

Data Quality in a Big Data Context

Franco Arolfo and Alejandro Vaisman^(✉)

Instituto Tecnológico de Buenos Aires, Buenos Aires, Argentina
{farolfo, avaisman}@itba.edu.ar

Abstract. In each of the phases of a Big Data analysis process, data quality (DQ) plays a key role. Given the particular characteristics of the data at hand, the traditional DQ methods used for relational databases, based on quality dimensions and metrics, must be adapted and extended, in order to capture the new characteristics that Big Data introduces. This paper dives into this problem, re-defining the DQ dimensions and metrics for a Big Data scenario, where data may arrive, for example, as unstructured documents in real time. This general scenario is instantiated to study the concrete case of Twitter feeds. Further, the paper also describes the implementation of a system that acquires tweets in real time, and computes the quality of each tweet, applying the quality metrics that are defined formally in the paper. The implementation includes a web user interface that allows filtering the tweets for example by keywords, and visualizing the quality of a data stream in many different ways. Experiments are performed and their results discussed.

Keywords: Data quality · Social networks · Big Data

1 Introduction and Motivation

The relevance of so-called Big Data has been acknowledged by researchers and practitioners even before the concept became widely popular through media coverage [1]. Although there is no precise and formal definition, it is accepted that Big Data refers to huge volumes of heterogeneous data that must be ingested at a speed that cannot be handled by traditional database systems tools. Big Data is characterized by the well-known “4 V’s” (volume, variety, velocity, and veracity), implying that not only the data volume is relevant, but also the different kinds of structured, semistructured and unstructured data, the speed at which data arrives (e.g., real time, near real time), and the reliability and usefulness of such data. However, it is also acknowledged that most of the promises and potential of Big Data are far from being realized. This gap between promise and reality is due to the many technical problems and challenges that are usually overlooked, although the database research community has warned about them, namely heterogeneity, scale, timeliness, complexity, and privacy, among other ones [3].

In each of the phases in a Big Data scenario, data quality (DQ) plays a key role, as the very nature of the “4 V’s” suggest. The largely studied concepts of DQ must be revisited and re-studied in a Big Data context, since, as it will be discussed in this paper, many new problems appear, which are not present

in traditional relational databases scenarios [8]. Intuitively, each of the “V’s” define a different context for data analysis, and therefore, for DQ. Thus, there is a strong relationship between the work about contexts in DQ (e.g., [7, 11]) and the problems of DQ in Big Data, since different notions of quality must be used for different types of Big Data. In particular, this paper deals with DQ in a real-time scenario, specifically Twitter feeds. This is a typical scenario where data come at high speed, highly unstructured, and with very volatile reliability and usefulness. All of these characteristics are the complete opposite of a relational database analytics scenario, where data are highly structured, and cleaned, transformed and analyzed offline. Therefore, DQ must be addressed considering these differences.

In spite of the relevance of the topic, there has been not much work so far, in particular regarding the implementation of quality processes over Big Data sources. This paper tackles this issue. More concretely, the contributions of this work are: (1) The definition of DQ dimensions and metrics in a Big Data scenario where data arrive as unstructured documents and in real time. Traditional DQ dimensions are redefined, to address those particular characteristics. This general scenario is instantiated to study the concrete case of Twitter feeds. (2) The implementation of a system that acquires tweets in real time, and computes the quality of each tweet, applying the quality metrics defined formally in the paper. The implementation includes a web user interface that allows filtering the tweets e.g., by keywords, computing their data quality, and visualizing the DQ, not only the overall one, but also along each dimension. (3) An experimental study of the quality of the feeds, using the tool described above. This study is aimed at showing how DQ can be used to determine the attributes that characterize the different quality of the tweets, filter out bad quality data, or validate the conclusions drawn in the data analysis phase.

The remainder of the paper is structured as follows. Section 2 discusses related work. In Sect. 3, the traditional DQ dimensions and metrics are presented, while Sect. 4 studies the DQ dimensions and metrics for Big Data. Section 5 presents the computation of DQ metrics to evaluate the quality of a tweet based on a collection of weighted metrics. Section 6 introduces the implementation of the system, and Sect. 7 presents an experimentation and reports and discusses the results. Section 8 concludes the paper.

2 Related Work

Ensuring the quality of data in databases has long been a research topic in the database community. Research in DQ has resulted in the definition of dimensions, metrics, and methods to assess the quality of a database [16], and also in an ISO standard specification¹. In spite of this, classic research considers DQ as a concept independent of the context in which data are produced and used, which is clearly not enough to solve complex problems, particularly in current times, when, among other facts, ubiquitous computing requires accounting for space and time when a query is being answered. Strong et al. [14] realized this problem, and claimed that data quality is highly dependent on the context,

¹ <http://iso25000.com/index.php/en/iso-25000-standards/iso-25012>.

which became an accepted fact thereon. The rationale for this conclusion was based on the fact that, similarly to quality in general, DQ cannot be assessed independently of the consumers who choose and use the products (e.g., [11]). The former proposes a system where contextual information allows evaluating the quality of blood pressure data. The latter proposes a framework that allows context-sensitive assessment of DQ, through the selection of dimensions for each particular decision-maker context and her information requirements.

There is a large corpus of work regarding data context management with a wide variety of uses. It is widely accepted that most modern applications, particularly over the web, are required to be context-aware. Bolchini et al. [6] presented a survey of context models, with a well-defined structure, that identified some important aspects of context models. In particular, they remark that models must account for *space*, *time*, *context history*, *subject*, and *user profile*. Preferences in databases have also been extensively studied [7, 13]. In the multi-dimensional databases domain, [10] proposes to define the context through the use of logic rules, representing the database as a first-order logic theory.

Recently, [4, 8] study the particularities of data quality in the context of Big data, that is, how the “4 V’s” mentioned in Sect. 1 impact on well-known DQ dimensions and metrics used in traditional structured databases [5]. The main message in [8] is that Big Data quality should be defined in source-specific terms and according to the specific dimension(s) under investigation. In some sense, this means that the context is again present in the Big Data scenario when quality is addressed. The present paper builds from these studies, as will be clear in the remaining sections.

3 A Short Background

Data Quality (DQ) is a multi-faceted concept, represented by different dimensions, each one referring to a different quality aspect [5, 14]. A *DQ dimension* captures a facet of DQ, while a *DQ metric* is a quantifiable instrument that defines the way in which a dimension is measured. Since a DQ dimension is in general a wide concept, an associated metric allows specifying a concrete meaning for the dimension. As a consequence, many different metrics can be associated to the same DQ dimension, and their application will measure several different aspects of the dimension. In a broader sense, the *quality of an object or service* represents how much this object or service fits the needs to solve a given problem. That is, quality is not absolute to the object or service *per se*, but relative to the problem to be solved. This is the approach followed in this work.

3.1 Dimensions and Metrics

While a large number of DQ dimensions were proposed in the literature, there is a basic set of them, which are generally acknowledged to be representative of the quality of data [5, 12]. This set includes accuracy, completeness, consistency, freshness (or timeliness), among other ones.

- *Accuracy*: Specifies how accurate data are, and involves the concepts of *semantic accuracy* and *syntactic accuracy*. The former refers to how close is

a real-world value to its representation in the database. The latter indicates if a value belongs to a valid domain. In other words, it describes the closeness between a value v and a value v' , considered as the correct representation of the real-life phenomenon that v aims at representing. For example, if someone wants to type the name “John” but typed “Jhn”, there is an accuracy issue.

- *Completeness*: Represents the extent to which data are of sufficient breadth, depth, and scope for the task at hand. For relational databases, this can be characterized as the presence/absence and meaning of null values, assuming that the schema is complete.
- *Redundancy*: Refers to the representation an aspect of the world with the minimal use of information resources.
- *Consistency*: Refers to the capability of the information to comply without contradictions with all the rules defined in a system. For example, in a relational database constraints are defined to guarantee consistency.
- *Readability*: Refers to the ease of understanding of information. This could be the case when, for example, a hand-written paragraph is scanned, and some of the characters are not well defined.
- *Accessibility*: Also called *availability*, is related to the ability of the user to access the information.
- *Trust*: Refers to how much the information source can be trusted, and therefore to what extent data are reliable. For example, people may rely on Facebook or Twitter posts to find out the quality of a movie, or check the IMDB site at <http://www.imdb.com>, which might provide more reliable data.
- *Usefulness* (cf. [8]): This is related to the benefits a user can obtain when using the data to produce information. For example, to observe technical details present in a picture of a painting, a user would choose the image with the highest contrast. Again, this is also a contextual quality dimension: a lower-quality picture may be enough for some users or for some kinds of requirements, while clearly not enough when the details are needed.

To quantify these dimensions and to be able to assess DQ according to them, the concept of *metrics* must be introduced. Mathematically, a DQ *metric* for a dimension D is a function that maps an entity to a value, such that this value, typically between 0 and 1, indicates the quality of a piece of data regarding the dimension D . For a given dimension, more than one metric could be defined and combined to obtain a concrete quality value. Note that metrics are highly context-dependent. For example, the readability of a hand-written text may be influenced not only by the text content, but also by the way the user writes. The same occurs with metrics for other DQ dimensions.

3.2 Big Data Quality

In a Big Data context, datasets are too large to store, analyze, handle or process, for the traditional database tools. As explained above, Big Data are characterized by the well-known “4 V’s”, namely *Volume* (size of the datasets), *Velocity* (speed of incoming data, e.g., the number of tweets per second (TPS)), *Variety* (refers to the type and nature of the data), and *Veracity* (the reliability of the data, which, in this context, is greatly volatile, even within the same data

stream). In the literature, many other “V’s” can be found, but only these four will be considered in the present work. According to the structure of data, they can be classified in: (a) *Structured*, where each piece of information has an associated fixed and formal structure, like in traditional relational databases; (b) *Semi Structured*, where the structure of the data has some degree of flexibility (e.g., an XML file with no associated schema, or a JSON response from an API, whose structure is not completely defined); (c) *Unstructured*, where no specific structure is defined. Further, the United Nations Economic Commission for Europe (UNECE) classifies Big Data according to the data sources in: *Human sourced*; *Process mediated*; and *Machine generated* [15]. These are explained next.

- *Human-sourced data*: Information people provide via text, photos or videos. Usually, this information lacks of a fixed structure, like the texts written in natural language. Therefore, the information streamed here is *loosely structured* and often ungoverned. Examples are social networks posts (Facebook, Twitter), YouTube videos or e-mails, and, in general, data coming from social networks.
- *Process-mediated data*: This is the information that concerns some business events of interest, like the purchase of a camera in an e-commerce site or the sign-up of clients in a system. This information is *highly structured*, such as relational databases, coming from traditional Business systems.
- *Machine-generated data*: Refers to the data resulting of the tracking of sensors of the physical world (e.g., temperature sensors, human health sensors, GPS coordinates, etc.). This type of source is associated with very large amounts of data, given the constant tracking of the sensors. In general, these are data coming from the so-called Internet of Things.

Given these characteristics of Big Data, the DQ along the dimensions explained in Sect. 3.1 must be quantified using metrics specific to such a context, therefore the typical quality metrics used for structured, process-mediated data must be adapted to this new situation. This is studied in the next section.

4 Data Quality Dimensions and Metrics in a Big Data Context

This section studies how the DQ dimensions can be used in a Big Data scenario. The study focuses in *human-sourced generated data*. The next sections describe how these dimensions can be applied to address the quality of Twitter² streams. Metrics for the dimensions defined here are presented later.

- *Readability (r)*. Given a dictionary D , and a collection of words considered as valid in a document x , the *Readability* of x , denoted $r(x)$ is defined as the quotient between the valid words in x and all the words in x , if any, otherwise it is zero. That is, given a set W of the words (valid and non-valid) that are present in the document x , the readability of x is

$$r(x) = \begin{cases} \frac{\#\{w \in W \wedge w \in D\}}{\#\{w \in W\}} & \text{if } W \neq \emptyset \\ 0 & \text{if } W = \emptyset \end{cases}$$

² <http://www.twitter.com>.

In the remainder, the problem to be addressed will refer to tweets in a Twitter stream, thus x will represent a tweet.

- *Completeness* (c). Consider an object x in domain, and an array $props_p$ that contains the names of the properties required to describe x for a given problem p ; assume that x is represented as a collection of $(property, value)$ pairs of the form $\{(p_1, v_1), \dots, (p_n, v_n)\}$, such that v_i is a value for p_i . If a property $p_i \in x$ has associated a non-null value v_i , it is called well-defined. There is also a function *validPropsOf* that, given an object x retrieves the set of well-defined properties in it. The *Completeness* of x , denoted $c(x)$ tells if all the properties in $props_p$ are well-defined in x . It is computed as:

$$c(x) = \begin{cases} 1 & \text{if } props_p \subset validPropsOf(x) \\ 0 & \text{otherwise} \end{cases}$$

Example 1. Given an object (a tweet) x , such that $x = \{text: "I like Bitcoin", user: null\}$, and an array $props_p = [text]$, it follows that $c(x) = 1$. Consider now $props_p = [text, user]$. In this case, $c(x) = 0$, since the user property is not well defined, because it has a null value.

- *Usefulness* (u). Since this paper deals with human-sourced datasets, it will be assumed that this property is directly related to the possibility of (among others):
 - (a) Detecting a sentiment, whether positive or negative, in an object x , say a tweet or post. Therefore, if x reflects a positive or negative feeling about a certain topic or person, x will be considered useful. If the sentiment is neutral, or no sentiment could be computed by a Natural Language Processing (NLP) tool, x will be considered not useful.
 - (b) Detecting the domain or topic of x , for example, politics, marketing, sports, and so on.

Many other ways of assessing usefulness could be considered, but this is outside the scope of this paper. In the remainder, Usefulness is defined as:

$$u(x) = \begin{cases} 1 & \text{if } (sentiment(x) = P \vee sentiment(x) = N) \\ 0 & \text{otherwise} \end{cases}$$

- *Trustworthiness* (t). In a social network (or, in general, for human-sourced datasets) anyone in general can publish any kind of information anywhere, whether truthful or not. Although this is mentioned here for completeness, validating the trustfulness of a post is outside the scope of this paper.

5 Computing Data Quality

As discussed above, the definition of DQ is not the same for all contexts and problems, but normally it depends on them. That is, DQ depends on the problem and the domain model at hand. This section provides a wide and general definition of DQ, that can be instantiated as needed.

Definition 1 (Problem (p)). A problem p is defined as a string or sequence of characters that defines the problem to be solved in a human-readable way.

Example 2. A problem can be defined as: Given this Twitter feeds stream, what are the best quality tweets for the hashtag #2020Elections?

Definition 2 (Domain Model (X)). The set of objects whose quality will be measured.

Example 3. For the case that will be studied in the next section, the domain model is defined as a set of Twitter feeds.

Definition 3 (Data Quality Metric (m_{Xp})). A Data Quality Metric is a function $m_{Xp} : X \rightarrow [0 \dots 1]$, such that, given $x \in X$, and a problem p , then $m_{Xp}(x) = 0$ if x contains data of very poor quality for the given problem p , and $m_{Xp}(x) = 1$ if x contains data of very good quality to fit the problem.

Definition 4 (Weight of a Data Quality Metric ($m_{Xp}.weight$)). Each DQ metric has an associated weight, which is a scalar value between 0 and 1, that measures the relevance of the metric for solving the problem p .

Definition 5 (Data Quality (Q_{Xp})). Consider a problem p , a domain X , and a set of metrics $M_{Xp} = \{m_1, m_2, \dots, m_n\}$. Each m_i is a DQ metric function. Note that n is an integer number greater than zero and the set M_{Xp} is finite. Data Quality (Q_{Xp}) is a function $Q_{Xp} : X \rightarrow [0 \dots 1]$ such that $Q_{Xp}(x) = g(m_1, m_2, \dots, m_n)(x)$, where g is a function $g : (X \rightarrow [0, 1])^n \rightarrow (X \rightarrow [0 \dots 1])$.

In this paper, the quality of a tweet x will be calculated as $Q(x) = g_{(r,c,u)}(x) = \sum_{m=\{r,c,u\}} m(x) * m.weight$, where r , c and u are the metrics for *Readability*, *Completeness* and *Usefulness* respectively, defined in Sect. 4.

Example 4. The Quality value of the following tweet x , using the weights values $r.weight = 0.5$, $c.weight = 0.25$ and $u.weight = 0.25$, is computed as follows.

```
- text: "I love Big Data Quality m#a!sc"
- id: 1
- coordinates: [48.864716, 2.349014]
```

– *Readability* (r)

$$r(x) = \frac{\#\{I, love, Big, Data, Quality\}}{\#\{I, love, Big, Data, Quality, m\#a!sc\}}$$

$$r(x) = \frac{5}{6} = 0.833$$

– *Completeness* (c). Consider that $props_p = \{text, id\}$. Then:

$$\{text, id\} \subset \{text, id, coordinates\}, \text{ and } c(x) = 1.$$

– *Usefulness* (u). The text provided expresses positive sentiment, thus:

$$sentiment(x) = P, \text{ and } u(x) = 1.$$

Finally, the quality value for x is $Q(x) = 0.83 * 0.5 + 1 * 0.25 + 1 * 0.25 = 0.915$.

6 Implementation

This section presents and describes the implementation of the concepts explained in previous sections, applying them to analyze the quality of Twitter feeds streams. The architecture is described first, detailing the technological components and how they interact with each other for capturing, filtering, and displaying the results. Finally, the user interface is described. The goal of the implementation is to develop a system that can let users to analyze a stream of Twitter feeds, based on a particular keyword-led search, and visualize the results to gain insight on the quality of the requested data.

The core of the system is an Apache Kafka³ cluster of three brokers. Kafka is a distributed streaming platform for capturing, processing and storing data streams. The implemented cluster has three nodes running Kafka. A Zookeeper⁴ service coordinates the cluster and manages the message topics. Besides the Kafka core, there are three components, one to produce data, one to consume data, and one to process and display data. Figure 1 illustrates these components and their orchestration. The components are briefly described next.

- *Kafka Producer Service*: A Java 8 program exposing a REST API using the Spring Boot⁵ framework. This API starts searches over the Twitter API and publishes the data to a particular *topic*.
- *User Interface (UI) Proxy*: In order to show the results, the UI needs a proxy that consumes the data from the producer and sends it to the UI via a persistent web socket connection, using the socket.io⁶ framework. Also, this Node.js service uses express.js framework in order to expose a REST API, through which the UI can request data in a new search.
- *Web UI*: The web UI is built on top of the dc.js library, which uses the Big Data processing framework crossfilter.js and the data visualization library d3.js. This UI is composed of five parts: (a) The general DQ results; (b) The DQ results per dimension; (c) The visualization of the presence of a dimension in the stream, that is, the percentage of tweets with values for each dimension; (d) The Tweets vs. re-tweets part, comparing DQ values considering or leaving out re-tweets, respectively; and (e) The verified vs. unverified users part, indicating the DQ for tweets coming from verified or unverified users. Of course, this UI can be easily extended according to the analysis needs, to gain insight on the DQ of the data streams.

The data flow works as follows. First, the UI performs a request to the UI Proxy via the REST API in order to start a search using some query, indicating the *topic* ID to use (just a randomly-generated name), and some optional advanced parameters (a list of *completeness* properties to analyze, and the weights of each DQ metric). Then, the Proxy performs a REST API call to a Kafka producer service instance in order to start the ingestion of the feeds using those parameters. At this point, the Kafka producer service creates the topic

³ <https://kafka.apache.org/>.

⁴ <https://zookeeper.apache.org/>.

⁵ <https://projects.spring.io/spring-boot/>.

⁶ <https://socket.io/>.

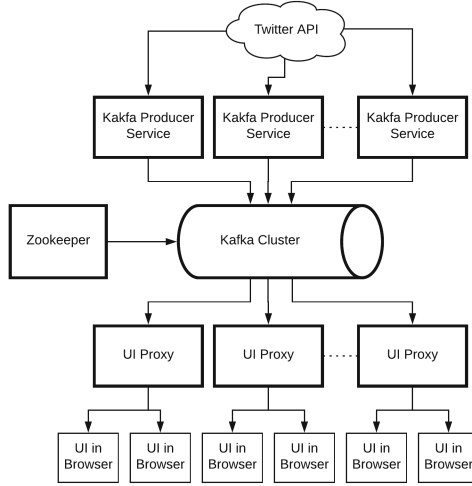


Fig. 1. Architecture diagram.

involved and starts a Kafka producer that fetches the feeds from the Twitter API and publishes its to the topic. After this initialization, the UI proxy starts a Kafka consumer peeking the feeds from the respective topic, and forwards the information to the UI. Finally, the UI processes the records using `crossfilter.js` and shows the data in real time using the `dc.js` library.

Remark 1. The system may scale horizontally as needed, just adding more Kafka producer services. To scale the Kafka Cluster more workers can be added, and Zookeeper will take care of their coordination.⁷

7 Experiments

This section describes the experiments performed over the implementation presented in Sect. 6, reports the results, and discusses them.

7.1 Use-Cases Description

The goal of the experiments is to illustrate how the quality of a Twitter stream can be measured using the dimensions and metrics presented above. Of course, there are countless ways in which the quality of data in tweets can be analyzed. The experiments presented here are just aimed at showing how the concepts discussed in previous sections can be studied using the tool presented in Sect. 6. The quality dimensions considered in all cases are: readability, completeness, and usefulness, with the metrics described in Sect. 4, and with the following weights: 0.5 for readability, and 0.25 for completeness and usefulness. Of course, the user can modify the weights according to her analysis needs. The dictionary used to

⁷ The system is available upon request to the authors.

check readability is given in [2] and contains 479,000 english words, including acronyms, abbreviations and even Internet slang. To compute sentiment (for usefulness), the Stanford CoreNLP software was used [9]. In all cases, the overall DQ of the stream is computed, as well as each DQ dimension individually, and the comparisons allowed by the UI (indicated in Sect. 6) are displayed. The experiments are aimed at:

- (a) Comparing the DQ of the whole stream of tweets, against the DQ of a stream filtered by a set of keywords related to some topic. The hypothesis is that the latter are more likely to have better quality than the former.
- (b) Performing the same comparisons above, but requesting the presence of different sets of properties (that is, changing the requested schema). This will give insight on which are the properties more likely to be present in a stream of tweets, and investigating the impact on this of the keyword filtering.
- (c) Determining if there is a correlation between the DQ of a stream, and the percentage of re-tweeted tweets that it contains. The hypothesis here is that a tweet with a high number of re-tweets is likely to be of high quality.
- (d) Comparing the quality of tweets from verified and not-verified users.

Next, the problems used to address the goals above, are described.

Problem 1. The first problem p_1 consists in analyzing the tweets in a stream, with no keyword filtering, i.e., all tweets provided by the Twitter API. The UI allows to indicate the set of properties considered for schema completeness. In this case, $props_{p_1} = \{\text{id}, \text{id_str}, \text{lang}, \text{retweet_count}, \text{usr}, \text{text}, \text{source}\}$.

Problem 2. For the second problem p_2 , the same set $props_{p_1}$ is used, but the stream is filtered using the keywords: $\{\text{Trump}, \text{Obama}, \text{Hillary}\}$.

Problem 3. The third problem p_3 consists in analyzing the tweets in a stream like in the previous problem, but considering a larger set of properties, namely all the ones supported by the UI. This allows to study how are the properties distributed in Twitter streams, that is, how many tweets contain the space coordinates or the language, for example.

Remark 2. Again, it must be clear that it is not intended here to draw conclusions on the DQ of Twitter feeds, but to suggest how the concepts and tools presented in this paper can be used to analyze such feeds.

7.2 Results and Discussion

Performance results were quite satisfactory. Tweets were captured and displayed at a rate of 2000 per minute (for non-filtered tweets), and at about 600 per minute, for filtered tweets, depending on how many tweets pass the filters.

The results obtained for the problems above are commented next, and illustrated by graphics. The figures show the status of some runs, such that after several thousands of consumed tweets, the results become stable, that is, the graphs do not change significantly.

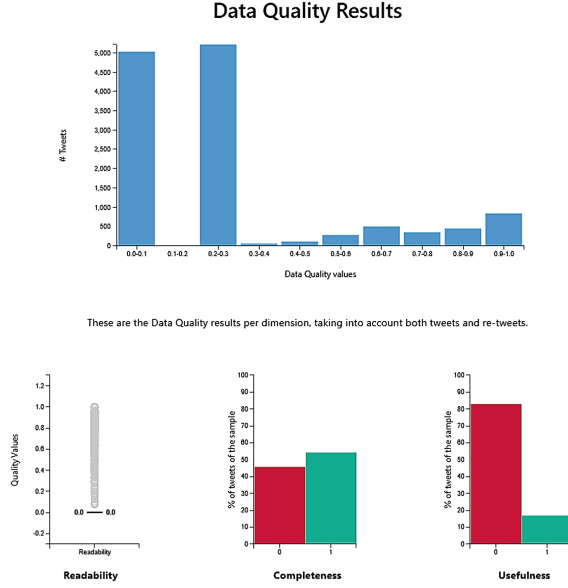


Fig. 2. Results for problem 1.

Figure 2 shows a portion of the UI, displaying the results obtained from running the system with the conditions of *Problem 1*. In the upper part, the Y-axis represents the number of tweets, and the X-axis the DQ values, in intervals of 0.1. It can be seen that the general DQ is low (most of the tweets fall on the left part of the graph). Readability, completeness and usefulness are shown in the lower part. The first one is displayed as a boxplot, while the other two are displayed as bar charts. Since Readability is low, most of the values in the boxplot are outliers. About 55% of the tweets have values in all the fields in $props_{p_1}$, and more that 70% of the tweets have usefulness = 0. Results also showed (graphics omitted for the sake of space) that all DQ values are better when only re-tweeted tweets are considered, and for tweets posted by verified users.

Figure 3 displays the results obtained under the conditions of *Problem 2*, using the keywords Trump, Obama, and Hillary (an “OR” condition). The intention is to capture tweets with political content, based on the hypothesis that their quality should be better than for non-filtered tweets. It can be seen that the general DQ is much better, and on the upper part of the right-hand side, most of the tweets fall on the right part of the graph. Readability, completeness and usefulness are much better than in Fig. 2, and all DQ values (like in Problem 1) are better when only re-tweeted tweets are considered, and are also better for tweets posted by verified users.

Figure 4 displays, for each property supported by the UI, the percentage of tweets that contains values for that property. Only the values for Problem 1 are displayed, due to space limitations. Anyway, the set of properties that were checked by the user is the same for Problems 1 and 2. The darker portion of the circles indicates the percentage of tweets containing no value for the property.



These are the Data Quality results per dimension, taking into account both tweets and re-tweets.

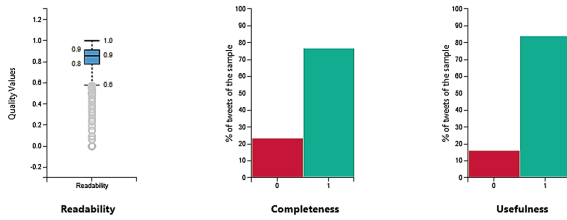


Fig. 3. Results for problem 2.



Fig. 4. Completeness for problem 1.

It can be seen that the `coordinates` field (at the top-right corner) has no value for any tweet, that is, no tweet is geo-referenced.

Figure 5, shows the results for Problem 3. Recall that in this problem, tweets are captured like in Problem 2, but the property set includes all properties supported by the system, for example, the spatial coordinates of the tweets.

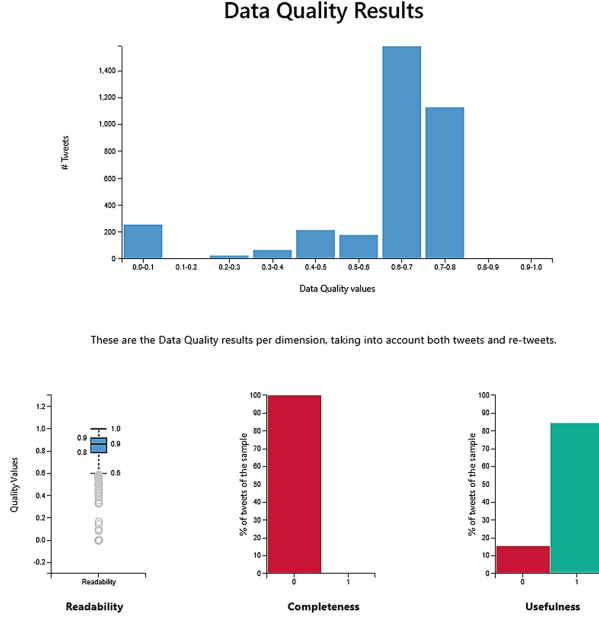


Fig. 5. Results for problem 3.

Figure 4 shows that this property is not satisfied by any tweet. Therefore, the completeness dimension for Problem 3 has value 0 (given that all properties are required to be present), which lowers the overall quality.

8 Conclusion and future work

This paper studied the particularities of assessing data quality in a Big Data context, and presented a system that allows analyzing such quality over Twitter streams in real time. Experiments performed over many different Twitter streams, showed how the concepts presented and tools developed could be applied in a real-world Big Data scenario. Still, there is plenty of room for further work. One line of research could be oriented to define more DQ dimensions and metrics for this or other settings (along the lines of [8]), since, as explained, this is typical context-dependent DQ. Also, new and more sophisticated visualization tools could extend and enhance the implemented framework. All of these will be part of future work.

Acknowledgments. Alejandro Vaisman was partially supported by the Argentinian Scientific Agency, PICT-2014 Project 0787.

References

1. Data, data everywhere (2008). <https://www.economist.com/node/15557443>
2. English-words project (2018). <https://github.com/dwyl/english-words>
3. Agrawal, D., Bernstein, P., Bertino, E., Davidson, S., Dayal, U.: Challenges and opportunities with big data (2011). <https://docs.lib.purdue.edu/cgi/viewcontent.cgi?referer=www.google.com.ar/&httpsredir=1&article=1000&context=cctech>
4. Batini, C., Rula, A., Scannapieco, M., Viscusi, G.: From data quality to big data quality. *J. Database Manag.* **26**(1), 60–82 (2015)
5. Batini, C., Scannapieco, M.: Data Quality: Concepts Methodologies and Techniques. DCSA. Springer, Heidelberg (2006). <https://doi.org/10.1007/3-540-33173-5>
6. Bolchini, C., Curino, C.A., Quintarelli, E., Schreiber, F.A., Tanca, L.: A data-oriented survey of context models. *SIGMOD Rec.* **36**(4), 19–26 (2007). <https://doi.org/10.1145/1361348.1361353>
7. Ciaccia, P., Torlone, R.: Modeling the propagation of user preferences. In: Jeusfeld, M., Delcambre, L., Ling, T.-W. (eds.) *ER 2011. LNCS*, vol. 6998, pp. 304–317. Springer, Heidelberg (2011). https://doi.org/10.1007/978-3-642-24606-7_23
8. Firmani, D., Mecella, M., Scannapieco, M., Batini, C.: On the meaningfulness of “big data quality” (invited paper). *Data Sci. Eng.* 1–15 (2015)
9. Manning, C.D., Surdeanu, M., Bauer, J., Finkel, J., Bethard, S.J., McClosky, D.: The stanford CoreNLP natural language processing toolkit. In: Association for Computational Linguistics (ACL) System Demonstrations, pp. 55–60 (2014). <http://www.aclweb.org/anthology/P/P14/P14-5010>
10. Marotta, A., Vaisman, A.: Rule-based multidimensional data quality assessment using contexts. In: Madria, S., Hara, T. (eds.) *DaWaK 2016. LNCS*, vol. 9829, pp. 299–313. Springer, Cham (2016). https://doi.org/10.1007/978-3-319-43946-4_20
11. Poeppelmann, D., Schultewolter, C.: Towards a data quality framework for decision support in a multidimensional context. *IJBIR* **3**(1), 17–29 (2012)
12. Scannapieco, M., Catarci, T.: Data quality under a computer science perspective. *Archivi Comput.* **2**, 1–15 (2002). <https://www.fing.edu.uy/inco/cursos/caldatos/articulos/ArchiviComputer2002.pdf>
13. Stefanidis, K., Pitoura, E., Vassiliadis, P.: Managing contextual preferences. *Inf. Syst.* **36**(8), 1158–1180 (2011)
14. Strong, D.M., Lee, Y.W., Wang, R.Y.: Data quality in context. *Commun. ACM* **40**(5), 103–110 (1997). <https://doi.org/10.1145/253769.253804>
15. Task Team on Big Data: Classification of types of big data (2007). <https://statswiki.unece.org/display/bigdata/Classification+of+Types+of+Big+Data>
16. Wang, R.Y., Strong, D.M.: Beyond accuracy: what data quality means to data consumers. *J. Manag. Inf. Syst.* **12**(4), 5–33 (1996)