

# Modelling Relational Statistics With Bayes Nets

Oliver Schulte, Hassan Khosravi, Arthur E. Kirkpatrick, Tong Man, Tianxiang Gao, and Yuke Zhu

School of Computing Science, Simon Fraser University,  
Burnaby-Vancouver, Canada

**Abstract.** Class-level dependencies model general relational statistics over attributes of linked objects and links. Class-level relationships are important in themselves, and they support applications like policy making, strategic planning, and query optimization. An example of a class-level query is “what is the percentage of friendship pairs where both friends are women?”. To represent class-level statistics, we utilize Parametrized Bayes nets (PBNs), a first-order logic extension of Bayes nets. The standard grounding semantics for PBNs is appropriate for answering queries about specific ground facts but not appropriate for answering queries about classes of individuals. We propose a *grounding-free* interpretation of PBNs that supports class-level queries, based on Halpern’s classic random selection semantics for probabilistic first-order logic [1]. Learning the parameters for this semantics can be done using the recent pseudo-likelihood measure [2] as the objective function. The parameter settings that maximize this objective function are the empirical frequencies in the relational data. A naive computation of the empirical frequencies of the relations is intractable due to the complexity imposed by negated relations. We render the computation tractable by using the fast Möbius transform. Evaluation on four benchmark datasets shows that maximum pseudo-likelihood provides accurate estimates at different sample sizes.

## 1 Introduction

Many applications store data in relational format, with different tables for entities and their links. Relational data introduces the machine learning problem of *class-level frequency estimation*: building a model that can answer generic statistical queries about classes of individuals in the database [3]. For example, a class-level query for a social network database may be “what is the percentage of friendship pairs where both are women”? A movie database example would be “what is the percentage of male users who have rated highly an action movie?”. A model of database statistics can be used for several applications:

**Statistical first-order Patterns.** AI research into combining first-order logic and probability investigated in depth the representation of statistical patterns in relational structures [1, 4]. Often such patterns can be expressed as *generic statements*, like “intelligent students tend to take difficult courses”.

**Policy making and strategic planning.** A university administrator may wish to know which program characteristics attract high-ranking students in general, rather than predict the rank of a specific student in a specific program. Maier *et al.* [5] describe several applications of causal-relational knowledge for decision making. These causal relations reflect generic correlations in the database.

**Query optimization.** A statistical model predicts a probability for given table join conditions that can be used to infer the size of a database query result [3]. Estimating join sizes (selectivity estimation) is used to minimize the size of intermediate join tables [6].

*Semantics.* We focus on building a Bayes net model for relational statistics, using the Parametrized Bayes nets (PBNs) of Poole [7]. The nodes in a PBN are constructed with functors and first-order variables (e.g., *gender(X)* may be a node). The original PBN semantics is a grounding semantics where the first-order Bayes net is instantiated with all possible groundings to obtain a directed graph whose nodes are functors with constants (e.g., *gender(sam)*). The ground graph can be used to answer queries *about individuals*, such as “if user *sam* has 3 friends, female *rozita*, males *ali* and *victor*, what is the probability that *sam* is a woman”? However, as pointed out by Getoor [8], the ground graph is not appropriate for answering class-level queries because these are about generic rates and percentages, not about any particular individuals.

We propose a new interpretation for Parametrized Bayes nets that supports class-level queries. The semantics is based on Halpern’s classic random selection semantics for probabilistic first-order logic [1, 4]. While we focus on PBNs, the random selection semantics can be applied to any statistical-relational model whose syntax is based on first-order logic. Halpern’s semantics views statements with first-order variables as expressing statistical information about classes (domains) of individuals. For instance, the claim “the percentage of friendship pairs where both are women is 60%” could be expressed by the formula

$$P(\text{Gender}(X) = \text{female}, \text{Gender}(Y) = \text{female} | \text{Friend}(X, Y)) = 60\%.$$

*Learning.* A standard Bayes net parameter learning method is maximum likelihood estimation, but this method is difficult to apply for Bayes nets that represent relational data because the cyclic data dependencies in relations violate the requirements of a traditional likelihood measure. We circumvent the limitations of classical likelihood measures by using a relational pseudo-likelihood measure for Bayes nets [2] that is well defined even in the presence of cyclic dependencies. In addition to this robustness, the relational pseudo-likelihood matches the random selection semantics because it is also based on the concept of random instantiations. An estimator that chooses the parameters that maximize this pseudo-likelihood function (MPLE), has a closed-form solution: the MPLE parameters are the empirical frequencies, as with classical i.i.d. maximum likelihood estimation. Since MPLE depends only on the generic event frequencies in the

data, it can be viewed as an instance of *lifted learning*. Computing the empirical frequencies for negated relationships is difficult, however, because enumerating the complement of a relationship table is computationally infeasible. We show that the fast Möbius transform [9] makes MPLE tractable, even in the case of negated relationships.

*Results.* We evaluate MPLE on four benchmark real-world datasets. On complete-population samples MPLE achieves near perfect accuracy in parameter estimates, and excellent performance on Bayes net queries. The accuracy of MPLE parameter values is high even on medium-size samples.

*Contributions.* Our main contributions for frequency modelling in relational data are the following:

1. A new class-level semantics for graphical first-order models, derived from the random selection semantics for probabilistic first-order logic.
2. Making the computation of frequency estimates tractable by computing database statistics using the fast Möbius transform. This transform is a general procedure for computing relational counts that involve negated links. It has application in Probabilistic Relational Models [10, Sec.5.8.4.2], multi-relational data mining, and inductive logic programming models with clauses that contain negated relationships.
3. We contribute to unification of instance-level/ground-level and class-level relational probabilities (defined in the next section) in two ways. (1) We show that the same first-order model can be used for both types of inference. (2) We show that the same objective function is suitable for learning models for both types of queries.

*Paper Organization.* We review related work, then background and notation. Section 4 presents the random selection semantics for Bayes nets. Section 5 presents the fast Möbius transform for relational data. Simulation results are presented in Section 6, showing the runtime cost of estimating parameters, and evaluations of their quality by (a) comparison with the true population parameter values, and (b) inference on random queries.

## 2 Related Work

*Class-level and Instance-level Relational Probabilities.* Classic AI research established a fundamental distinction between two types of probabilities associated with a relational structure [1, 4]. *Class-level probabilities*, also called type 1 probabilities are assigned to the rates, statistics, or frequencies of events in a database. These concern classes of entities (e.g., students, courses, users) rather than specific entities. *Instance-level probabilities*, also called type 2 probabilities are assigned to specific, non-repeatable events or the properties of specific entities. Syntactically, class-level probabilities are assigned to formulas that contain

first-order variables (e.g.,  $P(\textit{Flies}(X)|\textit{Bird}(X)) = 90\%$ , or “birds fly” with probability 90%), whereas instance-level probabilities are assigned to formulas that contain constants only (e.g.,  $P(\textit{Flies}(\textit{tweety})) = 90\%$ ). There has been much AI research on using Bayes nets for representing and reasoning both with class probabilities [11] and instance probabilities [12]. Most statistical-relational learning has been concerned with instance probabilities: For instance, Probabilistic Relational Models (PRMs) [10] and Markov Logic Networks (MLNs) [13] define probabilities for ground instances using a grounding semantics. A well-known problem for directed graphical models is that there may be cyclic dependencies among ground facts. For example, if the gender of a friend predicts the gender of an individual, then it may be the case that the gender of Sam predicts the gender of her friend Jim, and the gender of Jim predicts that of Sam, leading to a cyclic dependencies, which are not allowed in Bayes nets. Various approaches to the cyclicity problem have been considered, such as (1) applying Bayes nets only under the assumption that there are no cycles at the instance level, and (2) using another type of graphical model: undirected models like Markov nets, or a directed model that allows cycles such as dependency networks. For discussion please see [?]. In the grounding-free interpretation of Bayes nets that we propose for class-level probabilities, there is no issue with cycles among features of ground individuals because the Bayes net does not refer to ground individuals. Just as in propositional learning, any joint distribution over random variables can be represented by an acyclic Bayes net, in class-level modelling any joint distribution over class-level random variables can be represented by an acyclic Bayes net. Thus an advantage of class-level modelling from the point of view of Bayes net learning is that it avoids the cyclicity problem. We elaborate on this point after we define how Bayes nets represent class-level probabilities.

The notion of a parfactor is a key concept for efficient instance-level inference with first-order models [7]. The parfactor approach shares with our work an emphasis on reasoning about class of individuals (e.g., the number of male friends that Sam has). The counting algorithms in this paper can likely be applied to computing parfactors (e.g., parfactors that involve negated relations). The motivation for parfactors, however, is as a means to the end of computing instance-level predictions (e.g., predict the gender of Sam). In contrast, our motivation is to learn class frequencies as an end in itself. Because parfactors are used with instance-level inferences, they usually concern classes defined with reference to specific individuals, whereas the statistics we model in this paper are at the class-level only.

*Statistical Relational Models.* To our knowledge, Statistical Relational Models (SRMs) due to Getoor, Taskar and Koller [8], are the only prior statistical models with a class-level probability semantics. A direct empirical comparison is difficult as code has not been released (Getoor, personal communication). The syntax of SRMs differ from PBNs and other statistical-relational models in several respects. (1) The SRM syntax is not that of first-order logic, but is derived from a tuple semantics [8, Def.6.3], which is different from the random selection semantics we propose for PBNs. (2) SRMs are less expressive: They cannot ex-

press general combinations of positive and negative relationships [8, Def.6.11]. Basically, this restriction stems from the fact that the SRM semantics is based on randomly selecting tuples from *existing tables in the database*. Complements of relationship tables are usually not stored as existing tables (e.g., there is no table that lists the set of user pairs who are *not* friends). The expressive power of SRMs and PBNs becomes essentially equivalent, with respect to class-level probabilities, if the SRM semantics is extended so that random tuples can be drawn from complement tables as well as the original tables. Adding complement tables to SRMs leads to the challenge that SRM learning has to include learning with negated relations, which we address in this paper.<sup>1</sup>

*Unified Learning for Type 1 and Type 2 Probabilities.* Previous work used different models for the two basic types of probability query (SRMs for class-level, template models for instance-level). In this paper we employ PBNs and the pseudo likelihood to learn models that are accurate for class-level probabilities. Previous research employed the same model class and objective function for learning models that are accurate for instance-level probabilities [15, 2]. These findings indicate that learning accurate class-level probabilities leads to accurate instance-level predictions. We believe that a unified approach to learning for both relational probability types is an exciting research direction for statistical-relational learning.

add diagram about previous work?

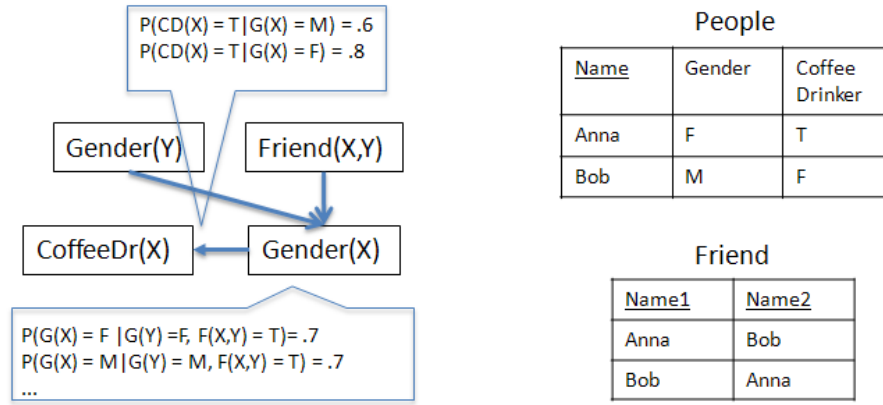
### 3 Background: Parametrized Bayes Nets

Our work combines concepts from relational databases and graphical models. As much as possible, we use standard notation in these different areas. Parametrized Bayes nets are a basic graphical model for relational data [7]. The syntax of PBNs is as follows. We assume familiarity with Bayes nets and concepts such as CP-table and I-map [16]. A **functor** is a function symbol or a predicate symbol. Each functor has a set of values (constants) called the **range** of the functor. There are two types of functor nodes: Boolean **relationship functors** that indicate whether a relationship holds (e.g., *Friend*), and **attribute functors** that correspond to the value of an attribute (e.g., *Gender*). Conforming to statistical terminology, Poole refers to first-order variables as population variables. A **population variable**  $X$  is associated with a population, a set of individuals; in logical terminology, a type, domain, or class. A **functor random variable** or **functor node** is of the form  $f(X_1, \dots, X_k)$ . In this paper we assume that functor nodes contain first-order variables only (no constants). A **Parametrized Bayes Net** is a Bayes net whose nodes are functor nodes. In the following we often omit the prefix “Parametrized” and speak simply of Bayes nets. Figure 1

<sup>1</sup> With complement tables included, the main difference is then that the SRM semantics randomly selects *tuples*, whereas Halpern’s semantics randomly selects *individuals*. This corresponds to the difference between the tuple relational calculus and the domain relational calculus, which are known to be equivalent in expressive power [14].

shows a PBN. The syntax of PBNs is similar to that of other directed relational graphical models (cf. [7]). An **instantiation** or **grounding** for a set of variables  $X_1, \dots, X_k$  assigns a constant  $c_i$  from the population of  $X_i$  to each variable  $X_i$ .

The functor formalism is rich enough to represent an entity-relationship schema via the following translation: Entity sets correspond to populations, descriptive attributes to functors, relationship tables to Boolean functors, and foreign key constraints to type constraints on the arguments of relationship predicates. Figure 1 shows a Parametrized Bayes net and a simple relational database instance.

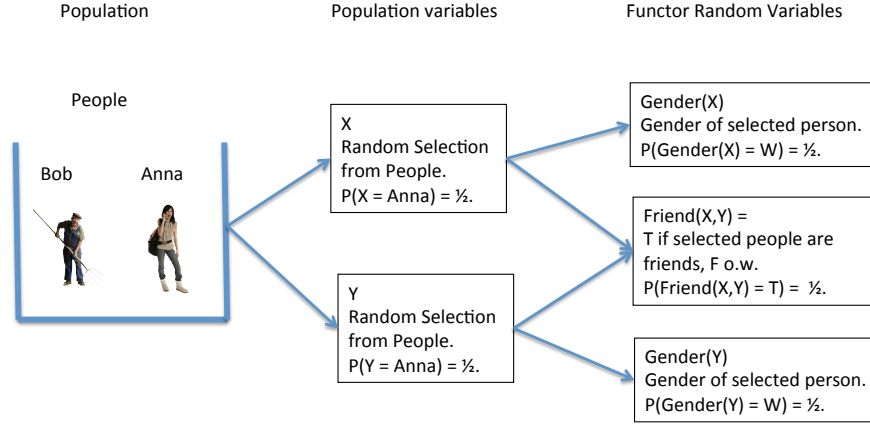


**Fig. 1.** Left: An illustrative Parametrized Bayes Net.  $Friend(X, Y)$  is a relationship node, the other three nodes are attribute nodes. Right: A simple relational database instance.

## 4 Random Selection Interpretation for Bayes Nets

We discuss how the random selection semantics can be applied to Bayes Nets, and relate this interpretation to database queries.

*Functor Node Interpretation.* The random selection semantics views first-order/population variables  $X_1, X_2, \dots, X_k$  as independent random variables that each sample an individual from the appropriate population [1, 4]. Then a functor of the form  $f(X_1, X_2, \dots, X_k)$  represents a function of a random  $k$ -tuple. Since a *function of a random variable is itself a random variable*, this means that *we can view functor nodes as random variables in their own right, without grounding the variables first*. Figure 2 illustrates the construction.



**Fig. 2.** Viewing population variables as random selections from the associated population, functor nodes are functions of random variables, hence themselves random variables.

*Bayes Net Joint Probabilities.* Via the standard product formula, a Bayes net  $B$  entails a probability value  $P_B$  for each joint assignment of values to functor nodes. In logical terms, the Bayes net applies probabilities to conjunctions of literals, where a literal is an assignment of a value to a functor node. For example, the Bayes net of Figure 1 entails a set of joint probabilities of the form  $P_B(G(Y) = a_1, F(X, Y) = a_2, G(X) = a_3, CDr(X) = a_4)$ , where we have used the obvious abbreviations for functors, and each  $a_i$  is a constant from the appropriate domain. To illustrate, with uniform priors on  $G(Y)$  and  $F(X, Y)$ , the BN of Figure 1 entails the probability

$$P_B(G(Y) = M, F(X, Y) = T, G(X) = W, CDr(X) = T) = 0.5 \times 0.5 \times 0.3 \times 0.8 = 0.06. \quad (1)$$

Using the axioms of probability calculus, we can derive further probability statements over the functor nodes, such as marginal and conditional probabilities. For instance, we have

$P_B(G(X) = W | G(Y) = W) \equiv P_B(G(X) = W, G(Y) = W) / P_B(G(Y) = W)$ . Random selection provides an interpretation for such probability statements. For instance, the meaning of Equation 1 is “if we randomly select two users  $X$  and  $Y$ , there is a 6% chance that they are friends, both are women, and one is a coffee drinker”. Halpern develops this interpretation as a formal semantics for first-order formulas that state probability assignments [1].

*Database Queries and the Database Distribution.* For each joint value assignment to nodes in a Bayes net there is a natural concept of the frequency with which this assignment holds in a given input database. The database frequency is the

move up to background

should we do this more rigorously in background like Chiang and Poole?

basis of the parameter learning algorithm we present below. It is connected to the cardinality of database queries, which is the basis for using Bayes net for selectivity estimation [3]. Consider a database query expressed in a *logical query language*, such as the domain relational calculus (DRC) [14]. An example would be

$$\{\langle X, Y \rangle | G(Y) = M, F(X, Y) = T, G(X) = W, CDr(X) = T\} \quad (2)$$

where  $X, Y$  are the *query variables* that occur free (unquantified) in the query formula. The query returns the set of all people pairs  $\langle x, y \rangle$  that satisfy the query formula. Logical queries are very powerful: a fundamental result in database theory states that a subset of DRC queries known as safe queries has expressive power equivalent to relational algebra [14]. The cardinality of a query is the number of tuples returned; the cardinality of Query 3 applied to the data of Figure 1 is 1. Assuming an assignment of types to the query variables, the largest possible query result is the cross-product of the type domains. Thus the **database frequency** of a query is the cardinality of the query, divided by the product of the domain sizes for the query variables. The database frequency for Query 3 in the database of Figure 1 is

$$P_{\mathcal{D}}(G(Y) = M, F(X, Y) = T, G(X) = W, CDr(X) = T) = \frac{1}{2 \cdot 2}. \quad (3)$$

We refer to  $P_{\mathcal{D}}$  as the **database distribution** over functor node assignments. Assuming knowledge of the type domain associated with the query variables, the database frequency  $P_{\mathcal{D}}$  immediately gives an estimate of the query cardinality. It is easy to see that the database distribution coincides with the random selection semantics when the random selection is uniform over the individuals listed for a given type in the database.

*Learning.* In learning, an observed database instance  $\mathcal{D}$  provides data only for a subpopulation (e.g., a web crawl). If  $U$  denotes the complete population data (e.g., the entire network), we can think of the goal of Bayes net learning as finding a model  $P_B$  that estimates the operating complete-population distribution  $P_U$  from a sample database  $\mathcal{D}$ . Since an acyclic Bayes net can represent any joint distribution over the random variables in its graph, for every complete population distribution  $P_U$  that defines a joint class-level distribution over the functor nodes, there is an acyclic Parametrized Bayes net that represents  $P_U$  using the random selection semantics.

The database distribution  $P_{\mathcal{D}}$  is analogous to the sampling distribution in a propositional i.i.d. setting. In propositional Bayes net learning, the basic parameter estimation method assigns the empirical conditional probability computed from a data table to the conditional probability parameters in a Bayes net. In this paper we use the analogous empirical conditional probabilities defined by the database distribution  $P_{\mathcal{D}}$ . A theoretical justification for this method is provided by the pseudo-likelihood measure recently defined for a PBN  $B$  and database

could consider a diagram



$\mathcal{D}$ . The pseudo log-likelihood for a database  $\mathcal{D}$  given a PBN  $B$  is the expected log-likelihood of a random instantiation of population variables in the PBN with individuals and values from the database  $\mathcal{D}$ . For a fixed database  $\mathcal{D}$  and Bayes net structure, the parameter values that maximize the pseudo-likelihood are the conditional empirical frequencies defined by the database distribution  $P_{\mathcal{D}}$  [2, Prop.3.1]. This result is exactly analogous to maximum likelihood estimation for i.i.d. data. In the remainder of the paper we evaluate database frequencies as parameter estimates. We begin with a procedure for computing them.

think about giving  
an example, copy  
from other paper

## 5 Computing Relational Frequencies

Initial work in SRL modelled the distribution of descriptive attributes given knowledge of existing links. Database statistics conditional on the *presence* of one or more relationships can be computed by table joins with SQL. More recent models represent *uncertainty about relationships* with link indicator variables. For instance, a Parametrized Bayes net includes relationship indicator variables such as *Friend*( $X, Y$ ). Learning with link uncertainty requires computing sufficient statistics that involve the *absence* of relationships. A naive approach would explicitly construct new data tables that enumerate tuples of objects that are *not* related. However, the number of unrelated tuples is too large to make this scalable (think about the number of user pairs who are *not* friends on Facebook). Can we instead reduce the computation of sufficient statistics that involve negated relationships to the computation of sufficient statistics that involve existing (positive) relationships only? The classic Möbius parametrization for binary random variables provides an affirmative answer [18, p.239]. Consider a set  $b_1, \dots, b_m$  of binary variables, where all marginal probabilities are available that involve only positive values. Thus we have available probabilities such as  $P(b_1 = 1)$ ;  $P(b_1 = 1, b_2 = 1)$ ;  $P(b_1 = 1, b_3 = 1, b_k = 1)$ ; etc. These joint probabilities are the **Möbius parameters** of the joint distribution. The Möbius inversion theorem entails that *all* joint probabilities, involving any number of 0 values, can be computed as an alternating sum of the Möbius parameters. We can apply this result for MPLE as follows. Consider a PBN family containing  $m$  relationship nodes. We wish to compute frequencies of the joint family assignments, from which conditional probabilities are easily derived. The Möbius inversion theorem entails that each joint frequency *can be computed from joint frequencies that involve existing relationships only*.

The **fast Möbius transform** (FMT) is an optimal algorithm for converting the Möbius parameters to a complete set of joint probabilities [9].<sup>2</sup> The FMT was originally described using category theory with lattice structures. Our version is adapted for **joint probability tables** (JP-tables). A JP-table is just like a CP-table whose rows correspond to joint probabilities rather than conditional probabilities. To represent a Möbius parameter, we allow relationship nodes to take on the value  $*$  for “unspecified”. For instance, suppose that the family nodes

<sup>2</sup> For the case of a single relationship, Getoor et al. [10] introduced a “1-minus trick”; the FMT generalizes this to the multi-relational case.

---

**Algorithm 1** The fast Möbius transform for parameter estimation in a Parametrized Bayes Net.

---

Input: database  $\mathcal{D}$ ; a set of functor nodes divided into attribute nodes  $A_1, \dots, A_j$  and relationship nodes  $R_1, \dots, R_m$ .

Output: Joint Probability specifying the data frequencies for each joint assignment to the input functor nodes.

```

1: for all attribute value assignments  $A_1 := a_1, \dots, A_j := a_j$  do
2:   initialize the JP-table with the Möbius parameters: set all relationship nodes to
     either  $T$  or  $*$ ; find joint frequencies with data queries.
3:   for  $i = 1$  to  $m$  do
4:     Change all occurrences of  $R_i = *$  to  $R_i = F$ .
5:     Update the joint frequencies using (4).
6:   end for
7: end for

```

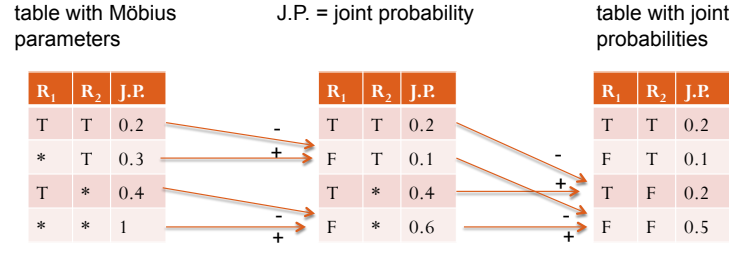
---

are  $Int(S)$ ,  $Reg(S, C)$ ,  $RA(S, P)$ . Then the Möbius parameter  $P(Int(S) = 1)$  is stored in the row where  $Int(S) = 1$ ,  $Registered(S, C) = *$ ,  $RA(S, P) = *$ . The FMT uses a local update operation corresponding to the simple probabilistic identity

$$P(\sigma, \mathbf{R}, R = F) := P(\sigma, \mathbf{R}) - P(\sigma, \mathbf{R}, R = T) \quad (4)$$

where  $\sigma$  is an attribute condition that does not involve relationships and  $\mathbf{R}$  specifies values for a list of relationship nodes. Eq. 4 shows how a probability that involves  $k + 1$  false relationships can be computed from two probabilities that each involve only  $k$  false relationships, for  $k \geq 0$ . The FMT initializes the JP-table with the Möbius parameters that do not contain negated relationships, that is, all relationship nodes have the value  $T$  or  $*$ . It then goes through the relationship nodes  $R_1, \dots, R_m$  in order, replaces at stage  $i$  all occurrences of  $R_i = *$  with  $R_i = F$ , and applies the local update equation for the probability value for the modified row. At termination, all  $*$  values have been replaced by  $F$  and the JP-table specifies all joint frequencies. Algorithm 1 gives pseudocode and Figure 3 illustrates the FMT in a schematic example with two relationship nodes.

*Complexity Analysis.* Kennes and Smets [9] provide a thorough theoretical analysis of the FMT. We summarize the main points. (1) The primary property of the FMT is that *it accesses data only about existing links*, never about non-existing links. A secondary but attractive property of FMT is that the number of additions performed is  $m2^{m-1}$ . A lower bound argument shows that this is optimal [9, Cor.1]. (3) Without a bound on  $m$ , computing sufficient statistics in a relational structure is #P-complete [13, Prop.12.4]. In practice, the number  $m$  of relationship nodes is small, 4 or less.



**Fig. 3.** The fast Möbius transform with  $m = 2$  relationship nodes. For simplicity we omit attribute conditions.

## 6 Evaluation

All experiments were done on a QUAD CPU Q6700 with a 2.66GHz CPU and 8GB of RAM. We evaluated the algorithm on real-world datasets. The datasets and our code are available on the Web [19].

### 6.1 Datasets

We used four benchmark real-world databases, with the modifications by [15], which contains details and references.

**Mondial Database.** A geography database. Mondial features a self-relationship, *Borders*, that indicates which countries border each other.

**Hepatitis Database.** A modified version of the PKDD'02 Discovery Challenge database.

**Financial** A dataset from the PKDD 1999 cup.

**MovieLens.** A dataset from the UC Irvine machine learning repository.

To obtain a Bayes net structure for each dataset, we applied the learn-and-join algorithm [15] to each database. This is the state-of-the-art structure learning algorithm for PBNs; for an objective function, it uses the pseudo-likelihood described in Section 4. We also conducted experiments with synthetic graphs and datasets. The results are similar to those on real-life datasets. We omit details for lack of space. Our implementation makes use of version 4.3.9-0 of CMU's Tetrad package [?].

### 6.2 Learning Times

Table 1 shows the runtimes for computing parameter values. The Complement method uses SQL queries that explicitly construct tables for the complement of relationships (tables that contain tuples of unrelated entities), while the FMT method uses the fast Möbius transform to compute the conditional probabilities. The FMT is faster by orders of magnitude, ranging from a factor of 15–237.

Figure5 is hard to read.

**Table 1.** Learning time results (sec) for the fast Möbius transform vs. constructing complement tables. For each database, we show the number of tuples, and the number of parameters (conditional probabilities) in the fixed Bayes net structure.

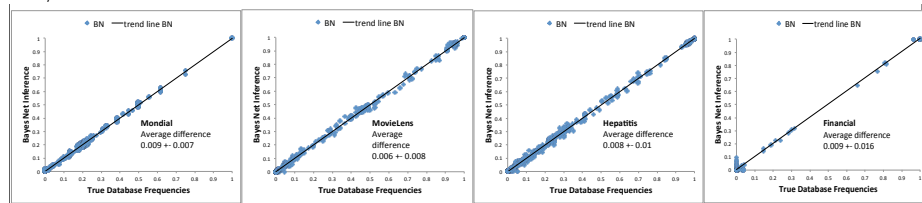
Database	Parameters	#tuples	Complement	FMT	Ratio
Mondial	1618	814	157	<b>7</b>	22
Hepatitis	1987	12,447	18,246	<b>77</b>	237
Financial	10926	17,912	228,114	<b>14,821</b>	15
MovieLens	326	82,623	2,070	<b>50</b>	41

### 6.3 Inference

The basic inference task for Bayes nets is answering probabilistic queries. If the given Bayes net structure is an I-map of the true distribution, then correct parameter values lead to correct predictions. Thus the performance on queries has been used to evaluate parameter learning [21]. We randomly generate queries for each dataset according to the following procedure. First, choose a target node  $V$  100 times, and go through each possible value  $a$  of  $V$  such that  $P(V = a)$  is the probability to be predicted. For each value  $a$ , choose the number  $k$  of conditioning variables, ranging from 1 to 3. Select  $k$  variables  $V_1, \dots, V_k$  and corresponding values  $a_1, \dots, a_k$ . The query to be answered is then  $P(V = a | V_1 = a_1, \dots, V_k = a_k)$ .

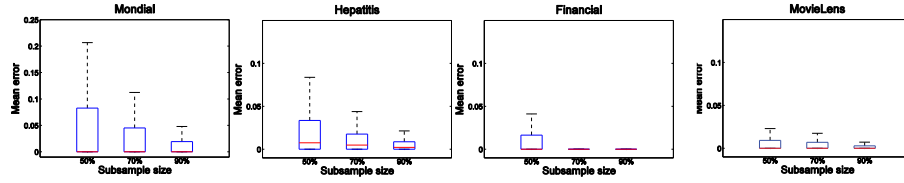
As in [3], we evaluate queries after learning parameter values on the entire database. Thus the BN is viewed as a statistical summary of the data rather than generalizing from a sample. BN inference is carried out using the Approximate Updater in CMU’s Tetrad program. Figure 4 shows the query performance for each database. A point  $(x, y)$  on a curve indicates that there is a query such that the true probability value in the database is  $x$  and the probability value estimated by the model is  $y$ . The Bayes net inference is close to the ideal identity line, with an average error of less than 1%.

**Fig. 4.** Query Performance: Estimated vs. true probability. The average error and standard deviation are shown as well. Number of queries/average inference time per query: Mondial, 506/0.08sec; MovieLens, 546/0.05sec; Hepatitis, 489/0.1sec; Financial, 140/0.02sec.



## 6.4 Conditional Probabilities

**Fig. 5.** Error in conditional probability estimates. The box plots show the distribution of the absolute difference between each true conditional probability value and the estimate produced by our method, averaged over 10 random subdatabases. The center line is the median error **box meaning? fix “mean” label**



In the previous inference results, the correct query answers are defined by the entire population (benchmark database), and the model is also trained on the entire database. To study parameter estimation at different sample sizes, we continue to define the correct value in terms of the entire population, but train the model on  $N\%$  of the data for varying  $N$ . Conceptually, we treated each benchmark database as specifying an entire population, and consider the problem of estimating the complete-population frequencies from partial-population data. The  $N\%$  parameter is uniform across tables and databases. We employed standard subgraph subsampling [20, 15], which selects entities from each entity table uniformly at random and restricts the relationship tuples in each subdatabase to those that involve only the selected entities. Subgraph sampling matches the random selection semantics which is based on random draws from a population. It is applicable when the observations include positive and negative link information (e.g., not listing two countries as neighbors implies that they are not neighbors). The subgraph method satisfies an ergodic law of large numbers in the sense that as the subsample size increases, the subsample relational frequencies approach the population relational frequencies.

Figure 5 plots the difference between the true conditional probabilities and the MPLE estimates. With increasing sample size, **MPLE** estimates approach the true value in all cases. Even for the smaller sample sizes, the median error is close to 0, confirming that most estimates are very close to correct. As the box plots show, the 3rd error quartile of estimates is bound within 10% on Mondial, the worst case, and within less than 5% on the other datasets.

## 7 Conclusion

We introduced a new interpretation of Parametrized Bayes nets as models of class-level statistics in a relational structure. For parameter learning we utilized

the empirical database frequencies, which can be feasibly computed using the fast Möbius transform, even for frequencies concerning negated links. In evaluation on four benchmark databases, the maximum pseudo-likelihood estimates approach the true conditional probabilities as observations increase. The fit is good even for medium data sizes.

A direction for future work is to adapt more techniques from propositional i.i.d. Bayes net parameter learning, such as smoothing frequencies and incorporating uncertainty in parameter estimates [21]. A theoretical understanding of estimator variance would be desirable: we may adapt the asymptotic approximations of [21], or apply graph estimator theory [20]. Halpern [1] showed that any instance-level inference model can be used for class-level inference by using ground queries that contain new constants only (e.g., random-student, random-course, and random-prof). We plan to use this scheme to evaluate instance-level models, such as Markov Logic Networks, for class-level queries.

## References

- [1] Halpern, J.Y.: An analysis of first-order logics of probability. *Artificial Intelligence* **46**(3) (1990) 311–350
- [2] Schulte, O.: A tractable pseudo-likelihood function for Bayes nets applied to relational data. In: *SIAM SDM*. (2011) 462–473
- [3] Getoor, L., Taskar, B., Koller, D.: Selectivity estimation using probabilistic models. *ACM SIGMOD Record* **30**(2) (2001) 461–472
- [4] Bacchus, F.: Representing and reasoning with probabilistic knowledge: a logical approach to probabilities. MIT Press, Cambridge, MA, USA (1990)
- [5] Maier, M., Taylor, B., Oktay, H., Jensen, D.: Learning causal models of relational domains. In: *AAAI*. (2010)
- [6] McMahan, B.J., Pan, G., Porter, P., Vardi, M.Y.: Projection pushing revisited. In: *EDBT*. (2004) 441–458
- [7] Poole, D.: First-order probabilistic inference. In: *IJCAI*. (2003) 985–991
- [8] Getoor, L.: Learning Statistical Models From Relational Data. PhD thesis, Department of Computer Science, Stanford University (2001)
- [9] Kennes, R., Smets, P.: Computational aspects of the Möbius transformation. In: *UAI*. (1990) 401–416
- [10] Getoor, L., Friedman, N., Koller, D., Pfeffer, A., Taskar, B.: Probabilistic relational models. [22] chapter 5 129–173
- [11] Bacchus, F.: Using first-order probability logic for the construction of Bayesian networks. In: *UAI*. (1993) 219–226
- [12] Ngo, L., Haddawy, P.: Answering queries from context-sensitive probabilistic knowledge bases. *Theor. Comput. Sci.* **171**(1-2) (1997) 147–177
- [13] Domingos, P., Richardson, M.: Markov logic: A unifying framework for statistical relational learning. [22]
- [14] Ullman, J.D.: Principles of database systems. 2. Computer Science Press (1982)
- [15] Khosravi, H., Schulte, O., Man, T., Xu, X., Bina, B.: Structure learning for Markov logic networks with many descriptive attributes. In: *AAAI*. (2010) 487–493
- [16] Pearl, J.: Probabilistic Reasoning in Intelligent Systems. Morgan Kaufmann (1988)

- [17] Schulte, O., Khosravi, H., Man, T.: Learning directed relational models with recursive dependencies. In: ILP. Volume 7207 of Springer LNAI. (2011) 39–44
- [18] Lauritzen, S.L.: Graphical Models (Oxford Statistical Science Series). Oxford University Press, USA (July 1996)
- [19] Khosravi, H., Man, T., Hu, J., Gao, E., Schulte, O.: Learn and join algorithm code. URL = <http://www.cs.sfu.ca/~oschulte/jbn/>.
- [20] Frank, O.: Estimation of graph totals. *Scandinavian Journal of Statistics* **4:2** (1977) 81–89
- [21] Allen, T.V., Singh, A., Greiner, R., Hooper, P.: Quantifying the uncertainty of a belief net response: Bayesian error-bars for belief net inference. *Artif. Intell.* **172**(4-5) (2008) 483–513
- [22] Getoor, L., Tasker, B.: Introduction to statistical relational learning. MIT Press (2007)