

An Approach to Probabilistic Data Integration for the Semantic Web

Andrea Calì¹ Thomas Lukasiewicz^{2 3}

¹ Facoltà di Scienze e Tecnologie Informatiche
Libera Università di Bolzano, Italy

² Dipartimento di Informatica e Sistemistica
Università di Roma "La Sapienza", Italy

³ Institut für Informationssysteme
Technische Universität Wien, Austria

Outline

- Probabilistic Description Logic Programs
- Probabilistic Data Integration

Probabilistic description logic programs (Lukasiewicz, ECSQARU-05, URSW-05, IJAR-07):

- Combine DLs, normal programs under the answer set resp. well-founded semantics, and probabilistic uncertainty.
- Generalize Poole's ICL, which in turn generalizes influence diagrams, Bayesian networks, MDPs, normal form games, and structural causal models.
- Consist of (i) a description logic knowledge base L , (ii) a normal program P involving queries to L , and (iii) a probability distribution on a set of total choices.
- Represent a set of probability distributions on a set of first-order interpretations.
- Instead of querying L in P , a variant of dl-programs also allows for using L to constrain the terms in P .

$PC \sqcup Camera \sqsubseteq Electronics$; $PC \sqcap Camera \sqsubseteq \perp$;
 $Book \sqcup Electronics \sqsubseteq Product$; $Book \sqcap Electronics \sqsubseteq \perp$;
 $Textbook \sqsubseteq Book$;

$Product \sqsubseteq \geq 1 \text{ related}$;
 $\geq 1 \text{ related} \sqcup \geq 1 \text{ related}^- \sqsubseteq Product$;

$Textbook(tb_ai)$; $Textbook(tb_lp)$;
 $PC(pc_ibm)$; $PC(pc_hp)$;

$related(tb_ai, tb_lp)$; $related(pc_ibm, pc_hp)$;
 $provides(ibm, pc_ibm)$; $provides(hp, pc_hp)$.

- Description logic knowledge base L (as above).

- Finite set L of dl-rules:

$avoid(X) \leftarrow DL[Camera](X), not\ offer(X), avoid_pos;$

$offer(X) \leftarrow DL[PC \uplus pc; Electronics](X), not\ brand_new(X), offer_pos;$

$buy(C, X) \leftarrow needs(C, X), view(X), not\ avoid(X), v_buy_pos;$

$buy(C, X) \leftarrow needs(C, X), buy(C, Y), also_buy(Y, X), a_buy_pos.$

- Probabilities associated with all atomic choices:

$\{avoid_pos: .9, avoid_neg: .1\}, \{offer_pos: .9, offer_neg: .1\},$

$\{v_buy_pos: .7, v_buy_neg: .3\}, \{a_buy_pos: .7, a_buy_neg: .3\}.$

- ⇒ Probabilities associated with all total choices:

$\{avoid_pos, offer_pos, v_buy_pos, a_buy_pos\} : .9 \times .9 \times .7 \times .7, \dots$

- ⇒ Probability intervals associated with all conditional events:

Probabilistic query: $?(buy(c, x) \mid needs(c, x) \wedge buy(c, y) \wedge$
 $also_buy(y, x) \wedge view(x) \wedge not\ avoid(x))[L, U]$

A **data integration system** is a triple $\langle G, S, M \rangle$, where

- G is the **global** or **mediated schema**, representing the domain of interest of the system,
- S is the **source schema**, representing the data sources that take part in the system, and
- M is a **mapping** that establishes a relation between the source schema and the global schema.

The vocabulary Φ is partitioned into:

- Φ_G : its symbols are of arity at least 1, and represent the (virtual) global predicates;
- Φ_S : its symbols are of arity at least 1, and represent the source predicates;
- Φ_C : its symbols are constants.

A probabilistic dl-program modeling a data integration system:

- **source dl-rules** (over Φ_S and Φ_C): they express properties and constraints of the data sources;
- **global dl-rules** (over Φ_G and Φ_C): they express properties and constraints on the global schema, which enhance its expressiveness to better fit the application domain;
- **mapping dl-rules**: they specify the mapping M between Φ_G and Φ_S ; they have only predicates in Φ_S and Φ_C in the body, and only predicates in Φ_G and Φ_C in the head.

G consists of Φ_G and the global dl-rules, S consists of Φ_S and the source dl-rules, and M consists of the mapping dl-rules.

The global predicate $buy(C, X)$ is derived from either the source predicate $s_1(C, X, Y)$ or the source predicates $s_2(C, D)$ and $s_3(D, X)$.

Suppose that C resp. X are restricted to customers resp. products from a description logic knowledge base L , and that there may be inconsistencies between the two different ways of deriving $buy(C, X)$.

To consistently integrate them, we assign to each derivation a total choice from $\{a, \bar{a}\}$ along with user-defined probabilities that depend on the reliability of the derivations (e.g., $\mu(a) = 0.7$ and $\mu(\bar{a}) = 0.3$).

The following dl-rules then realize the probabilistic data integration:

$$\begin{aligned} buy(C, X) &\leftarrow s_1(C, X, Y), DL[Customer](C), DL[Product](X), a; \\ buy(C, X) &\leftarrow s_2(C, D), s_3(D, X), DL[Customer](C), DL[Product](X), \bar{a}. \end{aligned}$$

Every fact that holds by the first resp. second dl-rule has the probability 0.7 resp. 0.3, while every fact that holds by both dl-rules has the probability 1.