

Ontology Learning from Interpretations in Lightweight Description Logics

Szymon Klarman^{1,2}(✉) and Katarina Britz^{2,3}

¹ Department of Computer Science, Brunel University London, Uxbridge, UK
szymon.klarman@gmail.com

² CSIR Centre for Artificial Intelligence Research, Pretoria, South Africa

³ Department of Information Science, Stellenbosch University,
Stellenbosch, South Africa

Abstract. Data-driven elicitation of ontologies from structured data is a well-recognized knowledge acquisition bottleneck. The development of efficient techniques for (semi-)automating this task is therefore practically vital — yet, hindered by the lack of robust theoretical foundations. In this paper, we study the problem of learning Description Logic TBoxes from interpretations, which naturally translates to the task of ontology learning from data. In the presented framework, the learner is provided with a set of positive interpretations (i.e., logical models) of the TBox adopted by the teacher. The goal is to correctly identify the TBox given this input. We characterize the key constraints on the models that warrant finite learnability of TBoxes expressed in selected fragments of the Description Logic \mathcal{EL} and define corresponding learning algorithms.

1 Introduction

In the advent of the Web of Data and various “e-” initiatives, such as e-science, e-health, e-governance, etc., the focus of the classical knowledge acquisition bottleneck becomes ever more concentrated around the problem of constructing rich and accurate ontologies enabling efficient management of the existing abundance of data [1]. Whereas the traditional understanding of this bottleneck has been associated with the necessity of developing ontologies *ex ante*, in a top-down, data-agnostic manner, this seems to be currently evolving into a new position, recently dubbed the knowledge reengineering bottleneck [2]. In this view, the contemporary challenge is to, conversely, enable data-driven approaches to ontology design — methods that can make use and make sense of the existing data, be it readily available on the web or crowdsourced, leading to elicitation of the ontological commitments implicitly present on the data-level. Even though the development of such techniques and tools, which could help (semi-)automate thus characterized ontology learning processes, becomes vital in practice, the robust theoretical foundations for the problem are still rather limited. This work

This work was funded in part by the National Research Foundation under Grant no. 85482.

is an attempt at establishing exactly such foundations and focuses on some key theoretical issues towards this goal.

We study the problem of learning *Description Logic* (DL) TBoxes from interpretations, which naturally translates to the task of ontology learning from data. DLs are a popular family of knowledge representation formalisms [3], which have risen to prominence as, among others, the logics underpinning different profiles of the Web Ontology Language OWL¹. In this paper, we focus on the lightweight DL \mathcal{EL} [4] and some of its more specific fragments. This choice is motivated, on the one hand, by the interesting applications of \mathcal{EL} , especially as the logic behind OWL 2 \mathcal{EL} profile, while on the other, by its relative complexity, which enables us to make interesting observations from the learning perspective. Our learning model is a variant of learning from positive interpretations (i.e., from models of the target theory) — a generally established framework in the field of inductive logic programming [5, 6]. In our scenario, the goal of the learner is to correctly identify the target TBox \mathcal{T} given a finite set of its finite models. Our overarching interest lies in algorithms warranting effective learnability in such setting with no or minimum supervision. Our key research questions and contributions are therefore concerned with the identification of specific languages and conditions on the learning input under which such algorithms can be in principle defined.

In the following two sections, we introduce DL preliminaries and discuss the adopted learning model. In Sect. 4, we identify two interesting fragments of \mathcal{EL} , called $\mathcal{EL}^{\text{rhs}}$ and $\mathcal{EL}^{\text{lhs}}$, which satisfy some basic necessary conditions enabling finite learnability, and at the same time, we show that full \mathcal{EL} does not meet that same requirement. In Sect. 5, we devise a generic algorithm which correctly identifies $\mathcal{EL}^{\text{rhs}}$ and $\mathcal{EL}^{\text{lhs}}$ TBoxes from finite data, employing a basic equivalence oracle. Further, in case of $\mathcal{EL}^{\text{rhs}}$, we significantly strengthen this result by defining an algorithm which makes no such calls to an oracle, and thus supports fully unsupervised learning. In Sect. 6, we compare our work to related contributions, in particular to the framework of learning TBoxes from entailment queries, by Konev et al. [7, 8]. We conclude in Sect. 7 with an overview of interesting open problems.

2 Description Logic Preliminaries

The language of the Description Logic (DL) \mathcal{EL} [4] is given by (1) a vocabulary $\Sigma = (N_C, N_R)$, where N_C is a set of concept names (i.e., unary predicates, e.g., *Father*, *Woman*) and N_R a set of role names (i.e., binary predicates, e.g., *hasChild*, *likes*), and (2) the following set of constructors for defining complex concepts, which shall be divided into two groups:

$$\begin{aligned} \mathcal{EL}: \quad C, D ::= \top \mid A \mid C \sqcap D \mid \exists r.C \\ \mathcal{L}^\square: \quad C, D ::= \top \mid A \mid C \sqcap D \end{aligned}$$

where $A \in N_C$ and $r \in N_R$. Concept \top denotes all individuals in the domain, $C \sqcap D$ the class of individuals that are instances of both C and D , and $\exists r.C$

¹ See <http://www.w3.org/TR/owl2-profiles/>.

describes all individuals that are related to some instance of C via the role r . The set of \mathcal{L}^\square concepts naturally captures the propositional part of \mathcal{EL} . The *depth* of a subconcept D in C is the number of existential restrictions within the scope of which D remains. The *depth of a concept* C is the depth of its subconcept with the greatest depth in C . Every \mathcal{L}^\square concept is trivially of depth 0.

A *concept inclusion* (or a *TBox axiom*) is an expression of the form $C \sqsubseteq D$, stating that all individuals of type C are D . We sometimes write $C \equiv D$ as an abbreviation for two inclusions: $C \sqsubseteq D$ and $D \sqsubseteq C$. For instance, axioms (i) and (ii) below state, respectively, that (i) the class of mothers consists of all and only those individuals who are women and have at least one child, (ii) while every individual of type `Father_of_boy` is a father and has at least one male child:

$$\begin{aligned} \text{Mother} &\equiv \text{Woman} \sqcap \exists \text{hasChild}.\top && (i) \\ \text{Father_of_boy} &\sqsubseteq \text{Father} \sqcap \exists \text{hasChild}.\text{Man} && (ii) \end{aligned}$$

A TBox (or *ontology*) is a finite set of such concept inclusions in a particular language fragment. The language fragments considered in this paper are classified according to the type of restrictions imposed on the syntax of concepts C and D in the concept inclusions $C \sqsubseteq D$ permitted in the TBoxes:

$$\begin{aligned} \mathcal{EL}: & \quad C \text{ and } D \text{ are both } \mathcal{EL} \text{ concepts;} \\ \mathcal{EL}^{\text{rhs}}: & \quad C \text{ is an } \mathcal{L}^\square \text{ concept and } D \text{ an } \mathcal{EL} \text{ concept;} \\ \mathcal{EL}^{\text{lhs}}: & \quad C \text{ is an } \mathcal{EL} \text{ concept and } D \text{ an } \mathcal{L}^\square \text{ concept;} \\ \mathcal{L}^\square: & \quad C \text{ and } D \text{ are both } \mathcal{L}^\square \text{ concepts} \end{aligned}$$

For instance, a TBox consisting of axioms (i) and (ii) above, belongs to language \mathcal{EL} , as it in fact contains some $\mathcal{EL}^{\text{rhs}}$ axioms ($\text{Mother} \sqsubseteq \text{Woman} \sqcap \exists \text{hasChild}.\top$ and (ii)) as well as one $\mathcal{EL}^{\text{lhs}}$ axiom ($\text{Woman} \sqcap \exists \text{hasChild}.\top \sqsubseteq \text{Mother}$).

The semantics of DL languages is defined through interpretations of the form $\mathcal{I} = (\Delta^\mathcal{I}, \cdot^\mathcal{I})$, where $\Delta^\mathcal{I}$ is a non-empty *domain of individuals* and $\cdot^\mathcal{I}$ is an *interpretation function* mapping each $A \in N_C$ to a subset $A^\mathcal{I} \subseteq \Delta^\mathcal{I}$ and each $r \in N_R$ to a binary relation $r^\mathcal{I} \subseteq \Delta^\mathcal{I} \times \Delta^\mathcal{I}$. The interpretation function is inductively extended over complex expressions according to the fixed semantics of the constructors:

$$\begin{aligned} \top^\mathcal{I} &= \Delta^\mathcal{I} \\ (C \sqcap D)^\mathcal{I} &= \{x \in \Delta^\mathcal{I} \mid x \in C^\mathcal{I} \cap D^\mathcal{I}\} \\ (\exists r.C)^\mathcal{I} &= \{x \in \Delta^\mathcal{I} \mid \exists y : (x, y) \in r^\mathcal{I} \wedge y \in C^\mathcal{I}\} \end{aligned}$$

An interpretation \mathcal{I} *satisfies* a concept inclusion $C \sqsubseteq D$ ($\mathcal{I} \models C \sqsubseteq D$) iff $C^\mathcal{I} \subseteq D^\mathcal{I}$. Whenever \mathcal{I} satisfies all axioms in a TBox \mathcal{T} ($\mathcal{I} \models \mathcal{T}$), we say that \mathcal{I} is a *model* of \mathcal{T} . Interpretations and models defined in this way are in fact usual Kripke structures, which can be naturally represented as labelled graphs, with nodes representing individuals in the domain, edges — roles, and labels — the interpretations of concept and role names, respectively. For instance, the three graphs in Fig. 1 all represent possible models of the TBox consisting of axioms (i) and (ii) above:

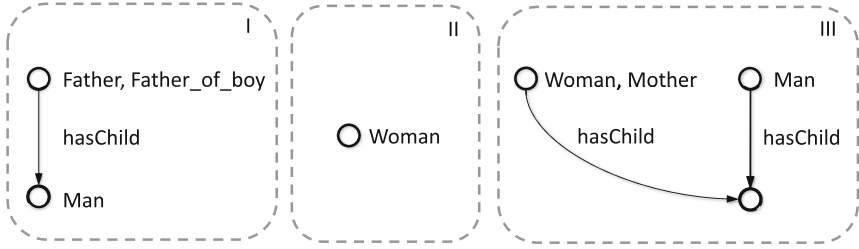


Fig. 1. Sample DL models.

Note that, as there are no a priori restrictions imposed on the number of domain individuals, a DL TBox might in general have infinitely many models of possibly infinite size. For a set of interpretations \mathcal{S} , we write $\mathcal{S} \models C \sqsubseteq D$ to denote that every interpretation in \mathcal{S} satisfies $C \sqsubseteq D$. We say that \mathcal{T} entails $C \sqsubseteq D$ ($\mathcal{T} \models C \sqsubseteq D$) iff every model of \mathcal{T} satisfies $C \sqsubseteq D$. Two TBoxes \mathcal{T} and \mathcal{H} are (logically) *equivalent* ($\mathcal{T} \equiv \mathcal{H}$) iff they have the same sets of models.

A *pointed interpretation* (\mathcal{I}, d) is a pair consisting of a DL interpretation $\mathcal{I} = (\Delta^{\mathcal{I}}, \cdot^{\mathcal{I}})$ and an individual $d \in \Delta^{\mathcal{I}}$, such that every $e \in \Delta^{\mathcal{I}}$ different from d is reachable from d through some role composition in \mathcal{I} . By a slight abuse of notation, given an arbitrary DL interpretation \mathcal{I} and an individual $d \in \Delta^{\mathcal{I}}$, we write (\mathcal{I}, d) to denote the largest subset \mathcal{I}' of \mathcal{I} such that (\mathcal{I}', d) is a pointed interpretation. If it is clear from the context, we refer to pointed interpretations and pointed models simply as interpretations and models. We say that (\mathcal{I}, d) is a model of a concept C iff $d \in C^{\mathcal{I}}$; it is a model of C w.r.t. \mathcal{T} whenever also $\mathcal{I} \models \mathcal{T}$.

An interpretation (\mathcal{I}, d) can be *homomorphically embedded* in an interpretation (\mathcal{J}, e) , denoted as $(\mathcal{I}, d) \mapsto (\mathcal{J}, e)$, iff there exists a mapping $h : \Delta^{\mathcal{I}} \mapsto \Delta^{\mathcal{J}}$, satisfying the following conditions:

- $h(d) = e$,
- if $(a, b) \in r^{\mathcal{I}}$ then $(h(a), h(b)) \in r^{\mathcal{J}}$, for every $a, b \in \Delta^{\mathcal{I}}$ and $r \in N_R$,
- if $a \in A^{\mathcal{I}}$ then $h(a) \in A^{\mathcal{J}}$, for every $a \in \Delta^{\mathcal{I}}$ and $A \in N_C$.

A model (\mathcal{I}, d) of C (w.r.t. \mathcal{T}) is called *minimal* iff it can be homomorphically embedded in every other model of C (w.r.t. \mathcal{T}). It is well-known that \mathcal{EL} concepts and TBoxes always have such minimal models (unique up to homomorphic embeddings) [9]. As in most modal logics, arbitrary \mathcal{EL} models can be unravelled into equivalent tree-shaped models. Finally, we observe that due to a tight relationship between the syntax and semantics of \mathcal{EL} , every tree-shaped interpretation (\mathcal{I}, d) can be viewed as an \mathcal{EL} concept $C_{\mathcal{I}}$, such that (\mathcal{I}, d) is a minimal model of $C_{\mathcal{I}}$. Formally, we set $C_{\mathcal{I}} = C(d)$, where for every $e \in \Delta^{\mathcal{I}}$ we let $C(e) = \top \sqcap A(e) \sqcap \exists(e)$, with $A(e) = \prod\{A \in N_C \mid e \in A^{\mathcal{I}}\}$ and $\exists(e) = \prod_{(r,f) \in N_R \times \Delta^{\mathcal{I}} \text{ s.t. } (e,f) \in r^{\mathcal{I}}} \exists r.C(f)$. In that case we call $C_{\mathcal{I}}$ the *covering concept* for (\mathcal{I}, d) . For instance, the covering concept for model I in Fig. 1 is $\top \sqcap \text{Father} \sqcap \text{Father_of_boy} \sqcap \exists \text{hasChild} . (\top \sqcap \text{Man})$, which can be simplified as $\text{Father} \sqcap \text{Father_of_boy} \sqcap \exists \text{hasChild} . (\text{Man})$.

3 Learning Model

The learning model studied in this paper is a variant of learning from positive interpretations [5,6]. In our setting, the teacher fixes a *target TBox* \mathcal{T} , whose set of all models is denoted by $\mathcal{M}(\mathcal{T})$. Further, the teacher presents a set of examples from $\mathcal{M}(\mathcal{T})$ to the learner, whose goal is to correctly identify \mathcal{T} based on this input. The learning process is conducted relative to a mutually known DL language \mathcal{L} and a finite vocabulary $\Sigma_{\mathcal{T}}$ used in \mathcal{T} .

In principle, $\mathcal{M}(\mathcal{T})$ contains sufficient information in order to enable correct identification of \mathcal{T} , as the following correspondence implies:

$$\mathcal{M}(\mathcal{T}) \models C \sqsubseteq D \text{ iff } \mathcal{T} \models C \sqsubseteq D, \text{ for every } C \sqsubseteq D \text{ in } \mathcal{L}.$$

However, as $\mathcal{M}(\mathcal{T})$ might consist of infinitely many models of possibly infinite size, the teacher cannot effectively present them all to the learner. Instead, the teacher must confine him- or herself to certain finitely presentable subset of $\mathcal{M}(\mathcal{T})$, called the *learning set*. For the sake of clarity, we focus here on the simplest case when learning sets consist of finitely many finite models.² Formally, we summarize the learning model with the following definitions.

Definition 1 (TIP). *A TBox Identification Problem (TIP) is a pair $(\mathcal{T}, \mathcal{S})$, where \mathcal{T} is a TBox in a DL language \mathcal{L} and \mathcal{S} , called the learning set, is a finite set of finite models of \mathcal{T} .*

Definition 2 (Learner, identification). *For a DL language \mathcal{L} , a learner is a computable function G , which for every set \mathcal{S} over $\Sigma_{\mathcal{T}}$ returns a TBox in \mathcal{L} over $\Sigma_{\mathcal{T}}$. Learner G correctly identifies \mathcal{T} on \mathcal{S} whenever $G(\mathcal{S}) \equiv \mathcal{T}$.*

Definition 3 (Learnability). *For a DL language \mathcal{L} , the class of TBoxes expressible in \mathcal{L} is learnable iff there exists a learner G such that for every TBox \mathcal{T} in \mathcal{L} there exists a learning set \mathcal{S} on which G correctly identifies \mathcal{T} . It is said to be finitely learnable whenever it is learnable from finite learning sets only.*

We are primarily interested here in the notion of finite learnability, as it provides a natural formal foundation for the task of ontology learning from data. By data, in the DL context, we understand collections of atomic concept and role assertions over domain individuals (e.g., `Father(john)`, `hasChild(john, mary)`), which under certain assumptions regarding their structuring with respect to the background ontology can be seen as models of that ontology and, consequently, as potentially valuable learning sets. Figure 2 presents an example of a TIP with a finite learning set, which consists of a single model of the assumed ontology. The key question is then what formal criteria must this set satisfy to warrant correct identification of the ontology constraining it. To this end we employ the

² An alternative, more general approach can be defined in terms of specific fragments of models. Such generalization, which lies beyond the scope of this paper, is essential when the learning problem concerns languages without finite model property.

$$\begin{aligned} \text{Mother} &\equiv \text{Woman} \sqcap \exists \text{hasChild}.\top \\ \text{Father} &\equiv \text{Man} \sqcap \exists \text{hasChild}.\top \\ \text{Father_of_boy} &\equiv \text{Father} \sqcap \exists \text{hasChild}.\text{Man} \end{aligned}$$

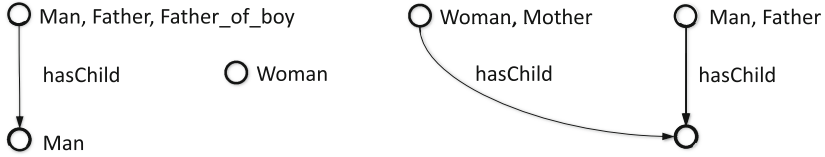


Fig. 2. A sample TIP with an \mathcal{EL} TBox and a finite learning set

basic *admissibility condition*, characteristic also of other learning frameworks [10], which ensures that the learning set is sufficiently rich to enable precise discrimination between the correct hypothesis and all the incorrect ones.

Definition 4 (Admissibility). A TIP $(\mathcal{T}, \mathcal{S})$ is admissible iff for every $C \sqsubseteq D$ in \mathcal{L} such that $\mathcal{T} \not\models C \sqsubseteq D$ there exists $\mathcal{I} \in \mathcal{S}$ such that $\mathcal{I} \not\models C \sqsubseteq D$.

For the target TBox \mathcal{T} , let \mathcal{T}^\neq be the set of all concept inclusions in \mathcal{L} that are not entailed by \mathcal{T} , i.e., $\mathcal{T}^\neq = \{C \sqsubseteq D \text{ in } \mathcal{L} \mid \mathcal{T} \not\models C \sqsubseteq D\}$. The admissibility condition requires that for every $C \sqsubseteq D \in \mathcal{T}^\neq$, the learning set \mathcal{S} must contain a “counterexample” for it, i.e., an individual $d \in \Delta^{\mathcal{I}}$, for some $\mathcal{I} \in \mathcal{S}$, such that $d \in C^{\mathcal{I}}$ and $d \notin D^{\mathcal{I}}$. Consequently, any learning set must contain such counterexamples to all elements of \mathcal{T}^\neq , or else, the learner might never be justified to exclude some of these concept inclusions from the hypothesis. If it was possible to represent them finitely we could expect that ultimately the learner can observe all of them and correctly identify the TBox. In the next section, we investigate this prospect formally in different fragments of \mathcal{EL} .

4 Finite Learning Sets

As argued in the previous section, to enable finite learnability of \mathcal{T} in a given language \mathcal{L} , the relevant counterexamples to all the concept inclusions not entailed by \mathcal{T} must be presentable within a finite learning set \mathcal{S} . Firstly, we can immediately observe that this requirement is trivially satisfied for \mathcal{L}^\square . Clearly, \mathcal{L}^\square can only induce finitely many different concept inclusions (up to logical equivalence) on finite vocabularies, such as $\Sigma_{\mathcal{T}}$. Hence, the set \mathcal{T}^\neq can always be finitely represented (up to logical equivalence) and it is straightforward to finitely present counterexamples to all its members. For more expressive fragments of \mathcal{EL} , however, this cannot be assumed in general, as the $\exists r.C$ constructor induces infinitely many concepts. One negative result comes with the case of \mathcal{EL} itself, as demonstrated in the next theorem.

Theorem 1 (Finite learning sets in \mathcal{EL}). *Let \mathcal{T} be a TBox in \mathcal{EL} . There exists no finite set \mathcal{S} such that $(\mathcal{T}, \mathcal{S})$ is admissible.*

The full proof of this and subsequent results is included in an online technical report [11]. The argument rests on the following lemma. Let $(\mathcal{T}, \mathcal{S})$ be an admissible TIP and C a concept. By $\mathcal{S}(C)$ we denote the set of all models (\mathcal{I}, d) of C w.r.t. \mathcal{T} such that $\mathcal{I} \in \mathcal{S}$. By $\bigcap \mathcal{S}(C)$ we denote the intersection of all these models, i.e., the model (\mathcal{J}, d) , such that:

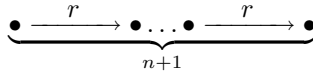
1. $(\mathcal{J}, d) \mapsto (\mathcal{I}, d)$ for every $(\mathcal{I}, d) \in \mathcal{S}(C)$,
2. for every other model (\mathcal{J}', d) such that $(\mathcal{J}', d) \mapsto (\mathcal{I}, d)$ for every $(\mathcal{I}, d) \in \mathcal{S}(C)$ and $(\mathcal{J}, d) \mapsto (\mathcal{J}', d)$, it is the case that $(\mathcal{J}', d) \mapsto (\mathcal{J}, d)$.

Lemma 1 (Minimal model lemma). *Let $(\mathcal{T}, \mathcal{S})$ be an admissible TIP for \mathcal{T} in \mathcal{EL} (resp. in \mathcal{EL}^{rhs}), and C be an \mathcal{EL} (resp. \mathcal{L}^\square) concept. Whenever $\mathcal{S}(C)$ is non-empty then $\bigcap \mathcal{S}(C)$ is a minimal model of C w.r.t. \mathcal{T} .*

Given the lemma, we consider a concept inclusion of type:

$$\tau_n := \underbrace{\exists r. \dots \exists r.}_{n} \top \sqsubseteq \underbrace{\exists r. \dots \exists r. \exists r.}_{n+1} \top$$

Suppose $\tau_n \in \mathcal{T}^\neq$ for some $n \in \mathbb{N}$. Since by the admissibility condition a counterexample to τ_n must be present in \mathcal{S} , it must be the case that $\mathcal{S}(C) \neq \emptyset$, where C is the left-hand-side concept in τ_n . By the lemma and the definition of a minimal model, it is easy to see that \mathcal{S} must contain a finite chain of individuals of length exactly $n + 1$, as depicted below:



Finally, since there can always exist some $n \in \mathbb{N}$, such that $\tau_m \in \mathcal{T}^\neq$ for every $m \geq n$, we see that the joint size of all necessary counterexamples in such cases must inevitably be also infinite. Consequently, for some \mathcal{EL} TBoxes admissible TIPs based on finite learning sets might not exist, and so finite learnability cannot be achieved in general.

One trivial way to tame this behavior is to “finitize” \mathcal{T}^\neq by delimiting the entire space of possible TBox axioms to a pre-defined, finite set. This can be achieved, for instance, by restricting the permitted depth of complex concepts or generally setting some a priori bound on the size of axioms. Such ad hoc solutions, though likely efficient in practice, are not very elegant. As a more interesting alternative, we are able to show that there exist at least two languages between \mathcal{L}^\square and \mathcal{EL} , namely \mathcal{EL}^{lhs} and \mathcal{EL}^{rhs} , for which finite learning sets are always guaranteed to exist, regardless of the fact that they permit infinitely many concept inclusions. In fact, we demonstrate that in both cases such learning sets might well consist of exactly one exemplary finite model.

We adopt the technique of so-called *types*, known from the area of modal logics [12]. Types are finite abstractions of possible individuals in the interpretation

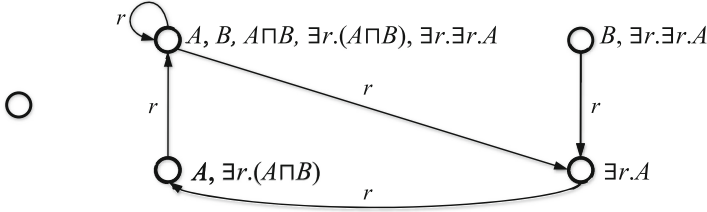


Fig. 3. A finite learning set for an $\mathcal{EL}^{\text{rhs}}$ TBox $\{A \sqsubseteq \exists r.(A \sqcap B), B \sqsubseteq \exists r.\exists r.A\}$. The figure includes type contents (in grey), as defined in the proof of Theorem 2.

domain, out of which arbitrary models can be constructed. Let $\text{con}(\mathcal{T})$ be the set of all concepts (and all their subconcepts) occurring in \mathcal{T} . A *type* over \mathcal{T} is a set $t \subseteq \text{con}(\mathcal{T})$, such that $C \sqcap D \in t$ iff $C \in t$ and $D \in t$, for every $C \sqcap D \in \text{con}(\mathcal{T})$. A type t is *saturated* for \mathcal{T} iff for every $C \sqsubseteq D \in \mathcal{T}$, if $C \in t$ then $D \in t$. For any $S \subseteq \text{con}(\mathcal{T})$, we write t_S to denote the smallest saturated type containing S . It is easy to see, that t_S must be unique for \mathcal{EL} .

The next theorem addresses the case of $\mathcal{EL}^{\text{rhs}}$. Figure 3 illustrates a finite learning set for a sample $\mathcal{EL}^{\text{rhs}}$ TBox, following the construction in the proof.

Theorem 2 (Finite learning sets in $\mathcal{EL}^{\text{rhs}}$). *Let \mathcal{T} be a TBox in $\mathcal{EL}^{\text{rhs}}$. There exists a finite set \mathcal{S} such that $(\mathcal{T}, \mathcal{S})$ is admissible.*

Proof sketch. Let Θ be the smallest set of types satisfying the following conditions:

- $t_S \in \Theta$, for every $S \subseteq N_C$ and for $S = \{\top\}$,
- if $t \in \Theta$ then $t_{\{C\}} \in \Theta$, for every $\exists r.C \in t$.

We define the interpretation $\mathcal{I} = (\Delta^{\mathcal{I}}, \cdot^{\mathcal{I}})$ as follows:

- $\Delta^{\mathcal{I}} := \Theta$,
- $t \in A^{\mathcal{I}}$ iff $A \in t$, for every $t \in \Theta$ and $A \in N_C$,
- $(t, t_{\{C\}}) \in r^{\mathcal{I}}$, for every $t \in \Theta$, whenever $\exists r.C \in t$.

Then $\mathcal{S} = \{\mathcal{I}\}$ is a finite learning set such that $(\mathcal{T}, \mathcal{S})$ is admissible. \square

A similar, though somewhat more complex construction demonstrates the existence of finite learning sets in $\mathcal{EL}^{\text{lhs}}$. Again, we illustrate the approach with an example in Fig. 4.

Theorem 3 (Finite learning sets in $\mathcal{EL}^{\text{lhs}}$). *Let \mathcal{T} be a TBox in $\mathcal{EL}^{\text{lhs}}$. There exists a finite set \mathcal{S} such that $(\mathcal{T}, \mathcal{S})$ is admissible.*

Proof Sketch. Let Θ be the set of all saturated types over \mathcal{T} , and Θ^* be its subset obtained by iteratively eliminating all those types t that violate the following condition: for every $r \in N_R$ and every existential restriction $\exists r.C \in t$ there is $u \in \Theta^*$ such that:

- $C \in u$,
- for every $\exists r.D \in \text{con}(\mathcal{T})$, if $D \in u$ then $\exists r.D \in t$.

Further, we define the interpretation $\mathcal{I} = (\Delta^{\mathcal{I}}, \cdot^{\mathcal{I}})$ as follows:

- $\Delta^{\mathcal{I}} := \Theta^*$,
- $t \in A^{\mathcal{I}}$ iff $A \in S_t$, for every $t \in \Theta^*$ and $A \in N_{\mathcal{C}}$,
- $(t, u) \in r^{\mathcal{I}}$ iff for every $\exists r.C \in \text{con}(\mathcal{T})$, if $C \in u$ then $\exists r.C \in t$.

Then $\mathcal{S} = \{\mathcal{I}\}$ is a finite learning set such that $(\mathcal{T}, \mathcal{S})$ is admissible. □

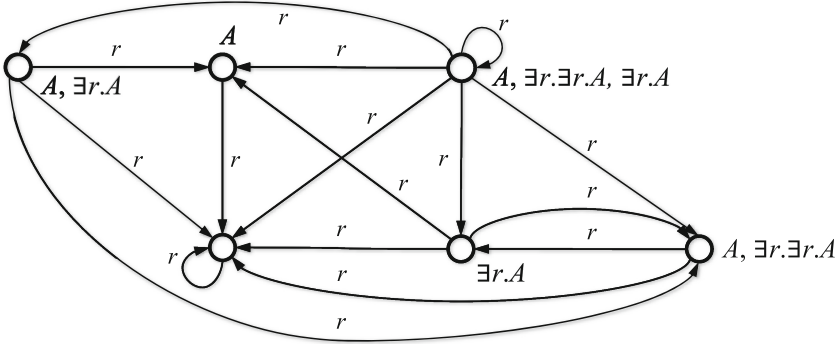


Fig. 4. A finite learning set for an $\mathcal{EL}^{\text{lhs}}$ TBox $\{\exists r.\exists r.A \sqsubseteq A\}$. The figure includes type contents (in grey), as defined in the proof of Theorem 3.

5 Learning Algorithms

In this section, we devise two learning algorithms for admissible TIPs with finite learning sets that correctly identify 1) $\mathcal{EL}^{\text{lhs}}$ and $\mathcal{EL}^{\text{rhs}}$ TBoxes using an equivalence oracle, and 2) $\mathcal{EL}^{\text{rhs}}$ TBoxes without such an oracle, i.e., in a fully unsupervised manner.

Since the set $\mathcal{T}^{\neq} = \{C \sqsubseteq D \text{ in } \mathcal{L} \mid \mathcal{T} \not\equiv C \sqsubseteq D\}$ can be in general infinite, our starting observation is that a learner cannot effectively eliminate concept inclusions from \mathcal{T}^{\neq} using a straightforward enumeration, thus arriving at the target TBox \mathcal{T} . The only feasible strategy is to try to identify the “good” candidate axioms to be included in \mathcal{T} , and possibly apply the elimination strategy only to finitely many incorrect guesses. One generic procedure to employ such heuristic, which we define as Algorithm 1, attempts to construct the hypothesis by extending it with consecutive axioms of systematically growing size that are satisfied by the learning set. There, by $\ell(C \sqsubseteq D)$ we denote the size of the axiom $C \sqsubseteq D$ measured in the total number of symbols used for expressing this axiom. At each step the algorithm makes use of a simple equivalence oracle, which informs whether the currently considered hypothesis is already equivalent to the learning target (in that case the identification succeeds) or whether some axioms are still missing. Theorem 4 demonstrates the correctness of this approach.

Algorithm 1. Learning $\mathcal{EL}^{\text{rhs}}/\mathcal{EL}^{\text{lhs}}$ TBoxes on finite inputs.

Input: a TIP $(\mathcal{T}, \mathcal{S})$

Output: a hypothesis TBox \mathcal{H}

```

1:  $n := 2$ 
2:  $\mathcal{H}_n := \emptyset$ 
3: while ' $\mathcal{H}_n \equiv \mathcal{T}$ '? is 'NO' (equivalence oracle querying) do
4:    $n := n + 1$ 
5:    $\text{Cand}_n := \{C \sqsubseteq D \in \mathcal{EL}^{\text{rhs}}/\mathcal{EL}^{\text{lhs}} \mid \ell(C \sqsubseteq D) = n\}$ 
6:    $\text{Accept}_n := \{C \sqsubseteq D \in \text{Cand}_n \mid \mathcal{S} \models C \sqsubseteq D\}$ 
7:    $\mathcal{H}_n := \mathcal{H}_{n-1} \cup \text{Accept}_n$ 
8: end while
9: return  $\mathcal{H}_n$ 

```

Theorem 4 (Correct identification in $\mathcal{EL}^{\text{rhs}}/\mathcal{EL}^{\text{lhs}}$). *Let $(\mathcal{T}, \mathcal{S})$ be an admissible TIP for \mathcal{T} in $\mathcal{EL}^{\text{rhs}}/\mathcal{EL}^{\text{lhs}}$. Then the hypothesis TBox \mathcal{H} generated by Algorithm 1 is equivalent to \mathcal{T} .*

Obviously the use of the oracle is essential to warrant termination of the algorithm. It is not difficult to see that without it, the algorithm must still converge on the correct TBox for some $n \in \mathbb{N}$, and consequently settle on it, i.e., $\mathcal{H}_m \equiv \mathcal{H}_n$ for every $m \geq n$. However, at no point of time can it guarantee that the convergence has been already achieved, and so it can only warrant learnability in the limit. This result is therefore not entirely satisfactory considering we aim at finite learnability from data in the unsupervised setting.

A major positive result, on the contrary, can be delivered for the case of $\mathcal{EL}^{\text{rhs}}$, for which we devise an effective learning algorithm making no reference to any oracle. It turns out that in $\mathcal{EL}^{\text{rhs}}$ the “good” candidate axioms can be directly extracted from the learning set, thus granting a proper unsupervised learning method. The essential insight is provided by Lemma 1, presented in the previous section. Given any \mathcal{L}^\square concept C such that $\mathcal{S}(C) \neq \emptyset$ we are able to identify a tree-shaped minimal model of C w.r.t. \mathcal{T} . Effectively, it suffices to retrieve only the initial part of this model, discarding its infinitely recurrent (cyclic) subtrees. Such an initial model $\mathcal{I}_{\text{init}}$ is constructed by Algorithm 2. The algorithm performs simultaneous unravelling of all models in $\mathcal{S}(C)$, while on the way, computing intersections of visited combinations of individuals, which are subsequently added to the model under construction. Whenever the same combination of individuals is about to be visited for the third time on the same branch it is skipped, as the cycle is evidently detected and further unravelling is unnecessary. The covering concept $C_{\mathcal{I}_{\text{init}}}$ for the resulting interpretation $\mathcal{I}_{\text{init}}$ is then included in the hypothesis within the axiom $C \sqsubseteq C_{\mathcal{I}_{\text{init}}}$. Meanwhile, all \mathcal{L}^\square concepts C such that $\mathcal{S}(C) = \emptyset$ are ensured to entail every \mathcal{EL} concept, as implied by the admissibility condition. The contents of the hypothesis TBox are formally specified in Definition 5.

Algorithm 2. Computing the initial part of the minimal model $\bigcap \mathcal{S}(C)$

Input: the set $\mathcal{S}(C) = \{(\mathcal{I}_i, d_i)\}_{0 \leq i \leq n}$, for some $n \in \mathbb{N}$

Output: a finite tree-shaped interpretation (\mathcal{J}, d) , where $\mathcal{J} = (\Delta^{\mathcal{J}}, \cdot^{\mathcal{J}})$

- 1: $\Delta^{\mathcal{J}} := \{f(d_0, \dots, d_n)\}$, for a “fresh” function symbol f
 - 2: $A^{\mathcal{J}} := \emptyset$, for every $A \in N_C$
 - 3: $r^{\mathcal{J}} := \emptyset$, for every $r \in N_R$
 - 4: **for** every $f(d_0, \dots, d_n) \in \Delta^{\mathcal{J}}$, $(e_0, \dots, e_n) \in \Delta^{\mathcal{I}_0} \times \dots \times \Delta^{\mathcal{I}_n}$, $r \in N_R$ **do**
 - 5: **if** $(d_i, e_i) \in r^{\mathcal{I}_i}$ for every $0 \leq i \leq n$ **and** there exists no function symbol g such that $g(e_0, \dots, e_n)$ is an ancestor of $f(d_0, \dots, d_n)$ in \mathcal{J} **then**
 - 6: $\Delta^{\mathcal{J}} := \Delta^{\mathcal{J}} \cup \{g(e_0, \dots, e_n)\}$, for a “fresh” function symbol g
 - 7: $r^{\mathcal{J}} := r^{\mathcal{J}} \cup \{(f(d_0, \dots, d_n), g(e_0, \dots, e_n))\}$
 - 8: **end if**
 - 9: **end for**
 - 10: **for** every $f(d_0, \dots, d_n) \in \Delta^{\mathcal{J}}$, $A \in N_C$ **do**
 - 11: **if** $d_i \in A^{\mathcal{I}_i}$ for every $0 \leq i \leq n$ **then**
 - 12: $A^{\mathcal{J}} := A^{\mathcal{J}} \cup \{f(d_0, \dots, d_n)\}$
 - 13: **end if**
 - 14: **end for**
 - 15: **return** $(\mathcal{J}, f(d_0, \dots, d_n))$, where $f(d_0, \dots, d_n)$ is the root of \mathcal{J} , created at step 1.
-

Definition 5 ($\mathcal{E}\mathcal{L}^{\text{rhs}}$ hypothesis TBox). *Let $(\mathcal{T}, \mathcal{S})$ be an admissible TIP for \mathcal{T} in $\mathcal{E}\mathcal{L}^{\text{rhs}}$ over the vocabulary $\Sigma_{\mathcal{T}}$. The hypothesis TBox \mathcal{H} is the set consisting of all the following axioms:*

- $C \sqsubseteq C_{\mathcal{I}_{\text{init}}}$ for every \mathcal{L}^{\square} concept C such that $\mathcal{S}(C) \neq \emptyset$, where $C_{\mathcal{I}_{\text{init}}}$ is the covering concept for the interpretation $(\mathcal{I}_{\text{init}}, d)$ generated by Algorithm 2 on $\mathcal{S}(C)$;
- $C \sqsubseteq \prod_{r \in N_R} \exists r. \prod N_C$ for every \mathcal{L}^{\square} concept C such that $\mathcal{S}(C) = \emptyset$.

To better illustrate the learning procedure, we consider a simple TIP consisting of an $\mathcal{E}\mathcal{L}^{\text{rhs}}$ TBox $\mathcal{T} = \{A \sqsubseteq \exists r.(A \sqcap B)\}$ and a finite learning set $\mathcal{S} = \{\mathcal{I}\}$, with \mathcal{I} as depicted in Fig. 5. The assumed vocabulary containing two concept names — A and B — induces four distinct \mathcal{L}^{\square} concepts, namely: \top , A , B and $A \sqcap B$. For every such concept C we identify the corresponding set of its all pointed models $\mathcal{S}(C)$ contained in the learning set. For instance, $\mathcal{S}(A) = \{(\mathcal{I}, e_2), (\mathcal{I}, e_3)\}$. Further, we use Algorithm 2 to compute the initial part of the minimal model $\bigcap \mathcal{S}(C)$, as illustrated in Fig. 6. Finally, based on these models we formulate the hypothesis TBox, as specified in Definition 5: $\mathcal{H} = \{\top \sqsubseteq \top, A \sqsubseteq \exists r.(A \sqcap B \sqcap \exists r.(A \sqcap B)), B \sqsubseteq B, A \sqcap B \sqsubseteq \exists r.(A \sqcap B)\}$. It is not difficult to verify that $\mathcal{H} \equiv \mathcal{T}$.

The correctness of the learning procedure is demonstrated in the following theorem.

Theorem 5 (Correct identification in $\mathcal{E}\mathcal{L}^{\text{rhs}}$). *Let $(\mathcal{T}, \mathcal{S})$ be an admissible TIP for \mathcal{T} in $\mathcal{E}\mathcal{L}^{\text{rhs}}$. Then the hypothesis TBox \mathcal{H} for \mathcal{S} is equivalent to \mathcal{T} .*

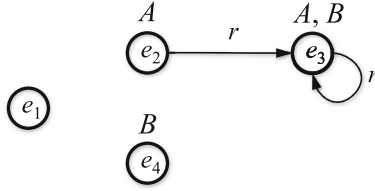


Fig. 5. A finite learning set for an $\mathcal{EL}^{\text{rhd}}$ TBox $\{A \sqsubseteq \exists r.(A \sqcap B)\}$.

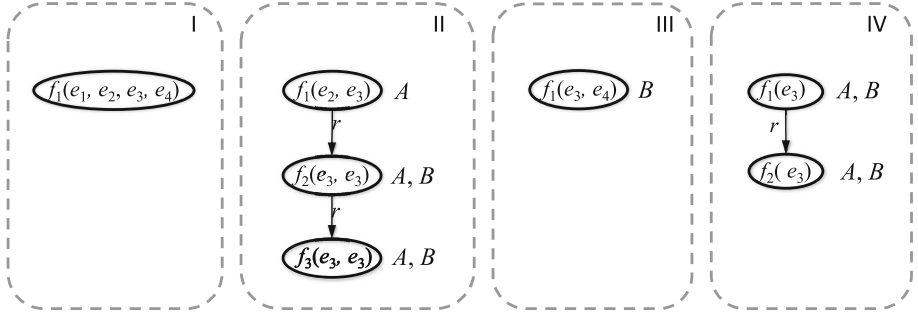


Fig. 6. The initial parts of the minimal models $\bigcap \mathcal{S}(C)$ computed with Algorithm 2 over the learning set in Fig. 5, where (I) $C = \top$, (II) $C = A$, (III) $C = B$, (IV) $C = A \sqcap B$.

Proof sketch. Let C be a concept in \mathcal{L}^\square such that $\mathcal{S}(C) \neq \emptyset$. By Lemma 1, the intersection $(\mathcal{J}, d) = \bigcap \mathcal{S}(C)$ is a minimal model of C w.r.t. \mathcal{T} . Without loss of generality, we assume that (\mathcal{J}, d) is a tree-shaped model. We also note, that every \mathcal{EL} concept C induces a syntactic tree, which corresponds directly to a minimal model of C . It is not difficult to see that Algorithm 2 indeed produces an initial part $(\mathcal{J}_{\text{init}}, d)$ of (\mathcal{J}, d) . By reconstructing the concept C_{init} from we in fact identify all minimal (i.e., necessary) consequences of C w.r.t. \mathcal{T} . However, certain infinite subtrees of (\mathcal{J}, d) are omitted in $(\mathcal{J}_{\text{init}}, d)$. This happens due to the condition at step 5 of Algorithm 2, which terminates the construction of certain branches whenever a cycle is detected. In the rest of the proof, we show that the covering concept $C_{\mathcal{J}_{\text{init}}}$ has the same minimal model w.r.t. \mathcal{H} as C has w.r.t. \mathcal{T} . Since this is demonstrated to hold for every \mathcal{L}^\square concept C , we can conclude that $\mathcal{H} \equiv \mathcal{T}$. \square

The learning algorithm runs in double exponential time in the worst case and generates TBoxes of double exponential size in the size of \mathcal{S} . This follows from the fact that the tree-shaped interpretations generated by Algorithm 2 might be of depth exponential in the number of individuals occurring in \mathcal{S} and have exponential branching factor. Importantly, however, there might exist solutions far closer to being optimal which we have not as far investigated.

It is our strong conjecture, which we leave as an open problem, that a similar learning strategy should also be applicable in the context of $\mathcal{EL}^{\text{lhs}}$.

6 Related Work

Ontology learning is an interdisciplinary research field drawing on techniques from Formal Concept Analysis [13, 14], Natural Language Processing [15, 16] and machine learning [6, 16, 17], to name a few. One classification of ontology learning techniques distinguishes between investigations of exact learnability, and approaches incorporating probabilistic, vague or fuzzy reasoning [18, 19]. Another classification is at the level at which learning takes place [15] — does the problem address learning of concepts, concept hierarchies, logical theories or rules? Lehmann and Völker [17] distinguishes between four types of ontology learning: learning from text, data mining, concept learning and crowdsourcing.

In this landscape, the present paper is on exact learnability and, within this framework, addresses the problem of learning logical theories. That is, we address the problem at the level of relationships between concepts, positing a logical theory, rather than at the concept level, learning concept descriptions [20–24]. Furthermore, the target theory is identified from interpretations, and is hence related to various contributions on learnability of different types of formal structures from data, e.g.: first-order theories from facts [10], finite automata descriptions from observations [25], and logic programs from interpretations [5, 6].

The model for exact learning of DL TBoxes which offers the most direct comparison to ours was introduced recently by Konev, et al. [8], and follows on prior research by the same authors based on Angluin’s model of learning from entailment [7, 26]. In their learning framework for learning from data retrieval queries, the learner identifies the TBox by posing two types of queries to an oracle: membership queries of the form “ $(\mathcal{T}, \mathcal{A}) \models q?$ ”, where \mathcal{A} is a given ABox and q is a query, and equivalence queries of the form “Does the hypothesis ontology \mathcal{H} entail the target ontology \mathcal{T} ?”. The authors study polynomial learnability in fragments of \mathcal{EL} and DL-Lite, and for queries ranging from atomic to conjunctive queries.

Essentially, given a finite learning set in an admissible TIP, a learner from interpretations can autonomously answer arbitrary membership queries, thus effectively simulating the membership oracle. However, it does not have by default access to an equivalence oracle. Once such an oracle is included, as in Algorithm 1, the learning power of both learners becomes comparable for the languages investigated in the present paper. In this sense, our Theorem 4 should be also indirectly derivable from the results by Konev et al. However, our stronger result for $\mathcal{EL}^{\text{rhs}}$ in Theorem 5 demonstrates that, at least in some cases, the learner from interpretations is able to succeed without employing any oracle. While learning from ABoxes and query answers makes sense in a semi-automated learning environment, learning from interpretations is in our view a more appropriate model in the context of fully autonomous learning.

7 Conclusions and Outlook

In this paper, we have delivered initial results on finite learnability of DL TBoxes from interpretations. We believe that this direction shows promise in establishing

formal foundations for the task of ontology learning from data. Some immediate problems that are left open with this work concern finite learnability of \mathcal{EL}^{hs} TBoxes in an unsupervised setting, and possibly of other lightweight fragments of DLs. Another set of very interesting research questions should deal, in our view, with the possibility of formulating alternative conditions on the learning sets and the corresponding learnability guarantees they would imply in different DL languages. In particular, some limited use of closed-world operator over the learning sets might allow to relax the practically restrictive admissibility condition. Finally, the development of practical learning algorithms, possibly building on existing inductive logic programming methods, is an obvious area to welcome further research efforts.

References

1. Maedche, A., Staab, S.: Ontology learning. In: Staab, S., Studer, R. (eds.) *Handbook on Ontologies*, pp. 173–189. Springer, New York (2004)
2. Hoekstra, R.: The knowledge reengineering bottleneck. *J. Semant. Web* **1**(1,2), 111–115 (2010)
3. Baader, F., Calvanese, D., McGuinness, D.L., Nardi, D., Patel-Schneider, P.F.: *The Description Logic Handbook: Theory, Implementation, and Applications*. Cambridge University Press, New York (2003)
4. Baader, F., Brandt, S., Lutz, C.: Pushing the \mathcal{EL} envelope. In: *Proceedings of IJCAI-05* (2005)
5. De Raedt, L., Lavrač, N.: The many faces of inductive logic programming. In: Komorowski, J., Raś, Z.W. (eds.) *ISMIS 1993*. LNCS, vol. 689, pp. 435–449. Springer, Heidelberg (1993)
6. De Raedt, L.: First order jk-clausal theories are PAC-learnable. *Artif. Intell.* **70**, 375–392 (1994)
7. Konev, B., Lutz, C., Ozaki, A., Wolter, F.: Exact learning of lightweight description logic ontologies. In: *Proceedings of Principles of Knowledge Representation and Reasoning (KR-14)* (2014)
8. Konev, B., Lutz, C., Wolter, F.: Exact learning of TBoxes in \mathcal{EL} and DL-Lite. In: *Proceedings of the 28th International Workshop on Description Logics* (2015)
9. Lutz, C., Piro, R., Wolter, F.: Enriching \mathcal{EL} -concepts with greatest fixpoints. In: *Proceedings of the 19th European Conference on Artificial Intelligence (ECAI 2010)*, pp. 41–46. IOS Press (2010)
10. Shapiro, E.Y.: *Inductive inference of theories from facts*. In: *Computational Logic: Essays in Honor of Alan Robinson* (1991). MIT Press (1981)
11. Klarman, S., Britz, K.: *Ontology learning from interpretations in lightweight description logics*. Technical report, CSIR Centre for Artificial Intelligence Research, South Africa (2015). <http://klarman.synthasite.com/resources/KlaBri-ILP15.pdf>
12. Pratt, V.: Models of program logics. In: *Proceedings of Foundations of Computer Science (FOCS 1979)* (1979)
13. Baader, F., Ganter, B., Sertkaya, B., Sattler, U.: Completing description logic knowledge bases using formal concept analysis. In: *Proceedings of the 20th International Joint Conference on Artificial Intelligence (IJCAI-07)* (2007)
14. Distel, F.: *Learning description logic knowledge bases from data using methods from formal concept analysis*. Ph.D. Thesis, TU Dresden (2011)

15. Buitelaar, P., Cimeano, P., Magnini, F. (eds.): *Ontology Learning from Text: Methods, Evaluation and Applications*. IOS Press, Amsterdam (2005)
16. Cimeano, P., Mädche, A., Staab, S., Völker, J.: *Ontology learning*. In: Staab, S., Studer, R. (eds.) *Handbook on Ontologies*. Springer, New York (2009)
17. Lehmann, J., Völker, J. (eds.): *Perspectives on Ontology Learning*. IOS Press, Amsterdam (2014)
18. Cohen, W., Hirsh, H.: The learnability of description logics with equality constraints. *Mach. Learn.* **17**(2–3), 169–199 (1994)
19. Lisi, F.A., Straccia, U.: A FOIL-like method for learning under incompleteness and vagueness. In: Zaverucha, G., Santos Costa, V., Paes, A. (eds.) *ILP 2013*. LNCS, vol. 8812, pp. 123–139. Springer, Heidelberg (2014)
20. Badea, L., Nienhuys-Cheng, S.-H.: A refinement operator for description logics. In: Cussens, J., Frisch, A.M. (eds.) *ILP 2000*. LNCS (LNAI), vol. 1866, pp. 40–59. Springer, Heidelberg (2000)
21. Lehmann, J., Hitzler, P.: A refinement operator based learning algorithm for the \mathcal{ALC} description logic. In: Blockeel, H., Ramon, J., Shavlik, J., Tadepalli, P. (eds.) *ILP 2007*. LNCS (LNAI), vol. 4894, pp. 147–160. Springer, Heidelberg (2008)
22. Fanizzi, N., d’Amato, C., Esposito, F.: DL-FOIL concept learning in description logics. In: Železný, F., Lavrač, N. (eds.) *ILP 2008*. LNCS (LNAI), vol. 5194, pp. 107–121. Springer, Heidelberg (2008)
23. Cohen, W.W., Hirsh, H.: Learning the classic description logic: Theoretical and experimental results. In: *Proceedings of Principles of Knowledge Representation and Reasoning (KR 1994)* (1994)
24. Chitsaz, M., Wang, K., Blumenstein, M., Qi, G.: Concept learning for $\mathcal{EL}++$ by refinement and reinforcement. In: Anthony, P., Ishizuka, M., Lukose, D. (eds.) *PRICAI 2012*. LNCS, vol. 7458, pp. 15–26. Springer, Heidelberg (2012)
25. Pitt, L.: Inductive inference, DFAs, and computational complexity. In: Jantke, K.P. (ed.) *All 1989*. LNCS, vol. 397, pp. 18–44. Springer, Heidelberg (1989)
26. Angluin, D.: Queries and concept learning. *Mach. Learn.* **2**(4), 319–342 (1988)