

# Data Search and Discovery in RDF Sources

Zoé Chevallier<sup>1,2</sup>, Zoubida Kedad<sup>1</sup>, Béatrice Finance<sup>1</sup>, and Frédéric Chaillan<sup>2</sup>

<sup>1</sup> David Lab. University of Versailles Paris-Sclay, France

`firstname.lastname@uvsq.fr`

<sup>2</sup> Grand Paris Sud, France

`f{firstname}.lastname@grandparissud.fr`

**Abstract.** The RDF data sources published on the web represent an unprecedented amount of knowledge. However, querying these sources to extract the relevant information for some specific needs represented by a target schema is a complex task, as the alignment between the target and the source schemas might not be provided or may be incomplete. This paper presents a system that aims to automatically populate the classes or a target schema from RDF data sources by identifying candidate instance patterns. This identification process relies on a semi-supervised learning algorithm and the systems automatically generates the SPARQL queries that populate the target schema.

**Keywords:** RDF data sources · Target Schema Instantiation · SPARQL Query Generation · Semi-supervised learning.

## 1 Introduction

The web represents a huge space of available data from which various applications can extract meaningful knowledge. However, finding relevant data for some specific need is not obvious, especially for irregular data sources. This problem has been addressed by dataset discovery approaches [1, 3], which aim at discovering the relevant datasets that could complement a given target dataset. These approaches are designed for structured datasets.

Considering that the specific needs of an application are described by a target schema, our problem is the identification and the extraction of relevant data to populate this target schema. A similar problem is the one of mapping generation which has been the subject of several works [2, 4], tarteting relational or XML data. Sacramento et al. [5] have addressed the problem of expressing an RDF data source in the terms defined by an ontology, which consists in generating a mapping between a source and this ontology. This requires the alignment between the source schema and the target schema, which is not always provided.

In this paper, we present a system that identifies candidate instances from RDF data sources to populate a given target schema. It relies on a semi-supervised learning algorithm to extract candidate instance patterns from an RDF data source and automatically generates the queries that extract these instances. This paper is organized as follows. Section 2 presents the architecture of our system. Section 3 details a use-case scenario, and section 4 presents some future works.

## 2 System Architecture

Figure 1 depicts the architecture of our system. We consider that the target schema is described in RDFS/OWL. The alignments between the source schemas and the target schema can be exploited if they are available, but they are not required.

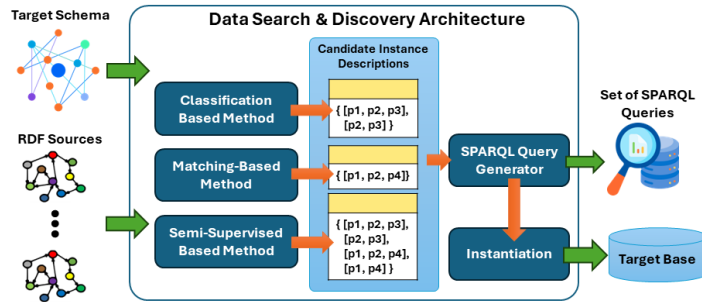


Fig. 1: The Data Search And Discovery Architecture

Our system identifies in RDF sources candidate instance patterns in order to populate the target schema. Each pattern is a property set which describes some candidate instances in the data source. Given a candidate pattern, a SPARQL Query is generated to retrieve all the candidate instances and automatically populate the target schema.

We assume that each class  $C$  of the target schema is described by a set of properties,  $Prop(C)$ , such that  $Prop(C) = \{p \mid \langle C, p, * \rangle \in T\}$ . An entity  $e$  in an RDF data source  $S$  is a resource that is neither a class, property, literal, or blank node, and is characterized by a set of properties,  $Prop(e)$ , such that:  $Prop(e) = \{p \mid \langle e, p, * \rangle \in S\}$ .

One straightforward way to identify candidate instance patterns is to compute the similarity between source entities and a target class  $C$ . If the similarity is higher than a threshold, then  $Prop(e)$  is a candidate pattern for  $C$ . We refer to this process as *classification-based*. If the output of some schema matching tool is provided, and if a target class  $C_T$  is equivalent to a source class  $C_S$ , then all the instances of  $C_S$  are candidate instances for  $C_T$ . We refer to this process as *matching-based* instantiation.

However, these correspondences are not always provided. Besides, the schema in RDF sources is only descriptive, and an instance can be characterized by a property set that is different from the one defined in the schema for its class. Therefore, new candidate instance patterns should be identified based on the class patterns in the target schema, but also based on the candidate instance descriptions already identified for a class. To do so, we propose an approach based on a semi-supervised learning algorithm [6] to identify candidate patterns.

An entity similar to a candidate pattern already identified could also be a candidate instance for  $C$ . We proposed an approach that iteratively computes the similarity of each entity to the candidate patterns already identified. If this similarity between  $e$  and the candidate patterns of  $C$  is higher than a threshold, then  $Prop(e)$  is considered as a candidate pattern of  $C$ . We introduced the notion of candidate instance descriptions, that represent the existing patterns of candidate instances for a target class  $C$  such that:

$$Desc_{CI}(C) = \{CID_i(C) \mid CID_i(C) = Prop(e), e \text{ candidate instance of } C\} \quad (1)$$

The similarity function based on the candidate instances between  $e$  and a target class  $C$  is such that:

$$Sim_I(e, C) = MAX\left(\frac{|Prop(e) \cap CID_i(C)|}{|Prop(e) \cup CID_i(C)|}, \text{ where } CID_i(C) \in Desc_{CI}(C)\right) \quad (2)$$

The semi-supervised algorithm first computes candidate instance descriptions based on matching-based and classification-based approaches. For a target class  $C$ , the similarity between each entity  $e$  and  $C$ ,  $Sim_I$  is computed. If the similarity is higher than a given threshold, then  $Prop(e)$  is a candidate instance description of  $C$ . If new candidate instance descriptions are identified, a new iteration is started. This identification process iterates until no more new candidate instance descriptions are found.

### 3 Demonstration Scenario

Our system is implemented in Java, using Apache Jena. An online presentation of the demonstration is available<sup>3</sup>. The demonstration will proceed in four phases:

1. **Configuration.** The process starts by selecting a project composed on a RDFS/OWL target schema, RDF sources, and optionally a set of correspondences between the target schema and the RDF sources. A graphical representation of both the target schema and the data sources is presented.
2. **Identification of Candidate Instance Descriptions.** In this phase, we will show the generation of candidate instance descriptions using different approaches. The candidate instance descriptions identified for each target class will be presented as well as their provenance, i.e. the data source from which it has been extracted. Moreover, a visualization tool allow the user to highlight the candidate instances identified in the source.
3. **Generation of Sparql Queries.** We will then demonstrates the query generation process, and show the queries generated for each candidate instance description and the execution result depicted in figure 2.
4. **Comparison of the Candidate Instance Sets.** In order to highlight the quality of the semi-supervised approach, we will compare its resulting candidate instance sets with the ones extracted using a baseline method, which consists in the union of the candidate instance descriptions obtained using both the matching-based approach and the classification-based approach.

<sup>3</sup> [https://mega.nz/folder/QK01RSha#dSW\\_77ur0QBPNp7ff4\\_qXg](https://mega.nz/folder/QK01RSha#dSW_77ur0QBPNp7ff4_qXg)

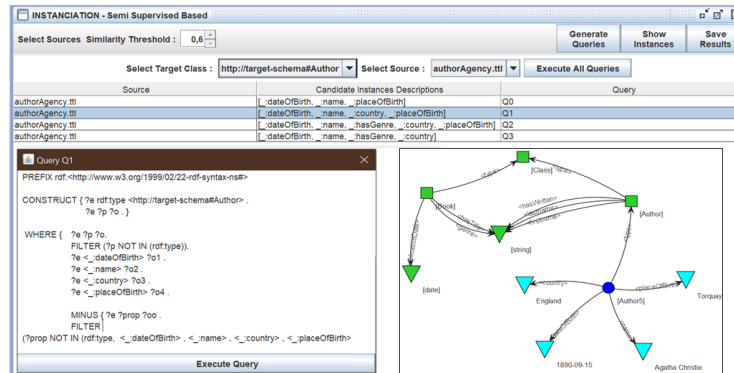


Fig. 2: A Screenshot of the Query Generated from a Candidate Description

## 4 Future Works

In future works, we will improve our data discovery system by taking into account some constraints defined on the target schema so as to extract only the source instances for which the constraints are verified. These constraints could be defined using the SHACL language<sup>4</sup>.

## References

1. Castro Fernandez, R., Abedjan, Z., Koko, F., Yuan, G., Madden, S., Stonebraker, M.: Aurum: A Data Discovery System. In: 2018 IEEE 34th International Conference on Data Engineering (ICDE). pp. 1001–1012. IEEE, Paris (Apr 2018)
2. Fagin, R., Haas, L.M., Hernandez, M., Miller, R.J., Popa, L., Velegarakis, Y.: Clio: Schema Mapping Creation and Data Exchange | SpringerLink. In: Conceptual Modeling: Foundations and Applications, Lecture Notes in Computer Science, vol. 5600, pp. 198–236 (2009)
3. Koutras, C., Siachamis, G., Ionescu, A., Psarakis, K., Brons, J., Fragkoulis, M., Lofi, C., Bonifati, A., Katsifodimos, A.: Valentine: Evaluating Matching Techniques for Dataset Discovery. In: 2021 IEEE 37th International Conference on Data Engineering (ICDE). pp. 468–479 (Apr 2021)
4. Mazilu, L., Paton, N.W., Fernandes, A.A., Koehler, M.: Schema mapping generation in the wild. *Information Systems* **104**, 101904 (Feb 2022)
5. Sacramento, E.R., Vidal, V.M.P., de Macêdo, J.A.F., Lóscio, B.F., Lopes, F.L.R., Casanova, M.A.: Towards automatic generation of application ontologies. *J. Inf. Data Manag.* **1**(3), 535–550 (2010)
6. Yarowsky, D.: Unsupervised word sense disambiguation rivaling supervised methods. In: Proceedings of the 33rd annual meeting. Association for Computational Linguistics, Cambridge, Massachusetts (1995)

<sup>4</sup> <https://www.w3.org/TR/shacl/>