

Exploiting Description Logic Reasoners in Inductive Logic Programming Systems: An Experience within the Semantic Web area

Francesca A. Lisi

Dipartimento di Informatica, Università degli Studi di Bari,
Via E. Orabona 4, 70125 Bari, Italy
lisi@di.uniba.it

Abstract. In spite of the increasing effort spent on building ontologies for the Semantic Web, little attention has been paid to the impact of these ontologies on knowledge-based intelligent systems such as Inductive Logic Programming (ILP) systems which were not conceived for dealing with DL knowledge bases. In this paper, we describe an extension of the ILP system \mathcal{AL} -QUIN to deal with a background knowledge in the form of OWL DL ontology. The extension consists of a preprocessing of the ontology that mainly relies on the services of the DL reasoner Pellet.

1 Introduction

Description Logics (DLs) are the most currently used among the logical formalisms proposed by Ontological Engineering [5]. Also the DL-based approach to Ontological Engineering is playing a relevant role in the definition of the Semantic Web. The Semantic Web is the vision of the World Wide Web enriched by machine-processable information which supports the user in his tasks [2]. Its architecture consists of several layers, each of which is equipped with an ad-hoc mark-up language. DLs, more precisely the very expressive DL \mathcal{SHIQ} , have guided the design of the mark-up language OWL for the *ontological layer* [6]. A DL reasoner, Pellet [16], has been recently proposed for OWL. In spite of the increasing effort spent on building ontologies for the Semantic Web, little attention has been paid to the impact of these ontologies on knowledge-based intelligent systems such as Inductive Logic Programming (ILP) systems which were not conceived for dealing with DL knowledge bases. Note that the use of background knowledge has been widely recognized as one of the strongest points of ILP when compared to other forms of concept learning and has been empirically studied in several application domains [11]. Yet the background knowledge in ILP systems is often not organized around a well-formed conceptual model and still ignores the latest developments in Knowledge Engineering such as ontologies and ontology languages based on DLs. In a recent position paper, Page and Srinivasan have pointed out that the use of special-purpose reasoners in ILP is among the pressing issues that have arisen from the most challenging ILP applications of today [12]. We think that this is the case for ILP applications in the

Semantic Web area. In this paper we report on an experience with DL reasoners in ILP within the Semantic Web application area. In particular, we describe an extension of the ILP system \mathcal{AL} -QUIN [10] to deal with a background knowledge in the form of OWL DL ontology. The extension consists of a preprocessing of the ontology that mainly relies on the reasoning services of Pellet.

The paper is structured as follows. Section 2 briefly describes \mathcal{AL} -QUIN. Section 3 illustrates the use of Pellet in \mathcal{AL} -QUIN. Section 4 concludes the paper.

2 The ILP system \mathcal{AL} -QuIn

The ILP system \mathcal{AL} -QUIN (\mathcal{AL} -log Query Induction) [10] supports a data mining task known under the name of *frequent pattern discovery*. In data mining a *pattern* is considered as an intensional description (expressed in a given language \mathcal{L}) of a subset of a given data set \mathbf{r} . The *support* of a pattern is the relative frequency of the pattern within \mathbf{r} and is computed with the evaluation function *supp*. The task of frequent pattern discovery aims at the extraction of all *frequent* patterns, i.e. all patterns whose support exceeds a user-defined threshold of *minimum support*. \mathcal{AL} -QUIN solves a variant of the frequent pattern discovery problem which takes concept hierarchies into account during the discovery process, thus yielding descriptions at multiple granularity levels up to a maximum level $maxG$. More formally, given

- a data set \mathbf{r} including a taxonomy \mathcal{T} where a reference concept C_{ref} and task-relevant concepts are designated,
- a multi-grained language $\{\mathcal{L}^l\}_{1 \leq l \leq maxG}$ of patterns
- a set $\{minsup^l\}_{1 \leq l \leq maxG}$ of user-defined minimum support thresholds

the problem of *frequent pattern discovery at l levels of description granularity*, $1 \leq l \leq maxG$, is to find the set \mathcal{F} of all the patterns $P \in \mathcal{L}^l$ that describe the reference concept w.r.t. the task-relevant concepts and turn out to be frequent in \mathbf{r} . Note that P 's with support s such that (i) $s \geq minsup^l$ and (ii) all ancestors of P w.r.t. \mathcal{T} are frequent in \mathbf{r} . Note that a pattern Q is considered to be an ancestor of P if it is a coarser-grained version of P .

Example 1. As a showcase we consider the task of finding frequent patterns that describe Middle East countries (reference concept) w.r.t. the religions believed and the languages spoken (task-relevant concepts) at three levels of granularity ($maxG = 3$). Minimum support thresholds are set to the following values: $minsup^1 = 20\%$, $minsup^2 = 13\%$, and $minsup^3 = 10\%$. The data set and the language of patterns will be illustrated in Example 2 and 3, respectively.

In \mathcal{AL} -QUIN data and patterns are represented according to the hybrid knowledge representation and reasoning system \mathcal{AL} -log [4]. In particular, the data set \mathbf{r} is represented as an \mathcal{AL} -log knowledge base \mathcal{B} , thus composed of a structural part and a relational part. The structural subsystem Σ is based on \mathcal{ALC} [14] whereas the relational subsystem Π is based on an extended form of DATALOG [3] that is obtained by using \mathcal{ALC} concept assertions essentially as type constraints on variables.

Example 2. For the task of interest, we consider an \mathcal{AL} -log knowledge base \mathcal{B}_{CIA} that integrates a \mathcal{ALC} component Σ_{CIA} containing taxonomies rooted into the concepts **Country**, **EthnicGroup**, **Language** and **Religion** and a DATALOG component Π_{CIA} containing facts¹ extracted from the on-line 1996 CIA World Fact Book². Note that Middle East countries have been defined as Asian countries that host at least one Middle Eastern ethnic group:

`MiddleEastCountry` \equiv `AsianCountry` \sqcap \exists `Hosts.MiddleEastEthnicGroup`.

In particular, Armenia ('ARM') and Iran ('IR') are classified as Middle East countries because the following membership assertions hold in Σ_{CIA} :

```
'ARM':AsianCountry.
'IR':AsianCountry.
'Arab':MiddleEastEthnicGroup.
'Armenian':MiddleEastEthnicGroup.
<'ARM','Armenian'>:Hosts.
<'IR','Arab'>:Hosts.
```

Also Π_{CIA} includes constrained DATALOG clauses such as:

```
believes(Code, Name)  $\leftarrow$ 
    religion(Code, Name, Percent) & Code:Country, Name:Religion.
speaks(Code, Name)  $\leftarrow$ 
    language(Code, Name, Percent) & Code:Country, Name:Language.
```

that define views on the relations `religion` and `language`, respectively.

The language $\mathcal{L} = \{\mathcal{L}^l\}_{1 \leq l \leq \max G}$ of patterns allows for the generation of \mathcal{AL} -log unary conjunctive queries, called \mathcal{O} -queries. Given a reference concept C_{ref} , an \mathcal{O} -query Q to an \mathcal{AL} -log knowledge base \mathcal{B} is a (linked and connected)³ constrained DATALOG clause of the form

$$Q = q(X) \leftarrow \alpha_1, \dots, \alpha_m \& X : C_{ref}, \gamma_1, \dots, \gamma_n$$

where X is the *distinguished variable* and the remaining variables occurring in the body of Q are the *existential variables*. Note that α_j , $1 \leq j \leq m$, is a DATALOG literal whereas γ_k , $1 \leq k \leq n$, is an assertion that constrains a variable already appearing in any of the α_j 's to vary in the range of individuals of a concept defined in \mathcal{B} . The \mathcal{O} -query

$$Q_t = q(X) \leftarrow \& X : C_{ref}$$

is called *trivial* for \mathcal{L} because it only contains the constraint for the *distinguished variable* X . Furthermore the language \mathcal{L} is *multi-grained*, i.e. it contains expressions at multiple levels of description granularity. Indeed it is implicitly defined

¹ <http://www.dbis.informatik.uni-goettingen.de/Mondial/mondial-rel-facts.flp>

² <http://www.odci.gov/cia/publications/factbook/>

³ For the definition of linkedness and connectedness see [11].

by a *declarative bias specification* which consists of a finite alphabet Δ of DATALOG predicate names and finite alphabets Γ^l (one for each level l of description granularity) of \mathcal{ALC} concept names. Note that the α_i 's are taken from \mathcal{A} and γ_j 's are taken from Γ^l . We impose \mathcal{L} to be finite by specifying some bounds, mainly $maxD$ for the maximum depth of search and $maxG$ for the maximum level of granularity.

Example 3. To accomplish the task of Example 1 we define \mathcal{L}_{CIA} as the set of \mathcal{O} -queries with $C_{ref} = \text{MiddleEastCountry}$ that can be generated from the alphabet $\Delta = \{\text{believes}/2, \text{speaks}/2\}$ of DATALOG binary predicate names, and the alphabets

$$\begin{aligned}\Gamma^1 &= \{\text{Language}, \text{Religion}\} \\ \Gamma^2 &= \{\text{IndoEuropeanLanguage}, \dots, \text{MonotheisticReligion}, \dots\} \\ \Gamma^3 &= \{\text{IndoIranianLanguage}, \dots, \text{MuslimReligion}, \dots\}\end{aligned}$$

of \mathcal{ALC} concept names for $1 \leq l \leq 3$, up to $maxD = 5$. Examples of \mathcal{O} -queries in \mathcal{L}_{CIA} are:

$$\begin{aligned}Q_t &= \text{q}(X) \leftarrow \& X:\text{MiddleEastCountry} \\ Q_1 &= \text{q}(X) \leftarrow \text{speaks}(X, Y) \& X:\text{MiddleEastCountry}, Y:\text{Language} \\ Q_2 &= \text{q}(X) \leftarrow \text{speaks}(X, Y) \& X:\text{MiddleEastCountry}, Y:\text{IndoEuropeanLanguage} \\ Q_3 &= \text{q}(X) \leftarrow \text{believes}(X, Y) \& X:\text{MiddleEastCountry}, Y:\text{MuslimReligion}\end{aligned}$$

where Q_t is the trivial \mathcal{O} -query for \mathcal{L}_{CIA} , $Q_1 \in \mathcal{L}_{CIA}^1$, $Q_2 \in \mathcal{L}_{CIA}^2$, and $Q_3 \in \mathcal{L}_{CIA}^3$. Note that Q_1 is an ancestor of Q_2 .

The *support* of an \mathcal{O} -query $Q \in \mathcal{L}^l$ w.r.t an \mathcal{AL} -log knowledge base \mathcal{B} is defined as

$$\text{supp}(Q, \mathcal{B}) = | \text{answerset}(Q, \mathcal{B}) | / | \text{answerset}(Q_t, \mathcal{B}) |$$

where $\text{answerset}(Q, \mathcal{B})$ is the set of correct answers to Q w.r.t. \mathcal{B} . An *answer* to Q is a ground substitution θ for the distinguished variable of Q . An answer θ to Q is a *correct (resp. computed) answer* w.r.t. \mathcal{B} if there exists at least one correct (resp. computed) answer to $\text{body}(Q)\theta$ w.r.t. \mathcal{B} . Thus the computation of support relies on query answering in \mathcal{AL} -log.

Example 4. The pattern Q_2 turns out to be frequent because it has support $\text{supp}(Q_2, \mathcal{B}_{CIA}) = (2/15)\% = 13.3\% (\geq \text{minsup}^2)$. It is to be read as '13.3 % of Middle East countries speak an Indoeuropean language'. The two correct answers to Q_2 w.r.t. \mathcal{B}_{CIA} are 'ARM' and 'IR'.

3 Exploiting Pellet in \mathcal{AL} -QuIn

3.1 Coverage of observations

In ILP the evaluation of inductive hypotheses (like candidate patterns in frequent pattern discovery) w.r.t. a set of observations (data units) is usually referred to as

the *coverage test* because it checks which observations satisfy (are covered by) the hypothesis. Since evaluation is the most computationally expensive step when inducing hypotheses expressed in (fragments of) first-order logic, an appropriate choice of representation for observations can help speeding up this step. In \mathcal{AL} -QUIN the extensional part of Π is partitioned into portions \mathcal{A}_i each of which refers to an individual a_i of C_{ref} . The link between \mathcal{A}_i and a_i is represented with the DATALOG literal $q(a_i)$. The pair $(q(a_i), \mathcal{A}_i)$ is called *observation*.

Example 5. By assuming `MiddleEastCountry` as reference concept, the observation \mathcal{A}_{ARM} contains DATALOG facts such as

```
language('ARM', 'Armenian', 96).
language('ARM', 'Russian', 2).
```

concerning the individual 'ARM' whereas \mathcal{A}_{IR} consists of facts like

```
language('IR', 'Turkish', 1).
language('IR', 'Kurdish', 9).
language('IR', 'Baloch', 1).
language('IR', 'Arabic', 1).
language('IR', 'Luri', 2).
language('IR', 'Persian', 58).
language('IR', 'Turkic', 26).
```

related to the individual 'IR'.

In ILP the coverage test must take the background knowledge into account. The portion \mathcal{K} of \mathcal{B} which encompasses the whole Σ and the intensional part (IDB) of Π is considered as *background knowledge* for \mathcal{AL} -QUIN. Therefore proving that an \mathcal{O} -query Q covers an observation $(q(a_i), \mathcal{A}_i)$ w.r.t. \mathcal{K} equals to proving that $\theta_i = \{X/a_i\}$ is a correct answer to Q w.r.t. $\mathcal{B}_i = \mathcal{K} \cup \mathcal{A}_i$.

Example 6. Checking whether Q_2 covers the observation $(q('ARM'), \mathcal{A}_{ARM})$ w.r.t. \mathcal{K}_{CIA} is equivalent to answering the query

$$Q_2^{(0)} = \leftarrow q('ARM')$$

w.r.t. $\mathcal{K}_{CIA} \cup \mathcal{A}_{ARM} \cup Q_2$. The coverage test for $(q('IR'), \mathcal{A}_{IR})$ is analogous.

A common practice in ILP is to use a reformulation operator, called *saturation* [13], to speed-up the coverage test. It enables ILP systems to make background knowledge explicit within the observations instead of implicit and apart from the observations. In the following we will discuss the implementation of the coverage test in \mathcal{AL} -QUIN and clarify the role of Pellet in supporting the saturation of observations w.r.t. a OWL-DL background knowledge Σ .

3.2 Saturation and instance retrieval

\mathcal{AL} -QUIN is implemented with Prolog as usual in ILP. Thus, the *actual* representation language in \mathcal{AL} -QUIN is a kind of DATALOG^{OI} [15], i.e. the subset of

DATALOG[≠] equipped with an equational theory that consists of the axioms of Clark's Equality Theory augmented with one rewriting rule that adds *inequality atoms* $s \neq t$ to any $P \in \mathcal{L}$ for each pair (s, t) of distinct terms occurring in P . Note that concept assertions are rendered as *membership atoms*, e.g. $a : C$ becomes $c.C(a)$.

Example 7. The following query

```
q(X) ← c.MiddleEastCountry(X), believes(X,Y), c.MonotheisticReligion(Y),
      believes(X,Z), Y≠Z
```

is the DATALOG^{OI} rewriting of:

```
q(X) ← believes(X,Y), believes(X,Z) &
      X:MiddleEastCountry, Y:MonotheisticReligion
```

where the absence of a \mathcal{ALC} constraint for the variable Z explains the need for the inequality atom.

When implementing the coverage test in \mathcal{AL} -QUIN, the goal has been to reduce constrained SLD-resolution of \mathcal{AL} -log to SLD-resolution on DATALOG^{OI}. A crucial issue in this mapping is to deal with the satisfiability tests of \mathcal{ALC} constraints w.r.t. Σ which are required by constrained SLD-resolution because they are performed by applying the tableau calculus for \mathcal{ALC} . The reasoning on the constraint part of \mathcal{O} -queries has been replaced by preliminary saturation steps of the observations w.r.t. the background knowledge Σ . By doing so, the observations are completed with concept assertions that can be derived from Σ by posing *instance retrieval* problems to a DL reasoner. Here, the retrieval is called *levelwise* because it follows the layering of \mathcal{T} : individuals of concepts belonging to the l -th layer \mathcal{T}^l of \mathcal{T} are retrieved all together. Conversely the retrieval for the *reference concept* is made only once at the beginning of the whole discovery process because it makes explicit knowledge of interest to all the levels of granularity. This makes SLD-refutations of queries in \mathcal{L}^l work only on extensional structural knowledge at the level l of description granularity.

A Java application, named OWL2DATALOG, has been developed to support the saturation of observations w.r.t. a OWL-DL background knowledge Σ in \mathcal{AL} -QUIN. To achieve this goal, it supplies the following functionalities:

- levelwise retrieval w.r.t. Σ
- DATALOG^{OI} rewriting of (asserted and derived) concept assertions of Σ

Note that the former is implemented by a client for the DIG server Pellet.

Example 8. The DATALOG^{OI} rewriting of the concept assertions derived for \mathcal{T}^2 produces facts like:

```
c_AfroAsiaticLanguage('Arabic').
...
c_IndoEuropeanLanguage('Armenian').
...
c_MonotheisticReligion('ShiaMuslim').
...
```

to be considered during coverage tests of \mathcal{O} -queries in \mathcal{L}^2 .

The concept assertions, once translated to DATALOG^{OI} , are added to the facts derived from the IDB of Π at the loading of each observation. The coverage test therefore concerns DATALOG^{OI} rewritings of both \mathcal{O} -queries and saturated observations.

Example 9. The DATALOG^{OI} rewriting

$q(X) \leftarrow c_MiddleEastCountry(X), \text{speaks}(X, Y), c_IndoEuropeanLanguage(Y)$

of Q_2 covers the DATALOG^{OI} rewriting:

$c_MiddleEastCountry('ARM')$.
 $\text{speaks}('ARM', 'Armenian')$.
 \dots
 $c_IndoEuropeanLanguage('Armenian')$.
 \dots

of the saturated observation $\hat{\mathcal{A}}_{ARM}$.

Note that the translation from OWL-DL to DATALOG^{OI} is possible because we assume that *all* the concepts are named. This means that an equivalence axiom is required for each complex concept in the knowledge base. Equivalence axioms help keeping concept names (used within constrained DATALOG clauses) independent from concept definitions.

4 Final remarks

In this paper we have shown how to exploit DL reasoners to make existing ILP systems compliant with the latest developments in Ontological Engineering. We would like to emphasize that \mathcal{AL} -QUIN was originally conceived to deal with background knowledge in the form of \mathcal{ALC} taxonomic ontologies but the implementation of this feature was still lacking⁴. Therefore, Pellet makes \mathcal{AL} -QUIN fulfill its design requirements. More precisely, the instance retrieval problems solved by Pellet support the saturation phase in \mathcal{AL} -QUIN. Saturation then compiles DL-based background knowledge down to the usual DATALOG -like formalisms of ILP systems. In this respect, the pre-processing method proposed in [8] to enable legacy ILP systems to work within the framework of CARIN [9] is related to ours but it lacks an implementation. Analogously, the method proposed in [7] for translating OWL-DL to disjunctive DATALOG is far too general with respect to the specific needs of our application. Rather, the proposal of interfacing existing reasoners to combine ontologies and rules [1] is more similar to ours in the spirit. For the future we intend to compare \mathcal{AL} -QUIN with other ILP systems able to deal with ontological background knowledge as soon as they are implemented and deployed.

⁴ \mathcal{AL} -QUIN could actually deal only with concept hierarchies in DATALOG^{OI} .

References

1. U. Assmann, J. Henriksson, and J. Maluszynski. Combining safe rules and ontologies by interfacing of reasoners. In J.J. Alferes, J. Bailey, W. May, and U. Schwerter, editors, *Principles and Practice of Semantic Web Reasoning*, volume 4187 of *Lecture Notes in Computer Science*, pages 33–47. Springer, 2006.
2. T. Berners-Lee, J. Hendler, and O. Lassila. The Semantic Web. *Scientific American*, May, 2001.
3. S. Ceri, G. Gottlob, and L. Tanca. *Logic Programming and Databases*. Springer, 1990.
4. F.M. Donini, M. Lenzerini, D. Nardi, and A. Schaerf. \mathcal{AL} -log: Integrating Datalog and Description Logics. *Journal of Intelligent Information Systems*, 10(3):227–252, 1998.
5. A. Gómez-Pérez, M. Fernández-López, and O. Corcho. *Ontological Engineering*. Springer, 2004.
6. I. Horrocks, P.F. Patel-Schneider, and F. van Harmelen. From \mathcal{SHIQ} and RDF to OWL: The making of a web ontology language. *Journal of Web Semantics*, 1(1):7–26, 2003.
7. U. Hustadt, B. Motik, and U. Sattler. Reducing \mathcal{SHIQ} -description logic to disjunctive datalog programs. In D. Dubois, C.A. Welty, and M.-A. Williams, editors, *Principles of Knowledge Representation and Reasoning: Proceedings of the Ninth International Conference (KR2004)*, pages 152–162. AAAI Press, 2004.
8. J.-U. Kietz. Learnability of description logic programs. In S. Matwin and C. Sammut, editors, *Inductive Logic Programming*, volume 2583 of *Lecture Notes in Artificial Intelligence*, pages 117–132. Springer, 2003.
9. A.Y. Levy and M.-C. Rousset. Combining Horn rules and description logics in CARIN. *Artificial Intelligence*, 104:165–209, 1998.
10. F.A. Lisi. Data Mining in Hybrid Languages with ILP. In D. Calvanese, G. De Giacomo, and E. Franconi, editors, *Proc. 2003 International Workshop on Description Logics*. <http://SunSITE.Informatik.RWTH-Aachen.de/Publications/CEUR-WS/Vol-81/lisi.pdf>, 2003.
11. S.-H. Nienhuys-Cheng and R. de Wolf. *Foundations of Inductive Logic Programming*, volume 1228 of *Lecture Notes in Artificial Intelligence*. Springer, 1997.
12. D. Page and A. Srinivasan. ILP: A short look back and a longer look forward. *Journal of Machine Learning Research*, 4:415–430, 2003.
13. C. Rouveirol. Flattening and saturation: Two representation changes for generalization. *Machine Learning*, 14(1):219–232, 1994.
14. M. Schmidt-Schauss and G. Smolka. Attributive concept descriptions with complements. *Artificial Intelligence*, 48(1):1–26, 1991.
15. G. Semeraro, F. Esposito, D. Malerba, N. Fanizzi, and S. Ferilli. A logic framework for the incremental inductive synthesis of Datalog theories. In N.E. Fuchs, editor, *Proceedings of 7th International Workshop on Logic Program Synthesis and Transformation*, volume 1463 of *Lecture Notes in Computer Science*, pages 300–321. Springer, 1998.
16. E. Sirin, B. Parsia, B. Cuenca Grau, A. Kalyanpur, and Y. Katz. Pellet: A practical OWL-DL reasoner. *Journal of Web Semantics*, 2006.