

PhD in Computer Science and Engineering
Bologna, April 2016

Machine Learning

Marco Lippi

marco.lippi3@unibo.it

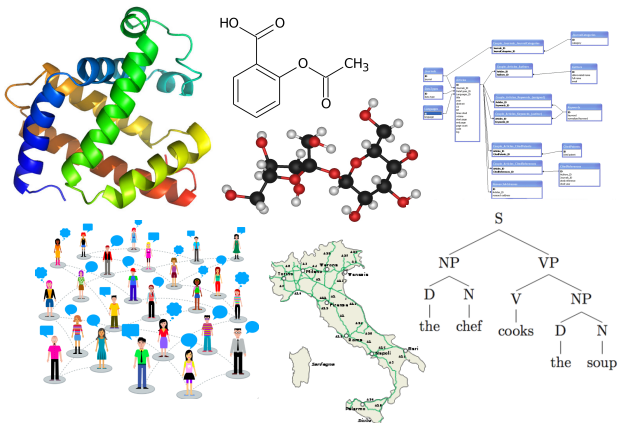


ALMA MATER STUDIORUM
UNIVERSITÀ DI BOLOGNA

Statistical Relational Learning

Structured and relational data

In real-world problems, machine learning has very often to deal with highly **structured** and **relational** data.



Why do we need SRL ?

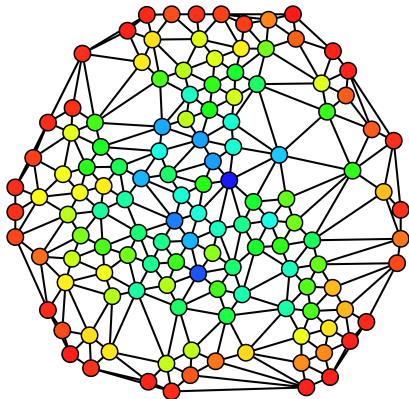
- Statistical learning assumes examples to be **independent and identically distributed** (i.i.d.)
- Traditional logic approaches (ILP) can handle relations but assume **no noise or uncertainty** in data

Many application domains need both **uncertainly handling** and **knowledge representation** !

General observations

- Improve **performance** on certain tasks 😊
- Improve **interpretability** of results 😊
- Maybe computationally more **expensive** 😞
- Learning could be much **harder** 😞

Collective classification



Node classification

- User profiles in a social network
- Gene functions in a regulatory network
- Congestions in a transportation network
- Service requests in p2p networks
- Fault diagnosis in sensor networks
- Hypertext categorization on the Internet
- ...

Node classification

- use attributes of each node
- use attributes of neighbor nodes
- use features coming from the graph structure
- **use labels of other nodes**

In social networks, two important phenomena can happen

- **homophily** → a link between individuals is correlated with those individuals being similar in nature
- **co-citation regularity** → similar individuals tend to be related/connected to the same things

Node classification

Iterative approaches:

- 1 label a few nodes, given existing labels and features
- 2 train using the newly labeled examples
- 3 repeat the procedure until all nodes are labeled

Random walk approaches:

- the probability that a node v receives label y is that of a random walk starting at v will end at a node labeled with y
- could be solved by **label propagation** or **graph regularization**

Edge classification

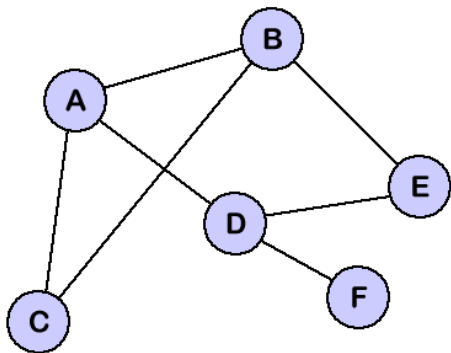
Classify properties of edges

- **viral marketing**: influence of a user on his neighbors
- **argumentation mining**: support/attack relationships
- **social networks**: types of interactions between users

In social networks there are two main theories:

- **balance** → the friend of my friend is my friend, the enemy of my friend is my enemy
- **status** → positive edge between v_i and v_j means v_j has a higher status than v_i

Link prediction (or Link mining)



Link prediction

- Friendship in a social network
- Recommendation in a customer-product network
- Interaction in a biological network
- Link congestion in a transportation network
- Link congestion in a p2p network
- ...

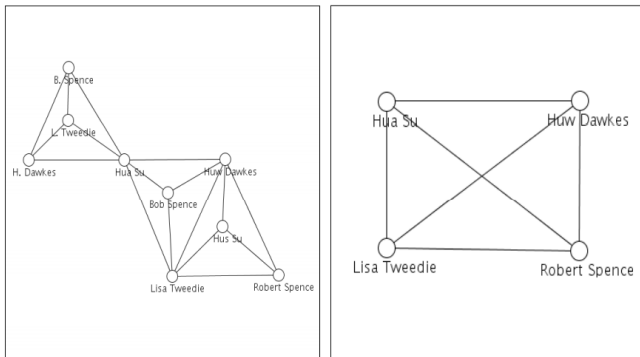
Networks are very often **dynamic**:

- nodes may change over time
- links may change over time
- node properties may change over time
- edge properties may change over time

Shall we predict the **evolution** of the network ?

E.g., given current links, which ones are likely to **change** ?

Entity resolution

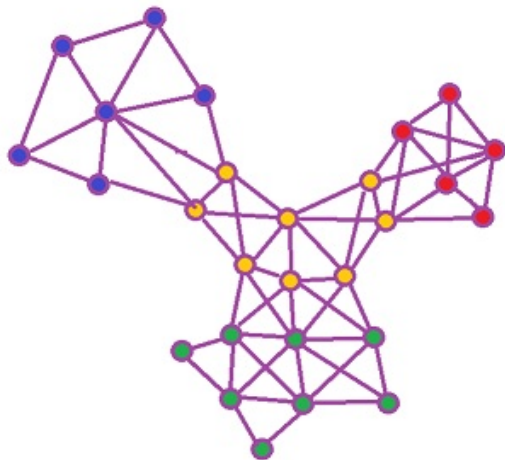


[Image from Entity Resolution Tutorial, VLDB2012, Lise Getoor]

Entity resolution

- Same paper in a bibliography corpus
- IP aliases over the Internet
- User matching across social networks
- Name disambiguation in a collection of documents
- Multiple records across databases

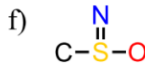
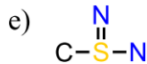
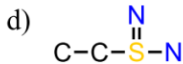
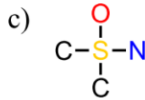
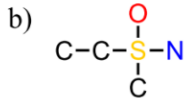
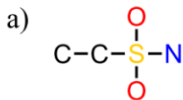
Group detection



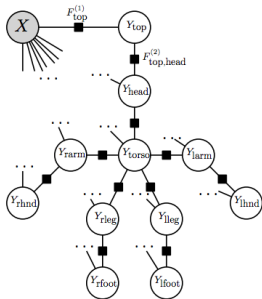
Group detection

- Communities in mobile networks
- Terrorist cells in social networks
- Correlation in international flight databases
- Document ranking and hypertext connectivity
- Information retrieval and document clustering
- ...

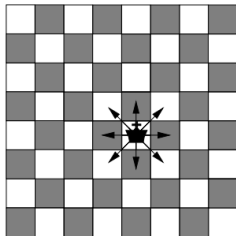
Frequent subgraph discovery



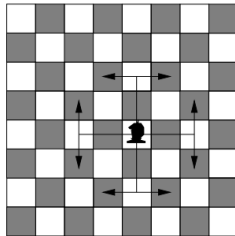
Structure learning



Predicate invention



a) The King



b) The Knight

[Image from Khan et al., ILP 1998]

The “relational revolution” of SRL stays at the intersection of Artificial Intelligence, Logic and Statistics

- **Logic**

- powerful and expressive formalism
- describes a domain in terms of quantified logic formulae
- allows to easily include background knowledge

- **Probabilistic graphical models**

- represent dependencies between random variables
- naturally handle uncertainty in data
- i.e., Bayesian Networks, Markov Random Fields

Propositional logic

- each **proposition** is associated to a symbol
- symbols have an associated **truth value** (true/false)
- logical operators (connectives) and inference at **propositional** level

Example

P = If it rains, Alice carries an umbrella

Q = It rains

R = Alice carries an umbrella

Can we **infer (entail)** R from P and Q ?

First-Order Logic (FOL)

- propositional logic is **not enough expressive**
- FOL described domains with objects, properties, relations
- it employs connectives as well as **quantifiers**

Constants: Alice, Italy, Snoopy, ...

Variables: x, y, z, \dots

Functions: father_of, ...

Predicates: blonde, friends, ...

Atoms: blonde(Alice), friends(x , father_of(y))

Literals: blonde(Alice), \neg blonde(Alice)

Formulas: $\forall x,y,z \text{ friends}(x,y) \wedge \text{friends}(y,z) \Rightarrow \text{friends}(x,z)$

Inference in first-order logic is **semi-decidable**

Remark

Semi-decidability of logic entailment $KB \models F$ implies that there is an algorithm which can always tell correctly when F is entailed by KB , but can provide either a negative answer or no answer at all in the case that F is not entailed by KB .

A research area that started **during the 70s**

- Logic programming (Prolog)

Key ideas:

- describe the domain of interest in terms of **logic facts**
- **background knowledge** is also encoded in the logic database
- **learn concepts** (rules) directly from data
- try to **entail** all positive examples and none of the negatives

Example

FOIL – First Order Inductive Learner

- **General-to-specific** learner
- Starts from an empty clause
- Iteratively adds literals
- Exploits **information gain**

Example:

strawberry(x)

strawberry(x) :- red(x)

strawberry(x) :- red(x), smaller_than_apple(x)

strawberry(x) :- red(x), smaller_than_apple(x), bigger_than_ribes(x)

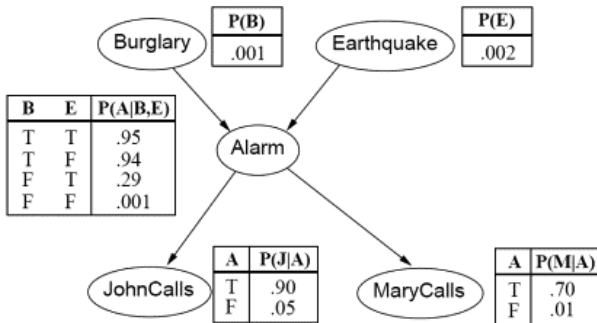
A graphical model is a **probabilistic model** where a graph encodes the **conditional dependence between random variables**

The model can be either directed or undirected:

- **directed** → Bayesian networks
- **undirected** → Markov networks

Bayesian networks

- The network structure is a **directed acyclic graph**
- It represents the **joint probability** of random variables
- Node have associated local **conditional probability tables**



[Image from Russel & Norvig]

Given variables X_1, \dots, X_n , the joint probability is computed as:

$$P(X_1, \dots, X_n) = \prod_{i=1}^n P(x_i | pa_i)$$

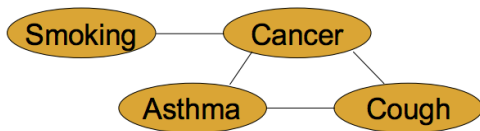
where pa_i are the parents of node i

Use **Bayes theorem** to find the values of other variables

Markov networks (or Markov random fields)

- The network structure is an **undirected graph**
- It represents the **joint probability** of random variables
- Random variables have to satisfy Markov properties
- Each clique has an associated **potential function**

$$P(X = x) = \frac{1}{Z} \prod_{C \in cl(G)} \phi_C(x_C)$$

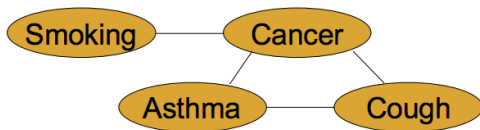


[Image from washington.edu]

Markov networks (or Markov random fields)

Markov properties:

- 1 Any two non-adjacent variables are conditionally independent given all other variables
- 2 A variable is conditionally independent of all other variables given its neighbors
- 3 Any two subsets of variables are conditionally independent given a separating subset



[Image from washington.edu]

Any Markov network can be written as a **log-linear model**:

$$P(X = x) = \frac{1}{Z} \exp \left(\sum_k w_k^T f(x_{\{k\}}) \right)$$

Remark

A Markov network can represent dependencies which a Bayesian network cannot (e.g., cyclic dependencies) and the other way round (induced dependencies)

Probabilistic inference: **inferring** the posterior distribution of **unobserved** variables, given observed ones

Exact methods (#P-complete):

- variable elimination
- clique tree propagation
- recursive conditioning
- ...

Approximate methods:

- loopy belief propagation
- Markov Chain Monte Carlo
- importance sampling
- ...

Example: **Markov Chain Monte Carlo with Gibbs sampling**

- 1 state \leftarrow random truth assignment
- 2 for $i \leftarrow 1$ to numSamples
- 3 for each variable x
- 4 sample x according to $P(x|\text{neighbors}(x))$
- 5 state \leftarrow state with new value of x
- 6 $P(F) \leftarrow$ fraction of states where F is true

[Slide from washington.edu]

The SRL alphabet soup:

- Stochastic Logic Programs (SLPs) [Muggleton, 1996]
- Relational Bayesian Networks (RBNs) [Jaeger, 1997]
- Bayesian Logic Programs (BLPs) [Kersting & De Raedt, 2001]
- Probabilistic Relational Models (PRMs) [Friedman et al., 2001]
- Relational Dependency Networks (RDNs) [Neville & Jensen, 2007]
- Problog [De Raedt et al., 2007]
- Type Extension Trees (TETs) [Jaeger & al., 2013]
- Learning From Constraints (LFC) [Gori & al., 2014]
- **Markov Logic Networks (MLNs)** [Domingos & Richardson, 2006]
- ...

Markov Logic Networks

An MLN is a set of FOL formulae with attached weights

An example

```
1.2 Friends(x,y) ∧ WatchedMovie(x,m) => WatchedMovie(y,m)
2.3 Friends(x,y) ∧ Friends(y,z) => Friends(x,z)
0.8 LikedMovie(x,m) ∧ Friends(x,y) => LikedMovie(y,m)
... ..
```

The **higher** the weight of a clause \rightarrow

\rightarrow The **lower** the probability for a world violating that clause

What is a **world** or **Herbrand interpretation** ?

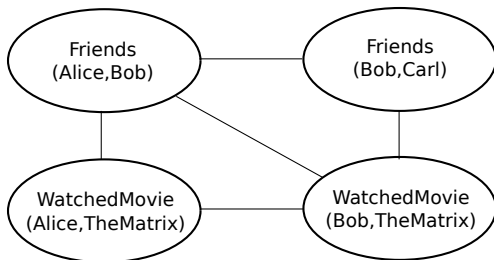
\rightarrow A truth assignment to all ground predicates

Together with a (finite) set of (unique and possibly typed) constants, an MLN defines a **Markov Network** which contains:

- 1 a binary node for each predicate grounding in the MLN, with value 0/1 if the atom is false/true
- 2 an edge between two nodes appearing together in (at least) one formula on the MLN
- 3 a feature for each formula grounding in the MLN, whose value is 0/1 if the formula is false/true, and whose weight is the weight of the formula

Set of constants

```
people = {Alice,Bob,Carl,David}
movie = {BladeRunner,ForrestGump,PulpFiction,TheMatrix}
```



The semantics of MLNs induces a probability distribution over all possible worlds. We indicate with X a set of random variables represented in the model, then we have:

$$P(X = x) = \frac{\exp\left(\sum_{F_i \in \mathcal{F}} w_i n_i(x)\right)}{Z}$$

being $n_i(x)$ the number of true groundings of formula i in world x .

The definition is similar to the joint probability distribution induced by a Markov network and expressed with a log-linear model:

$$P(X = x) = \frac{\exp\left(\sum_j w_j f_j(x)\right)}{Z}$$

- Typically, some atoms are always observed (**evidence X**), while others are unknown at prediction time (**query Y**)

EVIDENCE

Friends(Alice,Bob)

Friends(Bob,Carl)

WatchedMovie(Alice,PulpFiction)

WatchedMovie(David,BladeRunner)

...

QUERY

LikedMovie(Alice,BladeRunner) ?

LikedMovie(Alice,PulpFiction) ?

LikedMovie(Bob,BladeRunner) ?

LikedMovie(Bob,TheMatrix) ?

...

$$P(Y = y|X = x) = \frac{\exp\left(\sum_{F_i \in F_y} w_i n_i(x, y)\right)}{Z_x}$$

Three main MLN tasks:

- 1 Inference
- 2 Parameter Learning
- 3 Structure Learning

Popularity of Markov Logic:

→ **Alchemy**, an open-source software

In the discriminative setting, inference corresponds to finding **the most likely interpretation** (MAP – Maximum A Posteriori)

$$y^* = \underset{y}{\operatorname{argmax}} P(Y = y | X = x)$$

- **#P-complete** problem → **approximate** algorithms
- **MaxWalkSAT** [Kautz et al., 1996], stochastic local search
→ corresponds to minimizing the sum of weights of unsatisfied clauses

Note: there are also algorithms for estimating the **probability** of query atoms to be true

Main ideas:

- start with a **random** truth value assignment
- **flip** the atom giving the highest improvement (greedy)
- can get stuck in local minima
- sometimes perform a **random flip**
- stochastic algorithm (many runs often needed)
- need to build the whole **ground** network !

Main drawback

Memory explosion: with N constants and c as highest clause arity, the ground network requires $O(n^c)$ memory

One possible solution

Exploit sparseness by **lazily** grounding clauses (LazySAT) ... But

still not enough !

Key ideas:

- exploit **symmetries**
- reason at **first-order** level
- reason about **groups of objects**
- **scalable inference**

In many cases, one would prefer to have a **probability distribution** over query atoms, rather than just inferring the most likely state.

Several attempts:

- Markov Chain Monte Carlo over the Markov Network induced by the Markov Logic Networks, and the given set of constants
- Markov Chain SAT (MC-SAT), combining Markov Chain Monte Carlo and WalkSAT

Maximize the **conditional likelihood** of query predicates given evidence ones: requires **inference** as subroutine !

$$\frac{\partial}{\partial w_i} \log P(Y = y|X = x) = n_i - E_w[n_i]$$

Several different algorithms can be adopted to address the task:

- Voted Perceptron
- Contrastive Divergence
- Diagonal Newton
- (Preconditioned) Scaled Conjugate Gradient

Directly **infer** the rules from the data !

A classic task for Inductive Logic Programming (ILP), which for MLNs can be addressed together with parameter learning

- Modified ILP algorithms (e.g., Aleph)
- Bottom-Up Clause Learning
- Iterated Local Search
- Structural Motifs

Still an open problem !

Modified ILP algorithms

Two-step approach:

- 1 learn a set of clauses with FOIL
- 2 just learn the weights with Alchemy

It works quite well in practice, but it still relies on a **non-probabilistic** ILP framework !

Jointly learning clauses and weights should be better !

Bottom-Up Clause Learning, Iterated Local Search, Structural Motifs

- Start with **unit clauses** (single literal)
- Add/Remove literal
- Measure evaluation function (e.g. pseudo-likelihood)
- Learn weights by **counting groundings**

Hypertext Classification

Topic(page,topic)

HasWord(page,word)

Link(page,page)

HasWord(p,+w) => Topic(p,+t)

Topic(p,t) \wedge Link(p,q) => Topic(q,t)

Information Retrieval

InQuery(word)

HasWord(page,word)

Link(page,page)

Relevant(page)

$\text{HasWord}(p,+w) \wedge \text{InQuery}(w) \Rightarrow \text{Relevant}(p)$

$\text{Relevant}(p) \wedge \text{Link}(p,q) \Rightarrow \text{Relevant}(q)$

Entity Resolution

HasToken(token,field,record)

SameField(field,record,record)

SameRecord(record,record)

$\text{HasToken}(+t,+f,r1) \wedge \text{HasToken}(+t,+f,r2) \Rightarrow$

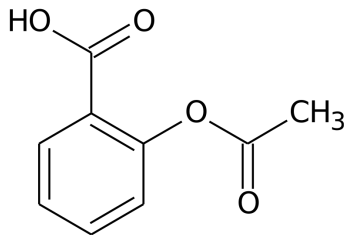
$\text{SameField}(f,r,s)$

$\text{SameField}(f,r1,r2) \Rightarrow \text{SameRecord}(r1,r2)$

$\text{SameRecord}(r1,r2) \wedge \text{SameRecord}(r2,r3) \Rightarrow$

$\text{SameRecord}(r1,r3)$

From images to structured data...



- IsOxygen(A1)
- IsCarbon(A2)
- IsCarbon(A3)
- DoubleBond(A1,A2)
- SingleBond(A1,A3)
- ...

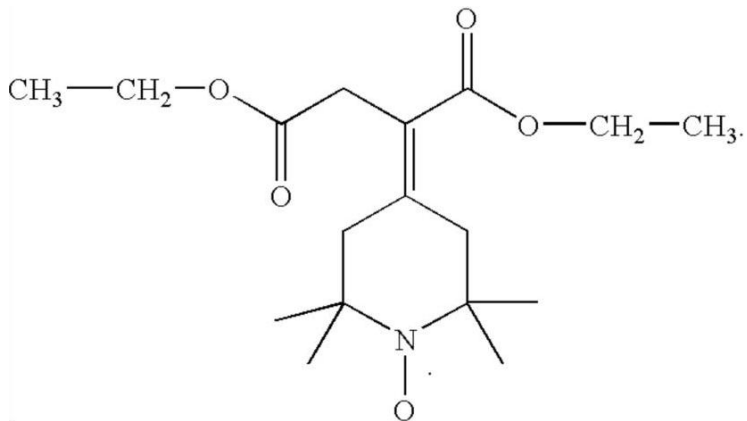
<http://mlocsr.dinfo.unifi.it>

Two-stage architecture

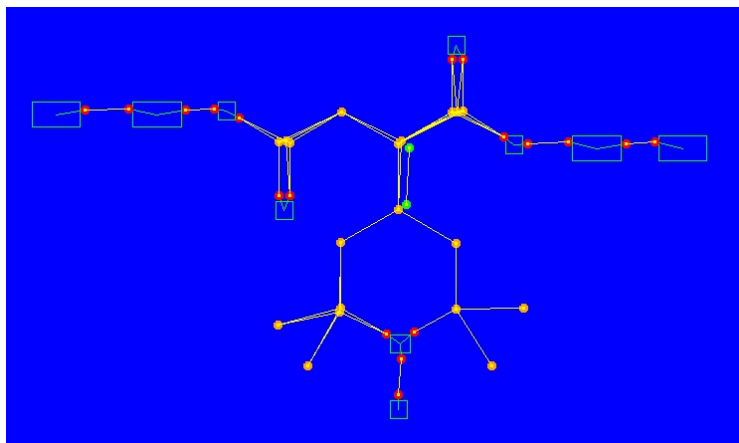
- 1 Image Processing level
- 2 Markov Logic level

Main idea

- extract low-level features
- translate them into Markov logic rules
- infer the final structure



Output of Image Processing stage



Output of Markov Logic stage



Opponent modeling in games

Games have always been a great opportunity for AI

- incomplete information
- uncertainty in data
- a lot of background knowledge

A nice domain also for SRL

Opponent modeling in games

Understand an adversarial's move given her past behavior:

- describe the domain in terms of **logic predicates**
- represent past played matches
- try to learn **strategies** through probabilistic rules

In many cases it is necessary to use **continuous features**

- classic machine learning classifiers naturally employ them
- they are often crucial in order to build accurate predictors

→ How to integrate them with first-order logic ?

Solution: grounding-specific weights [Lippi & Frasconi, 2009]

Instead of having a weight for a first-order logic rule, we allow different weights for different ground clauses.

Introducing vector of features

$$\text{Node}(N12, \$\text{Features_}N12) \wedge \text{Node}(N23, \$\text{Features_}N23) \\ \Rightarrow \text{Link}(N12, N23)$$
$$\text{Node}(N18, \$\text{Features_}N18) \wedge \text{Node}(N25, \$\text{Features_}N25) \\ \Rightarrow \text{Link}(N18, N25)$$

The two rules will have different weights, which will be computed using as feature vectors the “constants” with the dollar symbol (\$).

Re-parametrize the MLN by computing each weight w_i as a function of variables of each specific grounding c_{ij} :

Standard MLN

$$P(Y = y|X = x) = \frac{\exp\left(\sum_{F_i \in F_y} w_i n_i(x, y)\right)}{Z_x}$$

MLNs with grounding-specific weights

$$P(Y = y|X = x) = \frac{\exp\left(\sum_{F_i \in F_y} \sum_j w_i(c_{ij}, \theta_i) n_{ij}(x, y)\right)}{Z_x}$$

Instead of having a weight for a first-order logic rule, we allow different weights for different ground clauses.

The weights $w_i(c_{ij}, \theta_i)$ can be computed in several ways

- using **neural networks**, by taking as input an encoding of the grounding c_{ij} (in principle any feature can be used !)
- **Inference** algorithms do not change.
- **Learning** algorithm can implement gradient descent:

$$\frac{\partial P(y|x)}{\partial \theta_k} = \frac{\partial P(y|x)}{\partial w_i} \frac{\partial w_i}{\partial \theta_k}$$

where the **first** term is computed by MLN inference and the **second** one is computed by backpropagation.

Other frameworks

Generalization of:

- Hidden Markov Models
- Probabilistic Context-Free Grammars

Based on the concept of **stochastic clause**:

$$\lambda \quad A \leftarrow B_1, \dots, B_n$$

The weight λ represents the probability that a clause is used in a **derivation**, provided that **head is true**

If some rules have no λ , the SLP is called **impure**, and it can allow **both** stochastic and deterministic derivations

Directed acyclic graph:

- each node corresponds to a **relation** r in the domain
- attached **probability formula** $F_r(x_1, \dots, x_n)$ over the symbols in the parent nodes of r

Semantics

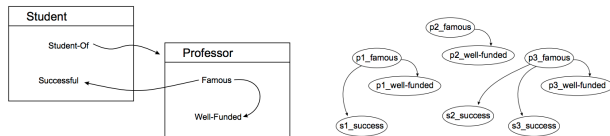
For every relation r , the RBN associates a probability distribution over interpretations of r , which is defined on the interpretations of $pa(r)$ within the underlying Bayesian network

Probability distribution over a **relational database**:

- description of the relational schema of the domain
- probabilistic dependencies among its objects

Given a set of ground objects, a PRM specifies a **joint probability distribution** over all instantiations of a relational schema.

Probabilities on unseen variables, given **partial** instantiations



[Image from mit.edu]

Probabilistic extension of Prolog

- logic clauses with associated **probability**
- probability of each clause is **independent** of the others
- distribution over **logic programs**

$$P(L|T) = \prod_{c_j \in L} p_j \prod_{c_j \notin L} (1 - p_j)$$

Initially designed to answer probabilistic queries

Representation formalism for **complex combinatorial features** in relational domains.

A TET is a tree which represents a feature or a set of features:

- nodes contain **literals**, or conjunctions of literals
- edges are labeled with **sets of variables** (possibly empty)

TETs are syntactically closely related to predicate logic formulae, where subtrees correspond to sub-formulae.

TETs have been used for feature discovery and for classification

Combining kernel machines and first-order logic

Translate logic clauses (background knowledge) into a **semantic regularization** term, in addition to classic regularization

$$\min_f \sum_{i=1}^N V(f) + R(f) + S(f)$$

$S(f)$ can encode **generic constraints**

- logic rules transformed through p -norms

SRL aims at combining logic and statistical learning
→ crucial when dealing with relational data

The problem is far from being solved:

- **scalability**
- **expressivity**
- more **efficient** approximation algorithms
- extend to many other **applications** and **contexts** !

... How about incorporating **deep learning** ?