



# Logical Detection of Invalid SameAs Statements in RDF Data

Laura Papaleo, Nathalie Pernelle, Fatiha Saïs, Cyril Dumont

► **To cite this version:**

Laura Papaleo, Nathalie Pernelle, Fatiha Saïs, Cyril Dumont. Logical Detection of Invalid SameAs Statements in RDF Data. In proceedings of the 19th International Conference on Knowledge Engineering and Knowledge Management, EKAW 2014, Nov 2014, Linköping, Sweden. 2014, <10.1007/978-3-319-13704-9\_29>. <hal-01275943>

**HAL Id: hal-01275943**

**<https://hal.inria.fr/hal-01275943>**

Submitted on 18 Feb 2016

**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.

# Logical Detection of Invalid SameAs Statements in RDF Data

Laura Papaleo<sup>1</sup>, Nathalie Pernelle<sup>1</sup>, Fatiha Saïs<sup>1</sup>, and Cyril Dumont<sup>1</sup>

Université de Paris-Sud, Laboratoire de Recherche en Informatique,  
Bâtiment 650, F-91405 Orsay Cedex, France  
firstname.lastname@lri.fr,  
home page: <http://www.lri.fr>

**Abstract.** In the last years, thanks to the standardization of Semantic Web technologies, we are experiencing an unprecedented production of data, published online as *Linked Data*. In this context, when a typed link is instantiated between two different resources referring to the same real world entity, the usage of *owl:sameAs* is generally predominant. However, recent research discussions within the Linked Data community have shown issues in the use of *owl:sameAs*. Problems arise both in cases in which sameAs is automatically discovered by a data linking tool erroneously, or when users declare it but meaning something less 'strict' than the semantics defined by OWL. In this work, we discuss further this issue and we present a method for logically detect invalid sameAs statements under specific circumstances. We report our experimental results, performed on OAEI datasets, to prove that the approach can be promising.

**Keywords:** RDF identity link, Linking validation, SameAs statement, Data and Link quality, Semantic Web, RDF

## 1 Introduction

As stated by the W3C, the Semantic Web provides a common framework that allows data to be shared and reused across applications, enterprises, and communities [29]. Generally speaking, the Semantic Web is a '*Web of Data*', where data can be processed by machines, extending the principles of the Web from *documents* to *data* [2]. In this context, resources (data) can be accessed using the conventional Web architecture (e.g. URIs) and it is possible to link resources using relations. In the Semantic Web framework, each relation (link) is named, improving the expressiveness of the connections and the semantic meaning of interrelations.

Today, we are experiencing an unprecedented production of data, published as *Linked Open Data* (LOD, for short). This is leading to the creation of a global data space containing billions of assertions [3]. RDF [14], which is one of the fundamental building blocks of the Semantic Web, provides formal ways to build these assertions.

Working in the LOD is basically about using the Web to create *typed links* between data from different sources: providers set RDF links from entities (described by URIs) to related entities (other URIs), improving the knowledge related to a specific resource and, thus, the global knowledge in the Web of Data. Most of the RDF links connecting resources coming from different data sources are *RDF identity links*, called also *sameAs statements*. They are defined using the *owl:sameAs* property, thus expressing that two URI references actually refer to the same thing. Unfortunately, many existing identity links do not reflect such genuine identity, as argued recently within the research community [7, 5].

So, as numerous independently developed data sources have been published over internet as Linked Data, the *problem of identity* is now casting a shadow over the shininess of the Semantic Web [13, 7]. Thus, it is becoming extremely important to develop means of *data and linking quality assurance*. The study of the quality of data and links in the LOD cloud may be particularly useful in applications that want to consume Linked Data as well as in Semantic Web frameworks dedicated to data linking or data integration.

In this work, we investigate and design a logical method to detect invalid sameAs statements in RDF data, by looking at the descriptions associated to the instances involved. We suppose that, in case of multiple data sources, mappings between properties are provided. Our approach is local, in the sense that, we build a contextual graph 'around' each one of the two resources involved in the sameAs statement and we study the descriptions provided in these contextual graphs. The construction of the contextual graph is based on properties that have specific characteristics (functional, local completeness). We claim that, when logical conflicts are encountered, the initial RDF identity link is 'inconsistent', meaning that it requires further investigation (supervised or automatic). We tested the approach on sameAs statements provided by linking tools that have been applied on Ontology Alignment Evaluation Initiative (OAEI) datasets, showing that our research direction is promising. We also introduce a set of rules which could eventually support the re-qualification of inconsistent sameAs statements when these statements involve resources that are not identical, in the *owl:sameAs* semantics, but similar as related to the same abstract concept.

The remainder of this paper is organized as follows. In Section 2 related works are described. In Section 3, we present the conceptual building blocks of our approach, while Sections 4 and 5 are dedicated to the formulation of the problem, the generation of the rules and the introduction of a possible re-qualification of inconsistent sameAs. In Sections 6 and 7 we present the logical method and the experimental results on sameAs links computed for OAEI datasets. Finally in Section 8 some concluding remarks are drawn.

## 2 Related Works

How to evaluate and assess the quality of data and links in the Linked Data Cloud is a generally novel problem, growing its importance in the last years, as the research community can, now, work with a massive quantity of data coming

from multiple data sources.

In [6] the authors present a 'global approach' where they analyzed the structural properties of large graphs of identity links focusing the attention on general network properties such as degree distributions and URI counts, without analyzing the quality. Recently, in [11] the authors describe another global approach in a framework dedicated to the assessment of Linked Data mappings using network metrics. Five different metrics have been performed on a set of known good and bad links concluding that most of these metrics are not meaningful with respect to the evaluation of the quality of an identity link. In [5], the author illustrates how to assess the quality of *owl:sameAs* links, using a constraint-based method. In the work, an interesting formalization of the problem in a graph-based fashion is presented, but the evaluation of the quality of the identity link is, in the end, performed using only one property, namely the name of each entity. The results are interesting, but as claimed but the author himself, it could be important to include advanced similarity measures and the evaluation of more properties. In [7, 12] the authors studied the problem of the quality of RDF identity links from a general point of view, making observations about the varying use of *owl:sameAs* in Linked Data. They proposed an ontology called the Similarity Ontology (SO) that aims at better classifying the different level of similarity between items in different data sources. However, the quality evaluation of the *owl:sameAs* links is performed manually, assessing around 250 *owl:sameAs* links in an Amazon Mechanical Turk experiment.

In this paper we propose a method which analyzes more information than simply the resource name, as opposite to [5]. Our approach is 'local', differently from [6, 11] as we assess the correctness of a sameAs statement by studying the information described in contextual sub-graphs built according to specific criteria. To the extent of our knowledge, there not exist similar logical methods in the literature.

To complete this Section, we need to recall that there exist a lot of interesting methods related to *owl:sameAs* link discovery (see [8, 9] as recent surveys). This is also referred as the 'coreference problem' in Semantic Web. The reader can, for example, see the works in [18, 19, 10, 23, 21] or, more recently those in [28, 26]. However, *sameAs statement quality assessment* is, generally, different from the coreference problem. From the 'coreference prospective', the goal is to analyze the knowledge related to two resources in order to decide if one new assertion can be added to the knowledge base. In various domains, there are generally accepted naming schemata [3]. If two resources in the knowledge base both support one or more of these identification schema, the implicit relationship between entities can be made explicit as identity link, automatically. This can be true, for example, in case of a unique code such as the italian 'Codice Fiscale' that can be derived through a deterministic algorithm from a person's name and his/her date and place of birth, or the International Standard Book Number (ISBN) which is a unique numeric commercial book identifier. When no shared naming schema exist, RDF identity links are usually generated by evaluating the similarity of entities using more or less complex similarity functions.

These functions generally take into account sub-parts of resource description that is known to be discriminative enough as, for example, inverse functional properties or composite keys.

Few linking tools are interested in generating *owl:differentFrom* links, as for example [24]. This idea of partitioning the resources into groups of 'different resources' is used also in blocking methods as [25]. Then, data linking tools will search for *sameAs* links only within a group. However, in such approaches, only direct data-type properties are taken into account.

Instead, once a *sameAs* statement exists in the knowledge base, it could be interesting to analyze different properties (not only inverse functional). To clarify this point, let us consider a very simple example: we have two resources (books)  $b_1$  and  $b_2$  both described using two data-type properties *isbn* and *pages*. We assume, for example, that the property *isbn* is inverse functional and *pages* is only functional. In order to infer *sameAs*( $b_1, b_2$ ), it is sufficient to check if the values of *isbn* are equal. Using the semantics of *owl:sameAs* it is possible to infer that the values of the property *pages* are equivalent. If they are not, one can detect a conflicting case contingent on the semantics of *sameAs*( $b_1, b_2$ ).

In conclusion, it is sure that the two problems are entailed and, as immediate future activities, we are planning to deepen the analysis of their interconnection, especially in the case of complex and hybrid linking methods that are recently emerging.

### 3 Preliminaries

In this Section we present the theoretical framework in order to define the building blocks for the logical invalidation approach.

**Definition 1. *RDF Graph.*** [17]

*An RDF graph is a set of RDF triples. The set of nodes of an RDF graph is the set of subjects and objects of triples in the graph.*

Given an infinite set  $U$  of URIs, an infinite set  $B$  of blank nodes and an infinite set  $L$  of literals, a RDF triple is a triple  $\langle s, p, o \rangle$  where the subject  $s \in (U \cup B)$ , the predicate  $p \in U$  and the object  $o \in (U \cup B \cup L)$ . A *RDF triple* represents an assertion a 'piece of knowledge', so if the triple  $\langle s, p, o \rangle$  exists, the logical assertion  $p(s, o)$  holds (is true).

An RDF graph  $G$  is simply a collection of RDF triples and it can be seen as a set of statements describing partially (or completely) a certain knowledge. Note that the knowledge described by an RDF graph  $G$  is contained in all the three elements in each RDF triple, since the possibility to create 'typed predicates' is one of the strength of the Web of Data.

**Definition 2. *SameAs Statement.*** [6]

*A SameAs statement *sameAs*( $s, o$ ) is an RDF triple  $\langle s, owl : sameAs, o \rangle$  in an RDF graph  $G$  which connects two RDF resources  $s$  and  $o$  by means of the *owl:sameAs* predicate.*

The built-in OWL property *owl:sameAs* links two resources in the RDF graph. Such an *owl:sameAs* statement indicates that two URI references refer to the same *thing* : the individuals have the same 'identity' [20, 16]. Given an RDF graph  $G$  as defined above, the following holds

$$\text{sameAs}(x, z) \wedge P(x, y) \Rightarrow P(z, y)$$

where the  $x, z \in (U \cup B)$ , the predicate  $P \in U$  and the object  $y \in (U \cup B \cup L)$ . *owl:sameAs* supports linking of data from diverse sources in a principled way. The underlying semantics provides support for inferences over this data but, as discussed before, this may yield unexpected or contradictory results. As explained in details in the following sections, scope of this work is to evaluate the correctness of a set of provided SameAs statements, using the notions of property-based walk and contextual graph we introduce now.

**Definition 3. Property-based walk of length  $n$**   $w_{\{n,s,P\}}$ .

Given an RDF graph  $G$ , a node  $s$  in  $G$ , given a set  $P$  of properties defined for  $G$ , a Property-based walk of length  $n$   $w_{\{n,s,P\}}$  is an alternating sequence of nodes and predicates  $\{v_0 \equiv s, p_0, v_1, p_1, v_2, \dots, v_{n-1}, p_{n-1}, v_n\}$ , such that

- $v_0, \dots, v_{n-1}$  are resources in  $G$ ,  $\forall i = 0, \dots, n-1$   $v_i \in U$ ,
- $v_n$  is a literal in  $G$ ,  $v_n \in L$
- each triple  $\{v_i, p_i, v_{i+1}\}$  in the sequence is an RDF triple in  $G$  such as  $p_i \in P$
- all the resources in the walk are distinct from one another. Thus, for each pair of resource  $\{v_i, v_j\}$ ,  $v_i$  and  $v_j$  are not the same resource, with  $\{i, j\} \in [0, \dots, n-1]$  (they have different URIs).

In the definition we suppose that each predicate in  $G$  has an associated weight 1 that expresses its existence (its length).  $w_{\{n,s,P\}}$  is basically a path in the RDF graph without cycle and of length  $n$ , involving  $n+1$  node,  $n$  resources defined by URIs and 1 node as a literal. It can be seen also as a collection of assertions selected according to specific conditions (the starting resource  $s$  and the set of properties  $P$ ).

In other words, with a walk  $w_{\{n,s,P\}}$  in the graph  $G$ , we select a sequence of assertions in some way related to the resource  $s$ . This means also that, for every RDF triple  $\langle v_i, p_i, v_{i+1} \rangle$  in  $w_{\{n,s,P\}}$  the fact  $p_i(v_i, v_{i+1})$  holds.

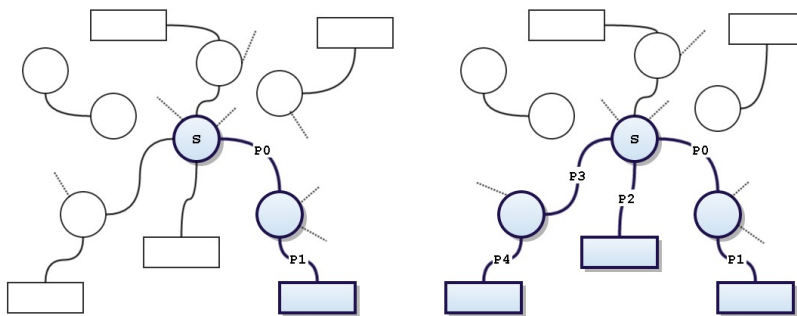
**Definition 4.  $m$ -degree Contextual Graph**  $G_{\{m,s,P\}}$

Given an RDF graph  $G$  and a node  $s \in G$ ,  $s \in U$ , an integer number  $m$  and a set  $P$  of properties defined for  $G$ , a  $m$ -degree Contextual Graph  $G_{\{m,s,P\}}$  for  $s$  is a sub-graph of  $G$  such that every node  $v_i \in G_{\{m,s,P\}}$  belong to a property-based walk of length  $n$ , with  $n \leq m$ .

A  $m$ -degree contextual graph for a resource  $s$  can be seen as a subset of knowledge pertinent to  $s$ , bounded by the set of predicates  $P$ .

In Figure 1 we depict two examples: a walk (on the left) and a contextual graph (on the right). Given an RDF graph  $G$ , in which circles identify resources with URI and rectangles represent literals, the Figure 1-(left) shows

a walk  $w_{\{2,s,P=\{P_0,P_1\}\}}$  for the resource  $s$ . The walk has length 2 and involves the properties  $P_0$  and  $P_1$ . Figure 1-(right) shows a 2-degree contextual graph  $G_{\{2,s,P=\{P_0,\dots,P_4\}\}}$  for the same resource  $s$ . It involves the properties  $P_0, P_1, P_2, P_3$  and  $P_4$ .



**Fig. 1.** (left) a walk of degree 2 for the resource  $s$ ,  $w_{\{2,s,P=\{P_0,P_1\}\}}$ . (right) The contextual graph  $G_{\{2,s,P=\{P_0,\dots,P_4\}\}}$  of degree 2 for the same resource  $s$ .

## 4 Problem Statement

This Section is dedicated to the explanation of the problem using the definitions provided before. The general situation is to have a portion of knowledge in the form of an RDF graph and an assertion declaring that two resources in the graph are the same (via the predicate *owl:sameAs*). The problem we are addressing is to check if a sameAs statement can be invalidated and eventually explain this deduction. We need to check if inconsistencies in the assertion  $sameAs(x, y)$  according to the knowledge provided in the RDF graph  $G$ .

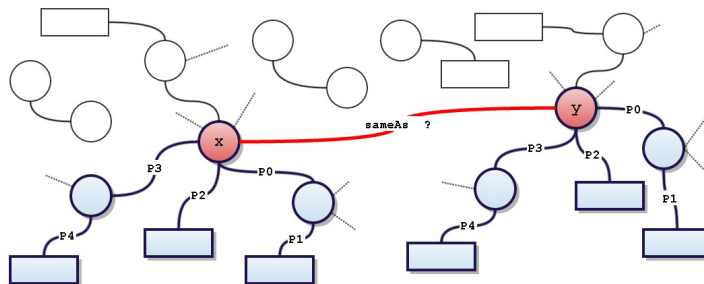
Our approach relies on building two contextual graphs (see Section 3), for  $x$  and  $y$  respectively and on reasoning on the assertions contained in these two graphs. The building blocks of the problem are the following:

- An RDF graph  $G$
- two resources  $x$  and  $y$ , such that  $x, y$  are resources in  $G$
- the triple  $\langle x, owl : sameAs, y \rangle$  (or  $sameAs(x, y)$ ) belonging to  $G$
- a set of properties  $P$  in  $G$
- a value  $n$  representing the depth of the contextual graphs
- the contextual graphs  $G_{\{n,x,P\}}$  and  $G'_{\{n,y,P\}}$  for  $x$  and  $y$

The problem becomes the evaluation of the following rule:

$$G_{\{n,x,P\}} \wedge G'_{\{n,y,P\}} \wedge sameAs(x, y) \Rightarrow \perp$$

The construction of the contextual graphs depends on the predicates (properties) we select and the value  $n$ . Indeed, in complex RDF graph, which can combine data coming from multiple data sources, limiting the depth of a contextual graph could be wise. The main reason is that long property-based walks can lead to not relevant piece of information which can eventually confuse the validation process.



**Fig. 2.** The statement  $sameAs(x, y)$  must be validated. The two 2-degree contextual graphs extracted for  $x$  and  $y$  are highlighted.

In Figure 2 we show an example of what we want to build. In this case, the statement  $sameAs(x, y)$  must be validated, and a value  $n = 2$  has been selected. The set of properties  $P$  has been defined as  $\{P_0, \dots, P_4\}$ . The image shows the two contextual graphs extracted for  $x$  and  $y$ . In the following Section we explain how we want to choose the predicates.

## 5 Properties Selection and Rules Generation

In this work, we chose to use functional properties and those properties declared as local complete. Here, we explain and motivate this choice, describing the logical rules we add in the resolution system. Additionally we propose a possible re-qualification of some particular  $sameAs$  statements.

### 5.1 Functional and Inverse Functional Properties

Let us suppose that  $p_1$  is a functional property. It can be expressed logically as follows [20, 16]:

$$p_1(r, v) \wedge p_1(r, v') \Rightarrow v \equiv v'$$

So if we want to validate  $sameAs(x, y)$  and we have a mapped functional property  $p_1$ , with  $p_1(x, w)$  and  $p_1(y, w_1)$ , and we can assert in some way that  $w \not\equiv w_1$  then:

$$sameAs(s, o) \wedge p_1(s, w) \wedge p_1(o, w_1) \wedge w \not\equiv w_1 \Rightarrow \perp$$



We have an inconsistency. In this situation, if we assume that the assertions already in the RDF graph are true and we have 'doubts' only on the sameAs statement, we can conclude that this latter has problems. In our approach, taking into consideration functional properties, we basically add the following rules for every property  $p_i, p_j, p_k$  in the contextual graphs we are considering.

$$R_{1_{FDP}} : \text{sameAs}(x, y) \wedge p_i(x, w_1) \wedge p_i(y, w_2) \rightarrow \text{synVals}(w_1, w_2)$$

$$R_{2_{FOP}} : \text{sameAs}(x, y) \wedge p_j(x, w_1) \wedge p_j(y, w_2) \rightarrow \text{sameAs}(w_1, w_2)$$

$$R_{3_{IFP}} : \text{sameAs}(x, y) \wedge p_k(w_1, x) \wedge p_k(w_2, y) \rightarrow \text{sameAs}(w_1, w_2)$$

Note that  $R_{1_{FDP}}$  is for data-type properties and  $R_{2_{FOP}}$  and  $R_{3_{IFP}}$  are for object-type properties.  $\text{synVals}$  and  $\neg\text{synVals}$  are further described in Section 6.

Given a property  $p$  in the graph  $G$ , the knowledge of  $p$  being a functional property can be already present among the assertions in  $G$  or derived after, collecting knowledge from experts or gathering it externally (existing ontologies, additional assertions on the Web and so on.)

## 5.2 Local completeness

The closed-world assumption is in general inappropriate for the Semantic Web due to its size and rate of change [15]. But in some domains and specific contexts, local-completeness for RDF predicates (properties) could be assured. In fact, there are predicates for which there is no need to express negative information because the available positive information about them is complete and, consequently, the negative information is simply the complement of the positive information [30]. A good example for a multi-valued local complete property could be one representing the authors of a publication.

When a predicate is like that, it should be declared closed in the specific knowledge base (making a local completeness assumption). A Local Completeness (LC rule) rule specifies that the resource is complete for a subset(s) of information (on a particular ontology). In other words, the information contained in the resource is all the information for the subset (specified by the rule) of the domain. In an RDF graph  $G$ , it is possible to declare a rule for each property that fulfills the local completeness in the form :

$$R_{4_{LC}} : \text{sameAs}(x, y) \wedge p(x, w_1) \rightarrow p(y, w_1)$$

where  $\text{sameAs}(x, y)$  is a sameAs statement,  $p$  is a predicate defined in the RDF graph  $G$ ,  $x$  and  $y$  are object-type resources in  $G$  ( $x, y \in U$ ) and  $w_1$  is a literal ( $w_1 \in L$ ).

Also in this case, given a property  $p$ , the knowledge of 'local completeness' for  $p$  can be asserted an expert or derived after gathering additional knowledge in a semi-automatic way.

### 5.3 Re-qualification of Erroneous sameAs

When two resources are not the same they do not necessarily are completely different. Indeed, it is possible that erroneous identity links could actually involve resources that, in some way, represent the same 'idea' (abstract concept) but at different levels of details. In this situation, a sameAs statement may not be an identity link, meaning that it is not correct to use the *owl:sameAs* predicate, but it can be a 'sameAs domain-dependent link'.

To clarify our idea, we provide a simple example. Let  $b_1$  and  $b_2$  be two books that represent the book entitled '*Madame Bovary*' written by '*Gustave Flaubert*'. The two books do not share the *ISBN*, the editor, the number of pages and the language, for example. In this case the relation that should be asserted between  $b_1$  and  $b_2$  should not be a *owl:sameAs* but a more abstract link, we call *sameAbsObject*. In this case this link expresses the *sameArtOfWork*.

To detect such abstract links, we propose to use a set of expert rules that express abstract key constraints that identify the abstract entities. For example, an expert can express that  $\{title, author\}$  is an abstract key of the class *book*.

The logical semantics of an abstract key  $k = \{p_1, \dots, p_n\}$  is:

$$c(x) \wedge c(y) \wedge \left( \bigwedge_{p \in k} \exists w p(x, w) \wedge p(y, w) \right) \rightarrow sameAbsObject(x, y)$$

where  $c(\cdot)$  indicates a class membership,  $p$  is a property and  $x, y$  are two resources with URI and  $w$  can be either a resource with URI or a literal.

The semantics of the construct *sameAbsObject* could be in some way similar to *owl:sameAs*, but only for the properties that can be 'related' to the abstract object. So, for each of such predicate  $p_{ao}$ , we have:

$$sameAbsObject(x, y) \wedge p_{ao}(x, w) \rightarrow p_{ao}(y, w)$$

Thus, in the case of books, a predicate like *has\_originalLanguage* is related to an abstract art of work. So, for the two books  $b_1$  and  $b_2$  described before, if we declare *sameAbsObject*( $b_1, b_2$ ), then the rule above holds for *has\_originalLanguage*. The same rule will not hold for properties not related to the abstract object, such as *pages* or *editor*, for example.

Other possible examples of *sameAbsObject* in domain dependent contexts can occur in music, events (such as festivals, concerts, conferences, and so on), even persons, depending on the social context.

Using the rules related to the definition of *sameAbsObject*, it is possible to propose a re-qualification of sameAs statements. In fact, given a sameAs statements with inconsistencies, if these inconsistencies do not involve properties that are related to the abstract object, the sameAs link could be replace by a *sameAbsObject* link.

## 6 The Invalidation Approach

In this Section we present our invalidation approach, on the basis of all the definitions and reasoning made so far. Given  $G$  the initial RDF graph with  $U$

the set of resources in  $G$  with URIs. Given  $sameAs(x, y)$  the input  $sameAs$  statement to validate, where  $x, y \in U$ . Let  $F$  be a set of facts, initially empty, and  $L$  the set of literals for  $G$ .

1. Build a set  $F_1$  of  $\neg synVals(w_1, w_2)$ , for each pair of semantically different  $w_1$  and  $w_2$ , with  $w_1, w_2 \in L$ .
2. Choose a value  $n$  indicating the depth of the contextual graphs
3. Build the contextual graphs for  $x$  and  $y$  considering functional properties and local complete properties
  - For all the functional properties  $p_{iFP}$  add the relative set of RDF facts to  $F$ , considering the rules  $R_{1FDP}, R_{2FOP}, R_{3IFP}$  in Section 5.
  - For each  $p_{iLC}$  that falls in the contextual graphs and fulfills the local completeness (i.e.  $R_{4LC}$  is declared), add to  $F$  a set of facts in the form  $\neg p_{iLC}(s, w)$  if  $w$  is different to all the  $w'$  s.t.  $p_{iLC}(s, w')$  belongs to  $F$ , using  $F_1$ . Note that  $w, w' \in L$ .
4. Apply iteratively unit resolution until saturation [22] using  $F \cup CNF^1\{R_{1FDP}, R_{2FOP}, R_{3IFP}, R_{4LC}\}$ .

Note also that disjointness of classes can be provided in input. So, for each pair of classes  $C$  and  $D$ , if a statement as  $DISJOINT(C, D)$  is declared, or such that their disjunction is inferred by inheritance, and  $x$  is of class  $C$  and  $y$  is of class  $D$ , we immediately infer that  $sameAs(x, y)$  is not correct.

The set of  $\neg synVals(w_1, w_2)$  with  $w_1, w_2 \in L$  can be obtained using different strategies. It is possible, for example, to perform a pre-processing step in which we build a clustering of the values according to specific criteria. To clarify, consider a simple example of names of cities in a specific domain: it is possible to pre-process all the possible values and assert that  $synVals('Paris', 'ParisCity')$  and that  $\neg synVals('Paris', 'Milan')$  and so on. Thus, the evaluation is based on determining if two values  $w_1, w_2$  belong to the same cluster. Another situation arises when the values are 'well defined' as in the case of enumeration, dates, years, geographical data or some types of measures. In these cases, the evaluation is again a simple syntactic comparison of the values. If they are the same, they are equivalent, otherwise they are not equivalent.

## 7 Experimental Results and Discussion

A prototype of our validation framework has been implemented in Java using the AIMA library [1] for the resolution. This first prototype does not include the re-qualification of  $sameAs$ .

In this Section we present the results of the experiments we performed for assessing the quality of the set of  $sameAs$  statements computed by different linking methods, respectively presented in [24], [27] and [31].

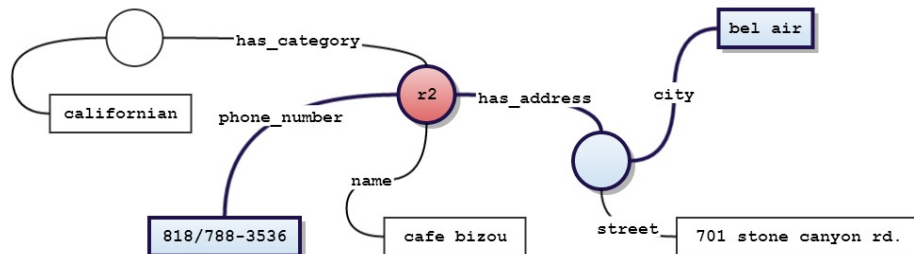
In [24] the  $sameAs$  statements are computed according to similarity measures over specific property descriptions, as in [31] where similarity between entities

---

<sup>1</sup> CNF: Conjunctive Normal Form

is iteratively calculated by analyzing specific features. In [27], instead sameAs statements are computed on the basis of a novel algorithm for key discovery. All the above methods have produced results on the Person-Restaurants (PR) data test available for the instance matching contest OAEI 2010 (IM@OAEI2010) [4]. For the key discovery method, we started from the links obtained by the method considering only the *name* of the restaurant as a key.

According to the knowledge base, we considered as 'meaningfully' functional the properties *phone\_number* and *has\_address* that describe a restaurant and *city* that describes an address<sup>2</sup>. Thus, given a sameAs statement in the form *sameAs(x, y)* we computed the contextual graph of degree 2 considering the three functional properties listed before. Figure 3 shows an example of the contextual graph computed for a restaurant in the first dataset 'restaurant1' (already mapped).



**Fig. 3.** An instance of restaurant in the dataset 'restaurant1' Given the functional properties *phone\_number*, *has\_address* and *city*, a contextual graph of degree 2 is depicted.

To build the  $\neg\text{synVals}$  (set  $F_1$ ) for the values of the properties selected, we did a normalization of the values. For *phone\_number*, we removed all the additional characters (e.g. '/', '-', and so on), leaving only the numbers. We note that the same number of digits are given for all the phone numbers. For *city*, we removed words which can be not meaningful such as 'city', the character '(' and so on. A  $\neg\text{synVals}$  is declared for each pair of syntactically different pairs of values.

To explain the results obtained, let us consider the answers collected by applying our invalidation approach on the sameAs statements computed by [27]. Note that, in this case, the analysis is performed using properties completely different from those used in the computation of the identity links.

Over the 90 sameAs statements computed, 4 were wrong with respect to the gold standard. We are able to detect 3 over 4 of these erroneous links. The

<sup>2</sup> Note that both the previous methods aligned the two initial datasets in order to compute the sameAs statements. We considered the same alignment in the explanation of the results.

only one we cannot detect is the one linking restaurant 91 defined in dataset 1 to restaurant 711 defined in dataset 2. By looking at the properties, the two restaurants share the same phone number and the same city. They even share the same street name. So an inconsistency cannot be detected. Most probably, they represent the same commercial site providing different services (they have different categories). In addition we classify as 'wrong sameAs' 5 statements which, with respect to the gold standard, are in fact good sameAs statements. The reason is that, in every statement  $sameAs(x, y)$ , the restaurants  $x$  and  $y$  have different phone number or different city (or both). This type of result can be seen dually. On the one hand this could mean, for example, that a restaurant can have two phone numbers, so maybe the property *phone\_number* is not functional. On the other hand, there can be errors in the data (for example 'los angeles' and 'los feliz') and the computation of the  $\neg synVals$  has been imprecise. In any case, the good idea is to highlight the inconsistency to the user (expert) and ask for confirmation or correction.

Linking Method	LM precision	TC	RG	TN	TP	FN	FP	accuracy	recall	IA precision	LM+IA precision
[27]	95.55%	90	4	81	3	1	5	93, 34%	75%	37%	98.85%
[24]	69.71%	142	43	94	38	5	5	92.9%	88.4%	88.4%	95.19%
[31]	90.17%	112	11	86	11	0	16	86.60%	100%	42.30%	100%

**Table 1.** Results of our approach on the sameAs links provided by the linking methods. We report the accuracy, recall and precision for the invalidation approach (IA) and the overall precision (LM+IA) in the last column.

Table 1 shows a tabular summary of our tests, including accuracy, recall and precision of the method. The table indicates as: (i) **TC**: total cases to be considered, namely the number of *sameAs* found by the linking algorithm. (ii) **RG**: the number of the sameAs statements really wrong, wrt the gold standard. (iii) **TN**: (true negative), the number of statements which we detected 'good' and were actually correct (wrt the gold standard). (iv) **TP**: (true positive), the number of statements which we detected 'wrong' and were actually wrong (wrt the gold standard). (v) **FP** (false positive), the number of statements which we detected 'wrong' but were actually correct. (vi) **FN**: (false negative), the number of real wrong statements which we could not detect.

By definition,  $accuracy = (TP + TN)/TC$ ,  $recall = TP/(TP + TN)$  and  $precision = TP/(TP + FP)$ . In other words, *accuracy* indicates the percentage of the predictions that were correct (where predictions for us are the wrong sameAs detected with inconsistencies or correct sameAs validated without inconsistencies), *precision* indicates the percentage of positive predictions which were correct (namely the numbers of wrong sameAs correctly detected over the total number of sameAs identified with inconsistencies) and *recall* indicates the percentage of positive cases caught (namely the number of wrong sameAs correctly detected over the total number of sameAs really wrong).

In conclusion, our results showed that, when our validation tool is applied after one of the linking tool, the precision of each tool can be improved, namely for [27] we pass from a precision of 95.55% to 98.85%, for [24] from a precision of 69.71% to 95.19% and finally for [31] from a precision of 90.17% to 100%.

## 8 Concluding remarks

In the last years, the amount of data published on online as Linked Data is growing significantly. In this context, the usage of *owl:sameAs* is generally predominant when linking resources from different data sources. Recent research discussions within the Linked Data community have shown that the use of *owl:sameAs* may be incorrect. The same community is demanding now innovative tools and methods to assure and validate the quality of data and links in RDF stores.

In this paper we argued on the problem of evaluating sameAs statements defined in RDF data. We designed a general logical evaluation method which relies on the descriptions associated to the resources involved in the sameAs statement to validate. Given a sameAs statement  $sameAs(x, y)$  in a RDF graph  $G$ , our method analyzes the functional properties of  $x$  and  $y$  and the properties defined as local complete. It builds a contextual graph for each resource and assesses the equality of each description involved. We formulated the necessary definitions and formally presented the approach, indicating the set of rules we use in the problem resolution. A possible re-qualification of a 'wrong' sameAs statement has been also introduced in cases in which two resources represent the same conceptual element but at different levels of details. We produced a prototype and experimented it with three datasets of sameAs statements produced by three different linking tools. The analysis of the results proved that, when our approach is applied after the linking, the precision is improved.

In the future, we are planning to explore different research directions. First of all we are planning to test the method on datasets more complex, with local completeness declared for some predicates, or eventually deducible in some way. Secondly, we want to complete the formalization of the re-qualification for wrong sameAs, studying the implication at both conceptual and instance level. We are also working on extending the approach in order to compute similarity measures on property values, thus relaxing specific conditions and allowing the method to work, for example, with data that can have typos errors.

As ultimate goal, we are aiming at designing an integration framework in which a specific knowledge base can be studied, assessed, enhanced and visualized, thanks also to suggestions coming from inference on the data and the links, including data fusion, identity links corrections, and organization of the knowledge and the data at different levels of abstraction.

## References

1. AIMA. The AIMA java library, June 2014.

2. Tim Berners-Lee, James Hendler, and Ora Lassila. The semantic web. *Scientific American*, 284(5):34–43, May 2001.
3. Christian Bizer, Tom Heath, and Tim Berners-Lee. Linked data - the story so far. *International Journal Semantic Web Information Systems*, 5(3):1–22, 2009.
4. OEAI Campaign. Im@oaei2010 - persons-restaurants (pr) dataset, April 2014.
5. Gerard de Melo. Not quite the same: Identity constraints for the Web of Linked Data. In Marie desJardins and Michael L. Littman, editors, *Proceedings of the Twenty-Seventh AAAI Conference on Artificial Intelligence*, Menlo Park, CA, USA, 2013. AAAI Press.
6. L. Ding, J. Shnavier, Z. Shangguan, and D.L. McGuinness. Sameas networks and beyond: Analyzing deployment status and implications of owl: sameas in linked data. In P.F. Patel-Schneider et al., editor, *Intern. Semantic Web Conference*, volume 6496 of *Lecture Notes in Computer Science*, pages 145–160. Springer, 2010.
7. Li Ding, Joshua Shnavier, Tim Finin, and Deborah L. McGuinness. owl:sameAs and Linked Data: An Empirical Study . In *Proceedings of the Second Web Science Conference*, Raleigh NC, USA, April 2010.
8. Alfio Ferrara, Andriy Nikolov, and François Scharffe. Data linking for the semantic web. *Int. J. Semantic Web Inf. Syst.*, 7(3):46–76, 2011.
9. Alfio Ferrara, Andriy Nikolov, and François Scharffe. Data linking. *Journal Web Semantics*, 23:1, 2013.
10. Hugh Glaser, Afraz Jaffri, and Ian Millard. Managing co-reference on the semantic web. In *WWW2009 Workshop: Linked Data on the Web (LDOW2009)*, April 2009.
11. C. Guret, P. Groth, C. Stadler, and J. Lehmann. Assessing linked data mappings using network measures. In E. et al. Simperl, editor, *The Semantic Web: Research and Applications*, volume 7295 of *Lecture Notes in Computer Science*, pages 87–102. Springer Berlin Heidelberg, 2012.
12. H. Halpin, P.J. Hayes, and H.S. Thompson. When owl: sameas isn’t the same redux: A preliminary theory of identity and inference on the semantic web. In *Workshop on Discovering Meaning On the Go in Large Heterogeneous Data 2011 (LHD)*, pages 25–30, 2011.
13. Harry Halpin and Patrick J. Hayes. When owl: sameas isn’t the same: An analysis of identity links on the semantic web. In C. Bizer, T. Heath, T. Berners-Lee, and M. Hausenblas, editors, *Proceedings of the WWW2010 Workshop on Linked Data on the Web*, volume 628 of *CEUR Workshop Proceedings*. CEUR-WS.org, 2010.
14. Tom Heath and Christian Bizer. *Linked Data: Evolving the Web into a Global Data Space*. Morgan & Claypool., 1st edition edition, 2011.
15. Jeff Heflin and Hector Muoz-avila. Lcw-based agent planning for the semantic web. In *In Ontologies and the Semantic Web. Papers from the 2002 AAAI Workshop WS-02-11*, pages 63–70. AAAI Press, 2002.
16. Pascal Hitzler, Markus Krötzsch, Bijan Parsia, Peter F. Patel-Schneider, and Sebastian Rudolph, editors. *OWL 2 Web Ontology Language: Primer*. W3C Recommendation, 27 October 2009. Available at <http://www.w3.org/TR/owl2-primer/>.
17. Ora Lassila, Ralph R. Swick, World Wide, and Web Consortium. Resource description framework (rdf) model and syntax specification, 2004.
18. J. Murdock M. Yatskevich, C. Welty. Coreference resolution on rdf graphs generated from information extraction: first results. In *Proceedings of Web Content Mining with Human Language Technologies Workshop*, 2006.
19. A. Nikolov, V.S. Uren, E. Motta, and A.N. De Roeck. Handling instance coreferencing in the knofuss architecture. In P. Bouquet, H. Halpin, H. Stoermer, and G. Tummarello, editors, *Identity and Reference on the Semantic Web International Workshop*, volume 422 of *CEUR Workshop Proceedings*. CEUR-WS.org, 2008.

20. W3C OWL Working Group. *OWL 2 Web Ontology Language: Document Overview*. W3C Recommendation, 27 October 2009. Available at <http://www.w3.org/TR/owl2-overview/>.
21. N. Pernelle, F. Saïs, B. Safar, M. Koutraki, and T. Ghosh. N2R-Part: Identity Link Discovery using Partially Aligned Ontologies. In *International Workshop on Open Data*, Paris, France, June 2013.
22. G. Robinson and L. Wos. Paramodulation and theorem-proving in first-order theories with equality. In J. Siekmann and G. Wrightson, editors, *Automation of Reasoning 2: Classical Papers on Computational Logic 1967-1970*, pages 298–313. Springer, Berlin, Heidelberg, 1969.
23. F. Saïs, N. Pernelle, and M.C. Rousset. Combining a logical and a numerical method for data reconciliation. *Journal of Data Semantics*, 12:66–94, 2009.
24. Fatiha Saïs, Nobal B. Niraula, Nathalie Pernelle, and Marie-Christine Rousset. Ln2r a knowledge based reference reconciliation system: Oaei 2010 results. In Pavel Shvaiko, Jérôme Euzenat, Fausto Giunchiglia, Heiner Stuckenschmidt, Ming Mao, and Isabel F. Cruz, editors, *OM*, volume 689 of *CEUR Workshop Proceedings*. CEUR-WS.org, 2010.
25. Dezhao Song and Jeff Heflin. Automatically generating data linkages using a domain-independent candidate selection approach. In Lora Aroyo, Chris Welty, Harith Alani, Jamie Taylor, Abraham Bernstein, Lalana Kagal, Natasha Fridman Noy, and Eva Blomqvist, editors, *International Semantic Web Conference (1)*, volume 7031 of *Lecture Notes in Computer Science*, pages 649–664. Springer, 2011.
26. Dezhao Song and Jeff Heflin. Domain-independent entity coreference for linking ontology instances. *J. Data and Information Quality*, 4(2):7:1–7:29, March 2013.
27. Danai Symeonidou, Vincent Armant, Nathalie Pernelle, and Fatiha Saïs. SAKey: Scalable Almost Key discovery in RDF data. In *Proceedings of the 13th International Semantic Web Conference (ISWC2014)*, ISWC '14. Springer Verlag, 2014.
28. Aynaz Taheri and Mehrnoush Shamsfard. Instance coreference resolution in multi-ontology linked data resources. In Hideaki Takeda, Yuzhong Qu, Riichiro Mizoguchi, and Yoshinobu Kitamura, editors, *JIST*, volume 7774 of *Lecture Notes in Computer Science*, pages 129–145. Springer, 2012.
29. World Wide Web Consortium (W3C). W3c semantic web activity. <http://www.w3.org/2001/sw/>. Accessed: 2014-01-30.
30. Gerd Wagner. Web rules need two kinds of negation. In François Bry, Nicola Henze, and Jan Mauszyski, editors, *Principles and Practice of Semantic Web Reasoning*, volume 2901 of *Lecture Notes in Computer Science*, pages 33–50. Springer Berlin Heidelberg, 2003.
31. JM R. Yves, E.P. Shironoshita, and M.R. Kabuka. Ontology matching with semantic verification. *Web Semantics*, 7(3):235–251, September 2009.