

Using Grammar-based Genetic Programming for Mining Disjointness Axioms Involving Complex Class Expressions

Thu Huong Nguyen Andrea G. B. Tettamanzi

Université Côte d'Azur, CNRS, Inria, I3S, Sophia Antipolis, France

Symposium MaDICS 2020, July 6, 2020



UNIVERSITÉ CÔTE D'AZUR



Motivation

- Linked Open Data (LOD) has made a huge number of interconnected RDF triples freely available for sharing and reuse
- Shared schemas and ontologies are needed to support reasoning
- Manual acquisition of axioms:
 - is exceedingly expensive and time-consuming
 - depends on the availability of domain specialists and knowledge engineers

Introduction: Ontology Learning

Top-down construction of ontologies has limitations

- aprioristic and dogmatic
- does not scale well
- does not lend itself to a collaborative effort

Bottom-up, *grass-roots* approach to ontology and KB creation

- start from RDF facts and learn OWL 2 axioms

Recent contributions towards OWL 2 ontology learning

- FOIL-like algorithms for learning concept definitions
- statistical schema induction via association rule mining
- light-weight schema enrichment (DL-Learner framework)

All these methods apply and extend ILP techniques.

Class Disjointness Axiom Learning

Problem:

We focus on learning OWL class disjointness axioms involving existential quantification ($\exists r.C$) and value restriction ($\forall r.C$)

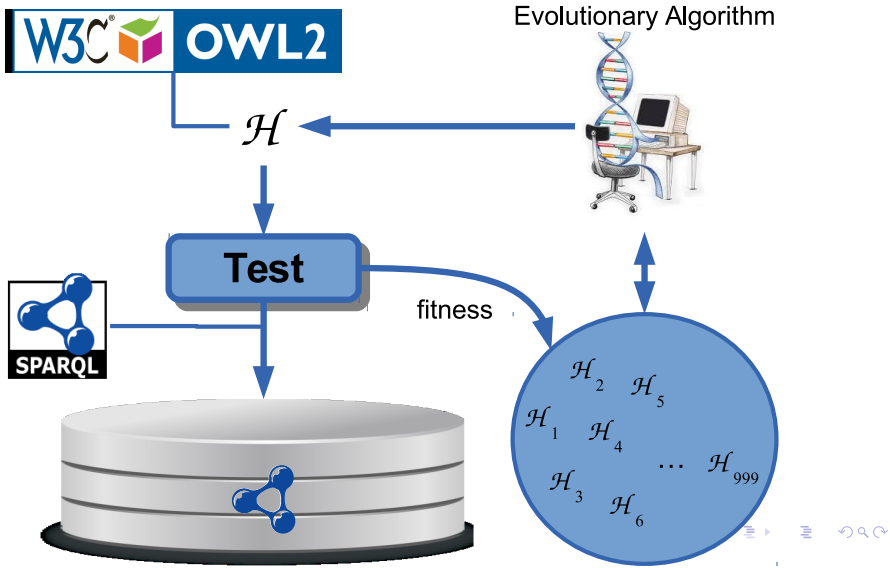
Class disjointness axioms are important

- to check the correctness of a knowledge base
- to derive new knowledge

Method:

Grammar-Based Genetic Programming

Synopsis



Grammatical Evolution

- A grammar-based form of Genetic Programming
- Search space distinguished from solution (= program) space
- Search space consists of variable-length bitstrings
- Bitstrings are mapped into programs thanks to a BNF *grammar*

In the mapping process, codons are used consecutively to choose production rules in the BNF grammar according to the function:

$$production = codon \bmod \left[\begin{array}{l} \text{Number of productions} \\ \text{for the current non-} \\ \text{terminal} \end{array} \right]$$

Static Part of the Grammar

Axiom := ClassAxiom

ClassAxiom := DisjointClasses

DisjointClasses := 'DisjointClasses' '(' ClassExpression1 ' ' ClassExpression2 ')'

ClassExpression1 := Class

| ObjectSomeValuesFrom

| ObjectAllValuesFrom

| ObjectIntersection

ClassExpression2 := ObjectSomeValuesFrom

| ObjectAllValuesFrom

ObjectIntersectionOf := 'ObjectIntersectionOf' '(' Class ' ' Class ')'

ObjectSomeValuesFrom := 'ObjectSomeValuesFrom' '(' ObjectPropertyOf ' ' Class ')'

ObjectAllValuesFrom := 'ObjectAllValuesFrom' '(' ObjectPropertyOf ' ' Class ')'

Dynamic Part of the Grammar

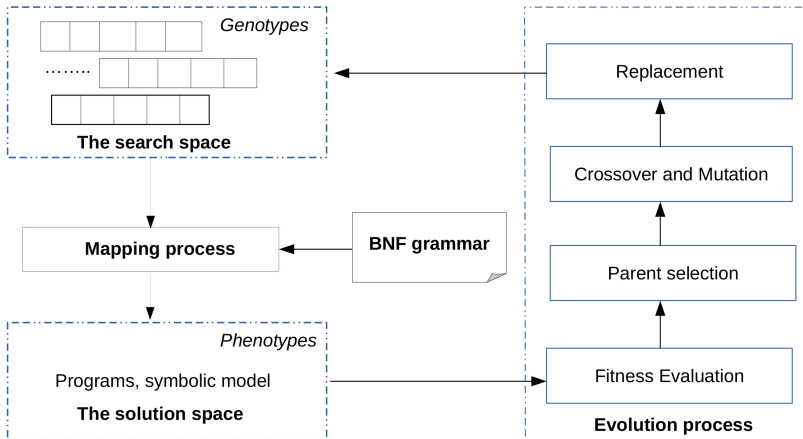
Class := ...

```
SELECT ?class
WHERE { ?instance rdf:type ?class . }
```

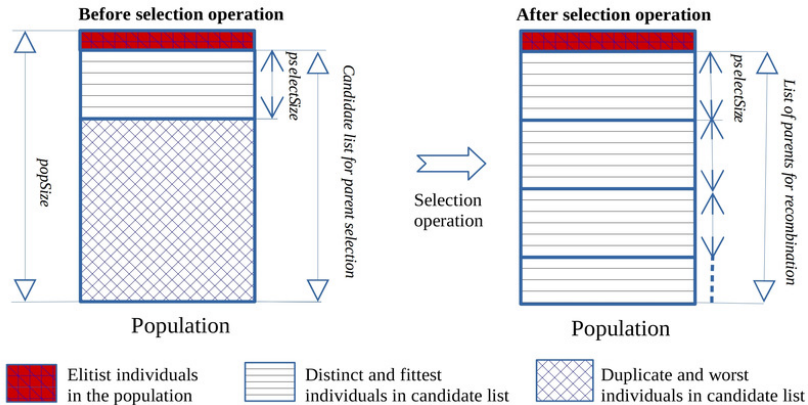
ObjectPropertyOf := ...

```
SELECT ?property
WHERE {
  ?subject ?property ?object .
  FILTER (isIRI(?object))
}
```


Grammatical Evolution Process



Parent Selection



Possibility Theory

Definition (Possibility Distribution)

$$\pi : \Omega \rightarrow [0, 1]$$

Definition (Possibility and Necessity Measures)

$$\Pi(A) = \max_{\omega \in A} \pi(\omega);$$

$$N(A) = 1 - \Pi(\bar{A}) = \min_{\omega \in \bar{A}} \{1 - \pi(\omega)\}.$$

For all subsets $A \subseteq \Omega$,

- ① $\Pi(\emptyset) = N(\emptyset) = 0$, $\Pi(\Omega) = N(\Omega) = 1$;
- ② $\Pi(A) = 1 - N(\bar{A})$ (duality);
- ③ $N(A) > 0$ implies $\Pi(A) = 1$, $\Pi(A) < 1$ implies $N(A) = 0$.

In case of complete ignorance on A , $\Pi(A) = \Pi(\bar{A}) = 1$.

Content of an Axiom

Definition (Content of Axiom ϕ)

Given an RDF dataset \mathcal{K} ,

$$\text{content}(\phi) = \{\psi : \phi \models_{\mathcal{K}} \psi\}$$

obtained through the instantiation of ψ to the vocabulary of \mathcal{K} .

Let $\phi = \text{Dis}(C, D)$

$$\text{content}(\phi) = \{\neg C(r) \vee \neg D(r) : r \text{ is a resource in } \mathcal{K}\}$$

Confirmation and Counterexample of an Axiom

Given $\psi \in \text{content}(\phi)$ and an RDF dataset \mathcal{K} , three cases:

- ① $\mathcal{K} \models \psi$: $\rightarrow \psi$ is a *confirmation* of ϕ ;
- ② $\mathcal{K} \models \neg\psi$: $\rightarrow \psi$ is a *counterexample* of ϕ ;
- ③ $\mathcal{K} \not\models \psi$ and $\mathcal{K} \not\models \neg\psi$: $\rightarrow \psi$ is neither of the above

Definition

Given axiom ϕ , let us define

$u_\phi = \|\text{content}(\phi)\|$ (a.k.a. the *support* of ϕ)

u_ϕ^+ = the number of confirmations of ϕ

u_ϕ^- = the number of counterexamples of ϕ

Axiom Evaluation

Definition (Generality)

$g_\phi = \min\{\| [C] \|, \| [D] \| \}$, where C, D are class expressions.

Definition (Possibility)

$$\Pi(\phi) = 1 - \sqrt{1 - \left(\frac{u_\phi - u_\phi^-}{u_\phi} \right)^2}$$

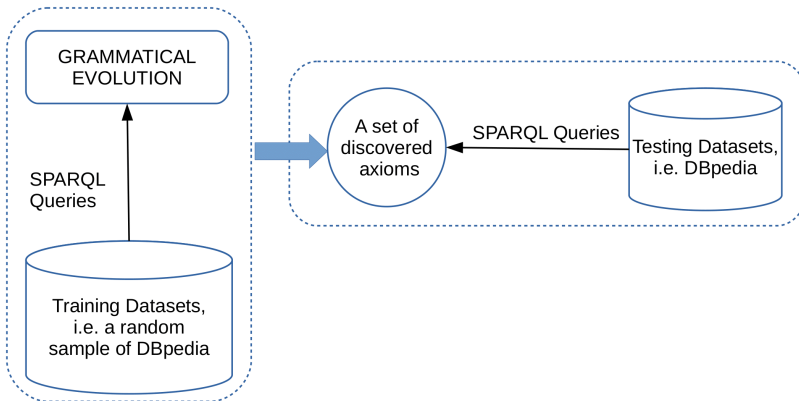
Definition (Fitness)

$$f(\phi) = g_\phi \cdot \Pi(\phi)$$

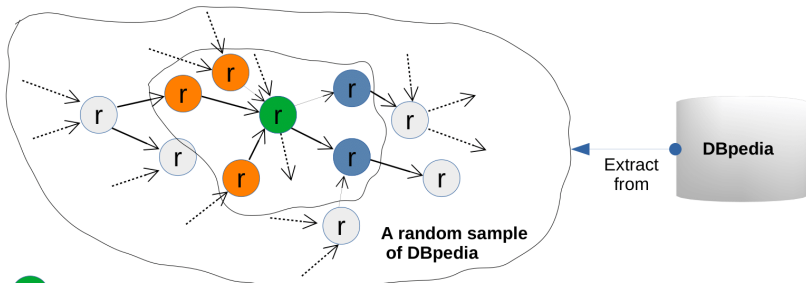
Experimental Protocol

- Training-Testing Model
- Experiments are divided into two phases:
 - ① mining class disjointness axioms with GE from a training RDF dataset, i.e., a 1% random sample of DBpedia 2015-04
 - ② testing the resulting axioms against the test dataset, i.e., the entire DBpedia 2015-04
- An objective benchmark to evaluate the effectiveness of the method.

Training-Testing Model



Training Set Construction



- r** Initial resource, e.g. <http://dbpedia.org/ontology/Plant>
- r** Resource as a subject in the triples relevant to **r**, i.e. **r** is the object of the extracted triple
- r** Resource as an object in the triples relevant to **r**, i.e. **r** is the subject of the extracted triple
- \longrightarrow Relation between resources as the predicate in the triples
- $\cdots \longrightarrow$

Experimental Setup

- 20 different runs on different parameter settings
- To allow fair comparisons, we define total effort

k = total number of fitness evaluations

- *maxGenerations* is set so that

$$\text{popSize} \cdot \text{maxGenerations} = k$$

Parameter	Value
Total effort k	100,000; 200,000; 300,000; 400,000
<i>initLenChrom</i>	6
<i>pCross</i>	80%
<i>pMut</i>	1%
<i>popSize</i>	1000; 2000; 5000; 10000

Measuring Accuracy

Since $\Pi(\phi)$ may be viewed as a fuzzy degree of membership, we use a fuzzy extension of the usual definition of *precision*, based on fuzzy set cardinality

$$\|F\| = \sum_{x \in \Delta} F(x),$$

The value of precision can thus be computed against the test dataset, i.e., DBpedia 2015-04, according to the formula

$$\text{precision} = \frac{\|\text{true positives}\|}{\|\text{discovered axioms}\|} = \frac{\sum_{\phi} \Pi_{\text{DBpedia}}(\phi)}{\sum_{\phi} \Pi_{\text{Training}}(\phi)}$$

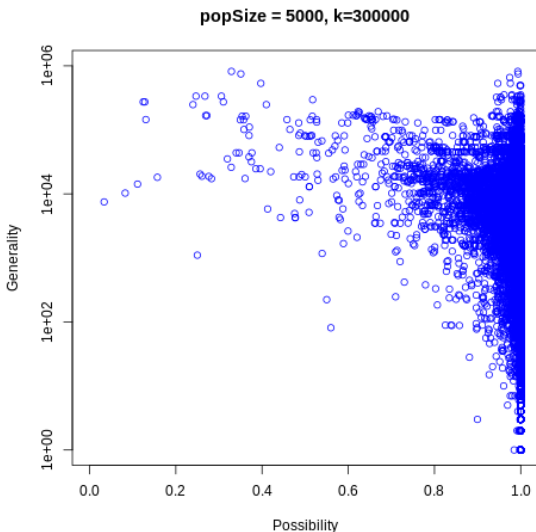
Axioms Discovered Over 20 Runs

popSize k	1000	2000	5000	10000
100000	8806	11389	4684	4788
200000	6204	13670	10632	9335
300000	5436	10541	53021	14590
400000	5085	9080	35102	21670

Average Precision per Run (\pm std)

$k \backslash popSize$	1,000	2,000	5,000	10,000
100,000	0.981 ± 0.019	0.999 ± 0.002	0.998 ± 0.002	0.998 ± 0.003
200,000	0.973 ± 0.024	0.979 ± 0.011	0.998 ± 0.001	0.998 ± 0.002
300,000	0.972 ± 0.024	0.973 ± 0.014	0.993 ± 0.007	0.998 ± 0.001
400,000	0.972 ± 0.024	0.969 ± 0.018	0.980 ± 0.008	0.998 ± 0.001

Π and g Distribution of Discovered Axioms



Examples of Discovered Axioms

Dis(\forall author.Place, \forall placeofBurial.Place) $\Pi(\phi) = 1.0$; $g_\phi = 4$

Dis(Writer, \forall writer.Agent) $\Pi(\phi) = 0.982$; $g_\phi = 79,464$

Dis(Journalist, \forall distributor.Agent) $\Pi(\phi) = 0.992$; $g_\phi = 32,533$

Dis(Stadium, \forall birthPlace.Place) $\Pi(\phi) = 1.0$; $g_\phi = 10,245$

Conclusions and Future Work

- Grammar-based GP method for mining disjointness axioms involving complex class expressions
- The use of a training-testing model allows to objectively validate the method, while also alleviating the computational bottleneck of SPARQL endpoints
- The experimental results confirm that the proposed method is capable of discovering highly accurate and general axioms
- Future Work:
 - Mining disjointness axioms involving operators such as `owl:hasValue` and `owl:OneOf`
 - Forbid atomic classes at the root of class expressions
 - Refining the evaluation of candidate axioms with some measure of their complexity

The End

Thank you for your attention!