

Basic Principles of Learning Bayesian Logic Programs^{*}

Kristian Kersting¹ and Luc De Raedt²

¹ CSAIL, Massachusetts Institute of Technologie,
32 Vassar Street, Cambridge, MA 02139-4307, USA
kersting@csail.mit.edu

² Departement Computerwetenschappen, K.U. Leuven,
Celestijnenlaan 200A - bus 2402, B-3001 Heverlee, Belgium
Luc.DeRaedt@cs.kuleuven.be

Abstract. Bayesian logic programs tightly integrate definite logic programs with Bayesian networks in order to incorporate the notions of objects and relations into Bayesian networks. They establish a one-to-one mapping between ground atoms and random variables, and between the *immediate consequence* operator and the *directly influenced by* relation. In doing so, they nicely separate the qualitative (i.e. logical) component from the quantitative (i.e. the probabilistic) one providing a natural framework to describe general, probabilistic dependencies among sets of random variables. In this chapter, we present results on combining Inductive Logic Programming with Bayesian networks to learn both the qualitative and the quantitative components of Bayesian logic programs from data. More precisely, we show how the qualitative components can be learned by combining the inductive logic programming setting learning from interpretations with score-based techniques for learning Bayesian networks. The estimation of the quantitative components is reduced to the corresponding problem of (dynamic) Bayesian networks.

1 Introduction

In recent years, there has been an increasing interest in integrating probability theory with first order logic. One of the research streams [42,40,24,19,29] concentrates on first order extensions of Bayesian networks [41]. The reason why this has attracted attention is, that even though Bayesian networks are one of the most important, efficient and elegant frameworks for representing and reasoning with probabilistic models, they suffer from an inherently propositional character. A single Bayesian network specifies a joint probability density over a finite set of random variables and consists of two components:

^{*} The is a slightly modified version of *Basic Principles of Learning Bayesian Logic Programs*, Technical Report No. 174, Institute for Computer Science, University of Freiburg, Germany, June 2002. The major change is an improved section on parameter estimation. For historical reasons, all other parts are left unchanged (next to minor editorial changes).

- a *qualitative* one that encodes the local influences among the random variables using a directed acyclic graph, and
- a *quantitative* one that encodes the probability densities over these local influences.

Imagine a Bayesian network modelling the localization of genes/proteins. Every gene would be a single random variable. There is no way of formulating general probabilistic regularities among the localizations of the genes such as *the protein P encoded by gene G has localization L if P interacts with another protein P' that has localization L*.

Bayesian logic programs are a language that overcomes this propositional character by tightly integrating definite logic programs with Bayesian networks to incorporate the notions of objects and relations. In doing so, they can naturally be used to do first order classification, clustering, and regression. Their underlying idea is to establish a one-to-one mapping between ground atoms and random variables, and between the *immediate consequence operator* and the *directly influences by* relation. In doing so, they nicely separate the qualitative (i.e. logical) component from the quantitative (i.e. the probabilistic) one providing a natural framework to describe general, probabilistic dependencies among sets of random variables such as the rule stated above.

It is, however, well-known that determining the structure of a Bayesian network, and therefore also of a Bayesian logic program, can be difficult and expensive. In 1997, Koller and Pfeffer [33] addressed the question “*where do the numbers come from?*” for similar frameworks. So far, this issue has not yet attracted much attention in the context of first order extensions of Bayesian networks (with the exception of [33,19]). In this context, we present for the first time how to calculate the *gradient* for a maximum likelihood estimation of the parameters of Bayesian logic programs. Together with the EM algorithm which we will present, this gives one a rich class of optimization techniques such as conjugate gradient and the possibility to speed up the EM algorithm, see e.g. [38].

Moreover, Koller and Pfeffer [33] rose the question whether techniques from inductive logic programming (ILP) could help to learn the logical component of first order probabilistic models. In [30], we suggested that the ILP setting *learning from interpretations* [13,14,6] is a good candidate for investigating this question. In this chapter we would like to make our suggestions more concrete. We present a novel scheme to learn intensional clauses within Bayesian logic programs [28,29]. It combines techniques from inductive logic programming with techniques for learning Bayesian networks. More precisely, we will show that *learning from interpretations* can indeed be integrated with score-based Bayesian network learning techniques in order to learn Bayesian logic programs. Thus, we answer Koller and Pfeffer’s question affirmatively.

We proceed as follows. After briefly reviewing the framework of Bayesian logic programs in Section 2, we define the learning problem in Section 3. Based on this, we then present a score-based greedy algorithm called SCOOPY solving the learning problem. More precisely, Section 4 presents SCOOPY first in the context of a special class of propositional Bayesian logic programs, i.e. Bayesian networks,

and then on general Bayesian logic programs. In Section 5, we formulate the likelihood of the parameters of a Bayesian logic program given some data and, based on this, we present a gradient-based and an EM method to find that parameters which maximize the likelihood. Section 6 reports on first experiments. Before concluding the paper, we touch upon related work.

We assume some familiarity with logic programming or Prolog (see e.g. [45,37]) as well as with Bayesian networks (see e.g. [41,10,27]).

2 Bayesian Logic Programs

Throughout the paper we will use an example from genetics which is inspired by Friedman et al. [19]: “it is a genetic model of the inheritance of a single gene that determines a person’s X blood type $\text{bt}(X)$. Each person X has two copies of the chromosome containing this gene, one, $\text{mc}(Y)$, inherited from her mother $\text{m}(Y, X)$, and one, $\text{pc}(Z)$, inherited from her father $\text{f}(Z, X)$.” We will use \mathbf{P} to denote a probability distribution, e.g. $\mathbf{P}(x)$, and the normal letter P to denote a probability value, e.g. $P(x = v)$, where v is a state of x .

2.1 Representation Language

The basic idea underlying our framework is that each Bayesian logic program specifies a (possibly infinite) Bayesian network, with one node for each (Bayesian) ground atom (see below). A Bayesian logic program B consist of two components:

- a *qualitative* or *logical* one, a set of Bayesian clauses (cf. below), and
- a *quantitative* one, a set of conditional probability distributions and combining rules (cf. below) corresponding to that logical structure.

Definition 1 (*Bayesian Clause*). A Bayesian (definite) clause c is an expression of the form

$$A \mid A_1, \dots, A_n$$

where $n \geq 0$, the A, A_1, \dots, A_n are Bayesian atoms and all Bayesian atoms are (implicitly) universally quantified. We define $\text{head}(c) = A$ and $\text{body}(c) = \{A_1, \dots, A_n\}$.

So, the differences between a *Bayesian clause* and a *logical* one are:

1. The atoms $p(t_1, \dots, t_n)$ and predicates p arising are Bayesian, which means that they have an associated (finite) domain¹ $\mathbf{S}(p)$, and
2. We use “ \mid ” instead of “ $:-$ ”.

For instance, consider the Bayesian clause c

$$\text{bt}(X) \mid \text{mc}(X), \text{pc}(X).$$

where $\mathbf{S}(\text{bt}) = \{a, b, ab, 0\}$ and $\mathbf{S}(\text{mc}) = \mathbf{S}(\text{pc}) = \{a, b, 0\}$. It says that the blood type of a person X depends on the inherited genetical information of X . Note

¹ For the sake of simplicity we consider finite random variables, i.e. random variables having a finite set \mathbf{S} of states. However, the ideas generalize to discrete and continuous random variables.

that the domain $\mathbf{S}(p)$ has nothing to do with the notion of a domain in the logical sense. The domain $\mathbf{S}(p)$ defines the states of random variables. Intuitively, a Bayesian predicate p generically represents a set of (finite) random variables. More precisely, each Bayesian ground atom g over p represents a (finite) random variable over the states $\mathbf{S}(g) := \mathbf{S}(p)$. E.g. $bt(ann)$ represents the blood type of a person named Ann as a random variable over the states $\{a, b, ab, 0\}$. Apart from that, most other *logical* notions carry over to Bayesian logic programs. So, we will speak of Bayesian predicates, terms, constants, substitutions, ground Bayesian clauses, Bayesian Herbrand interpretations etc. We will assume that all Bayesian clauses are range-restricted. A clause is *range-restricted* iff all variables occurring in the head also occur in the body. Range restriction is often imposed in the database literature; it allows one to avoid derivation of non-ground true facts.

In order to represent a probabilistic model we associate with each Bayesian clause c a conditional probability distribution $cpd(c)$ encoding $\mathbf{P}(head(c) \mid body(c))$. To keep the expositions simple, we will assume that $cpd(c)$ is represented as table, see Figure 1. More elaborate representations like decision trees or rules are also possible. The distribution $cpd(c)$ generically represents the conditional probability distributions of all ground instances $c\theta$ of the clause c . In general, one may have many clauses, e.g. clauses c_1 and the c_2

$$\begin{array}{l} bt(X) \mid mc(X) . \\ bt(X) \mid pc(X) . \end{array}$$

and corresponding substitutions θ_i that ground the clauses c_i such that $head(c_1\theta_1) = head(c_2\theta_2)$. They specify $cpd(c_1\theta_1)$ and $cpd(c_2\theta_2)$, but not the distribution required: $\mathbf{P}(head(c_1\theta_1) \mid body(c_1) \cup body(c_2))$. The standard solution to obtain the distribution required are so called *combining rules*.

Definition 2 (*Combining Rule*). A combining rule is a functions which maps finite sets of conditional probability distributions $\{\mathbf{P}(A \mid A_{i1}, \dots, A_{in_i}) \mid i = 1, \dots, m\}$ onto one (combined) conditional probability distribution $\mathbf{P}(A \mid B_1, \dots, B_k)$ with $\{B_1, \dots, B_k\} \subseteq \bigcup_{i=1}^m \{A_{i1}, \dots, A_{in_i}\}$.

We assume that for each Bayesian predicate p there is a corresponding combining rule cr , such as *noisy or* (see e.g. [27]) or *average*. The latter assumes $n_1 = \dots = n_m$ and $\mathbf{S}(A_{ij}) = \mathbf{S}(A_{kj})$, and computes the average of the distributions over $\mathbf{S}(A)$ for each joint state over $\bigotimes_j \mathbf{S}(A_{ij})$.

To summarize, we could define Bayesian logic program in the following way:

Definition 3 (*Bayesian Logic Program*). A Bayesian logic program B consists of a (finite) set of Bayesian clauses. To each Bayesian clause c there is exactly one conditional probability distribution $cpd(c)$ associated, and for each Bayesian predicate p there is exactly one associated combining rule $cr(p)$.

2.2 Declarative Semantics

The declarative semantics of Bayesian logic programs is given by the annotated *dependency graph*. The *dependency graph* $DG(B)$ is that directed graph whose nodes

```

m(ann,dorothy).
f(brian,dorothy).
pc(ann).
pc(brian).
mc(ann).
mc(brian).

```

```

mc(X) | m(Y,X),mc(Y),pc(Y).
pc(X) | f(Y,X),mc(Y),pc(Y).
bt(X) | mc(X),pc(X).

```

(1)

$mc(X)$	$pc(X)$	$\mathbf{P}(bt(X))$
a	a	(0.97, 0.01, 0.01, 0.01)
b	a	(0.01, 0.01, 0.97, 0.01)
...
0	0	(0.01, 0.01, 0.01, 0.97)

(2)

Fig. 1. (1) The Bayesian logic program *bloodtype* encoding our genetic domain. To each Bayesian predicate, the identity is associated as combining rule. (2) A conditional probability distribution associated to the Bayesian clause $bt(X) \mid mc(X), pc(X)$ represented as a table.

correspond to the ground atoms in the least Herbrand model $LH(B)$ (cf. below). It encodes the *directly influenced by* relation over the random variables in $LH(B)$:

there is an edge from a node x to a node y if and only if there exists a clause $c \in B$ and a substitution θ , s.t. $y = head(c\theta)$, $x \in body(c\theta)$ and for all atoms z in $c\theta$: $z \in LH(B)$.

The direct predecessors of a graph node x are called its parents, $\mathbf{Pa}(x)$. The Herbrand base $HB(B)$ is the set of all random variables we can talk about. It is defined as if B were a logic program (cf. [37]). The least Herbrand model $LH(B) \subseteq HB(B)$ consists of all *relevant* random variables, the random variables over which a probability distribution is well-defined by B , as we will see. It is the least fix point of the *immediate consequence operator* applied on the empty interpretation. Therefore, a ground atom which is true in the logical sense corresponds to a relevant random variables. Now, to each node x in $DG(B)$ we associate the combined conditional probability distribution which is the result of the combining rule $cr(p)$ of the corresponding Bayesian predicate p applied to the set of $cpd(c\theta)$'s where $head(c\theta) = x$ and $\{x\} \cup body(c\theta) \subseteq LH(B)$. Thus, if $DG(B)$ is acyclic and not empty then it encodes a (possibly infinite) Bayesian network, because the least Herbrand model always exists and is unique. Therefore, the following independence assumption holds:

Independence Assumption 1. *Each node x is independent of its non-descendants given a joint state of its parents $\mathbf{Pa}(x)$ in the dependency graph.*

E.g. in the program in Figure 1, the random variable $bt(dorothy)$ is independent from $pc(brian)$ given a joint state of $pc(dorothy), mc(dorothy)$. Using this assumption the following proposition holds:

Proposition 1. *Let B be a Bayesian logic program. If*

1. $LH(B) \neq \emptyset$,
2. $DG(B)$ is acyclic, and
3. each node in $DG(B)$ is influenced by a finite set of random variables

then B specifies a unique probability distribution \mathbf{P}_B over $LH(B)$.

```

m(ann,dorothy).
f(brian,dorothy).
pc(ann).
pc(brian).
mc(ann).
mc(brian).
mc(dorothy) | m(ann, dorothy),mc(ann),pc(ann).
pc(dorothy) | f(brian, dorothy),mc(brian),pc(brian).
bt(ann) | mc(ann), pc(ann).
bt(brian) | mc(brian), pc(brian).
bt(dorothy) | mc(dorothy),pc(dorothy).

```

Fig. 2. The grounded version of the Bayesian logic program of Figure 1. It (directly) encodes a Bayesian network.

The proof of the proposition can be sketched as follows (for a detailed proof see [29]). The least Herbrand $\text{LH}(B)$ always exists, is unique and countable. Thus, $DG(B)$ exists and is unique, and due to condition (3) the combined probability distribution for each node of $DG(B)$ is computable. Furthermore, because of condition (1) a total order π on $DG(B)$ exists, so that one can see B together with π as a stochastic process over $\text{LH}(B)$. An induction “along” π together with condition 2 shows that the family of finite-dimensional distribution of the process is projective (cf. [2]), i.e the joint probability density over each finite subset $s \subseteq \text{LH}(B)$ is uniquely defined and $\int_y p(s, x = y) dy = p(s)$. Thus, the preconditions of *Kolmogorov’s theorem* [2, page 307] hold, and it follows that B given π specifies a probability density function p over $\text{LH}(B)$. This proves the proposition because the total order π used for the induction is arbitrary.

A program B satisfying the conditions (1), (2) and (3) of proposition 1 is called *well-defined*. The program *bloodtype* in Figure 1 is an example of a well-defined Bayesian logic program. It encodes the regularities in our genetic example. Its grounded version, which is a Bayesian network, is shown in Figure 2. This illustrates that Bayesian networks [41] are well-defined propositional Bayesian logic programs. Each node-parents pair uniquely specifies a propositional Bayesian clause; we associate the identity as combining rule to each predicate; the conditional probability distributions are those of the Bayesian network.

Some interesting insights follow from the proof sketch. We interpreted a Bayesian logic program as a stochastic process. This places them in a wider context of what Cowell et al. call *highly structured stochastic systems* (HSSS), cf. [10], because Bayesian logic programs represent discrete-time stochastic processes in a more flexible manner. Well-known probabilistic frameworks such as dynamic Bayesian networks, first order hidden Markov models or Kalman filters are special cases of Bayesian logic programs. Moreover, the proof in [29] indicates the important *support network* concept. Support networks are a graphical representation of the finite-dimensional distribution, cf. [2], and are needed for the formulation of the likelihood function (see below) as well as for answering probabilistic queries in Bayesian logic programs.

Definition 4 (*Support Network*). The support network N of a variable $x \in \text{LH}(B)$ is defined as the induced subnetwork of $S = \{x\} \cup \{y \mid y \in \text{LH}(B) \text{ and } y \text{ is influencing } x\}$. The support network of a finite set $\{x_1, \dots, x_k\} \subseteq \text{LH}(B)$ is the union of the networks of each single x_i .

Because we consider well-defined Bayesian logic programs, each $x \in \text{LH}(B)$ is influenced by a finite subset of $\text{LH}(B)$. So, the support network N of a finite set $\{x_1, \dots, x_k\} \subseteq \text{LH}(B)$ of random variables is always a finite Bayesian network and computable in finite time. The distribution factorizes in the usual way, i.e. $\mathbf{P}_N(x_1 \dots, x_n) = \prod_{i=1}^n \mathbf{P}_N(x_i \mid \mathbf{Pa} x_i)$, where $\{x_1 \dots, x_n\} = S$, and $\mathbf{P}(x_i \mid \mathbf{Pa} x_i)$ is the combined conditional probability distribution associated to x_i . Because N models the finite-dimensional distribution specified by S , any interesting probability value over subsets of S is specified by N . For the proofs and an effective inference procedure (together with a Prolog implementation) we refer to [29].

3 The Learning Problem

So far, we have assumed that there is an expert who provides both the structure and the conditional probability distributions of the Bayesian logic program. This is not always easy. Often, there is no-one possessing necessary the expertise or knowledge. However, instead of an expert we may have access to data. In this section, we investigate and formally define the problem of learning Bayesian logic programs. While doing so, we exploit analogies with Bayesian network learning as well as with inductive logic programming.

3.1 Data Cases

In the last section, we have introduced Bayesian logic programs and argued that they contain two components, the quantitative (the combining rules and the conditional probability distributions) and the qualitative ones (the Bayesian clauses). Now, if we want to learn Bayesian logic programs, we need to employ data. Hence, we need to formally define the notions of a data case.

Let B be a Bayesian logic program consisting of the Bayesian clauses c_1, \dots, c_n , and let $\mathbf{D} = \{D_1, \dots, D_m\}$ be a set of data cases.

Definition 5 (*Data Case*). A data case $D_i \in \mathbf{D}$ for a Bayesian logic program B consists of a

Logical part: Which is a Herbrand interpretation $\text{Var}(D_i)$ such that $\text{Var}(D_i) = \text{LH}(B \cup \text{Var}(D_i))$, and a

Probabilistic part: Which is a partially observed joint state of some variables, i.e. an assignment of values to some of the facts in $\text{Var}(D_i)$.

Examples of data cases are

$$\begin{aligned}
 D_1 = \{ & m(\text{cecily}, \text{fred}) = \text{true}, f(\text{henry}, \text{fred}) = ?, pc(\text{cecily}) = a, \\
 & pc(\text{henry}) = b, pc(\text{fred}) = ?, mc(\text{cecily}) = b, mc(\text{henry}) = b, \\
 & mc(\text{fred}) = ?, bt(\text{cecily}) = ab, bt(\text{henry}) = b, bt(\text{fred}) = ?\},
 \end{aligned}$$

$$\begin{aligned}
D_2 = \{ & m(\text{ann}, \text{dorothy}) = \text{true}, f(\text{brian}, \text{dorothy}) = \text{true}, pc(\text{ann}) = b, \\
& mc(\text{ann}) = ?, mc(\text{brian}) = a, mc(\text{dorothy}) = a, \\
& pc(\text{dorothy}) = a, pc(\text{brian}) = ?, bt(\text{ann}) = ab, bt(\text{brian}) = ?, \\
& bt(\text{dorothy}) = a \},
\end{aligned}$$

where ‘?’ stands for an unobserved state. Notice that – for ease of writing – we merged the two components of a data case into one. Indeed, the *logical part* of a data case $D_i \in \mathbf{D}$, denoted as $Var(D_i)$, is a Herbrand interpretation, such as

$$\begin{aligned}
Var(D_1) = \{ & m(\text{cecily}, \text{fred}), f(\text{henry}, \text{fred}), pc(\text{cecily}), pc(\text{henry}), \\
& pc(\text{fred}), mc(\text{cecily}), mc(\text{henry}), mc(\text{fred}), bt(\text{cecily}), \\
& bt(\text{henry}), bt(\text{fred}) \}, \\
Var(D_2) = \{ & m(\text{ann}, \text{dorothy}), f(\text{brian}, \text{dorothy}), pc(\text{ann}), \\
& mc(\text{ann}), mc(\text{brian}), mc(\text{dorothy}), pc(\text{dorothy}), \\
& pc(\text{brian}), bt(\text{ann}), bt(\text{brian}), bt(\text{dorothy}) \},
\end{aligned}$$

satisfy this logical property w.r.t. the target Bayesian logic program B

$$\begin{array}{l|l}
mc(X) & | \ m(Y, X), mc(Y), pc(Y) . \\
pc(X) & | \ f(Y, X), mc(Y), pc(Y) . \\
bt(X) & | \ mc(X), pc(X) .
\end{array}$$

Indeed, $Var(B \cup Var(D_i)) = Var(D_i)$ for all $D_i \in \mathbf{D}$.

So, the logical components of the data cases should be seen as the least Herbrand models of the target Bayesian logic program. They specify different sets of *relevant* random variables, depending on the given “extensional context”. If we accept that the genetic laws are the same for both families then a learning algorithm should find regularities among the Herbrand interpretations that can be to compress the interpretations. The key assumption underlying any inductive technique is that the rules that are valid in one interpretation are likely to hold for any interpretation. This is exactly what the *learning from interpretations* in inductive logic programming [14,6] is doing. Thus, we will adapt this setting for learning the structure of the Bayesian logic program, cf. Section 4.

There is one further logical constraints to take into account while learning Bayesian logic programs. It is concerned with the acyclicity requirement (cf. property 2 in proposition 1) imposed on Bayesian logic programs. Thus, we require that for each $D_i \in \mathbf{D}$ the induced Bayesian network over $LH(B \cup Var(D_i))$ has to be acyclic.

At this point, the reader should also observe that we require that the logical part of a data case is a *complete* model of the target Bayesian logic program and not a *partial* one². This is motivated by 1) Bayesian network learning and 2) the problems with learning from partial models in inductive logic programming. First, data cases as they have been used in Bayesian network learning are the

² Partial models specify the truth-value (false or true) of *some* of the elements in the Herbrand Base.

propositional equivalent of the data cases that we introduced above. Indeed, if we have a Bayesian network B over the propositional Bayesian predicates $\{p_1, \dots, p_k\}$ then $\text{LH}(B) = \{p_1, \dots, p_k\}$ and a data case would assign values to some of the predicates in B . This also shows that the second component of a data case is pretty standard in the Bayesian network literature. Second, it is well-known that learning from partial models is harder than learning from complete models (cf. [12]). More specifically, learning from partial models is akin to multiple predicate learning, which is a very hard problem in general. These two points also clarify why the semantics of the set of relevant random variables coincided with the least Herbrand domain and at the same time why we do not restrict the domain of Bayesian predicates to $\{true, false\}$.

Before we are able to fully specify the problem of learning Bayesian logic programs, let us introduce the hypothesis space and scoring functions.

3.2 The Hypothesis Space

The *hypothesis space* \mathcal{H} explored consists of Bayesian logic programs, i.e. finite set of Bayesian clauses to which conditional probability distributions are associated. More formally, let \mathcal{L} be the language, which determines the set \mathcal{C} of clauses that can be part of a hypothesis. It is common to impose syntactic restrictions on the space \mathcal{H} of hypotheses.

Language Assumption: In this paper, we assume that the alphabet of \mathcal{L} only contains constant and predicate symbols that occur in one of the data cases, and we restrict \mathcal{C} to range-restricted, constant-free clauses containing maximum $k = 3$ atoms in the body. Furthermore, we assume that the combining rules associated to the Bayesian predicates are given.

E.g. given the data cases D_1 and D_2 , \mathcal{C} looks like

```

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

Not every element $H \in \mathcal{H}$ has to be a candidate. The logical parts of the data cases constraint the set of possible candidates. To be a candidate, H has to be

- (logically) valid on the data, and
- acyclic on the data i.e. the induced Bayesian network over $\text{LH}(H \cup \text{Var}(D_i))$ has to be acyclic.

E.g. given the data cases D_1 and D_2 , the Bayesian clause

```

mc(X) | mc(X)
    
```

is not included in any candidate, because the Bayesian network induced over the data cases would be cyclic.

3.3 Scoring Function

So far, we mainly exploit the logical part of the data cases. The probabilistic part of the data cases are partially observed joint states. They induce a joint distribution over the random variables of the logical parts of the data cases. A candidate $H \in \mathcal{H}$ should reflect this distribution. We assume that there is a scoring function $score_{\mathbf{D}} : \mathcal{H} \mapsto \mathbb{R}$ which expresses how well a given candidate H fits the data \mathbf{D} . Examples of scoring functions are the likelihood (see Section 5) or the *minimum description length* score (which bases on the likelihood).

Putting all together, we can define the basic learning problem as follows:

Definition 6 (*Learning Problem*). **Given** a set $\mathbf{D} = \{D_1, \dots, D_m\}$ of data cases, a set \mathcal{H} of sets of Bayesian clauses according to some language bias, and a scoring function $score_{\mathbf{D}} : \mathcal{H} \mapsto \mathbb{R}$, **find** a hypothesis $H^* \in \mathcal{H}$ such that for all $D_i \in \mathbf{D}$: $LH(H^* \cup Var(D_i)) = Var(D_i)$, H^* is acyclic on the data, and H^* maximizes $score_{\mathbf{D}}$.

As usual, we assume the all data cases are independently sampled from identical distributions. In the following section we will present an algorithm solving the learning problem.

4 Scooby: An Algorithm for Learning Intensional Bayesian Logic Programs

In this section, we introduce SCOOPY (structural learning of intensional Bayesian logic programs), cf. Algorithm 1. Roughly speaking, SCOOPY performs a heuristic search using traditional inductive logic programming refinement operators on clauses. The hypothesis currently under consideration is evaluated using some score as heuristic. The hypothesis that scores best is selected as the final hypothesis.

First, we will illustrate how SCOOPY works for the special case of Bayesian networks. As it will turn out, SCOOPY coincides in this case with well-known and effective score-based techniques for learning Bayesian networks [22]. Then, we will show that SCOOPY works for first-order Bayesian logic programs, too. For the sake of readability, we assume the existence of a method to compute the parameters maximizing the score given a candidate and data cases. Methods to do this will be discussed in Section 5. They assume that the combining rules are decomposable, a concept which we will introduce below. Furthermore we will discuss the basic framework only; extensions are possible.

4.1 The Propositional Case: Bayesian Networks

Let us first explain how SCOOPY works on Bayesian networks. and show that well-known score-based methods for structural learning of Bayesian networks are special cases of SCOOPY.

Let $\mathbf{x} = \{x_1, \dots, x_n\}$ be a fixed set of random variables. The set \mathbf{x} corresponds to a least Herbrand model of an unknown propositional Bayesian logic program

Algorithm 1. A simplified skeleton of a greedy algorithm for structural learning of intensional Bayesian logic programs (Scooby). Note that we have omitted the initialization of the conditional probability distributions associated with Bayesian clauses with random values. The operators ρ_g and ρ_s are generalization and specialization operators.

```

Let  $H$  be an initial (valid) hypothesis;
 $S(H) := score_{\mathbf{D}}(H)$ ;
repeat
   $H' := H$ ;
   $S(H') := S(H)$ ;
  foreach  $H'' \in \rho_g(H') \cup \rho_s(H')$  do
    if  $H''$  is (logically) valid on  $D$  then
      if the Bayesian networks induced by  $H''$  on the data are acyclic then
        if  $score_{\mathbf{D}}(H'') > S(H)$  then
           $H := H''$ ;
           $S(H) := S(H'')$ ;
until  $S(H') = S(H)$  ;
Return  $H$ ;

```

representing a Bayesian network. The probabilistic dependencies among the relevant random variables are not known, i.e. the propositional Bayesian clauses are unknown. Therefore, we have to select such a propositional Bayesian logic program as a *candidate* and estimate its parameters. Assume the data cases $\mathbf{D} = \{D_1, \dots, D_m\}$ look like

$$\begin{aligned}
 &\{m(ann, dorothy) = true, f(brian, dorothy) = true, pc(ann) = a, \\
 &mc(ann) = ?, mc(brian) = ?, mc(dorothy) = a, mc(dorothy) = a, \\
 &pc(brian) = b, bt(ann) = a, bt(brian) = ?, bt(dorothy) = a\}
 \end{aligned}$$

which is a data case for the Bayesian network in Figure 2. Note, that the atoms have to be interpreted as propositions. Each H in the hypothesis space \mathcal{H} is a Bayesian logic program consisting of n propositional clauses: for each $x_i \in \mathbf{x}$ a single clause c with $head(c) = x_i$ and $body(c) \subseteq \mathbf{x} \setminus \{x_i\}$. To traverse \mathcal{H} we specify two *refinement* operators $\rho_g : \mathcal{H} \mapsto 2^{\mathcal{H}}$ and $\rho_s : \mathcal{H} \mapsto 2^{\mathcal{H}}$, that take a hypothesis and modify it to produce a set of possible candidates. In the case of Bayesian networks the operator $\rho_g(H)$ deletes a Bayesian proposition from the body of a Bayesian clause $c_i \in H$, and the operator $\rho_s(H)$ adds a Bayesian proposition to the body of $c_i \in H$ (cf Figure 3). The search algorithm performs an greedy, informed search in \mathcal{H} based on $score_{\mathbf{D}}$.

As a simple illustration we consider a greedy hill-climbing algorithm incorporating $score_{\mathbf{D}}(H) := LL(\mathbf{D}, H)$, the log-likelihood of the data \mathbf{D} given a candidate structure H with the best parameters. We pick an initial candidate $S \in \mathcal{H}$

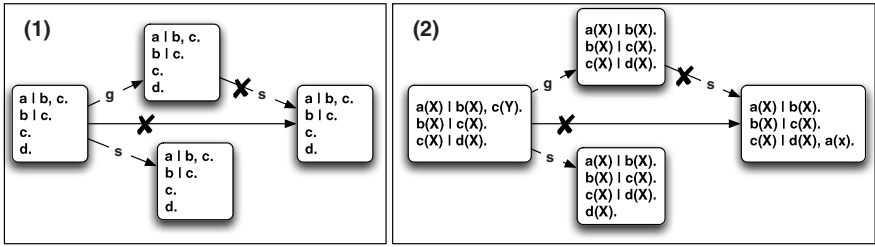


Fig. 3. (1) The use of refinement operators during structural search for Bayesian networks. We can add (ρ_s) a proposition to the body of a clause or delete (ρ_g) it from the body. (2) The use of refinement operators during structural search within the framework of Bayesian logic programs. We can add (ρ_s denoted as ‘s’) an atom to the body of a clause or delete (ρ_g denoted as ‘g’) it from the body. Candidates crossed out in (1) and (2) are illegal because they are cyclic.

as starting point (e.g. the set of all propositions) and compute the likelihood $LL(\mathbf{D}, S)$ with the best parameters. Then, we use $\rho(S)$ to compute the legal “neighbors” (candidates being acyclic) of S in \mathcal{H} and score them. All neighbors are valid (see below for a definition of validity). E.g. replacing $\text{pc}(\text{dorothy})$ with $\text{pc}(\text{dorothy}) \mid \text{pc}(\text{brian})$ gives such a “neighbor”. We take that $S' \in \rho(S)$ with the best improvements in the score. The process is continued until no improvements in score are obtained.

4.2 The First Order Case: Bayesian Logic Programs

Let us now explain how Scooby works in the first order case. The key differences with the propositional case are The key difference to the propositional case are

1. That some Bayesian logic programs will be logically invalid (see below for an example), and
2. That the traditional first order refinement operators must be used.

Difference 1 is the most important one, because it determines the hypotheses that are candidate Bayesian logic programs. To account for this difference, two modifications of the traditional Bayesian network algorithm are needed.

The first modification concerns the initialization phase where we have to choose a logically valid, acyclic Bayesian logic program. Such a program can be computed using a CLAUDIEN like procedure ([13,14,6]). CLAUDIEN is an ILP-program that computes a logically valid hypothesis H from a set of data cases. Furthermore, all clauses in H will be maximally general (w.r.t. θ -subsumption), and CLAUDIEN will compute all such clauses (within \mathcal{L}). This means none of the clauses in H can be generalized without violating the logical validity requirement (or leaving \mathcal{L}). Consider again the data cases

$$\begin{aligned}
 D_1 &= \{m(\text{cecily}, \text{fred}) = \text{true}, f(\text{henry}, \text{fred}) = ?, pc(\text{cecily}) = a, \\
 &\quad pc(\text{henry}) = b, pc(\text{fred}) = ?, mc(\text{cecily}) = b, mc(\text{henry}) = b, \\
 &\quad mc(\text{fred}) = ?, bt(\text{cecily}) = ab, bt(\text{henry}) = b, bt(\text{fred}) = ?\}, \\
 D_2 &= \{m(\text{ann}, \text{dorothy}) = \text{true}, f(\text{brian}, \text{dorothy}) = \text{true}, pc(\text{ann}) = b, \\
 &\quad mc(\text{ann}) = ?, mc(\text{brian}) = a, mc(\text{dorothy}) = a, \\
 &\quad pc(\text{dorothy}) = a, pc(\text{brian}) = ?, bt(\text{ann}) = ab, bt(\text{brian}) = ?, \\
 &\quad bt(\text{dorothy}) = a\},
 \end{aligned}$$

The clause $bt(X)$ is not a member of \mathcal{L} . The clause $bt(X) \mid mc(X), pc(X)$ is valid but not maximally general because the literal $pc(X)$ can be deleted without violating the logical validity requirement. Any hypothesis including $m(X, Y) \mid mc(X), pc(Y)$ would be logically invalid because *cecily* is not the mother of *henry*. Examples of maximally general clauses are

$$\begin{aligned}
 mc(X) &\mid m(Y, X). \\
 pc(X) &\mid f(Y, X). \\
 bt(X) &\mid mc(X). \\
 bt(X) &\mid pc(X). \\
 \dots &
 \end{aligned}$$

Roughly speaking, CLAUDIEN works as follows (for a detailed discussion we refer to [14]). It keeps track of a list of candidate clauses Q , which is initialized to the maximally general clause (in \mathcal{L}). It repeatedly deletes a clause c from Q , and tests whether c is valid on the data. If it is, c is added to the final hypothesis, otherwise, all maximally general specializations of c (in \mathcal{L}) are computed (using a so-called refinement operator ρ , see below) and added back to Q . This process continues until Q is empty and all relevant parts of the search space have been considered. The clauses generated by CLAUDIEN can be used as an initial hypothesis.

In the experiments, for each predicate, we selected one of the clause generated by CLAUDIEN for inclusion in the initial hypothesis such that the valid Bayesian logic program was also acyclic on the data cases (see below). An initial hypothesis is e.g.

$$\begin{aligned}
 mc(X) &\mid m(Y, X). \\
 pc(X) &\mid f(Y, X). \\
 bt(X) &\mid mc(X).
 \end{aligned}$$

The second modification concerns filtering out those Bayesian logic programs that are logically invalid during search. This is realized by the first **if**-condition in the loop. The second **if**-condition tests whether cyclic dependencies are induced on the data cases. This can be done in time $O(s \cdot r^3)$ where r is the number of random variables of the largest data case in \mathbf{D} and s is the number of clauses in H . To do so, we build the Bayesian networks induced by H over each $Var(D_i)$ by computing the ground instances for each clause $c \in H$ where the ground atoms are members of $Var(D_i)$. Thus, ground atoms, which are not appearing

```

m(ann,dorothy).   m(cecily,fred).
f(brian,dorothy). f(henry,fred).
pc(ann).   pc(brian).   pc(cecily).   pc(henry).
mc(ann).   mc(brian).   mc(cecily).   mc(henry).
mc(dorothy) | m(ann,dorothy).   mc(fred) | m(cecily,fred).
pc(dorothy) | f(brian,dorothy).   pc(fred) | f(cecily,fred).
bt(ann)      | mc(ann).      bt(brian)   | mc(brian).
bt(dorothy) | mc(dorothy).   bt(cecily) | mc(cecily).
bt(henry)   | mc(henry).     bt(fred)   | mc(fred).

```

Fig. 4. The support network induced by the initial hypothesis S (see text) over the the data cases D_1 and D_2

as a head atom of a valid ground instance, are a priori nodes, i.e. nodes with an empty parent set. This takes $O(s \cdot r_i^3)$. Then, we test in $O(r_i)$ for a topological order of the nodes in the induced Bayesian network. If it exists, then the Bayesian network is acyclic. Otherwise, it is cyclic. Figure 4 shows the support network induced by the initial hypothesis over D_1 and D_2 .

For Difference 2, i.e. the refinements operators, we employ the traditional ILP refinement operators. In our approach we use the two refinement operators $\rho_s : 2^{\mathcal{H}} \mapsto \mathcal{H}$ and $\rho_g : 2^{\mathcal{H}} \mapsto \mathcal{H}$. The operator $\rho_s(H)$ adds constant-free atoms to the body of a single clause $c \in H$, and $\rho_g(H)$ deletes constant-free atoms from the body of a single clause $c \in H$. Figure 3 shows the different refinement operators for the first order case and the propositional case for learning Bayesian networks. Instead of adding (deleting) propositions to (from) the body of a clause, they add (delete) according to our language assumption constant-free atoms. Furthermore, Figure 3 shows that using the refinement operators each hypothesis can in principle be reached.

Finally, we need to mention that whereas the maximal general clauses are the most interesting ones from the logical point of view, this is not necessarily the case from the probabilistic point of view. E.g. having data cases D_1 and D_2 (see Section 3.1), the initial candidate S

```

mc(X)   | m(Y, X).
pc(X)   | f(Y, X).
bt(X)   | mc(X).

```

is likely not to score maximally on the data cases. E.g. the blood type does not depend on the fatherly genetical information.

As a simple instantiation of Algorithm 1, we consider a greedy hill-climbing algorithm incorporating $score_{\mathbf{D}}(H) := LL(\mathbf{D}, H)$ with $\mathbf{D} = \{D_1, D_2\}$. It takes $S \in \mathcal{H}$ (see above) as starting point and computes $LL(\mathbf{D}, S)$ with the best parameters. Then, we use $\rho_s(S)$ and $\rho_g(S)$ to compute the legal “neighbors” of S in \mathcal{H} and score them. E.g. one such a “neighbor” is given by replacing $bt(X) \mid mc(X)$ with $bt(X) \mid mc(X), pc(X)$. Let S' be that valid and acyclic neighbor which scores best. If $LL(\mathbf{D}, S) < LL(\mathbf{D}, S')$, then we take S' as new hypothesis. The process is continued until no improvements in score are obtained.

4.3 Discussion

The algorithm presented serves as a basic, unifying framework. Several extensions and modifications based on ideas developed in both fields, inductive logic programming and Bayesian networks are possible. These include: lookaheads, background knowledge, mode declarations and improved scoring functions. Let us briefly address some of these.

Lookahead: In some cases, an atom might never be chosen by our algorithm because it will not – in itself – result in a better score. However, such an atom, while not useful in itself, might introduce new variables that make a better score possible by adding another atom later on. Within inductive logic programming this is solved by allowing the algorithm to *look ahead* in the search space. Immediately after refining a clause by putting some atom A into the body, the algorithm checks whether any other atom involving some variable of A results in a better score [5]. The same problem is encountered when learning Bayesian networks [47].

Background Knowledge: Inductive logic programming emphasizes background knowledge, i.e. predefined, fixed regularities which are common to all examples. Background knowledge can be incorporated into our approach in the following way. It is expressed as a fixed Bayesian logic program BK . Now, we search for a candidate H^* which is together with BK acyclic on the data such that for all $D_i \in \mathbf{D} : \text{LH}(BK \cup H^* \cup \text{Var}(D_i)) = \text{Var}(D_i)$, and $BK \cup H^*$ matches the data \mathbf{D} best according to $\text{score}_{\mathbf{D}}$. Therefore, all the Bayesian facts that can be derived from the background knowledge and an example are part of the corresponding “extended” example. This is particularly interesting to specify deterministic knowledge as in inductive logic programming. In [29], we showed how pure Prolog programs can be represented as Bayesian logic programs w.r.t. the conditions 1,2 and 3 of Proposition 1.

Improved Scoring Function: Using the likelihood directly as scoring function, score-based algorithm to learn Bayesian networks prefer fully connected networks. To overcome the problem advanced scoring functions were developed. One of these is the *minimum description length* (MDL) score which trades off the fit to the data with complexity. In the context of learning Bayesian networks, the whole Bayesian network is encoded to measure the compression [34]. In the context of learning clause programs, other compression measure were investigated such as the average length of proofs [44]. For Bayesian logic programs, a combination of both seems to be appropriate.

Finally, an extension for learning predicate definitions consisting of more than one clause is in principle possible. The refinement operators could be modified in such a way that for a clause $c \in H'$ with head predicate p another (valid) clause c' (e.g. computed by CLAUDIEN) with head predicate p is added or deleted.

5 Learning Probabilities in a Bayesian Logic Program

So far, we have assumed that there is a method estimating the parameters of an candidate program given data. In this section, we show how to learn the quantitative component of a Bayesian logic program, i.e. the conditional probability distributions. The learning problem can be stated as follows:

Definition 7 (*Parameter Estimation*). **Given** a set $\mathbf{D} = \{D_1, \dots, D_m\}$ of data cases³, a set H of Bayesian clauses according to some language bias, which is logically valid and acyclic on the data, and a scoring function $\text{score}_{\mathbf{D}} : \mathcal{H} \mapsto \mathbb{R}$, **find** the parameters of H maximizing $\text{score}_{\mathbf{D}}$.

We will concentrate on *maximum likelihood estimation* (MLE).

5.1 Maximum Likelihood Estimation

Maximum likelihood is a classical method for parameter estimation. The likelihood is the probability of the observed data as a function of the unknown parameters with respect to the current model. Let B be a Bayesian logic program consisting of the Bayesian clauses c_1, \dots, c_n , and let $\mathbf{D} = \{D_1, \dots, D_m\}$ be a set of data cases. The parameters $\text{cpd}(c_i)_{jk} = P(u_j \mid \mathbf{u}_k)$, where $u_j \in \mathbf{S}(\text{head}(c_i))$ and $\mathbf{u}_k \in \mathbf{S}(\text{body}(c_i))$, affecting the associated conditional probability distributions $\text{cpd}(c_i)$ constitute the set $\boldsymbol{\lambda} = \bigcup_{i=1}^n \text{cpd}(c_i)$. The version of B where the parameters are set to $\boldsymbol{\lambda}$ is denoted by $B(\boldsymbol{\lambda})$, and as long as no ambiguities occur we will not distinguish between the parameters $\boldsymbol{\lambda}$ themselves and a particular instance of them.

Now, the likelihood $L(\mathbf{D}, \boldsymbol{\lambda})$ is the probability of the data \mathbf{D} as a function of the unknown parameters $\boldsymbol{\lambda}$:

$$L(\mathbf{D}, \boldsymbol{\lambda}) := P_B(\mathbf{D} \mid \boldsymbol{\lambda}) = P_{B(\boldsymbol{\lambda})}(\mathbf{D}). \quad (1)$$

Thus, the search space \mathcal{H} is spanned by the product space over the possible values of $\lambda(c_i)$ and we seek to find the parameter values $\boldsymbol{\lambda}^*$ that maximize the likelihood, i.e.

$$\boldsymbol{\lambda}^* = \max_{\boldsymbol{\lambda} \in \mathcal{H}} P_{B(\boldsymbol{\lambda})}(\mathbf{D}).$$

Usually, B specifies a distribution over a (countably) infinite set of random variables namely $\text{LH}(B)$ and hence we cannot compute $P_{B(\boldsymbol{\lambda})}(\mathbf{D})$ by considering the whole dependency graph. But as we have argued it is sufficient to consider the support network $N(\boldsymbol{\lambda})$ of the random variables occurring in \mathbf{D} to compute $P_{B(\boldsymbol{\lambda})}(\mathbf{D})$. Thus, using the monotonicity of the logarithm, we seek to find

$$\boldsymbol{\lambda}^* = \max_{\boldsymbol{\lambda} \in \mathcal{H}} \log P_{N(\boldsymbol{\lambda})}(\mathbf{D}) \quad (2)$$

³ Given a well-defined Bayesian network B , the logical part of a data case D_i can also be a partial model only if we only estimate the parameters and do not learn the structure, i.e. $\text{RandVar}(D_i) \subseteq \text{LH}(B)$. The given Bayesian logic program will fill in the missing random variables.

where $P_{N(\lambda)}$ is the probability distribution specified by the support network $N(\lambda)$ of the random variables occurring in \mathbf{D} . Equation (2) expresses the original problem in terms of the maximum likelihood parameter estimation problem of Bayesian networks:

A Bayesian logic program together with data cases induces a Bayesian network over the variables of the data cases.

This is not surprising because the learning setting is an instance of the probabilistic learning from interpretations. More important, due to the reduction, all techniques for maximum likelihood parameter estimation within Bayesian networks are in principle applicable. We only need to take the following issues into account:

1. Some of the nodes in $N(\lambda)$ are hidden, i.e., their values are not observed in \mathbf{D} .
2. We are not interested in the conditional probability distributions associated to ground instances of Bayesian clauses, but in those associated to the Bayesian clauses themselves.
3. Not only $L(\mathbf{D}, \lambda)$ but also $N(\lambda)$ itself depends on the data, i.e. the data cases determine the subnetwork of $DG(B)$ that is sufficient to calculate the likelihood.

The available data cases may not be complete, i.e., some values may not be observed. For instance in medical domains, a patient rarely gets all of the possible tests. In presence of missing data, the maximum likelihood estimate typically cannot be written in closed form. Unfortunately, it is a numerical optimization problem, and all known algorithms involve nonlinear, iterative optimization and multiple calls to a Bayesian inference procedures as subroutines, which are typically computationally infeasible. For instance the inference within Bayesian network has been proven to be NP-hard [9]. Typical ML parameter estimation techniques (in the presence of missing data) are the Expectation-Maximization (EM) algorithm and gradient-based approaches. We will now discuss both approaches in turn.

5.2 Gradient-Based Approach

We will adapt Binder *et al.*'s solution for dynamic Bayesian networks based on the chain rule of differentiation [3]. For simplicity, we fix the current instantiation of the parameters λ and, hence, we write B and $N(\mathbf{D})$. Applying the chain rule to (2) yields

$$\frac{\partial \log P_N(\mathbf{D})}{\partial \text{cpd}(c_i)_{jk}} = \sum_{\substack{\text{subst. } \theta \text{ s.t.} \\ \text{sn}(c_i\theta)}} \frac{\partial \log P_N(\mathbf{D})}{\partial \text{cpd}(c_i\theta)_{jk}} \quad (3)$$

where θ refers to grounding substitutions and $\text{sn}(c_i\theta)$ is true iff $\{\text{head}(c_i\theta)\} \cup \text{body}(c_i\theta) \subset N$. Assuming that the data cases $D_l \in \mathbf{D}$ are independently sampled from the same distribution we can separate the contribution of the different data cases to the partial derivative of a single ground instance $c\theta$:

$$\begin{aligned}
 \frac{\partial \log P_N(\mathbf{D})}{\partial \text{cpd}(c_i\theta)_{jk}} &= \frac{\partial \log \prod_{l=1}^m P_N(D_l)}{\partial \text{cpd}(c_i\theta)_{jk}} && \text{by independence} \\
 &= \sum_{l=1}^m \frac{\partial \log P_N(D_l)}{\partial \text{cpd}(c_i\theta)_{jk}} && \text{by } \log \prod = \sum \log \\
 &= \sum_{l=1}^m \frac{\partial P_N(D_l) / \partial \text{cpd}(c_i\theta)_{jk}}{P_N(D_l)}. && (4)
 \end{aligned}$$

In order to obtain computations local to the parameter $\text{cpd}(c_i\theta)_{jk}$ we introduce the variables $\text{head}(c_i\theta)$ and $\text{body}(c_i\theta)$ into the numerator of the summand of (4) and average over their possible values, i.e.,

$$\frac{\partial P_N(D_l)}{\partial \text{cpd}(c_i\theta)_{jk}} = \frac{\partial}{\partial \text{cpd}(c_i\theta)_{jk}} \left(\sum_{j',k'} P_N(D_l, \text{head}(c_i\theta) = u_{j'}, \text{body}(c_i\theta) = \mathbf{u}_{k'}) \right)$$

Applying the chain rule yields

$$\begin{aligned}
 \frac{\partial P_N(D_l)}{\partial \text{cpd}(c_i\theta)_{jk}} &= \frac{\partial}{\partial \text{cpd}(c_i\theta)_{jk}} \left(\sum_{j',k'} P_N(D_l \mid \text{head}(c_i\theta) = u_{j'}, \text{body}(c_i\theta) = \mathbf{u}_{k'}) \right. \\
 &\quad \left. \cdot P_N(\text{head}(c_i\theta) = u_{j'}, \text{body}(c_i\theta) = \mathbf{u}_{k'}) \right) \\
 &= \frac{\partial}{\partial \text{cpd}(c_i\theta)_{jk}} \left(\sum_{j',k'} P_N(D_l \mid \text{head}(c_i\theta) = u_{j'}, \text{body}(c_i\theta) = \mathbf{u}_{k'}) \right. \\
 &\quad \cdot P_N(\text{head}(c_i\theta) = u_{j'} \mid \text{body}(c_i\theta) = \mathbf{u}_{k'}) \\
 &\quad \left. \cdot P_N(\text{body}(c_i\theta) = \mathbf{u}_{k'}) \right) && (5)
 \end{aligned}$$

where $u_j \in \mathbf{S}(\text{head}(c_i))$, $\mathbf{u}_k \in \mathbf{S}(\text{body}(c_i))$ and j, k refer to the corresponding entries in $\text{cpd}(c_i)$, respectively $\text{cpd}(c_i\theta)$. In (5), $\text{cpd}(c_i\theta)_{jk}$ appears only in linear form. Moreover, it appears only when $j' = j$, and $k' = k$. Therefore, (5) simplifies two

$$\frac{\partial P_N(D_l)}{\partial \text{cpd}(c_i\theta)_{jk}} = P_N(D_l \mid \text{head}(c_i\theta) = u_j, \text{body}(c_i\theta) = \mathbf{u}_k) \cdot P_N(\text{body}(c_i\theta) = \mathbf{u}_k). \tag{6}$$

Substituting (6) back into (4) yields

$$\begin{aligned}
 &\sum_{l=1}^m \frac{\partial \log P_N(D_l) / \partial \text{cpd}(c_i\theta)_{jk}}{P_N(D_l)} \\
 &= \sum_{l=1}^m \frac{P_N(D_l \mid \text{head}(c_i\theta) = u_j, \text{body}(c_i\theta) = \mathbf{u}_k) \cdot P_N(\text{body}(c_i\theta) = \mathbf{u}_k)}{P_N(D_l)}
 \end{aligned}$$

$$\begin{aligned}
 &= \sum_{l=1}^m \frac{P_N(\text{head}(c_i\theta) = u_j, \text{body}(c_i\theta) = \mathbf{u}_k \mid D_l) \cdot P_N(D_l) \cdot P_N(\text{body}(c_i\theta) = \mathbf{u}_k)}{P_N(\text{head}(c_i\theta) = u_j, \text{body}(c_i\theta) = \mathbf{u}_k) \cdot P_N(D_l)} \\
 &= \sum_{l=1}^m \frac{P_N(\text{head}(c_i\theta) = u_j, \text{body}(c_i\theta) = \mathbf{u}_k \mid D_l)}{P_N(\text{head}(c_i\theta) = u_j \mid \text{body}(c_i\theta) = \mathbf{u}_k)} \\
 &= \sum_{l=1}^m \frac{P_N(\text{head}(c_i\theta) = u_j, \text{body}(c_i\theta) = \mathbf{u}_k \mid D_l)}{\text{cpd}(c_i\theta)_{jk}}.
 \end{aligned}$$

Combining all these, (3) can be rewritten as

$$\frac{\partial \log P_N(\mathbf{D})}{\partial \text{cpd}(c_i)_{jk}} = \sum_{\substack{\text{subst. } \theta \text{ with} \\ \text{sn}(c_i\theta)}} \frac{\text{en}(c_{ijk} \mid \theta, \mathbf{D})}{\text{cpd}(c_i\theta)_{jk}} \quad (7)$$

where

$$\begin{aligned}
 \text{en}(c_{ijk} \mid \theta, \mathbf{D}) &:= \text{en}(\text{head}(c_i\theta) = u_j, \text{body}(c_i\theta) = \mathbf{u}_k \mid \mathbf{D}) \\
 &:= \sum_{l=1}^m P_N(\text{head}(c_i\theta) = u_j, \text{body}(c_i\theta) = \mathbf{u}_k \mid D_l)
 \end{aligned} \quad (8)$$

are the so-called *expected counts* of the joint state $\text{head}(c_i\theta) = u_j, \text{body}(c_i\theta) = \mathbf{u}_k$ given the data \mathbf{D} .

Equation (7) shows that $P_N(\text{head}(c_i\theta) = u_j, \text{body}(c_i\theta) = \mathbf{u}_k \mid D_l)$ is all what is needed. This can essentially be computed using any standard Bayesian network inference engine. This is not surprising because (7) differs from the one for Bayesian networks given in [3] only in that we sum over all ground instances of a Bayesian clause holding in the data. To stress this close relationship, we rewrite (7) in terms of expected counts of clauses instead of ground clauses. They are defined as follows:

Definition 8 (Expected Counts of Bayesian Clauses). *The expected counts of a Bayesian clauses c of a Bayesian logic program B for a data set \mathbf{D} are defined as*

$$\begin{aligned}
 \text{en}(c_{ijk} \mid \mathbf{D}) &:= \text{en}(\text{head}(c_i) = u_j, \text{body}(c_i) = \mathbf{u}_k \mid \mathbf{D}) \\
 &:= \sum_{\substack{\text{subst. } \theta \text{ with} \\ \text{sn}(c_i\theta)}} \text{en}(\text{head}(c_i\theta) = u_j, \text{body}(c_i\theta) = \mathbf{u}_k \mid \mathbf{D}). \quad (9)
 \end{aligned}$$

Reading (7) in terms of Definition 8 proves the following proposition:

Proposition 1 (Partial Derivative of Log-Likelihood). *Let B be a Bayesian logic program with parameter vector λ . The partial derivative of the log-likelihood of B with respect to $\text{cpd}(c_i)_{jk}$ for a given data set \mathbf{D} is*

$$\frac{\partial LL(\mathbf{D}, \lambda)}{\partial \text{cpd}(c_i)_{jk}} = \frac{\text{en}(c_{ijk} \mid \mathbf{D})}{\text{cpd}(c_i)_{jk}}. \quad (10)$$

Algorithm 2. A simplified skeleton of the algorithm for *adaptive Bayesian logic programs* estimating the parameters of a Bayesian logic program

input : B , a Bayesian logic program; associated cpds are parameterized by λ ;
 \mathbf{D} , a finite set of data cases
output: a modified Bayesian logic program

$\lambda \leftarrow \text{INITIALPARAMETERS}$
 $N \leftarrow \text{SUPPORTNETWORK}(B, \mathbf{D})$

repeat

- $\Delta\lambda \leftarrow 0$
- set associated conditional probability distribution of N according to λ
- foreach** $D_l \in \mathbf{D}$ **do**
 - set the evidence in N from D_l
 - foreach** *Bayesian clause* $c \in B$ **do**
 - foreach** *ground instance* $c\theta$ s.t. $\{\text{head}(c\theta)\} \cup \text{body}(c\theta) \subset N$ **do**
 - foreach** *single parameter* $\text{cpd}(c\theta)_{jk}$ **do**
 - $\Delta \text{cpd}(c)_{jk} \leftarrow \Delta \text{cpd}(c)_{jk} + (\partial \log P_N(D_l) / \partial \text{cpd}(c\theta)_{jk})$
- $\Delta\lambda \leftarrow \text{PROJECTIONONTOCONSTRAINTSURFACE}(\Delta\lambda)$
- $\lambda \leftarrow \lambda + \alpha \cdot \Delta\lambda$

until $\Delta\lambda \approx 0$
return B

Equation (10) can be viewed as the first-order logical equivalent of the Bayesian network formula. A simplified skeleton of a gradient-based algorithm employing (10) is shown in Algorithm 2.

Before showing how to adapt the EM algorithm, we have to explain two points, which we have left out so far for the sake of simplicity: Constraint satisfaction and decomposable combining rules.

In the problem at hand, the gradient ascent has to be modified to take into account the constraint that the parameter vector λ consists of probability values, i.e. $\text{cpd}(c_i)_{jk} \in [0, 1]$ and $\sum_j \text{cpd}(c_i)_{jk} = 1$. Following [3], there are two ways to enforce this:

1. Projecting the gradient onto the constraint surface (as used to formulate the Algorithm 2), and
2. Reparameterizing the problem.

In the experiments, we chose the reparameterization approach because the new parameters automatically respect the constraints on $\text{cpd}(c_i)_{jk}$ no matter what their values are. More precisely, we define the parameters β with $\beta_{ijk} \in \mathbb{R}$ such that

$$\text{cpd}(c_i)_{jk} = \frac{e^{\beta_{ijk}}}{\sum_l e^{\beta_{ilk}}} \quad (11)$$

where the β_{ijk} are indexed like $\text{cpd}(c_i)_{jk}$. This enforces the constraints given above, and a local maximum with respect to the β is also a local maximum with

respect to λ , and vice versa. The gradient with respect to β can be found by computing the gradient with respect to λ and then deriving the gradient with respect to β using the chain rule of derivatives. More precisely, the chain rule of derivatives yields

$$\frac{\partial LL(\mathbf{D}, \lambda)}{\partial \beta_{ijk}} = \sum_{i'j'k'} \frac{\partial LL(\mathbf{D}, \lambda)}{\partial \text{cpd}(c_{i'})_{j'k'}} \cdot \frac{\partial \text{cpd}(c_{i'})_{j'k'}}{\partial \beta_{ijk}} \quad (12)$$

Since $\partial \text{cpd}(c_{i'})_{j'k'}/\partial \beta_{ijk} = 0$ for all $i \neq i'$, and $k \neq k'$, (12) simplifies to

$$\frac{\partial LL(\mathbf{D}, \lambda)}{\partial \beta_{ijk}} = \sum_{j'} \frac{\partial LL(\mathbf{D}, \lambda)}{\partial \text{cpd}(c_i)_{j'k}} \cdot \frac{\partial \text{cpd}(c_i)_{j'k}}{\partial \beta_{ijk}}$$

The quotient rule yields

$$\begin{aligned} \frac{\partial LL(\mathbf{D}, \lambda)}{\partial \beta_{ijk}} &= \sum_{j'} \left\{ \frac{\partial LL(\mathbf{D}, \lambda)}{\partial \text{cpd}(c_i)_{j'k}} \cdot \frac{\left(\frac{\partial e^{\beta_{ij'k}}}{\partial \beta_{ijk}} \cdot \sum_l e^{\beta_{ilk}} \right) - \left(e^{\beta_{ij'k}} \cdot \frac{\partial \sum_l e^{\beta_{ilk}}}{\partial \beta_{ijk}} \right)}{\left(\sum_l e^{\beta_{ilk}} \right)^2} \right\} \\ &= \frac{\left\{ \sum_{j'} \left(\frac{\partial LL(\mathbf{D}, \lambda)}{\partial \text{cpd}(c_i)_{j'k}} \cdot \frac{\partial e^{\beta_{ij'k}}}{\partial \beta_{ijk}} \cdot \sum_l e^{\beta_{ilk}} \right) \right\}}{\left(\sum_l e^{\beta_{ilk}} \right)^2} - \frac{\left\{ \sum_{j'} \left(\frac{\partial LL(\mathbf{D}, \lambda)}{\partial \text{cpd}(c_i)_{j'k}} \cdot e^{\beta_{ij'k}} \cdot \frac{\partial \sum_l e^{\beta_{ilk}}}{\partial \beta_{ijk}} \right) \right\}}{\left(\sum_l e^{\beta_{ilk}} \right)^2} \end{aligned}$$

Because $\partial e^{\beta_{ij'k}}/\partial \beta_{ijk} = 0$ for $j' \neq j$ and $\partial e^{\beta_{ijk}}/\partial \beta_{ijk} = e^{\beta_{ijk}}$, this simplifies to

$$\begin{aligned} \frac{\partial LL(\mathbf{D}, \lambda)}{\partial \beta_{ijk}} &= \frac{\left(\frac{\partial LL(\mathbf{D}, \lambda)}{\partial \text{cpd}(c_i)_{jk}} \cdot e^{\beta_{ijk}} \cdot \sum_l e^{\beta_{ilk}} \right) - \sum_{j'} \left(\frac{\partial LL(\mathbf{D}, \lambda)}{\partial \text{cpd}(c_i)_{j'k}} \cdot e^{\beta_{ij'k}} \cdot e^{\beta_{ilk}} \right)}{\left(\sum_l e^{\beta_{ilk}} \right)^2} \\ &= \frac{e^{\beta_{ijk}}}{\left(\sum_l e^{\beta_{ilk}} \right)^2} \cdot \left\{ \frac{\partial LL(\mathbf{D}, \lambda)}{\partial \text{cpd}(c_i)_{jk}} \cdot \left(\sum_l e^{\beta_{ilk}} \right) - \sum_{j'} \frac{\partial LL(\mathbf{D}, \lambda)}{\partial \text{cpd}(c_i)_{j'k}} \cdot e^{\beta_{ij'k}} \right\} \quad (13) \end{aligned}$$

To further simplify the partial derivative, we note that $\partial LL(\mathbf{D}, \lambda)/\partial \text{cpd}(c_i)_{jk}$ can be rewritten as

$$\frac{\partial LL(\mathbf{D}, \lambda)}{\partial \text{cpd}(c_i)_{jk}} = \frac{\text{en}(c_{ijk} \mid \mathbf{D})}{\text{cpd}(c_i)_{jk}} = \frac{\text{en}(c_{ijk} \mid \mathbf{D})}{\frac{e^{ijk}}{\sum_l e^{\beta_{ijk}}}} = \frac{\text{en}(c_{ijk} \mid \mathbf{D})}{e^{\beta_{ijk}}} \cdot \left(\sum_l e^{\beta_{ijk}} \right)$$

by substituting (11) in (9). Using the last equation, (13) simplifies to

$$\begin{aligned} &\frac{\partial LL(\mathbf{D}, \lambda)}{\partial \beta_{ijk}} \\ &= \frac{e^{\beta_{ijk}}}{\left(\sum_l e^{\beta_{ilk}} \right)^2} \cdot \left\{ \frac{\text{en}(c_{ijk} \mid \mathbf{D})}{e^{\beta_{ijk}}} \cdot \left(\sum_l e^{\beta_{ilk}} \right)^2 - \sum_{j'} \frac{\text{en}(c_{ij'k} \mid \mathbf{D})}{e^{\beta_{ij'k}}} \cdot \left(\sum_l e^{\beta_{ilk}} \right) \cdot e^{\beta_{ij'k}} \right\} \\ &= \text{en}(c_{ijk} \mid \mathbf{D}) - \frac{e^{\beta_{ijk}}}{\sum_l e^{\beta_{ilk}}} \cdot \sum_{j'} \text{en}(c_{ij'k} \mid \mathbf{D}). \end{aligned}$$

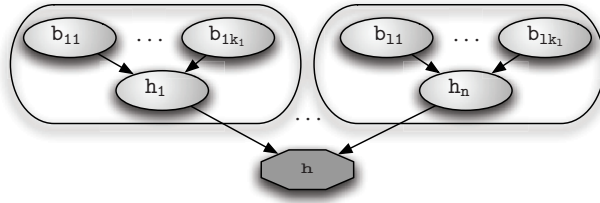


Fig. 5. The scheme of decomposable combining rules. Each rectangle corresponds to a ground instance $c\theta \equiv \mathbf{h}_i | \mathbf{b}_{11}, \dots, \mathbf{b}_{k_i}$ of a Bayesian clause $c \equiv \mathbf{h} | \mathbf{b}_1, \dots, \mathbf{b}_k$. The node \mathbf{h} is a deterministic node, i.e., its state is deterministic function of the parents joint state.

Using once more (11), the following proposition is proven:

Proposition 2 (Partial Derivative of Log-Likelihood of an Reparameterized BLP). *Let B be a Bayesian logic program reparameterized according to (11). The partial derivative of the log-likelihood of B with respect to β_{ijk} for a given data set \mathbf{D} is*

$$\frac{\partial LL(\mathbf{D}, \lambda)}{\partial \beta_{ijk}} = \text{en}(c_{ijk} | \mathbf{D}) - \text{cpd}(c_i)_{jk} \sum_{j'} \text{en}(c_{ij'k} | \mathbf{D}) . \quad (14)$$

Equation (14) shows that the partial derivative can be expressed solely in terms of expected counts and original parameters. Consequently, its computational complexity is linear in (10).

We assumed *decomposable* combining rules.

Definition 9 (Decomposable Combining Rule). *Decomposable combining rules can be expressed using a set of separate, deterministic nodes in the support network such that the family of every non-deterministic node uniquely corresponds to a ground Bayesian clause, as shown in Figure 5.*

Most combining rules commonly employed in Bayesian networks such as *noisy_or* or linear regression are decomposable (cp. [23]). The definition of decomposable combining rules directly imply the following proposition.

Proposition 3. *For each node x in the support network n there exist at most one clause c and a substitution θ such that $\text{body}(c\theta) \subset \text{LH}(B)$ and $\text{head}(c\theta) = x$.*

Thus, while the same clause c can induce more than one node in N , all of these nodes have identical local structure: the associated conditional probability distributions (and so the parameters) have to be identical, i.e.,

$$\forall \text{ substitutions } \theta : \text{cpd}(c\theta) = \text{cpd}(c) .$$

Example 1. Consider the nodes $\text{bt}(\text{ann})$, $\text{mc}(\text{ann})$, $\text{pc}(\text{ann})$ and $\text{bt}(\text{brain})$, $\text{mc}(\text{brain})$, $\text{pc}(\text{brian})$. Both families contribute to the conditional probability distribution associated with the clause defining $\text{bt}(X)$.

This is the same situation as for dynamic Bayesian networks where the parameters that encode the stochastic model of state evolution appear many times in the network. However, gradient methods might be applied to non-decomposable combining functions as well. In the general case, the partial derivatives of an inner function has to be computed. For instance, [3] derive the gradient for *noisy_or* when it is not expressed in the structure. This seems to be more difficult in the case of the EM algorithm, which we will now devise.

5.3 Expectation-Maximization (EM)

The Expectation-Maximization (EM) algorithm [15] is another classical approach to maximum likelihood parameter estimation in the presence of missing values. The basic observation of the Expectation-Maximization algorithm is as that *if the states of all random variables are observed, then learning would be easy*. Assuming that no value is missing, Lauritzen [36] showed that maximum likelihood estimation of Bayesian network parameters simply corresponds to frequency counting in the following way. Let $n(\mathbf{a} \mid \mathbf{D})$ denote the *counts* for a particular joint state \mathbf{a} of variables \mathbf{A} in the data, i.e. the number of cases in which the variables in \mathbf{A} are assigned the evidence \mathbf{a} . Then the maximum likelihood value for the conditional probability value $P(X = x \mid \mathbf{Pa}(X) = \mathbf{u})$ is the ratio

$$\frac{n(X = x, \mathbf{Pa}(X) = \mathbf{u}_k \mid D_l)}{n(\mathbf{Pa}(X) = \mathbf{u}_k \mid D_l)}. \quad (15)$$

However, in the presence of missing values, the maximum likelihood estimates typically cannot be written in closed form. Therefore, the Expectation-Maximization algorithm iteratively performs the following two steps:

- (E-Step)** Based on the current parameters λ and the observed data \mathbf{D} the algorithm computes a distribution over all possible completions of each partially observed data case. Each completion is then treated as a fully-observed data case weighted by its probability.
- (M-Step)** A new set of parameters is then computed based on Equation (15) taking the weights into accounts.

[36] showed that this idea leads to a modified Equation (15) where the *expected counts*

$$\text{en}(\mathbf{a} \mid \mathbf{D}) := \sum_{l=1}^m P_N(\mathbf{a} \mid D_l) \quad (16)$$

are used instead of *counts*. Again, essentially any Bayesian network engine can be used to compute $P(\mathbf{a} \mid D_l)$.

To apply the EM algorithm to parameter estimation of Bayesian logic programs, we assume decomposable combining rules. Thus,

- Each node in the support network was “produced” by exactly one Bayesian clause c , and
- Each node derived from c can be seen as a separate “experiment” for the conditional probability distribution $\text{cpd}(c)$.

Formally, due to the reduction of our problem at hand to parameter estimation within the support network N , the update rule becomes

$$cpd(c_i)_{jk} \leftarrow \frac{\text{en}(c_i|\mathbf{D})}{\text{en}(\text{body}(c_i)|\mathbf{D})} = \frac{\text{en}(\text{head}(c_i), \text{body}(c_i)|\mathbf{D})}{\text{en}(\text{body}(c_i)|\mathbf{D})} \quad (17)$$

where $\text{en}(\cdot|\mathbf{D})$ refers to the first order expected counts as defined in Equation (9). Note that the summation over data cases and ground instances is hidden in $\text{en}(\cdot|\mathbf{D})$. Equation (17) is similar to the one already encountered in Equation (10) for computing the gradient.

5.4 Gradient vs. EM

As one can see, the EM update rule in equation (17) and the corresponding equation (7) for the gradient ascent are very similar. Both rely on computing expected counts. The comparison between EM and (advanced) gradient techniques like conjugate gradient is not yet well understood. Both methods perform a greedy local search, which is guaranteed to converge to stationary points. They both exploit expected counts, i.e., sufficient statistics as their primary computation step. However, there are important differences.

The EM is easier to implement because it does not have to enforce the constraint that the parameters are probability distributions. It converges much faster (at least initially) than simple gradient, and is somewhat less sensitive to starting points. (Conjugate) gradients estimate the step size with a line search involving several additional Bayesian network inferences compared to EM. On the other hand, gradients are more flexible than EM, as they allow one to learn non-multinomial parameterizations using the chain rule for derivatives [3] or to choose other scoring functions than the likelihood [26]. Furthermore, although the EM algorithm is quite successful in practice due to its simplicity and fast initial progress, it has been argued (see e.g. [25,38] and references in there) that the EM convergence can be extremely slow in particular close to the solution, and that more advanced second-order methods should in general be favored to EM or one should switch to gradient-based method after a small number of initial iterations.

Finally, though we focused here on parameter estimation, methods for computing the gradient of the log-likelihood with respect to the parameters of a probabilistic model can also be used to employ generative models within discriminative learners such as SVMs. In the context of probabilistic ILP, this yields relational kernel methods [31,16].

6 Experiments

The presented learning algorithm for Bayesian logic programs is mainly meant as an overall and general framework. Indeed, it leaves several aspects open such as scoring functions. Nevertheless in this section, we report on experiments that show that the algorithm and its underlying principles work.

We implemented the score-based algorithm in Sicstus Prolog 3.9.0 on a Pentium-III 700 MHz Linux machine. The implementation has an interface to the Netica API (<http://www.norsys.com>) for Bayesian network inference and maximum likelihood estimation. To do the maximum likelihood estimation, we adapted the scaled conjugate gradient (SCG) as implemented in Bishop and Nabney's Netlab library (<http://www.ncrg.aston.ac.uk/netlab/>, see also [4]) with an upper bound on the scale parameter of $2 \cdot 10^6$. Parameters were initialized randomly. To avoid zero entries in the conditional probability tables, m-estimates were used.

6.1 Genetic Domain

The goal was to learn a global, descriptive model for our genetic domain, i.e. to learn the program *bloodtype*. We considered two totally independent families using the predicates given by *bloodtype* having 12 respectively 15 family members. For each least Herbrand model 1000 data cases from the induced Bayesian network were sampled with a fraction of 0.4 of missing at random values of the observed nodes making in total 2000 data cases.

Therefore, we first had a look at the (logical) hypotheses space. The space could be seen as the first order equivalent of the space for learning the structure of Bayesian networks (see Figure 3). The generating hypothesis is a member of it. In a further experiment, we fixed the definitions for $\mathbf{m}/2$ and $\mathbf{f}/2$. The hypothesis scored best included $\mathbf{bt}(X) \mid \mathbf{mc}(X), \mathbf{pc}(X)$, i.e. the algorithm re-discovered the intensional definition which was originally used to build the data cases. However, the definitions of $\mathbf{mc}/1$ and $\mathbf{pc}/1$ considered genetic information of the grandparents to be important. It failed to re-discover the original definitions for reasons explained above. The predicates $\mathbf{m}/2$ and $\mathbf{f}/2$ were not part of the learned model rendering them to be extensionally defined. Nevertheless, the founded global model had a slightly better likelihood than the original one.

6.2 Bongard Domain

The Bongard problems (due to the Russian scientist M. Bongard) are well-known problems within inductive logic programming. Consider Figure 6. Each example or scene consists of

- A variable number of geometrical objects such as triangles, rectangles and circles etc (predicate $\mathbf{obj}/2$ with $\mathbf{S}(obj) = \{triangle, circle\}$), each having a number of different properties such as color, orientation, size etc., and
- A variable number of relations between objects such as in (predicate $\mathbf{in}/3$ having states *true, false*), leftof, above etc.

The task is to find a set of rules which *discriminates* positive from negative examples (represented by $\mathbf{class}/1$ over the states *pos, neg*) by looking at the kind of objects they consists of. Though the Bongard problems are toy problems, they are very similar to real-world problems in e.g. the field of molecular biology where essentially the same representational problems arise. Data consists of

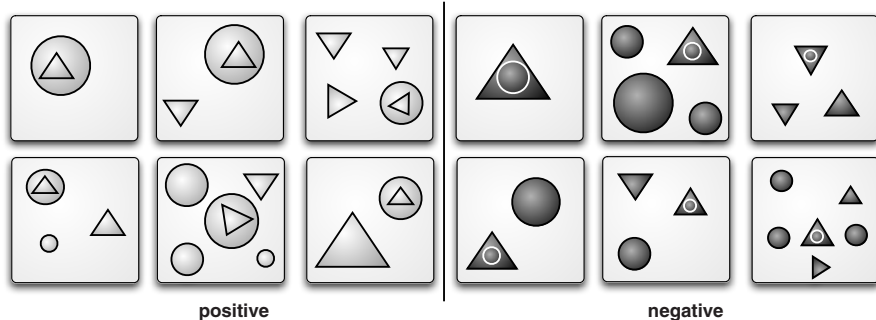


Fig. 6. A Bongard problem consisting of 12 scenes, six positive ones and six negative ones. The goal is to discriminate between the two classes.

molecules, each of which is composed of several atoms with specific properties such as charge. There exists a number of relations between atoms like e.g. bonds, structure etc.

In most real-world applications, the data is *noisy*. Class labels or object properties might be wrong or missing in the data cases. One extreme case concerns clustering where no class labels are given. Furthermore, we might be uncertain about relations among objects. Some positive examples might state that a triangle is not in a circle due to noise. In such cases, Bayesian logic programs are a natural method of choice. We conducted the following experiments.

First, we generated 20 positive and 20 negative examples of the concept “*there is a triangle in a circle.*” The number of objects varied from 2 to 8. We conducted three different experiments. We assumed the *in* relation to be deterministic and given as background knowledge $\text{in}(\text{Example}, \text{Obj1}, \text{Obj2})$, i.e. we assumed that there is no uncertainty about it. Due to that, no conditional probability distribution has to take $\text{in}/3$ into account. Because each scene is independent of the other, we represented the whole training data as one data case

$$\{ \text{class}(e1) = \text{pos}, \text{obj}(e1, o1) = \text{triangle}, \text{obj}(e1, o2) = \text{circle}, \\ \text{class}(e2) = \text{neg}, \text{obj}(e2, o1) = \text{triangle}, \text{size}(e2, o1) = \text{large}, \\ \text{obj}(e2, o2) = \text{'?'}, \dots \}$$

with the background knowledge $\text{in}(e1, o1, o2), \dots$ where $e1, e2, \dots$ are identifiers for examples and $o1, o2, \dots$ for objects. A fraction of 0.2 of the random variables were not observed. Our algorithm scored the hypothesis

$$\text{class}(\text{Ex}) \mid \text{obj}(\text{Ex}, \text{O1}), \text{in}(\text{Ex}, \text{O1}, \text{O2}), \text{obj}(\text{Ex}, \text{O2}).$$

best after specifying $\text{obj}(\text{Ex}, \text{O2})$ as a lookahead for $\text{in}(\text{Ex}, \text{O1}, \text{O2})$. The conditional probability distribution assigned *pos* a probability higher than 0.6 only if object O1 was a triangle and O2 a circle. Without the lookahead, adding $\text{in}(\text{Ex}, \text{O1}, \text{O2})$ yield no improvement in score, and the correct hypothesis was not considered. The hypothesis is not a well-defined Bayesian networks, but it

says that ground atoms over `obj/2` extensional defined. Therefore, we estimated the maximum likelihood parameters of

```
obj(Ex,0) | dom(Ex,0).
class(Ex) | obj(Ex,01), in(Ex,01,02), obj(Ex,02).
```

where `dom/2` ensured range-restriction and was part of the deterministic background knowledge. Using 0.6 as threshold, the learned Bayesian logic program had accuracy 1.0 on the training set and on an independently generated validation set consisting of 10 positive and 10 negative examples.

In a second experiments, we fixed the structure of the program learned in the first experiment, and estimated its parameters on a data set consisting of 20 positive and 20 negative examples of the disjunctive concept “*there is a (triangle or a circle) in a circle.*” The estimated conditional probability distribution gave almost equal probability for the object 01 to be a triangle or circle.

In third experiment, we assumed uncertainty about the *in* relation. We enriched the data case used for the first experiment in the following way

$$\{ \text{class}(e1) = \text{pos}, \text{obj}(e1, o1) = \text{triangle}, \text{obj}(e1, o2) = \text{circle}, \\ \text{in}(e1, o1, o2) = \text{true}, \text{class}(e2) = \text{neg}, \text{obj}(e2, o1) = \text{triangle}, \\ \text{size}(e2, o1) = \text{large}, \text{obj}(e2, o2) = \text{'?'}, \text{in}(e2, o1, o2) = \text{false}, \dots \},$$

i.e. for each pair of objects that could be related by *in* a ground atom over *in* was included. Note that the state need not to be observed. Here, the algorithm did not re-discovered the correct rule but

```
class(X) | obj(X,Y)
obj(X,Y) | in(X,Y,Z), obj(Z).
```

This is interesting, because when these rules are used to classify examples, only the first rule is needed. The class is independent of any information about `in/3` given full knowledge about the objects. The likelihood of the founded solution was close to the one of `class(Ex) | obj(Ex,01), in(Ex,01,02), obj(Ex,02)` on the data (absolute difference less than 0.001). However, the accuracy decreased (about 0.6 on an independently generated training set (10 pos / 10 neg)) for the reasons explained above: We learned a global model not focusing on the classification error. [18] showed for Bayesian network classifier that maximizing the conditional likelihood of the class variable comes close to minimizing the classification error. In all experiments we assumed *noisy* or as combining rule.

Finally, we conducted a simple clustering experiments. We generated 20 positive and 20 negative examples of the disjunctive concept “*there is a triangle.*” There were triangles, circles and squares. The number of objects varied from 2 to 8. All class labels were said to be observed, and 20% of the remaining stated were missing at random. The learned hypothesis was `class(X) | obj(X,Y)` totally separating the two classes.

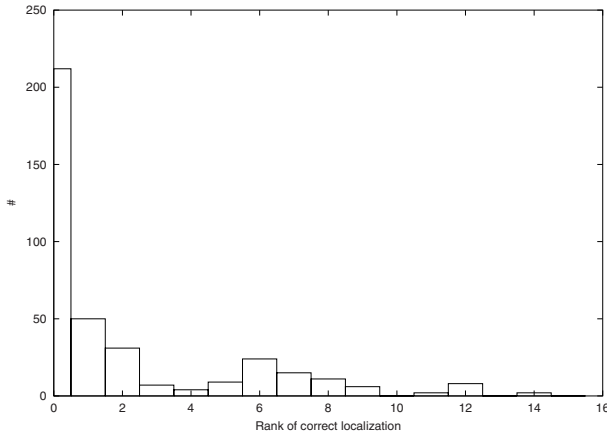


Fig. 7. An histogram of the ranks of the correct localization computed by an Bayesian logic program on the KDD Cup 2001 test set

6.3 KDD Cup 2001

We also started to conduct experiments on *large scale* data sets namely the KDD Cup 2001⁴ data sets, cf. [8]. Task 3 is to predict the localizations $\text{local}(\mathbf{G})$ of proteins encoded by the genes \mathbf{G} . This is a multi-class problem because there are 16 possible localizations. The training set consisted of 862 genes, the test set of 381 genes. The information we used included whether organisms with an mutation in this gene can survive $\text{ess}(\mathbf{G})$ (3 states), the class $\text{class}(\mathbf{G})$ of a gene/protein \mathbf{G} (24 states), the complex $\text{compl}(\mathbf{G})$ (56 states), and the other proteins \mathbf{G} with which each protein \mathbf{G} is known to interact encoded by $\text{inter}(\mathbf{G1}, \mathbf{G2})$. To avoid a large interaction table, we considered only those interactions with a correlation higher than 0.85. Furthermore, we introduced a hidden predicate $\text{hidden}/1$ with domain 0, 1, 2 to compress the representation size because e.g. the conditional probability table of $\text{local}(\mathbf{G1}) \mid \text{inter}(\mathbf{G1}, \mathbf{G2}), \text{local}(\mathbf{G2})$ would consists of 225 entries (instead of 45 using $\text{hidden}(\mathbf{G1})$). The ground atoms over $\text{hidden}/1$ were never observed in the data. Nevertheless, the naive Prolog representation of the support networks induced by some hypothesis (more than 4.400 random variables with more than 60.000 parameters) in our current implementation broke the memory capabilities of Sicstus Prolog. Due to that, we can only report on preliminary results. We only considered maximum likelihood parameter estimation on the training set. The (logical) structure is based on naive Bayes taking relational information into account:

```

local(G1) | gene(G1) .
hidden(G1) | local(G1) .
hidden(G1) | inter(G1,G2), local(G2) .

```

⁴ For details see <http://www.cs.wisc.edu/~dpage/kddcup2001/>

```

class(G1) | hidden(G1).
compl(G1) | local(G1).
ess(G1)   | local(G1).

```

As combining rule, we used for all predicates *average*. The given ground atoms over `inter/2` were used as pure logical background knowledge. Therefore, the conditional probability distribution associated to `hidden(G1) | inter(G1,G2)`, `local(G2)` had not to take it into account. The parameters were randomly initialized. Again, the training set was represented as one data case, so that no artificial independencies among genes were postulated. Estimate the parameters took 12 iteration (about 30 min). The learned Bayesian logic program achieved an accuracy of 0.57 (top 50% level of submitted models was 0.61, best predictive accuracy was 0.72). A learner predicting always the majority class would achieve an predictive accuracy of 0.44. Furthermore, when we rank for each test gene its possible localizations according to the probability computed by the program, then the correct localization was among the three highest ranked localizations in 293 out of 381 cases (77%) (cf. Figure 7). Not that it took 40 iterations to learn the corresponding grounded Bayesian logic program.

7 Related Work

The learning of Bayesian networks has been thoroughly investigated in the Uncertainty in AI community, see e.g. [22,7]. Binder *et al.* [3], whose approach we have adapted, present results for a gradient-based method. But so far – to the best of our knowledge – there has not been much work on learning within first order extensions of Bayesian networks. Koller and Pfeffer [33] adapt the EM algorithm for probabilistic logic programs [40], a framework which in contrast to Bayesian logic programs sees ground atoms as states of random variables. Although the framework seems to theoretically allow for continuous random variables there exists no (practical) query-answering procedure for this case; to the best of our knowledge, Ngo and Haddawy [40] give only a procedure for variables having finite domains. Furthermore, Koller and Pfeffer’s approach utilizes support networks, too, but requires the intersection of the support networks of the data cases to be empty. This could be in our opinion in some cases too restrictive, e.g. in the case of dynamic Bayesian networks. Friedman *et al.* [19,20] adapted the Structural-EM to learn the structure of probabilistic relational models. It applies the idea of the standard EM algorithm for maximum likelihood parameter estimation to the problem of learning the structure. If we know the values for all random variables, then the maximum likelihood estimate can be written in closed form. Based on the current hypothesis a distribution over all possible completions of each partially observed data case is computed. Then, new hypotheses are computed using a score-based method. However, the algorithm does not consider logical constraints on the space of hypotheses. Indeed, the considered clauses need not be logically valid on the data. Therefore, combining our approach with the structural EM seems to be reasonable and straightforward. Finally, there is work on learning object-oriented Bayesian networks [35,1].

There exist also methods for learning within first order probabilistic frameworks which do not build on Bayesian networks. Sato and Kameya [43] introduce an EM method for parameter estimation of PRISM programs, see also Chapter 5. Cussens [11] investigates EM like methods for estimating the parameters of stochastic logic programs (SLPs). As a reminder, SLPs lift probabilistic context-free grammars to the first order case by replacing production rules with probability values with clauses labeled with probability values. In turn, they define probability distributions over proofs. As discussed in Chapter 1, this is quite different from Bayesian logic programs, which lift Bayesian networks by defining probability distributions over an interpretation. Nevertheless, mappings between the two approaches exist as shown by Muggleton and Chen in Chapter 12. Cussens’ EM approach for SLPs has been successfully applied to protein fold discovery by Chen *et. al* as reported in Chapter 9. Muggleton [39] uses ILP techniques to learn the logical structure/program of stochastic logic programs. The used ILP setting is different to *learning from interpretations*, it is not based on learning Bayesian networks, and so far considers only for single predicates definitions.

To summarize, the related work on learning probabilistic relational models mainly differs in three points from ours:

- The underlying (logical) frameworks lack important knowledge representational features which Bayesian logic programs have.
- They adapt the EM algorithm to do parameter estimation which is particularly easy to implement. However, there are problematic issues both regarding speed of convergence as well as convergence towards a local (sub-optimal) maximum of the likelihood function. Different accelerations based on the gradient are discussed in [38]. Also, the EM algorithm is difficult to apply in the case of general probability density functions because it relies on computing the sufficient statistics (cf. [22]).
- No probabilistic extension of the *learning from interpretations* is established.

8 Conclusions

A new link between ILP and learning of Bayesian networks was established. We have proposed a scheme for learning both the probabilities and the structure of Bayesian logic programs. We addressed the question “where do the numbers come from?” by showing how to compute the gradient of the likelihood based on ideas known for (dynamic) Bayesian networks. The intensional representation of Bayesian logic programs, i.e. their compact representation should speed up learning and provide good generalization. The general learning setting built on the ILP setting *learning from interpretations*. We have argued that by adapting this setting score-based methods for structural learning of Bayesian networks could be updated to the first order case. The ILP setting is used to define and traverse the space of (logical) hypotheses.

The experiments proved the principle of the algorithm. Their results highlight that future work on improved scoring functions is needed. We plan to conduct experiments on real-world scale problems. The use of refinement operators adding

or deleting non constant-free atoms should be explored. Furthermore, it would be interesting to weaken the assumption that a data case corresponds to a complete interpretation. Not assuming all relevant random variables are known would be interesting for learning intensional rules like $\text{nat}(s(X)) \mid \text{nat}(X)$. Ideas for handling this within inductive logic programming might be adapted [14,6]. Furthermore, instead of traditional score-based greedy algorithm more advanced UAI methods such as Friedman's Structural-EM or structure search among equivalence classes of Bayesian logic programs may be adapted taking advantage of the logical constraints implied by the data cases. In any case, we believe that the proposed approach is a good point of departure for further research. The link established between ILP and Bayesian networks seems to be bi-directional. Can ideas developed in the UAI community be carried over to ILP?

Acknowledgements

The authors would like to thank Manfred Jaeger, Stefan Kramer and David Page for helpful discussions on the ideas of the paper. Furthermore, the authors would like to thank Jan Ramon and Hendrik Blockeel for making available their Bongard problems generators. This research was partly supported by the European Union IST programme under contract number IST-2001-33053 (Application of Probabilistic Inductive Logic Programming – APRIL).

References

1. Bangsø, O., Langseth, H., Nielsen, T.D.: Structural learning in object oriented domains. In: Russell, I., Kolen, J. (eds.) Proceedings of the Fourteenth International Florida Artificial Intelligence Research Society Conference (FLAIRS 2001), Key West, Florida, USA, pp. 340–344. AAAI Press, Menlo Park (2001)
2. Bauer, H.: Wahrscheinlichkeitstheorie, 4th edn., Walter de Gruyter, Berlin, New York (1991)
3. Binder, J., Koller, D., Russell, S., Kanazawa, K.: Adaptive Probabilistic Networks with Hidden Variables. *Machine Learning* 29(2–3), 213–244 (1997)
4. Bishop, C.M.: *Neural Networks for Pattern Recognition*. Oxford University Press, Oxford (1995)
5. Blockeel, H., De Raedt, L.: Lookahead and discretization in ilp. In: Džeroski, S., Lavrač, N. (eds.) *ILP 1997*. LNCS, vol. 1297, pp. 77–85. Springer, Heidelberg (1997)
6. Blockeel, H., De Raedt, L.: ISIDD: An Interactive System for Inductive Database Design. *Applied Artificial Intelligence* 12(5), 385 (1998)
7. Buntine, W.: A guide to the literature on learning probabilistic networks from data. *IEEE Transaction on Knowledge and Data Engineering* 8, 195–210 (1996)
8. Cheng, J., Hatzis, C., Krogel, M.-A., Morishita, S., Page, D., Sese, J.: KDD Cup 2002 Report. *SIGKDD Explorations* 3(2), 47–64 (2002)
9. Cooper, G.F.: The computational complexity of probabilistic inference using Bayesian belief networks. *Artificial Intelligence* 42, 393–405 (1990)
10. Cowell, R.G., Dawid, A.P., Lauritzen, S.L., Spiegelhalter, D.J.: Probabilistic networks and expert systems. In: *Statistics for engineering and information*, Springer, Heidelberg (1999)

11. Cussens, J.: Parameter estimation in stochastic logic programs. *Machine Learning* 44(3), 245–271 (2001)
12. De Raedt, L.: Logical settings for concept-learning. *Artificial Intelligence* 95(1), 197–201 (1997)
13. De Raedt, L., Bruynooghe, M.: A theory of clausal discovery. In: Bajcsy, R. (ed.) *Proceedings of the Thirteenth International Joint Conference on Artificial Intelligence (IJCAI 1993)*, Chambéry, France, pp. 1058–1063. Morgan Kaufmann, San Francisco (1993)
14. De Raedt, L., Dehaspe, L.: Clausal discovery. *Machine Learning* 26(2-3), 99–146 (1997)
15. Dempster, A.P., Laird, N.M., Rubin, D.B.: Maximum likelihood from incomplete data via the EM algorithm. *J. Royal Stat. Soc. B* 39, 1–39 (1977)
16. Dick, U., Kersting, K.: Fisher Kernels for relational data. In: Fürnkranz, J., Scheffer, T., Spiliopoulou, M. (eds.) *Proceedings of the 17th European Conference on Machine Learning (ECML 2006)*, Berlin, Germany, pp. 112–125 (2006)
17. Flach, P.A., Lachiche, N.: 1BC: A first-order Bayesian classifier. In: Džeroski, S., Flach, P.A. (eds.) *ILP 1999. LNCS (LNAI)*, vol. 1634, pp. 92–103. Springer, Heidelberg (1999)
18. Friedman, N., Geiger, D., Goldszmidt, M.: Bayesian network classifiers. *Machine Learning* 29, 131–163 (1997)
19. Friedman, N., Getoor, L., Koller, D., Pfeffer, A.: Learning probabilistic relational models. In: Dean, T. (ed.) *Proceedings of the Sixteenth International Joint Conferences on Artificial Intelligence (IJCAI 1999)*, Stockholm, Sweden, pp. 1300–1309. Morgan Kaufmann, San Francisco (1999)
20. Getoor, L., Koller, D., Taskar, B., Friedman, N.: Learning probabilistic relational models with structural uncertainty. In: Getoor, L., Jensen, D. (eds.) *Proceedings of the AAAI-2000 Workshop on Learning Statistical Models from Relational Data*, AAAI Press, Menlo Park (2000)
21. Gilks, W.R., Thomas, A., Spiegelhalter, D.J.: A language and program for complex bayesian modelling. *The Statistician* 43 (1994)
22. Heckerman, D.: A Tutorial on Learning with Bayesian Networks. Technical Report MSR-TR-95-06, Microsoft Research (1995)
23. Heckerman, D., Breese, J.: Causal Independence for Probability Assessment and Inference Using Bayesian Networks. Technical Report MSR-TR-94-08, Microsoft Research (1994)
24. Jaeger, M.: Relational Bayesian networks. In: Geiger, D., Shenoy, P.P. (eds.) *Proceedings of the Thirteenth Annual Conference on Uncertainty in Artificial Intelligence (UAI 1997)*, Providence, Rhode Island, USA, pp. 266–273. Morgan Kaufmann, San Francisco (1997)
25. Jamshidian, M., Jennrich, R.I.: Acceleration of the EM Algorithm by using Quasi-Newton Methods. *Journal of the Royal Statistical Society B* 59(3), 569–587 (1997)
26. Jensen, F.V.: Gradient descent training of bayesian networks. In: Hunter, A., Parsons, S. (eds.) *ECSQARU 1999. LNCS (LNAI)*, vol. 1638, pp. 190–200. Springer, Heidelberg (1999)
27. Jensen, F.V.: *Bayesian networks and decision graphs*. Springer, Heidelberg (2001)
28. Kersting, K., De Raedt, L.: Bayesian logic programs. In: Cussens, J., Frisch, A. (eds.) *Work-in-Progress Reports of the Tenth International Conference on Inductive Logic Programming (ILP 2000)* (2000), <http://SunSITE.Informatik.RWTH-Aachen.DE/Publications/CEUR-WS/Vol-35/>
29. Kersting, K., De Raedt, L.: Bayesian logic programs. Technical Report 151, University of Freiburg, Institute for Computer Science (submitted) (April 2001)

30. Kersting, K., De Raedt, L., Kramer, S.: Interpreting Bayesian Logic Programs. In: Getoor, L., Jensen, D. (eds.) Working Notes of the AAAI-2000 Workshop on Learning Statistical Models from Relational Data (SRL), Austin, Texas, AAAI Press, Menlo Park (2000)
31. Kersting, K., Gärtner, T.: Fisher Kernels for Logical Sequences. In: Boulicaut, J.-F., Esposito, F., Giannotti, F., Pedreschi, D. (eds.) ECML 2004. LNCS (LNAI), vol. 3201, p. 205. Springer, Heidelberg (2004)
32. Koller, D.: Probabilistic relational models. In: Džeroski, S., Flach, P.A. (eds.) ILP 1999. LNCS (LNAI), vol. 1634, pp. 3–13. Springer, Heidelberg (1999)
33. Koller, D., Pfeffer, A.: Learning probabilities for noisy first-order rules. In: Proceedings of the Fifteenth Joint Conference on Artificial Intelligence (IJCAI 1997), Nagoya, Japan, pp. 1316–1321 (1997)
34. Lam, W., Bacchus, F.: Learning Bayesian belief networks: An approach based on the MDL principle. *Computational Intelligence* 10(4) (1994)
35. Langseth, H., Bangsø, O.: Parameter learning in object oriented Bayesian networks. *Annals of Mathematics and Artificial Intelligence* 32(1-2), 221–243 (2001)
36. Lauritzen, S.L.: The EM algorithm for graphical association models with missing data. *Computational Statistics and Data Analysis* 19, 191–201 (1995)
37. Lloyd, J.W.: *Foundations of Logic Programming*, 2nd edn. Springer, Berlin (1989)
38. McKachlan, G.J., Krishnan, T.: *The EM Algorithm and Extensions*. John Eiley & Sons, Inc. (1997)
39. Muggleton, S.H.: Learning stochastic logic programs. In: Getoor, L., Jensen, D. (eds.) Working Notes of the AAAI-2000 Workshop on Learning Statistical Models from Relational Data (SRL), Austin, Texas, AAAI Press, Menlo Park (2000)
40. Ngo, L., Haddawy, P.: Answering queries from context-sensitive probabilistic knowledge bases. *Theoretical Computer Science* 171, 147–177 (1997)
41. Pearl, J.: *Reasoning in Intelligent Systems: Networks of Plausible Inference*, 2nd edn. Morgan Kaufmann, San Francisco (1991)
42. Poole, D.: Probabilistic Horn abduction and Bayesian networks. *Artificial Intelligence* 64, 81–129 (1993)
43. Sato, T., Kameya, Y.: Parameter learning of logic programs for symbolic-statistical modeling. *Journal of Artificial Intelligence Research* 15, 391–454 (2001)
44. Srinivasan, A., Muggleton, S., Bain, M.: The justification of logical theories based on data compression. In: Furukawa, K., Michie, D., Muggleton, S. (eds.) *Machine Intelligence*, vol. 13, Oxford University Press, Oxford (1994)
45. Sterling, L., Shapiro, E.: *The Art of Prolog: Advanced Programming Techniques*. MIT Press, Cambridge (1986)
46. Taskar, B., Segal, E., Koller, D.: Probabilistic clustering in relational data. In: Nebel, B. (ed.) *Seventeenth International Joint Conference on Artificial Intelligence (IJCAI 2001)*, Seattle, Washington, USA, pp. 870–887. Morgan Kaufmann, San Francisco (2001)
47. Xiang, Y., Wong, S.K.M., Cercone, N.: Critical remarks on single link search in learning belief networks. In: Horvitz, E., Jensen, F.V. (eds.) *Proceedings of the Twelfth Annual Conference on Uncertainty in Artificial Intelligence (UAI 1996)*, Portland, Oregon, USA, pp. 564–571. Morgan Kaufmann, San Francisco (1996)