



**HAL**  
open science

# A Semantic Measure for Outlier Detection in Knowledge Graph

Bara Diop, Cheikh Talibouya Diop, Lamine Diop

► **To cite this version:**

Bara Diop, Cheikh Talibouya Diop, Lamine Diop. A Semantic Measure for Outlier Detection in Knowledge Graph. *Revue Africaine de Recherche en Informatique et Mathématiques Appliquées*, 2022, Volume 35, Data Intelligibility, Business Intelligence and Semantic Web, 10.46298/arima.8679 . hal-03415728v3

**HAL Id: hal-03415728**

**<https://inria.hal.science/hal-03415728v3>**

Submitted on 9 Apr 2022

**HAL** is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.

## A Semantic Measure for Outlier Detection in Knowledge Graph

Bara DIOP<sup>1</sup>, Cheikh Talibouya DIOP<sup>1</sup>, Lamine DIOP<sup>2</sup>

<sup>1</sup>University Gaston Berger of Saint-Louis, Senegal

<sup>2</sup>University of Tours, France

\*E-mail : {diop.bara,cheikh-talibouya.diop}@ugb.edu.sn, lamine.diop@univ-tours.fr

DOI : [10.46298/arima.8679](https://doi.org/10.46298/arima.8679)

Submitted on November 5, 2021 - Published on March 20, 2022

Volume : 35 - Year : 2022

Special Issue : **Volume 35, Data Intelligibility, Business Intelligence and Semantic Web**

Editors : Ghislain Atemezing, Gaoussou Camara, Bruce Watson

---

### Abstract

Nowadays, there is a growing interest in data mining and information retrieval applications from Knowledge Graphs (KG). However, the latter (KG) suffers from several data quality problems such as accuracy, completeness, and different kinds of errors. In DBpedia, there are several issues related to data quality. Among them, we focus on the following: several entities are in classes they do not belong to. For instance, the query to get all the entities of the class Person also returns group entities, whereas these should be in the class Group. We call such entities “outliers.” The discovery of such outliers is crucial for class learning and understanding. This paper proposes a new outlier detection method that finds these entities. We define a semantic measure that favors the real entities of the class (inliers) with positive values while penalizing outliers with negative values and improving it with the discovery of frequent and rare itemsets. Our measure outperforms FPOF (Frequent Pattern Outlier Factor) ones. Experiments show the efficiency of our approach.

### Keywords

Knowledge graph; Pattern Mining; Itemset; Outlier Detection

---

## I INTRODUCTION

Nowadays, the use of different Knowledge Graphs is becoming more and more important. Knowledge Graph is defined as follows: “A knowledge graph (i) mainly describes real world entities and their interrelations, organized in a graph, (ii) defines possible classes and relations of entities in a schema, (iii) allows for potentially interrelating arbitrary entities with each other and (iv) covers various topical domains” [19]. Knowledge Graphs are constructed from various data sources using various ways. For example, DBpedia is built from the free encyclopedia Wikipedia. In the Semantic Web, the term Knowledge Graph is often used to designate the Knowledge Graphs of the Semantic Web [20]. In these, Semantic Web principles are used to

represent and publish publicly available data as Linked Open Data. Most Semantic Knowledge Graphs (KG) are defined as classes (types) hierarchy where entities may belong to several classes. For instance, in figure 2, we can see that entity *Bob\_Marley* belongs to classes *Artist*, *Person*, and *Agent*. Hence, the type prediction problem can be reformulated as a hierarchical multi-label classification one [18]. In this context, one solution for entity type prediction is the local classifier per node approach, which consists of training each node of the class hierarchy separately and determining the overall type. However, it has been observed that type information and, more generally, KG suffer from several issues related to data quality such as accuracy, like typing error [21], completeness, like schema completeness and property completeness, and different kinds of errors. In [22], a systematic literature review on Linked Data completeness and an approach for Linked Data completeness assessment have been proposed. Among them, we focus on the following: several entities are in classes they do not belong to. For instance, let us consider the following query (figure 1) to get all the entities of the class *Person*. It also returns group entities, whereas these should be in the class *Group* (see table 1). We call such entities "outliers." Discovering such outliers is very important for correct class learning and understanding. Outlier detection methods on semantic web data have been proposed to find these outliers [24, 26]. These methods use outlier factors based on frequent pattern discovery methods that have been proposed before [6, 13]. The idea behind this is that transactions that contain more frequent patterns will have a big value of FPOF measure [6] and are unlikely to be outliers. In contrast, transactions with small FPOF values are likely to be outliers. Let us consider the following three entities of class *Artist*: [https://dbpedia.org/resource/Masaba\\_Gupta](https://dbpedia.org/resource/Masaba_Gupta), [https://dbpedia.org/resource/Bertram\\_Goodman](https://dbpedia.org/resource/Bertram_Goodman), and <https://dbpedia.org/resource/Mystik>. Masaba Gupta and Bertram Goodman are real artists, while Mystik is an outlier because it is a song. However, as we will see in example 4, FPOF measure gives a value higher for Mystik than for Masaba Gupta and Bertram Goodman. FPOF does not consider semantics aspects of the entities, but these are very important in KG and specifically in DBpedia. The KG data is based on a well-designed schema which is the ontology. The application of the FPOF measure on semantic web data while completely ignoring the semantics behind the ontology has further motivated our work. In a knowledge graph, ontology plays a very important role. It is a reference schema for organizing data while defining the concepts and relationships between these concepts. Among its components, we have the classes, the attributes or properties, the relations, etc. The role of the ontology here is to provide the semantics behind the properties describing the entities. We propose a semantic measure based on properties (domain and range properties).

```

1 SELECT distinct ?entity
2 WHERE {
3     ?entity a dbo:Person.
4 }

```

Figure 1: Query to get all the entities of the class *Person*.

| Entities  | Person | Group |
|---|--------|-------|
| <a href="http://dbpedia.org/resource/Les_Twins">http://dbpedia.org/resource/Les_Twins</a>             | No     | Yes   |
| <a href="https://dbpedia.org/resource/Nelson_Mandela">https://dbpedia.org/resource/Nelson_Mandela</a> | Yes    | No    |
| <a href="https://dbpedia.org/resource/Barack_Obama">https://dbpedia.org/resource/Barack_Obama</a>     | Yes    | No    |
| ...   | ...    | ...   |

Table 1: Some entities of the result of the query of Figure 1.

The main contributions of our paper are as follows :

- We define a semantic measure that favors the real entities of the class (inliers) with positive values while penalizing outliers with negative values and improving it with the discovery of frequent and rare itemsets.
- We propose a generic algorithm for outlier detection based on the semantic measure in a knowledge graph. This algorithm can be used for frequent or rare patterns.
- We present a set of experiments on DBpedia classes showing that our method outperforms FPOF like ones.

The rest of the paper is organized as follows. Section II discusses related work. In section III, we propose the basic definitions and problem reformulation. Section IV proposes our measure and presents a generic algorithm for outlier detection in a knowledge graph while presenting a theoretical analysis. Section V presents the results of our experiments, and section VI concludes the paper.

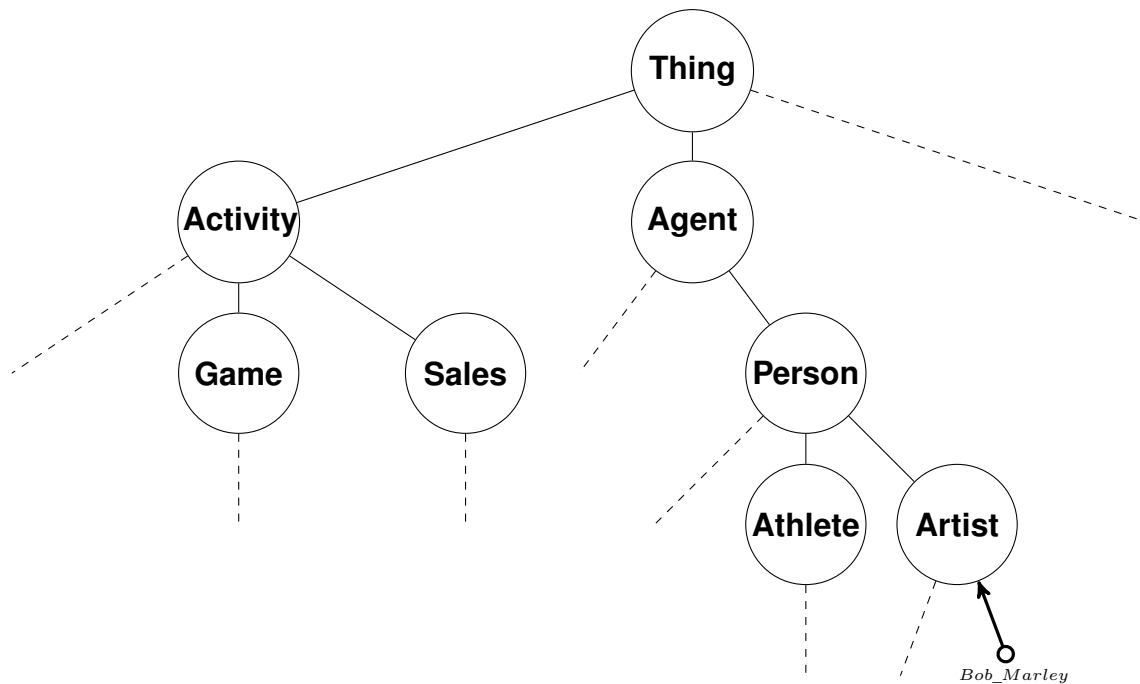


Figure 2: A sub-part of the DBpedia hierarchy

## II RELATED WORK

In this section, we will present the different dimensions of linked data quality and the outlier detection methods.

**Linked Data quality.** [25] cites three approaches that are used to construct a Knowledge Graph namely manual approaches, like Freebase [10] and Cyc [8], cooperative approaches, like DBpedia [17], and Yago [11] and automatic approaches, like NELL [14]. The manual approaches denote that the construction of the Knowledge Graph is done completely by a human, the automatic ones indicate that supervised or unsupervised learning techniques have been used and the cooperative ones indicate that most of the tasks are done by a human [25]. For the data quality of the Knowledge Graph, manual approaches tend to be more accurate compared to other approaches [21]. With the extensive growth of data and the many domains they covered it is impossible to build or maintain a knowledge graph by manual methods, automatic methods seem more suitable although they are more likely to contain several errors. For example on DBpedia several types of errors such as typing errors, incorrect numerical data, etc. have been noted. In [1] the dimensions of data quality were classified into four categories:

- Intrinsic Category which allows to assess the validity and consistent of the data. In this category, we have *Accuracy*, *Consistency* and *Trustworthiness*.
- Contextual category : “highlights the requirement that data quality must be considered within the context of the task at hand” [1]. In this category we have Relevancy, Completeness, and Timeliness
- Representational Category : "A concept that data quality is related to the “format of the data (concise and consistent representation) and meaning of data (interpretability and ease of understanding)” [1].
- Accessibility Category : "Accessibility dimensions are about how easily accessible and secure data is, such as availability and security." [22]

One of the solutions to the problem of data quality in KG is the language SHACL (Shapes Constraint Language)<sup>1</sup>. SHACL is a W3C Recommendation that defines a language for validating RDF graphs against a set of conditions. For its use, a good understanding of the ontology is necessary to define these conditions. Another problem is the scalability since each KG has its ontology.

It is almost impossible to define a generic approach to assessing data quality across all dimensions, and one dataset may be suitable for one purpose but not for another [22]. Depending on the application, it is imperative to address particular data quality issues. For class learning and understanding, the problem we face is misclassified entities (Accuracy dimension). It is very important to have suitable parameters to find misclassified entities in a given class.

**Outlier detection.** One solution for solving this problem of misclassified entities is outlier detection. Outlier detection has been the topic of several surveys and reviews. In these surveys, different classifications methodologies are presented. More cited categories are the following: Nearest Neighbour Based Outlier Detection techniques, Distance-Based Techniques, Density-Based Techniques, Cluster-Based Techniques, Statistical Approach Based Techniques. Distance-Based [2] and Nearest Neighbour Based Techniques [5] rely on the notion of distance. A distance-based outlier in a dataset  $\mathcal{D}$  is a data object with a given percentage of the objects in  $\mathcal{D}$  having a distance of more than  $d_{min}$  away from it [13]. Nearest Neighbour based outlier detection techniques require a distance or similarity measure between two data points. If a point  $x$  has a short distance to its  $k$  neighbors, it is considered as normal otherwise it is considered as outlier. Density-based techniques measure density of a point  $x$  within a small region by counting number of points within a neighborhood region. Breunig et al. [3] introduced the concept of Local Outlier Factor (LOF), a score which is assigned to every point based on its

<sup>1</sup><https://www.w3.org/TR/shacl/>

local density. All data points are sorted in decreasing order of LOF value. Points with high scores are detected as outliers. In Cluster-Based Outlier detection techniques, a cluster represents a collection of data objects similar to one another within the same cluster and dissimilar to the objects in other clusters. Inliers correspond to data in a cluster while outliers do not belong to any cluster. Statistical approach based technique assumes a distribution or probability model for the given data and then identifies outliers with respect to the model using a discordancy test [7]. Many of these techniques suffer from high dimensional space curse and high computational cost [16]. More recently, a new trend has appeared, using frequent pattern technique. Several measures were consequently proposed. FPOF (Frequent Pattern Outlier Factor) has been presented by [6]. The idea behind this measure is that transactions that contain more frequent patterns will have a big value of FPOF measure and are unlikely to be outliers. In contrast, transactions with small FPOF values are likely to be outliers. [13] proposed another measure WCFPOF (Weighted Closed Frequent Pattern Outlier Factor) in order to overcome the drawbacks of FPOF. According to this measure, transactions that contain more closed frequent patterns are more likely to be inliers and those that contain less closed frequent patterns are likely to be outliers. Another approach for overcoming the drawbacks of FPOF measurement was formulated by [15]. According to this one, transactions that contain longer frequent pattern (i.e. longer superset) are more likely to be inliers because they contain more subset frequent patterns, while transactions that contain short frequent patterns are likely to be outliers. By approximating the Frequent Pattern Outlier Factor (FPOF) with sampled patterns, [26] proposes a method for detecting misclassified entities. The observation is the same for these measures: they do not consider semantics.

### III PRELIMINARIES AND PROBLEM REFORMULATION

**Knowledge graph.** A *knowledge graph* is an RDF dataset described by a set of triples (*subject, predicate, object*). It can be represented as a tuple  $\mathcal{K} = (\mathcal{T}, \mathcal{A})$  where  $\mathcal{T}$  is called a Tbox formed by names and assertions about concepts (or classes) and roles (or predicates), and  $\mathcal{A}$  is called an ABox formed by assertions about individuals called entities and facts. This paper focuses on DBpedia[12] TBox and its entities. For example, (*dbo:Artist, rdfs:subClassOf, dbo:Person*) is an assertion in DBpedia TBox which means that the concept *Artist* is a subclass of the concept *Person* (or *Person* is a superclass of *Artist*). In other words, all artists are also persons. For example, (*Bob\_Marley, rdf:type, dbo:Artist*) means that *Bob\_Marley* is an artist, then an entity of the class *Artist*. The triple (*Bob\_Marley, spouse, Rita\_Marley*) is an assertion in DBpedia ABox which means that *Bob\_Marley* has spouse *Rita\_Marley*. Given a class  $C \in \mathcal{T}$  and an entity  $e \in \mathcal{A}$ , (*e, rdf:type, C*) implies that *e* is an instance of *C*. To better formalize and define all these notions and the others we need to formulate our problem, we use the Description Logics (DL) [9] formal notations. In that case, the assertion (*dbo:Artist, rdfs:subClassOf, dbo:Person*) is denoted by  $dbo:Artist \sqsubseteq dbo:Person$  which means that *dbo:Artist* is subsumed by *dbo:Person*.  $dbo:Artist(Bop\_Marley)$  denotes that *Bop\_Marley* is an instance (or entity) of the class *Artist*. The relation  $spouse(Bob\_Marley, Rita\_Marley)$  materializes that the predicate *spouse* links the subject *Bob\_Marley* to the object *Rita\_Marley*. We distinguish two types of predicates from these examples: outgoing and incoming. If we have the triple (*X, P, Y*), then we say that *P* is an outgoing predicate for *X* and an incoming predicate for *Y* whatever the type of *X* and *Y*. Then, there is a case where an outgoing predicate *P* is specific to subjects that instantiate a class *C*. In that case, we say that *C* is the domain of *P*, denoted by  $\exists P.\top \sqsubseteq C$ . Another case is when an incoming predicate is specific to objects of a class *C*. Here, we say that *C* is

the range of the predicate  $P$ , denoted by  $\top \sqsubseteq \forall P.C$ . The disjointness of two classes  $C_1$  and  $C_2$  is denoted by  $C_1 \sqsubseteq \neg C_2$ . For example, according to the ontology of DBpedia, it should not have entities belonging to both the class *Person* and the class *Organization*, then we have  $Person \sqsubseteq \neg Organization$ . There are several types of knowledge graphs, but for this work, we are interested in these where ABox and TBox are available and accessible by SPARQL queries like DBpedia.

In the following, we show the transformation from RDF data to the transactional database as done in [26].

**Transactional database.** Now, we will show how to represent a transaction from the predicates that describe an entity. Since we are interested in predicates having domains, we consider the transactional database defined on the set of items  $\mathcal{I} = \{P : (\exists C \in \mathcal{T})(\exists P.\top \sqsubseteq C)\}$ . An itemset (or pattern), denoted by  $\varphi$ , is a non empty subset of  $\mathcal{I}$ . Formally, we have  $\varphi \subseteq \mathcal{I}$ . The set of all patterns that can be generated from  $\mathcal{I}$  is called the pattern language  $\mathcal{L} = 2^{\mathcal{I}} \setminus \emptyset$ . In this paper, a transaction is a couple  $(e, \mathcal{I}^e) \in \mathcal{A} \times \mathcal{L}$  where  $\mathcal{I}^e$  is the set of all predicates describing the entity  $e$  that appears in  $\mathcal{I}$ . Formally, we have  $\mathcal{I}^e = \{P \in \mathcal{I} : (\exists C \in \mathcal{T})(C(e))(\exists P.\top \sqsubseteq C)\}$ . In the following, we denote such transaction as  $\mathcal{I}^e$ . So in our context, a transactional database  $\mathcal{D}_C$  is a multi-set of transactions defined in  $\mathcal{I}$  where all items (predicates) of  $\mathcal{I}$  describes an entity of class  $C$ . It means that we are making a restriction in class  $C$  in which we look for whether it contains outliers or not. For example, to find outliers in the class *Artist*, we only consider entities that instantiate *Artist*.

*Example 1:*

For instance,  $\mathcal{D}_{Artist}$  in Table 2 is a toy dataset of 20 transactions from the class Artist of DBpedia<sup>2</sup>(entities are prefixed by `dbr(http://dbpedia.org/resource/)` and predicates by `dbo(http://dbpedia.org/ontology/)`). For instance, to obtain the transaction  $\mathcal{I}^{Bob\_Marley}$ , we run the following SPARQL query on the Dbpedia endpoint (`https://dbpedia.org/sparql/`):

```

1      SELECT DISTINCT ?P WHERE {
2          <http://dbpedia.org/resource/Bob_Marley> ?P ?object.
3          ?P rdfs:domain ?C.
4      }

```

$\mathcal{D}_{Artist}$  is built from the set of 38 items  $\mathcal{I} = \{deathPlace, deathDate, birthPlace, birthDate, deathCause, partner, relative, child, parent, spouse, birthYear, deathYear, birthName, nationality, field, training, residence, restingPlacePosition, restingPlace, cinematography, director, producer, starring, writer, runtime, Work/runtime, artist, endingTheme, network, openingTheme, previousWork, starring, subsequentWork, completionDate, numberOfEpisodes, bandMember, formerBandMember, hometown\}$ . Table 3 gives the domain and the range (if it exists) of the set of items (predicates) in  $\mathcal{I}$ .

The transaction  $\mathcal{I}^{Bob\_Marley} = \{deathPlace, deathDate, birthPlace, birthDate, deathCause, partner, relative, child, parent, spouse, birthYear, deathYear\}$  contains 12 items. so, there are  $2^{12} - 1$  patterns that appear in  $\mathcal{L}(\mathcal{D}_{Artist})$ . For example,  $\varphi = \{spouse, birthYear\}$  is one among them. This pattern also belongs to the transaction  $\mathcal{I}^{Hank\_Williams}$ .

<sup>2</sup>Access in 18/10/2021

Table 2:  $\mathcal{D}_{Artist}$  : an example of a transactional database from the class Artist of DBpedia

| entity             | Itemset (set of predicates)   |
|--------------------|---|
| Bob_Marley         | {deathPlace, deathDate, birthPlace, birthDate, deathCause, partner, relative, child, parent, spouse, birthYear, deathYear}                        |
| Omar_Kiam          | {birthDate, birthName, birthPlace, deathDate, deathPlace, nationality}  |
| Bertram_Goodman    | {field, training}   |
| Giuliana_Camerino  | {birthDate, birthName, birthPlace, deathDate, deathPlace, nationality}  |
| Twice_as_Nice      | {cinematography, director, producer, starring, writer, runtime, Work/runtime}   |
| Robin_Harris       | {birthDate, birthName, birthPlace, deathDate, deathPlace, nationality}  |
| Brett_Newski       | {birthDate, hometown}   |
| LaWanda_Page       | {birthDate, birthName, birthPlace, deathDate, deathPlace, nationality}  |
| Frank_Suero        | {birthDate, birthName, birthPlace, deathDate, deathPlace, nationality}  |
| Josephus_Thimister | {birthDate, birthName, birthPlace, deathDate, deathPlace, nationality}  |
| Sid_James          | {birthDate, birthName, birthPlace, deathDate, deathPlace, nationality}  |
| Masaba_Gupta       | {residence, spouse}   |
| Hank_Williams      | {deathPlace, deathDate, birthPlace, birthDate, restingPlacePosition, deathCause, relative, restingPlace, spouse, birthName, birthYear, deathYear} |
| Children_of_Eve    | {cinematography, director, producer, writer, runtime, Work/runtime}   |
| Greg_Giraldo       | {birthDate, birthName, birthPlace, deathDate, deathPlace, nationality}  |
| Mystik             | {artist, previousWork, producer, writer, runtime, Work/runtime}   |
| Jerry_Clower       | {birthDate, birthName, birthPlace, deathDate, deathPlace, nationality}  |
| Mackenzie_Taylor   | {birthDate, birthName, birthPlace, deathDate, deathPlace, nationality}  |
| The_Lead           | {endingTheme, network, openingTheme, previousWork, starring, subsequentWork, completionDate, numberOfEpisodes, runtime, Work/runtime}             |
| 7icons             | {bandMember, formerBandMember, hometown}  |

Table 3: Domain and range (optional) of the predicates of our dataset  $\mathcal{D}_{Artist}$  (EduIns: EducationalInstitution, TShow: TelevisionShow, Bcaster: Broadcaster, nonNegInt: nonNegativeInteger)

| $P$        | dom( $P$ ) | range( $P$ ) | $P$                  | dom( $P$ ) | range( $P$ )           |
|------------|------------|--------------|----------------------|------------|------------------------|
| deathPlace | Person     | Place        | birthYear            | Person     | gYear                  |
| deathDate  | Person     | date         | deathYear            | Person     | gYear                  |
| birthPlace | Person     | Place        | birthName            | Person     | langString             |
| birthDate  | Person     | date         | nationality          | Person     | Country                |
| deathCause | Person     |              | field                | Artist     |                        |
| partner    | Person     | Person       | training             | Artist     | EducationalInstitution |
| relative   | Person     | Person       | residence            | Person     | Place                  |
| child      | Person     | Person       | restingPlacePosition | Person     | SpatialThing           |
| parent     | Person     | Person       | restingPlace         | Person     | Place                  |
| spouse     | Person     | Person       | cinematography       | Film       | Person                 |

| $P$          | dom( $P$ )  | range( $P$ ) | $P$              | dom( $P$ )     | range( $P$ ) |
|--------------|-------------|--------------|------------------|----------------|--------------|
| director     | Film        | Person       | openingTheme     | TelevisionShow | Work         |
| producer     | Work        | Agent        | previousWork     | Work           | Work         |
| starring     | Work        | Actor        | subsequentWork   | Work           | Work         |
| writer       | Work        | Person       | completionDate   | Work           | date         |
| runtime      | Work        | double       | numberOfEpisodes | TShow          | nonNegInt    |
| Work/runtime | Work        | minute       | bandMember       | Band           | Person       |
| artist       | MusicalWork | Agent        | formerBandMember | Band           | Person       |
| endingTheme  | TShow       | Work         | hometown         | Agent          | Settlement   |
| network      | Bcaster     | Bcaster      |                  |                |              |



We already remarked that, semantically, some of these entities such as *Twice\_as\_Nice*, *Children\_of\_Eve*, *Mystik*, and *The\_Lead*, are not artists, and yet they are instantiated in the Artist class. We call them outliers while the others are well instantiated in Artist; they are inliers.

*Definition 1: Outlier, Inlier*

Given two disjoint classes  $C$  and  $C'$  ( $C \sqcap C' = \emptyset$ ) and an entity  $e$ . If  $e$  is instantiated in  $C$  but really defined as an entity of  $C'$ , then  $e$  is an outlier. However, if really  $e$  is an instance of  $C$ , then it is called an inlier.

*Definition 2: Frequency of a pattern*

Given a transactional database  $\mathcal{D}_C = \{(e_1, \mathcal{I}^{e_1}), \dots, (e_n, \mathcal{I}^{e_n})\}$  of a class  $C$  and a pattern  $\varphi \in \mathcal{L}(\mathcal{D}_C)$ . The frequency of  $\varphi$  is the number of transactions of  $\mathcal{D}$  containing  $\varphi$ . Formally,

$$freq(\varphi, \mathcal{D}_C) = |\{(e_i, \mathcal{I}^{e_i}) \in \mathcal{D}_C : \varphi \subseteq \mathcal{I}^{e_i}\}|.$$

*Example 2:*

The frequency of  $\varphi = \{spouse, birthYear\}$  in  $\mathcal{D}_{Artist}$  is 2 because only  $\mathcal{I}^{Bob\_Marley}$  and  $\mathcal{I}^{Hank\_Williams}$  contain  $\varphi$ . The frequency of  $\varphi_1 = \{spouse\}$  is  $freq(\varphi_1, \mathcal{D}_{Artist}) = |\{Masaba\_Gupta, Bob\_Marley, Hank\_Williams\}| = 3$ .

According to the frequency, one can judge a pattern as frequent or rare (not frequent) in a dataset.

*Definition 3: Frequent and Rare pattern*

Let us consider a dataset  $\mathcal{D}_C$  of class  $C$ , a minimum support threshold  $\alpha \in [0, 1]$ , and  $\varphi$  a pattern of  $\mathcal{L}(\mathcal{D}_C)$ . We say that :

- $\varphi$  is frequent in  $\mathcal{D}_C$  if and only if  $freq(\varphi, \mathcal{D}_C) \geq \alpha \times |\mathcal{D}_C|$ .
- $\varphi$  is rare in  $\mathcal{D}_C$  if and only if  $0 < freq(\varphi, \mathcal{D}_C) < \alpha \times |\mathcal{D}_C|$ .

*Example 3:*

If we consider a minimum threshold  $\alpha = 4/20$ , then  $\varphi = \{spouse, birthYear\}$  is a rare pattern in  $\mathcal{D}_{Artist}$  because  $freq(\varphi, \mathcal{D}_{Artist}) = \frac{2}{20} < \frac{4}{20}$  but the pattern  $\varphi' = \{birthDate, birthName\}$  is frequent in  $\mathcal{D}_{Artist}$  because  $freq(\varphi', \mathcal{D}_{Artist}) = \frac{11}{20} > \frac{4}{20}$ .

In this paper, we focus on two interestingness measures *frequency* and *rare*. So, given an interestingness measure  $m$  it is possible to extract a set of interesting patterns. We denote by  $q(\cdot)$ , such that  $q(\cdot) \in \{true, false\}$ , the constraint that a pattern of  $\mathcal{D}_C$  should respect according to the interestingness measure  $m$ . Therefore, on the one hand, the set of all frequent patterns given a minimum threshold  $\alpha$  can be formulated by  $\mathcal{Pset}(\mathcal{D}_C, freq_\alpha) = \{\varphi \in \mathcal{L}(\mathcal{D}_C) : q(freq(\varphi, \mathcal{D}_C) \geq \alpha) = true\}$ . On the other hand, the set of all rare patterns given a minimum threshold  $\alpha$  can be formulated by  $\mathcal{Pset}(\mathcal{D}_C, rare_\alpha) = \{\varphi \in \mathcal{L}(\mathcal{D}_C) : q(0 < freq(\varphi, \mathcal{D}_C) < \alpha \times |\mathcal{D}_C|) = true\}$ .

Now we present FPOF in the knowledge graph as done in [26] to detect outliers.

**FPOF : Frequent Pattern Based Outlier Factor.** As it was introduced in Section II, FPOF is based on frequent patterns. The main idea behind this approach for transactional databases is that inliers will contain most of the frequent patterns while most of the patterns that outliers contain will have a low frequency. So, they formulate the following metric.

*Definition 4:*

Given a transactional database  $\mathcal{D}_C = \{(e_1, \mathcal{I}^{e_1}), \dots, (e_n, \mathcal{I}^{e_n})\}$  and a minimum frequency threshold  $\alpha > 0$ , the frequent pattern outlier factor of the transaction  $(e, \mathcal{I}^e)$  denoted by  $fprof(e, \mathcal{D}_C)$  is defined as follows:

$$fprof(e, \mathcal{D}_C) = \frac{1}{|\mathcal{Pset}(\mathcal{D}_C, freq_\alpha)|} \times \sum_{\varphi \in \mathcal{Pset}(\mathcal{D}_C, freq_\alpha) \wedge \varphi \subseteq \mathcal{I}^e} \frac{freq(\varphi, \mathcal{D}_C)}{|\mathcal{D}_C|}$$

*Example 4:*

Let us consider the database  $\mathcal{D}_{Artist}$  shown in Table 2. Each entity is described by a set of properties. For example, the entity *Sid\_James* is described by the properties birthDate, birthName, deathDate, birthPlace, and nationality. Considering the minsup 0.2 (4/20), the fprof values are calculated in Table 4 (column 3).

Table 4:  $\mathcal{D}_{Artist}$  : an example of transactional database of the class Artist of DBpedia with fprof values

| Artist             | Properties  | FPOF   |
|--------------------|---|--------|
| Sid_James          | {birthDate, birthName, birthPlace, deathDate, deathPlace, nationality}  | 0,5129 |
| Mackenzie_Taylor   | {birthDate, birthName, birthPlace, deathDate, deathPlace, nationality}  | 0,5129 |
| Giuliana_Camerino  | {birthDate, birthName, birthPlace, deathDate, deathPlace, nationality}  | 0,5129 |
| Jerry_Clower       | {birthDate, birthName, birthPlace, deathDate, deathPlace, nationality}  | 0,5129 |
| Robin_Harris       | {birthDate, birthName, birthPlace, deathDate, deathPlace, nationality}  | 0,5129 |
| Greg_Giraldo       | {birthDate, birthName, birthPlace, deathDate, deathPlace, nationality}  | 0,5129 |
| LaWanda_Page       | {birthDate, birthName, birthPlace, deathDate, deathPlace, nationality}  | 0,5129 |
| Frank_Suero        | birthDate, birthName, birthPlace, deathDate, deathPlace, nationality}   | 0,5129 |
| Josephus_Thimister | {birthDate, birthName, birthPlace, deathDate, deathPlace, nationality}  | 0,5129 |
| Omar_Kiam          | {birthDate, birthName, birthPlace, deathDate, deathPlace, nationality}  | 0,5129 |
| Hank_Williams      | {deathPlace, deathDate, birthPlace, birthDate, restingPlacePosition, deathCause, relative, restingPlace, spouse, birthName, birthYear, deathYear} | 0,2704 |
| Bob_Marley         | {deathPlace, deathDate, birthPlace, birthDate, deathCause, partner, relative, child, parent, spouse, birthYear, deathYear}                        | 0,1371 |
| Brett_Newski       | {birthDate, hometown}   | 0,0098 |
| Children_of_Eve    | {cinematography, director, producer, writer, runtime, Work/runtime}   | 0,0090 |
| Mystik             | {artist, previousWork, producer, writer, runtime, Work/runtime}   | 0,0091 |
| Twice_as_Nice      | {cinematography, director, producer, starring, writer, runtime, Work/runtime}   | 0,0091 |
| The_Lead           | {endingTheme, network, openingTheme, previousWork, starring, subsequentWork, completionDate, numberOfEpisodes, runtime, Work/runtime}             | 0,0091 |
| Masaba_Gupta       | {residence, spouse}   | 0      |
| Bertram_Goodman    | {field, training}   | 0      |
| 7icons             | {bandMember, formerBandMember, hometown}  | 0      |

Based on the frequency of the itemsets, the FPOF measure gives a higher score to the entity *Mystik* (which is an outlier) than to entities *Masaba\_Gupta* and *Bertram\_Goodman* (which are inliers). The FPOF measure is based only on the frequency of items that describe transactions while ignoring the semantics.

In knowledge graphs, the frequency of items only does not allow to find the outliers of a set because semantics are decisive for entity definition. For example, *birthName* is a predicate that defines a person, which is less frequent than the predicate *abstract*, which is generally used by the entities of all classes.

**The problem that we want to solve in this paper can be formulated as follows:**

**Given a class  $C$  of knowledge graph  $\mathcal{K} = (\mathcal{T}, \mathcal{A})$ , let  $\mathcal{D}_C$  be its transactional database, and  $\mathcal{Pset}(\mathcal{D}_C, m_\alpha)$  a set of interesting patterns built with an interestingness measure  $m$ ,**

**Q1: which metric can be used to favor the inliers and penalize the outliers?**

**Q2: how to benefit from the semantic of the TBox  $\mathcal{T}$  to improve the metrics?**

Table 5: Notations

| Symbol                                    | Definition   |
|---|--|
| $e$                                       | An entity  |
| $\mathcal{I}$                             | The set of items of the database   |
| $\mathcal{I}^e$                           | All the items of the entity $e$  |
| $\varphi$                                 | A pattern of items (a non empty subset of $\mathcal{I}$ )  |
| $\mathcal{A}$                             | The ABox   |
| $\mathcal{T}$                             | The TBox   |
| $\mathcal{K}$                             | The Knowledge Graph  |
| $C$                                       | A class  |
| $\mathcal{D}_C$                           | A transactional database from the class $C$  |
| $m$                                       | An interestingness measure   |
| $\alpha$                                  | A minimum frequency threshold  |
| $\mathcal{Pset}(\mathcal{D}_C, m_\alpha)$ | The set of all frequent (rare) patterns given a minimum threshold $\alpha$                                   |
| $k$                                       | The number of super-classes  |
| $\mathcal{C}$                             | The set of concepts  |
| $\mathcal{C}^{\leq k}(C)$                 | The set of concepts containing $C$ , its super-classes at a level at most equal to $k$ and its sub-classes.  |
| $\mathcal{SCl}$                           | The set of super-classes   |
| $\mathcal{SCl}^{\leq k}(C)$               | The set of super-classes of $C$ which are located at a level less than or equal to $k$ with respect to $C$ . |
| $setOut$                                  | The set of outliers  |
| $setIn$                                   | The set of inliers   |

#### IV ONTOLOGY AND PATTERN-BASED OUTLIER DETECTION IN KNOWLEDGE GRAPH

This section shows how to take advantage of semantic relationships that arise in ontology to improve the mined pattern according to a measure of interest. However, let us introduce the idea behind the rare patterns for detecting outliers.

In a similar way to the FPOF metric, the rare patterns can be used to detect outliers from the transactional database. A naive method is to sum, for any transaction, the frequency of the

rare patterns it contains. In that case, the outliers will have high scores while the inliers will have low scores. But, the main problem with this basis metric follows from the fact that many inliers will have null scores. This is why in this paper, we will introduce semantics on metrics, particularly on items forming patterns that are nothing other than properties that appear in the ontology. It is important to note that all the properties that we use to make up the patterns have domains, and some of them also have ranges. In the rest of this paper, the metrics as well as the algorithm that we are going to introduce are valid for both rare and frequent patterns.

**Intuition behind this method.** Our approach is based on the fact that a property  $P$  with domain  $C$  must necessarily be an outgoing property of an entity of the class  $C$  when it is used. It is very important to note that  $P$  can also be an outgoing property of any entity of a subclass of  $C$  ( $C' \sqsubseteq C$ ) or of a super-class of  $C$  located at a certain level  $k \geq 0$ . We say that a super-class  $C'$  of the class  $C$  is at level  $k$  with refer to  $C$  if there are  $k - 1$  classes  $\{C_1, \dots, C_{k-1}\}$ , with  $C_i \neq C_j$  for all  $i \neq j$ , in  $\mathcal{T}$  such that  $C \sqsubseteq C_1 \sqsubseteq \dots \sqsubseteq C_{k-1} \sqsubseteq C'$ . Indeed, the fact of favoring all the super-classes of  $C$  will distort the calculations because there will be outliers who will benefit from it. For example, *Artist* and *Athlete* are two classes in the same level that have super-classes Person (of level  $k = 1$ ), Agent (of level  $k = 2$ ) and Thing (of level  $k = 3$ ). When searching for outliers in *Artist*, it is not interesting to consider properties that have Thing as its domain, because the latter is the parent class of all classes in DBpedia. However, properties that have Agent as their domain are only interesting if the outliers are in a twin class  $C'$  (*Activity* for instance) or one of sub-classes of  $C'$  (*Game*, *Sales*, ...). The same is true for properties that have Person as their domain, if some outliers really belong in class *Organization*, then they are useful for detecting outliers. If, on the other hand, the real class of the outliers is *Athlete*, then the properties having Person as their domain are not interesting to detect the outliers. So, the properties which can allow us to verify if an entity is an outlier or not, are not obvious. It is why retrieving these properties is a real problem of efficiency for our method because we do not know apriori at what level  $k \geq 0$  we must stop to detect outliers since we ignore their real classes. To solve this problem, we propose to vary the value of  $k$  to find the most relevant properties to detect outliers. In practice, the value of  $k$  is small. Therefore, we denote by  $\mathcal{SC}^{\leq k}(C) = \{C' \in \mathcal{T} : (\exists i, 1 \leq i \leq k)(C_i \in \mathcal{T})(C \sqsubseteq C' \sqsubseteq C_i)\}$  the set of super-classes of  $C$  which are located at a level less than or equal to  $k$  with respect to  $C$ . So,  $\mathcal{SC}^{\leq \infty}(C) = \{C' \in \mathcal{T} : C \sqsubseteq C'\}$ .

Let us denote by  $\mathcal{C}^{\leq k}(C)$  the set of concepts containing  $C$ , its super-classes at a level at most equal to  $k$  and its sub-classes:  $\mathcal{C}^{\leq k}(C) = \{C\} \cup \{C' \in \mathcal{T} : C' \sqsubseteq C\} \cup \{C' : C' \in \mathcal{SC}^{\leq k}(C)\}$ . The semantic judgment we have on a property  $P$  that appears in a pattern  $\varphi$  can be formulated as follows:

- (a) If a property  $P$  is interesting enough to well describe an inlier  $e$  of a class  $C$  in a certain level  $k$  then it must necessarily have a domain which is part of  $\mathcal{C}^{\leq k}(C)$ . The set of predicates in the pattern  $\varphi$  that meet this intuition is defined by  $\{P \in \varphi : (\exists C' \in \mathcal{C}^{\leq k}(C))(\exists P.T \sqsubseteq C')\}$ .
- (b) An entity found in class  $C$  which has an outgoing property  $P$  whose domain does not belong to  $\mathcal{C}^{\leq \infty}(C)$  is likely to be poorly described, and looks like an outlier if most of the properties it contains have this tendency. The set of predicates in the pattern  $\varphi$  that meet this intuition is defined by  $\{P \in \varphi : (\nexists C' \in \mathcal{C}^{\leq \infty}(C))(\exists P.T \sqsubseteq C')\}$ .
- (c) If a property  $P$  that matches the intuition in (a) additionally has a range that is in  $\mathcal{C}^{\leq k}(C)$ , then it is determinant for class  $C$ . The set of predicates in the pattern  $\varphi$  that meet this intuition is defined by  $\{P \in \varphi : (\exists C' \in \mathcal{C}^{\leq k}(C))(\exists P.T \sqsubseteq C')(\exists C'' \in \mathcal{C}^{\leq k}(C))(T \sqsubseteq P.C'')\}$ .

- (d) If a property  $P$  that matches the intuition in (b) additionally has a range that is part of  $\mathcal{C}^{\leq k}(C)$ , then the use of  $P$  is ambiguous. The set of predicates in the pattern  $\varphi$  that meet this intuition is defined by  $\{P \in \varphi : (\nexists C' \in \mathcal{C}^{\leq \infty}(C))(\exists P. \top \sqsubseteq C')(\exists C'' \in \mathcal{C}^{\leq k}(C))(\top \sqsubseteq \forall P.C'')\}$ .

Based on these intuitions, we can now define our algorithm which combines the mined patterns and the semantic relationships from ontology to detect outliers.

**A generic algorithm for outlier detection in knowledge graph.** We propose an algorithm based on the four intuitions presented in the previous section. To be in agreement with these intuitions, we say that a predicate that verifies intuition (a) brings a bonus of 1 to the score of the entity that it describes while that which verifies (b) brings a penalty of 1 to the entity's score. On the other hand, a predicate that verifies intuition (c) brings a bonus of 2 (bonus of 1 for the domain and of 1 for the range) while a predicate that verifies intuition (d) brings a penalty of 2 for the entity it describes (penalty of 1 for the domain and 1 for the range). Thereby, given a pattern we introduce the notion of reliability, a metric to compute the total score a pattern brings to an entity for a given class.

*Definition 5: Reliability*

Let  $e$  be an entity of a class  $C$  described by a certain predicates that belong to  $\mathcal{I}$ , a pattern  $\varphi \subseteq \mathcal{I}$  and  $k$  a positive integer. The reliability of the pattern  $\varphi$  to the entity  $e$  with respect to the class  $C$  is defined as follows:

$$reliability(\varphi, e, C, k) = \begin{cases} 0 & \text{if } \varphi \not\subseteq \mathcal{I}^e \\ |\{P \in \varphi : (\exists C' \in \mathcal{C}^{\leq k}(C))(\exists P. \top \sqsubseteq C')\}| \\ +|\{P \in \varphi : (\exists C' \in \mathcal{C}^{\leq k}(C))(\exists P. \top \sqsubseteq C')(\top \sqsubseteq P.C')\}| & (1) \\ -|\{P \in \varphi : (\nexists C' \in \mathcal{C}^{\leq \infty}(C))(\exists P. \top \sqsubseteq C')\}| \\ -|\{P \in \varphi : (\nexists C' \in \mathcal{C}^{\leq \infty}(C))(\exists P. \top \sqsubseteq C') \\ (\exists C'' \in \mathcal{C}^{\leq k}(C))(\top \sqsubseteq \forall P.C'')\}| & \text{otherwise} \end{cases}$$

It is clear that the reliability of a pattern is the sum of the contribution of all the predicates it contains. Now, we can present the algorithm exploring this metric to capture the outliers that disturb the veracity of all the instances in a given class. Let us recall that the goal of this paper is not to find frequent or rare patterns, we just reuse the existing methods.

*Example 5:*

Let us consider this two patterns  $\varphi_1 = \{birthDate, birthName\}$  and  $\varphi_2 = \{bandMember, hometown\}$ . We are now going to compute the reliability of the transactions  $(Robin\_Harris, \mathcal{I}^{Robin\_Harris})$  and  $(7icons, \mathcal{I}^{7icons})$  with refer to the class *Artist* and with different values of  $k$ . We have :

$reliability(\varphi_1, Robin\_Harris, Artist, 0) = 0$  because none of the properties in  $\varphi_1$  has a domain that belongs in  $\mathcal{C}^{\leq 0}(Artist) = \{Artist\}$  but in  $\mathcal{C}^{\leq \infty}(Artist)$ . These predicates are neutral since they have a domain that is a superclass of *Artist*. However,  $reliability(\varphi_1, Robin\_Harris, Artist, 1) = (+1) + (+1) = +2$  because *birthDate* and *birthName* have a domain that belongs in  $\mathcal{C}^{\leq 1}(Artist) = \{Artist, Person\}$  then  $(+1)$ .  $reliability(\varphi_2, 7icons, Artist, 0) = -1$  because, in one hand, the domain of the property *bandMember* doesn't belong in  $\mathcal{C}^{\leq \infty}(Artist)$  while its range doesn't belong in  $\mathcal{C}^{\leq 0}(Artist)$ . In other hand, the domain of the property *hometown* belongs in the set  $\mathcal{C}^{\leq \infty}(Artist)$  but not in

$\mathcal{C}^{\leq 0}(\text{Artist})$ . Let us now compute  $\text{reliability}(\varphi_2, 7icons, \text{Artist}, 1)$ . We know that the range of `bandMember` belongs in  $\mathcal{C}^{\leq 1}(\text{Artist})$  while its domain doesn't belong in  $\mathcal{C}^{\leq \infty}(\text{Artist})$ . So, this property gives a score of  $((-1)+(-1))$  to the entity `7icons`. The property `hometown` gives the same score like in  $k = 0$ . So, we have  $\text{reliability}(\varphi_2, 7icons, \text{Artist}, 1) = (-1) + (-1) = -2$ .

---

**Algorithm 1** ONTOPOD (Ontology and Pattern-based Outlier Detection)

---

**Input:** A set of patterns  $\mathcal{P}_{set}(\mathcal{D}_C, m_\alpha)$  obtained from the transactional database  $\mathcal{D}_C$  of a class  $C$  of a knowledge graph  $\mathcal{K} = (\mathcal{T}, \mathcal{A})$  and an interestingness measure  $m_\alpha$  with  $\alpha \in [0, 1]$  and a positive integer  $k$

**Output:** A set of entities weighted by their scores  $setOut$  (for outliers) and  $setIn$  (for inliers) obtained from  $\mathcal{D}_C$

```

1:  $setOut \leftarrow \emptyset$  ▷ The set of outliers
2:  $setIn \leftarrow \emptyset$  ▷ The set of inliers
3: for  $(e, \mathcal{I}^e) \in \mathcal{D}_C$  do
4:    $score(e, C, k) \leftarrow \text{reliability}(\mathcal{I}^e, e, C, k) + \sum_{\varphi \in \mathcal{P}_{set}(\mathcal{D}_C, m_\alpha) \wedge \varphi \subset \mathcal{I}^e} \text{reliability}(\varphi, e, C, k)$ 
5:   if  $score(e, C, k) < 0$  then
6:      $setOut \leftarrow setOut \cup \{(e, score(e, C, k))\}$ 
7:   else
8:      $setIn \leftarrow setIn \cup \{(e, score(e, C, k))\}$ 
9: return  $setOut$ 

```

---

Algorithm 1, named ONTOPOD, computes and returns a set of entities likely to be outliers. It takes as input the transactional database obtained from the entities of a class  $C$  of a given knowledge graph  $\mathcal{K} = (\mathcal{T}, \mathcal{A})$ , an interestingness measure  $m_\alpha$ , with  $\alpha$  the minimum threshold, and a positive integer  $k$  for the maximum level of superclasses to be considered for  $C$ . In the end, it returns all the entities that have a negative score, and therefore likely to be outliers. Then, it computes the score of each entity  $e$  according to its transaction  $\mathcal{I}^e$  and the extracted patterns  $\mathcal{P}_{set}(\mathcal{D}_C, m_\alpha)$  if it is not empty (line 4). Finally, all the entities having a negative score are retrieved and stored in the  $setOut$  variable (line 6) which is finally returned at output (line 9). Interestingly, it is also possible to easily flip inliers (those with positive scores) in the  $setIn$  variable (line 8). We maintain the idea that the method fails to judge entities that have a score equal to zero (0). We specifically use three variants of ONTOPOD. The first is ONTOPOD – *No* which means that  $\mathcal{P}_{set}(\mathcal{D}_C, m_\alpha)$  is empty. The second is ONTOPOD – *Freq* that corresponds to the case where the interestingness measure is the frequency,  $\mathcal{P}_{set}(\mathcal{D}_C, freq_\alpha)$ . The last is ONTOPOD – *Rare* which means that the rare patterns are used  $\mathcal{P}_{set}(\mathcal{D}_C, rare_\alpha)$ .

*Example 6:*

Table 6 gives the scores of ONTOPOD applied in the dataset of Table 7 that corresponds to the Artist class by promoting the super-class at level 1 (Person) ( $m_\alpha = no$  means that there is no interestingness measure that has been provided). As we can see it, our method manages to find outliers that the FPOF method failed to judge (Masaba\_Gupta, Bertram\_Goodman, 7icons). It is also important to note that without an interestingness measure, certain entities risk to have a score equal to zero (Brett\_Newski), hence the interest in finding frequent or rare patterns for deciding. Moreover, with frequent patterns, ONTOPOD finds all the outliers while with rare patterns, it adds an inlier (Brett\_Newski) to the outliers.

**Theoretical analysis of Algorithm 1.** In this section, we give the theoretical analysis of our algorithm ONTOPOD. As we said earlier, we do not consider the complexity of pattern

Table 6:  $\mathcal{D}_{Artist}$  : an example of transactional database from the Artist class of DBpedia with fpop and ONTOPOD values

| Artist             | FPOF   | ONTOPOD with $\mathcal{Pset}(\mathcal{D}_{Artist}, m_\alpha)$ and $k = 1$ |                         |                         |
|--------------------|--------|---|-------------------------|-------------------------|
|                    |        | $m_\alpha = no$   | $m_\alpha = freq_{0.2}$ | $m_\alpha = rare_{0.2}$ |
| Sid_James          | 0,5129 | 6   | 198                     | 6                       |
| Mackenzie_Taylor   | 0,5129 | 6   | 198                     | 6                       |
| Giuliana_Camerino  | 0,5129 | 6   | 198                     | 6                       |
| Jerry_Clower       | 0,5129 | 6   | 198                     | 6                       |
| Robin_Harris       | 0,5129 | 6   | 198                     | 6                       |
| Greg_Giraldo       | 0,5129 | 6   | 198                     | 6                       |
| LaWanda_Page       | 0,5129 | 6   | 198                     | 6                       |
| Frank_Suero        | 0,5129 | 6   | 198                     | 6                       |
| Josephus_Thimister | 0,5129 | 6   | 198                     | 6                       |
| Omar_Kiam          | 0,5129 | 6   | 198                     | 6                       |
| Hank_Williams      | 0,2704 | 14  | 94                      | 28606                   |
| Bob_Marley         | 0,1371 | 17  | 49                      | 34801                   |
| Brett_Newski       | 0,0098 | 1   | 2                       | 2                       |
| Children_of_Eve    | 0,0090 | -9  | -13                     | -293                    |
| Mystik             | 0,0091 | -7  | -11                     | -227                    |
| Twice_as_Nice      | 0,0091 | -11   | -15                     | -711                    |
| The_Lead           | 0,0091 | -11   | -15                     | -5639                   |
| Masaba_Gupta       | 0      | 3   | 3                       | 9                       |
| Bertram_Goodman    | 0      | 2   | 2                       | 6                       |
| 7icons             | 0      | -4  | -4                      | -20                     |

mining, which depends on the used pattern mining algorithm. However, the complexity of our method also depends on the size of the set of mined patterns. Generally, to compute the score of a transaction  $(e, \mathcal{I}^e)$ , we first compute the reliability of the itemset  $\mathcal{I}^e$  with refer to the corresponding entity  $e$  in  $O(|\mathcal{I}|)$ . Then, we compute the reliability of each pattern  $\varphi$  in  $\mathcal{Pset}(\mathcal{D}_C, m_\alpha)$  to a transaction  $(e, \mathcal{I}^e)$  with a complexity in  $O(|\mathcal{I}| \times |\mathcal{Pset}(\mathcal{D}_C, m_\alpha)|)$ . Thereby, the complexity of ONTOPOD is in  $O(|\mathcal{I}| \times |\mathcal{D}_C|) + O(|\mathcal{I}| \times |\mathcal{D}_C| \times |\mathcal{Pset}(\mathcal{D}_C, m_\alpha)|)$ . So, the complexity of ONTOPOD is generally in  $O(|\mathcal{I}| \times |\mathcal{D}_C| \times (1 + |\mathcal{Pset}(\mathcal{D}_C, m_\alpha)|))$ .

In the case of ONTOPOD – *No*, where the set of patterns  $\mathcal{Pset}(\mathcal{D}_C, m_\alpha)$  is empty, the complexity is only in  $O(|\mathcal{I}| \times |\mathcal{D}_C|)$ .

## V EXPERIMENTS

We will show in this section the effectiveness of our method by presenting the dataset used and the different measures to evaluate our measure.

**Datasets.** We use two DBpedia datasets for our experiments: online and benchmark datasets [23]. The benchmark dataset lacks information on the properties of the entities, which is why we use the online dataset.

*DBpedia Online:* DBpedia Online is a set of human-understandable online RDF schema descriptions accessible to applications. The data is organized in a hierarchical structure. In this dataset, we will extract, for a class  $C$ , its entities (name, properties, domain, and range of properties), its super-classes, and sub-classes. As already stated these classes contain outliers. We will use the benchmark dataset to test the performance of our methods.

*Benchmark dataset:* This benchmark dataset is a reference database containing 342,781 data instances that can be used for hierarchical classification tasks. It has 3 levels 11, 12, 13, with respectively 9, 70, and 219 classes. We will focus on some classes of level 2. This dataset contains the real classes of the entities but does not cover the data of the online database.

**Extracting information from DBpedia online.** We have chosen to work with some DBpedia classes. These are the classes «Animal,» «Artist,» «Athlete,» «Company,» «EducationalInstitution,» «Group,» «Politician,» and «NaturalPlace.» We first choose a class C of these classes for the online extraction and then use SPARQL queries to extract its entities, super-classes, and sub-classes. The entities of the sub-classes of C are also part of the entities of C. The properties of each entity and the domains and ranges are also extracted. For instance, for the Artist class, the following SPARQL queries were executed:

```
1 SELECT distinct ?e ?P ?domain ?range
2 WHERE {
3   ?e a dbo:Artist.
4   ?e ?P ?objet.
5   ?P rdfs:domain ?domain.
6   OPTIONAL { ?P rdfs:range ?range }
7 }
```

```
1 SELECT distinct ?C
2 WHERE {
3   dbo:Artist rdfs:subClassOf* ?C.
4 }
```

```
1 SELECT distinct ?C
2 WHERE {
3   ?C rdfs:subClassOf+ dbo:Artist.
4 }
```

After the extraction of the online information, for a class C, we keep only the resources having at least two properties that are in the gold standard since it is these resources whose real class is known.

In Table 7, we give the number of entities, the number of outliers, and the number of inliers in each class. A set of properties describes each entity.



Table 7: Number of outlier and inlier entities in classes

| $\mathcal{D}$                 | $ Entities $ | $ Outliers $ | $ Inliers $ |
|-------------------------------|--------------|--------------|-------------|
| <i>Animal</i>                 | 1991         | 20           | 1971        |
| <i>Artist</i>                 | 2484         | 96           | 2388        |
| <i>Athlete</i>                | 3040         | 16           | 3024        |
| <i>Company</i>                | 7848         | 1013         | 6835        |
| <i>EducationalInstitution</i> | 3477         | 15           | 3462        |
| <i>Group</i>                  | 1519         | 386          | 1133        |
| <i>Politician</i>             | 6920         | 210          | 6710        |
| <i>NaturalPlace</i>           | 8195         | 12           | 8183        |

**Frequent and rare patterns extraction and Evaluation metrics.** To detect frequent and rare patterns, we used the *fpgrowth* [4] algorithm to find all supports of the patterns in the database. Then we separated the database into frequent and rare patterns using the average support as a separation threshold.

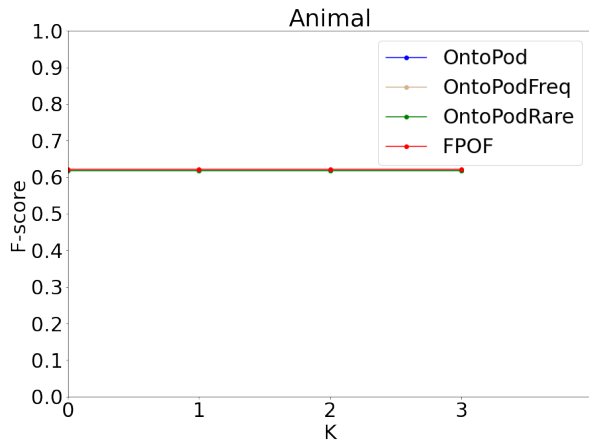
For the evaluation of the methods, we used three metrics: the inliers rate *IR*, the outliers rate *OR* and the harmonic mean of the *IR* and the *OR* called *F-score*. For each method *B*, we calculate  $IR(B)$ ,  $OR(B)$  and  $F-score(B)$ .

$$IR(B) = \frac{ni_B}{ni_T} \quad OR(B) = \frac{no_B}{no_T} \quad F-score = 2 * \frac{IR * OR}{IR + OR}$$

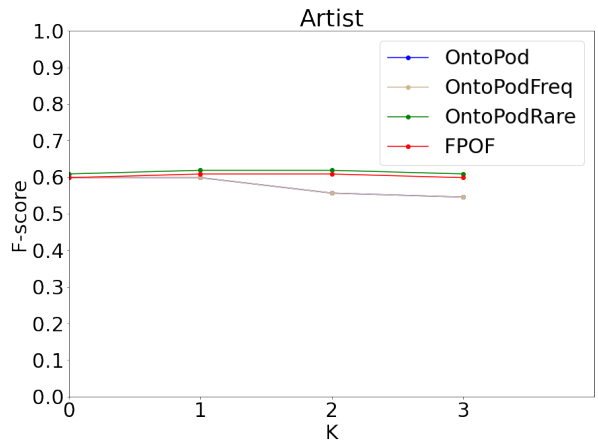
Where  $ni_B$  is the number of inliers found by the method *B*, and  $ni_T$  is the total number of inliers in the database, and  $no_B$  is the number of outliers found by the method *B* and  $no_T$  is the total number of outliers in the database.

There are two sets of experiments conducted in each class to evaluate the performance of the methods. The first one is the variation of the f-score of the methods according to the value of *k* and the second one is the outlier rate and the f-score of the methods for the best value of *k*.

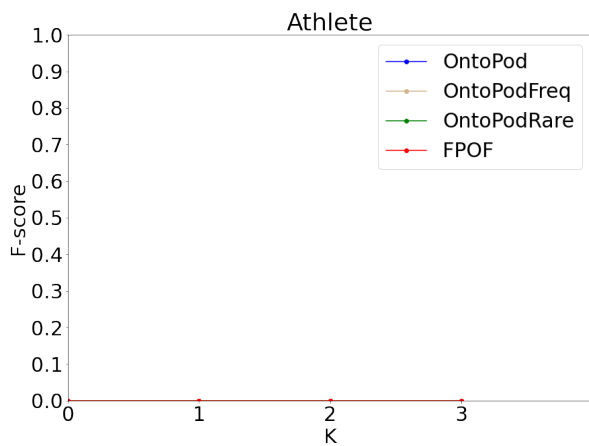
**F-scores evolution according to *k*.** Figure 3 shows the evolution of the F-scores of the methods according to the value of *k*. We notice that, in general, the best values of the F-score correspond to  $k = 0$ . The figures show that our method outperforms the FPOF method. For some classes, the FPOF method cannot even find outliers. Most outliers are entities of superclasses, and their properties are very frequent.



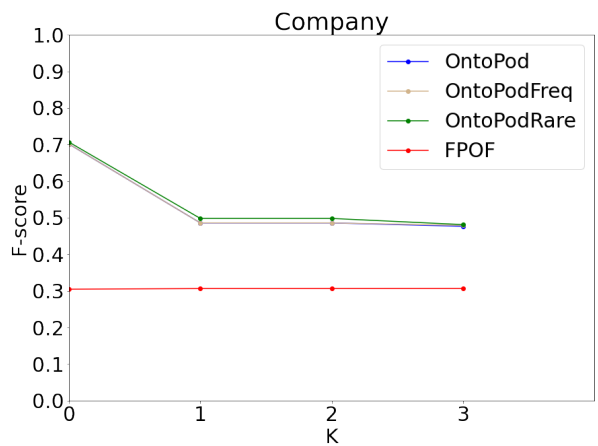
(a)



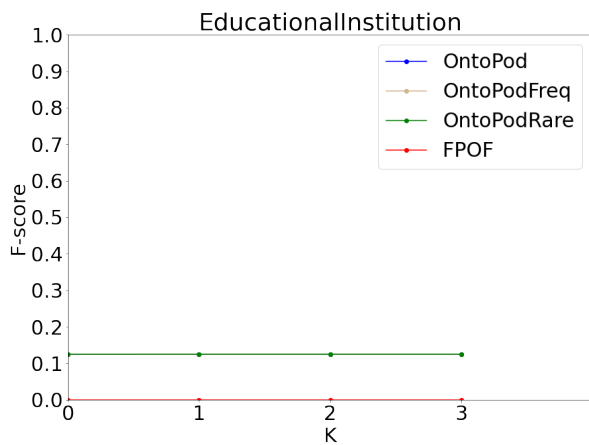
(b)



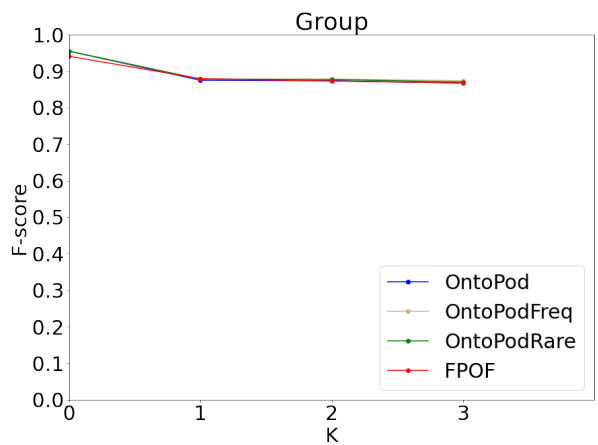
(c)



(d)



(e)



(f)

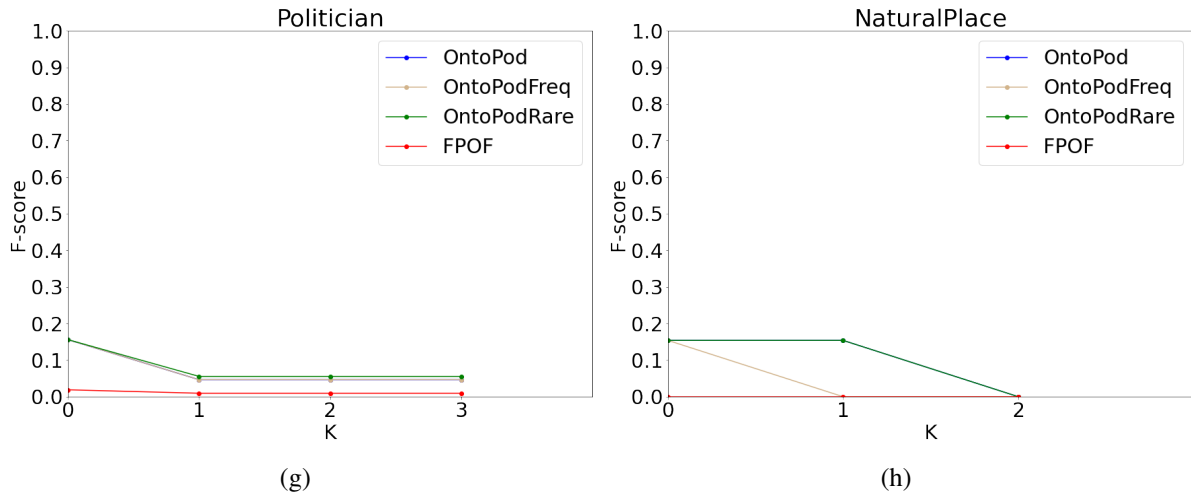
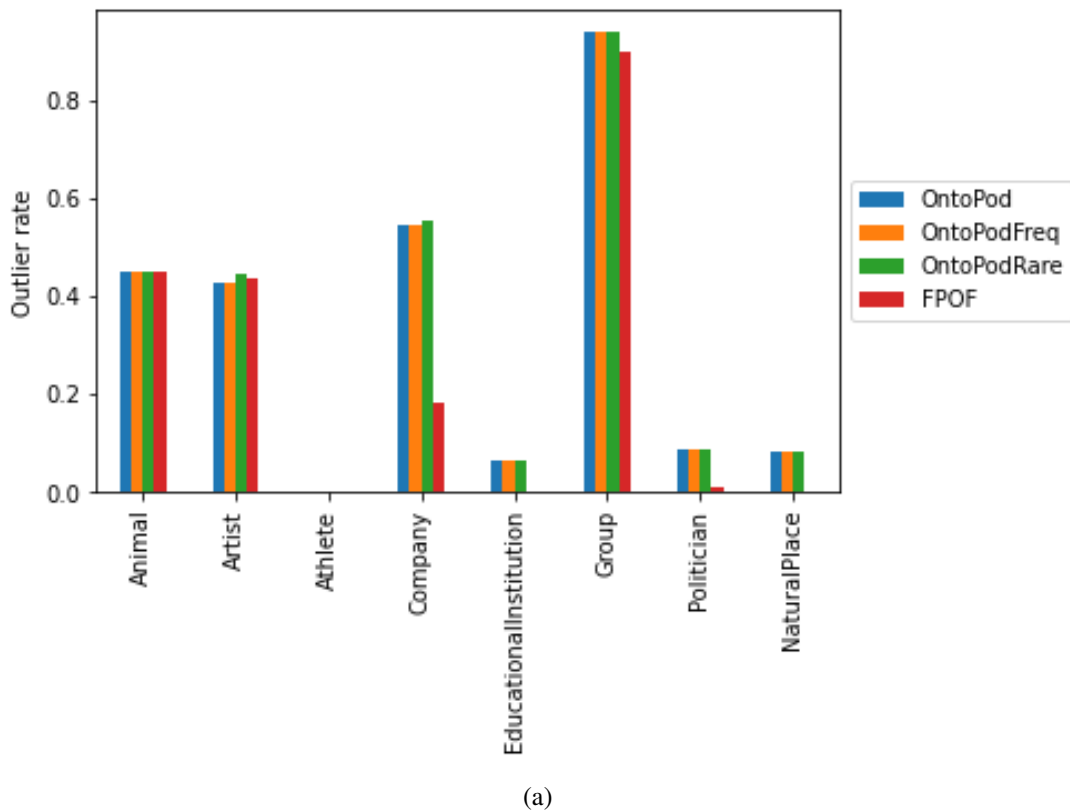
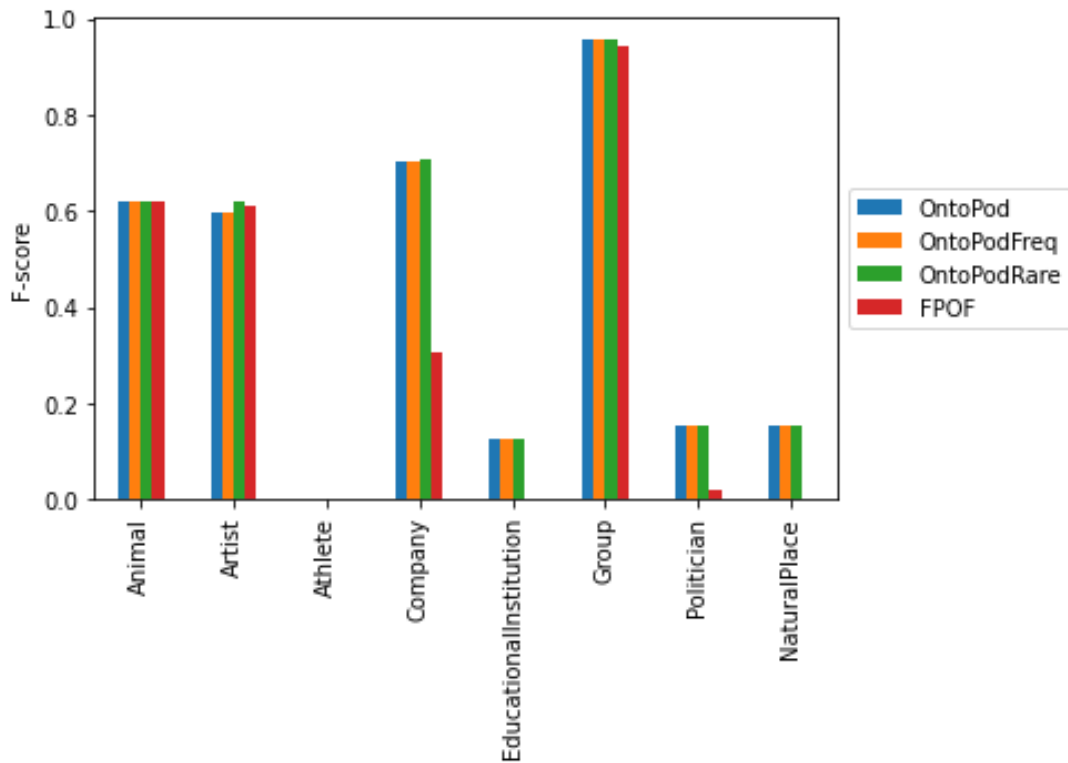


Figure 3: F-score of methods for each class

**Performance of the methods on the different classes.** Figure 4 shows the performance of the methods on the different classes. Figure 4(b) is the result of figure 3 for the maximum f-score value (which generally corresponds to  $k = 0$ ). Figure 4(a) shows the outlier rate for each method in each class. For the classes Company, EducationalInstitution, Politician, and NaturalPlace our method greatly outperforms the FPOF method. The outliers have frequent properties that are why the FPOF method fails to find them.





(b)

Figure 4: Outlier rate and F-score per class

## VI CONCLUSION

This paper proposed a generic semantic measure of outlier detection that favors some entities and disfavors others from the semantics behind the ontology (especially the properties and class hierarchy). We have tested our measure on a gold standard dataset of DBpedia. The experiments showed the effectiveness of our approach. In semantic web knowledge graphs, ontology plays a vital role that should not be neglected. Our measure is based on the domains and ranges of properties, the hierarchy of classes, the frequency of properties which are excellent parameters to find the entities outliers. If there are instances classified in class C that have (frequent) properties which have for a domain (and range) the class C (or its super-classes), we can reasonably suppose that this increases the precision of their membership in C. It is clear that if the ontology is not complete (i.e., it does not contain all the classes and properties of the data) or the real entities are poorly described (for instance, they have few properties), the performance of our measure decreases.

In our future work, we will add a preprocessing step using this outlier detection measure in the automatic classification process of DBpedia resources. Indeed, to build the learning base, it is necessary to select positive and negative examples. The measure will allow us to find true positive examples. Finally, we will apply this method to other Knowledge Graphs like Yago and Wikidata.

## REFERENCES

### Publications

- [1] Y. Wand and R. Y. Wang. “Anchoring data quality dimensions in ontological foundations”. In: *Commun. ACM* 39 (1996), pages 86–95.
- [2] E. M. Knorr and R. T. Ng. “Algorithms for Mining Distance-Based Outliers in Large Datasets”. In: *VLDB*. 1998.
- [3] M. M. Breunig, H.-P. Kriegel, R. T. Ng, and J. Sander. “**LOF: Identifying Density-Based Local Outliers**”. In: *Proceedings of the 2000 ACM SIGMOD International Conference on Management of Data*. SIGMOD '00. Dallas, Texas, USA: Association for Computing Machinery, 2000, pages 93–104. ISBN: 1581132174.
- [4] J. Han, J. Pei, and Y. Yin. “**Mining Frequent Patterns without Candidate Generation**”. In: *SIGMOD Rec.* 29.2 (May 2000), pages 1–12. ISSN: 0163-5808.
- [5] S. Ramaswamy, R. Rastogi, and K. Shim. “**Efficient Algorithms for Mining Outliers from Large Data Sets**”. In: *SIGMOD Rec.* 29.2 (May 2000), pages 427–438. ISSN: 0163-5808.
- [6] Z. He, X. Xu, J. Z. Huang, and S. Deng. “FP-outlier: Frequent pattern based outlier detection”. In: *Comput. Sci. Inf. Syst.* 2 (2005), pages 103–118.
- [7] Z. abu bakar, R. Mohamad, A. Ahmad, and M. Mat Deris. “**A Comparative Study for Outlier Detection Techniques in Data Mining**”. In: July 2006, pages 1–6.
- [8] C. Matuszek, J. Cabral, M. Witbrock, and J. DeOliveira. “An Introduction to the Syntax and Content of Cyc.” In: Jan. 2006, pages 44–49.
- [9] F. Baader, D. Calvanese, D. McGuinness, D. Nardi, and P. Patel-Schneider. *The Description Logic Handbook: Theory, Implementation, and Applications*. Jan. 2007.
- [10] K. Bollacker, R. Cook, and P. Tufts. “Freebase: A Shared Database of Structured General Human Knowledge.” In: Jan. 2007, pages 1962–1963.
- [11] F. Suchanek, G. Kasneci, and G. Weikum. “**YAGO: a core of semantic knowledge**”. In: Jan. 2007, pages 697–706.
- [12] C. Bizer, J. Lehmann, G. Kobilarov, S. Auer, C. Becker, R. Cyganiak, and S. Hellmann. “**DBpedia - A crystallization point for the Web of Data**”. In: *Journal of Web Semantics* 7.3 (2009). The Web of Data, pages 154–165. ISSN: 1570-8268.
- [13] J. Ren, Q. Wu, C. Hu, and K. Wang. “**An Approach for Analyzing Infrequent Software Faults Based on Outlier Detection**”. In: *2009 International Conference on Artificial Intelligence and Computational Intelligence*. Volume 4. 2009, pages 302–306.
- [14] A. Carlson, J. Betteridge, B. Kisiel, B. Settles, E. Hruschka, and T. Mitchell. “Toward an Architecture for Never-Ending Language Learning.” In: volume 3. Jan. 2010.
- [15] W. Zhang, J. Wu, and J. Yu. “**An Improved Method of Outlier Detection Based on Frequent Pattern**”. In: *2010 WASE International Conference on Information Engineering*. Volume 2. 2010, pages 3–6.
- [16] A. M. Said, D. D. Dominic, and B. B. Samir. “Outlier Detection Scoring Measurements Based on Frequent Pattern Technique”. In: *Research Journal of Applied Sciences, Engineering and Technology* 6 (2013), pages 1340–1347.
- [17] J. Lehmann, R. Isele, M. Jakob, A. Jentzsch, D. Kontokostas, P. Mendes, S. Hellmann, M. Morsey, P. Van Kleef, S. Auer, and C. Bizer. “**DBpedia - A Large-scale, Multilingual Knowledge Base Extracted from Wikipedia**”. In: *Semantic Web Journal* 6 (Jan. 2014).
- [18] A. Melo, H. Paulheim, and J. Völker. “Type Prediction in RDF Knowledge Bases Using Hierarchical Multilabel Classification”. In: *Proceedings of the 6th International Conference on Web Intelligence, Mining and Semantics* (2016).

- [19] H. Paulheim. “Knowledge graph refinement: A survey of approaches and evaluation methods”. In: *Semantic Web* 8 (Dec. 2016), pages 489–508.
- [20] P. Ristoski. “Exploiting semantic web knowledge graphs in data mining”. In: *Studies on the Semantic Web*. 2018.
- [21] D. Caminhas, D. Cones, N. Hervieux, and D. Barbosa. “Detecting and Correcting Typing Errors in DBpedia.” In: *DI2KG@ KDD*. 2019.
- [22] S. Issa. “Linked data quality : completeness and conciseness”. Theses. Conservatoire national des arts et metiers - CNAM, Dec. 2019.
- [23] D. Ofer. *DBPedia Classes : Hierarchical Taxonomy of Wikipedia article classes*. <https://www.kaggle.com/danofer/dbpedia-classes>. [Online; accessed 02-November-2021]. 2019.
- [24] L. Diop, C. Diop, A. Giacometti, and A. Soulet. “Pattern Sampling in Distributed Databases”. In: Aug. 2020, pages 60–74. ISBN: 978-3-030-54832-2.
- [25] F. Al-Aswadi, S. Mishra Tiwari, and D. Gaurav. “Recent trends in knowledge graphs: theory and practice”. In: *Soft Computing* 25 (July 2021).
- [26] L. Diop, C. Diop, A. Giacometti, and A. Soulet. “Pattern on demand in transactional distributed databases”. In: *Information Systems* 104 (Oct. 2021), page 101908.

## A ACKNOWLEDGEMENTS

This research was supported by the Partnership for skills in Applied Sciences, Engineering and Technology (PASET) - Regional Scholarship and Innovation Fund (RSIF).