }
...

```



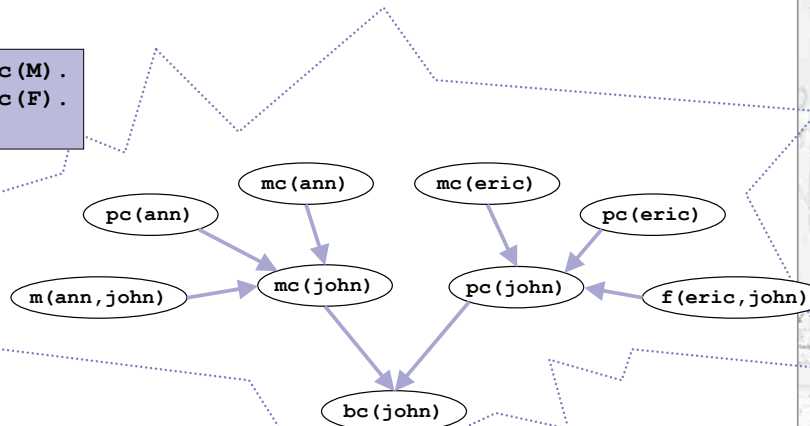
# Example

## 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).

```



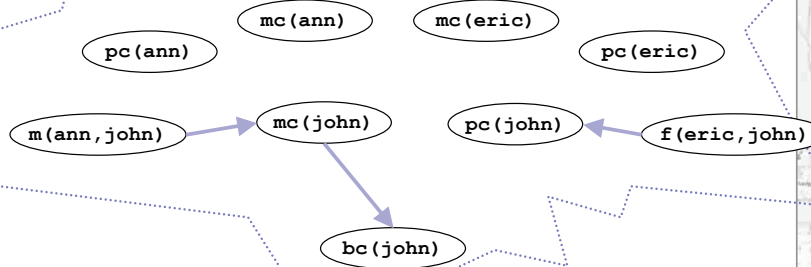
## Example

### 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).
```



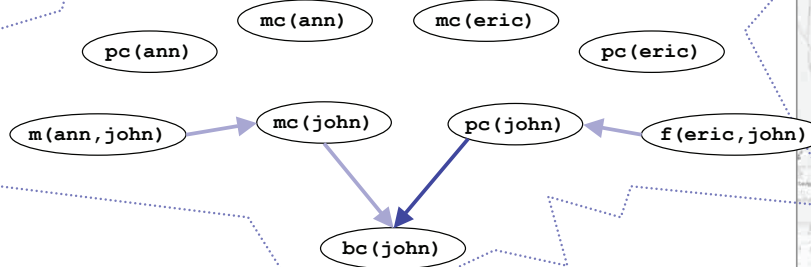
## Example

### 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).
```



# Example

## 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).

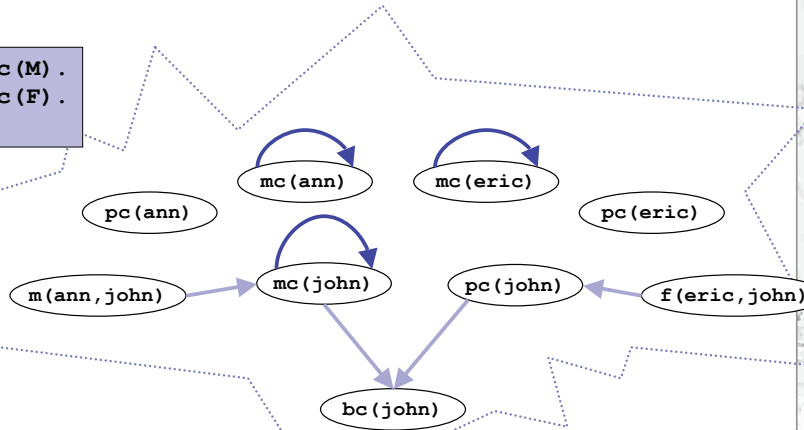
```

## Refinement

```


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


```



# Example

## 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).

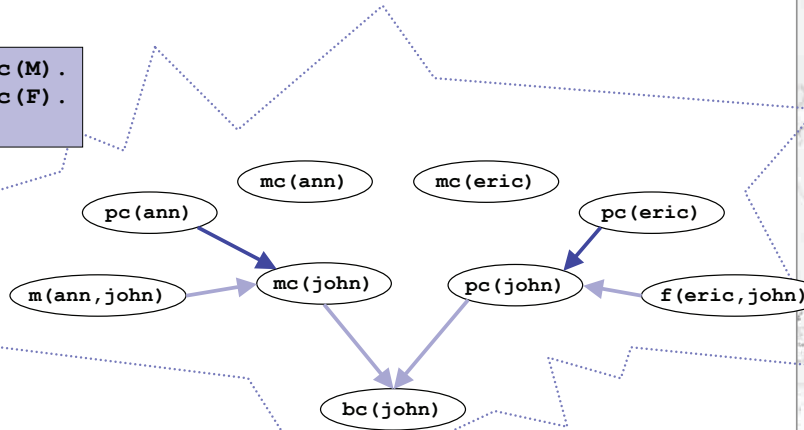
```

## Refinement

```

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

```





## Example

### 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

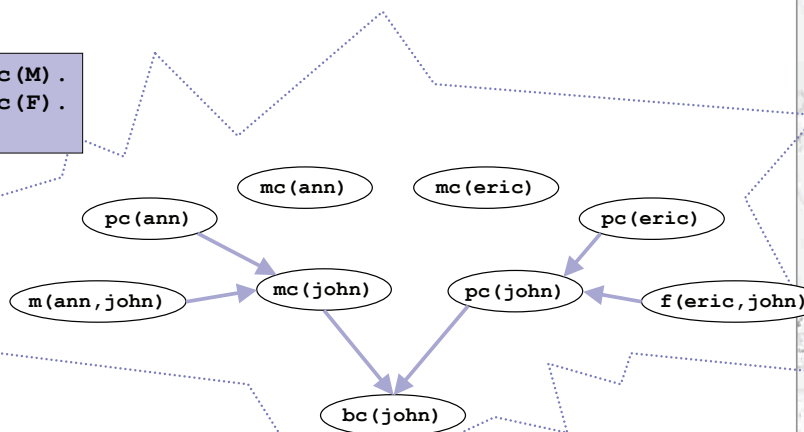
```
mc(X) | m(M,X).  
pc(X) | f(F,X).  
bt(X) | mc(X).
```

### 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).
```



## Undirected Probabilistic Relational Models

- So far, **directed** graphical models only
- Impose **acyclicity constraint**
- **Undirected** graphical models do not impose the acyclicity constraint

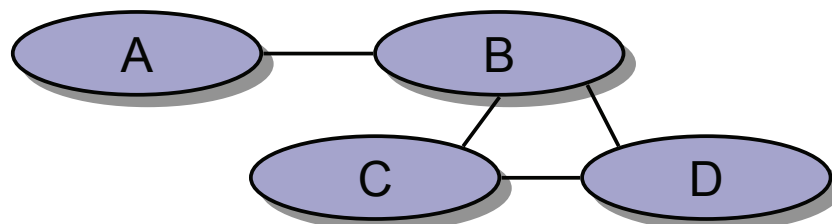


## Undirected Probabilistic Relational Models

- 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



## Undirected Graphical Models/Markov Networks



- To each clique  $c$ , a potential  $\phi_c$  is associated
- Given the values  $\mathbf{v}$  of all nodes in the Markov Network

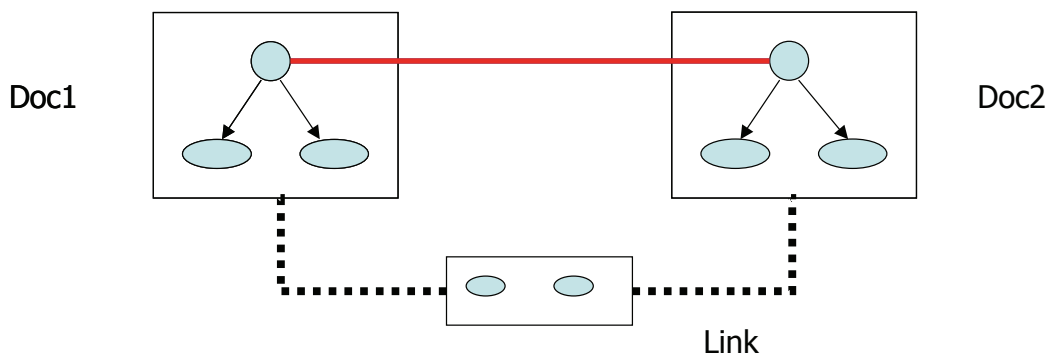
$$P(\mathbf{v}) = \frac{1}{Z} \prod_{c \in \mathcal{C}(G)} \phi_c(\mathbf{v}_c) \quad Z = \sum_{\mathbf{v}} \prod_{c \in \mathcal{C}(G)} \phi_c(\mathbf{v}_c)$$

$$\log P(\mathbf{v}) = \sum_c \mathbf{w}_c \cdot \mathbf{f}_c(\mathbf{v}_c) - \log Z = \mathbf{w} \cdot \mathbf{f}(\mathbf{v}) - \log Z$$



## Relational Markov Networks

```
SELECT doc1.Category,doc2.Category
FROM doc1,doc2,Link link
WHERE link.From=doc1.key and link.To=doc2.key
```



16th ECML, 9th PKDD, Porto, Portugal, Oct. 7th, 2005



## Markov Logic Networks

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)



slides by Pedro Domingos

16th ECML, 9th PKDD, Porto, Portugal, Oct. 7th, 2005



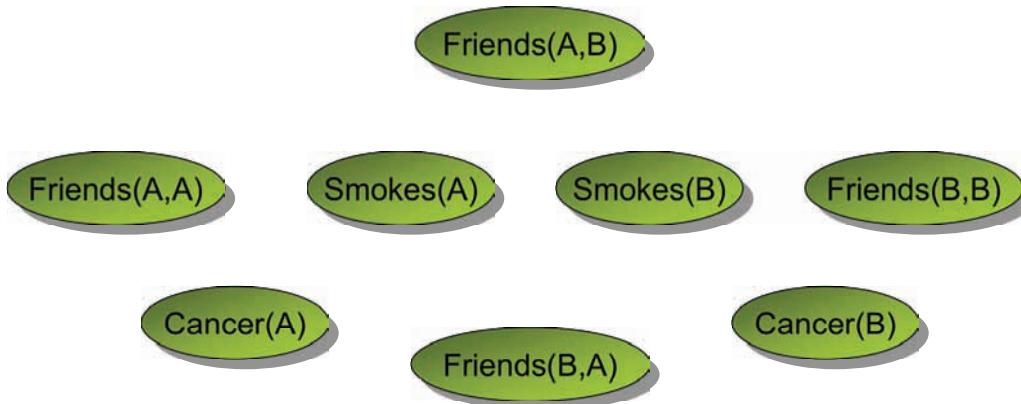


# Markov Logic Networks

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

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

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



slides by Pedro Domingos

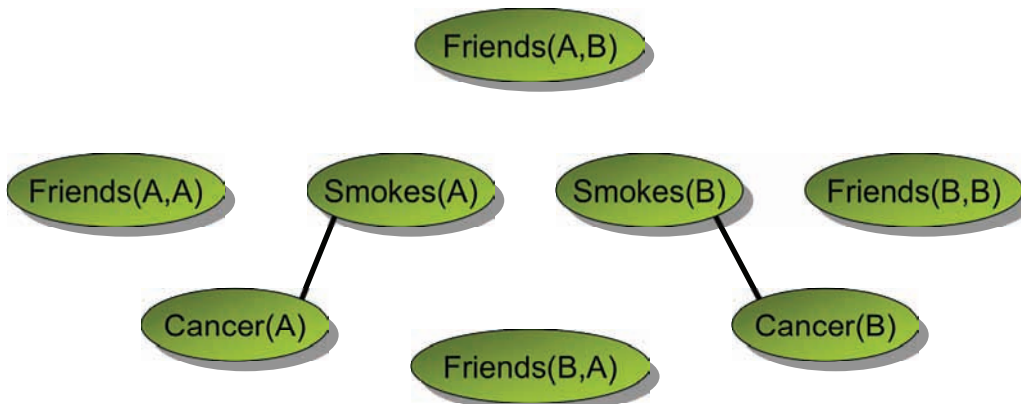


# Markov Logic Networks

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

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

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



slides by Pedro Domingos

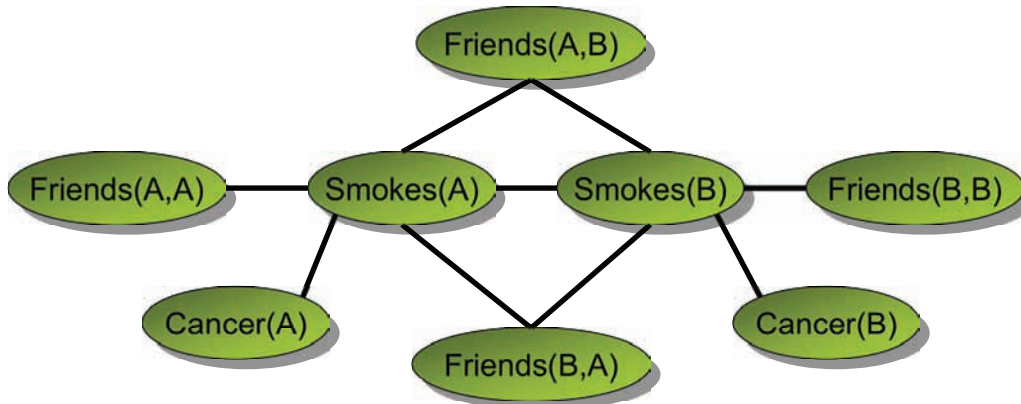


## Markov Logic Networks

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)



slides by Pedro Domingos



## Learning Undirected Probabilistic Relational Models

- Parameter estimation
  - discriminative (gradient, max-margin)
  - generative setting using pseudo-likelihood
- Structure learning
  - Similar to (Probabilistic) Inductive Logic Programming



## Applications

- Computer Vision (Taskar et al.)
  - Collective classification of 3D scan data
- Citation Analysis (Taskar et al., Singla&Domingos)
- Activity Recognition (Liao et al.)



## Conclusions Learning from Interpretations

- Data cases are herbrand interpretations
- Most prob. ILP approaches incorporates objects and relations among the objects into Bayesian and Markov networks
- Learning includes principles from
  - Inductive logic programming / multi-relational data mining
    - Refinement operators, Background knowledge, Bias
  - Statistical learning
    - Likelihood, Independencies, Priors



## Outline

1. Motivation / Introduction
2. Inductive Logic Programming (ILP)
  - Logic
  - Learning setting, cover relation
  - Learning from entailment, interpretations, and traces/proofs
3. Probabilistic ILP
  - Learning setting, probabilistic cover relation
4. Probabilistic Learning from
  - Interpretations, entailment, and traces/proofs
5. Discriminative ILP
6. Conclusions



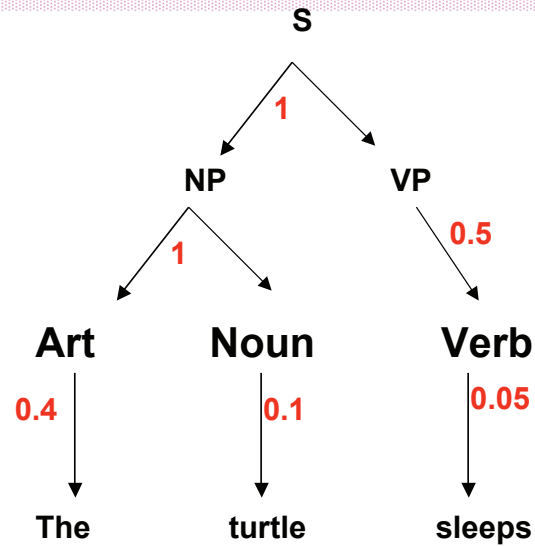
## Learning from entailment and from proofs

- Stochastic Logic Programs
  - Derived from Probabilistic Context Free Grammars by Eisele and Muggleton
  - Closely related to Sato's PRISM and Poole's ICL
- Learning from entailment
  - Parameter estimation (Cussens' FAM)
- Learning from proofs
  - Structure learning



# 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 label of rule  $i$

and  $c_i$  the number of times it was applied

$$P(\text{sentence}) = \sum_{i \text{ is a parse tree for sentence}} P(\text{tree}_i)$$

Observe: **all derivation/rewriting steps succeed**

i.e. S->T,Q

T->R,U

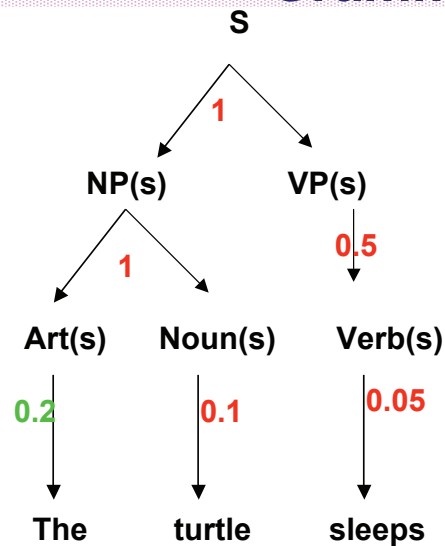
always gives

S-> R,U,Q



# Probabilistic Definite Clause Grammar

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



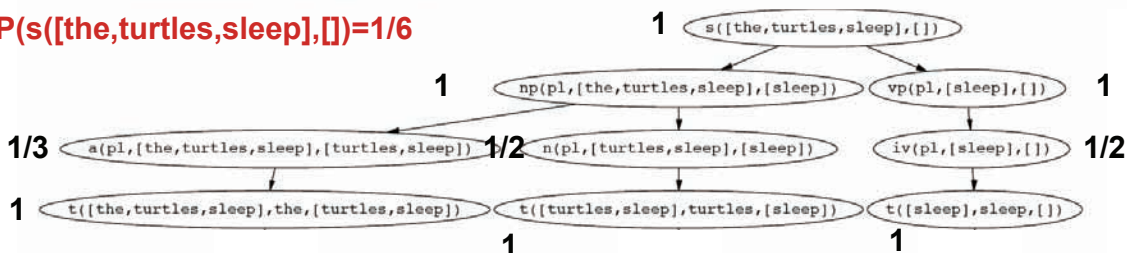
$P(\text{derivation tree}) = 1 \times 1 \times 0.5 \times 1 \times 0.2 \times 0.05$



## In SLP notation

- 1 sentence(A, B) :- noun\_phrase(C, A, D), verb\_phrase(C, D, B).
- 1 noun\_phrase(A, B, C) :- article(A, B, D), noun(A, D, C).
- verb\_phrase(A, B, C) :- intransitive\_verb(A, B, C).
- 1/3 article(singular, A, B) :- terminal(A, a, B).
- 1/3 article(singular, A, B) :- terminal(A, the, B).
- 1/3 article(plural, A, B) :- terminal(A, the, B).
- 1/2 noun(singular, A, B) :- terminal(A, turtle, B).
- 1/2 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).
- 1 terminal([A|B], A, B).

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





# Probabilistic Definite Clause Grammar

1.0 : S -> NP(Num), VP(Num)

1.0 NP(Num) -> Art(Num),  
Noun(Num)

0.6 Art(sing) -> a

0.2 Art(sing) -> the

0.2 Art(sing) -> the

0.1 Noun(sing) -> turtle

0.1 Noun(plur) -> turtles

...

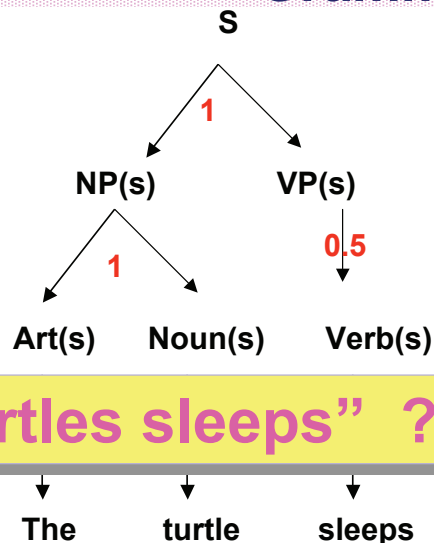
0.5 VP(Num) -> verb(Num)

0.5 VP(Num) -> Verb(Num), NP(Num)

0.05 Verb(sing) -> sleep

0.05 Verb(plur) -> sleeps

....



**What about "A turtles sleeps" ?**

$$P(\text{derivation tree}) = 1 \times 1 \times 0.5 \times 1 \times 0.2 \times 0.05$$



## SLPs

$$P_D(\text{derivation for goal } g(X_1, \dots, X_n)) = \prod_i p_i^{c_i}$$

Observe: **some derivations/resolution steps fail**

e.g. NP(Num) -> Art(Num), Noun(Num)

and Art(sing) -> a and Noun(plur) -> turtles

np(Num, S1, S2) :- art(Num, S1, S3), noun(Num, S3, S2)

and art(sing, [a|S], S) and noun(plur, [turtles|S], S)

Interest in **successful** derivations/proofs/refutations

-> **normalization necessary**

$$P_S(\text{proof}) = \frac{P_D(\text{proof})}{\sum_i P_D(\text{proof}_i)}$$

$$P_A(\text{ground atom } g(X_1, \dots, X_n)\theta) = \sum_{i \text{ is a proof tree for } g(X_1, \dots, X_n)\theta} P_S(\text{tree}_i)$$

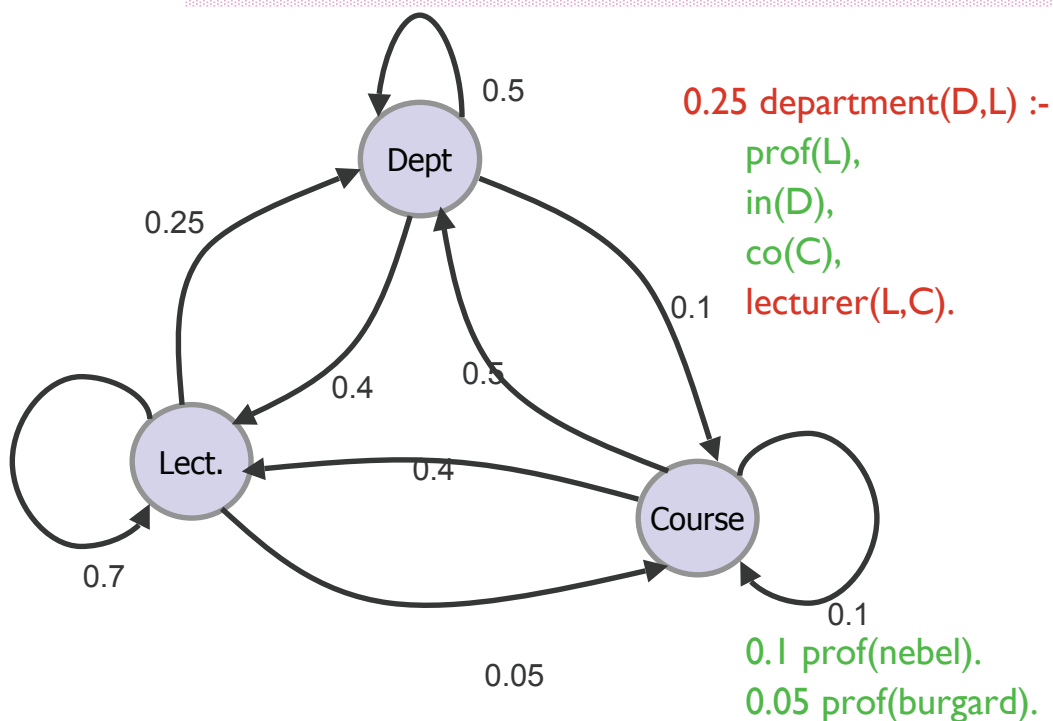


## Example Application

- Consider traversing a university website
- Pages are characterized by predicates  
`department(cs,nebel)` denotes the page of cs following the link to nebel
- Rules applied would be of the form  
`department(cs,nebel) :-`  
    `prof(nebel), in(cs), co(ai), lecturer(nebel,ai).`  
`pagetype1(t1,t2) :-`  
    `type1(t1), type2(t2), type3(t3), pagetype2(t2,t3)`
- SLP models probabilities over traces / proofs / web logs  
`department(cs,nebel), lecturer(nebel,ai007),`  
`course(ai007,burgard), ...`
- This is actually a Logical Markov Model
  - Logical **Hidden** Markov Model (cf. Kersting et al. JAIR)
  - Includes also structured observations and abstraction



## Logical Markov Model







## PRISM (Sato) / ICL (Poole)

- A logic program in which probability labels are attached to facts;
- Clauses carry no probability label (or equiv.  $P = 1$ )  
 $disjoint(h_1 : p_1; \dots, h_n : p_n)$   
statements, facts  $h_i; \sum_i p_i = 1$ 
  - Disjoint(head(C) : 0.5; tail(C): 0.5)
- Probability distributions can be defined in a related/similar fashion on proofs
  - on explanations
  - on atoms
  - though some differences with SLP



## PRISM (Sato) / ICL (Poole)

### Logical Abduction

mortal(X) :- human(X) and mortal(X) :- animal(X)

From mortal(socrates) infer human(socrates)

OR

infer animal(socrates)

### Probabilistic Abduction (ICL Poole)

From mortal(socrates) infer

probabilities for human(socrates)

and animal(socrates)

### Backward reasoning



## 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 + prob switch*

Infer e.g.  $P(\text{btype}('A'))$  from probabilities of proofs - facts

Infer e.g. explanations and prob.

$$P(\text{gene}(\text{father},a), \text{gene}(\text{mother},a)) = p_2 * p_2$$

...

$$P(\text{gene}(\text{father},a), \text{gene}(\text{mother},a) \mid \text{btype}('A')) : \text{normalize}$$



## Parameter Estimation for CFGs

- Given
  - A set of **sentences** for the startsymbol
    - E.g. **the, turtles, sleep**
  - The structure of the CFG
    - (i.e. the rules - not the probability labels)
- Find
  - The maximum likelihood parameters of the CFG
- Approach
  - Inside - Outside Algorithm
    - Parse trees are unobserved
    - Instance of EM
    - dynamic programming (using CYK)



## Parameter Estimation for SLPs

- Given
  - A set of **ground facts** for a predicate  $g$ 
    - E.g.  $s([the, turtles, sleep], [])$
  - The structure of the stochastic logic program
    - (i.e. the logic part - not the probability labels)
- Find
  - The maximum likelihood parameters of the SLP
- Approach
  - Failure Adjusted Maximisation Algorithm (Cussens)
    - Parse trees and **failures** are unobserved
    - Instance of EM
    - Ako dynamic programming using tabling (Sato)
    - Very efficient



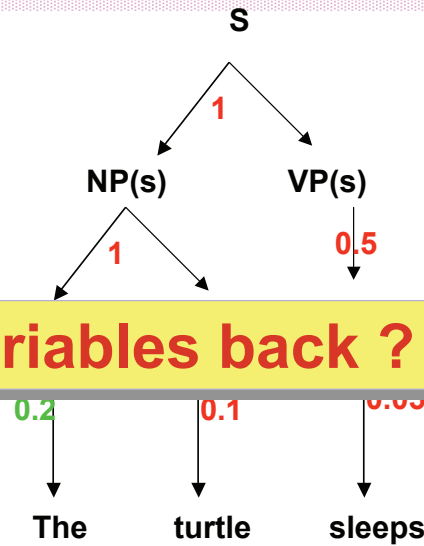
## Structure Learning

- From entailment : Muggleton ILP 02
  - Learns a single clause at a time from facts only
  - Hard problem, requires one to solve the full inductive logic programming problem
- From proof trees : De Raedt et al AAI 05
  - Learn from **proof-trees** instead of from ground facts
  - Proof-trees carry **much more** information
  - Upgrade idea of tree bank grammars
- Given
  - A set of proof trees
- Find
  - An SLP that maximizes the likelihood



## Initial Rule Set DCG

$S \rightarrow NP(s), VP(s)$   
 $NP(s) \rightarrow Art(s), Noun(s)$   
 $VP(s) \rightarrow Verb(s)$   
 $Art(s) \rightarrow the$   
 $Noun(s) \rightarrow turtle$   
 $Verb(s) \rightarrow sleeps$



$P(\text{derivation tree}) = 1 \times 1 \times 0.2 \times 0.1 \times 0.5 \times 0.05$



## Learning SLPs from Proof Trees

- Based on Tree-Bank Grammar idea, e.g. Penn Tree Bank
- **Key algorithm**
  - Let S be the set of all (instantiated) rules that occur in an example proof tree
  - Initialize parameters
  - repeat as long as the score of S improves
    - **Generalize** S
    - **Estimate** the parameters of S using Cussens' FAM
      - (which can be simplified - proofs are now observed)
  - Output S



## Generalizing Rules in SLPs

- Generalization in ILP
  - Take two clauses for same predicate and replace them by the lgg under  $\theta$  - subsumption (Plotkin)
  - Example
    - 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 r1 and r2 by the lgg **should preserve the proofs !**
- So, two rules r1 and r2 should only be generalized when
  - There is a one to one mapping (with corresponding substitutions) between literals in r1, r2 and lgg(r1,r2)
- Exclude
  - father(j,a) :- m(j),f(a),parent(j,a)
  - father(j,t) :- m(j),m(t), parent(j,t)
- Gives
  - father(j,P) :- m(j),m(X),parent(j,P)



[Anderson et al.]

## Web Log Data

- Log data of web sides
- KDDCup 200 ([www.gazelle.com](http://www.gazelle.com))
- RMM over
  - 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)

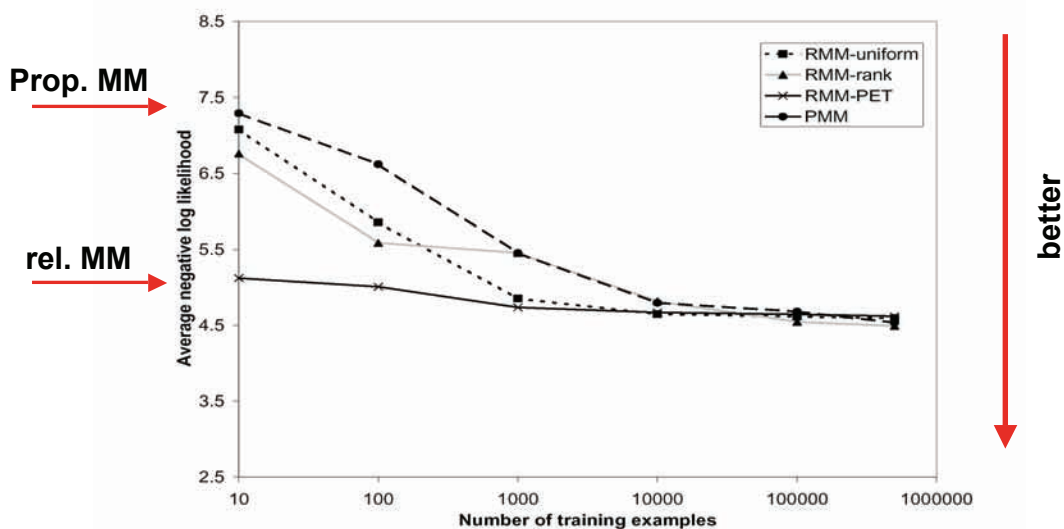
16th ECML, 9th PKDD, Porto, Portugal, Oct. 7th, 2005



[Anderson et al.]

## User Log Data

### Relational representation pays off



16th ECML, 9th PKDD, Porto, Portugal, Oct. 7th, 2005



## Protein Fold Recognition

[Kersting et al.; Kersting, Gaertner]

- Comparison of protein structure is fundamental to biology, e.g. function prediction
- Two proteins show sufficient sequence similarity = essentially adopt the same structure.

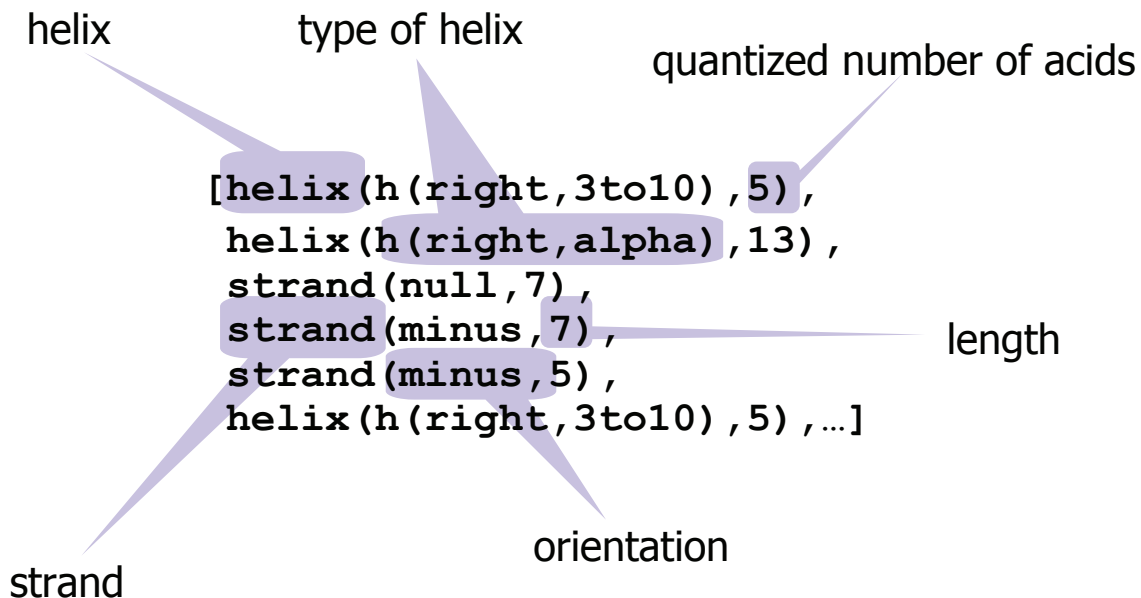


- If one of the two similar proteins has a known structure, can build a rough model of the protein of unknown structure.



## Protein Secondary Structure

[Kersting et al.; Kersting, Gaertner]

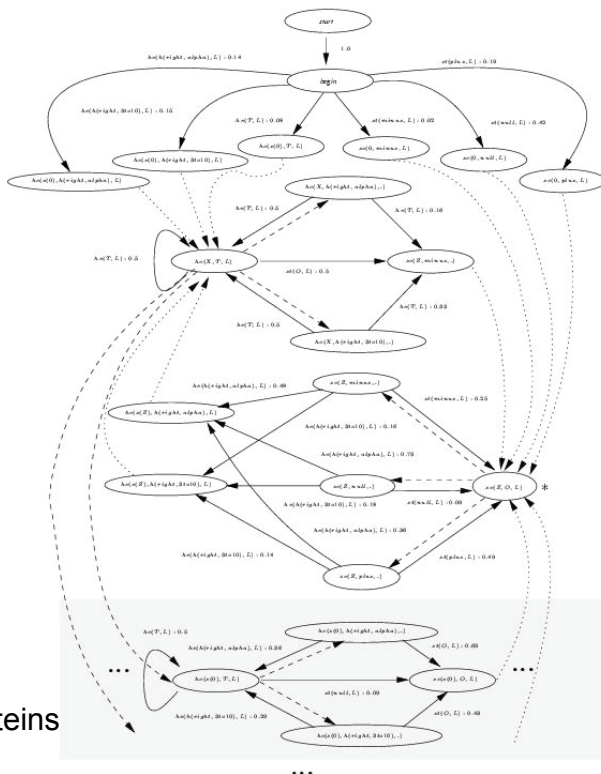




# Model

[Kersting et al.]

~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)



# Conclusions 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
- Learning from proof trees is a new setting for inductive logic programming/statistical relational learning
  - Generalizes learning from traces
- Strong logical constraints at structure level
- Allows one also to elegantly model and study RMMs and LOHMMs
  - Sequential relational / logical traces.
- A lot of further research questions
  - Most of all : experiments on real data





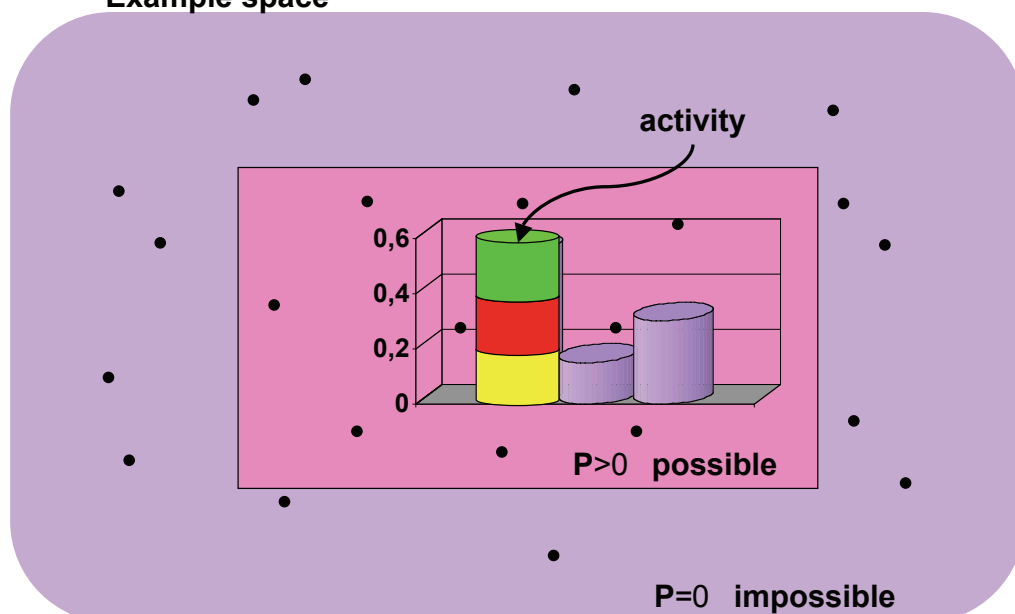
## Outline

1. Motivation / Introduction
2. Inductive Logic Programming (ILP)
  - Logic
  - Learning setting, cover relation
  - Learning from entailment, interpretations, and traces/proofs
3. Probabilistic ILP
  - Learning setting, probabilistic cover relation
4. Probabilistic Learning from
  - Interpretations, entailment, and traces/proofs
5. Discriminative ILP
6. Conclusions



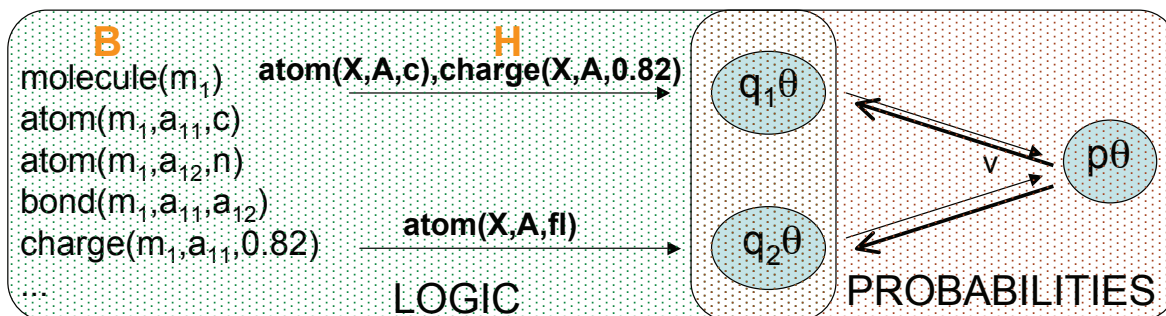
## Probabilistic ILP Problem

### Example space





## nFOIL = naive Bayes + FOIL



- Clause set + simple probabilistic model
- Idea: Clauses are independent
- Success/failure of a query is random variable in a Naive Bayes model



## The nFOIL model

- Naive Bayes assumption translates into

$$\begin{aligned} P(p\theta \mid H, B) &= P(p\theta \mid q_1\theta, \dots, q_k\theta) \\ &= \frac{P(q_1\theta, \dots, q_k\theta \mid p\theta) * P(p\theta)}{P(q_1\theta, \dots, q_k\theta)} \\ &= \frac{\prod_i P(q_i\theta \mid p\theta) * P(p\theta)}{P(q_1\theta, \dots, q_k\theta)} \end{aligned}$$

- Model consists of clauses  $q_1, \dots, q_k$  and parameters  $P(q_i\theta \mid p\theta), P(p\theta)$
- Classify positive if  $P(p\theta \mid H, B) > 0.5$



## Learning: nFOIL

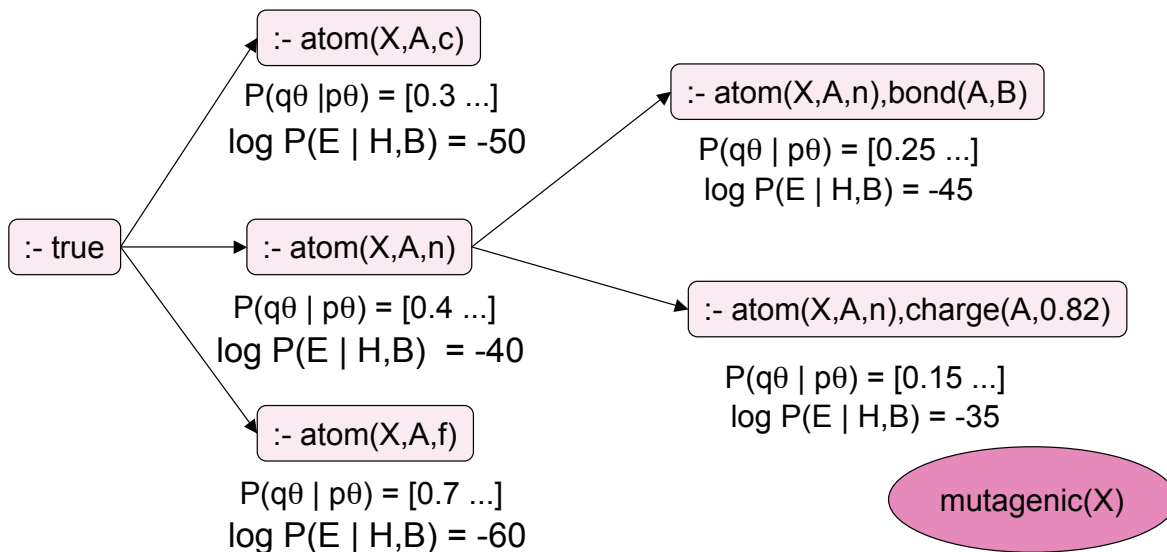
- Modified FOIL: search guided by cond. likelihood
- FOIL
  - a clause is scored by its coverage
  - Covered positive examples are removed
- nFOIL
  - score a **set of clauses**  $\{q_1, \dots, q_k\}$  by conditional likelihood:

$$P(E | H, B) = \prod_{e \in E} \frac{\prod_i P(q_i | p_\theta) * P(p_\theta)}{P(q_1, \theta, \dots, q_k | \theta)}$$

where  $P(q_i | p_\theta) = \frac{\text{count}(q_i | p_\theta)}{\text{count}(p_\theta)}$

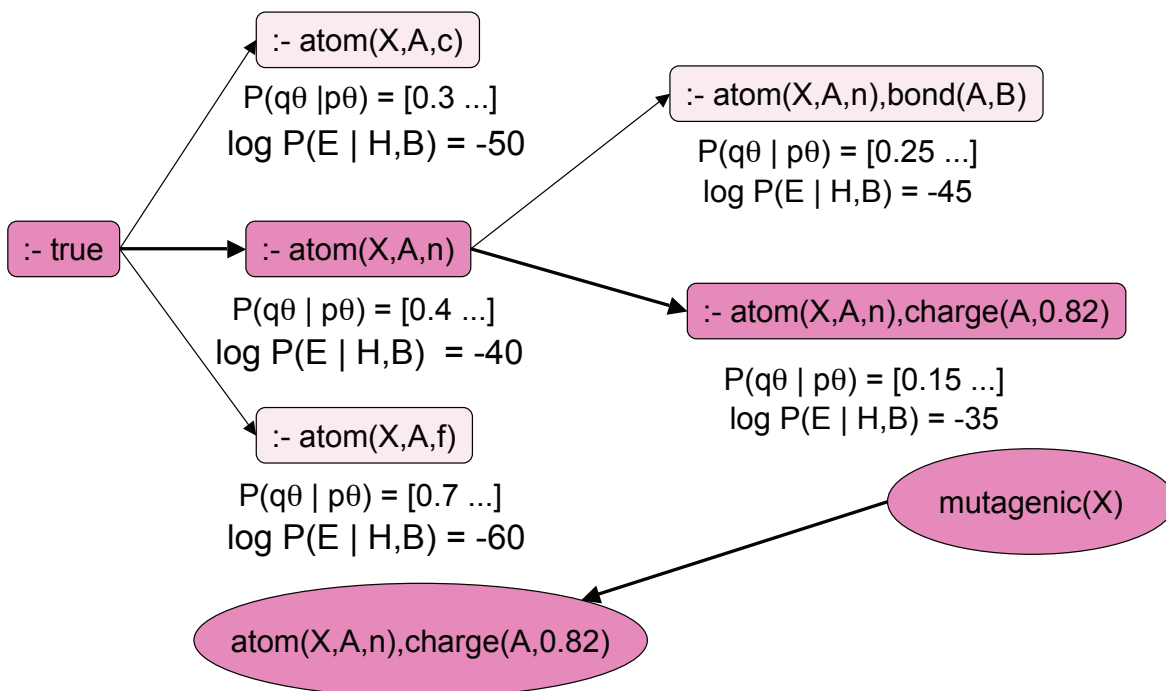


## Learning (Example)

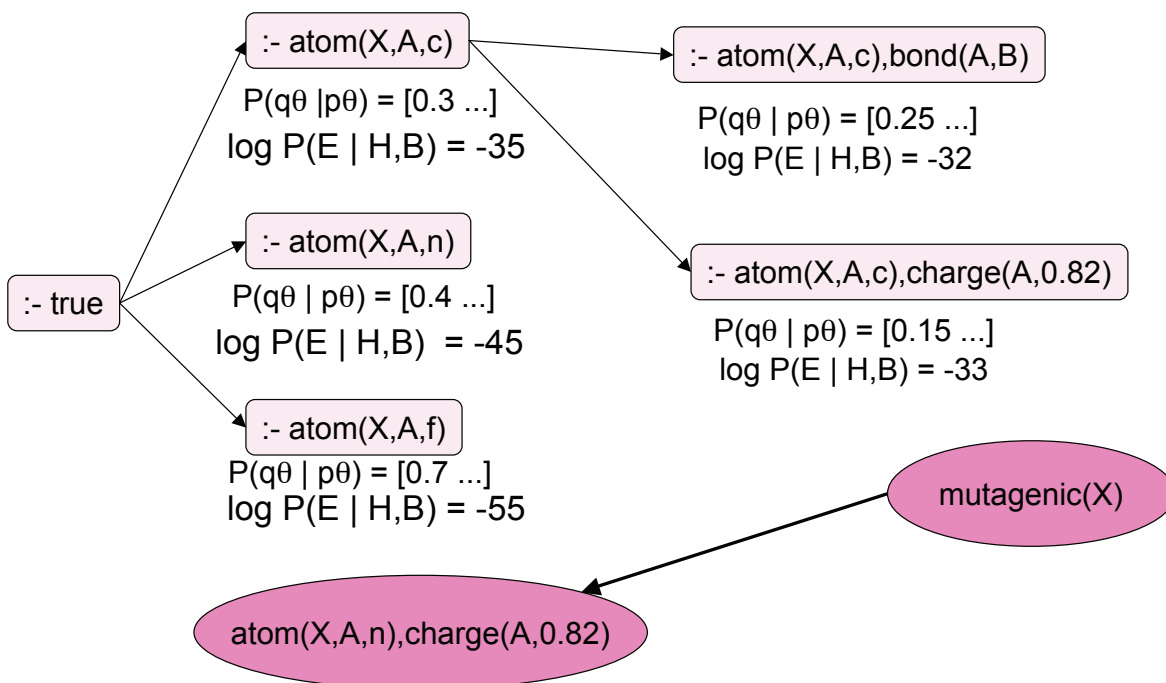




## Learning (Example)

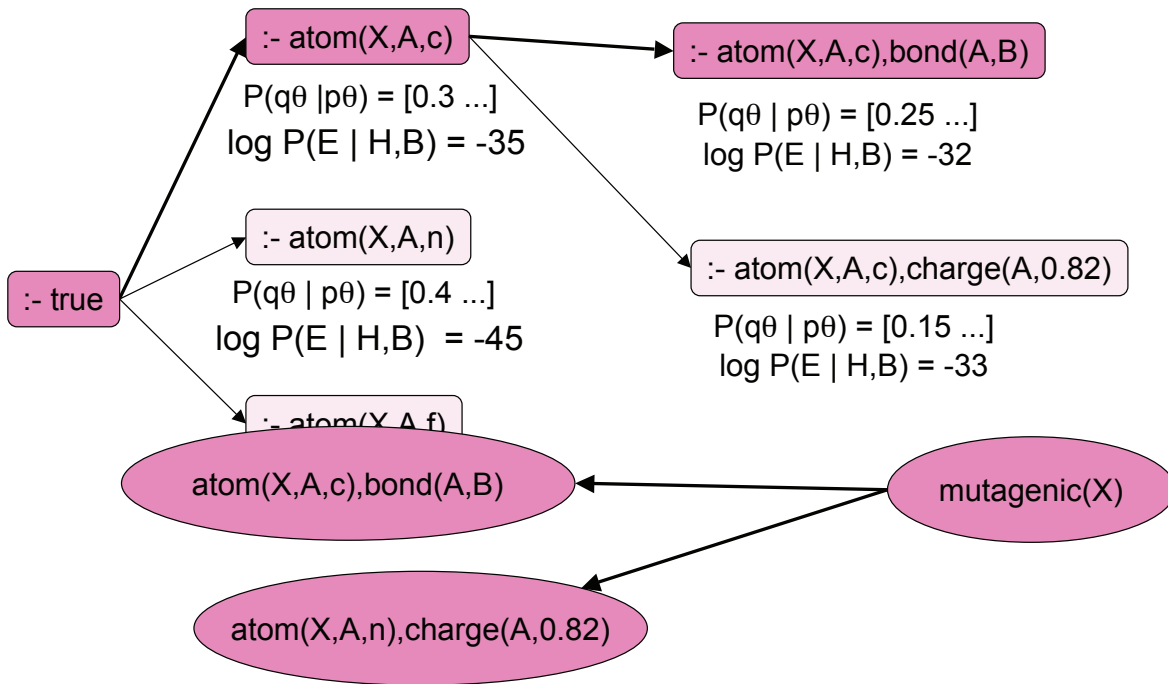


## Learning (Example)

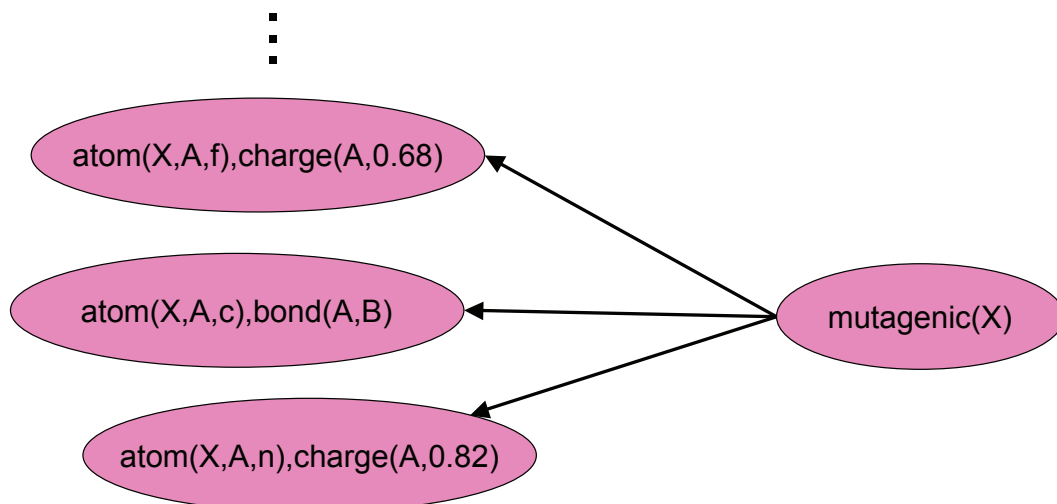




## Learning (Example)

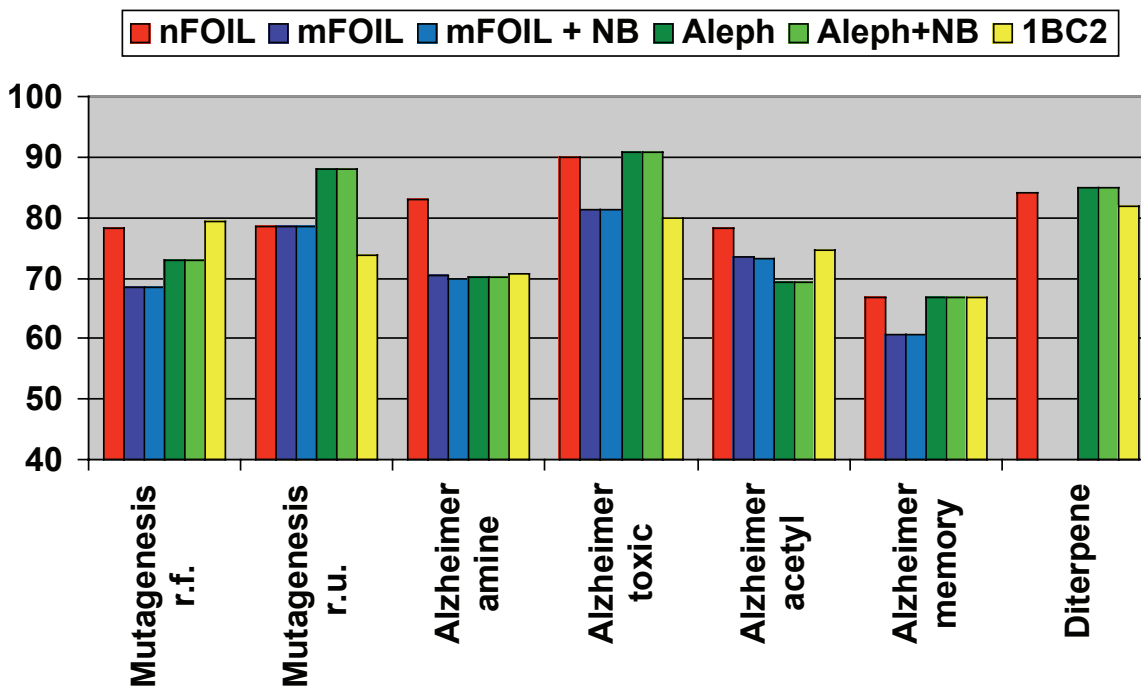


## Learning (Example)





## Experimental Results



## Fisher Kernels

[Jaakkola, Haussler NIPS'99]

- Kernels based on probability models
- Idea: capture the change in parameters to accommodate new data point

Fisher score of observed data  $x$

$$U_x = \nabla_{\theta} \log P(x | \theta^*, M) = \left( \frac{\partial \log P(x | \theta^*, M)}{\partial \theta_1}, \dots, \frac{\partial \log P(x | \theta^*, M)}{\partial \theta_n} \right)^T$$

Fisher information matrix  $J_{\theta} = E_x[U_x U_x^T]$

Fisher kernel  $k(x, x') = U_x^T J_{\theta^*}^{-1} U_{x'}$

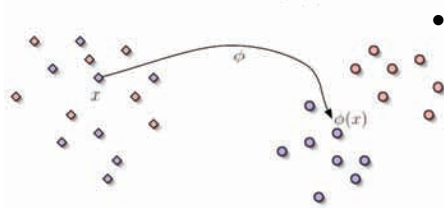
In practice:  $k(x, x') = U_x^T U_{x'}$



## SVMs in a Nutshell

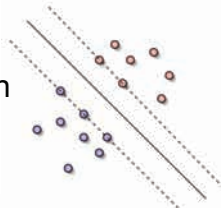
[Boser, Guyon, Vapnik COLT'92]

$$x \rightarrow \phi(x)$$

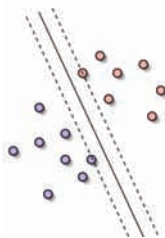


- With an appropriate mapping  $\phi$  to a sufficiently high dimension, data from 2 classes can always be separated by a hyperplane
- Find optimal separating hyperplane by maximizing the margin of separation
- Computationally effective algorithms exist to solve this optimization problem (in higher dimensional feature space).

Larger Margin



Smaller Margin



## Protein Fold Recognition

[Kersting et al.; Kersting, Gaertner]

- Comparison of protein structure is fundamental to biology, e.g. function prediction
- Two proteins show sufficient sequence similarity = essentially adopt the same structure.



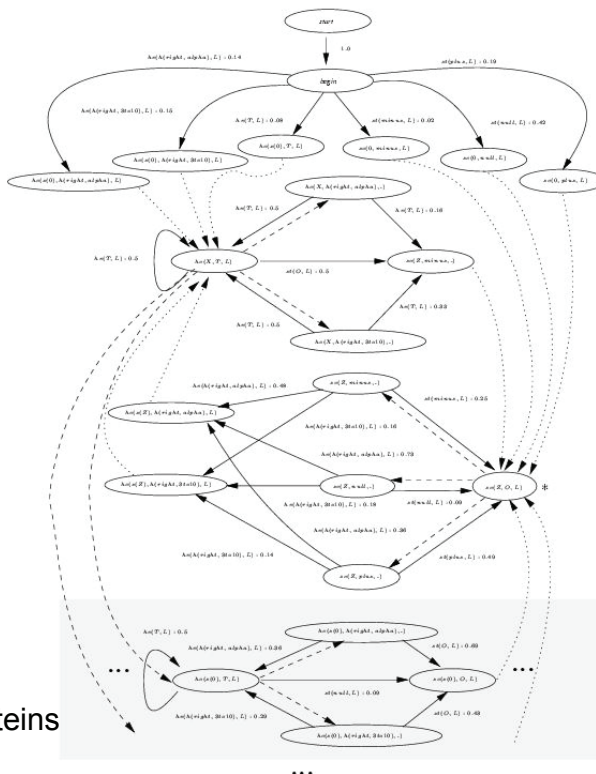
- If one of the two similar proteins has a known structure, can build a rough model of the protein of unknown structure.



# Model

[Kersting et al.]

~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]

- 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

## New class of probabilistic relational Kernels





[Kersting et al.; Kersting, Gaertner]

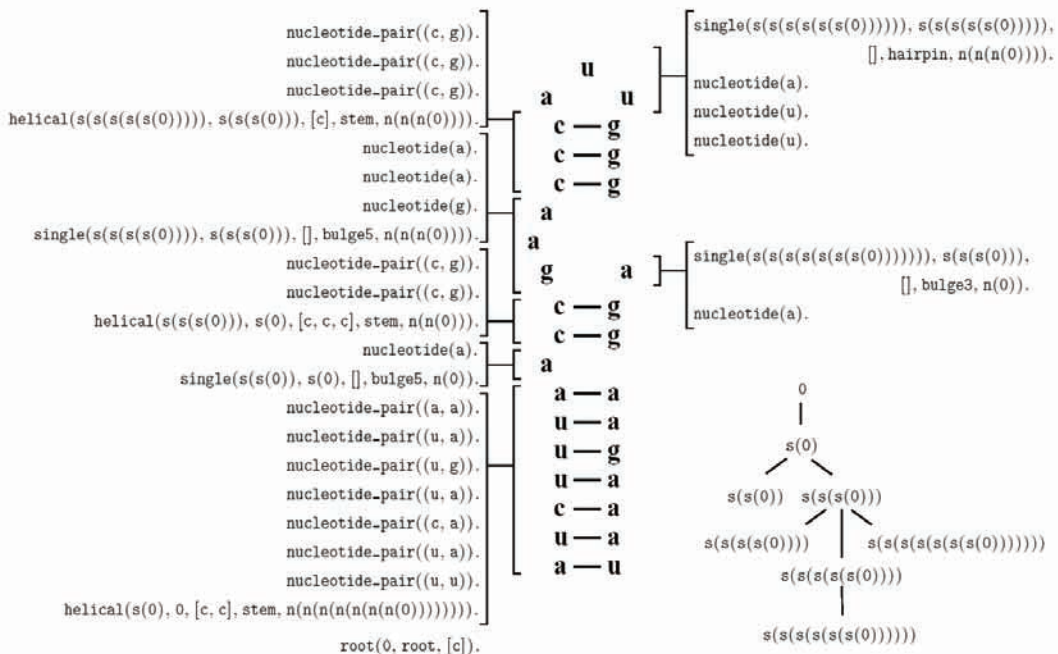
# mRNA

- Identifying subsequences in mRNA that are responsible for biological functions.
- Secondary structures of mRNAs form tree structures: not easily for HMMs



[Kersting et al.; Kersting, Gaertner]

# mRNA





[Kersting et al.; Kersting, Gaertner]

## mRNA

93 logical sequences (in total 3122 ground atoms)

- 15 and 5 SECIS (Selenocysteine Insertion Sequence),
- 27 IRE (Iron Responsive Element),
- 36 TAR (Trans Activating Region) and
- 10 histone stemloops.

Leave-one-out crossvalidation:

Plug-In Estimates: 4.3 % error

Fisher kernels SVM: 2.2 % error



## Outline

1. Motivation / Introduction
2. Inductive Logic Programming (ILP)
  - Logic
  - Learning setting, cover relation
  - Learning from entailment, interpretations, and traces/proofs
3. Probabilistic ILP
  - Learning setting, probabilistic cover relation
4. Probabilistic Learning from
  - Interpretations, entailment, and traces/proofs
5. Discriminative ILP
6. Conclusions



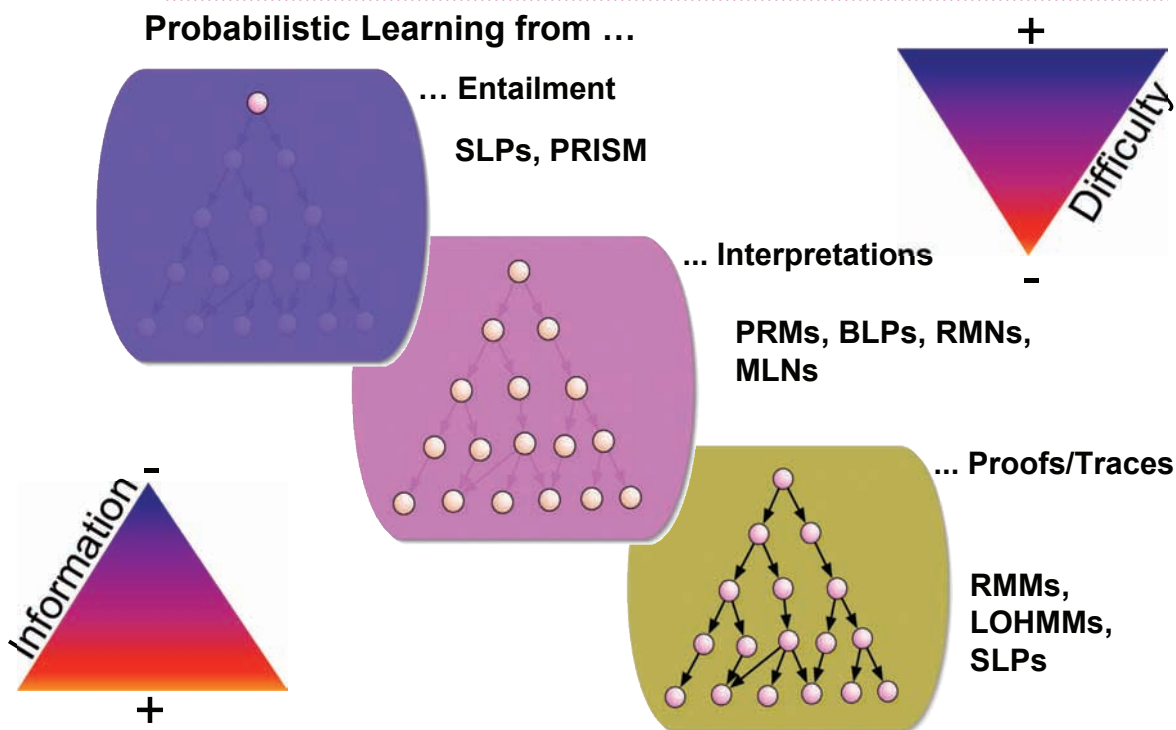
## Conclusions

- 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 settings (as for ILP)
  - **Learning from entailment**: Parameter Est. SLPs/Prism
  - **Learning from interpretations**: BLPs, PRMs, RMNs, MLNs
  - **Learning from traces or proofs**: SLPs, RMMs, LOHMMs
  - Discriminative ILP: nFOIL
- Different settings have different complexities



## Use of different Prob. ILP Settings

Probabilistic Learning from ...





## Conclusions

Many interesting problems left !

# Thank you for your attention!

## Please join PILP / SRL !

<http://www.aprill.org>

<http://www.informatik.uni-freiburg.de/~kersting/plmr>



## Special thanks to the APRIL II consortium

- „Application of Probabilistic ILP“
- 3 years EU project
- 5 institutes
- [www.aprill.org](http://www.aprill.org)

