

# Online Research @ Cardiff

This is an Open Access document downloaded from ORCA, Cardiff University's institutional repository: <http://orca.cf.ac.uk/106737/>

This is the author's version of a work that was submitted to / accepted for publication.

Citation for final published version:

De Raedt, Luc and Kimmig, Angelika 2015. Probabilistic (logic) programming concepts. Machine Learning 100 (1) , pp. 5-47. 10.1007/s10994-015-5494-z file

Publishers page: <http://dx.doi.org/10.1007/s10994-015-5494-z> <<http://dx.doi.org/10.1007/s10994-015-5494-z>>

Please note:

Changes made as a result of publishing processes such as copy-editing, formatting and page numbers may not be reflected in this version. For the definitive version of this publication, please refer to the published source. You are advised to consult the publisher's version if you wish to cite this paper.

This version is being made available in accordance with publisher policies. See <http://orca.cf.ac.uk/policies.html> for usage policies. Copyright and moral rights for publications made available in ORCA are retained by the copyright holders.



# Probabilistic (Logic) Programming Concepts

Luc De Raedt · Angelika Kimmig

Received: date / Accepted: date

**Abstract** A multitude of different probabilistic programming languages exists today, all extending a traditional programming language with primitives to support modeling of complex, structured probability distributions. Each of these languages employs its own probabilistic primitives, and comes with a particular syntax, semantics and inference procedure. This makes it hard to understand the underlying programming concepts and appreciate the differences between the different languages.

To obtain a better understanding of probabilistic programming, we identify a number of core programming concepts underlying the primitives used by various probabilistic languages, discuss the execution mechanisms that they require and use these to position and survey state-of-the-art probabilistic languages and their implementation.

While doing so, we focus on probabilistic extensions of *logic* programming languages such as Prolog, which have been considered for over 20 years.

## 1 Introduction

The substantial interest in statistical relational learning (Getoor and Taskar, 2007), probabilistic (inductive) logic programming (De Raedt et al, 2008) and probabilistic programming languages (Roy et al, 2008) has resulted in a wide variety of different formalisms, models and languages, with applications in structured, uncertain domains such as natural language processing, bioinformatics, and activity recognition. The multitude of probabilistic languages that exists today provides evidence for the richness and maturity of the field, but on the other hand, makes it hard to get an appreciation and understanding of the relationships and differences between the different languages. Furthermore, most arguments in the literature about the relationship amongst these languages are about the expressiveness of

---

L. De Raedt · A. Kimmig  
Department of Computer Science, KU Leuven, Celestijnenlaan 200A - bus 2402, 3001 Heverlee, Belgium  
E-mail: {firstname.lastname}@cs.kuleuven.be

these languages, that is, they state (often in an informal way) that one language is more expressive than another one (implying that the former could be used to emulate the latter). By now, it is commonly accepted that the more interesting question is concerned with the underlying concepts that these languages employ and their effect on the inference mechanisms, as their expressive power is often very similar. However, a multitude of different probabilistic primitives exists, which makes it hard to appreciate their relationships.<sup>1</sup>

To alleviate these difficulties and obtain a better understanding of the field we identify a number of core probabilistic programming concepts and relate them to one another. We cover the basic concepts representing different types of random variables, but also general modeling concepts such as negation or time and dynamics, and programming constructs such as meta-calls and ways to handle sets. While doing so, we focus on probabilistic extensions of logic programming languages because this is (arguably) the first and best studied probabilistic programming paradigm. It has been studied for over 20 years starting with the seminal work of David Poole (1992) and Taisuke Sato (1995), and now includes languages such as the independent choice logic (ICL) (Poole, 1995, 2008), stochastic logic programs (SLPs) (Muggleton, 1996), PRISM (Sato and Kameya, 1997, 2001), Bayesian logic programs (BLPs) (Kersting and De Raedt, 2001, 2008), CLP( $\mathcal{BN}$ ) (Santos Costa et al, 2003, 2008), logic programs with annotated disjunctions (LPADs) (Vennekens et al, 2004), P-log (Baral et al, 2004, 2009), Dyna (Eisner et al, 2005), CP-logic (Vennekens et al, 2006, 2009), ProbLog (De Raedt et al, 2007), and programming with personalized Pagerank (PROPPR) (Wang et al, 2013). Another reason for focussing on probabilistic extensions of logic programming languages is that the concepts are all embedded within the same host language, so we can focus on semantics rather than syntax. At the same time, we also relate the concepts to alternative probabilistic programming languages such as IBAL (Pfeffer, 2001), Bayesian logic (BLOG) (Milch et al, 2005), Church (Goodman et al, 2008) and Figaro (Pfeffer, 2009) and to some extent also to statistical relational learning models such as relational Bayesian networks (RBNs) (Jaeger, 1997, 2008), probabilistic relational models (PRMs) (Koller and Pfeffer, 1998; Getoor et al, 2007) and Markov logic (Richardson and Domingos, 2006). Most statistical relational learning approaches employ a knowledge-based model construction approach (Wellman et al, 1992), in which the logic is used as a template for constructing a graphical model. Typical probabilistic programming languages, on the other hand, employ a variant of Sato’s distribution semantics (Sato, 1995), in which random variables directly correspond to ground facts and a traditional program specifies how to deduce further knowledge from these facts. This difference explains why we introduce the concepts in the context of the distribution semantics, and discuss approaches to knowledge-based model construction separately.

Inference, that is, evaluating the probability distribution defined by a program or model, is a key challenge in probabilistic programming and statistical relational learning. Furthermore, the choice of inference approach often influences which probabilistic primitives can be supported. Enormous progress has been made in the past few years w.r.t. probabilistic inference and numerous inference procedures

<sup>1</sup> Throughout the paper we use the term *primitive* to denote a particular syntactic and semantic construct that is available in a particular probabilistic programming language, and the term *concept* to denote the underlying notion. Different primitives may hence realize the same concept.

have been contributed. Therefore, we also identify some core classes of inference mechanisms for probabilistic programming and discuss which ones to use for which probabilistic concept. Inference in probabilistic languages also is an important building block of approaches that learn the structure and/or parameters of such models from data. Given the variety of approaches that exist today, a discussion of learning is beyond the scope of this paper.

To summarize, the key contributions of this paper are (1) the identification of a number of core concepts that are used by various probabilistic languages, (2) a discussion of the execution mechanisms that they require, and (3) a positioning of state-of-the-art probabilistic languages and implementations w.r.t. these concepts. Although many of the concepts we discuss are well-described in the literature, some even in survey papers (De Raedt and Kersting, 2003; Poole, 2008), we believe a new and up-to-date survey is warranted due to the rapid developments of the field which rapidly renders existing surveys incomplete and even outdated. To the best of our knowledge, this is also the first time that such a wide variety of probabilistic programming concepts and languages, also in connection to inference, is discussed in a single paper.

We expect the reader to be familiar with basic language concepts and terms of Prolog (Lloyd, 1989; Flach, 1994); a quick summary can be found in Appendix A.

This paper is organized as follows. We first discuss the distribution semantics (Section 2) and classify corresponding inference approaches according to their logical and probabilistic components (Section 3). Section 4 identifies the probabilistic programming concepts. In Section 5, we discuss the relation with statistical relational modeling approaches rooted in graphical models. Section 6 relates the different inference approaches to the probabilistic programming concepts. We touch upon applications of probabilistic logic programming in Section 7. Section 8 concludes the survey and summarizes a number of key open questions.

## 2 Distribution Semantics

Sato's distribution semantics (Sato, 1995) is a well-known semantics for probabilistic logics that has been considered many times and under varying names, cf. (Dantsin, 1991; Poole, 1993; Fuhr, 2000; Poole, 2000; Sato and Kameya, 2001; Dalvi and Suciu, 2004; De Raedt et al, 2007). It generalizes the least model semantics of logic programming to a distribution over (the least models of) a set of logic programs that share their definite clauses, but differ in the set of facts. This is achieved based on a joint distribution over truth value assignments to these facts, cf. Equation (4) below. Prominent examples of Prolog-based languages using this semantics include ICL (Poole, 2008), PRISM (Sato and Kameya, 2001) and ProbLog (De Raedt et al, 2007; Kimmig et al, 2011b), even though there exist subtle differences between these languages as we will illustrate later. Sato has defined the distribution semantics for a countably infinite set of random variables and a general class of distributions. We focus on the finite case here, discussing the two most popular instances of the semantics, based on a set of independent random variables and independent probabilistic choices, respectively, and refer to Sato (1995) for details on the general case.

## 2.1 Probabilistic Facts

The arguably most basic instance of the distribution semantics uses a finite set of Boolean random variables that are all assumed to be independent. We use the following running example inspired by the well-known alarm Bayesian network (Pearl, 1988):

$$\begin{aligned}
0.1 &:: \text{burglary}. & 0.7 &:: \text{hears\_alarm}(\text{mary}). \\
0.2 &:: \text{earthquake}. & 0.4 &:: \text{hears\_alarm}(\text{john}). \\
& \text{alarm} :- \text{earthquake}. & & \\
& \text{alarm} :- \text{burglary}. & & \\
& \text{calls}(\text{X}) :- \text{alarm}, \text{hears\_alarm}(\text{X}). & & \\
& \text{call} :- \text{calls}(\text{X}). & & 
\end{aligned} \tag{1}$$

The program consists of a set  $R$  of definite clauses or *rules*,<sup>2</sup> and a set  $F$  of ground facts  $f$ , each of them labeled with a probability  $p$ , written as  $p :: f$ . We call such labeled facts *probabilistic facts*. Each probabilistic fact corresponds to a Boolean random variable that is *true* with probability  $p$  and *false* with probability  $1 - p$ . We use  $b$ ,  $e$ ,  $hm$  and  $hj$  to denote the random variables corresponding to `burglary`, `earthquake`, `hears_alarm(mary)` and `hears_alarm(john)`, respectively. Assuming that all these random variables are independent, we obtain the following probability distribution  $P_F$  over truth value assignments to these random variables and their corresponding sets of ground facts  $F' \subseteq F$ :

$$P_F(F') = \prod_{f_i \in F'} p_i \cdot \prod_{f_i \in F \setminus F'} (1 - p_i) \tag{2}$$

For instance, the truth value assignment `burglary = true`, `earthquake = false`, `hears_alarm(mary) = true`, `hears_alarm(john) = false`, which we will abbreviate as  $b \wedge \neg e \wedge hm \wedge \neg hj$ , corresponds to the set of facts  $\{\text{burglary}, \text{hears\_alarm}(\text{mary})\}$ , and has probability  $0.1 \cdot (1 - 0.2) \cdot 0.7 \cdot (1 - 0.6) = 0.0336$ . The corresponding logic program obtained by adding the set of rules  $R$  to the set of facts, also called a *possible world*, is

$$\begin{aligned}
& \text{burglary}. \\
& \text{hears\_alarm}(\text{mary}). \\
& \text{alarm} :- \text{earthquake}. \\
& \text{alarm} :- \text{burglary}. \\
& \text{calls}(\text{X}) :- \text{alarm}, \text{hears\_alarm}(\text{X}). \\
& \text{call} :- \text{calls}(\text{X}).
\end{aligned} \tag{3}$$

As each logic program obtained by choosing a truth value for every probabilistic fact has a unique least Herbrand model (i.e., a unique least model using only symbols from the program; cf. Appendix A),  $P_F$  can be used to define the *success*

<sup>2</sup> Rules are deterministic and no further constraints on  $R$  are imposed in general; an exception is discussed on page 12. Probabilistic rules can be modeled in this setting by adding a fresh probabilistic fact to the body; an alternative is presented in Section 2.2.

probability  $P(q)$  of a query  $q$ , that is, the probability that  $q$  is true in a randomly chosen such program, as the sum over all programs that entail  $q$ :

$$P(q) := \sum_{\substack{F' \subseteq F \\ \exists \theta F' \cup R \models q\theta}} P_F(F') \quad (4)$$

$$= \sum_{\substack{F' \subseteq F \\ \exists \theta F' \cup R \models q\theta}} \prod_{f_i \in F'} p_i \cdot \prod_{f_i \in F \setminus F'} (1 - p_i). \quad (5)$$

Naively, the success probability can thus be computed by enumerating all sets  $F' \subseteq F$ , for each of them checking whether the corresponding possible world entails the query, and summing the probabilities of those that do. As fixing the set of facts yields an ordinary logic program, the entailment check can use any reasoning technique for such programs.

For instance, *forward reasoning*, also known as applying the  $T_P$  operator, starts from the set of facts and repeatedly uses rules to derive additional facts until no more facts can be derived. In our example possible world (3), we thus start from  $\{\text{burglary}, \text{hears\_alarm}(\text{mary})\}$ , and first add **alarm** due to the second rule based on **burglary**. This in turn makes it possible to add **calls(mary)** using the third rule and substitution  $X = \text{mary}$ , and finally, **call** is added using the last rule, resulting in the least Herbrand model  $\{\text{burglary}, \text{hears\_alarm}(\text{mary}), \text{alarm}, \text{calls}(\text{mary}), \text{call}\}$ . This possible world thus contributes to the success probabilities of **alarm**, **calls(mary)** and **call**, but not to the one of **calls(john)**. Similarly, starting from the world where all probabilistic facts are false, we obtain the empty set as the least Herbrand model, and this world thus does not contribute to the probability of any atom.

An alternative to forward reasoning is *backward reasoning*, also known as SLD-resolution or proving, which we again illustrate for our example possible world (3). It starts from a given query, e.g., **call**, and uses the rules in the opposite direction: in order to prove a fact appearing in the head of a clause, we have to prove all literals in the clause's body. For instance, based on the last rule, to prove **call**, we need to prove **calls(X)** for some instantiation of  $X$ . Using the third rule, this means proving **alarm**, **hears\_alarm(X)**. To prove **alarm**, we could use the first rule and prove **earthquake**, but this fails for our choice of facts, as there is no rule (or fact) for the latter. We thus *backtrack* to the second rule for **alarm**, which requires proving **burglary**, which is proven by the corresponding fact. Finally, we prove **hears\_alarm(X)** using the fact **hears\_alarm(mary)**, substituting **mary** for  $X$ , which completes the proof for **call**.

Going over all possible worlds in this way, we obtain the success probability of **calls(mary)**,  $P(\text{calls}(\text{mary})) = 0.196$ , as the sum of the probabilities of six possible worlds (listed in Table 1).

Clearly, enumerating all possible worlds is infeasible for larger programs; we will discuss alternative inference techniques from the literature in Section 3.

For ease of modeling (and to allow for countably infinite sets of probabilistic facts), probabilistic languages such as ICL and ProbLog use *non-ground probabilistic facts* to define sets of random variables. All ground instances of such a fact are mutually independent and share the same probability value. As an example, consider a simple coin game which can be won either by throwing two times heads

| world                                      | calls(john) | probability  |
|--|-------------|--|
| $b \wedge \neg e \wedge hm \wedge \neg hj$ | false       | $0.1 \cdot (1 - 0.2) \cdot 0.7 \cdot (1 - 0.4) = 0.0336$ |
| $b \wedge \neg e \wedge hm \wedge hj$      | true        | $0.1 \cdot (1 - 0.2) \cdot 0.7 \cdot 0.4 = 0.0224$       |
| $b \wedge e \wedge hm \wedge \neg hj$      | false       | $0.1 \cdot 0.2 \cdot 0.7 \cdot (1 - 0.4) = 0.0084$       |
| $b \wedge e \wedge hm \wedge hj$           | true        | $0.1 \cdot 0.2 \cdot 0.7 \cdot 0.4 = 0.0056$             |
| $\neg b \wedge e \wedge hm \wedge \neg hj$ | false       | $(1 - 0.1) \cdot 0.2 \cdot 0.7 \cdot (1 - 0.4) = 0.0756$ |
| $\neg b \wedge e \wedge hm \wedge hj$      | true        | $(1 - 0.1) \cdot 0.2 \cdot 0.7 \cdot 0.4 = 0.0504$       |

**Table 1** The possible worlds of program (1) where `calls(mary)` is true.

or by cheating. This game can be modeled by the program below. The probability to win the game is then defined by the success probability  $P(\text{win})$ .

```

0.5 :: heads(X).           0.2 :: cheat_successfully.
win :- cheat_successfully.
win :- heads(1), heads(2).

```

Legal groundings of such facts can also be restricted by providing a domain, as in the following variant of our alarm example where all persons have the same probability of independently hearing the alarm:

```

0.1 :: burglary.  0.2 :: earthquake
0.7 :: hears_alarm(X) :- person(X).
person(mary). person(john). person(bob). person(ann).
alarm :- earthquake.
alarm :- burglary.
calls(X) :- alarm, hears_alarm(X).
call :- calls(X).

```

If such domains are defined purely logically, without using probabilistic facts, the basic distribution is still well defined.

It is often assumed that probabilistic facts do not unify with other probabilistic facts or heads of rules. This ensures that the label of a probabilistic fact equals the fact's success probability, and achieves a clean separation between the facts  $F$  used to define the distribution  $P_F$  and their logical consequences given by the set of rules  $R$ . We discuss dropping this assumption by using independent causation below.

## 2.2 Probabilistic Choices

As already noted by Sato (1995), probabilistic facts (or binary switches) are expressive enough to represent a wide range of models, including Bayesian networks, Markov chains and hidden Markov models. However, for ease of modeling, it is often more convenient to use multi-valued random variables instead of binary ones. The concept commonly used to realize such variables in the distribution semantics is a probabilistic choice, that is, a finite set of ground atoms exactly one of which is true in any possible world. Examples of primitives implementing the concept of a probabilistic choice are the *probabilistic alternatives* of the independent choice logic (ICL) (Poole, 2000) and probabilistic Horn abduction (PHA) (Poole, 1993), the

*multi-ary random switches* of PRISM (Sato and Kameya, 2001), the *probabilistic clauses* of stochastic logic programs (SLPs) (Muggleton, 1996), and the *annotated disjunctions* of logic programs with annotated disjunctions (LPADs) (Vennekens et al, 2004), or the *CP-events* of CP-logic (Vennekens, 2007). We restrict the following discussion to annotated disjunctions (Vennekens et al, 2004), using the notation introduced below, and return to the relation between these languages in Section 2.3.

An *annotated disjunction* (AD) is an expression of the form

$$p_1 :: h_1; \dots ; p_N :: h_N :- b_1, \dots, b_M.$$

where  $b_1, \dots, b_M$  is a possibly empty conjunction of literals, the  $p_i$  are probabilities and  $\sum_{i=1}^N p_i \leq 1$ . Considered in isolation, an annotated disjunction states that if the body  $b_1, \dots, b_M$  is true at most one of the  $h_i$  is true as well, where the choice is governed by the probabilities (see below for interactions between multiple ADs with unifying atoms in the head). If the  $p_i$  in an annotated disjunction do not sum to 1, there is also the case that nothing is chosen. The probability of this event is  $1 - \sum_{i=1}^n p_i$ . A probabilistic fact is thus a special case of an AD with a single head atom and empty body.

For instance, consider the following program:

$$\begin{aligned} 0.4 &:: \text{draw}. \\ \frac{1}{3} &:: \text{color}(\text{green}); \frac{1}{3} :: \text{color}(\text{red}); \frac{1}{3} :: \text{color}(\text{blue}) :- \text{draw}. \end{aligned}$$

The probabilistic fact states that we draw a ball from an urn with probability 0.4, and the annotated disjunction states that if we draw a ball, the color is picked uniformly among **green**, **red** and **blue**. The program thus has four possible worlds, the empty one (with probability 0.6), and three that each contain **draw** and one of the **color** atoms (each with probability 0.4/3). As for probabilistic facts, a non-ground AD denotes the set of all its groundings, and for each such grounding, choosing one of its head atoms to be true is seen as an independent random event. That is, the annotated disjunction

$$\frac{1}{3} :: \text{color}(\text{B}, \text{green}); \frac{1}{3} :: \text{color}(\text{B}, \text{red}); \frac{1}{3} :: \text{color}(\text{B}, \text{blue}) :- \text{ball}(\text{B}).$$

defines an independent probabilistic choice of color for each ball B.

As noted already by Vennekens et al (2004), the probabilistic choice over head atoms in an annotated disjunction can equivalently be expressed using a set of logical clauses, one for each head, and a probabilistic choice over facts added to the bodies of these clauses, e.g.

$$\begin{aligned} \text{color}(\text{B}, \text{green}) &:- \text{ball}(\text{B}), \text{choice}(\text{B}, \text{green}). \\ \text{color}(\text{B}, \text{red}) &:- \text{ball}(\text{B}), \text{choice}(\text{B}, \text{red}). \\ \text{color}(\text{B}, \text{blue}) &:- \text{ball}(\text{B}), \text{choice}(\text{B}, \text{blue}). \\ \frac{1}{3} &:: \text{choice}(\text{B}, \text{green}); \frac{1}{3} :: \text{choice}(\text{B}, \text{red}); \frac{1}{3} :: \text{choice}(\text{B}, \text{blue}). \end{aligned}$$

This example illustrates that annotated disjunctions define a distribution  $P_F$  over basic facts as required in the distribution semantics, but can simplify modeling by directly expressing probabilistic consequences.



*Independent Causes* Some languages, e.g. ICL (Poole, 2008), assume that head atoms in the same or different annotated disjunctions cannot unify with one another, while others, e.g., LPADs (Vennekens et al, 2004), do not make this restriction, but instead view each annotated disjunction as an independent cause for the conclusions to hold. In that case, the structure of the program defines the combined effect of these causes, similarly to how the two clauses for `alarm` in our earlier example (1) combine the two causes `burglary` and `earthquake`. We illustrate this on the Russian roulette example by Vennekens et al (2009), which involves two guns.

$$\begin{aligned} \frac{1}{6} &:: \text{death}:- \text{pull\_trigger}(\text{left\_gun}). \\ \frac{1}{6} &:: \text{death}:- \text{pull\_trigger}(\text{right\_gun}). \end{aligned}$$

Each gun is an independent cause for death. Pulling both triggers will result in death being true with a probability of  $1 - (1 - \frac{1}{6})^2$ , which exactly corresponds to the probability of `death` being proven via the first or via the second annotated disjunction (or both). Assuming independent causes closely corresponds to the noisy-or combining rule that is often employed in the Bayesian network literature, cf. Section 5.

### 2.3 Discussion

While we have distinguished probabilistic facts and probabilistic choices here for ease of exposition, both views are closely connected and exchangeable from the perspective of expressivity. Indeed, as mentioned above, a probabilistic fact  $p :: f$  directly corresponds to an annotated disjunction  $p :: f :- \text{true}$  with a single atom in the head and an empty (or true) body. Conversely, each annotated disjunction can – for the purpose of calculating success probabilities – be equivalently represented using a set of probabilistic facts and deterministic clauses, which together simulate a sequential choice mechanism; we refer to Appendix B for details.

Annotated disjunctions are one of many primitives that implement probabilistic choices. In ICL, a probabilistic choice is implemented as a probabilistic alternative of the form

$$\text{prob } a_1 : p_1, \dots, a_n : p_n.$$

As pointed out by Vennekens et al (2004), the probabilistic alternatives of ICL map onto annotated disjunctions (and vice versa), that is, the alternative above rewrites as

$$p_1 :: a_1; \dots; p_n :: a_n :- \text{true}.$$

Similarly, for PRISM, a multi-ary random switch  $\text{msw}(\text{term}, V)$  with identifier `term`, values  $v_1, \dots, v_n$  and probabilities  $p_1, \dots, p_n$  directly corresponds to an annotated disjunction

$$p_1 :: \text{msw}(\text{term}, v_1); \dots; p_n :: \text{msw}(\text{term}, v_p) :- \text{true}.$$

However, a key distinguishing feature of PRISM is its use of *stochastic memoization*, that is, the fact that different occurrences of the same `msw` atom denote independent random variables; we will discuss this aspect in more detail in Section 4.4. Finally, it is well-known that the probabilistic clauses of SLPs map onto the switches of PRISM (Cussens, 2005). The correspondence is direct in the case

of a predicate in an SLP being defined by a set of probabilistic clauses that all have an empty body, and uses the same idea of explicitly representing a choice in the body of a clause as illustrated for annotated disjunctions above (page 7) else. An example can be found in Section 4.4.

## 2.4 Inference Tasks

In probabilistic programming and statistical relational learning, the following inference tasks have been considered:

- In the  $SUCC(q)$  task, a ground query  $q$  is given, and the task is to compute

$$SUCC(q) = P(q),$$

the success probability of the query as specified in Equation (4).<sup>3</sup>

- In the  $MARG(Q | e)$  task, a set  $Q$  of ground atoms of interest, the query atoms, and a ground query  $e$ , the evidence, are given. The task is to compute the marginal probability distribution of each atom  $q \in Q$  given the evidence,

$$MARG(Q | e) = P(q|e) = \frac{P(q \wedge e)}{P(e)}.$$

The  $SUCC(q)$  task corresponds to the special case of the  $MARG(Q | e)$  task with  $Q = \{q\}$  and  $e = \text{true}$  (and thus  $P(e) = 1$ ).

- The  $MAP(Q | e)$  task is to find the most likely truth-assignment  $v$  to the atoms in  $Q$  given the evidence  $e$ , that is, to compute

$$MAP(Q | e) = \arg \max_v P(Q = v|e)$$

- The  $MPE(e)$  task is to find the most likely world where the given evidence query  $e$  holds. Let  $U$  be the set of all atoms in the Herbrand base that do not occur in  $e$ . Then, the task is to compute the most likely truth-assignment  $u$  to the atoms in  $U$ ,

$$MPE(e) = MAP(U | e).$$

- In the  $VIT(q)$  task, a query  $q$  (but no evidence) is given, and the task is to find a Viterbi proof of  $q$ . Let  $E(q)$  be the set of all explanations or proofs of  $q$ , that is, of all sets  $F'$  of ground probabilistic atoms for which  $q$  is true in the corresponding possible world. Then, the task is to compute

$$VIT(q) = \arg \max_{X \in E(q)} P\left(\bigwedge_{f \in X} f\right).$$

To illustrate, consider our initial alarm example (1) with  $e = \text{calls}(\text{mary})$  and  $Q = \{\text{burglary}, \text{calls}(\text{john})\}$ . The worlds where the evidence holds are listed in Table 1, together with their probabilities. The answer to the MARG task is

$$\begin{aligned} P(\text{burglary}|\text{calls}(\text{mary})) &= 0.07/0.196 = 0.357 \\ P(\text{calls}(\text{john})|\text{calls}(\text{mary})) &= 0.0784/0.196 = 0.4 \end{aligned}$$

<sup>3</sup> Non-ground queries have been considered as well, in which case the success probability corresponds to the probability that  $q\theta$  is true for some grounding substitution  $\theta$ .

The answer to the MAP task is `burglary=false, calls(john)=false`, as its probability  $0.0756/0.196$  is higher than  $0.028/0.196$  (for true, true),  $0.042/0.196$  (for true, false) and  $0.0504/0.196$  (for false, true). The world returned by MPE is the one corresponding to the set of facts `{earthquake, hears_alarm(mary)}`. Finally, the Viterbi proof of query `calls(john)`, which does not take into account evidence, is  $e \wedge hj$ , as  $0.2 \cdot 0.4 > 0.1 \cdot 0.4$  (for  $b \wedge hj$ ), whereas the Viterbi proof for query `burglary` is its only proof  $b$ .

### 3 Inference

We now provide an overview of existing inference approaches in probabilistic (logic) programming. As most existing work addresses the SUCC task of computing success probabilities, cf. Equation (4), we focus on this task here, and mention other tasks in passing where appropriate. For simplicity, we assume probabilistic facts as basic building blocks. Computing marginals under the distribution semantics has to take into account both *probabilistic* and *logical* aspects. We therefore distinguish between *exact* inference and approximation using either *bounds* or *sampling* on the probabilistic side, and between methods based on *forward* and *backward reasoning* and *grounding to CNF* on the logical side. Systems implementing (some of) these approaches include the ICL system AILog2,<sup>4</sup> the PRISM system,<sup>5</sup> the ProbLog implementations ProbLog1<sup>6</sup> and ProbLog2,<sup>7</sup> and the LPAD implementations cplint<sup>8</sup> and PITA.<sup>9</sup> General statements about systems in the following refer to these six systems.

#### 3.1 Exact Inference

As most methods for exact inference can be viewed as operating (implicitly or explicitly) on a propositional logic representation of all possible worlds that entail the query  $q$  of interest, we first note that this set of possible worlds is given by the following formula in disjunctive normal form (DNF)

$$DNF(q) = \bigvee_{\substack{F' \subseteq F \\ \exists \theta F' \cup R \models q\theta}} \left( \bigwedge_{f_i \in F'} f_i \wedge \bigwedge_{f_i \in F \setminus F'} \neg f_i \right) \quad (6)$$

and that the structure of this formula exactly mirrors that of Equation (5) defining the success probability in the case of probabilistic facts, where we replace summation by disjunction, multiplication by conjunction, and probabilities by truth values of random variables (or facts).

<sup>4</sup> <http://artint.info/code/ailog/ailog2.html>

<sup>5</sup> <http://sato-www.cs.titech.ac.jp/prism/>

<sup>6</sup> included in YAP Prolog, <http://www.dcc.fc.up.pt/~vsc/Yap/>

<sup>7</sup> <http://dtai.cs.kuleuven.be/problog/>

<sup>8</sup> included in YAP Prolog, <http://www.dcc.fc.up.pt/~vsc/Yap/>

<sup>9</sup> included in XSB Prolog, <http://xsb.sourceforge.net/>

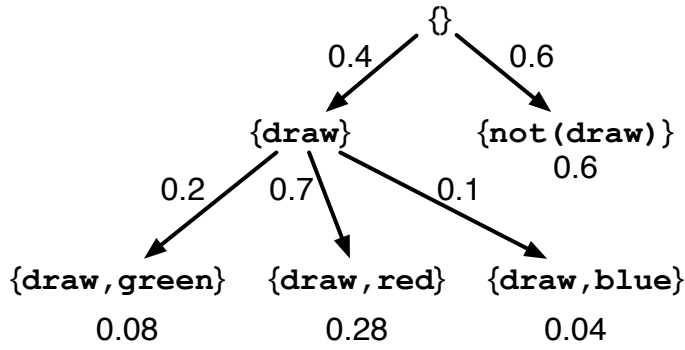


Fig. 1 Forward reasoning example.

In our initial alarm example (1), the DNF corresponding to `calls(mary)` contains the worlds shown in Table 1, and thus is

$$\begin{aligned}
 & (b \wedge e \wedge hm \wedge hj) \vee (b \wedge e \wedge hm \wedge \neg hj) \vee (b \wedge \neg e \wedge hm \wedge hj) \\
 & \vee (b \wedge \neg e \wedge hm \wedge \neg hj) \vee (\neg b \wedge e \wedge hm \wedge hj) \vee (\neg b \wedge e \wedge hm \wedge \neg hj).
 \end{aligned} \tag{7}$$

*Forward Reasoning:* Following the definition of the semantics of CP-logic (Vennekens et al, 2009), forward reasoning can be used to build a tree whose leaves correspond to possible worlds, on which success probabilities can be calculated. Specifically, the root of the tree is the empty set, and in each node, one step of forward reasoning is executed, creating a child for each possible outcome in the case of probabilistic facts or annotated disjunctions. For instance, consider the program

```

0.4 :: draw.
0.2 :: green; 0.7 :: red; 0.1 :: blue :- draw.
    
```

As illustrated in Figure 1, the first step using the probabilistic fact `draw` adds two children to the root, one containing `draw`, and one containing `not(draw)`. In the latter case, the body of the AD is false and thus no further reasoning steps are possible. For the world where `draw` is true, the AD introduces three children, adding `green`, `red` and `blue`, respectively, and no further reasoning steps are possible in the resulting worlds. Thus, each path from the root to a leaf constructs one possible world, whose probability is the product of assignments made along the path. Domains for non-ground facts have to be explicitly provided to ensure termination. While this approach clearly illustrates the semantics, even in the finite case, it suffers from having to enumerate all possible worlds, and is therefore not used in practice. A possible solution to this could be based on the magic set transformation (Bancilhon et al, 1986), which restricts forward reasoning to atoms relevant for deriving a query.

*Backward Reasoning:* Probably the most common inference strategy in probabilistic logic programming is to collect all possible *proofs* or *explanations* of a given query using backward reasoning, represent them in a suitable data structure, and

compute the probability on that structure. As discussed in Section 2.4, an explanation is a partial truth value assignment to probabilistic facts that is sufficient to prove the query via SLD-resolution.

For instance,  $b \wedge hm$  is the explanation for `calls(mary)` given by the derivation discussed in Section 2.1 (page 5), as it depends on `burglary` and `hears_alarm(mary)` being true, but not on any particular truth values of `earthquake` and `hears_alarm(john)`. This query has a second proof,  $e \wedge hm$ , obtained by using the first clause for `alarm` during backward reasoning. We can describe the set of possible worlds where `calls(mary)` is true by the disjunction of all proofs of the query,

$$(b \wedge hm) \vee (e \wedge hm)$$

which is more compact than the disjunction (7) explicitly listing the six possible worlds.

We cannot, however, calculate the probability of this more compact DNF by simply replacing conjunction by multiplication and disjunction by addition as we did for the longer DNF above. The reason is that the two proofs are *not mutually exclusive*, that is, they can be true in the same possible world. Specifically, in our example this holds for the two worlds  $b \wedge e \wedge hm \wedge hj$  and  $b \wedge e \wedge hm \wedge \neg hj$ , and the probability of these worlds,

$$0.1 \cdot 0.2 \cdot 0.7 \cdot 0.4 + 0.1 \cdot 0.2 \cdot 0.7 \cdot (1 - 0.4) = 0.014$$

is exactly the difference between 0.21 as obtained by the direct sum of products  $0.1 \cdot 0.7 + 0.2 \cdot 0.7$  and the true probability 0.196. This is also known as the *disjoint-sum-problem*, which is #P-complete (Valiant, 1979).

Existing languages and systems approach the problem from different angles. PHA (Poole, 1992) and PRISM (Sato and Kameya, 2001) rely on the *exclusive explanation assumption*, that is, they assume that the structure of the program guarantees mutual exclusiveness of all conjunctions in the DNF, which allows one to evaluate it as a direct sum of products (as done in the PRISM system). This assumption allows for natural modeling of many models, including e.g., probabilistic grammars and Bayesian networks, but prevents direct modeling of e.g., connection problems over uncertain graphs where each edge independently exists with a certain probability, or simple variations of Bayesian network models such as our running example.

ICL (Poole, 2000) is closely related to PHA, but does not assume exclusive explanations. Poole instead suggests symbolic disjoining techniques to split explanations into mutually exclusive ones (implemented in AILog2). The ProbLog1 implementation of ProbLog (De Raedt et al, 2007; Kimmig et al, 2011b) has been the first probabilistic programming system representing DNFs as Binary Decision Diagrams (BDDs), an advanced data structure that disjoins explanations. This technique has subsequently also been adopted for ICL and LPADs in the cplint and PITA systems (Riguzzi, 2009; Riguzzi and Swift, 2011). AILog2 and cplint also support computing conditional probabilities.

Riguzzi (2014) has introduced an approach called PITA(OPT) that automatically recognizes certain independencies that allow one to avoid the use of disjoining techniques when computing marginal probabilities.

Given its focus on proofs, backward reasoning can easily be adapted to solve the VIT task of finding most likely proofs, as done in the PRISM, ProbLog1 and PITA systems.

*Reduction to Weighted Model Counting:* A third way to approach the logic side of inference in probabilistic logic programming has been suggested by Fierens et al (2011, 2014), who use the propositional logic semantics of logic programming to reduce MARG inference to *weighted model counting* (WMC) and MPE inference to weighted MAX-SAT. The first step again builds a Boolean formula representing all models where the query is true, but this time, using conjunctive normal form (CNF), and associating a weight with every literal in the formula. More specifically, it grounds the parts of the logic program relevant to the query (that is, the rule groundings contributing to a proof of the query, as determined using backward reasoning), similar to what happens in answer set programming, transforms this ground program into an equivalent CNF based on the semantics of logic programming, and defines the weight function for the second step using the given probabilities. The second step can then use any existing approach to WMC or weighted MAX-SAT, such as representing the CNF as an *sd-DNNF*, a data structure on which WMC can be performed efficiently.

For instance, the relevant ground program for `calls(mary)` in our initial alarm example (1) is

```

0.1 :: burglary.    0.7 :: hears_alarm(mary).
0.2 :: earthquake.
                alarm :- earthquake.
                alarm :- burglary.
calls(mary) :- alarm,hears_alarm(mary).

```

Next, the rules in the ground program are translated to equivalent formulas in propositional logic, taking into account that their head atoms can only be true if a corresponding body is true:

$$\begin{aligned}
 \text{alarm} &\leftrightarrow \text{earthquake} \vee \text{burglary} \\
 \text{calls}(\text{mary}) &\leftrightarrow \text{alarm} \wedge \text{hears\_alarm}(\text{mary})
 \end{aligned}$$

The conjunction of these formulas is then transformed into CNF as usual in propositional logic. The weight function assigns the corresponding probabilities to literals of probabilistic facts, e.g.,  $w(\text{burglary}) = 0.1$ ,  $w(\neg \text{burglary}) = 0.9$ , and 1.0 to all other literals, e.g.,  $w(\text{calls}(\text{mary})) = w(\neg \text{calls}(\text{mary})) = 1.0$ . The weight of a model is the product of the weights of the literals in the model, and the WMC of a formula the sum of weights of all its models. As a full truth value assignment  $F'$  to probabilistic facts  $F$  uniquely determines the truth values of all other literals, there is a single model extending  $F'$ , with weight  $P_F(F') \cdot 1.0 \cdot \dots \cdot 1.0 = P_F(F')$ . The WMC of a formula thus exactly corresponds to the success probability. Evidence can directly be incorporated by conjoining it with the CNF. Exact MARG inference using this approach is implemented in ProbLog2.

*Lifted Inference* is a central research topic in statistical relational learning today (Kersting, 2012; Poole, 2003). Lifted inference wants to realize probabilistic logic inference at the lifted, that is, non-grounded level in the same way that resolution realizes this for logical inference. The problem of lifted inference can be illustrated on the following example (cf. also Poole (2008)):

```

p :: famous(Y).
popular(X) :- friends(X, Y), famous(Y).

```

In this case  $P(\text{popular}(\text{john})) = 1 - (1 - p)^m$  where  $m$  is the number of friends of `john`, that is, to determine the probability that `john` is popular, it suffices to know how many friends `john` has. We do not need to know the identities of these friends, and hence, need not ground the clauses.

Various techniques for lifted inference have been obtained over the past decade. For instance, Poole (2003) shows how variable elimination, a standard approach to probabilistic inference in graphical models, can be lifted and Van den Broeck et al (2011) studied weighted model counting for first order probabilistic logic using a generalization of d-DNNFs for first order logic. Lifted inference techniques are – to the best of our knowledge – not yet supported by current probabilistic logic programming language implementations, which explains why we do not provide more details in this paper. It remains a challenge for further work; but see Van den Broeck et al (2014); Bellodi et al (2014) for recent progress. A recent survey on lifted inference is provided by Kersting (2012).

### 3.2 Approximate Inference using Bounds

As the probability of a set of possible worlds monotonically increases if more models are added, hard lower and upper bounds on the success probability can be obtained by considering a subset or a superset of all possible worlds where a query is true. For instance, let  $W$  be the set of possible worlds where a query  $q$  holds. The success probability of  $q$  thus is the sum of the probabilities of all worlds in  $W$ . If we restrict this sum to a subset of  $W$ , we obtain a lower bound, and an upper bound if we sum over a superset of  $W$ . In our example, as `calls(mary)` is true in  $b \wedge e \wedge hm \wedge hj$ , but false in  $b \wedge e \wedge \neg hm \wedge hj$ , we have

$$0.1 \cdot 0.2 \cdot 0.7 \cdot 0.4 \leq P(\text{calls}(\text{mary})) \leq 1 - (0.1 \cdot 0.2 \cdot (1 - 0.7) \cdot 0.4).$$

In practice, this approach is typically used with the DNF obtained by backward reasoning, that is, the set of proofs of the query, rather than with the possible worlds directly. This has initially been suggested for PHA by Poole (1992), and later also been adapted for ProbLog (De Raedt et al, 2007; Kimmig et al, 2008) and LPADs (Bragaglia and Riguzzi, 2010). The idea is to maintain a set of partial derivations during backward reasoning, which allows one to, at any point, obtain a lower bound based on all complete explanations or proofs found so far, and an upper bound based on those together with all partial ones (based on the assumption that those will become proofs with probability one). For instance,  $(e \wedge hm) \vee b$  provides an upper bound of 0.226 for the probability of `calls(mary)` based on the proof  $e \wedge hm$  (which provides the corresponding lower bound 0.14) and the partial derivation  $b$  (which still requires to prove `hears_alarm(mary)`). Different search strategies are possible here, including e.g., iterative deepening or best first search. Lower bounds based on a fixed number of proofs have been proposed as well, either using the  $k$  explanations with highest individual probabilities (Kimmig et al, 2011b), or the  $k$  explanations chosen by a greedy procedure that maximizes

the probability an explanation adds to the one of the current set (Renkens et al, 2012). Approximate inference using bounds is available in ProbLog1, cplint, and ProbLog2.

### 3.3 Approximate Inference by Sampling

While probabilistic logic programming often focuses on exact inference, approximate inference by sampling is probably the most popular approach to inference in many other probabilistic languages. Sampling uses a large number of random executions or randomly generated possible worlds, from which the probability of a query is estimated as the fraction of samples where the query holds. For instance, samples can be generated by randomly choosing truth values of probabilistic facts as needed during backward reasoning, until either a proof is found or all options are exhausted (Kimmig et al, 2008; Bragaglia and Riguzzi, 2010; Riguzzi, 2013b). Fierens et al (2014) have used MC-SAT (Poon and Domingos, 2006) to perform approximate WMC on the CNF representing all models. Systems for languages that specify generative models, such as BLOG (Milch et al, 2005) and distributional clauses (Gutmann et al, 2011), cf. Sec. 4.2, often use forward reasoning to generate samples. A popular approach to sampling are MCMC algorithms, which, rather than generating each sample from scratch, generate a sequence of samples by making random modifications to the previous sample based on a so-called proposal distribution. This approach has been used e.g., for the probabilistic functional programming language Church (Goodman et al, 2008), for BLOG (Arora et al, 2010), and for the probabilistic logic programming languages PRISM (Sato, 2011) and ProbLog (Moldovan et al, 2013). ProbLog1 and cplint provide inference techniques based on backward sampling, and the PRISM system includes MCMC inference.

## 4 Probabilistic Programming Concepts

While probabilistic programming languages based on the distribution semantics as discussed so far are expressive enough for a wide range of models, an important part of their power is their support for additional programming concepts. Based on primitives used in a variety of probabilistic languages, we discuss a range of such concepts next, also touching upon their implications for inference.

### 4.1 Flexible Probabilities

A probabilistic fact with *flexible probability* is of the form  $P :: atom$  where  $atom$  contains the logical variable  $P$  that has to be instantiated to a probability when using the fact. The following example models drawing a red ball from an urn with  $R$  red and  $G$  green balls, where each ball is drawn with uniform probability from the urn:

```
Prob :: red(Prob).
draw_red(R, G):- Prob is R/(R + G),
                 red(Prob).
```



The combination of flexible probabilities and Prolog code offers a powerful tool to compute probabilities on-the-fly, cf. e.g., (Poole, 2008). Flexible probabilities have also been used in extended SLPs (Angelopoulos and Cussens, 2004), and are supported by the probabilistic logic programming systems AILog2, ProbLog1, cplint and ProbLog2. For such facts to be meaningful, their probabilities have to be bound at inference time. Probabilistic facts with flexible probabilities are thus easily supported by backward inference as long as these facts are ground on calling. In the example, this holds for ground queries such as `draw_red(3,1)`, which binds `Prob` to 0.75. They however cannot directly be used with exact forward inference, as they abbreviate an infinite set of ground facts and thus would create an infinite tree of possible worlds.<sup>10</sup>

## 4.2 Distributional Clauses

Annotated disjunctions – as specified in Section 2.2 – are of limited expressivity, as they can only define distributions over a fixed, finite number of head elements. While more flexible discrete distributions can be expressed using a combination of flexible probabilities and Prolog code, this may require significant programming effort. Gutmann et al (2010) introduce Hybrid ProbLog, an extension of ProbLog to continuous distributions, but their inference approach based on exact backward reasoning and discretization severely limits the use of such distributions. To alleviate these problems, *distributional clauses* were introduced by Gutmann et al (2011), whom we closely follow.

A *distributional clause* is a clause of the form

$$h \sim \mathcal{D} :- b_1, \dots, b_n.$$

where  $\sim$  is a binary predicate used in infix notation. Similarly to annotated disjunctions, the head  $(h \sim \mathcal{D})\theta$  of a distributional clause is defined for a grounding substitution  $\theta$  whenever  $(b_1, \dots, b_n)\theta$  is true in the semantics of the logic program. Then the distributional clause defines the random variable  $h\theta$  as being distributed according to the associated distribution  $\mathcal{D}\theta$ . Possible distributions include finite discrete distributions such as a uniform distribution, discrete distributions over infinitely many values, such as a Poisson distribution, and continuous distributions such as Gaussian or Gamma distributions. The outcome of a random variable  $h$  is represented by the term  $\simeq(h)$ . Both random variables  $h$  and their outcome  $\simeq(h)$  can be used as other terms in the program. However, the typical use of terms  $\simeq(h)$  is inside comparison predicates such as `equal/2` or `lessthan/2`.<sup>11</sup> In this case these predicates act in the same way as probabilistic facts in Sato’s distribution semantics. Indeed, depending on the value of  $\simeq(h)$  (which is determined probabilistically) they will be true or false.

<sup>10</sup> If only finitely many different instances of such a fact are relevant for any possible world of a given program, a mechanism similarly to the magic set transformation (Bancilhon et al, 1986) may circumvent this problem.

<sup>11</sup>  $p/n$  denotes a predicate  $p$  with  $n$  arguments, cf. Appendix A

Consider the following distributional clause program.

```
color(B) ~ discrete((0.7 : green), (0.3 : blue)) :- ball(B).
    diameter(B, MD) ~ gamma(MD1, 20) :- mean_diameter(≃(color(B)), MD),
        MD1 is 1/20 * MD.
    mean_diameter(green, 15).
    mean_diameter(blue, 25).
ball(1). ball(2). ... ball(k).
```

The first clause states that for every ball  $B$ , there is a random variable  $\text{color}(B)$  whose value is either **green** (with probability 0.7) or **blue** (with probability 0.3). This discrete distribution directly corresponds to the one given by the annotated disjunction

```
0.7 :: color(B, green); 0.3 :: color(B, blue) :- ball(B).
```

The second distributional clause in the example defines a random variable  $\text{diameter}(B, MD)$  for each ball  $B$ . This random variable follows a Gamma distribution with parameters  $MD/20$  and  $20$ , where the mean diameter  $MD$  depends on the color of the ball.

Distributional clauses are the logic programming equivalent of the mechanisms employed in statistical relational languages such as Bayesian Logic (BLOG) (Milch et al, 2005), Church (Goodman et al, 2008) and IBAL (Pfeffer, 2001), which also use programming constructs to define generative process that can define new variables in terms of existing one.

As we have seen in the example, annotated disjunctions can easily be represented as distributional clauses with finite, discrete distributions. However, distributional clauses are more expressive than annotated disjunctions (and the standard distribution semantics) as they can also represent continuous distributions.

Performing inference with distributional clauses raises some extra difficulties (see (Gutmann et al, 2011) for more details). The reason for this is that continuous distributions (such as a Gaussian or a Gamma-distribution) have uncountable domains. Typical inference with constructs such as distributional clauses will therefore resort to sampling approaches in order to avoid the need for evaluating complex integrals. It is quite natural to combine sampling for distributional clauses with forward reasoning,<sup>12</sup> realizing a kind of generative process, though more complex strategies are also possible, cf. (Gutmann et al, 2011).

### 4.3 Unknown Objects

One of the key contributions of Bayesian Logic (BLOG) (Milch et al, 2005) is that it allows one to drop two common assumptions, namely the *closed world assumption* (all objects in the world are known in advance) and the *unique names assumption* (different terms denote different objects), which makes it possible to define probability distributions over outcomes with varying sets of objects. This

<sup>12</sup> Valid distributional clause programs are required to have finite support, which ensures termination.

is achieved by defining generative processes that construct possible worlds, where the existence and the properties of objects can depend on objects created earlier in the process.

As already shown by Poole (2008), such generative processes with an unknown number of objects can often be modeled using flexible probabilities and Prolog code to specify a distribution over the number of objects as done in BLOG. Distributional clauses simplify this modeling task, as they make introducing a random variable corresponding to this number straightforward. We can then use the `between/3` predicate to enumerate the objects in definitions of predicates that refer to them, cf. also Poole (2008). Below, the random variable `nballs` stands for the number of balls, which is Poisson distributed with  $\lambda = 6$ . For each possible value  $\simeq(\text{nballs})$ , the corresponding number of balls are generated which are identified by the numbers  $1, 2, \dots, \simeq(\text{nballs})$ .

```
nballs ~ poisson(6).
ball(N) : -between(1, \simeq(nballs), N).
```

#### 4.4 Stochastic Memoization

A key concept in the probabilistic functional programming language Church (Goodman et al, 2008) is *stochastic memoization*. If a random variable in Church is memoized, subsequent calls to it simply look up the result of the first call, similarly to *tabling* in logic programming (Warren, 1992). On the other hand, for random variables that are not memoized, each reference to the variable corresponds to an independent draw of an outcome. In contrast to Church, probabilistic logic programming languages and their implementations typically do not leave this choice to the user. In ICL, ProbLog, LPADs and the basic distribution semantics as introduced in (Sato, 1995), each ground probabilistic fact directly corresponds to a random variable, i.e., within a possible world, each occurrence of such a fact has the same truth value, and the fact is thus memoized. Furthermore, the probability of the fact is taken into account once when calculating the probability of a proof, independently of the number of times it occurs in that proof. While early versions of PRISM (Sato, 1995; Sato and Kameya, 1997) used binary or n-ary probabilistic choices with an argument that explicitly distinguished between different calls, this argument has been made implicit later on (Sato and Kameya, 2001), meaning that the PRISM implementation never memoizes the outcome of a random variable.

The difference between the two approaches can be explained using the following example. For the AD

$$\frac{1}{3} :: \text{color}(\text{green}); \frac{1}{3} :: \text{color}(\text{red}); \frac{1}{3} :: \text{color}(\text{blue}),$$

there are three answers to the goal `color(X), color(Y)`, one answer  $X = Y = c$  for each color  $c$  with probability  $\frac{1}{3}$ , as exactly one of the facts `color(c)` is true in each possible world when memoizing `color` (as in ProbLog and ICL). Asking the same question when `color` is not memoized (as in PRISM) results in 9 possible answers with probability  $\frac{1}{9}$  each. The query then – implicitly – corresponds to an

ICL or ProbLog query (`color(X, id1)`, `color(Y, id2)`), where the original AD is replaced by a non-ground variant

$$\frac{1}{3} :: \text{color}(\text{green}, \text{ID}); \frac{1}{3} :: \text{color}(\text{red}, \text{ID}); \frac{1}{3} :: \text{color}(\text{blue}, \text{ID})$$

and `id1` and `id2` are identifiers that are unique to the call.

Avoiding the memoization of probabilistic facts is necessary in order to model stochastic automata, probabilistic grammars, or stochastic logic programs (Muggleton, 1996) under the distribution semantics. There, a new rule is chosen randomly for each occurrence of the same nonterminal state/symbol/predicate within a derivation, and each such choice contributes to the probability of the derivation. The rules for a nonterminal thus form a family of independent identically distributed random variables, and each choice is automatically associated with one variable from this family.

Consider the following stochastic logic program. It is in fact a fragment of a stochastic definite clause grammar; the rules essentially encode the probabilistic context free grammar rules defining  $0.3 : vp \rightarrow verb$ ,  $0.5 : vp \rightarrow verb, np$  and  $0.2 : vp \rightarrow verb, pp$ . There are three rules for the non-terminal `vp` and each of them is chosen with an associated probability. Furthermore, the sum of the probabilities for these rules equals 1.

```
0.3 : vp(H, T) :- verb(H, T).
0.5 : vp(H, T) :- verb(H, H1), np(H1, T).
0.2 : vp(H, T) :- verb(H, H1), pp(H1, T).
```

This type of stochastic grammar can easily be simulated in the distribution semantics using one *dememoized* AD (or switch) for each non-terminal, a rule calling the AD to make the selection, and a set of rules linking the selection to the SLP rules:<sup>13</sup>

```
dememoize 0.3 :: vp_sel(rule1); 0.5 :: vp_sel(rule2); 0.2 :: vp_sel(rule3).
vp(H, T) :- vp_sel(Rule), vp_rule(Rule, H, T).
vp_rule(rule1, H, T) :- verb(H, T).
vp_rule(rule2, H, T) :- verb(H, H1), np(H1, T).
vp_rule(rule3, H, T) :- verb(H, H1), pp(H1, T).
```

All inference approaches discussed here naturally support stochastic memoization; this includes the ones implemented in AILog2, ProbLog1, ProbLog2, cplint and PITA. The PRISM system uses exact inference based on backward reasoning in the setting without stochastic memoization. In principle, stochastic memoization can be disabled in backward reasoning by automatically adding a unique identifier to each occurrence of the same random variable. However, for techniques that build propositional representations different from mutually exclusive DNFs

<sup>13</sup> The `dememoize` keyword is used for clarity here; it is not supported by existing systems.

(such as the DNFs of BDD-based methods and the CNFs when reducing to WMC), care is needed to ensure that these identifiers are correctly shared among different explanations when manipulating these formulas. Backward sampling can easily deal with both memoized and dememoized random variables. As only one possible world is considered at any point, each repeated occurrence of the same dememoized variable is simply sampled independently, whereas the first result sampled within the current world is reused for memoized ones. Forward sampling cannot be used without stochastic memoization, as it is unclear up front how many instances are needed. MCMC methods have been developed both for ProbLog (with memoization; implementation not available) and PRISM (without memoization; included in the PRISM system).

#### 4.5 Constraints

In knowledge representation, answer set programming and databases, it is common to allow the user to specify constraints on the possible models of a theory. In knowledge representation, one sometimes distinguishes inductive definitions (such as the definite clauses used in logic programming) from constraints. The former are used to define predicates, the latter impose constraints on possible worlds. While the use of constraints is still uncommon in probabilistic logic programming it is conceptually easy to accommodate this when working with the distribution semantics, cf. Fierens et al (2012). While such constraints can in principle be any first-order logic formula, we will employ clausal constraints here.

A *clausal constraint* is an expression of the form

$$h_1; \dots; h_N :- b_1, \dots, b_M.$$

where the  $h_i$  and  $b_j$  are literals. The constraint specifies that whenever  $(b_1 \dots b_M)\theta$  is true for a substitution  $\theta$  grounding the clause at least one of the  $h_i\theta$  must also be true. All worlds in which a constraint is violated become impossible, that is, their probability becomes 0. Constraints are very useful for specifying complex properties that possible worlds must satisfy.

To illustrate constraints, reconsider the alarm example and assume that it models a situation in the 1930s where there is only one phone available in the neighborhood implying that at most one person can call. This could be represented by the constraint

$$X = Y :- \text{calls}(X), \text{calls}(Y).$$

Imposing this constraint would exclude all worlds in which both Mary and John hear the alarm and call. The total probability mass for such worlds is  $0.4 \cdot 0.8 = 0.32$ . By excluding these worlds, one loses probability mass and thus has to normalize the probabilities of the remaining possible worlds. For instance, the possible world corresponding to the truth value assignment `burglary=true`, `earthquake=false`, `hears_alarm(mary)=true`, `hears_alarm(john)=false` yielded a probability mass of  $0.1 \cdot (1 - 0.2) \cdot 0.7 \cdot (1 - 0.6) = 0.0336$  without constraints. Now, when enforcing the constraint, one obtains  $0.0336 / (1 - 0.32)$ . Thus the semantics of constraints correspond to computing conditional probabilities where one conditions on the constraints being satisfied.

Handling constraints during inference has not been a focus of inference in probabilistic logic programming, and – to the best of our knowledge – no current system provides explicit support for both logic programming (or inductive definitions) and constraints.

Nevertheless, the constraints discussed here are related to Markov Logic (Richardson and Domingos, 2006), where first order logic formulas express soft and hard constraints. In Markov Logic, possible worlds or interpretations (i.e., truth value assignments to ground atoms) become less likely as they violate more groundings of soft constraints, and have probability zero if they violate some grounding of a hard constraint. It is well known that the transitive closure of a binary relation cannot be represented in first order logic, but requires second order constructs; see Huth and Ryan (2004) for a detailed formal discussion. Thus, the hard constraints in a Markov Logic network (MLN), which form a first order logic theory, cannot enforce probability zero for all worlds that do not respect the transitive closure of a binary relation.

On the other hand, the least Herbrand semantics of definite clause logic (i.e., pure Prolog) naturally represents such transitive closures. For instance, under the least Herbrand semantics,

$$\begin{aligned} \text{path}(A, C) &:- \text{edge}(A, C). \\ \text{path}(A, C) &:- \text{edge}(A, B), \text{path}(B, C). \end{aligned}$$

inductively defines `path` as the transitive closure of the relation `edge`, that is, a ground atom `path(a, c)` is true if and only if there is a sequence of true `edge` atoms connecting `a` and `c`. As an example, consider the case of two nodes 1 and 2 and a single edge pointing from 1 to 2, i.e., the edge relation is fully given by  $\{\text{edge}(1, 2)\}$ . Under least Herbrand semantics, there is a single model  $\{\text{edge}(1, 2), \text{path}(1, 2)\}$ , as the first clause requires that `path(1, 2)` is true, and no other facts can be derived. Thus, the probability of `path(1, 2)` is one, and the other three ground `path` atoms all have a probability of zero.

Note that an MLN that maps the definition above to the hard constraints

$$\begin{aligned} \text{edge}(A, C) &\rightarrow \text{path}(A, C) \\ \text{edge}(A, B) \wedge \text{path}(B, C) &\rightarrow \text{path}(A, C) \end{aligned}$$

enforces the transitivity property, as these rules are violated if there is a sequence of edges connecting two nodes, but the corresponding path atom is false. Still, these hard constraints do not correspond to the transitive closure, as they can for instance be satisfied by setting all ground path atoms to true, independently of the truth values of edge atoms. For our example with the single edge, the only ground MLN constraints that are not trivially satisfied based on the edge relation alone are

$$\begin{aligned} \text{true} &\rightarrow \text{path}(1, 2) \\ \text{true} \wedge \text{path}(2, 1) &\rightarrow \text{path}(1, 1) \\ \text{true} \wedge \text{path}(2, 2) &\rightarrow \text{path}(1, 2) \end{aligned}$$

Any model of the first constraint has to contain `path(1, 2)`, and thus trivially satisfies the third constraint as well. The second constraint then rules out all interpretations where `path(2, 1)` is true, but `path(1, 1)` is false, leading to a total of six possible models of the hard constraints:

```

{edge(1, 2), path(1, 2)}
{edge(1, 2), path(1, 2), path(1, 1)}
{edge(1, 2), path(1, 2), path(2, 1), path(1, 1)}
{edge(1, 2), path(1, 2), path(2, 2)}
{edge(1, 2), path(1, 2), path(2, 2), path(1, 1)}
{edge(1, 2), path(1, 2), path(2, 2), path(2, 1), path(1, 1)}

```

The only difference with the first order logic case is that an MLN assigns a probability to each of the models of its hard constraints based on which soft constraints hold in the model, and a probability of zero to all other interpretations. As there are no soft constraints in our example MLN, it assigns the same probability to each of the six models. Each ground path atom appears in a different number of models, and thus has a different, non-zero probability according to the MLN, whereas under Prolog's least Herbrand semantics, where transitive closure puts the full probability mass on the first of the MLN models, three such atoms have probability zero.

#### 4.6 Negation as Failure

So far, we have only considered probabilistic programs using definite clauses, that is, programs that only use positive literals in clause bodies, as those are guaranteed to have a unique model for any truth value assignment to basic probabilistic events. It is however possible to adopt Prolog's *negation as failure* on ground literals under the distribution semantics, as long as all truth values of derived atoms are still uniquely determined by those of the basic facts, cf., e.g., (Poole, 2000; Sato et al, 2005; Kimmig et al, 2009; Riguzzi, 2009; Fierens et al, 2014). Then, in each possible world, any ground query  $q$  either succeeds or fails, and its negation  $\text{not}(q)$  succeeds in exactly those worlds where  $q$  fails. Thus, the probability of a ground query  $\text{not}(q)$  is the sum of the probabilities of all possible worlds that do *not* entail  $q$ . Consider the following variant of our alarm example, where people also call if there is no alarm, but they have gossip to share:

```

0.1 :: burglary.    0.7 :: hears_alarm(mary).
0.2 :: earthquake. 0.4 :: hears_alarm(john).
0.3 :: has_gossip(mary). 0.6 :: has_gossip(john).

alarm :- earthquake.
alarm :- burglary.
calls(X) :- alarm, hears_alarm(X).
calls(X) :- not(alarm), has_gossip(X).
call :- calls(X).

```

The new rule for `calls(X)` can only possibly apply in worlds where `not(alarm)` succeeds, that is, `alarm` fails, which are exactly those containing neither `burglary` nor `earthquake`. Using *gm* as shorthand for `has_gossip(mary) = true`, we obtain the additional explanation  $\neg e \wedge \neg b \wedge gm$  for `calls(mary)`. Thus, in the presence of negation, explanations no longer correspond to sets of probabilistic facts as in the case of definite clause programs, but to sets of positive and negative literals

for probabilistic facts. While `not(alarm)` has a single explanation in this simple example, in general, explanations for negative literals can be much more complex, as they have to falsify every possible explanation of the corresponding positive literal by flipping the truth value of at least one probabilistic fact included in the explanation.

Negation as failure can be handled in forward and backward reasoning both for exact inference and for sampling, though forward reasoning has to ensure to proceed in the right order. Exact inference with backward reasoning often benefits from tabling. Negation as failure complicates approximate inference using bounds, as explanations for failing goals have to be considered (Renkens et al, 2014). AILog2, ProbLog1, ProbLog2, cplint and PITA all support negation as failure in their exact and sampling based approaches. The PRISM system follows the approach proposed by Sato et al (2005) and compiles negation into a definite clause program with unification constraints. Current MCMC approaches in probabilistic logic programming do not support negation beyond that of probabilistic facts.

#### 4.7 Second Order Predicates

When modeling relational domains, it is often convenient to reason over sets of objects that fulfill certain conditions, for instance, to aggregate certain values over them. In logic programming, this is supported by second order predicates such as `findall/3`, which collects all answer substitutions for a given query in a list. In the following example, the query `sum(S)` will first collect all arguments of `f/1` into a list and then sum the values using predicate `sum_list/2`, thus returning `S=3`.

```
f(1).
f(2).
sum(Sum) :- findall(X, f(X), L), sum_list(L, Sum).
```

Note that in Prolog, the list returned by `findall/3` is unique. Under the distribution semantics, however, this list will be different depending on which possible world is considered. To illustrate this, we replace the definition of `f/1` in our example with probabilistic facts:

```
0.1 :: f(1).
0.2 :: f(2).
sum(Sum) :- findall(X, f(X), L), sum_list(L, Sum).
```

We now have four sets of facts – `{f(1), f(2)}`, `{f(1)}`, `{f(2)}`, and `{}` – leading to the four possible worlds `{f(1), f(2), sum(3)}`, `{f(1), sum(1)}`, `{f(2), sum(2)}`, and `{sum(0)}`, as the answer list `L` is different in each case.

This behavior of second order predicates in the probabilistic setting can pose a challenge to inference. In principle, all inference approaches could deal with second order predicates. However, exact approaches would suffer from a blow-up, as they have to consider all possible lists of elements – and thus all possible worlds – explicitly, whereas in sampling, each sample only considers one such list. As far as we know, the only systems with some support for second order predicates are



`cpint`, which allows `bagof` and `setof` with one of its backward reasoning modules (Riguzzi, 2013a), and `ProbLog1`, whose backward sampling technique supports the second order predicates of the underlying YAP Prolog engine.

#### 4.8 Meta-Calls

One of the distinct features of programming languages such as Prolog and Lisp is the possibility to use programs as objects within programs, which enables meta-level programming. For their probabilistic extensions, this means reasoning about the probabilities of queries within a probabilistic program, a concept that is central to the probabilistic programming language Church, which builds upon a Lisp dialect (Goodman et al, 2008), and has also been considered with `ProbLog` (Mantadelis and Janssens, 2011). Possible uses of such a feature include filtering of proofs based on the probability of subqueries, or the dynamic definition of probabilities using queries, e.g., to implement simple forms of combining rules as in the following example, where `max_true(G1,G2)` succeeds with the success probability of the more likely argument.

```
P :: p(P).
max_true(G1, G2) :- prob(G1, P1), prob(G2, P2), max(P1, P2, P), p(P).
% rest of program (omitted)
```

In this section, we will use `prob(Goal, Prob)` to refer to an atom returning the success probability `Prob` of goal `Goal`, that is, implementing Equation (4). Note that such atoms are independent queries, that is, they do not share truth values of probabilistic facts with other atoms occurring in a derivation they are part of. Finally, if the second argument is a free variable upon calling, the success probability of `prob(goal, Prob)` is 1. For the sake of simplicity, we will assume here that the second argument will always be free upon calling.<sup>14</sup>

We extend the example above with the following program.

```
0.5 :: a. 0.7 :: b. 0.2 :: c.
d :- a, not(b).
e :- b, c.
```

Querying for `max_true(d,e)` using backward reasoning will execute two calls to `prob/2` in sequence: `prob(d,P1)` and `prob(e,P2)`. Note that if multiple calls to `prob/2` atoms occur in a proof, they are independent, i.e., even if they use the same probabilistic facts, those will (implicitly) correspond to different copies of the corresponding random variables local to that specific `prob/2` call. Put differently, `prob/2` encapsulates part of our possible worlds. In the example, `b` is thus a different random variable in `prob(d,P1)` and `prob(e,P2)`. The reason for this encapsulation is twofold: first, the probability of a goal is not influenced by calculating the probability of another (or even the same) event before, and second, as `prob/2` summarizes a set of possible worlds, the value of a random variable cannot be made visible to the outside world, as it may be different in different

<sup>14</sup> This is not a restriction, as `prob(Goal, c)` is equivalent to `prob(Goal, P), P=c`.

internal worlds. Indeed, in our example, `b` needs to be false to prove `d`, but true to prove `e`, so using the same random variable would force the top level query to be unprovable. We thus obtain a kind of hierarchically organized world: some probabilistic facts are used in the top level query, others are encapsulated in `prob/2` atoms, whose queries might in turn rely on both directly called probabilistic facts and further calls to `prob/2`. In our example, `prob(d,P1)` uses random variables corresponding to probabilistic facts `a` and `b`, returning  $P1 = 0.5 \cdot (1 - 0.7) = 0.15$ , `prob(e,P2)` uses random variables corresponding to probabilistic facts `b` and `c`, returning  $P2 = 0.7 \cdot 0.2 = 0.14$ , and the top level query `max_true(d,e)` uses probabilistic fact `p(0.15)` and has probability  $P(\text{more\_likely\_is\_true}(d,e)) = 0.15$ .

The probability of a derivation is determined by the probabilities of the probabilistic facts it uses outside all `prob/2` calls. Those facts define the possible worlds from the point of view of the top level query. In those worlds, the random variables of the encapsulated parts are hidden, as they have been aggregated by `prob/2`. Returning to our example and abstracting from the concrete remainder of the program, we observe that for any given pair of goals `g1,g2` and suitable program defining those goals, `max_true(g1,g2)` has exactly one proof: the first two body atoms always succeed and return the probabilities of the goals, the third atom deterministically finds the maximum  $m$  of the two probabilities, and the proof finally uses a single random variable `p(m)` with probability  $m$ . Thus, the query indeed succeeds with the probability of the more likely goal.

Another example for the use of `prob/2` is filtering goals based on their probability:

```
almost_always_false(G) :- prob(G,P), P < 0.00001.
% rest of program (omitted)
```

Note that in contrast to the previous example, this is a purely logical decision, that is, the success probability will be either 0 or 1 depending on the goal  $G$ .

To summarize, using meta-calls to turn probabilities into usable objects in probabilistic logic programming is slightly different from the other probabilistic programming concepts considered in this paper: it requires a notion of encapsulation or hierarchical world structure and cannot be interpreted directly on the level of individual possible worlds for the entire program.

Mantadelis and Janssens (2011) introduce MetaProbLog,<sup>15</sup> a prototype implementation for ProbLog supporting nested meta-calls based on exact backward inference. As they discuss, meta-calls can be supported by any inference mechanism that can be suspended to perform inference for the query inside the meta-call. Such suspending is natural in backward reasoning, where the proof of a subgoal becomes a call to inference rather than a continuation of backward reasoning. With forward reasoning, such non-ground `prob(goal,P)` goals raise the same issues as other non-ground facts. Meta-calls of the form `prob(goal,P)` compute the grounding of  $P$  as the goal's probability, and using approximate inference to compute the latter will thus influence the grounding of such a fact, and therefore potentially also the consequences of this fact. This may affect the result of inference in unexpected ways, and it is thus unclear in how far approximation approaches are suitable for meta-calls. Goodman et al (2008) state that supporting meta-calls (or nested

<sup>15</sup> <http://people.cs.kuleuven.be/~theoфраstos.mantadelis/tools/metaproblog.tar.gz>, also supports flexible probabilities, stochastic memoization, and negation as failure

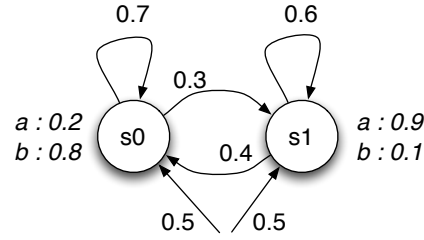


Fig. 2 Example HMM

queries) in MCMC inference in Church is expected to be straightforward, but do not provide details. AILog2, PRISM, ProbLog1, ProbLog2, cplint and PITA do not support nested meta-calls, i.e., querying for probabilities is only possible at the top level.

#### 4.9 Time and Dynamics

Among the most popular probabilistic models are those that deal with dynamics and time such as Hidden Markov Models (HMMs) and Dynamic Bayesian Networks. Dynamic models have received quite some attention within probabilistic logic programming. They can naturally be represented using logic programs through the addition of an extra "time" argument to each of the predicates. We illustrate this by giving two encodings of the Hidden Markov Model shown in Figure 2, where we restrict sequences to a given length (10 in the example). Following Vennekens et al (2004), this model can be written as a set of annotated disjunctions:

```

length(10).
0.5 :: state(s0, 0); 0.5 :: state(s1, 0).
0.7 :: state(s0, T1); 0.3 :: state(s1, T1) :- state(s0, T), length(L), T < L, T1 is T + 1.
0.4 :: state(s0, T1); 0.6 :: state(s1, T1) :- state(s1, T), length(L), T < L, T1 is T + 1.
0.2 :: out(a, T); 0.8 :: out(b, T) :- state(s0, T).
0.9 :: out(a, T); 0.1 :: out(b, T) :- state(s1, T).

```

Alternatively, following Sato and Kameya (1997), but writing PRISM's multi-valued switches as unconditional annotated disjunctions,<sup>16</sup> the model can be writ-

<sup>16</sup> In this example, the program structure causes the time argument to act as a unique identifier for different calls to the same AD, thus making memoized ADs and dememoized switches equivalent.



**Fig. 3** Example graph illustrating generalized labels: Boolean case (left), shortest path (right).

ten as follows:

```

0.2 :: output(s0, a, T) ; 0.8 :: output(s0, b, T).
0.9 :: output(s1, a, T) ; 0.1 :: output(s1, b, T).

0.5 :: init(s0) ; 0.5 :: init(s1).

0.7 :: trans(s0, s0, T) ; 0.3 :: trans(s0, s1, T).
0.4 :: trans(s1, s0, T) ; 0.6 :: trans(s1, s1, T).

length(10).

hmm(List) :- init(S), hmm(1, S, List).

% last time T :
hmm(T, S, [Obs]) :- length(T), output(S, Obs, T).

% earlier time T : output Obs in state S, transit from S to Next
hmm(T, S, [Obs|R]) :- length(L), T < L,
    output(S, Obs, T), trans(S, Next, T),
    T1 is T + 1, hmm(T1, Next, R).
    
```

Forward and backward sampling naturally deal with a time argument (provided time is bounded in the case of forward reasoning). Naively using such a time argument with exact inference results in exponential running times (in the number of time steps), though this can often be avoided using dynamic programming approaches and principles, as shown by the PRISM system, which achieves the same time complexity for HMMs as corresponding special-purpose algorithms (Sato and Kameya, 2001).

Other approaches that have devoted special attention to modeling and inference for dynamics include Logical HMMs (Kersting et al, 2006), a language for modeling HMMs with structured states, CPT-L (Thon et al, 2011), a dynamic version of CP-logic, and the work on a particle filter for dynamic distributional clauses (Nitti et al, 2013).

#### 4.10 Generalized Labels for Facts and Queries

As we have seen in Section 3, computing success probabilities in probabilistic logic programming is closely related to evaluating the truth value of a logical formula.

Weighted logic programming languages such as Dyna (Eisner et al, 2005)<sup>17</sup> and aProbLog (Kimmig et al, 2011a) take this observation a step further and replace probabilities (or Boolean truth values) by elements from a semiring and corresponding combination operators.<sup>18</sup>

More specifically, Dyna assigns labels to ground facts in a logic program and computes weights of atoms in the heads of clauses as follows: conjunction ( $\wedge$ ) in clause bodies is replaced by semiring multiplication  $\otimes$ , that is, the weight of a body is the  $\otimes$ -product of the weights of its atoms, and if multiple clauses share the same head atom, this atom's weight is the  $\oplus$ -sum of the corresponding bodies, that is,  $:-$  is replaced by semiring addition  $\oplus$ . We illustrate the idea with a logic program defining reachability in a directed graph adapted from Cohen et al (2008):

```
reachable(S) :- initial(S).
reachable(S) :- reachable(R), edge(R, S).
```

which in Dyna is interpreted as a system of (recursive) semiring equations

$$\begin{aligned} \text{reachable}(S) \oplus &= \text{initial}(S). \\ \text{reachable}(S) \oplus &= \text{reachable}(R) \otimes \text{edge}(R, S). \end{aligned}$$

To get the usual logic programming semantics, we can combine this program with facts labeled with values from the Boolean semiring (with  $\otimes = \wedge$  and  $\oplus = \vee$ ), as illustrated on the left of Figure 3:

```
initial(a) = T
edge(a, b) = T  edge(a, d) = T  edge(b, c) = T  edge(d, b) = T  edge(d, c) = T
```

which means that the weights of `reachable` atoms are computed as follows:

$$\begin{aligned} \text{reachable}(a) &= \text{initial}(a) = T \\ \text{reachable}(d) &= \text{reachable}(a) \wedge \text{edge}(a, d) = T \\ \text{reachable}(b) &= \text{reachable}(a) \wedge \text{edge}(a, b) \vee \text{reachable}(d) \wedge \text{edge}(d, b) = T \\ \text{reachable}(c) &= \text{reachable}(b) \wedge \text{edge}(b, c) \vee \text{reachable}(d) \wedge \text{edge}(d, c) = T \end{aligned}$$

Alternatively, one can label facts with non-negative numbers denoting costs, as illustrated on the right of Figure 3, and use  $\otimes = +$  and  $\oplus = \min$  to describe single-source shortest paths:

```
initial(a) = 0
edge(a, b) = 7  edge(a, d) = 5  edge(b, c) = 13  edge(d, b) = 4  edge(d, c) = 9
```

<sup>17</sup> Dyna is currently being extended into a more general language (Eisner and Filardo, 2011), but we consider the initial version here, as that one is more closely related to the probabilistic programming languages we discuss.

<sup>18</sup> A *semiring* is a structure  $(\mathcal{A}, \oplus, \otimes, e^\oplus, e^\otimes)$ , where *addition*  $\oplus$  is an associative and commutative binary operation over the set  $\mathcal{A}$ , *multiplication*  $\otimes$  is an associative binary operation over the set  $\mathcal{A}$ ,  $\otimes$  distributes over  $\oplus$ ,  $e^\oplus \in \mathcal{A}$  is the neutral element of  $\oplus$ , i.e., for all  $a \in \mathcal{A}$ ,  $a \oplus e^\oplus = a$ ,  $e^\otimes \in \mathcal{A}$  is the neutral element of  $\otimes$ , i.e., for all  $a \in \mathcal{A}$ ,  $a \otimes e^\otimes = a$ , and for all  $a \in \mathcal{A}$ ,  $e^\oplus \otimes a = a \otimes e^\oplus = e^\oplus$ . In a *commutative semiring*,  $\otimes$  is commutative as well.

resulting in evaluation

$$\text{reachable}(a) = \text{initial}(a) = 0$$

$$\text{reachable}(d) = \text{reachable}(a) + \text{edge}(a, d) = 5$$

$$\text{reachable}(b) = \min(\text{reachable}(a) + \text{edge}(a, b), \text{reachable}(d) + \text{edge}(d, b)) = 7$$

$$\text{reachable}(c) = \min(\text{reachable}(b) + \text{edge}(b, c), \text{reachable}(d) + \text{edge}(d, c)) = 14$$

That is, the values of `reachable` atoms now correspond to the length of the shortest path rather than the existence of a path.

Given its origins in natural language processing, Dyna is closely related to PRISM in two aspects. First, it does not memoize labeled facts, but takes into account their weights each time they appear in a derivation, generalizing how each use of a rule in a probabilistic grammar contributes to a derivation. Second, again as in probabilistic grammars, it sums the weights of all derivations, but in contrast to PRISM or grammars does not require them to be mutually exclusive to do so.

The inference algorithm of basic Dyna as given by Eisner et al (2005)<sup>19</sup> computes weights by forward reasoning, keeping intermediate results in an agenda and updating them until a fixpoint is reached, though other execution strategies could be used as well, cf. (Eisner and Filardo, 2011).

As Dyna, aProbLog (Kimmig et al, 2011a) replaces probabilistic facts by semiring-labeled facts, with the key difference that it bases the labels of derived facts on the labels of their models rather than those of their derivations. It thus directly generalizes the success probability (5) and the possible world DNF (6). As, in contrast to derivations, models do not provide an order in which labels of facts have to be multiplied, aProbLog requires semirings to be commutative. This restriction ensures that labels of derived facts are uniquely defined, and allows one to use inference approaches based on BDDs or sd-DNNFs, which may reorder facts when constructing the efficient representation. ProbLog inference algorithms based on BDDs have been directly adapted to aProbLog.<sup>20</sup>

Rather than replacing probabilities with semiring labels, one can also combine them with utilities or costs, and use the resulting language for decision making under uncertainty, as done in DTProbLog (Van den Broeck et al, 2010).<sup>21</sup>

## 5 Knowledge-Based Model Construction

So far, we have focused on probabilistic logic languages with strong roots in logic, where the key concepts of logic and probability are unified, that is, a random variable corresponds to a ground fact (or sometimes a ground term, as in distributional clauses), and standard logic programs are used to specify knowledge that can be derived from these facts. In this section, we discuss a second important group of probabilistic logic languages with strong roots in probabilistic graphical models, such as Bayesian or Markov networks. These formalisms typically use logic as a templating language for graphical models in relational domains,

<sup>19</sup> Implementation available at <http://dyna.org/>

<sup>20</sup> A prototype implementation of aProbLog is included in ProbLog1, cf. Footnote 6.

<sup>21</sup> An implementation of DTProbLog is included in ProbLog1 and ProbLog2, cf. Footnotes 6 and 7.

and thus take a quite different approach to combine logic and probabilities, also known as *knowledge-based model construction* (KBMC). Important representatives of this stream of research include probabilistic logic programs (PLPs) (Haddawy, 1994), relational Bayesian networks (RBNs) (Jaeger, 1997), probabilistic relational models (PRMs) (Koller and Pfeffer, 1998; Getoor et al, 2007), Bayesian logic programs (BLPs) (Kersting and De Raedt, 2001, 2008),  $\text{CLP}(\mathcal{BN})$  (Santos Costa et al, 2003, 2008), logical Bayesian networks (LBNs) (Fierens et al, 2005), Markov Logic (Richardson and Domingos, 2006), chain logic (Hommersom et al, 2009), and probabilistic soft logic (PSL) (Bröcheler et al, 2010). A recent survey of this field is provided by Kimmig et al (2014).

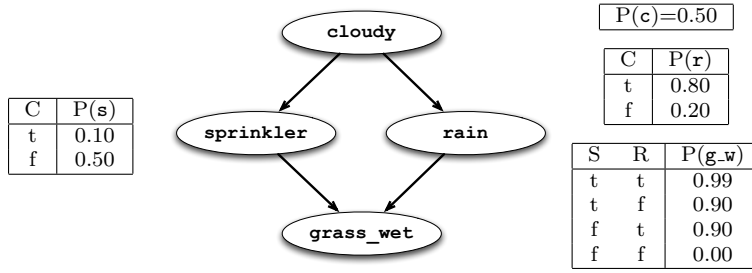
In the following, we relate the key concepts underlying the knowledge-based model construction approach to those discussed in the rest of this article. We again focus on languages based on logic programming, such as PLPs, BLPs, LBNs, chain logic, and  $\text{CLP}(\mathcal{BN})$ , but mostly abstract from the specific language. These representation languages are typically designed so that implication in logic (“:–”) corresponds to the direct influence relation in Bayesian networks. The logical knowledge base is then used to construct a Bayesian network. So inference proceeds in two steps: the logical step, in which one constructs the network, and the probabilistic step, in which one performs probabilistic inference on the resulting network. We first discuss modeling Bayesian networks and their relational counterpart in the context of the distribution semantics, and then focus on  $\text{CLP}(\mathcal{BN})$  as an example of a KBMC approach whose primitives clearly expose the separation between model construction via logic programming and probabilistic inference on the propositional model.

### 5.1 Bayesian Networks and Conditional Probability Tables

A Bayesian network (BN) defines a joint probability distribution over a set of random variables  $\mathcal{V} = \{V_1, \dots, V_m\}$  by factoring it into a product of conditional probability distributions, one for each variable  $V_i$  given its parents  $\text{par}(V_i) \subseteq \mathcal{V}$ . The parent relation is given by an acyclic directed graph (cf. Figure 4), where the random variables are the nodes and an edge  $V_i \rightarrow V_j$  indicates that  $V_i$  is a parent of  $V_j$ . The conditional probability distributions are typically specified as *conditional probability tables* (CPTs), which form the key probabilistic concept of BNs. For instance, the CPT on the left of Figure 4 specifies that the random variable `sprinkler` takes value true with probability 0.1 (and false with 0.9) if its parent `cloudy` is true, and with probability 0.5 if `cloudy` is false. Formally, a CPT contains a row for each possible assignment  $x_1, \dots, x_n$  to the parent variables  $X_1, \dots, X_n$  specifying the distribution  $P(X|x_1, \dots, x_n)$ . As has been shown earlier, e.g., by Poole (1993) and Vennekens et al (2004), any Bayesian network can be modeled in languages based on the distribution semantics by representing every row in a CPT as an annotated disjunction

$$p_1 :: X(w_1); \dots ; p_k :: X(w_k) :- X_1(v_1), \dots, X_n(v_n)$$

where  $X(v)$  is true when  $v$  is the value of  $X$ . The body of this AD is true if the parent nodes have the values specified in the corresponding row of the CPT, in which case the AD chooses a value for the child from the corresponding distribution. As an example, consider the sprinkler network shown in Figure 4. The CPT



**Fig. 4** The sprinkler network is a Bayesian network modeling an environment where both the sprinkler and the rain can cause the grass getting wet (Russell and Norvig, 2003).

for the root node `cloudy` corresponds to an AD with empty body

$$0.5 :: \text{cloudy}(\text{t}); 0.5 :: \text{cloudy}(\text{f}).$$

whereas the CPTs for `sprinkler` and `rain` require the state of their parent node `cloudy` to be present in the body of the ADs

$$\begin{aligned} 0.1 :: \text{sprinkler}(\text{t}); 0.9 :: \text{sprinkler}(\text{f}) :- \text{cloudy}(\text{t}). \\ 0.5 :: \text{sprinkler}(\text{t}); 0.5 :: \text{sprinkler}(\text{f}) :- \text{cloudy}(\text{f}). \\ 0.8 :: \text{rain}(\text{t}); 0.2 :: \text{rain}(\text{f}) :- \text{cloudy}(\text{t}). \\ 0.2 :: \text{rain}(\text{t}); 0.8 :: \text{rain}(\text{f}) :- \text{cloudy}(\text{f}). \end{aligned}$$

The translation for the CPT of `grass_wet` is analogous.

## 5.2 Relational Dependencies

Statistical relational learning formalisms such as BLPs, PLPs, LBNs and  $\text{CLP}(\mathcal{BN})$  essentially replace the specific random variables in the CPTs of Bayesian networks by logically defined random variable templates, commonly referred to as *parameterized* random variables or par-RVs for short (Poole, 2003), though the actual syntax amongst these systems differs significantly. We here use annotated disjunctions to illustrate the key idea. For instance, in a propositional setting, the following annotated disjunctions express that a specific student's grade in a specific course probabilistically depends on whether he has read the corresponding textbook or not:

$$\begin{aligned} 0.6 :: \text{grade}(\text{high}); 0.4 :: \text{grade}(\text{low}) :- \text{reads}(\text{true}). \\ 0.1 :: \text{grade}(\text{high}); 0.9 :: \text{grade}(\text{low}) :- \text{reads}(\text{false}). \end{aligned}$$

Using logical variables, this dependency can directly be expressed for many students, courses, and books:

$$\begin{aligned} 0.6 :: \text{grade}(\text{S}, \text{C}, \text{high}); 0.4 :: \text{grade}(\text{S}, \text{C}, \text{low}) :- \text{book}(\text{C}, \text{B}), \text{reads}(\text{S}, \text{B}). \\ 0.1 :: \text{grade}(\text{S}, \text{C}, \text{high}); 0.9 :: \text{grade}(\text{S}, \text{C}, \text{low}) :- \text{book}(\text{C}, \text{B}), \text{not}(\text{reads}(\text{S}, \text{B})). \end{aligned}$$



More concretely, the annotated disjunctions express that  $P(\text{grade}(S, C) = \text{high}) = 0.6$  if the student has read the book of the course and  $P(\text{grade}(S, C) = \text{high}) = 0.1$  otherwise. Thus the predicate `grade` depends on `book/2` and `reads/2`. The dependency holds for all instantiations of the rule, that is, it acts as a template for all persons, courses, and books. This is what knowledge-based model construction approaches all share: the logic acts as a template to generate dependencies (here CPTs) in the graphical model. This also introduces a complication that is not encountered in propositional Bayesian networks or their translation to annotated disjunctions. To illustrate this, let us assume the predicate `book/2` is deterministic and known. Then the propositional case arises when for each course there is exactly one book. The annotated disjunctions then effectively encode the conditional probability table  $P(\text{Grade}|\text{Reads})$ . However, if there are multiple books, say two, for one course, then the above template would specify two CPTs: one for the first book,  $P(\text{Grade}|\text{Reads}_1)$ , and one for the second,  $P(\text{Grade}|\text{Reads}_2)$ . In Bayesian networks, these CPTs need to be combined and there are essentially two ways for realizing this.

The first is to use a so-called *combining rule*, that is, a function that maps these CPTs into a single CPT of the form  $P(\text{Grade}|\text{Reads}_1, \text{Reads}_2)$ . The most popular combining rule is noisy-or, for which

$$P(\text{Grade} = \text{high}|\text{Reads}_1, \dots, \text{Reads}_n) = 1 - \prod_{i=1}^n (1 - P(\text{Grade} = \text{high}|\text{Reads}_i = \text{true}))$$

where  $n$  is the number of books for the course. Using annotated disjunctions, this combining rule is obtained automatically, cf. Section 2.2. In the statistical relational learning literature, this approach is followed for instance in RBNs and BLPs, and several other combining rules exist, cf., e.g., (Jaeger, 1997; Kersting and De Raedt, 2008; Natarajan et al, 2005). While combining rules are an important concept in KBMC, using them in their general form under the distribution semantics requires one to change the underlying logic, which is non-trivial. Hommersom and Lucas (2011) introduce an approach that models these interactions by combining the distribution semantics with default logic. Alternatively, one could use meta-calls, cf. Section 4.8.

The second way of dealing with the two distributions uses aggregation. In this way, the random variable upon which one conditions grade is the number of books the person read, rather than the reading of the individual books. This approach is taken for instance in PRMs and  $\text{CLP}(\mathcal{BN})$ . In the context of the distribution semantics, aggregation can be realized within the logic program using second order predicates, cf. Section 4.7. For instance, the following program makes a distinction between reading more than two, two, one, or none of the books:

```
0.9 :: grade(S, C, high); 0.1 :: grade(S, C, low) :- nofbooksread(S, C, N), N > 2.
0.8 :: grade(S, C, high); 0.2 :: grade(S, C, low) :- nofbooksread(S, C, 2).
0.6 :: grade(S, C, high); 0.4 :: grade(S, C, low) :- nofbooksread(S, C, 1).
0.1 :: grade(S, C, high); 0.9 :: grade(S, C, low) :- nofbooksread(S, C, 0).
nofbooksread(S, C, N) :- findall(B, (book(C, B), reads(S, B)), List), length(List, N).
```

### 5.3 Example: $\text{CLP}(\mathcal{BN})$

An example of a KBMC approach that clearly exposes the separation between model construction and probabilistic inference in the resulting model is  $\text{CLP}(\mathcal{BN})$  (Santos Costa et al, 2008), which we now discuss in more detail.<sup>22</sup>  $\text{CLP}(\mathcal{BN})$  uses constraint programming principles to construct Bayesian networks. The key inference task in  $\text{CLP}(\mathcal{BN})$  is to compute marginal distributions of query variables, conditioned on evidence if available. Syntactically,  $\text{CLP}(\mathcal{BN})$  extends logic programming with *constraint atoms* that (a) define random variables together with their CPTs and (b) establish constraints linking these random variables to logical variables used in the logic program.

The first phase of inference in  $\text{CLP}(\mathcal{BN})$  uses backward reasoning in the logic program to collect all relevant constraints in a constraint store. These constraints define the relevant Bayesian network, on which the second phase computes the required marginals. Conditioning on evidence is straightforward, as it only requires to add the corresponding constraints to the store.<sup>23</sup>

We first illustrate this for the propositional case, using the following model<sup>24</sup> of the sprinkler Bayesian network as given in Figure 4:<sup>25</sup>

```
cloudy(C) :-
    { C = cloudy with p([f,t], [0.5,0.5], []) }.

sprinkler(S) :-
    cloudy(C),
    { S = sprinkler with p([f,t], [0.5,0.9,
                                0.5,0.1],
                          [C])
      % C = f , t
      % S = f
      % S = t
    }.

rain(R) :-
    cloudy(C),
    { R = rain with p([f,t], [0.8,0.2,
                            0.2,0.8],
                      [C])
      % C = f , t
      % R = f
      % R = t
    }.

wet_grass(W) :-
    sprinkler(S),
    rain(R),
    { W = wet with p([f,t],
                    /* S/R = f/f, f/t, t/f, t/t */
                    [1.0, 0.1, 0.1, 0.01,
                     0.0, 0.9, 0.9, 0.99],
                    [C])
      % W = f
      % W = t
    }.
```

<sup>22</sup> implementation included in YAP Prolog, <http://www.dcc.fc.up.pt/~vsc/Yap/>

<sup>23</sup> The implementation adds evidence declared in the input program to the store at compile time.

<sup>24</sup> taken from the examples in the  $\text{CLP}(\mathcal{BN})$  system

<sup>25</sup> We include comments for better readability, as  $\text{CLP}(\mathcal{BN})$  swaps rows and columns of CPTs compared to the notation in Figure 4.

```

    [S,R])
  }.

```

In the clause for the top node `cloudy`, the body consists of a single constraint atom (delimited by curly braces) that constrains the logical variable `C` to the value of the random variable `cloudy`. The term `p([f,t],[0.5,0.5],[ ])` specifies that this random variable takes values `f` or `t` with probability 0.5 each, and has an empty parent list. Note that within constraint atoms, the `=` sign does not denote Prolog unification, but an equality constraint between a logical variable and the value of a random variable. The clause for `sprinkler` first calls `cloudy(C)`, which as discussed sets up a constraint between `C` and the `cloudy` random variable, and then uses `C` as the only parent of the random variable `sprinkler` it defines. The first column of the CPT corresponds to the first parent value, the first row to the first child value, and so on, i.e., in case of `cloudy=f`, the probability of `sprinkler=f` is 0.5, whereas for `cloudy=t`, it is 0.9. The remaining two random variables `rain` and `wet` are defined analogously, with their clauses again first calling the predicates for the parent variables to include the corresponding constraints. To answer the query `sprinkler(S)`, which asks for the marginal of the random variable `sprinkler`,  $CLP(\mathcal{BN})$  performs backward reasoning to find all constraints in the proof of the query, and thus the part of the Bayesian network relevant to compute the marginal. This first calls `cloudy(C)`, adding the constraint `C=cloudy` to the store (and thus the `cloudy` node to the BN), and then adds the constraint `S=sprinkler` to the store, and the `sprinkler` node with parent `cloudy` to the BN. Any BN inference algorithm can be used to compute the marginal in the second phase.

In general, a  $CLP(\mathcal{BN})$  clause (in canonical form) is either a standard Prolog clause, or has the following structure:

$$h(A_1, \dots, A_n, V) :- \text{body}, \{V = \text{sk}(C_1, \dots, C_t) \text{ with CPT}\}.$$

Here, `body` is a possibly empty conjunction of logical atoms, and the part in curly braces is a constraint atom. The term `sk(C1, ..., Ct)` is a Skolem term not occurring in any other clause of the program. Its arguments `Ci` are given via the input variables `Aj` and the logical `body`. `CPT` is a term of the form `p(Values, Table, Parents)`, where `Values` is a list of possible values for `sk(C1, ..., Ct)`, `Parents` is a list of logical variables specifying the parent nodes, and `Table` the probability table given as a list of probabilities. The order of entries in this list corresponds to the valuations obtained by backtracking over the parents' values in the order given in the corresponding definitions. This CPT term can be given either directly (as in the example above) or via the use of logical variables and unification (see below).

When defining relational models, random variables can be parameterized by logical variables as in the following clause from the school example included in the implementation:

```

registration_grade(R, Grade) :-
  registration(R, C, S),
  course_difficulty(C, Dif),
  student_intelligence(S, Int),
  grade_table(Int, Dif, Table),
  { Grade = grade(R) with Table }.

```

```

grade_table(I, D,
  p([a,b,c,d],
/* I,D = h h h m h l m h m m m l l h l m l l */
  [ 0.20, 0.70, 0.85, 0.10, 0.20, 0.50, 0.01, 0.05, 0.10, %a
    0.60, 0.25, 0.12, 0.30, 0.60, 0.35, 0.04, 0.15, 0.40, %b
    0.15, 0.04, 0.02, 0.40, 0.15, 0.12, 0.50, 0.60, 0.40, %c
    0.05, 0.01, 0.01, 0.20, 0.05, 0.03, 0.45, 0.20, 0.10 ],%d
  [I,D])).

```

Here, `registration/3` is a purely logical predicate linking a registration `R` to a course `C` and a student `S`. We omit the clauses for `course.difficulty` and `student.intelligence`; these define distributions over possible values `h(igh)`, `m(edium)`, and `l(ow)` for the difficulty `Dif` of course `C` and the intelligence `Int` of student `S`, respectively. For each grounding `r` of the variable `R` in the database of registrations, the clause above defines a random variable `grade(r)` with values `a`, `b`, `c` and `d` that depends on the difficulty of the corresponding course and the intelligence of the corresponding student. In this case, the CPT itself is not defined within the constraint atom, but obtained from the Prolog predicate `grade_table` via unification.

Defining aggregation using second order predicates is straightforward in  $CLP(\mathcal{BN})$ , as random variables and constraints are part of the object level vocabulary. For instance, the following clause defines the performance level of a student based on the average of his grades:

```

student_level(S,L) :-
  findall(G,(registration(R,_,S),registration_grade(R,G)),Grades),
  avg_grade(Grades,Avg),
  level_table(T),
  { L = level(S) with p([h,m,l],T,[Avg])}.

```

First, the list `Grades` of all grade random variables for student `S` is obtained using the Prolog predicate `findall`. Then, `avg_grade/2` constrains `Avg` to a new random variable defined as the average of these grades (with a deterministic CPT). Finally, the CPT specifying how the performance level depends on this average is obtained from the deterministic predicate `level_table`, and the corresponding random variable and constraint are set up in the constraint atom. We refer to Santos Costa et al (2008) for a discussion of the inference challenges aggregates raise.

Despite the differences in syntax, probabilistic primitives, and inference between  $CLP(\mathcal{BN})$  and probabilistic extensions of Prolog following the distribution semantics, there are also many commonalities between those. As we discussed above, conditional probability tables can be represented using annotated disjunctions, and it is thus possible to transform  $CLP(\mathcal{BN})$  clauses into Prolog programs using annotated disjunctions. On the other hand, Santos Costa and Paes (2009) discuss the relation between PRISM and  $CLP(\mathcal{BN})$  based on a number of PRISM programs that they map into  $CLP(\mathcal{BN})$  programs.

## 6 Probabilistic Programming Concepts and Inference

We complete this survey by summarizing the relations between the dimensions of SUCC inference as discussed in Section 3 and the probabilistic programming con-

|                   | Flexible Probabilities | Continuous Distributions | Stochastic Memoization     | Negation as Failure | 2nd Order Predicates | Meta-Calls | Time and Dynamics | Generalized Labels (aProbLog) |
|-------------------|------------------------|--------------------------|----------------------------|---------------------|----------------------|------------|-------------------|-------------------------------|
| forward exact     | no                     | no                       | with                       | yes                 | yes*                 | no         | yes*              | yes                           |
| backward exact    | yes                    | limited                  | with<br>or<br>with-<br>out | yes                 | yes*                 | yes        | yes*              | yes                           |
| forward sampling  | no                     | yes                      | with                       | yes                 | yes                  | no         | yes               | n.a.                          |
| backward sampling | yes                    | yes                      | with<br>or<br>with-<br>out | yes                 | yes                  | yes        | yes               | n.a.                          |

**Table 2** Relation between key probabilistic programming concepts and main dimensions of inference; see Section 6 for details. (\* number of proofs/worlds exponential in length of answer list or time sequence)

cepts identified in Section 4. On the probabilistic side, we focus on *exact inference* versus *sampling*, as conclusions for exact inference carry over to approximate inference with bounds in most cases. On the logical side, we focus on *forward* versus *backward* reasoning, as conclusions for backward reasoning carry over to the approach using weighted model counting. We provide an overview in Table 2, where we omit the concepts *unknown objects*, as those are typically simulated via flexible probabilities and/or continuous distributions, and *constraints*, as those have not yet been considered during inference. For *generalized labels*, we focus on aProbLog, as it is closer to the distribution semantics than Dyna, due to its semantics based on worlds rather than derivations. We do not include MCMC here, as existing MCMC approaches in the context of the distribution semantics are limited to the basic case of definite clause programs without additional concepts.

*Dimensions of inference:* The main difference between exact inference and sampling is that the former has to consider *all* possible worlds or all proofs of the query, whereas the latter always considers one possible world or proof in isolation. As *second order predicates* and *time and dynamics* can increase the number of proofs exponentially (in the length of the answer list or the number of time steps), they are more easily handled by sampling-based approaches, though tabling can significantly improve performance of exact inference in dynamic domains. Sampling based approaches do not directly apply for *generalized labels*, as sampling exploits the probabilistic semantics of fact labels.

The main difference between forward and backward reasoning is that the former generates all consequences of the probabilistic logic program, whereas the latter is query-driven and only considers relevant consequences, which can drastically improve efficiency. This difference is well-known in logic programming, and becomes even more important in the probabilistic setting, where we are interested in not just a single world or proof, but in all possible worlds or all proofs. The fact

that backward reasoning is query-driven makes it well-suited for *flexible probabilities* and *meta-calls*, which cannot directly be handled in forward reasoning. The reason is that the corresponding subgoals have an infinite number of groundings, among which backward reasoning easily picks the relevant ones, which forward reasoning cannot do. The same effect makes it necessary to use *stochastic memoization* in forward reasoning, while backward reasoning can support dememoization (as in PRISM) as well as memoization (as in the various ICL, ProbLog and LPAD systems).

The roots of the distribution semantics in logic programming become apparent when considering inference for the two remaining key concepts, *negation as failure* and continuous distributions as provided by *distributional clauses*. While the logic concept of negation as failure is naturally supported in all combinations of exact inference or sampling and forward or backward reasoning, the probabilistic concept of continuous distributions is much more challenging, and only practical in sampling-based approaches.

*Inference approaches:* More specifically, *exact inference using forward reasoning* in the form discussed in Section 3.1 can be used for all programs with finitely many finite worlds, which (a) excludes the use of non-ground facts without explicitly given domains, flexible probabilities, meta-calls and continuous probabilities, and (b) requires stochastic memoization. As this approach additionally suffers from having to enumerate all possible worlds, it is not used in practice.<sup>26</sup>

*Exact inference using backward reasoning* is the most widely supported inference technique in probabilistic logic programming, provided by AILog2, PRISM, ProbLog1, cplint, PITA and MetaProbLog. PRISM never uses stochastic memoization, whereas the other systems always use it. Only very limited forms of continuous distributions can be supported, cf. the work on Hybrid ProbLog (Gutmann et al, 2010). All other concepts can be supported, but implementations differ in the ones they cover. Negation as failure is supported in all implementations. In addition, AILog2 and cplint support flexible probabilities, MetaProbLog supports flexible probabilities and meta-calls, and ProbLog1 supports flexible probabilities, limited use of continuous distributions (Hybrid ProbLog) and generalized labels (aProbLog). Approximate inference with bounds using backward reasoning is available in ProbLog1 and cplint, but restricted to definite clause programs, as the use of negation as failure complicates proof finding (as discussed in Section 4.6). As the WMC approach as implemented in ProbLog2 uses backward inference to determine the relevant grounding, that is, the groundings of clauses that appear in some proof of a query, the same observations as for exact backward inference apply in this case as well. ProbLog2 supports flexible probabilities and negation as failure.

*Forward sampling* in its simplest form as discussed in Section 3.1 can be used with programs whose worlds are all finite, which excludes the use of non-ground facts without explicitly given domains, flexible probabilities, and meta-calls, and requires stochastic memoization. In contrast to exact forward inference, forward sampling does support continuous distributions, as only one value is considered

<sup>26</sup> Dyna’s exact inference is based on forward reasoning, but uses a different type of algorithm that propagates value updates using forward reasoning based on an agenda of pending updates.

at a time. None of the probabilistic logic programming systems discussed here implement forward sampling.

*Backward sampling* is the most flexible approach and can in principle deal with all concepts except generalized labels. Backward sampling approaches are provided by ProbLog1 and cplint, which both support flexible probabilities and negation as failure. PRISM has a builtin for sampling the outcome of a query using backward reasoning, but does not use it for probability estimation.

## 7 Applications

Probabilistic programming has been applied in many domains, including natural language processing, bioinformatics, action- and activity recognition, robotics, web analysis, and music analysis. We now present a selection of such applications, again focusing on probabilistic logic programming languages; examples using other paradigms include the works of Mansinghka et al (2013), Gerstenberg and Goodman (2012), and Stuhlmüller et al (2010).

Several probabilistic logic languages, including PRISM, SLPs and Dyna, have their roots in natural language processing and sequence analysis; they directly upgrade representations and techniques from probabilistic grammars or were designed with natural language applications in mind; cf. (Sato et al, 2008; Eisner et al, 2005). Sneyers et al (2010) developed CHR(PRISM), a constraint-based version of PRISM that has been applied to music analysis.

By far the most popular and successful application area of probabilistic logic programming is bioinformatics. For instance, PRISM has been used for biological sequence analysis (Mørk and Holmes, 2012) in the LOST project, ProbLog for interpreting gene interaction networks (De Maeyer et al, 2013), Stochastic Logic Programs for inferring metabolic pathways (Lodhi and Muggleton, 2005), and a combination of CP-logic and Allen’s interval calculus for modeling progress of diseases (van der Heijden and Lucas, 2013).

Another key application area of probabilistic programs are action- and activity recognition and planning. For instance, Thon et al (2011) build CPT-L models of Travian, a popular multi-player online strategy game, whereas Skarlatidis et al (2015) apply their probabilistic event calculus based on ProbLog to event recognition in surveillance videos. In the closely related area of robotics, distributional clauses have been applied to, for instance, tracking (Nitti et al, 2013).

Probabilistic logic programming is also strongly connected to (probabilistic) databases (Suciu et al, 2011) and reasoning over the web. For instance, Fuhr (1995) introduced probabilistic Datalog in the context of information retrieval and the web. A variation called WHIRL was developed by Cohen (2000) as an information representation language that synergistically combines properties of logic-based and text-based representation systems. More recently, the work on Never-Ending Language Learning (Nell) (Carlson et al, 2010) associates probabilities or confidences with instances of relations automatically extracted from the web, and the probabilistic logic programming system ProPPR (Wang et al, 2013) has been designed for reasoning and learning with the resulting large, uncertain knowledge bases.

## 8 Conclusions and Open Questions

Probabilistic programming is a rapidly developing field of research as witnessed by the many probabilistic programming languages and primitives that have been introduced over the past few years. In this paper, we have attempted to provide a gentle introduction to this field by focussing on probabilistic *logic* programming languages and identifying the underlying probabilistic concepts that these languages support. The same concept (e.g., probabilistic choice) can be realized using different syntactic primitives (e.g., switches, annotated disjunctions, etc.) leading to differences in the probabilistic programming languages. Probabilistic programming implementations not only differ in the primitives they provide but also in the way they perform probabilistic inference. Inference is a central concern in these languages, as probabilistic inference is computationally expensive. We have therefore also presented various probabilistic inference mechanisms and discussed their suitability for supporting the probabilistic programming concepts. This in turn allowed us to position different languages and implementations, leading to a broad survey of the state-of-the-art in probabilistic logic programming.

This work also reveals a number of challenges and directions for future work. First, there is an ongoing quest for efficient inference approaches for languages that support a broad range of programming concepts. Promising directions include lifted inference, which aims at exploiting symmetries and abstraction over individuals to speed up inference, knowledge compilation, which has contributed many data structures for compactly representing and efficiently querying various types of knowledge, and approximate methods such as MCMC, which is used in many probabilistic programming languages, but still requires proposal functions to be custom made for the program at hand. There also is a need for a clear understanding of the relative computational complexity of the various probabilistic languages and concepts that exist to date. Another question that has only seen partial answers so far is how to efficiently deal with evidence and constraints in different inference techniques. Adapting and extending program transformation and analysis techniques to the probabilistic setting promises opportunities to recognize and exploit program parts that are amenable to more efficient inference. Concepts such as time and dynamics require inference approaches that on the one hand exploit repeated structure, but on the other hand can also deal with changing structure over time. Last but not least, it still is a challenge to learn probabilistic programs, although a wide variety of learning techniques for probabilistic programming has already been developed. Many key challenges for both parameter and structure learning remain, many of which are related to efficient inference, as learning requires inference.

### Acknowledgements

The authors are indebted to Bernd Gutmann and Ingo Thon for participating in many discussions, and contributing several ideas during the early stages of the research that finally led to this paper. Angelika Kimmig is supported by the Flemish Research Foundation (FWO-Vlaanderen).



## References

- Angelopoulos N, Cussens J (2004) On the implementation of MCMC proposals over stochastic logic programs. In: Carro M, Morales JF (eds) Proceedings of the Colloquium on Implementation of Constraint and Logic Programming Systems (CICLOPS-04)
- Arora NS, de Salvo Braz R, Sudderth EB, Russell SJ (2010) Gibbs sampling in open-universe stochastic languages. In: Grünwald P, Spirtes P (eds) Proceedings of the 26th Conference on Uncertainty in Artificial Intelligence (UAI-10), AUA Press, pp 30–39
- Bancilhon F, Maier D, Sagiv Y, Ullman JD (1986) Magic sets and other strange ways to implement logic programs (extended abstract). In: Silberschatz A (ed) Proceedings of the 5th ACM SIGACT-SIGMOD symposium on Principles of Database Systems (PODS-86), ACM, pp 1–15
- Baral C, Gelfond M, Rushton JN (2004) Probabilistic reasoning with answer sets. In: Lifschitz V, Niemelä I (eds) Proceedings of the 7th International Conference on Logic Programming and Nonmonotonic Reasoning (LPNMR-04), Springer, Lecture Notes in Computer Science, vol 2923, pp 21–33
- Baral C, Gelfond M, Rushton N (2009) Probabilistic reasoning with answer sets. *Theory and Practice of Logic Programming (TPLP)* 9(1):57–144
- Bellodi E, Lamma E, Riguzzi F, Santos Costa V, Zese R (2014) Lifted variable elimination for probabilistic logic programming. *Theory and Practice of Logic Programming (TPLP)* 14(4–5):681–695
- Bragaglia S, Riguzzi F (2010) Approximate inference for logic programs with annotated disjunctions. In: Frasconi P, Lisi FA (eds) Revised papers of the 20th International Conference on Inductive Logic Programming (ILP-10), Springer, Lecture Notes in Computer Science, vol 6489, pp 30–37
- Bröcheler M, Mihalkova L, Getoor L (2010) Probabilistic similarity logic. In: Grünwald P, Spirtes P (eds) Proceedings of the 26th Conference on Uncertainty in Artificial Intelligence (UAI-10), AUA Press, pp 73–82
- Carlson A, Betteridge J, Kisiel B, Settles B, Hruschka Jr ER, Mitchell TM (2010) Toward an architecture for never-ending language learning. In: Fox M, Poole D (eds) Proceedings of the 24th AAAI Conference on Artificial Intelligence (AAAI-10), AAAI Press, pp 1306–1313
- Cohen SB, Simmons RJ, Smith NA (2008) Dynamic programming algorithms as products of weighted logic programs. In: Garcia de la Banda M, Pontelli E (eds) Proceedings of the 24th International Conference on Logic Programming (ICLP-08), Springer, Lecture Notes in Computer Science, vol 5366, pp 114–129
- Cohen WW (2000) WHIRL: A word-based information representation language. *Artificial Intelligence* 118(1):163–196
- Cussens J (2005) Integrating by separating: Combining probability and logic with ICL, PRISM and SLPs. APRIL II project report
- Dalvi NN, Suciu D (2004) Efficient query evaluation on probabilistic databases. In: Nascimento MA, Özsu MT, Kossmann D, Miller RJ, Blakeley JA, Schiefer KB (eds) Proceedings of the 30th International Conference on Very Large Databases (VLDB-04), Morgan Kaufmann Publishers, pp 864–875
- Dantsin E (1991) Probabilistic logic programs and their semantics. In: Voronkov A (ed) Proceedings of the First Russian Conference on Logic Programming, Springer, Lecture Notes in Computer Science, vol 592, pp 152–164

- De Maeyer D, Renkens J, Cloots L, De Raedt L, Marchal K (2013) Phenetic: Network-based interpretation of unstructured gene lists in *E. coli*. *Molecular BioSystems* 9(7):1594–1603
- De Raedt L, Kersting K (2003) Probabilistic logic learning. *SIGKDD Explorations* 5(1):31–48
- De Raedt L, Kimmig A, Toivonen H (2007) ProbLog: A probabilistic Prolog and its application in link discovery. In: Veloso MM (ed) *Proceedings of the 20th International Joint Conference on Artificial Intelligence (IJCAI-07)*, Morgan Kaufmann Publishers, pp 2462–2467
- De Raedt L, Frasconi P, Kersting K, Muggleton S (eds) (2008) *Probabilistic Inductive Logic Programming — Theory and Applications*, *Lecture Notes in Artificial Intelligence*, vol 4911. Springer
- Eisner J, Filardo NW (2011) Dyna: Extending Datalog for modern AI. In: de Moor O, Gottlob G, Furche T, Sellers A (eds) *Datalog Reloaded*, *Lecture Notes in Computer Science*, vol 6702, Springer, pp 181–220
- Eisner J, Goldlust E, Smith N (2005) Compiling Comp Ling: Weighted dynamic programming and the Dyna language. In: Brew C, Chien LF, Kirchhoff K (eds) *Proceedings of the Human Language Technology Conference and Conference on Empirical Methods in Natural Language Processing (HLT/EMNLP-05)*, The Association for Computational Linguistics, pp 281–290
- Fierens D, Blockeel H, Bruynooghe M, Ramon J (2005) Logical Bayesian networks and their relation to other probabilistic logical models. In: Kramer S, Pfahringer B (eds) *Proceedings of the 15th International Conference on Inductive Logic Programming (ILP-05)*, Springer, *Lecture Notes in Computer Science*, vol 3625, pp 121–135
- Fierens D, Van den Broeck G, Thon I, Gutmann B, De Raedt L (2011) Inference in probabilistic logic programs using weighted CNFs. In: Cozman FG, Pfeffer A (eds) *Proceedings of the 27th Conference on Uncertainty in Artificial Intelligence (UAI-11)*, AUAI Press, pp 211–220
- Fierens D, Van den Broeck G, Bruynooghe M, De Raedt L (2012) Constraints for probabilistic logic programming. In: Mansinghka V, Roy D, Goodman N (eds) *Proceedings of the NIPS Probabilistic Programming Workshop*
- Fierens D, Van den Broeck G, Renkens J, Shterionov D, Gutmann B, Thon I, Janssens G, De Raedt L (2014) Inference and learning in probabilistic logic programs using weighted Boolean formulas. *Theory and Practice of Logic Programming (TPLP)* FirstView
- Flach PA (1994) *Simply Logical: Intelligent Reasoning by Example*. John Wiley
- Fuhr N (1995) Probabilistic Datalog—a logic for powerful retrieval methods. In: Fox EA, Ingwersen P, Fidel R (eds) *Proceedings of the 18th Annual International ACM SIGIR Conference on Research and Development in Information Retrieval (SIGIR-95)*, ACM Press, pp 282–290
- Fuhr N (2000) Probabilistic Datalog: Implementing logical information retrieval for advanced applications. *Journal of the American Society for Information Science (JASIS)* 51(2):95–110
- Gerstenberg T, Goodman ND (2012) Ping pong in Church: Productive use of concepts in human probabilistic inference. In: Miyake N, Peebles D, Cooper RP (eds) *Proceedings of the 34th Annual Conference of the Cognitive Science Society*, Cognitive Science Society, pp 1590–1595

- Getoor L, Taskar B (eds) (2007) *An Introduction to Statistical Relational Learning*. MIT Press
- Getoor L, Friedman N, Koller D, Pfeffer A, Taskar B (2007) Probabilistic relational models. In: Getoor and Taskar (2007), pp 129–174
- Goodman ND, Mansinghka VK, Roy DM, Bonawitz K, Tenenbaum JB (2008) Church: A language for generative models. In: McAllester DA, Myllymäki P (eds) *Proceedings of the 24th Conference on Uncertainty in Artificial Intelligence (UAI-08)*, AUAI Press, pp 220–229
- Gutmann B, Jaeger M, De Raedt L (2010) Extending ProbLog with continuous distributions. In: Frasconi P, Lisi FA (eds) *Proceedings of the 20th International Conference on Inductive Logic Programming (ILP-10)*, Springer, *Lecture Notes in Computer Science*, vol 6489, pp 76–91
- Gutmann B, Thon I, Kimmig A, Bruynooghe M, De Raedt L (2011) The magic of logical inference in probabilistic programming. *Theory and Practice of Logic Programming (TPLP)* 11(4–5):663–680
- Haddawy P (1994) Generating Bayesian networks from probabilistic logic knowledge bases. In: de Mántaras RL, Poole D (eds) *Proceedings of the 10th Annual Conference on Uncertainty in Artificial Intelligence (UAI-94)*, Morgan Kaufmann Publishers, pp 262–269
- van der Heijden M, Lucas PJF (2013) Describing disease processes using a probabilistic logic of qualitative time. *Artificial Intelligence in Medicine* 59(3):143–155
- Hommersom A, Lucas PJF (2011) Generalising the interaction rules in probabilistic logic. In: Walsh T (ed) *Proceedings of the 22nd International Joint Conference on Artificial Intelligence (IJCAI-11)*, AAAI Press, pp 912–917
- Hommersom A, de Carvalho Ferreira N, Lucas PJF (2009) Integrating logical reasoning and probabilistic chain graphs. In: Buntine WL, Grobelnik M, Mladenic D, Shawe-Taylor J (eds) *Proceedings of the European Conference on Machine Learning and Principles and Practice of Knowledge Discovery in Databases (ECML/PKDD-09)*, Springer, *Lecture Notes in Computer Science*, vol 5781, pp 548–563
- Huth M, Ryan M (2004) *Logic in Computer Science: Modelling and Reasoning About Systems*. Cambridge University Press
- Jaeger M (1997) Relational Bayesian networks. In: Geiger D, Shenoy PP (eds) *Proceedings of the 13th Conference on Uncertainty in Artificial Intelligence (UAI-97)*, Morgan Kaufmann Publishers, pp 266–273
- Jaeger M (2008) Model-theoretic expressivity analysis. In: De Raedt et al (2008), pp 325–339
- Kersting K (2012) Lifted probabilistic inference. In: De Raedt L, Bessière C, Dubois D, Doherty P, Frasconi P, Heintz F, Lucas PJF (eds) *Proceedings of the 20th European Conference on Artificial Intelligence (ECAI-12)*, IOS Press, *Frontiers in Artificial Intelligence and Applications*, vol 242, pp 33–38
- Kersting K, De Raedt L (2001) Bayesian logic programs. CoRR cs.AI/0111058
- Kersting K, De Raedt L (2008) Basic principles of learning Bayesian logic programs. In: De Raedt et al (2008), pp 189–221
- Kersting K, De Raedt L, Raiko T (2006) Logical hidden Markov models. *Journal of Artificial Intelligence Research (JAIR)* 25:425–456
- Kimmig A, Santos Costa V, Rocha R, Demoen B, De Raedt L (2008) On the efficient execution of ProbLog programs. In: Garcia de la Banda M, Pontelli E (eds) *Proceedings of the 24th International Conference on Logic Programming*

- (ICLP-08), Springer, Lecture Notes in Computer Science, vol 5366, pp 175–189
- Kimmig A, Gutmann B, Santos Costa V (2009) Trading memory for answers: Towards tabling ProbLog. In: Domingos P, Kersting K (eds) Proceedings of the International Workshop on Statistical Relational Learning (SRL-2009)
- Kimmig A, Van den Broeck G, De Raedt L (2011a) An algebraic Prolog for reasoning about possible worlds. In: Burgard W, Roth D (eds) Proceedings of the 25th AAAI Conference on Artificial Intelligence (AAAI-11), AAAI Press, pp 209–214
- Kimmig A, Demoen B, De Raedt L, Santos Costa V, Rocha R (2011b) On the implementation of the probabilistic logic programming language ProbLog. *Theory and Practice of Logic Programming (TPLP)* 11(2–3):235–262
- Kimmig A, Mihalkova L, Getoor L (2014) Lifted graphical models: A survey. *Machine Learning* (published online)
- Koller D, Pfeffer A (1998) Probabilistic frame-based systems. In: Mostow J, Rich C (eds) Proceedings of the 15th National Conference on Artificial Intelligence (AAAI-98), AAAI Press / MIT Press, pp 580–587
- Lloyd JW (1989) *Foundations of Logic Programming*, 2nd edn. Springer
- Lodhi H, Muggleton S (2005) Modelling metabolic pathways using stochastic logic programs-based ensemble methods. In: Danos V, Schächter V (eds) Revised Selected Papers of the International Conference on Computational Methods in Systems Biology (CMSB-04), Springer, Lecture Notes in Computer Science, vol 3082, pp 119–133
- Mansinghka VK, Kulkarni TD, Perov YN, Tenenbaum JB (2013) Approximate Bayesian image interpretation using generative probabilistic graphics programs. In: Burges CJC, Bottou L, Ghahramani Z, Weinberger KQ (eds) Proceedings of the 27th Annual Conference on Neural Information Processing Systems (NIPS-13), Neural Information Processing Systems, pp 1520–1528
- Mantadelis T, Janssens G (2011) Nesting probabilistic inference. In: Abreu S, Santos Costa V (eds) Proceedings of the Colloquium on Implementation of Constraint and Logic Programming Systems (CICLOPS-11)
- Milch B, Marthi B, Russell SJ, Sontag D, Ong DL, Kolobov A (2005) BLOG: Probabilistic models with unknown objects. In: Kaelbling LP, Saffiotti A (eds) Proceedings of the 19th International Joint Conference on Artificial Intelligence (IJCAI-05), Professional Book Center, pp 1352–1359
- Moldovan B, Thon I, Davis J, De Raedt L (2013) MCMC estimation of conditional probabilities in probabilistic programming languages. In: van der Gaag LC (ed) Proceedings of the 12th European Conference on Symbolic and Quantitative Approaches to Reasoning with Uncertainty (ECSQARU-13), Springer, Lecture Notes in Computer Science, vol 7958, pp 436–448
- Mørk S, Holmes I (2012) Evaluating bacterial gene-finding HMM structures as probabilistic logic programs. *Bioinformatics* 28(5):636–642
- Muggleton S (1996) Stochastic logic programs. In: De Raedt L (ed) *Advances in Inductive Logic Programming*, IOS Press, pp 254–264
- Natarajan S, Tadepalli P, Altendorf E, Dietterich TG, Fern A, Restificar AC (2005) Learning first-order probabilistic models with combining rules. In: De Raedt L, Wrobel S (eds) Proceedings of the 22nd International Conference on Machine Learning (ICML-05), ACM, ACM International Conference Proceeding Series, vol 119, pp 609–616

- Nitti D, De Laet T, De Raedt L (2013) A particle filter for hybrid relational domains. In: Amato N (ed) Proceedings of the IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS-13), IEEE, pp 2764–2771
- Pearl J (1988) Probabilistic reasoning in intelligent systems: Networks of plausible inference. Morgan Kaufmann Publishers
- Pfeffer A (2001) IBAL: A probabilistic rational programming language. In: Nebel B (ed) Proceedings of the 17th International Joint Conference on Artificial Intelligence (IJCAI-01), Morgan Kaufmann Publishers, pp 733–740
- Pfeffer A (2009) Figaro: An object-oriented probabilistic programming language. Tech. rep., Charles River Analytics
- Poole D (1992) Logic programming, abduction and probability. In: Proceedings of the International Conference on Fifth Generation Computing Systems (FGCS-92), IOS Press, pp 530–538
- Poole D (1993) Probabilistic Horn abduction and Bayesian networks. *Artificial Intelligence* 64:81–129
- Poole D (1995) Exploiting the rule structure for decision making within the independent choice logic. In: Besnard P, Hanks S (eds) Proceedings of the Eleventh Annual Conference on Uncertainty in Artificial Intelligence (UAI-95), Morgan Kaufmann Publishers, pp 454–463
- Poole D (2000) Abducing through negation as failure: Stable models within the independent choice logic. *Journal of Logic Programming* 44(1–3):5–35
- Poole D (2003) First-order probabilistic inference. In: Gottlob G, Walsh T (eds) Proceedings of the 18th International Joint Conference on Artificial Intelligence (IJCAI-03), Morgan Kaufmann Publishers, pp 985–991
- Poole D (2008) The independent choice logic and beyond. In: De Raedt et al (2008), pp 222–243
- Poole D (2010) Probabilistic programming languages: Independent choices and deterministic systems. In: R Dechter HG, Halpern J (eds) *Heuristics, Probability and Causality: A Tribute to Judea Pearl*, College Publications, pp 253–269
- Poon H, Domingos P (2006) Sound and efficient inference with probabilistic and deterministic dependencies. In: Gil Y, Mooney RJ (eds) Proceedings of the 21st National Conference on Artificial Intelligence (AAAI-06), AAAI Press, pp 458–463
- Renkens J, Van den Broeck G, Nijssen S (2012) k-optimal: A novel approximate inference algorithm for ProbLog. *Machine Learning* 89(3):215–231
- Renkens J, Kimmig A, Van den Broeck G, De Raedt L (2014) Explanation-based approximate weighted model counting for probabilistic logics. In: Brodley CE, Stone P (eds) Proceedings of the 28th AAAI Conference on Artificial Intelligence (AAAI-14), AAAI Press, pp 2490–2496
- Richardson M, Domingos P (2006) Markov logic networks. *Machine Learning* 62(1–2):107–136
- Riguzzi F (2009) Extended semantics and inference for the Independent Choice Logic. *Logic Journal of the IGPL* 17(6):589–629
- Riguzzi F (2013a) cplint Manual. <https://sites.google.com/a/unife.it/ml/cplint/cplint-manual>
- Riguzzi F (2013b) MCINTYRE: A Monte Carlo system for probabilistic logic programming. *Fundamenta Informaticae* 124(4):521–541
- Riguzzi F (2014) Speeding up inference for probabilistic logic programs. *The Computer Journal* 57(3):347–363

- Riguzzi F, Swift T (2011) The PITA system: Tabling and answer subsumption for reasoning under uncertainty. *Theory and Practice of Logic Programming (TPLP)* 11(4–5):433–449
- Roy D, Mansinghka V, Winn J, McAllester D, Tenenbaum D (eds) (2008) *Probabilistic Programming*, NIPS Workshop
- Russell SJ, Norvig (2003) *Artificial Intelligence: A Modern Approach*, 2nd edn. Prentice Hall
- Santos Costa V, Paes A (2009) On the relationship between PRISM and  $CLP(BN)$ . In: Domingos P, Kersting K (eds) *Proceedings of the International Workshop on Statistical Relational Learning (SRL-2009)*
- Santos Costa V, Page D, Qazi M, Cussens J (2003)  $CLP(BN)$ : constraint logic programming for probabilistic knowledge. In: Meek C, Kjærulff U (eds) *Proceedings of the 19th Conference on Uncertainty in Artificial Intelligence (UAI-03)*, Morgan Kaufmann Publishers, pp 517–524
- Santos Costa V, Page D, Cussens J (2008)  $CLP(BN)$ : Constraint logic programming for probabilistic knowledge. In: De Raedt et al (2008), pp 156–188
- Sato T (1995) A statistical learning method for logic programs with distribution semantics. In: Sterling L (ed) *Proceedings of the 12th International Conference on Logic Programming (ICLP-95)*, MIT Press, pp 715–729
- Sato T (2011) A general MCMC method for Bayesian inference in logic-based probabilistic modeling. In: Walsh T (ed) *Proceedings of the 22nd International Joint Conference on Artificial Intelligence (IJCAI-11)*, AAAI Press, pp 1472–1477
- Sato T, Kameya Y (1997) PRISM: A language for symbolic-statistical modeling. In: Pollack ME (ed) *Proceedings of the 15th International Joint Conference on Artificial Intelligence (IJCAI-97)*, Morgan Kaufmann Publishers, pp 1330–1339
- Sato T, Kameya Y (2001) Parameter learning of logic programs for symbolic-statistical modeling. *Journal of Artificial Intelligence Research (JAIR)* 15:391–454
- Sato T, Kameya Y, Zhou NF (2005) Generative modeling with failure in PRISM. In: Kaelbling LP, Saffiotti A (eds) *Proceedings of the 19th International Joint Conference on Artificial Intelligence (IJCAI-05)*, Professional Book Center, pp 847–852
- Sato T, Kameya Y, Kurihara K (2008) Variational Bayes via propositionalized probability computation in PRISM. *Annals of Mathematics and Artificial Intelligence* 54(1–3):135–158
- Skarlatidis A, Artikis A, Filippou J, Paliouras G (2015) A probabilistic logic programming event calculus. *Theory and Practice of Logic Programming (TPLP)* 15(2):213–245
- Sneyers J, Meert W, Vennekens J, Kameya Y, Sato T (2010) CHR(PRISM)-based probabilistic logic learning. *Theory and Practice of Logic Programming (TPLP)* 10(4–6):433–447
- Stuhlmüller A, Tenenbaum JB, Goodman ND (2010) Learning structured generative concepts. In: Ohlsson S, Catrambone R (eds) *Proceedings of the 32nd Annual Conference of the Cognitive Science Society*, Cognitive Science Society, pp 2296–2301
- Suciu D, Olteanu D, Ré C, Koch C (2011) *Probabilistic Databases*, Synthesis Lectures on Data Management, vol 16. Morgan & Claypool Publishers

- Thon I, Landwehr N, De Raedt L (2011) Stochastic relational processes: Efficient inference and applications. *Machine Learning* 82(2):239–272
- Valiant LG (1979) The complexity of enumeration and reliability problems. *SIAM Journal on Computing* 8(3):410–421
- Van den Broeck G, Thon I, van Otterlo M, De Raedt L (2010) DTProbLog: A decision-theoretic probabilistic Prolog. In: Fox M, Poole D (eds) *Proceedings of the 24th AAAI Conference on Artificial Intelligence (AAAI-10)*, AAAI Press, pp 1217–1222
- Van den Broeck G, Taghipour N, Meert W, Davis J, De Raedt L (2011) Lifted probabilistic inference by first-order knowledge compilation. In: Walsh T (ed) *Proceedings of the 22nd International Joint Conference on Artificial Intelligence (IJCAI-11)*, AAAI Press, pp 2178–2185
- Van den Broeck G, Meert W, Darwiche A (2014) Skolemization for weighted first-order model counting. In: Baral C, Giacomo GD, Eiter T (eds) *Proceedings of the 14th International Conference on Principles of Knowledge Representation and Reasoning (KR-14)*, AAAI Press, pp 111–120
- Vennekens J (2007) Algebraic and logical study of constructive processes in knowledge representation. PhD thesis, K.U. Leuven
- Vennekens J, Verbaeten S, Bruynooghe M (2004) Logic programs with annotated disjunctions. In: Demoen B, Lifschitz V (eds) *Proceedings of the 20th International Conference on Logic Programming (ICLP-04)*, Springer, *Lecture Notes in Computer Science*, vol 3132, pp 431–445
- Vennekens J, Denecker M, Bruynooghe M (2006) Representing causal information about a probabilistic process. In: Fisher M, van der Hoek W, Konev B, Lisitsa A (eds) *Proceedings of the 10th European Conference on Logics in Artificial Intelligence (JELIA-06)*, Springer, *Lecture Notes in Computer Science*, vol 4160, pp 452–464
- Vennekens J, Denecker M, Bruynooghe M (2009) CP-logic: A language of causal probabilistic events and its relation to logic programming. *Theory and Practice of Logic Programming (TPLP)* 9(3):245–308
- Wang WY, Mazaitis K, Cohen WW (2013) Programming with personalized PageRank: A locally groundable first-order probabilistic logic. In: He Q, Iyengar A, Nejdil W, Pei J, Rastogi R (eds) *Proceedings of the 22nd ACM International Conference on Information and Knowledge Management (CIKM-13)*, ACM, pp 2129–2138
- Warren DS (1992) Memoing for logic programs. *Communications of the ACM (CACM)* 35(3):93–111
- Wellman MP, Breese JS, Goldman RP (1992) From knowledge bases to decision models. *The Knowledge Engineering Review* 7(1):35–53

## A Logic Programming Basics

The basic building blocks of logic programs are *variables* (denoted by strings starting with upper case letters), *constants*, *functors* and *predicates* (all denoted by strings starting with lower case letters). A *term* is a variable, a constant, or a functor  $f$  of *arity*  $n$  followed by  $n$  terms  $t_i$ , i.e.,  $f(t_1, \dots, t_n)$ . An *atom* is a predicate  $p$  of arity  $n$  followed by  $n$  terms  $t_i$ , i.e.,  $p(t_1, \dots, t_n)$ . A predicate  $p$  of arity  $n$  is also written as  $p/n$ . A *literal* is an atom or a negated atom  $\text{not}(p(t_1, \dots, t_n))$ . A *definite clause* is a universally quantified expression of the form  $h :- b_1, \dots, b_n$  where  $h$  and the  $b_i$  are atoms.  $h$  is called the *head* of the clause, and  $b_1, \dots, b_n$  its *body*. Informally, the meaning of such a clause is that if all the  $b_i$  are true,  $h$  has to be true as well. A *normal clause* is a universally quantified expression of the form  $h :- b_1, \dots, b_n$  where  $h$  is an atom and the  $b_i$  are literals. If  $n = 0$ , a clause is called *fact* and simply written as  $h$ . A *definite clause program* or *logic program* for short is a finite set of definite clauses. A *normal logic program* is a finite set of normal clauses. A *substitution*  $\theta$  is an expression of the form  $\{V_1/t_1, \dots, V_m/t_m\}$  where the  $V_i$  are different variables and the  $t_i$  are terms. Applying a substitution  $\theta$  to an expression  $e$  (term or clause) yields the *instantiated* expression  $e\theta$  where all variables  $V_i$  in  $e$  have been simultaneously replaced by their corresponding terms  $t_i$  in  $\theta$ . If an expression does not contain variables it is *ground*. Two expressions  $e_1$  and  $e_2$  can be *unified* if and only if there are substitutions  $\theta_1$  and  $\theta_2$  such that  $e_1\theta_1 = e_2\theta_2$ . In Prolog, unification is written using  $=$  as an infix predicate.

The *Herbrand base* of a logic program is the set of ground atoms that can be constructed using the predicates, functors and constants occurring in the program.<sup>27</sup> Subsets of the Herbrand base are called *Herbrand interpretations*. A Herbrand interpretation is a *model* of a clause  $h :- b_1, \dots, b_n$  if for every substitution  $\theta$  such that all  $b_i\theta$  are in the interpretation,  $h\theta$  is in the interpretation as well. It is a model of a logic program if it is a model of all clauses in the program. The model-theoretic semantics of a definite clause program is given by its smallest Herbrand model with respect to set inclusion, the so-called *least Herbrand model* (which is unique). We say that a logic program  $P$  *entails* an atom  $a$ , denoted  $P \models a$ , if and only if  $a$  is true in the least Herbrand model of  $P$ .

The main inference task of a logic programming system is to determine whether a given atom, also called *query* (or *goal*), is true in the least Herbrand model of a logic program. If the answer is yes (or no), we also say that the query *succeeds* (or *fails*). If such a query is not ground, inference asks for the existence of an *answer substitution*, that is, a substitution that grounds the query into an atom that is part of the least Herbrand model.

Normal logic programs use the notion of *negation as failure*, that is, for a ground atom  $a$ ,  $\text{not}(a)$  is true exactly if  $a$  cannot be proven in the program. They are not guaranteed to have a unique minimal Herbrand model. Various ways to define the canonical model of such programs have been studied; see, e.g., (Lloyd, 1989, Chapter 3) for an overview.

## B Annotated Disjunctions and Probabilistic Facts

As mentioned in Section 2.2, each annotated disjunction can be equivalently represented using a set of probabilistic facts and deterministic clauses. Using probabilistic facts is not sufficient, as those correspond to independent random variables. For instance, using probabilistic facts

$$\frac{1}{3} :: \text{color}(\text{green}). \quad \frac{1}{3} :: \text{color}(\text{red}). \quad \frac{1}{3} :: \text{color}(\text{blue}).$$

the probability of  $\text{color}(\text{green})$ ,  $\text{color}(\text{red})$  and  $\text{color}(\text{blue})$  all being true is  $1/27$ , whereas it is 0 for the annotated disjunction

$$\frac{1}{3} :: \text{color}(\text{green}); \frac{1}{3} :: \text{color}(\text{red}); \frac{1}{3} :: \text{color}(\text{blue}).$$

On the other hand, we can exploit the fact that negation of probabilistic facts is easily handled under the distribution semantics<sup>28</sup> to encode an AD by simulating a sequential choice mecha-

<sup>27</sup> If the program does not contain constants, one arbitrary constant is added.

<sup>28</sup> For a probabilistic fact  $p :: f$ ,  $\text{not}(f)$  succeeds in a possible world exactly if  $f$  is not among the probabilistic facts included in that world; cf. Section 4.6 for a more general discussion of negation.



nism.<sup>29</sup> With this encoding, the three possible outcomes are mutually exclusive as in the AD and exactly one will be true in any possible world:

```

1/3 :: sw_1(color(green)).    1/2 :: sw_1(color(red)).    1 :: sw_1(color(blue)).
color(green) :- sw_1(color(green)).
color(red)    :- not(sw_1(color(green))), sw_1(color(red)).
color(blue)   :- not(sw_1(color(green))), not(sw_1(color(red))), sw_1(color(blue)).

```

Note that the probabilities have been adapted to reproduce the probabilities of the different head atoms; we discuss the details of this adaptation below.<sup>30</sup> This mapping follows the general idea of representing a probabilistic model in an *augmented space* where random variables can be assumed independent, while capturing the dependencies in the deterministic part of the program (Poole, 2010).

For non-ground ADs, all logical variables have to be included in the probabilistic facts to ensure that all groundings correspond to independent random events. For instance, the AD

$$\frac{1}{2} :: \text{color}(\text{green}); \frac{1}{2} :: \text{color}(\text{red}) :- \text{ball}(\text{Ball})$$

would be represented as

```

1/2 :: sw_1(color(green), Ball).    1 :: sw_1(color(red), Ball).
color(green) :- ball(Ball), sw_1(color(green), Ball).
color(red)   :- ball(Ball), not(sw_1(color(green), Ball)), sw_1(color(red), Ball).

```

As this example suggests, annotated disjunctions can be expressed using probabilistic facts by representing each annotated disjunction using the set of probabilistic facts  $\tilde{p}_i :: \text{sw\_id}(h_i, v_1, \dots, v_f)$  and the following clauses

$$h_i :- b_1, \dots, b_m, \text{not}(\text{sw\_id}(h_1, v_1, \dots, v_f)), \dots, \text{not}(\text{sw\_id}(h_{i-1}, v_1, \dots, v_f)), \text{sw\_id}(h_i, v_1, \dots, v_f) \quad (8)$$

where *id* is a unique identifier for a particular AD and  $v_1, \dots, v_f$  are the free variables in the body of the AD. The probability  $\tilde{p}_1$  is defined as  $p_1$  and for  $i > 1$  it is

$$\tilde{p}_i := \begin{cases} p_i \cdot \left(1 - \sum_{j=1}^{i-1} p_j\right)^{-1} & \text{if } p_i > 0 \\ 0 & \text{if } p_i = 0 \end{cases} \quad (9)$$

One can recover the original probabilities from  $\tilde{p}$  by setting  $p_1 := \tilde{p}_1$  and iteratively applying the following transformation for  $i = 2, 3, \dots, n$

$$p_i := \tilde{p}_i \cdot \left(1 - \sum_{j=1}^{i-1} p_j\right) \quad (10)$$

Equation (9) and (10) together define a bijection between  $p$  and  $\tilde{p}$  which allows one to use parameter learning in either representation and map learned probabilities onto the other representation. If the  $p_i$  sum to 1, it is possible to drop the last probabilistic fact  $\text{sw\_id}(h_n)$  since its probability  $\tilde{p}_n$  is 1.

<sup>29</sup> used, e.g., by Sato and Kameya (1997) with parameters learned from data

<sup>30</sup> This transformation is correct for computing success probabilities, but care has to be taken to accommodate for the additional random variables in MPE inference.