

UNIVERSITY OF SÃO PAULO
INSTITUTE OF MATHEMATICS AND STATISTICS
BACHELOR OF COMPUTER SCIENCE

Tractable Probabilistic Description Logic
Algorithms and Implementation

Andrew Ijano Lopes

FINAL ESSAY
MAC 499 — CAPSTONE PROJECT

Program: Computer Science

Advisor: Prof. Dr. Marcelo Finger

São Paulo
January 20th, 2021

Tractable Probabilistic Description Logic
Algorithms and Implementation

Andrew Ijano Lopes

This is the original version of the
capstone project report prepared by
the candidate Andrew Ijano Lopes, as
submitted to the Examining Committee.

Autorizo a reprodução e divulgação total ou parcial deste trabalho, por qualquer meio convencional ou eletrônico, para fins de estudo e pesquisa, desde que citada a fonte.

Resumo

Andrew Ijano Lopes. **Lógica de Descrição Probabilística Tratável: *Algoritmos e Implementação***. Monografia (Bacharelado). Instituto de Matemática e Estatística, Universidade de São Paulo, São Paulo, 2020.

Elemento obrigatório, constituído de uma sequência de frases concisas e objetivas, em forma de texto. Deve apresentar os objetivos, métodos empregados, resultados e conclusões. O resumo deve ser redigido em parágrafo único, conter no máximo 500 palavras e ser seguido dos termos representativos do conteúdo do trabalho (palavras-chave). Deve ser precedido da referência do documento.

Palavras-chave: Lógicas de descrição. Palavra-chave2. Palavra-chave3.

Abstract

Andrew Ijano Lopes. **Tractable Probabilistic Description Logic: *Algorithms and Implementation***. Capstone Project Report (Bachelor). Institute of Mathematics and Statistics, University of São Paulo, São Paulo, 2020.

Elemento obrigatório, elaborado com as mesmas características do resumo em língua portuguesa. De acordo com o Regimento da Pós-Graduação da USP (Artigo 99), deve ser redigido em inglês para fins de divulgação. É uma boa ideia usar o sítio www.grammarly.com na preparação de textos em inglês.

Keywords: Description logics. Keyword2. Keyword3.

Lista de Abreviaturas

DL Description Logic

Lista de Símbolos

List of Figures

List of Tables

2.1 Syntax and semantics of \mathcal{EL}^{++} without concrete domains	4
--	---

List of Programs

Contents

1	Introduction	1
2	Background	3
2.1	The description logic \mathcal{EL}^{++}	3
2.1.1	Syntax	3
2.1.2	Semantics	4
2.1.3	Normal form	4
2.2	Graphic \mathcal{EL}^{++} (\mathcal{GEL}^{++})	5
2.3	MaxSAT for \mathcal{GEL}^{++}	5
2.4	Probabilistic \mathcal{GEL}^{++}	5
2.5	Related Work	5
3	Development	7
3.1	Input and output format	7
3.2	OWL parser	7
3.3	Knowledge Base	7
3.4	GEL-MaxSAT	7
3.5	Linear solver	7
3.6	PGEL-SAT reasoner	7
4	Experiments	9
5	Results	11
6	Conclusion	13
	References	15

Chapter 1

Introduction

Description logics are a family of formal knowledge representation languages, being of particular importance in providing a logical formalism for ontologies and the Semantic Web. Also, they are notable in biomedical informatics for assisting the codification of biomedical knowledge. Due to these uses, there is a great demand to find tractable (i.e., polynomial-time decidable) description logics.

One of them, the logic \mathcal{EL}^{++} , is one of the most expressive description logics in which the complexity of inferential reasoning is tractable (BAADER *et al.*, 2005). Even though it is expressive enough to deal with several practical applications, there was also a need to model uncertain knowledge.

Example 1.1

Consider a medical situation, in which a patient may have non-specific symptoms, such as high fever, cough and headache. Also, COVID-19, a severe acute respiratory syndrome caused by the SARS-CoV-2 virus, is a disease that can account for those symptoms, but not all patients present all symptoms. Such an uncertain situation is suitable for probabilistic modeling.

In a certain hospital, a patient with high fever has some probability of having COVID-19, but that probability is 20% larger if the patient has cough too. On the other hand, COVID-19 is not very prevalent and is not observed in the hospital 90% of the time. If those probabilistic constraints are satisfiable, one can also ask the minimum and maximum probability that a hospital patient Mary, with fever and cough, is a suspect of suffering from COVID-19.

For classical propositional formulas, this problem, called *probabilistic satisfiability* (PSAT), has already been presented with tractable fragments (ANDERSEN and PRETOLANI, 2001). On the other hand, in description logics, most studies result in intractable reasoning; moreover, by adding probabilistic reasoning capabilities to \mathcal{EL}^{++} , in order to model such situation, the complexity reaches NP-completeness (FINGER, 2019).

To solve this problem, probabilistic constraints can be applied to axioms and its probabilistic satisfaction can be seen in a linear algebraic view. Furthermore, it can be reduced to an optimization problem, which can be solved by an adaptation of the simplex method with

column generation (FINGER, 2019). Thus, it is possible to reduce the column generation problem to the *weighted partial maximum satisfiability*.

Moreover, recent studies show that it is necessary to focus on a fragment of \mathcal{EL}^{++} for obtain a tractable probabilistic reasoning, which will be called Graphic \mathcal{EL}^{++} (\mathcal{GEL}^{++}) (FINGER, n.d.). Therefore, this fragment allows axioms to be seen as edges in a graph, as opposed to hyperedges in a hypergraph, which is the case of \mathcal{EL}^{++} . This allows the use of graph-based machinery to develop tractable algorithm for the *weighted partial Maximum SATisfiability* for \mathcal{GEL}^{++} (Max \mathcal{GEL}^{++} -SAT) and, as a result, a tractable probabilistic description logic.

Then, the objective of this project is to propose and implement tractable algorithms for weighted partial Max-SAT and Probabilistic SAT for a fragment of \mathcal{EL}^{++} description logic.

In this paper, we describe the implementation of these algorithms¹ and is organized as follows: Section ?? highlights related results in the literature. The basic definition of \mathcal{GEL}^{++} with its algorithms for MaxSAT and PSAT are described in Section ?? and followed by Section ??, which presents details about the implementation and its experimental evaluation.

¹Available at <https://github.com/AndrewIjano/pgel-sat>

Chapter 2

Background

In this chapter, we present the theoretical background of DLs and some definitions.

Add initial description of the chapter

2.1 The description logic \mathcal{EL}^{++}

\mathcal{EL}^{++} is an extension of the DL \mathcal{EL} (BAADER *et al.*, 2005). It was created with large biohealth ontologies in mind, such as SNOMED-CT, the NCI thesaurus and Galen, and became an official OWL 2 profile (HITZLER *et al.*, 2009). We concentrate on presenting \mathcal{EL}^{++} without concrete domains.

2.1.1 Syntax

In \mathcal{EL}^{++} , *concept descriptions* are defined inductively from a set N_C of *concept names*, a set N_R of *role names* and set N_I of *individual names* as follows:

- \top , \perp and concept names in N_C are concept descriptions;
- if C and D are concept descriptions, $C \sqcap D$ is a concept description;
- if C is a concept description and $r \in N_R$, $\exists r.C$ is a concept description;
- if $a \in N_I$, $\{a\}$ is a concept description.

An *axiom*, or a *general concept inclusion* (GCI), is an expression of the form $C \sqsubseteq D$, where C and D are concept inclusions. A *role inclusion* (RI) is an expression of the form $r_1 \circ \dots \circ r_k \sqsubseteq r$, where $r_1, \dots, r_k, r \in N_R$. The symbol “ \circ ” denotes composition of binary relations. We call a general *terminology box* (TBox) a finite set of GCIs. Also, a *constraint box* (CBox) is a finite set of GCIs and a finite set of RIs.

Improve the structure of this paragraph, making it more readable

Similarly, a *concept assertion* is an expression of the form $C(a)$ and a *role assertion*, $r(a, b)$, where C is a concept description, $a, b \in N_I$ and $r \in N_R$. A finite set of concept assertions and role assertions is an *assertional box* (ABox).

Then, an \mathcal{EL}^{++} *knowledge base* \mathcal{K} (KB) is a pair (C, \mathcal{A}) , where C is a CBox and \mathcal{A} is an ABox.

2.1.2 Semantics

The semantics of \mathcal{EL}^{++} are given by *interpretations* $\mathcal{I} = (\Delta^{\mathcal{I}}, \cdot^{\mathcal{I}})$. The *domain* $\Delta^{\mathcal{I}}$ is a non-empty set of individuals, and the *interpretation function* $\cdot^{\mathcal{I}}$ maps each concept name $A \in \mathbf{N}_C$ to a subset $A^{\mathcal{I}}$ of $\Delta^{\mathcal{I}}$, each role name r to a binary relation $r^{\mathcal{I}} \subseteq \Delta^{\mathcal{I}} \times \Delta^{\mathcal{I}}$, and each individual name a to an element $a^{\mathcal{I}} \in \Delta^{\mathcal{I}}$. The extension of $\cdot^{\mathcal{I}}$ for an arbitrary concept description is inductively defined by the third column of Table 2.1.

Name	Syntax	Semantics
top	\top	$\Delta^{\mathcal{I}}$
bottom	\perp	\emptyset
nominal	$\{a\}$	$\{a^{\mathcal{I}}\}$
conjunction	$C \sqcap D$	$C^{\mathcal{I}} \cap D^{\mathcal{I}}$
existential restriction	$\exists r.C$	$\{x \in \Delta^{\mathcal{I}} \mid \exists y \in \Delta^{\mathcal{I}} : (x, y) \in r^{\mathcal{I}} \wedge y \in C^{\mathcal{I}}\}$
GCI	$C \sqsubseteq D$	$C^{\mathcal{I}} \subseteq D^{\mathcal{I}}$
RI	$r_1 \circ \dots \circ r_k \sqsubseteq r$	$r_1^{\mathcal{I}} \circ \dots \circ r_k^{\mathcal{I}} \subseteq r^{\mathcal{I}}$
concept assertion	$C(a)$	$a^{\mathcal{I}} \in C^{\mathcal{I}}$
role assertion	$r(a, b)$	$(a^{\mathcal{I}}, b^{\mathcal{I}}) \in r^{\mathcal{I}}$

Table 2.1: Syntax and semantics of \mathcal{EL}^{++} without concrete domains

The interpretation \mathcal{I} satisfies:

- an axiom $C \sqsubseteq D$ if $C^{\mathcal{I}} \subseteq D^{\mathcal{I}}$ (represented as $\mathcal{I} \models C \sqsubseteq D$);
- a RI $r_1 \circ \dots \circ r_k \sqsubseteq r$ if $r_1^{\mathcal{I}} \circ \dots \circ r_k^{\mathcal{I}} \subseteq r^{\mathcal{I}}$ (represented as $\mathcal{I} \models r_1 \circ \dots \circ r_k \sqsubseteq r$);
- an assertion $C(a)$ if $a^{\mathcal{I}} \in C^{\mathcal{I}}$ (represented as $\mathcal{I} \models C(a)$);
- an assertion $r(a, b)$ if $(a^{\mathcal{I}}, b^{\mathcal{I}}) \in r^{\mathcal{I}}$ (represented as $\mathcal{I} \models r(a, b)$).

Also, we say that \mathcal{I} is a *model* of:

- a CBox \mathcal{C} if it satisfies every axiom and RI in \mathcal{C} (represented as $\mathcal{I} \models \mathcal{C}$);
- an ABox \mathcal{A} if it satisfies every assertion in \mathcal{A} (represented as $\mathcal{I} \models \mathcal{A}$);

Then, an important problem in \mathcal{EL}^{++} is to determine its *consistency*, that is if \mathcal{A} and \mathcal{C} have a common model, which is in PTime (BAADER *et al.*, 2005).

2.1.3 Normal form

We can convert an \mathcal{EL}^{++} knowledge base into a normal form, in polynomial time, by introducing new concept and role names (BAADER *et al.*, 2005).

First, there is no need of explicit ABox, because $\mathcal{I} \models C(a) \iff \mathcal{I} \models \{a\} \sqsubseteq C$ and $\mathcal{I} \models r(a, b) \iff \{a\} \sqsubseteq \exists r.\{b\}$. In other words, a knowledge base can be represented by just a CBox, by transforming assertions in axioms.

In addition, given a CBox \mathcal{C} , consider the set $\text{BC}_{\mathcal{C}}$ of *basic concept descriptions*, which is the smallest set of concept descriptions that contains the top concept \top , all concept names used in \mathcal{C} and all concepts of the form $\{a\}$ used in \mathcal{C} .

Then, every axiom can be represented in the following normal form, where $C_1, C_2 \in \text{BC}_{\mathcal{C}}$, $D \in \text{BC}_{\mathcal{C}} \cup \{\perp\}$:

1. $C_1 \sqsubseteq D$
2. $C_1 \sqsubseteq \exists r.C_2$
3. $C_2 \sqsubseteq D$
4. $C_1 \sqcap C_2 \sqsubseteq D$

And every RI are of the form $r \sqsubseteq s$ or $r_1 \circ r_2 \sqsubseteq s$.

Example 2.1

Consider the following CBox representing the situation in [Example 1.1](#). In the left we have basic knowledge on diseases and, in the right, the specific knowledge about Mary. Note that, for simplicity, it is not in normal form.

Fever \sqsubseteq Symptom	
Cough \sqsubseteq Symptom	
COVID-19 \sqsubseteq Disease	$\{\text{mary}\} \sqsubseteq \text{Patient}$
Symptom $\sqsubseteq \exists \text{hasCause.Disease}$	$\{s1\} \sqsubseteq \text{Fever} \sqcap \text{Cough}$
Patient $\sqsubseteq \exists \text{hasSymptom.Symptom}$	$\{\text{mary}\} \sqsubseteq \exists \text{hasSymptom}.\{s1\}$
CovidPatient $\equiv \exists \text{suspectOf.COVID-19}$	
$\text{hasSymptom} \circ \text{hasCause} \sqsubseteq \text{suspectOf}$	

Because CBoxes can only represent facts, there is no way to describe uncertain knowledge. Even though, $A_{x_1} := \text{Fever} \sqsubseteq \exists \text{hasCause.COVID-19}$, $A_{x_1} := \text{Fever} \sqcap \text{Cough} \sqsubseteq \exists \text{hasCause.COVID-19}$ and $\text{COVID-19} \sqsubseteq \perp$

We want to model uncertain information using DLs. However, it was been proved that, by adding probabilistic reasoning capabilities to \mathcal{EL}^{++} , the complexity reaches NP-completeness ([FINGER, 2019](#)). Then, it is necessary to reduce the expressiveness of this language.

2.2 Graphic \mathcal{EL}^{++} (\mathcal{GEL}^{++})

2.3 MaxSAT for \mathcal{GEL}^{++}

2.4 Probabilistic \mathcal{GEL}^{++}

2.5 Related Work

Chapter 3

Development

3.1 Input and output format

3.2 OWL parser

3.3 Knowledge Base

3.4 GEL-MaxSAT

3.5 Linear solver

3.6 PGEL-SAT reasoner

Chapter 4

Experiments

Chapter 5

Results

Chapter 6

Conclusion

References

- [ANDERSEN and PRETOLANI 2001] Kim Allan ANDERSEN and Daniele PRETOLANI. “Easy cases of probabilistic satisfiability”. In: *Annals of Mathematics and Artificial Intelligence* 33.1 (2001), pp. 69–91 (cit. on p. 1).
- [BAADER *et al.* 2005] Franz BAADER, Sebastian BRANDT, and Carsten LUTZ. “Pushing the EL Envelope”. In: *Proceedings of the 19th International Joint Conference on Artificial Intelligence*. IJCAI’05. 2005, pp. 364–369 (cit. on pp. 1, 3, 4).
- [FINGER 2019] Marcelo FINGER. “Extending EL++ with Linear Constraints on the Probability of Axioms”. In: *Description Logic, Theory Combination, and All That. Essays Dedicated to Franz Baader on the Occasion of His 60th Birthday*. Ed. by Carsten LUTZ, Uli SATTLER, Cesare TINELLI, Anni-Yasmin TURHAN, and Frank WOLTER. Vol. LLNCS 11560. Theoretical Computer Science and General Issues. Springer International Publishing, 2019. ISBN: 978-3-030-22101-0. DOI: [10.1007/978-3-030-22102-7](https://doi.org/10.1007/978-3-030-22102-7) (cit. on pp. 1, 2, 5).
- [FINGER n.d.] Marcelo FINGER. “Tractable Max-SAT in Graphic Description Logics”. In preparation (cit. on p. 2).
- [HITZLER *et al.* 2009] Pascal HITZLER, Markus KRÖTZSCH, Bijan PARSIA, Peter F PATEL-SCHNEIDER, Sebastian RUDOLPH, *et al.* “OWL 2 web ontology language primer”. In: *W3C recommendation* 27.1 (2009), p. 123 (cit. on p. 3).