

# Markov Logic Networks

Matthew Richardson (mattr@cs.washington.edu) and

Pedro Domingos (pedrod@cs.washington.edu)

*Department of Computer Science and Engineering, University of Washington, Seattle, WA 98195-250, U.S.A.*

**Abstract.** We propose a simple approach to combining first-order logic and probabilistic graphical models in a single representation. A Markov logic network (MLN) is a first-order knowledge base with a weight attached to each formula (or clause). Together with a set of constants representing objects in the domain, it specifies a ground Markov network containing one feature for each possible grounding of a first-order formula in the KB, with the corresponding weight. Inference in MLNs is performed by MCMC over the minimal subset of the ground network required for answering the query. Weights are efficiently learned from relational databases by iteratively optimizing a pseudo-likelihood measure. Optionally, additional clauses are learned using inductive logic programming techniques. Experiments with a real-world database and knowledge base in a university domain illustrate the promise of this approach.

**Keywords:** Statistical relational learning, Markov networks, Markov random fields, log-linear models, graphical models, first-order logic, satisfiability, inductive logic programming, knowledge-based model construction, Markov chain Monte Carlo, pseudo-likelihood, link prediction

## 1. Introduction

Combining probability and first-order logic in a single representation has long been a goal of AI. Probabilistic graphical models enable us to efficiently handle uncertainty. First-order logic enables us to compactly represent a wide variety of knowledge. Many (if not most) applications require both. Interest in this problem has grown in recent years due to its relevance to statistical relational learning (Getoor & Jensen, 2000; Getoor & Jensen, 2003; Dietterich et al., 2003), also known as multi-relational data mining (Džeroski & De Raedt, 2003; Džeroski et al., 2002; Džeroski et al., 2003; Džeroski & Blockeel, 2004). Current proposals typically focus on combining probability with restricted subsets of first-order logic, like Horn clauses (e.g., Wellman et al. (1992); Poole (1993); Muggleton (1996); Ngo and Haddawy (1997); Sato and Kameya (1997); Cussens (1999); Kersting and De Raedt (2001); Santos Costa et al. (2003)), frame-based systems (e.g., Friedman et al. (1999); Pasula and Russell (2001); Cumby and Roth (2003)), or database query languages (e.g., Taskar et al. (2002); Popescul and Ungar (2003)). They are often quite complex. In this paper, we introduce *Markov logic networks (MLNs)*, a representation that is quite simple, yet combines probability and first-order logic with no restrictions other than finiteness of the domain. We develop

efficient algorithms for inference and learning in MLNs, and evaluate them in a real-world domain.

A Markov logic network is a first-order knowledge base with a weight attached to each formula, and can be viewed as a template for constructing Markov networks. From the point of view of probability, MLNs provide a compact language to specify very large Markov networks, and the ability to flexibly and modularly incorporate a wide range of domain knowledge into them. From the point of view of first-order logic, MLNs add the ability to soundly handle uncertainty, tolerate imperfect and contradictory knowledge, and reduce brittleness. Many important tasks in statistical relational learning, like collective classification, link prediction, link-based clustering, social network modeling, and object identification, are naturally formulated as instances of MLN learning and inference.

Experiments with a real-world database and knowledge base illustrate the benefits of using MLNs over purely logical and purely probabilistic approaches. We begin the paper by briefly reviewing the fundamentals of Markov networks (Section 2) and first-order logic (Section 3). The core of the paper introduces Markov logic networks and algorithms for inference and learning in them (Sections 4–6). We then report our experimental results (Section 7). Finally, we show how a variety of SRL tasks can be cast as MLNs (Section 8), discuss how MLNs relate to previous approaches (Section 9) and list directions for future work (Section 10).

## 2. Markov Networks

A *Markov network* (also known as *Markov random field*) is a model for the joint distribution of a set of variables  $X = (X_1, X_2, \dots, X_n) \in \mathcal{X}$  (Pearl, 1988). It is composed of an undirected graph  $G$  and a set of potential functions  $\phi_k$ . The graph has a node for each variable, and the model has a potential function for each clique in the graph. A potential function is a non-negative real-valued function of the state of the corresponding clique. The joint distribution represented by a Markov network is given by

$$P(X=x) = \frac{1}{Z} \prod_k \phi_k(x_{\{k\}}) \quad (1)$$

where  $x_{\{k\}}$  is the state of the  $k$ th clique (i.e., the state of the variables that appear in that clique).  $Z$ , known as the *partition function*, is given by  $Z = \sum_{x \in \mathcal{X}} \prod_k \phi_k(x_{\{k\}})$ . Markov networks are often conveniently represented as *log-linear models*, with each clique potential replaced by an exponentiated weighted sum of features of the state, leading to

$$P(X=x) = \frac{1}{Z} \exp \left( \sum_j w_j f_j(x) \right) \quad (2)$$

A feature may be any real-valued function of the state. This paper will focus on binary features,  $f_j(x) \in \{0, 1\}$ . In the most direct translation from the potential-function form (Equation 1), there is one feature corresponding to each possible state  $x_{\{k\}}$  of each clique, with its weight being  $\log \phi_k(x_{\{k\}})$ . This representation is exponential in the size of the cliques. However, we are free to specify a much smaller number of features (e.g., logical functions of the state of the clique), allowing for a more compact representation than the potential-function form, particularly when large cliques are present. MLNs will take advantage of this.

Inference in Markov networks is #P-complete (Roth, 1996). The most widely used method for approximate inference in Markov networks is Markov chain Monte Carlo (MCMC) (Gilks et al., 1996), and in particular Gibbs sampling, which proceeds by sampling each variable in turn given its Markov blanket. (The Markov blanket of a node is the minimal set of nodes that renders it independent of the remaining network; in a Markov network, this is simply the node's neighbors in the graph.) Marginal probabilities are computed by counting over these samples; conditional probabilities are computed by running the Gibbs sampler with the conditioning variables clamped to their given values. Another popular method for inference in Markov networks is belief propagation (Yedidia et al., 2001).

Maximum-likelihood or MAP estimates of Markov network weights cannot be computed in closed form, but, because the log-likelihood is a concave function of the weights, they can be found efficiently using standard gradient-based or quasi-Newton optimization methods (Nocedal & Wright, 1999). Another alternative is iterative scaling (Della Pietra et al., 1997). Features can also be learned from data, for example by greedily constructing conjunctions of atomic features (Della Pietra et al., 1997).

### 3. First-Order Logic

A *first-order knowledge base (KB)* is a set of sentences or formulas in first-order logic (Genesereth & Nilsson, 1987). Formulas are constructed using four types of symbols: constants, variables, functions, and predicates. Constant symbols represent objects in the domain of interest (e.g., people: Anna, Bob, Chris, etc.). Variable symbols range over the objects in the domain. Function symbols (e.g., `MotherOf`) represent mappings from tuples of objects to objects. Predicate symbols represent relations among objects in the domain (e.g., `Friends`) or attributes of objects (e.g., `Smokes`). An *inter-*

*pretation* specifies which objects, functions and relations in the domain are represented by which symbols. Variables and constants may be *typed*, in which case variables range only over objects of the corresponding type, and constants can only represent objects of the corresponding type. For example, the variable  $x$  might range over people (e.g., Anna, Bob, etc.), and the constant  $C$  might represent a city (e.g., Seattle).

A *term* is any expression representing an object in the domain. It can be a constant, a variable, or a function applied to a tuple of terms. For example,  $Anna$ ,  $x$ , and  $\text{GreatestCommonDivisor}(x, y)$  are terms. An *atomic formula* or *atom* is a predicate symbol applied to a tuple of terms (e.g.,  $\text{Friends}(x, \text{MotherOf}(Anna))$ ). Formulas are recursively constructed from atomic formulas using logical connectives and quantifiers. If  $F_1$  and  $F_2$  are formulas, the following are also formulas:  $\neg F_1$  (negation), which is true iff  $F_1$  is false;  $F_1 \wedge F_2$  (conjunction), which is true iff both  $F_1$  and  $F_2$  are true;  $F_1 \vee F_2$  (disjunction), which is true iff  $F_1$  or  $F_2$  is true;  $F_1 \Rightarrow F_2$  (implication), which is true iff  $F_1$  is false or  $F_2$  is true;  $F_1 \Leftrightarrow F_2$  (equivalence), which is true iff  $F_1$  and  $F_2$  have the same truth value;  $\forall x F_1$  (universal quantification), which is true iff  $F_1$  is true for every object  $x$  in the domain; and  $\exists x F_1$  (existential quantification), which is true iff  $F_1$  is true for at least one object  $x$  in the domain. Parentheses may be used to enforce precedence. A *positive literal* is an atomic formula; a *negative literal* is a negated atomic formula. The formulas in a KB are implicitly conjoined, and thus a KB can be viewed as a single large formula. A *ground term* is a term containing no variables. A *ground atom* or *ground predicate* is an atomic formula all of whose arguments are ground terms. A *possible world* or *Herbrand interpretation* assigns a truth value to each possible ground atom.

A formula is *satisfiable* iff there exists at least one world in which it is true. The basic inference problem in first-order logic is to determine whether a knowledge base  $KB$  *entails* a formula  $F$ , i.e., if  $F$  is true in all worlds where  $KB$  is true (denoted by  $KB \models F$ ). This is often done by *refutation*:  $KB$  entails  $F$  iff  $KB \cup \neg F$  is unsatisfiable. (Thus, if a KB contains a contradiction, all formulas trivially follow from it, which makes painstaking knowledge engineering a necessity.) For automated inference, it is often convenient to convert formulas to a more regular form, typically *clausal form* (also known as *conjunctive normal form (CNF)*). A KB in clausal form is a conjunction of *clauses*, a clause being a disjunction of literals. Every KB in first-order logic can be converted to clausal form using a mechanical sequence of steps.<sup>1</sup> Clausal form is used in resolution, a sound and refutation-complete inference procedure for first-order logic (Robinson, 1965).

---

<sup>1</sup> This conversion includes the removal of existential quantifiers by Skolemization, which is not sound in general. However, in finite domains an existentially quantified formula can simply be replaced by a disjunction of its groundings.

Inference in first-order logic is only semidecidable. Because of this, knowledge bases are often constructed using a restricted subset of first-order logic with more desirable properties. The most widely-used restriction is to *Horn clauses*, which are clauses containing at most one positive literal. The Prolog programming language is based on Horn clause logic (Lloyd, 1987). Prolog programs can be learned from databases by searching for Horn clauses that (approximately) hold in the data; this is studied in the field of inductive logic programming (ILP) (Lavrač & Džeroski, 1994).

Table I shows a simple KB and its conversion to clausal form. Notice that, while these formulas may be *typically* true in the real world, they are not *always* true. In most domains it is very difficult to come up with non-trivial formulas that are always true, and such formulas capture only a fraction of the relevant knowledge. Thus, despite its expressiveness, pure first-order logic has limited applicability to practical AI problems. Many *ad hoc* extensions to address this have been proposed. In the more limited case of propositional logic, the problem is well solved by probabilistic graphical models. The next section describes a way to generalize these models to the first-order case.

#### 4. Markov Logic Networks

A first-order KB can be seen as a set of hard constraints on the set of possible worlds: if a world violates even one formula, it has zero probability. The basic idea in MLNs is to soften these constraints: when a world violates one formula in the KB it is less probable, but not impossible. The fewer formulas a world violates, the more probable it is. Each formula has an associated weight that reflects how strong a constraint it is: the higher the weight, the greater the difference in log probability between a world that satisfies the formula and one that does not, other things being equal.

**DEFINITION 4.1.** *A Markov logic network  $L$  is a set of pairs  $(F_i, w_i)$ , where  $F_i$  is a formula in first-order logic and  $w_i$  is a real number. Together with a finite set of constants  $C = \{c_1, c_2, \dots, c_{|C|}\}$ , it defines a Markov network  $M_{L,C}$  (Equations 1 and 2) as follows:*

1.  $M_{L,C}$  contains one binary node for each possible grounding of each predicate appearing in  $L$ . The value of the node is 1 if the ground atom is true, and 0 otherwise.
2.  $M_{L,C}$  contains one feature for each possible grounding of each formula  $F_i$  in  $L$ . The value of this feature is 1 if the ground formula is true, and 0 otherwise. The weight of the feature is the  $w_i$  associated with  $F_i$  in  $L$ .

Table I. Example of a first-order knowledge base and MLN.  $\text{Fr}()$  is short for  $\text{Friends}()$ ,  $\text{Sm}()$  for  $\text{Smokes}()$ , and  $\text{Ca}()$  for  $\text{Cancer}()$ .

English	First-Order Logic	Clausal Form	Weight
Friends of friends are friends.	$\forall x \forall y \forall z \text{Fr}(x, y) \wedge \text{Fr}(y, z) \Rightarrow \text{Fr}(x, z)$	$\neg \text{Fr}(x, y) \vee \neg \text{Fr}(y, z) \vee \text{Fr}(x, z)$	0.7
Friendless people smoke.	$\forall x (\neg(\exists y \text{Fr}(x, y)) \Rightarrow \text{Sm}(x))$	$\text{Fr}(x, g(x)) \vee \text{Sm}(x)$	2.3
Smoking causes cancer.	$\forall x \text{Sm}(x) \Rightarrow \text{Ca}(x)$	$\neg \text{Sm}(x) \vee \text{Ca}(x)$	1.5
If two people are friends, either	$\forall x \forall y \text{Fr}(x, y) \Rightarrow (\text{Sm}(x) \Leftrightarrow \text{Sm}(y))$	$\neg \text{Fr}(x, y) \vee \text{Sm}(x) \vee \neg \text{Sm}(y),$	1.1
both smoke or neither does.		$\neg \text{Fr}(x, y) \vee \neg \text{Sm}(x) \vee \text{Sm}(y)$	1.1

The syntax of the formulas in an MLN is the standard syntax of first-order logic (Genesereth & Nilsson, 1987). Free (unquantified) variables are treated as universally quantified at the outermost level of the formula.

An MLN can be viewed as a *template* for constructing Markov networks. Given different sets of constants, it will produce different networks, and these may be of widely varying size, but all will have certain regularities in structure and parameters, given by the MLN (e.g., all groundings of the same formula will have the same weight). We call each of these networks a *ground Markov network* to distinguish it from the first-order MLN. From Definition 4.1 and Equations 1 and 2, the probability distribution over possible worlds  $x$  specified by the ground Markov network  $M_{L,C}$  is given by

$$P(X=x) = \frac{1}{Z} \exp\left(\sum_i w_i n_i(x)\right) = \frac{1}{Z} \prod_i \phi_i(x_{\{i\}})^{n_i(x)} \quad (3)$$

where  $n_i(x)$  is the number of true groundings of  $F_i$  in  $x$ ,  $x_{\{i\}}$  is the state (truth values) of the atoms appearing in  $F_i$ , and  $\phi_i(x_{\{i\}}) = e^{w_i}$ . Notice that, although we defined MLNs as loglinear models, they could equally well be defined as products of potential functions, as the second equality above shows. This will be the most convenient approach in domains with a mixture of hard and soft constraints (i.e., where some formulas hold with certainty, leading to zero probabilities for some worlds).

The graphical structure of  $M_{L,C}$  follows from Definition 4.1: there is an edge between two nodes of  $M_{L,C}$  iff the corresponding ground atoms appear together in at least one grounding of one formula in  $L$ . Thus, the atoms in each ground formula form a (not necessarily maximal) clique in  $M_{L,C}$ . Figure 1 shows the graph of the ground Markov network defined by the last two formulas in Table I and the constants Anna and Bob. Each node in this graph is a ground atom (e.g., Friends(Anna, Bob)). The graph contains an arc between each pair of atoms that appear together in some grounding of one of the formulas.  $M_{L,C}$  can now be used to infer the probability that Anna and Bob are friends given their smoking habits, the probability that Bob has cancer given his friendship with Anna and whether she has cancer, etc.

Each state of  $M_{L,C}$  represents a possible world. A possible world is a set of objects, a set of functions (mappings from tuples of objects to objects), and a set of relations that hold between those objects; together with an interpretation, they determine the truth value of each ground atom. The following assumptions ensure that the set of possible worlds for  $(L, C)$  is finite, and that  $M_{L,C}$  represents a unique, well-defined probability distribution over those worlds, irrespective of the interpretation and domain. These assumptions are quite reasonable in most practical applications, and greatly simplify the use of MLNs. For the remaining cases, we discuss below the extent to which each one can be relaxed.

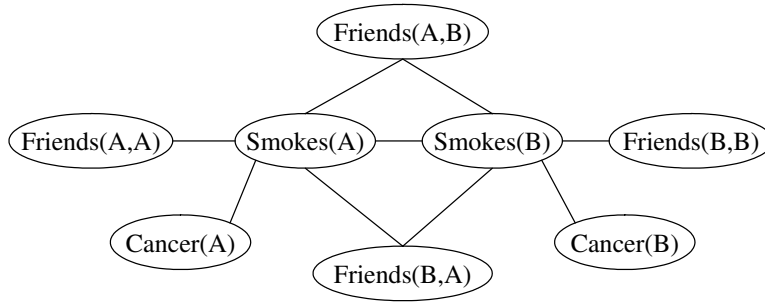


Figure 1. Ground Markov network obtained by applying the last two formulas in Table I to the constants `Anna(A)` and `Bob(B)`.

**ASSUMPTION 1. Unique names.** *Different constants refer to different objects (Genesereth & Nilsson, 1987).*

**ASSUMPTION 2. Domain closure.** *The only objects in the domain are those representable using the constant and function symbols in  $(L, C)$  (Genesereth & Nilsson, 1987).*

**ASSUMPTION 3. Known functions.** *For each function appearing in  $L$ , the value of that function applied to every possible tuple of arguments is known, and is an element of  $C$ .*

This last assumption allows us to replace functions by their values when grounding formulas. Thus the only ground atoms that need to be considered are those having constants as arguments. The infinite number of terms constructible from all functions and constants in  $(L, C)$  (the “Herbrand universe” of  $(L, C)$ ) can be ignored, because each of those terms corresponds to a known constant in  $C$ , and atoms involving them are already represented as the atoms involving the corresponding constants. The possible groundings of a predicate in Definition 4.1 are thus obtained simply by replacing each variable in the predicate with each constant in  $C$ , and replacing each function term in the predicate by the corresponding constant. Table II shows how the groundings of a formula are obtained given Assumptions 1–3.

Assumption 1 (unique names) can be removed by introducing the equality predicate (`Equals(x, y)`, or `x = y` for short) and adding the necessary axioms to the MLN: equality is reflexive, symmetric and transitive; for each unary predicate  $P$ ,  $\forall x \forall y \ x = y \Rightarrow (P(x) \Leftrightarrow P(y))$ ; and similarly for higher-order predicates and functions (Genesereth & Nilsson, 1987). The resulting MLN will have a node for each pair of constants, whose value is 1 if the constants represent the same object and 0 otherwise; these nodes will be connected to each other and to the rest of the network by arcs representing the axioms above. Notice that this allows us to make probabilistic inferences about the



Table II. Construction of all groundings of a first-order formula under Assumptions 1–3.

---

```

function Ground( $F$ )
  input:  $F$ , a formula in first-order logic
  output:  $G_F$ , a set of ground formulas
  for each existentially quantified subformula  $\exists x S(x)$  in  $F$ 
     $F \leftarrow F$  with  $\exists x S(x)$  replaced by  $S(c_1) \vee S(c_2) \vee \dots \vee S(c_{|C|})$ ,
    where  $S(c_i)$  is  $S(x)$  with  $x$  replaced by  $c_i$ 
   $G_F \leftarrow \{F\}$ 
  for each universally quantified variable  $x$ 
    for each formula  $F_j(x)$  in  $G_F$ 
       $G_F \leftarrow (G_F \setminus F_j(x)) \cup \{F_j(c_1), F_j(c_2), \dots, F_j(c_{|C|})\}$ ,
      where  $F_j(c_i)$  is  $F_j(x)$  with  $x$  replaced by  $c_i$ 
    for each formula  $F_j \in G_F$ 
      repeat
        for each function  $f(a_1, a_2, \dots)$  all of whose arguments are constants
           $F_j \leftarrow F_j$  with  $f(a_1, a_2, \dots)$  replaced by  $c$ , where  $c = f(a_1, a_2, \dots)$ 
        until  $F_j$  contains no functions
      return  $G_F$ 

```

---

equality of two constants. We have successfully used this as the basis of an approach to object identification (see Subsection 8.5).

If the number  $u$  of unknown objects is known, Assumption 2 (domain closure) can be removed simply by introducing  $u$  arbitrary new constants. If  $u$  is unknown but finite, Assumption 2 can be removed by introducing a distribution over  $u$ , grounding the MLN with each number of unknown objects, and computing the probability of a formula  $F$  as  $P(F) = \sum_{u=0}^{u_{max}} P(u)P(F|M_{L,C}^u)$ , where  $M_{L,C}^u$  is the ground MLN with  $u$  unknown objects. An infinite  $u$  requires extending MLNs to the case  $|C| = \infty$ .

Let  $H_{L,C}$  be the set of all ground terms constructible from the function symbols in  $L$  and the constants in  $L$  and  $C$  (the ‘‘Herbrand universe’’ of  $(L, C)$ ). Assumption 3 (known functions) can be removed by treating each element of  $H_{L,C}$  as an additional constant and applying the same procedure used to remove the unique names assumption. For example, with a function  $G(\mathbf{x})$  and constants  $A$  and  $B$ , the MLN will now contain nodes for  $G(A) = A$ ,  $G(A) = B$ , etc. This leads to an infinite number of new constants, requiring the corresponding extension of MLNs. However, if we restrict the level of nesting to some maximum, the resulting MLN is still finite.

To summarize, Assumptions 1–3 can be removed as long the domain is finite. We believe it is possible to extend MLNs to infinite domains (see Jaeger (1998)), but this is an issue of chiefly theoretical interest, and we leave it for future work. In the remainder of this paper we proceed under Assumptions 1–3, except where noted.

A first-order KB can be transformed into an MLN simply by assigning a weight to each formula. For example, the clauses and weights in the last two columns of Table I constitute an MLN. According to this MLN, other things being equal, a world where  $n$  friendless people are non-smokers is  $e^{(2.3)n}$  times less probable than a world where all friendless people smoke. Notice that all the formulas in Table I are false in the real world as universally quantified logical statements, but capture useful information on friendships and smoking habits, when viewed as features of a Markov network. For example, it is well known that teenage friends tend to have similar smoking habits (Lloyd-Richardson et al., 2002). In fact, an MLN like the one in Table I succinctly represents a type of model that is a staple of social network analysis (Wasserman & Faust, 1994).

It is easy to see that MLNs subsume essentially all propositional probabilistic models, as detailed below.

**PROPOSITION 4.2.** *Every probability distribution over discrete or finite-precision numeric variables can be represented as a Markov logic network.*

**Proof.** Consider first the case of Boolean variables  $(X_1, X_2, \dots, X_n)$ . Define a predicate of zero arity  $R_h$  for each variable  $X_h$ , and include in the MLN  $L$  a formula for each possible state of  $(X_1, X_2, \dots, X_n)$ . This formula is a conjunction of  $n$  literals, with the  $h$ th literal being  $R_h()$  if  $X_h$  is true in the state, and  $\neg R_h()$  otherwise. The formula's weight is  $\log P(X_1, X_2, \dots, X_n)$ . (If some states have zero probability, use instead the product form (see Equation 3), with  $\phi_i()$  equal to the probability of the  $i$ th state.) Since all predicates in  $L$  have zero arity,  $L$  defines the same Markov network  $M_{L,C}$  irrespective of  $C$ , with one node for each variable  $X_h$ . For any state, the corresponding formula is true and all others are false, and thus Equation 3 represents the original distribution (notice that  $Z = 1$ ). The generalization to arbitrary discrete variables is straightforward, by defining a zero-arity predicate for each value of each variable. Similarly for finite-precision numeric variables, by noting that they can be represented as Boolean vectors.  $\square$

Of course, compact factored models like Markov networks and Bayesian networks can still be represented compactly by MLNs, by defining formulas for the corresponding factors (arbitrary features in Markov networks, and states of a node and its parents in Bayesian networks).<sup>2</sup>

First-order logic (with Assumptions 1–3 above) is the special case of MLNs obtained when all weights are equal and tend to infinity, as described below.

<sup>2</sup> While some conditional independence structures can be compactly represented with directed graphs but not with undirected ones, they still lead to compact models in the form of Equation 3 (i.e., as products of potential functions).

**PROPOSITION 4.3.** *Let  $KB$  be a satisfiable knowledge base,  $L$  be the MLN obtained by assigning weight  $w$  to every formula in  $KB$ ,  $C$  be the set of constants appearing in  $KB$ ,  $P_w(x)$  be the probability assigned to a (set of) possible world(s)  $x$  by  $M_{L,C}$ ,  $\mathcal{X}_{KB}$  be the set of worlds that satisfy  $KB$ , and  $F$  be an arbitrary formula in first-order logic. Then:*

1.  $\forall x \in \mathcal{X}_{KB} \lim_{w \rightarrow \infty} P_w(x) = |\mathcal{X}_{KB}|^{-1}$   
 $\forall x \notin \mathcal{X}_{KB} \lim_{w \rightarrow \infty} P_w(x) = 0$
2. For all  $F$ ,  $KB \models F$  iff  $\lim_{w \rightarrow \infty} P_w(F) = 1$ .

**Proof.** Let  $k$  be the number of ground formulas in  $M_{L,C}$ . By Equation 3, if  $x \in \mathcal{X}_{KB}$  then  $P_w(x) = e^{kw}/Z$ , and if  $x \notin \mathcal{X}_{KB}$  then  $P_w(x) \leq e^{(k-1)w}/Z$ . Thus all  $x \in \mathcal{X}_{KB}$  are equiprobable and  $\lim_{w \rightarrow \infty} P(\mathcal{X} \setminus \mathcal{X}_{KB})/P(\mathcal{X}_{KB}) \leq \lim_{w \rightarrow \infty} (|\mathcal{X} \setminus \mathcal{X}_{KB}|/|\mathcal{X}_{KB}|)e^{-w} = 0$ , proving Part 1. By definition of entailment,  $KB \models F$  iff every world that satisfies  $KB$  also satisfies  $F$ . Therefore, letting  $\mathcal{X}_F$  be the set of worlds that satisfy  $F$ , if  $KB \models F$  then  $\mathcal{X}_{KB} \subseteq \mathcal{X}_F$  and  $P_w(F) = \sum_{x \in \mathcal{X}_F} P_w(x) \geq P_w(\mathcal{X}_{KB})$ . Since, from Part 1,  $\lim_{w \rightarrow \infty} P_w(\mathcal{X}_{KB}) = 1$ , this implies that if  $KB \models F$  then  $\lim_{w \rightarrow \infty} P_w(F) = 1$ . The inverse direction of Part 2 is proved by noting that if  $\lim_{w \rightarrow \infty} P_w(F) = 1$  then every world with non-zero probability in the limit must satisfy  $F$ , and this includes every world in  $\mathcal{X}_{KB}$ .  $\square$

In other words, in the limit of all equal infinite weights, the MLN represents a uniform distribution over the worlds that satisfy the KB, and all entailment queries can be answered by computing the probability of the query formula and checking whether it is 1. Even when weights are finite, first-order logic is “embedded” in MLNs in the following sense. Assume without loss of generality that all weights are non-negative. (A formula with a negative weight  $w$  can be replaced by its negation with weight  $-w$ .) If the knowledge base composed of the formulas in an MLN  $L$  (negated, if their weight is negative) is satisfiable, then, for any  $C$ , the satisfying assignments are the modes of the distribution represented by  $M_{L,C}$ . This is because the modes are the worlds  $x$  with maximum  $\sum_i w_i n_i(x)$  (see Equation 3), and this expression is maximized when all groundings of all formulas are true (i.e., the KB is satisfied). Unlike an ordinary first-order KB, however, an MLN can produce useful results even when it contains contradictions. An MLN can also be obtained by merging several KBs, even if they are partly incompatible. This is potentially useful in areas like the Semantic Web (Berners-Lee et al., 2001) and mass collaboration (Richardson & Domingos, 2003).

It is interesting to see a simple example of how MLNs generalize first-order logic. Consider an MLN containing the single formula  $\forall x R(x) \Rightarrow S(x)$  with weight  $w$ , and  $C = \{A\}$ . This leads to four possible worlds:

$\{\neg R(A), \neg S(A)\}$ ,  $\{\neg R(A), S(A)\}$ ,  $\{R(A), \neg S(A)\}$ , and  $\{R(A), S(A)\}$ . From Equation 3 we obtain that  $P(\{R(A), \neg S(A)\}) = 1/(3e^w + 1)$  and the probability of each of the other three worlds is  $e^w/(3e^w + 1)$ . (The denominator is the partition function  $Z$ ; see Section 2.) Thus, if  $w > 0$ , the effect of the MLN is to make the world that is inconsistent with  $\forall x R(x) \Rightarrow S(x)$  less likely than the other three. From the probabilities above we obtain that  $P(S(A)|R(A)) = 1/(1 + e^{-w})$ . When  $w \rightarrow \infty$ ,  $P(S(A)|R(A)) \rightarrow 1$ , recovering the logical entailment.

A first-order KB partitions the set of possible worlds into two subsets: those that satisfy the KB and those that do not. An MLN has many more degrees of freedom: it can partition the set of possible worlds into many more subsets, and assign a different probability to each. How to use this freedom is a key decision for both knowledge engineering and learning. At one extreme, the MLN can add little to logic, treating the whole knowledge base as a single formula, and assigning one probability to the worlds that satisfy it and another to the worlds that do not. At the other extreme, each formula in the KB can be converted into clausal form, and a weight associated with each clause.<sup>3</sup> The more finely divided into subformulas a KB is, the more gradual the dropoff in probability as a world violates more of those subformulas, and the greater the flexibility in specifying distributions over worlds. From a knowledge engineering point of view, the decision about which formulas constitute indivisible constraints should reflect domain knowledge and the goals of modeling. From a learning point of view, dividing the KB into more formulas increases the number of parameters, with the corresponding tradeoff in bias and variance.

In practice, we have found it useful to add each predicate to the MLN as a unit clause. In other words, for each predicate  $R(x_1, x_2, \dots)$  appearing in the MLN, we add the formula  $\forall x_1, x_2, \dots R(x_1, x_2, \dots)$  with some weight  $w_R$ . The weight of a unit clause can (roughly speaking) capture the marginal distribution of the corresponding predicate, leaving the weights of the non-unit clauses free to model only dependencies between predicates.

When manually constructing an MLN or interpreting a learned one, it is useful to have an intuitive understanding of the weights. The weight of a formula  $F$  is simply the log odds between a world where  $F$  is true and a world where  $F$  is false, other things being equal. However, if  $F$  shares variables with other formulas, as will typically be the case, it may not be possible to keep those formulas's truth values unchanged while reversing  $F$ 's. In this case there is no longer a one-to-one correspondence between

---

<sup>3</sup> This conversion can be done in the standard way (Genesereth & Nilsson, 1987), except that, instead of introducing Skolem functions, existentially quantified formulas should be replaced by disjunctions, as in Table II.

weights and probabilities of formulas.<sup>4</sup> Nevertheless, the probabilities of all formulas collectively determine all weights, if we view them as constraints on a maximum entropy distribution, or treat them as empirical probabilities and learn the maximum likelihood weights (the two are equivalent) (Della Pietra et al., 1997). Thus a good way to set the weights of an MLN is to write down the probability with which each formula should hold, treat these as empirical frequencies, and learn the weights from them using the algorithm in Section 6. Conversely, the weights in a learned MLN can be viewed as collectively encoding the empirical formula probabilities.

The size of ground Markov networks can be vastly reduced by having typed constants and variables, and only grounding variables to constants of the same type. However, even in this case the size of the network may be extremely large. Fortunately, many inferences do not require grounding the entire network, as we see in the next section.

## 5. Inference

MLNs can answer arbitrary queries of the form “What is the probability that formula  $F_1$  holds given that formula  $F_2$  does?” If  $F_1$  and  $F_2$  are two formulas in first-order logic,  $C$  is a finite set of constants including any constants that appear in  $F_1$  or  $F_2$ , and  $L$  is an MLN, then

$$\begin{aligned}
 P(F_1|F_2, L, C) &= P(F_1|F_2, M_{L,C}) \\
 &= \frac{P(F_1 \wedge F_2|M_{L,C})}{P(F_2|M_{L,C})} \\
 &= \frac{\sum_{x \in \mathcal{X}_{F_1} \cap \mathcal{X}_{F_2}} P(X=x|M_{L,C})}{\sum_{x \in \mathcal{X}_{F_2}} P(X=x|M_{L,C})} \quad (4)
 \end{aligned}$$

where  $\mathcal{X}_{F_i}$  is the set of worlds where  $F_i$  holds, and  $P(x|M_{L,C})$  is given by Equation 3. Ordinary conditional queries in graphical models are the special case of Equation 4 where all predicates in  $F_1$ ,  $F_2$  and  $L$  are zero-arity and the formulas are conjunctions. The question of whether a knowledge base  $KB$  entails a formula  $F$  in first-order logic is the question of whether  $P(F|L_{KB}, C_{KB,F}) = 1$ , where  $L_{KB}$  is the MLN obtained by assigning infinite weight to all the formulas in  $KB$ , and  $C_{KB,F}$  is the set of all constants

---

<sup>4</sup> This is an unavoidable side-effect of the power and flexibility of Markov networks. In Bayesian networks, parameters are probabilities, but at the cost of greatly restricting the ways in which the distribution may be factored. In particular, potential functions must be conditional probabilities, and the directed graph must have no cycles. The latter condition is particularly troublesome to enforce in relational extensions (Taskar et al., 2002).

appearing in  $KB$  or  $F$ . The question is answered by computing  $P(F|L_{KB}, C_{KB,F})$  by Equation 4, with  $F_2 = \text{True}$ .

Computing Equation 4 directly will be intractable in all but the smallest domains. Since MLN inference subsumes probabilistic inference, which is  $\#P$ -complete, and logical inference, which is NP-complete even in finite domains, no better results can be expected. However, many of the large number of techniques for efficient inference in either case are applicable to MLNs. Because MLNs allow fine-grained encoding of knowledge, including context-specific independences, inference in them may in some cases be more efficient than inference in an ordinary graphical model for the same domain. On the logic side, the probabilistic semantics of MLNs facilitates approximate inference, with the corresponding potential gains in efficiency.

In principle,  $P(F_1|F_2, L, C)$  can be approximated using an MCMC algorithm that rejects all moves to states where  $F_2$  does not hold, and counts the number of samples in which  $F_1$  holds. However, even this is likely to be too slow for arbitrary formulas. Instead, we provide an inference algorithm for the case where  $F_1$  and  $F_2$  are conjunctions of ground literals. While less general than Equation 4, this is the most frequent type of query in practice, and the algorithm we provide answers it far more efficiently than a direct application of Equation 4. Investigating lifted inference (where queries containing variables are answered without grounding them) is an important direction for future work (see Jaeger (2000) and Poole (2003) for initial results). The algorithm proceeds in two phases, analogous to knowledge-based model construction (Wellman et al., 1992). The first phase returns the minimal subset  $M$  of the ground Markov network required to compute  $P(F_1|F_2, L, C)$ . The algorithm for this is shown in Table III. The size of the network returned may be further reduced, and the algorithm sped up, by noticing that any ground formula which is made true by the evidence can be ignored, and the corresponding arcs removed from the network. In the worst case, the network contains  $O(|C|^a)$  nodes, where  $a$  is the largest predicate arity in the domain, but in practice it may be much smaller.

The second phase performs inference on this network, with the nodes in  $F_2$  set to their values in  $F_2$ . Our implementation uses Gibbs sampling, but any inference method may be employed. The basic Gibbs step consists of sampling one ground atom given its Markov blanket. The Markov blanket of a ground atom is the set of ground atoms that appear in some grounding of a formula with it. The probability of a ground atom  $X_l$  when its Markov blanket  $B_l$  is in state  $b_l$  is

$$\begin{aligned}
 P(X_l = x_l | B_l = b_l) & \\
 &= \frac{\exp(\sum_{f_i \in F_l} w_i f_i(X_l = x_l, B_l = b_l))}{\exp(\sum_{f_i \in F_l} w_i f_i(X_l = 0, B_l = b_l)) + \exp(\sum_{f_i \in F_l} w_i f_i(X_l = 1, B_l = b_l))}
 \end{aligned} \tag{5}$$

Table III. Network construction for inference in MLNs.

---

```

function ConstructNetwork( $F_1, F_2, L, C$ )
  inputs:  $F_1$ , a set of ground atoms with unknown truth values (the “query”)
            $F_2$ , a set of ground atoms with known truth values (the “evidence”)
            $L$ , a Markov logic network
            $C$ , a set of constants
  output:  $M$ , a ground Markov network
  calls:  $MB(q)$ , the Markov blanket of  $q$  in  $M_{L,C}$ 
   $G \leftarrow F_1$ 
  while  $F_1 \neq \emptyset$ 
    for all  $q \in F_1$ 
      if  $q \notin F_2$ 
         $F_1 \leftarrow F_1 \cup (MB(q) \setminus G)$ 
         $G \leftarrow G \cup MB(q)$ 
         $F_1 \leftarrow F_1 \setminus \{q\}$ 
  return  $M$ , the ground Markov network composed of all nodes in  $G$ , all arcs between them
  in  $M_{L,C}$ , and the features and weights on the corresponding cliques.

```

---

where  $F_l$  is the set of ground formulas that  $X_l$  appears in, and  $f_i(X_l = x_l, B_l = b_l)$  is the value (0 or 1) of the feature corresponding to the  $i$ th ground formula when  $X_l = x_l$  and  $B_l = b_l$ . For sets of atoms of which exactly one is true in any given world (e.g., the possible values of an attribute), blocking can be used (i.e., one atom is set to true and the others to false in one step, by sampling conditioned on their collective Markov blanket). The estimated probability of a conjunction of ground literals is simply the fraction of samples in which the ground literals are true, after the Markov chain has converged. Because the distribution is likely to have many modes, we run the Markov chain multiple times. When the MLN is in clausal form, we minimize burn-in time by starting each run from a mode found using MaxWalkSat, a local search algorithm for the weighted satisfiability problem (i.e., finding a truth assignment that maximizes the sum of weights of satisfied clauses) (Kautz et al., 1997). When there are hard constraints (clauses with infinite weight), MaxWalkSat finds regions that satisfy them, and the Gibbs sampler then samples from these regions to obtain probability estimates.

## 6. Learning

We learn MLN weights from one or more relational databases. (For brevity, the treatment below is for one database, but the generalization to many is trivial.) We make a closed world assumption (Genesereth & Nilsson, 1987): if a ground atom is not in the database, it is assumed to be false. If there are  $n$  possible ground atoms, a database is effectively a vector  $x = (x_1, \dots, x_l, \dots, x_n)$

where  $x_l$  is the truth value of the  $l$ th ground atom ( $x_l = 1$  if the atom appears in the database, and  $x_l = 0$  otherwise). Given a database, MLN weights can in principle be learned using standard methods, as follows. If the  $i$ th formula has  $n_i(x)$  true groundings in the data  $x$ , then by Equation 3 the derivative of the log-likelihood with respect to its weight is

$$\frac{\partial}{\partial w_i} \log P_w(X=x) = n_i(x) - \sum_{x'} P_w(X=x') n_i(x') \quad (6)$$

where the sum is over all possible databases  $x'$ , and  $P_w(X=x')$  is  $P(X=x')$  computed using the current weight vector  $w = (w_1, \dots, w_i, \dots)$ . In other words, the  $i$ th component of the gradient is simply the difference between the number of true groundings of the  $i$ th formula in the data and its expectation according to the current model. Unfortunately, counting the number of true groundings of a formula in a database is intractable, even when the formula is a single clause, as stated in the following proposition (due to Dan Suciu).

**PROPOSITION 6.1.** *Counting the number of true groundings of a first-order clause in a database is #P-complete in the length of the clause.*

**Proof.** Counting satisfying assignments of propositional monotone 2-CNF is #P-complete (Roth, 1996). This problem can be reduced to counting the number of true groundings of a first-order clause in a database as follows. Consider a database composed of the ground atoms  $R(0,1)$ ,  $R(1,0)$  and  $R(1,1)$ . Given a monotone 2-CNF formula, construct a formula  $\Phi$  that is a conjunction of predicates of the form  $R(x_i, x_j)$ , one for each disjunct  $x_i \vee x_j$  appearing in the CNF formula. (For example,  $(x_1 \vee x_2) \wedge (x_3 \vee x_4)$  would yield  $R(x_1, x_2) \wedge R(x_3, x_4)$ .) There is a one-to-one correspondence between the satisfying assignments of the 2-CNF and the true groundings of  $\Phi$ . The latter are the false groundings of the clause formed by disjoining the negations of all the  $R(x_i, x_j)$ , and thus can be counted by counting the number of true groundings of this clause and subtracting it from the total number of groundings.  $\square$

In large domains, the number of true groundings of a formula may be counted approximately, by uniformly sampling groundings of the formula and checking whether they are true in the data. In smaller domains, and in our experiments below, we use an efficient recursive algorithm to find the exact count.

A second problem with Equation 6 is that computing the expected number of true groundings is also intractable, requiring inference over the model. Further, efficient optimization methods also require computing the log-likelihood itself (Equation 3), and thus the partition function  $Z$ . This can be done approximately using a Monte Carlo maximum likelihood estimator (MC-MLE)



(Geyer & Thompson, 1992). However, in our experiments the Gibbs sampling used to compute the MC-MLEs and gradients did not converge in reasonable time, and using the samples from the unconverged chains yielded poor results.

A more efficient alternative, widely used in areas like spatial statistics, social network modeling and language processing, is to optimize instead the pseudo-likelihood (Besag, 1975)

$$P_w^*(X=x) = \prod_{l=1}^n P_w(X_l=x_l|MB_x(X_l)) \quad (7)$$

where  $MB_x(X_l)$  is the state of the Markov blanket of  $X_l$  in the data. The gradient of the pseudo-log-likelihood is

$$\frac{\partial}{\partial w_i} \log P_w^*(X=x) = \sum_{l=1}^n [n_i(x) - P_w(X_l=0|MB_x(X_l)) n_i(x_{[X_l=0]}) - P_w(X_l=1|MB_x(X_l)) n_i(x_{[X_l=1]})] \quad (8)$$

where  $n_i(x_{[X_l=0]})$  is the number of true groundings of the  $i$ th formula when we force  $X_l = 0$  and leave the remaining data unchanged, and similarly for  $n_i(x_{[X_l=1]})$ . Computing this expression (or Equation 7) does not require inference over the model. We optimize the pseudo-log-likelihood using the limited-memory BFGS algorithm (Liu & Nocedal, 1989). The computation can be made more efficient in several ways:

- The sum in Equation 8 can be greatly sped up by ignoring predicates that do not appear in the  $i$ th formula.
- The counts  $n_i(x)$ ,  $n_i(x_{[X_l=0]})$  and  $n_i(x_{[X_l=1]})$  do not change with the weights, and need only be computed once (as opposed to in every iteration of BFGS).
- Ground formulas whose truth value is unaffected by changing the truth value of any single literal may be ignored, since then  $n_i(x) = n_i(x_{[X_l=0]}) = n_i(x_{[X_l=1]})$ . In particular, this holds for any clause which contains at least two true literals. This can often be the great majority of ground clauses.

To combat overfitting, we penalize the pseudo-likelihood with a Gaussian prior on each weight.

Inductive logic programming (ILP) techniques can be used to learn additional clauses, refine the ones already in the MLN, or learn an MLN from scratch. We use the CLAUDIEN system for this purpose (De Raedt & Dehaspe, 1997). Unlike most other ILP systems, which learn only Horn clauses, CLAUDIEN is able to learn arbitrary first-order clauses, making it well suited

to MLNs. Also, by constructing a particular language bias, we are able to direct CLAUDIEN to search for refinements of the MLN structure. In the future we plan to more fully integrate structure learning into MLNs, by generalizing techniques like Della Pietra et al.'s (1997) to the first-order realm, as done by MACCENT for classification problems (Dehaspe, 1997).

## 7. Experiments

We tested MLNs using a database describing the Department of Computer Science and Engineering at the University of Washington (UW-CSE). The domain consists of 12 predicates and 2707 constants divided into 10 types. Types include: publication (342 constants), person (442), course (176), project (153), academic quarter (20), etc. Predicates include: Professor(person), Student(person), Area(x, area) (with x ranging over publications, persons, courses and projects), AuthorOf(publication, person), AdvisedBy(person, person), YearsInProgram(person, years), CourseLevel(course, level), TaughtBy(course, person, quarter), TeachingAssistant(course, person, quarter), etc. Additionally, there are 10 equality predicates: SamePerson(person, person), SameCourse(course, course), etc. which always have known, fixed values that are true iff the two arguments are the same constant.

Using typed variables, the total number of possible ground atoms ( $n$  in Section 6) was 4,106,841. The database contained a total of 3380 tuples (i.e., there were 3380 true ground atoms). We obtained this database by scraping pages in the department's Web site ([www.cs.washington.edu](http://www.cs.washington.edu)). Publications and AuthorOf relations were obtained by extracting from the Bib-Serv database ([www.bibserv.org](http://www.bibserv.org)) all records with author fields containing the names of at least two department members (in the form "last name, first name" or "last name, first initial").

We obtained a knowledge base by asking four volunteers to each provide a set of formulas in first-order logic describing the domain. (The volunteers were not shown the database of tuples, but were members of the department who thus had a general understanding about it.) Merging these yielded a KB of 96 formulas. The complete KB, volunteer instructions, database, and algorithm parameter settings are online at <http://www.cs.washington.edu/ai/mln>. Formulas in the KB include statements like: students are not professors; each student has at most one advisor; if a student is an author of a paper, so is her advisor; advanced students only TA courses taught by their advisors; at most one author of a given publication is a professor; students in Phase I of the Ph.D. program have no advisor; etc. Notice that these statements are not always true, but are typically true.

For training and testing purposes, we divided the database into five sub-databases, one for each area: AI, graphics, programming languages, systems, and theory. Professors and courses were manually assigned to areas, and other constants were iteratively assigned to the most frequent area among other constants they appeared in some tuple with. Each tuple was then assigned to the area of the constants in it. Tuples involving constants of more than one area were discarded, to avoid train-test contamination. The sub-databases contained, on average, 521 true ground atoms out of a possible 58457.

We performed leave-one-out testing by area, testing on each area in turn using the model trained from the remaining four. The test task was to predict the  $\text{AdvisedBy}(x, y)$  predicate given (a) all others (All Info) and (b) all others except  $\text{Student}(x)$  and  $\text{Professor}(x)$  (Partial Info). In both cases, we measured the average conditional log-likelihood of all possible groundings of  $\text{AdvisedBy}(x, y)$  over all areas, drew precision/recall curves, and computed the area under the curve. This task is an instance of link prediction, a problem that has been the object of much interest in statistical relational learning (see Section 8). All KBs were converted to clausal form. Timing results are on a 2.8Ghz Pentium 4 machine.

## 7.1. SYSTEMS

In order to evaluate MLNs, which use logic and probability for inference, we wished to compare with methods that use only logic or only probability. We were also interested in automatic induction of clauses using ILP techniques. This subsection gives details of the comparison systems used.

### 7.1.1. Logic

One important question we aimed to answer with the experiments is whether adding probability to a logical knowledge base improves its ability to model the domain. Doing this requires observing the results of answering queries using only logical inference, but this is complicated by the fact that computing log-likelihood and the area under the precision/recall curve requires real-valued probabilities, or at least some measure of “confidence” in the truth of each ground atom being tested. We thus used the following approach. For a given knowledge base  $KB$  and set of evidence atoms  $E$ , let  $\mathcal{X}_{KB \cup E}$  be the set of worlds that satisfy  $KB \cup E$ . The probability of a query atom  $q$  is then defined as  $P(q) = \frac{|\mathcal{X}_{KB \cup E \cup q}|}{|\mathcal{X}_{KB \cup E}|}$ , the fraction of  $\mathcal{X}_{KB \cup E}$  in which  $q$  is true.

A more serious problem arises if the KB is inconsistent (which was indeed the case with the KB we collected from volunteers). In this case the denominator of  $P(q)$  is zero. (Also, recall that an inconsistent knowledge base trivially entails any arbitrary formula). To address this, we redefine  $\mathcal{X}_{KB \cup E}$  to be the set of worlds which satisfies the maximum possible number of ground clauses. We use Gibbs sampling to sample from this set, with each chain

initialized to a mode using WalkSat. At each Gibbs step, the step is taken with probability: 1 if the new state satisfies more clauses than the current one (since that means the current state should have 0 probability), 0.5 if the new state satisfies the same number of clauses (since the new and old state then have equal probability), and 0 if the new state satisfies fewer clauses. We then use only the states with maximum number of satisfied clauses to compute probabilities. Notice that this is equivalent to using an MLN built from the KB and with all infinite equal weights.

### 7.1.2. Probability

The other question we wanted to answer with these experiments is whether existing (propositional) probabilistic models are already powerful enough to be used in relational domains without the need for the additional representational power provided by MLNs. In order to use such models, the domain must first be propositionalized by defining features that capture useful information about it. Creating good attributes for propositional learners in this highly relational domain is a difficult problem. Nevertheless, as a tradeoff between incorporating as much potentially relevant information as possible and avoiding extremely long feature vectors, we defined two sets of propositional attributes: order-1 and order-2. The former involves characteristics of individual constants in the query predicate, and the latter involves characteristics of relations between the constants in the query predicate.

For the order-1 attributes, we defined one variable for each  $(a, b)$  pair, where  $a$  is an argument of the query predicate and  $b$  is an argument of some predicate with the same value as  $a$ . The variable is the fraction of true groundings of this predicate in the data. Some examples of first-order attributes for `AdvisedBy(Matt, Pedro)` are: whether Pedro is a student, the fraction of publications that are published by Pedro, the fraction of courses for which Matt was a teaching assistant, etc.

The order-2 attributes were defined as follows: for a given (ground) query predicate  $Q(q_1, q_2, \dots, q_k)$ , consider all sets of  $k$  predicates and all assignments of constants  $q_1, q_2, \dots, q_k$  as arguments to the  $k$  predicates, with exactly one constant per predicate (in any order). For instance, if  $Q$  is `AdvisedBy(Matt, Pedro)` then one such possible set would be  $\{\text{TeachingAssistant}(\_, \text{Matt}, \_), \text{TaughtBy}(\_, \text{Pedro}, \_)\}$ . This forms  $2^k$  attributes of the example, each corresponding to a particular truth assignment to the  $k$  predicates. The value of an attribute is the number of times, in the training data, the set of predicates have that particular truth assignment, when their unassigned arguments are all filled with the same constants. For example, consider filling the above empty arguments with “CSE546” and “Autumn\_0304”. The resulting set,  $\{\text{TeachingAssistant}(\text{CSE546}, \text{Matt}, \text{Autumn\_0304}), \text{TaughtBy}(\text{CSE546}, \text{Pedro}, \text{Autumn\_0304})\}$  has some truth assignment in the training data (e.g.,  $\{\text{True}, \text{True}\}$ ,  $\{\text{True}, \text{False}\}$ ,  $\dots$ ). One attribute is the number of

such sets of constants that create the truth assignment  $\{\text{True}, \text{True}\}$ , another for  $\{\text{True}, \text{False}\}$  and so on. Some examples of second-order attributes generated for the query `AdvisedBy(Matt, Pedro)` are: how often `Matt` is a teaching assistant for a course that `Pedro` taught (as well as how often he is not), how many publications `Pedro` and `Matt` have coauthored, etc.

The resulting 28 order-1 attributes and 120 order-2 attributes (for the All Info case) were discretized into five equal-frequency bins (based on the training set). We used two propositional learners: Naive Bayes (Domingos & Elkan, 1997) and Bayesian networks (Heckerman et al., 1995) with structure and parameters learned using the VFBN2 algorithm (Hulten & Domingos, 2002) with a maximum of four parents per node. The order-2 attributes helped the naive Bayes classifier but hurt the performance of the Bayesian network classifier, so below we report results using the order-1 and order-2 attributes for naive Bayes, and only the order-1 attributes for Bayesian networks.

### 7.1.3. *Inductive logic programming*

Our original knowledge base was acquired from volunteers, but we were also interested in whether it could have been developed automatically using inductive logic programming methods. As mentioned earlier, we used CLAUDIEN to induce a knowledge base from data. CLAUDIEN was run with: local scope; minimum accuracy of 0.1; minimum coverage of 1; maximum complexity of 10; and breadth-first search. CLAUDIEN's search space is defined by its language bias. We constructed a language bias which allowed: a maximum of 3 variables in a clause; unlimited predicates in a clause; up to 2 non-negated appearances of a predicate in a clause, and 2 negated ones; and use of knowledge of predicate argument types. To minimize search, the equality predicates (e.g., `SamePerson`) were not used in CLAUDIEN, and this improved its results.

Besides inducing clauses from the training data, we were also interested in using data to automatically refine the knowledge base provided by our volunteers. CLAUDIEN does not support this feature directly, but it can be emulated by an appropriately constructed language bias. We did this by, for each clause in the KB, allowing CLAUDIEN to (1) remove any number of the literals, (2) add up to  $v$  new variables, and (3) add up to  $l$  new literals. We ran CLAUDIEN for 24 hours on a Sun-Blade 1000 for each  $(v, l)$  in the set  $\{(1, 2), (2, 3), (3, 4)\}$ . All three gave nearly identical results; we report the results with  $v = 3$  and  $l = 4$ .

### 7.1.4. *MLNs*

Our results compare the above systems to Markov logic networks. The MLNs were trained using a Gaussian weight prior with zero mean and unit variance, and with the weights initialized at the mode of the prior (zero). For optimization, we used the Fortran implementation of L-BFGS from Zhu et al. (1997)

and Byrd et al. (1995), leaving all parameters at their default values, and with a convergence criterion (*ftol*) of  $10^{-5}$ . Inference was performed using Gibbs sampling as described in Section 5, with ten parallel Markov chains, each initialized to a mode of the distribution using MaxWalkSat. The number of Gibbs steps was determined using the criterion of DeGroot and Schervish (2002, pp. 707 and 740-741). Sampling continued until we reached a confidence of 95% that the probability estimate was within 1% of the true value in at least 95% of the nodes (ignoring nodes which are always true or false). A minimum of 1000 and maximum of 500,000 samples was used, with one sample per complete Gibbs pass through the variables. Typically, inference converged within 5000 to 100,000 passes. The results were insensitive to variation in the convergence thresholds.

## 7.2. RESULTS

### 7.2.1. Training with MC-MLE

Our initial system used MC-MLE to train MLNs, with ten Gibbs chains, and each ground atom being initialized to true with the corresponding first-order predicate's probability of being true in the data. Gibbs steps may be taken quite quickly by noting that few counts of satisfied clauses will change on any given step. On the UW-CSE domain, our implementation took 4-5 ms per step. We used the maximum across all predicates of the Gelman criterion  $R$  (Gilks et al., 1996) to determine when the chains had reached their stationary distribution. In order to speed convergence, our Gibbs sampler preferentially samples atoms that were true in either the data or the initial state of the chain. The intuition behind this is that most atoms are always false, and sampling repeatedly from them is inefficient. This improved convergence by approximately an order of magnitude over uniform selection of atoms. Despite these optimizations, the Gibbs sampler took a prohibitively long time to reach a reasonable convergence threshold (e.g.,  $R = 1.01$ ). After running for 24 hours (approximately 2 million Gibbs steps per chain), the average  $R$  value across training sets was 3.04, with no one training set having reached an  $R$  value less than 2 (other than briefly dipping to 1.5 in the early stages of the process). Considering this must be done iteratively as L-BFGS searches for the minimum, we estimate it would take anywhere from 20 to 400 days to complete the training, even with a weak convergence threshold such as  $R = 2.0$ . Experiments confirmed the poor quality of the models that resulted if we ignored the convergence threshold and limited the training process to less than ten hours. With a better choice of initial state, approximate counting, and improved MCMC techniques such as the Swendsen-Wang algorithm (Edwards & Sokal, 1988), MC-MLE may become practical, but it is not a viable option for training in the current version. (Notice that during learning

MCMC is performed over the full ground network, which is too large to apply MaxWalkSat to.)

### 7.2.2. *Training with pseudo-likelihood*

In contrast to MC-MLE, pseudo-likelihood training was quite fast. As discussed in Section 6, each iteration of training may be done quite quickly once the initial clause and ground atom satisfiability counts are complete. On average (over the five test sets), finding these counts took 2.5 minutes. From there, training took, on average, 255 iterations of L-BFGS, for a total of 16 minutes.

### 7.2.3. *Inference*

Inference was also quite quick. Inferring the probability of all `AdvisedBy(x, y)` atoms in the All Info case took 3.3 minutes in the AI test set (4624 atoms), 24.4 in graphics (3721), 1.8 in programming languages (784), 10.4 in systems (5476), and 1.6 in theory (2704). The number of Gibbs passes ranged from 4270 to 500,000, and averaged 124,000. This amounts to 18 ms per Gibbs pass and approximately 200,000–500,000 Gibbs steps per second. The average time to perform inference in the Partial Info case was 14.8 minutes (vs. 8.3 in the All Info case).

### 7.2.4. *Comparison of systems*

We compared twelve systems: the original KB (KB); CLAUDIEN (CL); CLAUDIEN with the original KB as language bias (CLB); the union of the original KB and CLAUDIEN’s output in both cases (KB+CL and KB+CLB); an MLN with each of the above KBs (MLN(KB), MLN(CL), MLN(KB+CL), and MLN(KB+CLB)); naive Bayes (NB); and a Bayesian network learner (BN). Add-one smoothing of probabilities was used in all cases.

Table IV summarizes the results. Figure 2 shows precision/recall curves for all areas (i.e., averaged over all `AdvisedBy(x, y)` atoms), and Figures 3 to 7 show precision/recall curves for the five individual areas. MLNs are clearly more accurate than the alternatives, showing the promise of this approach. The purely logical and purely probabilistic methods often suffer when intermediate predicates have to be inferred, while MLNs are largely unaffected. Naive Bayes performs well in AUC in some test sets, but very poorly in others; its CLLs are uniformly poor. CLAUDIEN performs poorly on its own, and produces no improvement when added to the KB in the MLN. Using CLAUDIEN to refine the KB typically performs worse in AUC but better in CLL than using CLAUDIEN from scratch; overall, the best-performing logical method is KB+CLB, but its results fall well short of the best MLNs’. The general drop-off in precision around 50% recall is attributable to the fact that the database is very incomplete, and only allows identifying a minority of the `AdvisedBy` relations. Inspection reveals that the occasional smaller

Table IV. Experimental results for predicting `AdvisedBy(x, y)` when all other predicates are known (All Info) and when `Student(x)` and `Professor(x)` are unknown (Partial Info). CLL is the average conditional log-likelihood, and AUC is the area under the precision-recall curve. The results are averages over all atoms in the five test sets and their standard deviations. (See <http://www.cs.washington.edu/ai/mln> for details on how the standard deviations of the AUCs were computed.)

System	All Info		Partial Info	
	AUC	CLL	AUC	CLL
MLN(KB)	0.215±0.0172	-0.052±0.004	0.224±0.0185	-0.048±0.004
MLN(KB+CL)	0.152±0.0165	-0.058±0.005	0.203±0.0196	-0.045±0.004
MLN(KB+CLB)	0.011±0.0003	-3.905±0.048	0.011±0.0003	-3.958±0.048
MLN(CL)	0.035±0.0008	-2.315±0.030	0.032±0.0009	-2.478±0.030
MLN(CLB)	0.003±0.0000	-0.052±0.005	0.023±0.0003	-0.338±0.002
KB	0.059±0.0081	-0.135±0.005	0.048±0.0058	-0.063±0.004
KB+CL	0.037±0.0012	-0.202±0.008	0.028±0.0012	-0.122±0.006
KB+CLB	0.084±0.0100	-0.056±0.004	0.044±0.0064	-0.051±0.005
CL	0.048±0.0009	-0.434±0.012	0.037±0.0001	-0.836±0.017
CLB	0.003±0.0000	-0.052±0.005	0.010±0.0001	-0.598±0.003
NB	0.054±0.0006	-1.214±0.036	0.044±0.0009	-1.140±0.031
BN	0.015±0.0006	-0.072±0.003	0.015±0.0007	-0.215±0.003

drop-offs in precision at very low recalls are due to students who graduated or changed advisors after co-authoring many publications with them.

## 8. Statistical Relational Learning Tasks

Many SRL tasks can be concisely formulated in the language of MLNs, allowing the algorithms introduced in this paper to be directly applied to them. In this section we exemplify this with five key tasks: collective classification, link prediction, link-based clustering, social network modeling, and object identification.

### 8.1. COLLECTIVE CLASSIFICATION

The goal of ordinary classification is to predict the class of an object given its attributes. Collective classification also takes into account the classes of related objects (e.g., Chakrabarti et al. (1998); Taskar et al. (2002); Neville and Jensen (2003)). Attributes can be represented in MLNs as predicates of the form  $A(x, v)$ , where  $A$  is an attribute,  $x$  is an object, and  $v$  is the value of  $A$  in  $x$ . The class is a designated attribute  $C$ , representable by  $C(x, v)$ , where



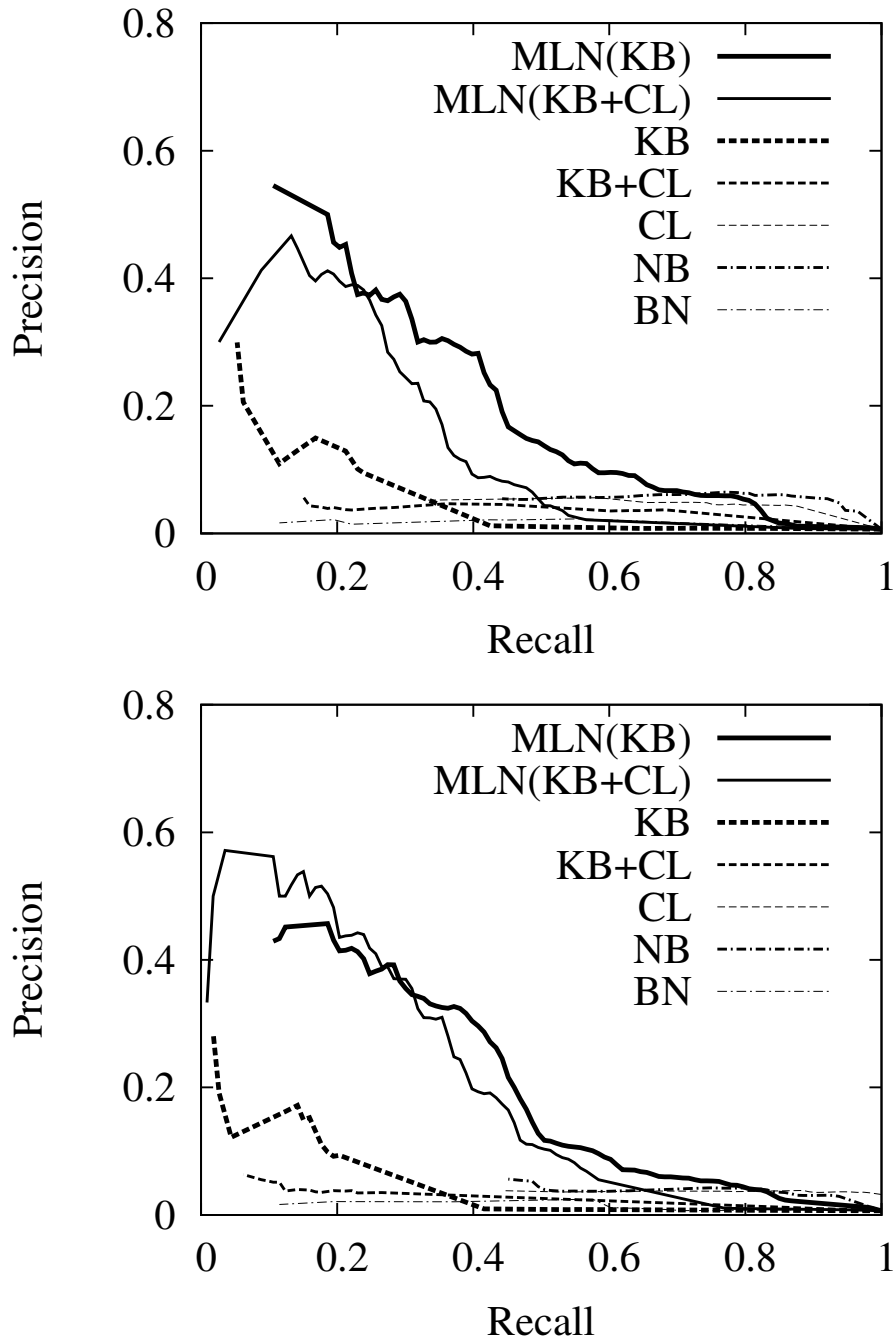


Figure 2. Precision and recall for all areas: All Info (upper graph) and Partial Info (lower graph).

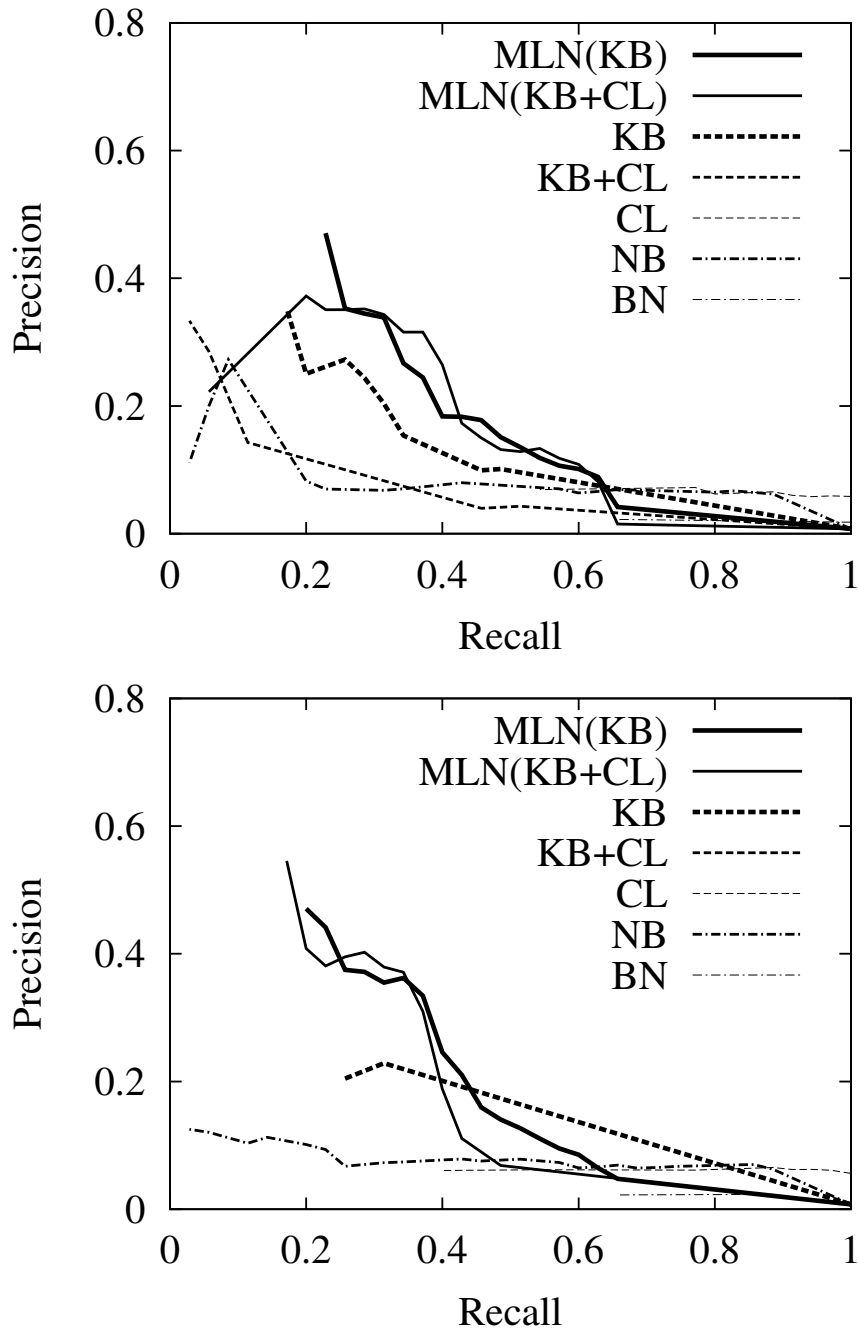


Figure 3. Precision and recall for the AI area: All Info (upper graph) and Partial Info (lower graph).

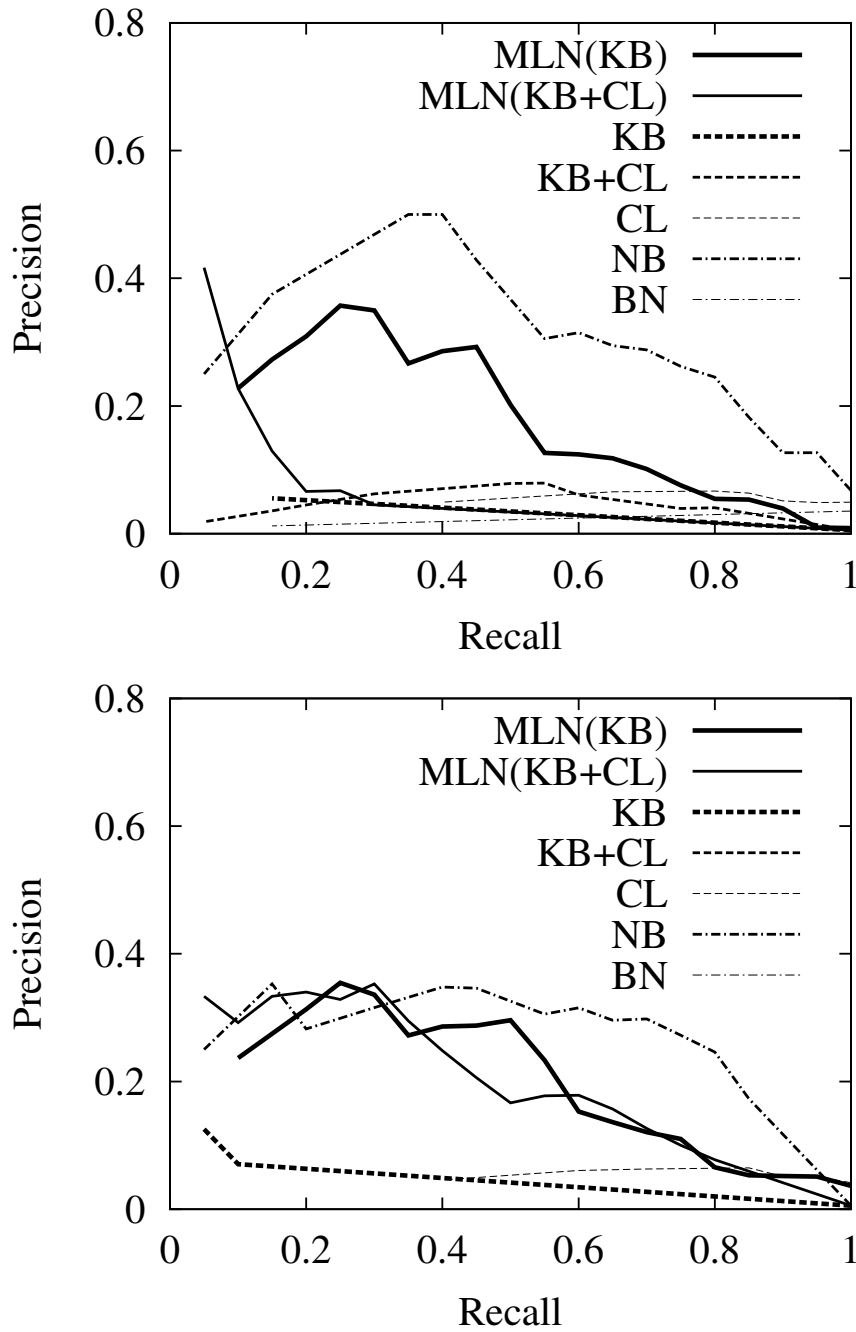


Figure 4. Precision and recall for the graphics area: All Info (upper graph) and Partial Info (lower graph).

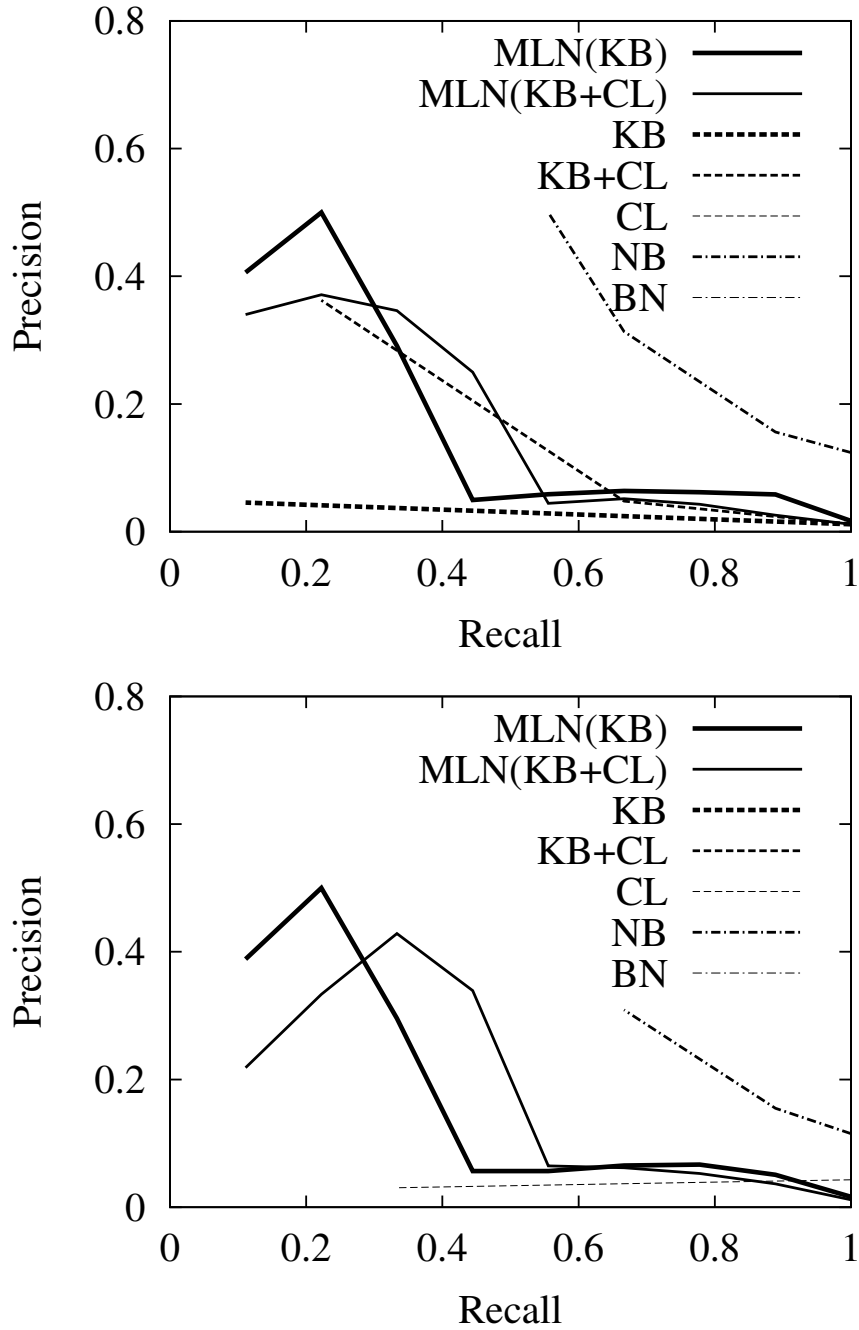


Figure 5. Precision and recall for the programming languages area: All Info (upper graph) and Partial Info (lower graph).

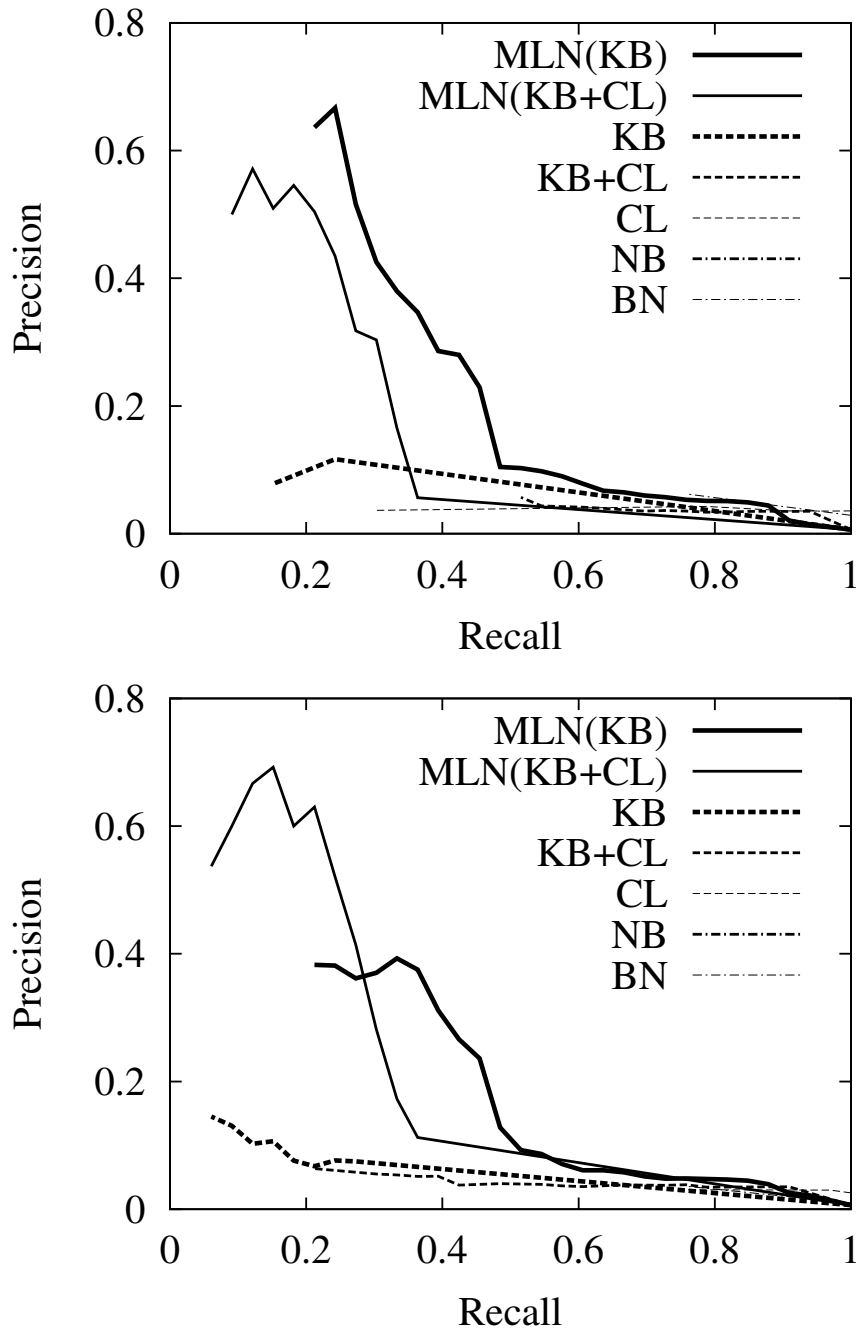


Figure 6. Precision and recall for the systems area: All Info (upper graph) and Partial Info (lower graph). The curves for naive Bayes are indistinguishable from the X axis.

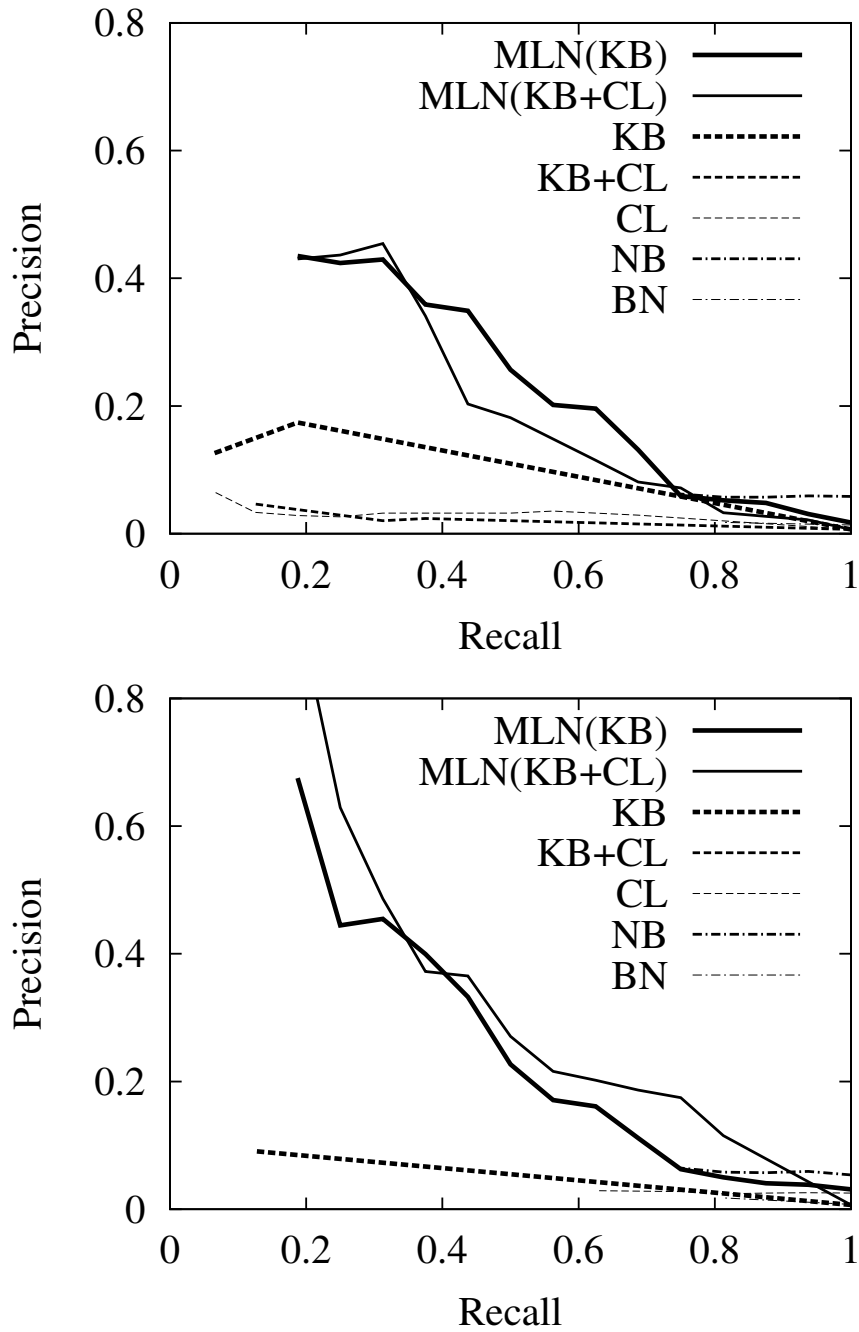


Figure 7. Precision and recall for the theory area: All Info (upper graph) and Partial Info (lower graph).

$v$  is  $x$ 's class. Classification is now simply the problem of inferring the truth value of  $C(x, v)$  for all  $x$  and  $v$  of interest given all known  $A(x, v)$ . Ordinary classification is the special case where  $C(x_i, v)$  and  $C(x_j, v)$  are independent for all  $x_i$  and  $x_j$  given the known  $A(x, v)$ . In collective classification, the Markov blanket of  $C(x_i, v)$  includes other  $C(x_j, v)$ , even after conditioning on the known  $A(x, v)$ . Relations between objects are represented by predicates of the form  $R(x_i, x_j)$ . A number of interesting generalizations are readily apparent, for example  $C(x_i, v)$  and  $C(x_j, v)$  may be indirectly dependent via unknown predicates, possibly including the  $R(x_i, x_j)$  predicates themselves.

## 8.2. LINK PREDICTION

The goal of link prediction is to determine whether a relation exists between two objects of interest (e.g., whether Anna is Bob's Ph.D. advisor) from the properties of those objects and possibly other known relations (e.g., Popescul and Ungar (2003)). The formulation of this problem in MLNs is identical to that of collective classification, with the only difference that the goal is now to infer the value of  $R(x_i, x_j)$  for all object pairs of interest, instead of  $C(x, v)$ . The task used in our experiments was an example of link prediction.

## 8.3. LINK-BASED CLUSTERING

The goal of clustering is to group together objects with similar attributes. In model-based clustering, we assume a generative model  $P(X) = \sum_C P(C) P(X|C)$ , where  $X$  is an object,  $C$  ranges over clusters, and  $P(C|X)$  is  $X$ 's degree of membership in cluster  $C$ . In link-based clustering, objects are clustered according to their links (e.g., objects that are more closely related are more likely to belong to the same cluster), and possibly according to their attributes as well (e.g., Flake et al. (2000)). This problem can be formulated in MLNs by postulating an unobserved predicate  $C(x, v)$  with the meaning " $x$  belongs to cluster  $v$ ," and having formulas in the MLN involving this predicate and the observed ones (e.g.,  $R(x_i, x_j)$  for links and  $A(x, v)$  for attributes). Link-based clustering can now be performed by learning the parameters of the MLN, and cluster memberships are given by the probabilities of the  $C(x, v)$  atoms conditioned on the observed ones.

## 8.4. SOCIAL NETWORK MODELING

Social networks are graphs where nodes represent social actors (e.g., people) and arcs represent relations between them (e.g., friendship). Social network analysis (Wasserman & Faust, 1994) is concerned with building models relating actors' properties and their links. For example, the probability of two actors forming a link may depend on the similarity of their attributes, and conversely two linked actors may be more likely to have certain properties.

These models are typically Markov networks, and can be concisely represented by formulas like  $\forall x \forall y \forall v R(x, y) \Rightarrow (A(x, v) \Leftrightarrow A(y, v))$ , where  $x$  and  $y$  are actors,  $R(x, y)$  is a relation between them,  $A(x, v)$  represents an attribute of  $x$ , and the weight of the formula captures the strength of the correlation between the relation and the attribute similarity. For example, a model stating that friends tend to have similar smoking habits can be represented by the formula  $\forall x \forall y \text{Friends}(x, y) \Rightarrow (\text{Smokes}(x) \Leftrightarrow \text{Smokes}(y))$  (Table I). As well as encompassing existing social network models, MLNs allow richer ones to be easily stated (e.g., by writing formulas involving multiple types of relations and multiple attributes, as well as more complex dependencies between them).

## 8.5. OBJECT IDENTIFICATION

Object identification (also known as record linkage, de-duplication, and others) is the problem of determining which records in a database refer to the same real-world entity (e.g., which entries in a bibliographic database represent the same publication) (Winkler, 1999). This problem is of crucial importance to many companies, government agencies, and large-scale scientific projects. One way to represent it in MLNs is by removing the unique names assumption as described in Section 4, i.e., by defining a predicate  $\text{Equals}(x, y)$  (or  $x = y$  for short) with the meaning “ $x$  represents the same real-world entity as  $y$ .” This predicate is applied both to records and their fields (e.g., “ICML” = “Intl. Conf. on Mach. Learn.”). The dependencies between record matches and field matches can then be represented by formulas like  $\forall x \forall y x = y \Leftrightarrow f_i(x) = f_i(y)$ , where  $x$  and  $y$  are records and  $f_i(x)$  is a function returning the value of the  $i$ th field of record  $x$ . We have successfully applied this approach to de-duplicating the Cora database of computer science papers (Parag & Domingos, 2004). Because it allows information to propagate from one match decision (i.e., one grounding of  $x = y$ ) to another via fields that appear in both pairs of records, it effectively performs collective object identification, and in our experiments outperformed the traditional method of making each match decision independently of all others. For example, matching two references may allow us to determine that “ICML” and “MLC” represent the same conference, which in turn may help us to match another pair of references where one contains “ICML” and the other “MLC.” MLNs also allow additional information to be incorporated into a de-duplication system easily, modularly and uniformly. For example, transitive closure is incorporated by adding the formula  $\forall x \forall y \forall z x = y \wedge y = z \Rightarrow x = z$ , with a weight that can be learned from data.



## 9. Related Work

There is a very large literature relating logic and probability; here we will focus only on the approaches most relevant to statistical relational learning, and discuss how they relate to MLNs.

### 9.1. EARLY WORK

Attempts to combine logic and probability in AI date back to at least Nilsson (1986). Bacchus (1990), Halpern (1990) and coworkers (e.g., Bacchus et al. (1996)) studied the problem in detail from a theoretical standpoint. They made a distinction between statistical statements (e.g., “65% of the students in our department are undergraduate”) and statements about possible worlds (e.g., “The probability that Anna is an undergraduate is 65%”), and provided methods for computing the latter from the former. In their approach, a KB did not specify a complete and unique distribution over possible worlds, requiring additional assumptions to obtain one. Bacchus et al. considered a number of alternatives, all of them quite restrictive (e.g., all worlds compatible with the KB should be equally likely). In contrast, by viewing KBs as Markov network templates, MLNs can represent arbitrary distributions over possible worlds.

Paskin (2002) extended the work of Bacchus et al. by associating a probability with each first-order formula, and taking the maximum entropy distribution compatible with those probabilities. This representation was still quite brittle, with a world that violates a single grounding of a universally quantified formula being considered as unlikely as a world that violates all of them. In contrast, in MLNs a rule like  $\forall x \text{Smokes}(x) \Rightarrow \text{Cancer}(x)$  causes the probability of a world to decrease gradually as the number of cancer-free smokers in it increases.

### 9.2. KNOWLEDGE-BASED MODEL CONSTRUCTION

Knowledge-based model construction (KBMC) is a combination of logic programming and Bayesian networks (Wellman et al., 1992; Ngo & Haddawy, 1997; Kersting & De Raedt, 2001). As in MLNs, nodes in KBMC represent ground atoms. Given a Horn KB, KBMC answers a query by finding all possible backward-chaining proofs of the query and evidence atoms from each other, constructing a Bayesian network over all atoms in the proofs, and performing inference over this network. The parents of an atom in the network are deterministic AND nodes representing the bodies of the clauses that have that node as head. The conditional probability of the node given these is specified by a combination function (e.g., noisy OR, logistic regression, arbitrary CPT). MLNs have several advantages compared to KBMC: they allow arbitrary formulas (not just Horn clauses) and inference in any direction, they sidestep the thorny problem of avoiding cycles in the Bayesian

networks constructed by KBMC, and they do not require the introduction of *ad hoc* combination functions for clauses with the same consequent.

A KBMC model can be translated into an MLN by writing down a set of formulas for each first-order predicate  $P_k(\dots)$  in the domain. Each formula is a conjunction containing  $P_k(\dots)$  and one literal per parent of  $P_k(\dots)$  (i.e., per first-order predicate appearing in a Horn clause having  $P_k(\dots)$  as the consequent). A subset of these literals are negated; there is one formula for each possible combination of positive and negative literals. The weight of the formula is  $w = \log[p/(1-p)]$ , where  $p$  is the conditional probability of the child predicate when the corresponding conjunction of parent literals is true, according to the combination function used. If the combination function is logistic regression, it can be represented using only a linear number of formulas, taking advantage of the fact that a logistic regression model is a (conditional) Markov network with a binary clique between each predictor and the response. Noisy OR can similarly be represented with a linear number of parents.

### 9.3. OTHER LOGIC PROGRAMMING APPROACHES

Stochastic logic programs (SLPs) (Muggleton, 1996; Cussens, 1999) are a combination of logic programming and log-linear models. Puech and Muggleton (2003) showed that SLPs are a special case of KBMC, and thus they can be converted into MLNs in the same way. Like MLNs, SLPs have one coefficient per clause, but they represent distributions over Prolog proof trees rather than over predicates; the latter have to be obtained by marginalization. Similar remarks apply to a number of other representations that are essentially equivalent to SLPs, like independent choice logic (Poole, 1993) and PRISM (Sato & Kameya, 1997).

MACCENT (Dehaspe, 1997) is a system that learns log-linear models with first-order features; each feature is a conjunction of a class and a Prolog query (clause with empty head). A key difference between MACCENT and MLNs is that MACCENT is a classification system (i.e., it predicts the conditional distribution of an object's class given its properties), while an MLN represents the full joint distribution of a set of predicates. Like any probability estimation approach, MLNs can be used for classification simply by issuing the appropriate conditional queries.<sup>5</sup> In particular, a MACCENT model can be converted into an MLN simply by defining a class predicate (as in Subsection 8.1), adding the corresponding features and their weights to the MLN, and adding a formula with infinite weight stating that each object must have exactly one class. (This fails to model the marginal distribution of the non-class predicates, which is not a problem if only classification

<sup>5</sup> Conversely, joint distributions can be built up from classifiers (e.g., (Heckerman et al., 2000)), but this would be a significant extension of MACCENT.

queries will be issued.) MACCENT can make use of deterministic background knowledge in the form of Prolog clauses; these can be added to the MLN as formulas with infinite weight. In addition, MLNs allow uncertain background knowledge (via formulas with finite weights). As described in Subsection 8.1, MLNs can be used for collective classification, where the classes of different objects can depend on each other; MACCENT, which requires that each object be represented in a separate Prolog knowledge base, does not have this capability.

Constraint logic programming (CLP) is an extension of logic programming where variables are constrained instead of being bound to specific values during inference (Laffar & Lassez, 1987). Probabilistic CLP generalizes SLPs to CLP (Riezler, 1998), and  $\text{CLP}(\mathcal{BN})$  combines CLP with Bayesian networks (Santos Costa et al., 2003). Unlike in MLNs, constraints in  $\text{CLP}(\mathcal{BN})$  are hard (i.e., they cannot be violated; rather, they define the form of the probability distribution).

#### 9.4. PROBABILISTIC RELATIONAL MODELS

Probabilistic relational models (PRMs) (Friedman et al., 1999) are a combination of frame-based systems and Bayesian networks. PRMs can be converted into MLNs by defining a predicate  $S(x, v)$  for each (propositional or relational) attribute of each class, where  $S(x, v)$  means “The value of attribute  $S$  in object  $x$  is  $v$ .” A PRM is then translated into an MLN by writing down a formula for each line of each (class-level) conditional probability table (CPT) and value of the child attribute. The formula is a conjunction of literals stating the parent values and a literal stating the child value, and its weight is the logarithm of  $P(x|Parents(x))$ , the corresponding entry in the CPT. In addition, the MLN contains formulas with infinite weight stating that each attribute must take exactly one value. This approach handles all types of uncertainty in PRMs (attribute, reference and existence uncertainty).

As Taskar et al. (2002) point out, the need to avoid cycles in PRMs causes significant representational and computational difficulties. Inference in PRMs is done by creating the complete ground network, which limits their scalability. PRMs require specifying a complete conditional model for each attribute of each class, which in large complex domains can be quite burdensome. In contrast, MLNs create a complete joint distribution from whatever number of first-order features the user chooses to specify.

#### 9.5. RELATIONAL MARKOV NETWORKS

Relational Markov networks (RMNs) use database queries as clique templates, and have a feature for each state of a clique (Taskar et al., 2002). MLNs generalize RMNs by providing a more powerful language for constructing features (first-order logic instead of conjunctive queries), and by

allowing uncertainty over arbitrary relations (not just attributes of individual objects). RMNs are exponential in clique size, while MLNs allow the user (or learner) to determine the number of features, making it possible to scale to much larger clique sizes. RMNs are trained discriminatively, and do not specify a complete joint distribution for the variables in the model. Discriminative training of MLNs is straightforward (in fact, easier than the generative training used in this paper), and we have carried out successful preliminary experiments using a voted perceptron algorithm (Collins, 2002). RMNs use MAP estimation with belief propagation for inference, which makes learning quite slow, despite the simplified discriminative setting; maximizing the pseudo-likelihood of the query variables may be a more effective alternative.

#### 9.6. STRUCTURAL LOGISTIC REGRESSION

In structural logistic regression (SLR) (Popescul & Ungar, 2003), the predictors are the output of SQL queries over the input data. Just as a logistic regression model is a discriminatively-trained Markov network, an SLR model is a discriminatively-trained MLN.<sup>6</sup>

#### 9.7. RELATIONAL DEPENDENCY NETWORKS

In a relational dependency network (RDN), each node's probability conditioned on its Markov blanket is given by a decision tree (Neville & Jensen, 2003). Every RDN has a corresponding MLN in the same way that every dependency network has a corresponding Markov network, given by the stationary distribution of a Gibbs sampler operating on it (Heckerman et al., 2000).

#### 9.8. PLATES AND PROBABILISTIC ER MODELS

Large graphical models with repeated structure are often compactly represented using plates (Buntine, 1994). MLNs subsume plates as a representation language. In addition, they allow individuals and their relations to be explicitly represented (see Cussens (2003)), and context-specific independencies to be compactly written down, instead of left implicit in the node models. More recently, Heckerman et al. (2004) have proposed a language based on entity-relationship models that combines the features of plates and PRMs; this language is a special case of MLNs in the same way that ER models are a special case of logic. Probabilistic ER models allow logical expressions as constraints on how ground networks are constructed, but the truth values of these expressions have to be known in advance; MLNs allow uncertainty over all logical expressions.

---

<sup>6</sup> Use of SQL aggregates requires that their definitions be imported into the MLN.

## 9.9. BLOG

Milch et al. (2004) have proposed a language, called BLOG, designed to avoid making the unique names and domain closure assumptions. A BLOG program specifies procedurally how to generate a possible world, and does not allow arbitrary first-order knowledge to be easily incorporated. Also, it only specifies the structure of the model, leaving the parameters to be specified by external calls. BLOG models are directed graphs and need to avoid cycles, which substantially complicates their design. We saw in Section 4 how to remove the unique names and domain closure assumptions in MLNs. (When there are unknown objects of multiple types, a random variable for the number of each type is introduced.) Inference about an object's attributes, rather than those of its observations, can be done simply by having variables for objects as well as for their observations (e.g., for books as well as citations to them). To our knowledge, BLOG has not yet been implemented and evaluated.

## 9.10. OTHER WORK

There are many more approaches to statistical relational learning than we can possibly cover here. This section briefly considers some additional works that are potentially relevant to MLNs.

Pasula and Russell (2001), Poole (2003) and Sanghai et al. (2003) have studied efficient inference in first-order probabilistic models. While they focus on directed graphical models, some of the ideas (e.g., different MCMC steps for different types of predicates, combining unification with variable elimination, abstraction hierarchies) may be applicable to MLNs.

MLNs have some interesting similarities with the KBANN system, which converts a propositional Horn KB into a neural network and uses backpropagation to learn the network's weights (Towell & Shavlik, 1994). More generally, MLNs can be viewed as an extension to probability estimation of a long line of work on knowledge-intensive learning (e.g., Bergadano and Giordana (1988); Pazzani and Kibler (1992); Ourston and Mooney (1994)).

## 10. Future Work

MLNs are potentially a tool of choice for many AI problems, but much remains to be done. Directions for future work fall into three main areas:

**Inference:** We plan to develop more efficient forms of MCMC for MLNs, study the use of belief propagation, identify and exploit useful special cases, and investigate the possibility of lifted inference.

**Learning:** We plan to develop algorithms for learning and revising the structure of MLNs by directly optimizing (pseudo) likelihood, study alternate approaches to weight learning, train MLNs discriminatively, learn MLNs from incomplete data, use MLNs for link-based clustering, and develop methods for probabilistic predicate discovery.

**Applications:** We would like to apply MLNs in a variety of domains, including information extraction and integration, natural language processing, vision, social network analysis, computational biology, etc.

## 11. Conclusion

Markov logic networks (MLNs) are a simple way to combine probability and first-order logic in finite domains. An MLN is obtained by attaching weights to the formulas (or clauses) in a first-order knowledge base, and can be viewed as a template for constructing ordinary Markov networks. Each possible grounding of a formula in the KB yields a feature in the constructed network. Inference is performed by grounding the minimal subset of the network required for answering the query and running a Gibbs sampler over this subnetwork, with initial states found by MaxWalkSat. Weights are learned by optimizing a pseudo-likelihood measure using the L-BFGS algorithm, and clauses are learned using the CLAUDIEN system. Empirical tests with real-world data and knowledge in a university domain illustrate the promise of MLNs. Source code for learning and inference in MLNs will be made available at <http://www.cs.washington.edu/ai/mln>.

## Acknowledgements

We are grateful to Julian Besag, James Cussens, Nilesch Dalvi, Alon Halevy, Mark Handcock, Henry Kautz, Kristian Kersting, Tian Sang, Bart Selman, Dan Suci, Jeremy Tantrum, and Wei Wei for helpful discussions. This research was partly supported by ONR grant N00014-02-1-0408 and by a Sloan Fellowship awarded to the second author. We used the VFML library in our experiments (<http://www.cs.washington.edu/dm/vfml/>).

## References

Bacchus, F. (1990). *Representing and reasoning with probabilistic knowledge*. Cambridge, MA: MIT Press.

- Bacchus, F., Grove, A. J., Halpern, J. Y., & Koller, D. (1996). From statistical knowledge bases to degrees of belief. *Artificial Intelligence*, 87, 75–143.
- Bergadano, F., & Giordana, A. (1988). A knowledge-intensive approach to concept induction. *Proceedings of the Fifth International Conference on Machine Learning* (pp. 305–317). Ann Arbor, MI: Morgan Kaufmann.
- Berners-Lee, T., Hendler, J., & Lassila, O. (2001). The Semantic Web. *Scientific American*, 284 (5), 34–43.
- Besag, J. (1975). Statistical analysis of non-lattice data. *The Statistician*, 24, 179–195.
- Buntine, W. (1994). Operations for learning with graphical models. *Journal of Artificial Intelligence Research*, 2, 159–225.
- Byrd, R. H., Lu, P., & Nocedal, J. (1995). A limited memory algorithm for bound constrained optimization. *SIAM Journal on Scientific and Statistical Computing*, 16, 1190–1208.
- Chakrabarti, S., Dom, B., & Indyk, P. (1998). Enhanced hypertext categorization using hyperlinks. *Proceedings of the 1998 ACM SIGMOD International Conference on Management of Data* (pp. 307–318). Seattle, WA: ACM Press.
- Collins, M. (2002). Discriminative training methods for hidden Markov models: Theory and experiments with perceptron algorithms. *Proceedings of the 2002 Conference on Empirical Methods in Natural Language Processing*. Philadelphia, PA.
- Cumby, C., & Roth, D. (2003). Feature extraction languages for propositionalized relational learning. *Proceedings of the IJCAI-2003 Workshop on Learning Statistical Models from Relational Data* (pp. 24–31). Acapulco, Mexico: IJCAI.
- Cussens, J. (1999). Loglinear models for first-order probabilistic reasoning. *Proceedings of the Fifteenth Conference on Uncertainty in Artificial Intelligence* (pp. 126–133). Stockholm, Sweden: Morgan Kaufmann.
- Cussens, J. (2003). Individuals, relations and structures in probabilistic models. *Proceedings of the IJCAI-2003 Workshop on Learning Statistical Models from Relational Data* (pp. 32–36). Acapulco, Mexico: IJCAI.
- De Raedt, L., & Dehaspe, L. (1997). Clausal discovery. *Machine Learning*, 26, 99–146.
- DeGroot, M. H., & Schervish, M. J. (2002). *Probability and statistics*. Boston, MA: Addison Wesley, 3rd edition.
- Dehaspe, L. (1997). Maximum entropy modeling with clausal constraints. *Proceedings of the Seventh International Workshop on Inductive Logic Programming* (pp. 109–125). Prague, Czech Republic: Springer.
- Della Pietra, S., Della Pietra, V., & Lafferty, J. (1997). Inducing features of random fields. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 19, 380–392.
- Dietterich, T., Getoor, L., & Murphy, K. (Eds.). (2003). *Proceedings of the ICML-2004 Workshop on Statistical Relational Learning and its Connections to Other Fields*. Banff, Canada: IMLS.

- Domingos, P., & Pazzani, M. (1997). On the optimality of the simple Bayesian classifier under zero-one loss. *Machine Learning*, 29, 103–130.
- Džeroski, S., & Blockeel, H. (Eds.). (2004). *Proceedings of the Third International Workshop on Multi-Relational Data Mining*. Seattle, WA: ACM Press.
- Džeroski, S., & De Raedt, L. (2003). Special issue on multi-relational data mining: The current frontiers. *SIGKDD Explorations*, 5.
- Džeroski, S., De Raedt, L., & Wrobel, S. (Eds.). (2002). *Proceedings of the First International Workshop on Multi-Relational Data Mining*. Edmonton, Canada: ACM Press.
- Džeroski, S., De Raedt, L., & Wrobel, S. (Eds.). (2003). *Proceedings of the Second International Workshop on Multi-Relational Data Mining*. Washington, DC: ACM Press.
- Edwards, R., & Sokal, A. (1988). Generalization of the Fortuin-Kasteleyn-Swendsen-Wang representation and Monte Carlo algorithm. *Physics Review D* (pp. 2009–2012).
- Flake, G. W., Lawrence, S., & Giles, C. L. (2000). Efficient identification of Web communities. *Proceedings of the Sixth ACM SIGKDD International Conference on Knowledge Discovery and Data Mining* (pp. 150–160). Boston, MA: ACM Press.
- Friedman, N., Getoor, L., Koller, D., & Pfeffer, A. (1999). Learning probabilistic relational models. *Proceedings of the Sixteenth International Joint Conference on Artificial Intelligence* (pp. 1300–1307). Stockholm, Sweden: Morgan Kaufmann.
- Genesereth, M. R., & Nilsson, N. J. (1987). *Logical foundations of artificial intelligence*. San Mateo, CA: Morgan Kaufmann.
- Getoor, L., & Jensen, D. (Eds.). (2000). *Proceedings of the AAAI-2000 Workshop on Learning Statistical Models from Relational Data*. Austin, TX: AAAI Press.
- Getoor, L., & Jensen, D. (Eds.). (2003). *Proceedings of the IJCAI-2003 Workshop on Learning Statistical Models from Relational Data*. Acapulco, Mexico: IJCAI.
- Geyer, C. J., & Thompson, E. A. (1992). Constrained Monte Carlo maximum likelihood for dependent data. *Journal of the Royal Statistical Society, Series B*, 54, 657–699.
- Gilks, W. R., Richardson, S., & Spiegelhalter, D. J. (Eds.). (1996). *Markov chain Monte Carlo in practice*. London, UK: Chapman and Hall.
- Halpern, J. (1990). An analysis of first-order logics of probability. *Artificial Intelligence*, 46, 311–350.
- Heckerman, D., Chickering, D. M., Meek, C., Rounthwaite, R., & Kadie, C. (2000). Dependency networks for inference, collaborative filtering, and data visualization. *Journal of Machine Learning Research*, 1, 49–75.
- Heckerman, D., Geiger, D., & Chickering, D. M. (1995). Learning Bayesian networks: The combination of knowledge and statistical data. *Machine Learning*, 20, 197–243.



- Heckerman, D., Meek, C., & Koller, D. (2004). Probabilistic entity-relationship models, PRMs, and plate models. *Proceedings of the ICML-2004 Workshop on Statistical Relational Learning and its Connections to Other Fields* (pp. 55–60). Banff, Canada: IMLS.
- Hulten, G., & Domingos, P. (2002). Mining complex models from arbitrarily large databases in constant time. *Proceedings of the Eighth ACM SIGKDD International Conference on Knowledge Discovery and Data Mining* (pp. 525–531). Edmonton, Canada: ACM Press.
- Jaeger, M. (1998). Reasoning about infinite random structures with relational Bayesian networks. *Proceedings of the Sixth International Conference on Principles of Knowledge Representation and Reasoning*. Trento, Italy: Morgan Kaufmann.
- Jaeger, M. (2000). On the complexity of inference about probabilistic relational models. *Artificial Intelligence*, 117, 297–308.
- Kautz, H., Selman, B., & Jiang, Y. (1997). A general stochastic approach to solving problems with hard and soft constraints. In D. Gu, J. Du and P. Pardalos (Eds.), *The satisfiability problem: Theory and applications*, 573–586. New York, NY: American Mathematical Society.
- Kersting, K., & De Raedt, L. (2001). Towards combining inductive logic programming with Bayesian networks. *Proceedings of the Eleventh International Conference on Inductive Logic Programming* (pp. 118–131). Strasbourg, France: Springer.
- Laffar, J., & Lassez, J. (1987). Constraint logic programming. *Proceedings of the Fourteenth ACM Conference on Principles of Programming Languages* (pp. 111–119). Munich, Germany: ACM Press.
- Lavrač, N., & Džeroski, S. (1994). *Inductive logic programming: Techniques and applications*. Chichester, UK: Ellis Horwood.
- Liu, D. C., & Nocedal, J. (1989). On the limited memory BFGS method for large scale optimization. *Mathematical Programming*, 45, 503–528.
- Lloyd, J. W. (1987). *Foundations of logic programming*. Berlin, Germany: Springer.
- Lloyd-Richardson, E., Kazura, A., Stanton, C., Niaura, R., & Papandonatos, G. (2002). Differentiating stages of smoking intensity among adolescents: Stage-specific psychological and social influences. *Journal of Consulting and Clinical Psychology*, 70.
- Milch, B., Marthi, B., & Russell, S. (2004). BLOG: Relational modeling with unknown objects. *Proceedings of the ICML-2004 Workshop on Statistical Relational Learning and its Connections to Other Fields* (pp. 67–73). Banff, Canada: IMLS.
- Muggleton, S. (1996). Stochastic logic programs. In L. De Raedt (Ed.), *Advances in inductive logic programming*, 254–264. Amsterdam, Netherlands: IOS Press.
- Neville, J., & Jensen, D. (2003). Collective classification with relational dependency networks. *Proceedings of the Second International Workshop on Multi-Relational Data Mining* (pp. 77–91). Washington, DC: ACM Press.

- Ngo, L., & Haddawy, P. (1997). Answering queries from context-sensitive probabilistic knowledge bases. *Theoretical Computer Science*, 171, 147–177.
- Nilsson, N. (1986). Probabilistic logic. *Artificial Intelligence*, 28, 71–87.
- Nocedal, J., & Wright, S. J. (1999). *Numerical optimization*. New York, NY: Springer.
- Ourston, D., & Mooney, R. J. (1994). Theory refinement combining analytical and empirical methods. *Artificial Intelligence*, 66, 273–309.
- Parag, & Domingos, P. (2004). Multi-relational record linkage. *Proceedings of the Third International Workshop on Multi-Relational Data Mining*. Seattle, WA: ACM Press.
- Paskin, M. (2002). *Maximum entropy probabilistic logic* (Technical Report UCB/CSD-01-1161). Computer Science Division, University of California, Berkeley, CA.
- Pasula, H., & Russell, S. (2001). Approximate inference for first-order probabilistic languages. *Proceedings of the Seventeenth International Joint Conference on Artificial Intelligence* (pp. 741–748). Seattle, WA: Morgan Kaufmann.
- Pazzani, M., & Kibler, D. (1992). The utility of knowledge in inductive learning. *Machine Learning*, 9, 57–94.
- Pearl, J. (1988). *Probabilistic reasoning in intelligent systems: Networks of plausible inference*. San Francisco, CA: Morgan Kaufmann.
- Poole, D. (1993). Probabilistic Horn abduction and Bayesian networks. *Artificial Intelligence*, 64, 81–129.
- Poole, D. (2003). First-order probabilistic inference. *Proceedings of the Eighteenth International Joint Conference on Artificial Intelligence* (pp. 985–991). Acapulco, Mexico: Morgan Kaufmann.
- Popescul, A., & Ungar, L. H. (2003). Structural logistic regression for link analysis. *Proceedings of the Second International Workshop on Multi-Relational Data Mining* (pp. 92–106). Washington, DC: ACM Press.
- Puech, A., & Muggleton, S. (2003). A comparison of stochastic logic programs and Bayesian logic programs. *Proceedings of the IJCAI-2003 Workshop on Learning Statistical Models from Relational Data* (pp. 121–129). Acapulco, Mexico: IJCAI.
- Richardson, M., & Domingos, P. (2003). Building large knowledge bases by mass collaboration. *Proceedings of the Second International Conference on Knowledge Capture* (pp. 129–137). Sanibel Island, FL: ACM Press.
- Riezler, S. (1998). *Probabilistic constraint logic programming*. Doctoral dissertation, University of Tübingen, Tübingen, Germany.
- Robinson, J. A. (1965). A machine-oriented logic based on the resolution principle. *Journal of the ACM*, 12, 23–41.
- Roth, D. (1996). On the hardness of approximate reasoning. *Artificial Intelligence*, 82, 273–302.

- Sanghai, S., Domingos, P., & Weld, D. (2003). Dynamic probabilistic relational models. *Proceedings of the Eighteenth International Joint Conference on Artificial Intelligence* (pp. 992–997). Acapulco, Mexico: Morgan Kaufmann.
- Santos Costa, V., Page, D., Qazi, M., , & Cussens, J. (2003). CLP(BN): Constraint logic programming for probabilistic knowledge. *Proceedings of the Nineteenth Conference on Uncertainty in Artificial Intelligence* (pp. 517–524). Acapulco, Mexico: Morgan Kaufmann.
- Sato, T., & Kameya, Y. (1997). PRISM: A symbolic-statistical modeling language. *Proceedings of the Fifteenth International Joint Conference on Artificial Intelligence* (pp. 1330–1335). Nagoya, Japan: Morgan Kaufmann.
- Taskar, B., Abbeel, P., & Koller, D. (2002). Discriminative probabilistic models for relational data. *Proceedings of the Eighteenth Conference on Uncertainty in Artificial Intelligence* (pp. 485–492). Edmonton, Canada: Morgan Kaufmann.
- Towell, G. G., & Shavlik, J. W. (1994). Knowledge-based artificial neural networks. *Artificial Intelligence*, 70, 119–165.
- Wasserman, S., & Faust, K. (1994). *Social network analysis: Methods and applications*. Cambridge, UK: Cambridge University Press.
- Wellman, M., Breese, J. S., & Goldman, R. P. (1992). From knowledge bases to decision models. *Knowledge Engineering Review*, 7.
- Winkler, W. (1999). *The state of record linkage and current research problems* (Technical Report). Statistical Research Division, U.S. Census Bureau.
- Yedidia, J. S., Freeman, W. T., & Weiss, Y. (2001). Generalized belief propagation. In T. Leen, T. Dietterich and V. Tresp (Eds.), *Advances in neural information processing systems 13*, 689–695. Cambridge, MA: MIT Press.
- Zhu, C., Byrd, R. H., Lu, P., & Nocedal, J. (1997). Algorithm 778: L-BFGS-B, FORTRAN routines for large scale bound constrained optimization. *ACM Transactions on Mathematical Software*, 23, 550–560.

*Address for Offprints:*

Department of Computer Science and Engineering  
University of Washington  
Seattle, WA 98195-2350, U.S.A.

