

Partial Matchmaking using Approximate Subsumption

Heiner Stuckenschmidt

KR and KM Research Group
University of Mannheim
A5, 6 68159 Mannheim, Germany
heiner@informatik.uni-mannheim.de

Abstract

Description Logics, and in particular the web ontology language OWL has been proposed as an appropriate basis for computing matches between structured objects for the sake of information integration and service discovery. A drawback of the direct use of subsumption as a matching criterion is the inability to compute partial matches and qualify the degree of mismatch. In this paper, we describe a method for overcoming these problems that is based on approximate logical reasoning. In particular, we approximate the subsumption relation by defining the notion of subsumption with respect to a certain subset of the concept and relation names. We present the formal semantics of this relation, describe a sound and complete algorithm for computing approximate subsumption and discuss its application to matching tasks.

Introduction

Description Logics are becoming more and more popular as a formalism for representing and reasoning about conceptual knowledge in different areas such as databases and semantic web technologies. In particular, subsumption reasoning for expressive ontologies has been used to compute matches between conceptual descriptions in the context of different real world tasks including information integration (Stuckenschmidt and van Harmelen 2004), product and service matching (Li and Horrocks 2004) and data retrieval (Bechhofer *et al.* 2005). In practical situations, however, it often turns out that logical reasoning is inadequate in many cases, because it does not leave any room for *partial matches*.

Recently, there are some efforts that try to address this problem by combining description logics with numerical techniques for uncertain reasoning in OWL, in particular with techniques for probabilistic (Giugno and Lukasiewicz 2002) and fuzzy reasoning (Straccia 2005). These approaches are able to compute partial matches by assigning an assessment of the degree of matching to the subsumption relation. This degree of matching normally is a real number or an interval between zero and one and therefore allows some ordering of the solutions. Although, in principle this is a solution to the problem of computing the best partial

match but defining an interpreting numerical assessments of uncertainty is a difficult problem. Further, the reduction to a single numerical assessment of the mismatch does not allow different users to discriminate between different kinds of mismatches.

In this paper, we propose a notion of approximate subsumption that supports the computation of partial matches between complex concept expressions without relying on a single number to represent the degree of mismatch. Instead, *the approach describes the degree of matching in terms of a subset of the aspects of the request that are met by the solution*. This approach allows the user to decide whether to accept a partial match based on whether important aspects are missed or not. In order to implement this approach we borrow from the area of approximate deduction. In particular, we extend the notion of S-Interpretations of propositional logic proposed in (Schaerf and Cadoli 1995) to description logics and use the result notion of a non-standard interpretation of concept expressions to define an approximate subsumption operator that computes subsumption with respect to a particular subset of the vocabulary used.

Our approach is similar to the notion of approximate entailment for description logics proposed in (Cadoli and Schaerf 1992). Our work extends this work different ways:

- Previous work was restricted to rather inexpressive description logics, in particular *ALC*. We extend this to expressive description logics. In particular, our approach includes qualified number restrictions.
- We provide a more natural way of approximating concept descriptions based on the set of concept and relation names.
- While the work of Cadoli and Schaerf relied on a rather complicated formalization in terms of the Herbrand Universe of the first-order translation of description logics, we provide a direct model theoretic semantics for approximate subsumption and show that it has the required properties¹

Approximation based on Sub-Vocabularies

In propositional logic, the vocabulary of a formula consists of a set of propositional letters. A formula consists of a

Copyright © 2007, American Association for Artificial Intelligence (www.aaai.org). All rights reserved.

¹Proofs for the theorems are omitted due to lack of space.

Boolean expression over these letters. A classical interpretation I assigns to each letter either the value *true* or *false*. The semantics of negation now implies that a letter and its negation cannot have the same truth value, in particular, for all propositional letters p one of the following :

$$\begin{aligned} I(p \wedge \neg p) &= false \\ I(p \vee \neg p) &= true \end{aligned} \quad (1)$$

Checking satisfiability of a formula relies on showing that there is no assignment of truth values that satisfies this condition and makes the whole formula true. A possible way for approximating satisfiability testing for propositional logic is now to restrict the condition above to a subset of the propositional letters. This subset is denoted as S and the corresponding interpretation is called an S-interpretation of the formula (Schaerf and Cadoli 1995).

Depending on how the letters not in S are treated, an S-Interpretation is sound or complete with respect to the classical interpretation. One kind of non-standard interpretation called S-3 Interpretation assigns both, a letter and its negation to *true*.

$$I(p \wedge \neg p) = true, p \notin S \quad (2)$$

When applying this interpretation to the satisfiability problem, we observe that formulas that were unsatisfiable before now become satisfiable. This means that the resulting calculus is sound, but incomplete, because some results that could be proven using the principle of proof by refutation can not be proven any more, because the conjunction of the knowledge base with the negation of the result to be proven becomes satisfiable under the new interpretation. The counterpart of S-3 interpretation are S-1 Interpretations that assign *false* to both a letters and their negation if the letters are not in the set S .

$$I(p \vee \neg p) = false, p \notin S \quad (3)$$

Following the same argument as above, S-1 Interpretations define a complete but unsound calculus for propositional logic. In both cases, the advantage of the approach is that we can decide which parts of the problem to approximate by selecting an appropriate set of letters S . Therefore the approach provides a potential solution to the problem of partial matching described above.

The idea of our approach is now to apply the underlying idea of S-Interpretations to the Description Logic \mathcal{SHIQ} which covers most of the expressive power of OWL in order to support approximate subsumption reasoning where parts of the vocabulary are interpreted in the classical way and other parts are approximated. In fact, Cadoli and Schaerf do propose an extension of S-Interpretations to Description logics, but they define S not in terms of a subset of the vocabulary, but in terms of the structure of the concept expression (Cadoli and Schaerf 1992). In (Groot *et al.* 2005) it has been shown that this way of applying S-Interpretations to description logics does not produce satisfying results on real data. In this paper, we therefore propose an alternative way of defining S-Interpretations for description logics which is closer to the notion of S-Interpretations in propositional logic. The idea is to interpret description logics as an extension of propositional logic, where class names correspond to propositional letters². As for propositional logic, we select a subset of the class names that is interpreted in the

classical way and approximate class names not in this set. In particular, a classical interpretation $(\Delta^{\mathcal{I}}, \mathcal{I})$ of class names requires that a concept name and its negation form a disjoint partition of the domain:

$$\begin{aligned} C^{\mathcal{I}} \cap (\neg C)^{\mathcal{I}} &= \emptyset \\ C^{\mathcal{I}} \cup (\neg C)^{\mathcal{I}} &= \Delta^{\mathcal{I}} \end{aligned} \quad (4)$$

We can now define approximations for description logics by relaxing these requirements for a subset of the concept names. The corresponding S-3 and S-1 Interpretations are very similar to the ones for propositional logic. In particular, for S-3 Interpretations we have.

$$C^{\mathcal{I}} \cap (\neg C)^{\mathcal{I}} = \Delta^{\mathcal{I}}, C \notin S \quad (5)$$

This means that both, C and $\neg C$ are mapped to $\Delta^{\mathcal{I}}$ by the interpretation. As a consequence, the concept name C cannot cause a clash in a tableaux proof and therefore, constraints that force a certain value to be of type C will be ignored in a subsumption proof. The resulting subsumption operator is sound, but incomplete. For S-1 Interpretations, we have

$$C^{\mathcal{I}} \cup (\neg C)^{\mathcal{I}} = \emptyset, C \notin S \quad (6)$$

which means that both C and $\neg C$ are mapped to the empty set. In a tableaux proof, all attempts to construct a model that involves a variable of type C will fail. The corresponding subsumption operator is complete, but unsound with respect to classical subsumption.

While approximation based on concept names is a straightforward application of the notion of S-1 and S-3 interpretations, things become more complicated if we want to extend the approach to relation names. In Description Logics relations are used to formulate constraints that apply to all members of a certain class. The most general formulation of these constraints is in terms of qualified number restrictions. Qualified number restrictions have the following form $(\leq nr.C)$ or $(\geq nr.C)$ where n is a positive natural number (including zero), r is the name of a binary relation and C is a concept expression. In a tableaux these qualified number restrictions are a second potential source of inconsistency besides the negation operator. In particular, we have

$$(\leq nr.C)^{\mathcal{I}} \cap (\geq mr.C)^{\mathcal{I}} = \emptyset \text{ for all } n < m$$

on the other hand, we have

$$(\leq nr.C)^{\mathcal{I}} \cup (\geq mr.C)^{\mathcal{I}} = \Delta^{\mathcal{I}} \text{ for all } n \geq m$$

We can use this analogy to extend the notion of S-1 and S-3 interpretations to qualified number restrictions in the following way. For S-3 Interpretations we define that

$$(\leq nr.C)^{\mathcal{I}} \cup (\geq mr.C)^{\mathcal{I}} = \Delta^{\mathcal{I}} \text{ for all } r \notin S \quad (7)$$

In particular, we weaken the condition for the expression to become the universal concept by making it independent of the values for m and n . Further, we claim that the conjunction of qualified number expressions can never be the empty concept, i.e.

$$(\leq nr.C)^{\mathcal{I}} \cap (\geq mr.C)^{\mathcal{I}} \neq \emptyset \text{ for all } r \notin S \quad (8)$$

This leaves us with a weaker interpretation, because inconsistencies arising from the relations not in the set S cannot be detected. For S-1 interpretations, we make analogous claim by demanding that the union of two qualified number restrictions can never be the universal concept

$$(\leq nr.C)^{\mathcal{I}} \cup (\geq mr.C)^{\mathcal{I}} \neq \Delta^{\mathcal{I}} \text{ for all } r \notin S \quad (9)$$

²In fact, a description logic that just contains the Boolean operators is equivalent to propositional logic.

Further, we strengthen the interpretation by claiming that the intersection of the two qualified number restrictions on the same relation and concept is always inconsistent

$$(\leq n r.C)^{\mathcal{I}} \cap (\geq m r.C)^{\mathcal{I}} = \emptyset \text{ for all } r \notin S \quad (10)$$

This gives us a stronger version of the semantics, because any two assertions using this relation in combination with the same concept expression C leads to an inconsistency³. The result is a complete but unsound subsumption operator. This unsound approximation operator is exactly what we need for specifying the notion of a partial match, because it forces a match on the constraints involving class names from S and treats constraints involving classes not in S as optional. Using subsumption operators with different sets S , we can focus on different aspects of the matching task and also rank results based on the number of requirements met. In the following, we will therefore concentrate on complete, but unsound approximations of subsumption reasoning for concept expressions based on the idea described above. In particular, we will formally specify non-standard interpretations and define a family of approximate subsumption operators that can be used to compute partial matches.

Non-Standard Semantics

In the following, we introduce a non-standard interpretation for concept expressions in the logic \mathcal{SHIQ} . Details about the language can be found in (Tobies 2001). A limited vocabulary is a subset $S \subseteq \mathcal{V}$ of the concept and relation names occurring in a concept expression. Our aim is to define approximate reasoning in Description Logics based on such a subset of the vocabulary. For this purpose, we define an upper and a lower approximation of an interpretation \mathcal{I} with respect to a set S referred to as \mathcal{I}_S^+ and \mathcal{I}_S^- respectively. We call \mathcal{I}_S^+ an upper approximation and \mathcal{I}_S^- a lower approximation of \mathcal{I} with respect to S .

Definition 1 (Lower Approximation) *A lower approximation of an interpretation \mathcal{I} with respect to S is a non standard interpretation $(\Delta^{\mathcal{I}}, \mathcal{I}_S^-)$ such that:*

$$\begin{aligned} A^{\mathcal{I}_S^-} &= \begin{cases} A^{\mathcal{I}} & A \in S \\ \emptyset & \text{otherwise} \end{cases} \\ (\neg C)^{\mathcal{I}_S^-} &= \Delta^{\mathcal{I}} - C^{\mathcal{I}_S^+} \\ (C \sqcap D)^{\mathcal{I}_S^-} &= C^{\mathcal{I}_S^-} \cap D^{\mathcal{I}_S^-} \\ (C \sqcup D)^{\mathcal{I}_S^-} &= C^{\mathcal{I}_S^-} \cup D^{\mathcal{I}_S^-} \\ (\geq n r.C)^{\mathcal{I}_S^-} &= \begin{cases} \{x | \#\{y.(x, y) \in r^{\mathcal{I}} \wedge y \in C^{\mathcal{I}_S^-}\} \geq n\} & r \in S \\ \{x | \#\{y.(x, y) \in r^{\mathcal{I}} \wedge y \in C^{\mathcal{I}_S^-}\} \geq \infty\} & r \notin S \end{cases} \\ (\leq n r.C)^{\mathcal{I}_S^-} &= \begin{cases} \{x | \#\{y | (x, y) \in r^{\mathcal{I}} \wedge y \in C^{\mathcal{I}_S^+}\} \leq n\} & r \in S \\ \{x | \#\{y | (x, y) \in r^{\mathcal{I}} \wedge y \in C^{\mathcal{I}_S^+}\} \leq 0\} & r \notin S \end{cases} \end{aligned}$$

where $(\Delta^{\mathcal{I}}, \mathcal{I}_S^+)$ is an upper approximation as defined in definition 2

Definition 2 (Upper Approximation) *An upper approximation of an interpretation \mathcal{I} with respect to S is a non stan-*

³As we will see later, it is sufficient if the two restrictions use concept expressions that are logically equivalent

ard interpretation $(\Delta^{\mathcal{I}}, \mathcal{I}_S^+)$ such that:

$$\begin{aligned} A^{\mathcal{I}_S^+} &= \begin{cases} A^{\mathcal{I}} & A \in S \\ \Delta^{\mathcal{I}} & \text{otherwise} \end{cases} \\ (\neg C)^{\mathcal{I}_S^+} &= \Delta^{\mathcal{I}} - C^{\mathcal{I}_S^-} \\ (C \sqcap D)^{\mathcal{I}_S^+} &= C^{\mathcal{I}_S^+} \cap D^{\mathcal{I}_S^+} \\ (C \sqcup D)^{\mathcal{I}_S^+} &= C^{\mathcal{I}_S^+} \cup D^{\mathcal{I}_S^+} \\ (\geq n r.C)^{\mathcal{I}_S^+} &= \begin{cases} \{x | \#\{y.(x, y) \in r^{\mathcal{I}} \wedge y \in C^{\mathcal{I}_S^+}\} \geq n\} & r \in S \\ \{x | \#\{y.(x, y) \in r^{\mathcal{I}} \wedge y \in C^{\mathcal{I}_S^+}\} > 0\} & r \notin S \end{cases} \\ (\leq n r.C)^{\mathcal{I}_S^+} &= \begin{cases} \{x | \#\{y | (x, y) \in r^{\mathcal{I}} \wedge y \in C^{\mathcal{I}_S^-}\} \leq n\} & r \in S \\ \{x | \#\{y | (x, y) \in r^{\mathcal{I}} \wedge y \in C^{\mathcal{I}_S^-}\} < \infty\} & r \notin S \end{cases} \end{aligned}$$

where $(\Delta^{\mathcal{I}}, \mathcal{I}_S^-)$ is a lower approximation as defined in definition 1

Negation normal forms play an important role with respect to tableaux-based algorithms for description logics. The ability to compute the negation normal form as a basis for a tableaux proof relies on the following transformation rules.

Lemma 1 (Equivalent Transformations) *The following equivalence rules also hold under the non-standard interpretations \mathcal{I}_S^+ and \mathcal{I}_S^- .*

$$\neg(C \sqcap D) \equiv (\neg C) \sqcup (\neg D) \quad (11)$$

$$\neg(C \sqcup D) \equiv (\neg C) \sqcap (\neg D) \quad (12)$$

$$\neg(\geq n r.C) \equiv (\leq n - 1 r.C) \quad (13)$$

$$\neg(\leq n r.C) \equiv (\geq n + 1 r.C) \quad (14)$$

The fact that these equivalences also hold under the non-standard interpretations implies that we can translate every concept expression into its negation normal form without changing the non-standard interpretation. This result is formalized in the following theorem.

Corollary 1 (Negation Normal Form) *For every concept expression C there is an expression $nnf(C)$ in negation normal form such that $nnf(C)^{\mathcal{I}_S^-} = C^{\mathcal{I}_S^-}$ and $nnf(C)^{\mathcal{I}_S^+} = C^{\mathcal{I}_S^+}$*

Based on the negation normal form, we can define a simplified version of the semantics with respect to negation. Instead of the general definition of negation, we can use a special rule for negation with respect to negation of atomic concept names. In particular, for a concept expression in negation normal form equation 11 can be replaced by the following equation:

$$(\neg A)^{\mathcal{I}_S^-} = \begin{cases} A^{\mathcal{I}} & A \in S \\ \emptyset & \text{otherwise} \end{cases} \quad (15)$$

Analogously, equation 11 can be replaced by the following equation:

$$(\neg A)^{\mathcal{I}_S^+} = \begin{cases} A^{\mathcal{I}} & A \in S \\ \Delta^{\mathcal{I}} & \text{otherwise} \end{cases} \quad (16)$$

The main implication of the existence of an equivalent negation normal form is the possibility to define inference mechanisms that work on this normal form. In particular, this allows us to use tableaux-style proof procedures to determine the satisfiability of a concept expression under the non-standard semantics.

From the negation normal form it is easy to see that the upper and lower approximation of the interpretation satisfies

equations 5 and 6. We can also easily show that the other intuitions about the generalization of S-Interpretations are satisfied.

Theorem 1 (Properties of approximate interpretations)
The upper approximation of an interpretation \mathcal{I} with respect to a sub-vocabulary S satisfies equations 5, 7 and 8. The lower approximation satisfies equations 6, 9 and 10.

Another useful property of the non standard interpretation is that it makes concept expressions strictly more general for upper and strictly more specific for lower approximations. This property which we call monotonicity is important in order to be able to guarantee formal properties of approximation methods defined based on this interpretation. Therefore the following theorem describes a central property of approximation in description logics.

Lemma 2 (Monotonicity) *Given a non-standard interpretation as defined above, the following equation holds for all concept expressions C :*

$$c^{\mathcal{I}_{\overline{S}}} \subseteq c^{\mathcal{I}} \subseteq c^{\mathcal{I}_{\overline{S}^+}} \quad (17)$$

We can generalize the theorem by observing that the standard interpretation is an extreme case of the non-standard interpretation with $S = \mathcal{V}$. In particular, the general version of monotonicity says that for upper approximations removing names from the set S will make concepts expressions strictly more general. Conversely, for lower approximations concept expressions become less general when we remove concept or relation names from the set S . The corresponding general property is defined in the following theorem:

Lemma 3 (Generalized Monotonicity) *Given a non-standard interpretation as defined above and two sub-vocabularies S_1 and S_2 with $S_1 \subseteq S_2$, the following equations hold for all concept expressions C :*

$$c^{\mathcal{I}_{\overline{S_1}}} \supseteq c^{\mathcal{I}_{\overline{S_2}}} \quad c^{\mathcal{I}_{\overline{S_1}^+}} \subseteq c^{\mathcal{I}_{\overline{S_2}^+}} \quad (18)$$

The generalized monotonicity property is interesting, because it allows us to successively compute more precise upper and lower approximations of a concept by adding names to the set S . This is convenient in cases where users provide a preference order over the vocabulary indicating the relative importance of different aspects of a concept. In this case, use the preference relation provided by the user to determine a sequence of approximations to be used in the matching process.

An Approximate Subsumption Operator

Up to now, we have only considered interpretations as such. As our aim is to develop approximate notions of subsumption as a basis for approximate matching, we now have to define the notion of approximate subsumption based on the non-standard interpretation defined above. It turns out, that this can be done in a straightforward way using the standard definition of the subsumption operator as:

$$\forall \mathcal{I} : \mathcal{I} \models C \sqsubseteq D \Leftrightarrow (C \sqcap \neg D)^{\mathcal{I}} = \emptyset$$

The idea is now to use this definition and replace the standard interpretation \mathcal{I} by the lower approximation $\mathcal{I}_{\overline{S}}$ with

respect to a certain sub-vocabulary S . Based on the choice of S , this defines different subsumption operators with certain formal properties that will be discussed in the following.

Definition 3 (Approximate Subsumption) *Let $S \subseteq \mathcal{V}$ be a subset of the concept names and $(\Delta^{\mathcal{I}}, \mathcal{I}_{\overline{S}}^-)$ a lower approximation, then the corresponding approximate subsumption relation $\sqsubseteq_{\overline{S}}$ is defined as follows*

$$\forall \mathcal{I} : \mathcal{I} \models (C \sqsubseteq_{\overline{S}} D) \Leftrightarrow_{def} (C \sqcap \neg D)^{\mathcal{I}_{\overline{S}}^-} = \emptyset \quad (19)$$

We say that C is subsumed by D with respect to sub-vocabulary S .

The monotonicity of the non-standard interpretation has an impact on the formal properties of the approximate subsumption operator. In particular, we can establish a relation between the subset of the vocabulary considered and the strength of the subsumption operator. The more concepts we exclude from the set S the weaker the subsumption operator as well as the matches we can compute get. This implies that if we can prove subsumption with respect to a particular set S the subsumption relation also holds for all subsets of S . Conversely, if we fail to prove subsumption with respect to a set S , we can be sure that the subsumption relation does also not hold with respect to any superset of S . These properties are stated formally in the following theorem.

Theorem 2 (Properties of Approximate Subsumption)
Let f be a lower approximation, then the following equation holds:

$$(C \sqsubseteq_{\overline{S_2}} D) \Rightarrow (C \sqsubseteq_{\overline{S_1}} D) \text{ for } S_1 \subseteq S_2 \quad (20)$$

$$(C \not\sqsubseteq_{\overline{S_1}} D) \Rightarrow (C \not\sqsubseteq_{\overline{S_2}} D) \text{ for } S_1 \subseteq S_2 \quad (21)$$

These properties allow us to develop approximation strategies by successively selecting smaller subsets of concepts to be considered for matching and trying to compute the corresponding subsumption relation until we succeed.

Applying the Approximation

A nice feature of our approach is that it can actually be implemented by simply performing syntactic modifications on concept expressions. In particular, in order to check whether a statement $C \sqsubseteq_{\overline{S}} D$ holds, we take the expression $(C \sqcap \neg D)$

and transform it into a concept expression that simulates the non-standard interpretation. For the lower approximation, the corresponding transformation $(\cdot)^-$ is defined as follows

$$\begin{aligned} (A)^- &\rightarrow \perp \text{ if } A \in S \\ (\neg A)^- &\rightarrow \perp \text{ if } A \in S \\ (\neg C)^- &\rightarrow \neg(C)^+ \\ (C \sqcap D)^- &\rightarrow (C)^- \sqcap (D)^- \\ (C \sqcup D)^- &\rightarrow (C)^- \sqcup (D)^- \\ (\leq n r.C)^- &\rightarrow (\leq 0 r.(C)^+) \text{ if } r \in S \\ (\leq n r.C)^- &\rightarrow (\leq n r.(C)^+) \text{ if } r \notin S \\ (\geq n r.C)^- &\rightarrow (\geq \max r.(C)^-) \text{ if } r \in S \\ (\geq n r.C)^- &\rightarrow (\geq n r.(C)^-) \text{ if } r \notin S \end{aligned}$$

Here \max is an integer number that is larger than any other number occurring in any qualified number restriction

in the concept expression. This is sufficient to model the interpretation that requires less than an infinite number of r-successors. Analogously, we define a transformation function $(\cdot)^+$ that creates a concept expression that simulates the upper approximation of a concept expression. This transformation is defined as follows:

$$\begin{aligned}
(A)^+ &\rightarrow \top \text{ if } A \in S \\
(\neg A)^+ &\rightarrow \top \text{ if } A \in S \\
(\neg C)^+ &\rightarrow \neg(C)^- \\
(C \sqcap D)^+ &\rightarrow (C)^+ \sqcap (D)^+ \\
(C \sqcup D)^+ &\rightarrow (C)^+ \sqcup (D)^+ \\
(\leq n r.C)^+ &\rightarrow (\leq \max - 1 r.(C)^-) \text{ if } r \in S \\
(\leq n r.C)^+ &\rightarrow (\leq n r.(C)^-) \text{ if } r \notin S \\
(\geq n r.C)^+ &\rightarrow (\geq 1 r.(C)^+) \text{ if } r \in S \\
(\geq n r.C)^+ &\rightarrow (\geq n r.(C)^+) \text{ if } r \notin S
\end{aligned}$$

We again use the number \max for modeling an infinite number of r-successors. Further, we have to use the condition ≥ 1 instead of < 0 which is equivalent. It can be shown that these rewriting rules provide a way for computing approximate subsumption as stated by the following theorem.

Theorem 3 (Syntactic approximation I) *Let C and D be concept expressions in $SHIQ$, then $\mathcal{I} \models C \sqsubseteq_S D$ if and only if $(C \sqcap \neg D)^-$ is unsatisfiable.*

It turns out that the equivalence of a concept expression and its normal form and the symmetry of upper and lower approximation with respect to negation can be used to define an alternative way of computing approximate subsumption based on the syntactic manipulations shown above

Theorem 4 (Syntactic approximation II) *Let C and D be concept expressions in $SHIQ$, then $\mathcal{I} \models C \sqsubseteq_S D$ if and only if $\mathcal{I} \models (C)^- \sqsubseteq (D)^+$*

This means that we have two rather straightforward ways of computing approximate subsumption using standard DL reasoners.

Information Integration

Our first example of the potential use of the approximate subsumption operator for partial Matchmaking is taken from chapter 4 of (Stuckenschmidt and van Harmelen 2004). The problem addressed is the integration of different land-use classification schemes (ATKIS and CORINE) for supporting the automatic update of official registry records with satellite image data. As an example, we take the land-use class 'Mixed-Forest' from the ATKIS catalogue which is defined as a region that has a vegetation composed of coniferous and broad-leaved plants that are all trees or shrubs. The corresponding concept expression is the following⁴.

$$\begin{aligned}
Mixed - Forest &\equiv Region \sqcap \\
&\geq 1 \text{ vegetation.Magnoliophyta} \sqcap \\
&\geq 1 \text{ vegetation.Coniferophyta} \sqcap \\
&\leq 0 \text{ vegetation.}\neg(\text{trees} \sqcup \text{shrubs})
\end{aligned}$$

⁴We already transformed existential and universal quantifiers into qualified number restrictions

When matching this description with the CORINE classification using standard subsumption, the most specific CORINE concept that subsumes Mixed Forest is the concept Vegetation Area which is defined by the presence of some vegetation:

$$Vegetation \equiv Region \sqcap (\geq 1 \text{ vegetation.}\top)$$

This result is somewhat disappointing as there are more specific classes in CORINE that we would have expected to also match. In particular, there is a concept 'Forest' that would qualify as the correct solution to the integration problem from a commonsense point of view. Looking at the definition of the Forest Concept in CORINE reveals that the problem is caused by the fact, that the definition does not mention shrubs as a possible form of vegetation of forest areas

$$Forests \equiv Vegetation \sqcap (\leq 0 \text{ vegetation.}\neg(\text{trees}))$$

Using the algorithm given above, we can show that $Mixed - Forest \sqsubseteq_{V-\{\text{shrubs}\}} Forests$ holds. This also

shows the main advantage of our approach over existing proposals for partial matching: the result of our method clearly states what aspects of the two concepts did not match; in our case the presence of shrubs. The importance of this feature becomes clear when we look at the other concepts in the CORINE classification that are candidates for a match with the concept of mixed forest. In particular, there is a concept that specifies herbaceous and shrub areas. This concept is defined as follows:

$$\begin{aligned}
Herbaceous &\equiv Vegetation \sqcap \\
&(\leq 0 \text{ vegetation.}\neg(\text{shrubs} \sqcup \text{herbs}))
\end{aligned} \tag{22}$$

If we apply our partial matching approach with respect to this concept expression, we can find out that $Mixed - Forest \sqsubseteq_{V-\{\text{trees}\}} Herbaceous$. In order to determine the

best match, we can ask the user which concept she considers to be more important with respect to determining a match for mixed forest. In this case most users will decide that the concept tree is more significant with respect to matching mixed forest than shrub and therefore should not be excluded from S and therefore lead the system to prefer the 'natural' match between mixed forest and forest.

Service Discovery

Our second example for the use of approximate subsumption as a basis for partial matchmaking is the problem of service discovery based on matchmaking between service profiles and service requests. We base our example on the description logic framework for service discovery proposed in (Li and Horrocks 2004) where numerical attributes are encoded as qualified number restrictions.

We assume a request asking for a Sales service that offers PCs or Laptops with at least 512 MB main memory, at least 256 MB Cache Memory and a price of at most 500 Dollars. The corresponding request can be formulated using the following concept expression:

$$\begin{aligned}
request &\equiv Sales \sqcap \\
&(\forall item.(PC \sqcup Laptop \sqcap (\geq 512 \text{ has} - \text{memory.Main}) \\
&\sqcap (\geq 256 \text{ has} - \text{memory.Cache}) \\
&\sqcap (\leq 500 \text{ price.Dollar})))
\end{aligned}$$

We further assume a sales service offering PCs with 256 MB main memory and 256 MB Cache Memory at a price of 450 Dollars and Laptops with 512 MB Main Memory and 256 MB Cache Memory at a price of 650 Dollar. This service can be described using the following concept expression:

$$\begin{aligned} advert1 &\equiv Sales \sqcap \\ &(\forall item.(PC \sqcap (\geq 256 has - memory.Main) \\ &\sqcap (\geq 256 has - memory.Cache) \\ &\sqcap (\leq 450 price.Dollar))) \\ \\ advert2 &\equiv Sales \sqcap \\ &(\forall item.(Laptop \sqcap (\geq 512 has - memory.Main) \\ &\sqcap (\geq 256 has - memory.Cache) \\ &\sqcap (\leq 650 price.Dollar))) \end{aligned}$$

It is easy to see that neither $advert1 \sqsubseteq request$ nor $advert2 \sqsubseteq request$ holds. The both adverts satisfy the condition of being a sales service, it also offers the right kinds of items - PCs or Laptops, but each of the items offered fails to satisfy one of the requirements. While the PCs do not have enough main memory, the Laptops are too expensive. The reasoner is unable to detect a clash in the expression $advert \sqcap \neg request$. In particular, it fails to detect a clash between the expressions $(\leq 511 has - memory.Main)$ and $(\geq 256 has - memory.Main)$ in the case of PCs and between $(\leq 650 price.Dollar)$ and $(\geq 501 price.Dollar)$ for the case of Laptops. Existing approaches for approximate deduction cannot solve this problem, as the problem involves qualified number restrictions. Using our approach, excluding $has - memory$ from the set S leads to a re-formulation of the problem where $(\leq max - 1 has - memory.Main)$ and $(\geq max has - memory.Main)$ are compared. As a consequence

$$advert1 \sqsubseteq_{\nu - \{has - memory\}} request$$

holds. In the second case, using our approach will re-write the problem to $(\leq 0 price.Dollar)$ and $(\geq 1 price.Dollar)$ which also leads to a clash. Thus we also have

$$advert2 \sqsubseteq_{\nu - \{price\}} request$$

for the Laptop case. As in the case of information integration, the subset of the vocabulary chosen provides the user with valuable feedback with respect to the relevance of the different matches. In particular, the user can decide whether the original constraint on the price or on the memory should be relaxed.

Discussion

We presented an approach for computing approximate subsumption between concept expressions in *SHIQ* based on a subset of the vocabulary used in the expressions. The approach solves some of the problems of classical reasoning in description logics, in particular, the inability to accept imperfect matches between concepts without having to leave the realms of formal logic. As a side-effect, the subset of the vocabulary also provides us with a qualitative characterization of the mismatch between the expressions, which is clearly an advantage over numerical approaches for dealing with imperfect matches. An approach for partial matching in description logics that is more similar to ours is reported

in (Di Noia *et al.* 2003). This approach, however, cannot deal with disjunction and qualified number restrictions.

Another advantage of your approach is the fact, that it does not only cover subsumption between concept expressions, but that it also provides us with the possibility to approximate subsumption with respect to a background terminology. In particular, SHIQ allows for the internalization of TBoxes into a single concept using a universal role, which is a transitive super-role of all roles (Tobies 2001). Internalization allows reduce subsumption with respect to general TBoxes and role hierarchies to unsatisfiability of concepts with respect to role hierarchies. This means that we can approximate subsumption reasoning with general TBoxes by approximating the satisfiability of the resulting concept expression.

References

- Sean Bechhofer, Ian Horrocks, and Daniele Turi. The owl instance store: System description. In *Proceedings CADE-20*, Lecture Notes in Computer Science. Springer-Verlag, 2005.
- Marco Cadoli and Marco Schaerf. Approximation in concept description languages. In *Proceedings of the International Conference on Knowledge Representation and Reasoning*, pages 330–341, 1992.
- Tommaso Di Noia, Francesco Donini Eugenio Di Sciascio, and Marina Mongiello. A system for principled match-making in an electronic marketplace. In *Proceedings of the Twelfth International World Wide Web Conference*, pages 321–330, 2003.
- Rosalba Giugno and Thomas Lukasiewicz. P-shoq(d): A probabilistic extension of shoq(d) for probabilistic ontologies in the semantic web. In *Proceedings of the 8th European Conference on Logics in Artificial Intelligence (JELIA'02)*, 2002.
- Perry Groot, Heiner Stuckenschmidt, and Holger Wache. Approximating description logic classification for semantic web reasoning. In *Proceedings of the 2nd European Semantic Web Conference*, Heraklion, Crete, 2005.
- Lei Li and Ian Horrocks. A software framework for match-making based on semantic web technology. *International Journal of Electronic Commerce*, 8(4):39 – 60, 2004.
- Marco Schaerf and Marco Cadoli. Tractable reasoning via approximation. *Artificial Intelligence*, 74(2):249–310, 1995.
- Umberto Straccia. Towards a fuzzy description logic for the semantic web (preliminary report). In *Proceedings of the 2nd European Semantic Web Conference ESWC-05*, pages 167–181, 2005.
- H. Stuckenschmidt and F. van Harmelen. *Information Sharing on the Semantic Web*. Advanced Information Processing. Springer Verlag, Berlin, Heidelberg, 2004. to appear.
- Stephan Tobies. *Complexity results and practical algorithms for logics in Knowledge Representation*. Phd thesis, LuFG Theoretical Computer Science, RWTH-Aachen, Germany, 2001.