

# How Redundant Is It? - An Empirical Analysis on Linked Datasets

Honghan Wu<sup>1</sup>, Boris Villazon-Terrazas<sup>2</sup>, Jeff Z. Pan<sup>1</sup>, and Jose Manuel Gomez-Perez<sup>2</sup>

<sup>1</sup> Department of Computing Science, University of Aberdeen, UK

<sup>2</sup> iSOCO, Intelligent Software Components S.A., Spain

**Abstract.** While there are some popular vocabularies widely used across linked open data, many linked data sets do not have T-Box axioms at all. How does such fact, i.e., the usage patterns on T-Boxes, affect the linked data consumption? This might be an interesting question to be asked by linked data consumers. In this paper, we analysis one particular aspect of this question i.e., redundancy analysis, on several popular datasets and vocabularies in web of data. We start with the analysis on semantic redundancies of datasets given their T-Boxes. Then, we propose useful t-box axioms which are helpful for consumption but are absent in the dataset in question. Finally, we reveal how the linkages of the web of data, both in concept-level and instance-level, affects data set redundancies.

## 1 Introduction

### 1.1 What is data redundancy with linked data?

### 1.2 Why is it of special interest to linked data consumption?

The Good

The Bad

The Ugly

## 2 Linked Data Redundancy Categorisation

Before analysing the redundancy in Linked Data datasets, it makes sense to have a good understanding about the redundancy of Linked Data. In this section, we briefly discuss the categorisation of Linked Data Redundancy. Given the fact that most Linked Datasets are represented in RDF data model, in the first part of this section, we categorise the redundancies in RDF data and point out the main focuses of this paper. In addition to RDF representation, the other important characteristic of Linked Data is its linked aspect. What are the different aspects of the linked nature which are relevant to the redundancy of the data and how can they be effecting the redundancies? In the second subsection we try to answer these questions.

## 2.1 Redundancy in RDF Data

The representation of RDF datasets can be broken down into two levels. In the data model level, an RDF dataset is essentially a set of triples. To share or consume an RDF dataset, e.g. for storage, transmission or query-answering, it needs to be represented in the second level, i.e. the serialisation level, where it has to be serialised as a sequence of bits. In this level, an RDF dataset usually takes the form of textual or binary files by using predefined syntaxes, e.g. RDF/XML, N-Triples or even sophisticated compression format (e.g. HDT [2]). Accordingly, the redundancies of RDF data reside in both levels of RDF data representation, i.e. data model level and serialisation level.

In the data model level, the size of data can be calculated by the number of triples. Hence, in this level, the data redundancy exists if less triples can be used to represent the same semantic meanings of the original data. In the serialisation level, the data is represented as a sequence of bits. Given a fix set of triples, one serialisation is said to be more redundant than the other if it uses more bits than its counterpart.

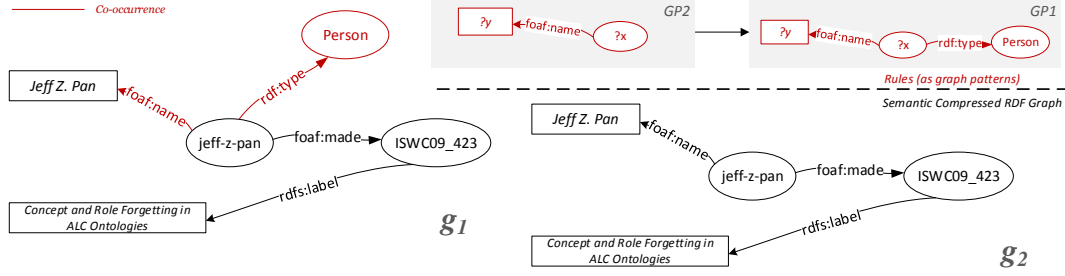
In [5], we proposed a fine-grained categorisation of RDF data redundancies. Table 1 illustrates such categorisation of RDF redundancies and also puts it in the dimension of RDF representation levels. In this paper we mainly focus on the semantic redundancy and the second type of syntactic redundancy, i.e. the inter-structural redundancy. Both are highlighted with a grey background in Table 1. Semantic redundancy is selected because it can be generated or removed by the T-Box axioms of a dataset. Hence, it is of interest to most of the Linked Data consumption tasks like inference computation and ontology based data access. Syntactic redundancy is also important to data consumption because a concise serialisation is beneficial not only to data transmission but also to query answering tasks [2]. In this category, a recent study [5] points out that most existing compression techniques, e.g. [2], cannot deal with the inter-structural one [5]. Hence, it is particularly interesting to analyse inter-structural redundancies in Linked Open Data.

**Table 1.** RDF Data Redundancy Categorisation

Types	Semantic Redundancy	Syntactic Redundancy		Symbolic Redundancy
		Intra-structural	Inter-structural	
Data Model Level	✓	-	-	-
Serialisation Level	-	✓	✓	✓

**Semantic Redundancy** An RDF dataset is said to be semantically redundant if some triples can be removed without leading to any changes in its meaning. In most cases the removal of these triples requires additional rules to be added in the dataset so that it is possible to re-generate the removed triples when needed. The usual form of such rules is the T-Box, i.e., the concept level statements in an Ontology. For example, in Fig. 1, we have an RDF graph of  $g_1$ . In the FOAF ontology <sup>3</sup>, there is a rule

<sup>3</sup> <http://xmlns.com/foaf/spec/>



**Fig. 1.** Original v.s. Semantically Compressed

of `<foaf:name, rdfs:domain, foaf:Person>`. Based on the `rdfs2` rule in the RDF specification, the type assertion in  $g_1$ , i.e. `<jeff-z-pan, rdf:type, foaf:Person>`, can be removed, given the presence of `<jeff-z-pan, foaf:name, Jeff Z. Pan>`.

In a more general perspective, the semantic redundancy can be identified by co-occurrences of triple patterns (cf. the red part of  $g_1$  in Fig. 1). Hence, it is not necessary to restrain the ability of semantic redundancy identification by the limitation of T-Box rules in representing triple pattern co-occurrences. In this regard, Joshi et al. [4] applied association rule mining approach to identify the co-occurrences in terms of association rules. We apply a graph pattern based rule system to represent the semantic redundancies. For example, in Fig. 1,  $g_2$  has less triples than  $g_1$  but the two are semantically equivalent. The redundancy in  $g_1$  is represented by the graph pattern rule in the upper part of  $g_2$ . Obviously, graph pattern based rules can represent more complex triple co-occurrences than T-Box rules.

**Inter-structure Syntactic Redundancy** As shown in Table 1, the other two types of redundancies, i.e. syntactic and symbolic ones, both reside in the serialisation level. Their volume can be evaluated using the number of bits used in the serialisation. To separate the two, we can use a simple formula  $|F| = n \times r$ , where  $F$  is the serialisation file,  $n$  is the number of resource occurrences and  $r$  is the average bits needed to represent a resource. The syntactic redundancies reside in the component of  $n$ , while the symbolic ones are in  $r$ .

Most existing serialisation approaches apply syntaxes to reduce  $n$ . The idea is to group triples by subjects or objects so that the multiple occurrences of the same resource only need to be serialised once. For example, the RDF/XML serialisation standard provides abbreviation and striping syntaxes. However, such syntaxes only work on concrete graph structures. Similar graph structures (i.e. graph patterns) which repeatedly occur in the data are not taken into account. In the following example, the graph pattern  $GP$  has two instances of  $Inst_1$  and  $Inst_2$ . The structure of  $GP$  appears twice in both instances. This means that each of the two predicate resources of  $GP$ , i.e. `foaf:name` and `foaf:made`, has two occurrences in its instances, which is avoidable when  $GP$  structure (stored wherever) is referred instead of duplicated in both instances.

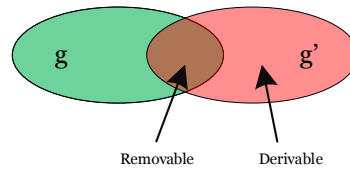
$GP : \langle ?U, foaf:name, ?Y \rangle, \langle ?U, foaf:made, ?Z \rangle$   
 $Inst_1 : \langle jeff-z-pan, foaf:name, Jeff Z. Pan \rangle, \langle ?U, foaf:made, ISWC09_423 \rangle$   
 $Inst_2 : \langle jose, foaf:name, Jose Manuel Gomez Perez \rangle, \langle ?U, foaf:made, ISWC13_1xx \rangle$

We use the term of intra-structural redundancy to denote the unnecessary resource occurrences in concrete graph structure, while the term of inter-structural redundancy is used for the unnecessary resource occurrences of graph patterns in its instances. Most of existing serialisation approaches do not provide facilities to identify graph patterns. Hence, they leave the inter-structural redundancy untouched.

## 2.2 Redundancy in Linked Data

One of the most important characteristics of Linked Data is its linking capability, which is to create connections from one information source to the other. The connection can be created in the concept level, where individuals of one dataset are described using concepts from another dataset or vocabulary. We denote this type of connections as T-Box Reuse. The second type of connections is the linkage in individual level, where individuals in one dataset are specified to be the same as their counterparts in the other, e.g. by using *owl:sameAs* assertions. This type of connections is called as A-Box Linkage in this paper. Both types of connections have implications in changing the semantics of the original dataset. These implications might change the redundancies of the dataset in question. In rest part of this subsection, we briefly discuss the possible changes on data redundancies caused by both types of connections.

**T-Box Reuse** When an individual is described using concepts defined in another T-Box. The axioms in that T-Box will be applied to infer more assertions not only on this individual but also potentially on other individuals which are related to it. For example, in LinkedMDB <sup>4</sup>, all individuals of *movie:actor* are also asserted to be instances of *foaf:Person*. Given the axiom of  $\langle foaf:Person, rdfs:subClassOf, foaf:Agent \rangle$  in FOAF vocabulary, all actors in LinkedMDB are also individuals of *foaf:Agent*. Essentially T-Box Reuse will bring new rules to the original dataset. These rules are usually from a relevant subset of the axioms of the materialised version of the reused T-Box.



**Fig. 2.** Data Redundancies affected by T-Box Reuse

<sup>4</sup> Linked Movie Database <http://data.linkedmdb.org/>

These new rules can infer a set of new triples <sup>5</sup> to be added in the reusing RDF dataset. Fig. 2 illustrates a typical case where the new triples overlaps with the existing ones. In Fig. 2, the original RDF dataset, labelled as  $g$ , is denoted by the green oval, and the new set of inferred triples is denoted by the pink oval with a label of  $g'$ . Depending on whether they are in the overlap with existing triples, new triples can be divided into two parts. The two parts will lead to two different consequences to the redundancies of the original data. The overlap part, labelled *removable* in Fig. 2, contains those triples in  $g$ , which can be inferred by new rules. This means that they are turned to be semantically redundant by the reused T-Box. As a consequence, the dataset turns to have more redundancies. On the contrary, the other part ( $g' \setminus \text{removable}$ ), labelled as *derivable* in Fig. 2, will decrease the data redundancy by means of increasing data compression ratio <sup>6</sup>, it can be calculated as  $\text{compression ratio} = \frac{|g| + |\text{derivable}|}{|g|}$ , in which the bigger  $|\text{derivable}|$  component becomes, the larger the compression ratio will be. In a word, *removable* and *derivable* affect the data redundancies in two opposite ways. Analyses on them will reveal the effects on data redundancy of concept-level connections between Linked Data.

**A-Box Linkage** Individual level linkages by *owl:sameAs* assertions are among the most advocated best practices in the Linked Data community. Data publishers are encouraged to provide such linkages to break the information silos. However, its semantic implications will lead to substantial inference computations across different datasets, some of which might be unexpectedly expensive. In [3], Halpin et al. pointed out that *quick-and-dirty use of owl:sameAs will almost always lead to OWL Full*. The other issue is that transitive closures of *owl:sameAs* in the Web scale can easily get too large to be manageable.

Despite the practical issues, we briefly discuss how the semantic implications of *owl:sameAs* links might lead to the changes in data redundancies. As supporting efficient query answering is the basis of many Linked Data consumption tasks, our discussion in this paper is based on OWL2 QL profile [1]. When an individual  $i$  in dataset  $D$  is asserted to be *sameAs* another individual  $i_1$  in  $D_{i_1}$ , all materialised assertions of  $i_1$  in  $D_{i_1}$ , denoted as  $A(i_1)$ , are immediately true to  $D$ . Hence, the data of  $D$  will be  $g_D \cup A(i_1)$ , where  $g_D$  is the original RDF graph of  $D$ . Considering OWL2 QL semantics, it is reasonable that only type assertions of  $i_1$  in  $D_{i_1}$  are included in  $A(i_1)$  directly. Regarding other triples of  $i_1$  in  $D_{i_1}$ , they will be simplified as existential assertions to be added.

Similar to T-Box Reuse case, the adding of new data has two consequences to the data redundancy of  $D$ . The first one is some triples turned to be redundant, while they were not originally. We also use the term *removable* to denote these triples. Let  $M$  be an A-Box materialisation function. The *removable* can be calculated by  $\text{removable} = g_D \cap (M(g_D \cup A(i_1)) \setminus M(g_D))$ . The second consequence is that  $g_D$ 's data compression ratio can be increased, which means less redundancy. This is caused by new triples

<sup>5</sup> Some triples inferred by new rules might be inferred by the dataset's own T-Box as well. To make discussion easier, we use the term *new triples* to denote those triples which can NOT be inferred by the dataset's own T-Box but the new rules.

<sup>6</sup> [http://en.wikipedia.org/wiki/Data\\_compression\\_ratio](http://en.wikipedia.org/wiki/Data_compression_ratio)

which are not in the original data and turn to be derivable after adding  $A(i_1)$ , while they could not be derived before that. The term *derivable* is used to denote such triples, which can be calculated by  $derivable = (M(g_D \cup A(i_1)) \setminus M(g_D)) \setminus g_D$ .

### 3 Two Dimension Analysis

According to above discussions, the redundancies in Linked Data can be analysed from two dimensions (cf. Fig. 3). The first dimension is that of the RDF data redundancy, where the focus is to reveal different categories of redundancies from data model level to serialisation level. The second dimension is from the linked semantic point of view, where the focus is to analyse the data redundancy based on different types of semantics. The *A-Box* semantics is to analyse the data redundancy only from data level without any T-Box axioms. The other three types are considering both data and T-Box information. The *No Linkage* semantics is to do the analysis based on the dataset's main T-Box. The main T-Box is the one which is defined and published by the data provider. It is not necessary that every dataset has a main T-Box. In such cases, this analysis is not applicable. The *T-Box Reuse* semantics is the type of analysis focusing on redundancy changes caused by concept level connections, i.e. reusing T-Box, while the *A-Box Linkage* is to reveal the changes of data redundancies caused by the individual level connections.

		RDF Redundancy Dimension		
		Semantic	Syntactic	Symbolic
Linked Semantic Dimension	A-Box	✓	✓	
	No Linkage	✓	-	-
	T-Box Reuse	✓	-	-
	A-Box Linkage		-	-

**Fig. 3.** Two Dimension Analysis on Linked Data Redundancy

In the matrix of Fig. 3, there are total 6 valid types of analyses. In this paper, we focus on four of them, which are marked with ticks in the figure. For *A-Box* semantics, we propose a graph pattern based approach to reveal the semantic and syntactic (inter-structure only) redundancies. Based on the graph pattern approach, we propose a virtual materialisation approach to analyse data redundancies in *No Linkage* and *T-Box Reuse* semantics. The other two types of analyses are left for future work.

#### 3.1 Graph Pattern Based Analysis Method

As discussed in section 2.1, we focus on semantic redundancy and the inter-structural syntactic redundancy, both of which will be benefited from the ability to identify frequent graph patterns. For semantic redundancy, identified graph patterns can be used to

identify possible rules for removing redundant triples. Combined with the knowledge of instance numbers of these graph patterns, these rules can be used to calculate the volume of semantic redundancy as number of removable triples. Similarly, for inter-structural redundancy, the structure of graph patterns and their instances numbers can give us a way to calculate the volume of syntactic redundancy in the data. In this subsection, we describe an entity description pattern based approach for redundancy analysis.

In an RDF graph, we call its non-literal nodes as entities. For an entity  $e$  in an RDF graph  $G$ , we can get a data block for it by extracting triples in  $G$  each of which has  $e$  as its subject or object. We call such kind of data blocks as Entity Description Blocks (EDBs for short).

For an EDB, it can be summarised by a notion of entity description pattern which is defined in Definition 1. EDP, the short name for entity description pattern, is the building block of our analysis approach.

**Definition 1.** (*Entity Description Pattern*) Given an entity description block  $B_e$ , its description pattern is a tuple  $P_e = (C_e, A_e, R_e, V_e)$ , where

- $C_e = \{c_i | \langle e, rdf:type, c_i \rangle \in G\}$  is called as the class component;
- $A_e = \{p_i | \langle e, p_i, l_i \rangle \in G \text{ and } l_i \text{ is a literal}\}$  is called as the attribute component;
- $R_e = \{r_i | \langle e, r_i, o_i \rangle \in G \text{ and } o_i \text{ is a URI resource or blank node}\}$  is called as the relation component;
- $V_e = \{v_i | \langle s_i, v_i, e \rangle \in G\}$  is called as the inverse relation component.

Taking the RDF graph  $g_1$  in Fig. 1 for example, the EDP of entity *jeff-z-pan* is  $(\{foaf:Person\}, \{foaf:name\}, \{foaf:made\}, \emptyset)$ , and the one of *ISWC09\_423* is  $(\emptyset, \{rdfs:label\}, \emptyset, \{foaf:made\})$ .

By generating EDPs for all entities in an RDF graph  $G$ , we can get a EDP representation  $EDP_G$  of  $G$ , which is a set of EDPs. As we will show in later section, the number of EDPs in an RDF dataset is usually much less than the number of entities. This is because many entities are sharing same EDP. The more entities are sharing one EDP; the more times its structure is duplicated. Such duplications are the source of data redundancies, which we break down to semantic redundancy and syntactic redundancy.

**Semantic Redundancy Identified By EDP** In the definition of EDP, the class component  $C_e$  is a set of constant class names. Hence, it is straightforward to generate a graph pattern based substitution rule for removing these type assertions from the original data. In particular, the rule can be defined as  $(\emptyset, A_e, R_e, V_2) \rightarrow (C_e, A_e, R_e, V_e)$ . The number of triples can be removed is  $|C_e| \times f_{P_e}$ , where  $f_{P_e}$  is the number of instances of EDP  $P_e$ . For the example of *jeff-z-pan* in Fig. 1, the rule will be  $(\emptyset, \{foaf:name\}, \{foaf:made\}, \emptyset) \rightarrow (\{foaf:Person\}, \{foaf:name\}, \{foaf:made\}, \emptyset)$ . In this particular case, this rule is equivalent to the rule of  $\langle foaf:name, rdfs:domain, foaf:Person \rangle$  from *FOAF* vocabulary. Although EDP based rules are more expressive than T-Box axioms, it is interesting to know the overlap of identifiable semantic redundancies between the two. In next section, we will provide results in this regard.

**Inter-structural Syntactic Redundancy Identified By EDP** The inter-structural redundancy denotes the unnecessary structure recurrences in EDBs of graph pattern instances. Given an EDP  $P_e$ , it is straightforward to calculate its recurrences as  $(|A_e| + |R_e| + |V_e|) \times f_{P_e}$ , the unit of which is resource occurrence.

**Virtual A-Box Materialisation on EDP** When considering T-Box of an RDF dataset, the data redundancy might be changed due to the above-mentioned *removable* and *derivable* triple sets. To compute the two triple sets, inferences need to be performed on A-Boxes, which might be too expensive for large scale data analyses. For redundancy analysis purpose, we only need to know the size of *removable* and *derivable* sets instead of getting the exact triples. In this subsection, we propose a virtual materialisation approach based on EDP for computing the triple sizes for the two triple sets.

Given an EDP  $P_e$ , according OWL2 QL profile, the four components of  $P_e$  are sufficient for inference computation on EDP. After applying inference rules defined in [1], we will get a new EDP:  $P_e^M = (C_e^M, A_e^M, R_e^M, V_e^M)$ . The number of *derivable* triples can be calculated as follows.

$$|derivable| = f_{P_e} \times \sum_{c \in N} f(c), \quad (1)$$

where  $N = (C_e^M \setminus C_e) \cup (A_e^M \setminus A_e) \cup (R_e^M \setminus R_e) \cup (V_e^M \setminus V_e)$  and  $f$  is an auxiliary dictionary of  $P_e$  which stores the instantiated times of each concept in  $P_e$ .

The number of *removable* triples can be calculated as follows.

$$|removable| = f_{P_e} \times \sum_{c \in E} f(c), \quad (2)$$

where  $E = (C_e^M \cap C_e) \cup (A_e^M \cap A_e) \cup (R_e^M \cap R_e) \cup (V_e^M \cap V_e)$ .

### 3.2 The Metrics

## 4 Related Work

The analysis papers; The compression papers; The summary paper

## 5 Linked Dataset Analysis

## 6 Conclusion

## References

1. D. Calvanese, G. De Giacomo, D. Lembo, M. Lenzerini, and R. Rosati. Tractable reasoning and efficient query answering in description logics: The dl-lite family. Journal of Automated reasoning, 39(3):385–429, 2007.



2. J. D. Fernández, M. A. Martínez-Prieto, C. Gutiérrez, A. Polleres, and M. Arias. Binary rdf representation for publication and exchange (HDT). Web Semantics: Science, Services and Agents on the World Wide Web, 19:22–41, 2013.
3. H. Halpin, P. J. Hayes, J. P. McCusker, D. L. McGuinness, and H. S. Thompson. When owl: sameas isn't the same: An analysis of identity in linked data. In The Semantic Web–ISWC 2010, pages 305–320. Springer, 2010.
4. A. K. Joshi, P. Hitzler, and G. Dong. Logical linked data compression. In The Semantic Web: Semantics and Big Data, pages 170–184. Springer, 2013.
5. J. Z. Pan, J. M. G. Pérez, Y. Ren, H. Wu, and M. Zhu. SSP: Compressing RDF data by summarisation, serialisation and predictive encoding. Manuscript submitted for CIKM, 2014.