# Data Accuracy: What does it mean to LOD?

Behshid Behkamal
Computer Engineering Department
Ferdowsi university of Mashhad
Mashhad, Iran
behkamal@um.ac.ir

Mohsen Kahani
Computer Engineering Department
Ferdowsi university of Mashhad
Mashhad, Iran
kahani@um.ac.ir

Ebrahim Bagheri
Department of Electrical and Computer Engineering
Ryerson University
Toronto, Canada
bagheri@ryesron.ca

Majid Sazvar
Computer Engineering Department
Ferdowsi university of Mashhad
Mashhad, Iran
sazvar@alumni.um.ac.ir

*Abstract* - **Linked Open Data provides a distributed model for the semantic web to create knowledge by publishing public available data and meaningfully interlinking dispersed data sources. It is undeniable that the realization of this goal depends strongly on the quality of the published data. Since, data quality is a multi-dimensional concept which is defined by a number of quality factors, in order to study data quality in depth; it is necessary to study each quality factor separately as well as the properties of its environment. The main objective of this work is to propose a set of metrics that enable the assessment of the accuracy of data sets from both semantic and syntactic accuracy viewpoints.**

*General Terms - Measurement, Experimentation*

*Keywords - Accuracy, Metrics, Quality Assessment, LOD*

## I. INTRODUCTION

Linked Open Data (LOD) allows for any data provider to publish its publicly available data and meaningfully link them with other information sources over the Web. The main goal of the LOD initiative is to create knowledge by interlinking dispersed data. It is undeniable that realization of this goal depends strongly on the quality of the published data.

In recent years, researchers have already made several proposals to evaluate data quality using different assessment methods such as user experience; expert judgment, sampling, parsing, continuous assessment and cleansing techniques[1, 2], but those cannot be applied directly to the Web of Data. Although data quality is an important issue for the successful organic growth of LOD, there are only a very limited number of research initiatives that focus on data quality specifically for LOD. Based on our practical experience in publishing linked data[3] we have observed that many of the published datasets suffer from quality issues such as syntax errors, redundant instances, and incorrect/uncompleted attribute values. We believe that the assessment of data quality before publishing to the LOD cloud can help publishers filter out low-quality data based on the quality assessment results. So, data owners/providers can evaluate their data before publishing as linked data as

well as data consumers can make better and more informed decisions when deciding which data to use. The high number of quality factors and their inter-relationship makes quality evaluation to be a complex problem and it is difficult to consider all quality factors at once. In order to study data quality in depth, it is necessary to study each quality factor separately as well as the properties of the environment that affect it. Since, most of the quality problems of dataset, studied in our experiment [3]were caused by the inaccuracy of their sources; this work focuses on the data accuracy before they are published.

The rest of this paper is organized as follows: first, data quality research in the area of LOD is reviewed in Section 2. Then, Section 3 defines data accuracy and discusses quality sub-factors that characterize our notion of data accuracy in the context of LOD. The process of metrics definition is presented in Section 4. The results of our experiments and some discussions are then provided in Sections 5. Finally, the paper is concluded by presenting future works in Section 6.

## II. RELATED WORKS

Despite its importance, data quality has not yet received a lot of attention by the researchers in the area of LOD. In [4] a framework is proposed to assess the information quality of Web data sources based on the provenance information. Also, Bizer describes a framework, called WIQA, to filter poor information in Web-based information systems according to user defined quality requirements [5]. Other approaches have used Semantic Web technologies to identify and correct data quality issues [6, 7].

Pedantic Web Group classifies quality problems of the published linked datasets and discusses common errors in RDF publishing, their consequences for applications, along with possible publisher-oriented approaches to improve the quality of machine-readable and open data on the Web[8]. In other work, Furber proposes an approach to evaluate the quality of datasets using SPARQL queries to identify some quality problems of already available datasets e.g. Geonames [9].

Generally, all of these related works focus on data quality problems in the published datasets and neither of them concentrate on a special quality dimension, nor propose a method for systematically evaluating data quality. In this paper, we focus on the assessment of data accuracy in the context of LOD by proposing a set of metrics for both aspects of semantic and syntactic accuracy.

## III. OUR NOTION OF DATA ACCURACY

To evaluate the quality of any dataset, it is imperative to define the quality dimensions based on the domain of use. Here, by touching upon the definitions of data accuracy in the literature, we define data accuracy in the context of LOD by proposing two quality sub-factors, namely semantic accuracy and syntactic accuracy.

Most data quality research includes accuracy as a key dimension of data quality for different domains. In the past, accuracy was known as "data quality". For that reason, several accuracy definitions include other quality aspects such as completeness or freshness, i.e. incomplete and expired data have been considered inaccurate [10, 11]. For example, ISO-25012 defines accuracy as correctly representation of a concept or event in a specific context of use[12]. In [13] it is characterized as the percentage of objects without data errors such as misspellings and out-of-range values, while in [14], accuracy is described as the degree of agreement between a collection of data values and a source agreed to be correct. All of the definitions are concerning different concepts and metrics, which are mainly due to the different objectives of the systems that they are used. For Better understanding of accuracy in the context of LOD and covering related quality metrics, we define accuracy from two perspectives of semantic accuracy and syntactic accuracy based on the classification presented in our previous work [15].

### A. Semantic Accuracy

Semantic accuracy mainly relates to the correctness of a data values in comparison to the actual real world values. In [10], this aspect of accuracy is described as "semantic correctness factor" and concerns the degree of correctness and validity of the data in comparison to the real world or with the reference data agreed to be correct. In [14],semantic accuracy is defined as the closeness of the data values to a set of values defined in a domain considered semantically correct. In the context of LOD, it means that every entity described in a dataset should represent a real world situation. Therefore, resources referencing a wrong real world correspondent and entities with erroneous attribute values are examples of quality deficiencies related to this sub-factor.

### B. Syntactic Accuracy

Syntactic accuracy commonly expresses the degree to which a set of data is free of syntactic errors such as misspellings. According to [10], data is syntactically correct, if it satisfies syntactic rules and constraints imposed by users. In the context of LOD, there are some tools for checking syntax validity of RDF documents, each with its own error-checking functionalities. Some that are available online accept an RDF/XML document as input and check if the document is syntactically valid, such as W3C Markup[1], W3C RDF/XML[2]. Other kind of online validators check the dereferencability of a given URI and determine whether the given URI is an information resource or a non-information resource, such as URIDebugger[3]and Vapour[4]; and some are command line tools designed for larger jobs; such as Eyeball[5] and VRP[6].

In this paper, we define syntactic accuracy as the validity of RDF documents and will propose a set of automated metrics to measure aspects of a given dataset that cannot be checked by the mentioned syntax validators.

## IV. ACCURACY ASSESSMENT METRICS

In order to make accuracy quantifiable, we define a set of metrics to measure both aspects of accuracy. Our approach for data accuracy assessment involves the measurement of quality aspects which cannot be assessed either by available validators, or by experts. The employed approach for metric definition is Goad-Question-Metric (GQM) [16]. In GQM, the goals are gradually refined into several questions and each question is then refined into metrics. Also, one metric can be used to answer multiple questions. Although the GQM was initially proposed in the software engineering field, it has since been widely applied in a variety of other domains as well [17].

We define the primary goal of our metrics as "the assessment of the accuracy of a dataset from the users' point of view in the context of LOD". As mentioned, we consider two sub-characteristics for the accuracy. Therefore, the main goal is decomposed into two sub-goals corresponding to semantic accuracy and syntactic accuracy. Based on this classification, we address these sub-goals by developing appropriate questions, which in turn substantiate the definition of related metrics. A data quality metric is a procedure for measuring an information quality characteristic [5]. Considering the fact that only few studies have been conducted which define quality metrics for LOD [2, 18-20], we undertake an exploratory analysis of the previous and current researches on data accuracy in the database community [2, 10, 11, 21].

### A. Semantic Accuracy Metrics

Semantic Accuracy relates to the correctness of a data value in comparison to its actual real world value. ISO 25012 defines semantic accuracy as the 'closeness of the data values to a set of values defined in a domain considered semantically correct'[12]; while [10] and [13] characterize semantic accuracy as the percentage of objects without data errors such as misspellings, out-of-range values, etc. In our

---

[1] http://validator.w3.org

[2] http://www.w3.org/RDF/Validator

[3] http://linkeddata.informatik.hu-berlin.de/uridbg

[4] http://validator.linkeddata.org/vapour

[5] http://jena.sourceforge.net/Eyeball

[6] http://139.91.183.30:9090/RDF/VRP

study, semantic accuracy focuses on the correctness of data presented in a dataset. To assess semantic accuracy, we need to illustrate that all of the attributes used to describe the entities contain correct values. Therefore, two main questions are developed in the context of GQM: 1) Is all the required information for each entity present? 2) Are the entities described with the appropriate/correct values? In the following, we propose six metrics to answer these questions.

- *Missing Properties Values (Miss_Prp_Vlu)*

Miss_Prp_Vlu measures the ratio of the properties defined in the schema, but not presented in a given dataset. It is calculated as:

$$Miss\_Prp\_Vlu = 1 - \frac{No.of\ defined\ properties, but\ not\ presented}{No.\ of\ classes * No.of\ properties} \quad (1)$$

By this metric, we measure the presence of required properties for the instances according to defined properties in the schema. We assume that all of the properties defined for a class, should be presented for all of the instances of that class. Although in specific cases, some of the defined properties for a class may not be applicable for all instances, we assume that if a property is not used for an instance, we consider it as missing property.

- *Average Missing Properties Values (Avg_MPV)*

Avg_MPV measures the average missing properties per instance. It is calculated as:

$$Avg\_MPV = 1 - \frac{Sum\ of\ the\ ratio\ of\ the\ missing\ properties\ per\ instance}{No.\ of\ instances} \quad (2)$$

This metric measures the presence of all properties for each instance, based on the defined properties for corresponding class in the schema. It is similar to the first metric, but *Avg_MPV* measures the ratio of missing properties per instance, while in*Miss_Prp_Vlu,* we measure the ratio of missing properties in the dataset.

- *Misspelled Property Values (Msspl_Prp_Vlu)*

Msspl_Prp_Vlu measures the ratio of the properties of a dataset which contain misspelled values. It is computed as:

$$Msspl\_prp\_Vlu = 1 - \frac{No.of\ triples\ which\ contain\ misspelled\ properties}{No.of\ triples} \quad (3)$$

This metric is defined to measure the misspelling errors of the values of data type properties. To this end, we have used Lucene spell checker [22] in our implementation. This spell checker includes different languages, including English, Danish, Dutch and Spanish.

- *Misspelled classes (Msspl_Cls)*

Msspl_Cls measures the ratio of the classes defined in the schema having misspelling errors in their names. It is computed as:

$$Msspl\_cls = 1 - \frac{No.of\ classes\ with\ misspelled\ names}{No.of\ classes\ defined\ in\ the\ schema} \quad (4)$$

This metric is defined to measure the misspelling errors in the classes' names. As mentioned in 4.1.3, we have used Lucene spell checker for this purpose.

- *Misspelled properties (Msspl_Prp)*

Msspl_Prp measures the ratio of the properties defined in the schema having misspelling errors in their names as:

$$Msspl\_prp = 1 - \frac{No.of\ properties\ with\ misspelled\ names}{No.of\ properties\ defined\ in\ the\ schema} \quad (5)$$

This metric is defined to measure the misspelling errors in the names of properties.

- *Out of range properties (Out_Prp_Vlu)*

Out_Prp_Vlu measures the ratio of the triples of dataset that contain properties with out of range values. It is calculated as:

$$Out\_Prp\_Vlu = 1 - \frac{No.of\ triples\ containing\ out\ of\ range\ properties\ values}{No.of\ triples} \quad (6)$$

Based on this definition, *Out_Prp_Vlu* measures the ratio of triples containing out of range properties, both data type properties and object properties.

*B.    Syntactic Accuracy Metrics*

ISO 25012 defines syntactic accuracy as the closeness of the data values to a set of values defined in a domain considered syntactically correct[12]. In another definition, [10] states that data is argued to be syntactically correct, if it satisfies syntactic rules and constraints imposed by the users. Furthermore, syntactic accuracy can additionally be defined as the structural validity of a dataset, such as compliance with RDF/XML standard. In this study, we focus on the syntactic accuracy of entities as well as the appropriateness of the properties which are used for describing the entities. To this end, the following questions are developed in the framework of GQM: 1) Have the resources been described with appropriate properties? 2) Are there formal definitions in the schema for all of the classes and properties used in the dataset? 3) What is the degree of inconsistency in terms of using classes, properties and data types in the dataset? To answer these questions, a set of metrics are proposed as follows.

- *Improper Data Types (Imp_DT)*

Imp_DT measures the ratio of the triples of a dataset that contain data type properties with inappropriate data types as:

$$Imp\_DT = 1 - \frac{No.of\ triples\ containing\ inappropriate\ data\ type}{No.of\ triples} \quad (7)$$

This metricconcerns the incorrect usage of data types, which is a relatively common error in the Web of Data. In RDF, a subset of well-defined XML data types is used to provide structure and semantics to literal values. For example, date values can be specified using the xsd:date data type, which provides a lexical syntax for date strings and a mapping from date strings to date values interpretable by an application [8].

- *Undefined Classes (Und_Cls)*

*Und_Cls* measures the ratio of the triples of a given dataset that have used classes without any formal definition as:

$$Und\_CLs = 1 - \frac{No.\ of\ triples\ using\ undefined\ classes}{No.of\ triples} \qquad (8)$$

It is defined to detect the classes used in a dataset, but not defined in the schema. In some published datasets, properties and classes are used without any formal definition. The use of ad-hoc undefined classes and properties makes automatic integration of data less effective and foregoes the possibility of making inferences through reasoning [8].

- *Undefined properties (Und_Prp)*

*Und_Prp* measures the ratio of the triples of a given dataset that have used properties without any formal definition, calculated as:

$$Und\_Prp = 1 - \frac{No.\ of\ triples\ using\ undefined\ properties}{No.of\ triples} \qquad (9)$$

This metric is defined to detect the properties used, but not defined in the schema. Thus, all of the properties which are not user-defined are considered as undefined properties.

- *Membership of disjoint classes (Dsj_Cls)*

*Dsj_Cls* measures the ratio of the instances of a dataset being members of disjoint classes. It is calculated as:

$$Und_{CLs} = 1 - \frac{No.\ of\ instances\ being members\ of\ disjoint\ classes}{No.of\ instances} \qquad (10)$$

Based on this formula, it is understood that *Dsj_Cls* is related to the members of disjoint classes either asserted directly by the publisher, or inferred through reasoning. For example, the instances of classes which were defined as complements of each other (using owl:complementOf), or the instances of foaf:Person and foaf:Document classes in FOAF, which are defined as being disjoint.

- *Usage of disjoint properties (Dsj_Prp)*

*Dsj_Prp* measures the ratio of the instances of a dataset that have used disjoint properties. It is calculated as:

$$Und\_Prp = 1 - \frac{No.\ of\ instances\ using\ disjoint\ properties}{No.of\ instances} \qquad (11)$$

The example of *Dsj_Prp* is similar to *Dsj_Cls*, where an instance has used two properties which are defined as being disjoint in the schema.

- *Functional properties with inconsistent values (FP)*

*FP* measures the ratio of triples with functional properties which contain inconsistent values. It is calculated as:

$$FP = 1 - \frac{No.\ of\ triples\ with\ inconsistent\ values\ for\ functional\ properties}{No.of\ triples} \qquad (12)$$

According to this definition, FP counts the triples in which their predicates are a specific functional property with the same subjects, but different objects.

- *Invalid usage of inverse-functional properties (IFP)*

*IFP* measures the ratio of triples that contain invalid usage of inverse-functional properties. Aside from URIs, resources are identified by the values of properties which uniquely identify them, named "inverse-functional property". IFP metric is calculated as:

$$IFP = 1 - \frac{No.\ of\ triples\ with\ inconsistent\ values\ for\ functional\ properties}{No.of\ triples} \qquad (13)$$

The definition of IFP is similar to FP, where IFP counts the triples in which their predicates are the same inverse functional property with the same objects, but different subjects. If two resources share a common value for one of these properties, reasoning will view these resources as equivalent (referring to the same resource). An example of this issue is presented in [8], where the FOAF ontology has defined foaf:mbox for email addresses to identify people, but there are a lot of void values for this property; and as a result all of these people are interpreted as equivalent and represent the same real-world person. The issue can easily be avoided by validating user input and also, it can automatically be resolved by checking the validity of inverse-functional values.

- *Misusage of Properties (Msusg_Prp)*

*Msusg_Prp* measures the ratio of triples which have used data type properties instead of object properties or vice versa.

$$Msusg_{Prp} = 1 - \frac{No.\ of\ triples\ using\ DT\ prp\ instead\ of\ object\ prp\ or\ vice\ versa}{No.of\ triples} \qquad (14)$$

A data-type property describes properties, which relate some resource to a literal value, while an object property describes properties, which relate one resource to another. In some cases of published datasets, data-type properties are used between two resources or conversely, the object properties are used with literal values. This metric is defined to measure these issues.

- *Misplaced Classes and Properties (Misplc_Cls_Prp)*

*Misplc_Cls_Prp* is defined to measure the ratio of triples with misplaced classes or properties.

$$Misplc\_Cls\_Prp = 1 - \frac{No.\ of\ triples\ with\ misplaced\ classes\ or\ properties}{No.of\ triples}$$
$$(15)$$

This metric is related to the usage of classes as properties, or conversely the usage of properties as a class. According to the examples presented in [8], rdfs:range is a core RDFS property, but is sometimes defined in a document as a class. In this section, fifteen metrics have been defined to assess the accuracy of a given dataset from two viewpoints: semantic and syntactic viewpoints. In the next section, these metrics are used in practice and the results of our observations are discussed.

## V. EXPERIMENTS

To show the applicability of the proposed metrics and observe their behaviors over different datasets, it is necessary to place them under empirical evaluation. Here, we report the results of our observations with regards to the calculation of the proposed metrics for several real world datasets. We have selected four datasets from across a variety of LOD domains. We also made sure that these datasets were of different sizes as shown in Table 1. In order to put the proposed metrics into practice, we have implemented a tool that is able to automatically compute the values of the metrics for any given input dataset. The code is implemented in the Java programming language (JDK 7 Update 25 x64) using Jena 2.6.3 semantic web library and is publicly accessible [23].

Table 1.The details of the datasets used in the experiment

| Dataset | Number of triples | Number of instances | Domain |
|---|---|---|---|
| Geonames[7] | 6,590 | 699 | Geography |
| IMDB[8] | 866 | 291 | Movie |
| Anatomy[9] | 6,449 | 6449 | Anatomy |
| Citeseer[10] | 948,770 | 173963 | Publication |

Table 2 presents all of the collected values of the metrics for each of the datasets.

Table 2. Observations for Metrics

| No | Metrics | Geonames | IMDB | Anatomy | Citeseer |
|---|---|---|---|---|---|
| 1 | **Miss_Prp_Vlu** | 0.28 | 0.00 | 0.00 | 0.05 |
| 2 | **Avg_MPV** | 0.29 | 0.00 | 0.00 | 0.05 |
| 3 | **Msspl_Prp_Vlu** | 0.40 | 1.00 | 1.00 | 1.00 |
| 4 | **Msspl_Cls** | 0.76 | 0.79 | 0.42 | 0.05 |
| 5 | **Msspl_Prp** | 0.57 | 0.86 | 0.00 | 0.27 |
| 6 | **Out_Prp_Vlu** | 0.21 | 1.00 | 1.00 | 0.73 |
| 7 | *Im_DT* | 1.00 | 1.00 | 1.00 | 1.00 |
| 8 | *Und_Cls* | 1.00 | 0.00 | 1.00 | 1.00 |
| 9 | *Und_Prp* | 1.00 | 1.00 | 1.00 | 0.54 |
| 10 | *Dsj_Cls* | 1.00 | 1.00 | 1.00 | 1.00 |
| 11 | *Dsj_prp* | 1.00 | 1.00 | 1.00 | 1.00 |
| 12 | *FP* | 1.00 | 1.00 | 1.00 | 1.00 |
| 13 | *IFP* | 1.00 | 1.00 | 1.00 | 1.00 |
| 14 | *Misplc_cls_Prp* | 1.00 | 1.00 | 1.00 | 1.00 |
| 15 | *Msusg_Prp* | 1.00 | 1.00 | 1.00 | 1.00 |

[7]http://www.geonames.org/ontology/ontology_v3.1.rdf

[8]https://babbage.inf.unibz.it/trac/obdapublic/raw-attachment/wiki/Example_MovieOntology/movieontology.owl

[9]http://oaei.ontologymatching.org/2013/anatomy/anatomy-dataset.zip

[10]http://citeseer.rkbexplorer.com/models/dump.tgz

The values of fifteen metrics are reported in Table 2. The first six metrics are related to semantic accuracy and highlighted with **'bold'**, while the last nine rows of the table refer to semantic accuracy metrics. According the metric definitions presented in Section 4, all of the metrics are defined as the ratio of the desired outcomes to total outcomes. This adheres to the convention where value '1' for a given metric represents the highest quality in terms of a quality deficiency and value '0' denotes poorest quality regarding the same deficiency. In our study, a preferred way for computation of metrics values is to calculate the ratio of the quality deficiencies and then subtract the result from '1'. In this way, all of the values of quality-driven metric are in the range of [0, 1], where the value '1' for a specific metric means that there is no quality deficiency measured by that metric. For example, the first metric is Miss_Prp_Vlu which is defined in order to measure the ratio of the properties defined in the schema, but not present in the dataset. As shown in the first row of Table 2, the value of this metric for the first dataset (*Geonames*) is 0.28. As '1' represents the most desirable value for this metric, it means that 28% of the properties are presented for the instances and 72% are missing. Similarly for the sixth metric, it is reported that only 21% of triples used in the *Geonames* ontology, are not out of range, but the value of this metric in the second column shows that none of the properties values of *IMDB* are out-of-range. The other values of these metrics can be interpreted similarly. For better representation of metrics behavior of the metrics over experimented datasets, a radar chart is depicted in Figure 1.
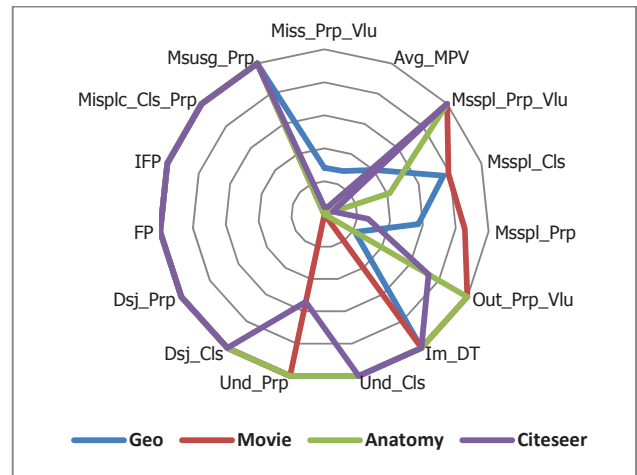


**Figure 1.The behavior of the proposed metrics**

As shown in the Figure, four experimented datasets are in different levels of semantic accuracy, while they are similar in terms of syntactic accuracy. In light of the values of syntactic accuracy metrics (rows 7-15 of Table 2), it is obvious that most of them have the value '1' for all experimented datasets. For example, *FP (Metric 12)* is measured by the ratio of triples which contain functional properties with inconsistent values. The metric value '1'in

our observations indicates that there is no functional property with inconsistent values in any of the datasets. Similarly, the value '1' for all of syntax accuracy metrics indicates that there is no problem in all datasets in terms of syntax errors. One possible explanation can be that in most cases, datasets are built using tool support, which ensures that the datasets are syntactically correct.

According to the reported metric values, it is clear that selected datasets are not in the same level of quality in terms of fifteen metrics. In addition, the trends of metrics values show the appropriate behavior of the proposed metrics.

## VI. CONCLUSION AND FUTURE WORKS

The main objective of this work is to propose a set of metrics for the assessment of the accuracy of datasets from two viewpoints of semantic and syntactic accuracy. Thus, by characterizing data accuracy in the context of LOD, a set of fifteen metrics are proposed. Finally, by putting the metrics under empirical validation, the results of our observation are discussed.

In the future work, we are going to find relations between the metrics values and perceived quality by collecting the opinions of the experts in LOD domain. If the proposed metrics are shown to have meaningful correlation with the quality, then we are able to predict the quality of any dataset once it is integrated into the LOD. The results will also help publishers to filter out low-quality data, which in turn enables data consumers to make better and more informed decisions when using the shared datasets.

## REFERENCES

[1]     F. Naumann and C. Rolker, "Assessment Methods for Information Quality Criteria," in *5'th Conference on Information Quality* Boston, Mass, USA, 2000, pp. 148-162.
[2]     C. Batini and M. Scannapieca, *Data quality: concepts, methodologies and techniques*, 1.0 ed.: Springer, 2006.
[3]     B. Behkamal, M. Kahani, S. Paydar, M. Dadkhah, and E. Sekhavaty, "Publishing Persian linked data; challenges and lessons learned," in *5th International Symposium on Telecommunications (IST)*, 2010, pp. 732-737.
[4]     O. Hartig and J. Zhao, "Using Web Data Provenance for Quality Assessment," *SWPM*, vol. 526, 2009.
[5]     C. Bizer and R. Cyganiak, "Quality-driven information filtering using the WIQA policy framework," *Web Semantics: Science, Services and Agents on the World Wide Web,* vol. 7, pp. 1-10, 2009.
[6]     Y. Lei, A. Nikolov, V. Uren, and E. Motta, "Detecting Quality Problems in Semantic Metadata without the Presence of a Gold Standard," in *5th International EON Workshop at International Semantic Web Conference (ISWC'07)*, Busan, Korea, 2007, pp. 51-60.
[7]     S. Brüggemann and F. Grüning, "Using ontologies providing domain knowledge for data quality management," in *Networked Knowledge-Networked Media*, ed: Springer, 2009, pp. 187-203.
[8]     A. Hogan, A. Harth, A. Passant, S. Decker, and A. Polleres, "Weaving the pedantic web," in *3rd International Workshop on Linked Data on the Web (LDOW2010)*, Raleigh, North Carolina, 2010.
[9]     C. Fürber and M. Hepp, "Using semantic web resources for data quality management," in *Knowledge Engineering and Management by the Masses*, ed: Springer, 2010, pp. 211-225.
[10]     V. Peralta, "Data freshness and data accuracy: A state of the art," Instituto de Computacion, Facultad de Ingenieria, Universidad de la Republica2006.
[11]     R. Y. Wang, D. M. Strong, and L. M. Guarascio, "Beyond accuracy: What data quality means to data consumers," *Journal of Management Information Systems,* vol. 12, pp. 5-33, 1996.
[12]     ISO, "ISO/IEC 25012- Software engineering - Software product Quality Requirements and Evaluation (SQuaRE)," in *Data quality model*, ed, 2008.
[13]     F. Naumann, U. Leser, and J. C. Freytag, "Quality-driven integration of heterogeneous information systems," presented at the 25th International Conference on Very Large Data Bases (VLDB'99), Edinburgh, Scotland, UK, 1999.
[14]     T. C. Redman and A. Blanton, *Data quality for the information age*: Artech House, Inc., 1997.
[15]     B. Behkamal, M. Kahani, E. Bagheri, and Z. Jeremic, "A Metrics-Driven approach for quality Assessment of Linked open Data," *Journal of Theoritical and Applied Electronic Commerce Research* vol. 9, pp. 64-79, 2014.
[16]     V. R. Basili, G. Caldiera, and H. D. Rombach, "The goal question metric approach," in *Encyclopedia of software engineering*, ed: John Wiley & Sons, 1994, pp. 528-532.
[17]     S. A. Sarcia, "Is GQM+ Strategies really applicable as is to non-software development domains?," in *ACM-IEEE International Symposium on Empirical Software Engineering and Measurement*, 2010, p. 45.
[18]     A. Zaveri, A. Rula, A. Maurino, R. Pietrobon, J. Lehmann, S. Auer*, et al.*, "Quality Assessment Methodologies for Linked Open Data," *Submitted to Semantic Web Journal,* 2013.
[19]     O. Hartig, "Trustworthiness of data on the web," in *Proceedings of the STI Berlin & CSW PhD Workshop*, 2008.
[20]     A. Zaveri, D. Kontokostas, M. A. Sherif, L. Bühmann, M. Morsey, S. Auer*, et al.*, "User-driven quality evaluation of dbpedia," in *Proceedings of the 9th International Conference on Semantic Systems*, 2013, pp. 97-104.
[21]     A. F. Karr, A. P. Sanil, and D. L. Banks, "Data quality: A statistical perspective," *Statistical Methodology,* vol. 3, pp. 137-173, 2006.
[22]     J. Ashraf, O. K. Hussain, and F. K. Hussain, "A framework for measuring ontology usage on the web," *The Computer Journal,* vol. 56, pp. 1083-1101, 2012.
[23]     B. Behkamal. (2013). *The code of metrics calculation tool (1.0 ed.)*. Available: https://bitbucket.org/behkamal/new-metrics-codes/src on 2014-10-18