

Scalable data cleansing based on qualitative attributes

Mika Huttunen
Helsinki University

Keywords

data cleansing, big data, scalability

1. INTRODUCTION

Nowadays companies are gathering large amounts of data in different ways such as via user input, and all kinds of sensors. It's also becoming easier and easier for them to use data in all kinds of decision making, analytics, and for example in automating human-involved tasks. Problems arise when the used data is *dirty*, or invalid, and thus can lead into making incorrect decisions. These can sometimes cause serious issues - especially in health care and financing sides [3].

Humans often make data input errors via misspelling, and sensors may grab unwanted noise along the data they're designed to catch. In fact, over 25% of critical data in the world's top companies is flawed [4]. Not to mention that today, the variety of data is also large which leads into collecting data of different formats together. *Data cleansing* is the solution to the above-mentioned problem. For data processors, understanding data cleansing, and the problems it tries to solve is thus naturally important.

Data cleansing is a research field that explores ways to improve *quality* of dirty data. If data isn't of high quality, it means that it has usually both *schema*, and *instance* level problems [5]. Another used definition for dirty data is that it doesn't meet its usage needs. A simple example of that would be such that a front developer of a service can't provide end-users a good UI because the server cannot deliver all kinds of useful data because there are faults in the way the data is stored in the database.

Figure 1 shows how *single-source problems* with data quality can yet be divided into schema and instance level data problems. Single-source problems simply stand for such that can occur in data that is stored in a single site with possible data replication on the database side. The same problems can also arise in *multi-source* systems where the data may also be scattered across several sites.

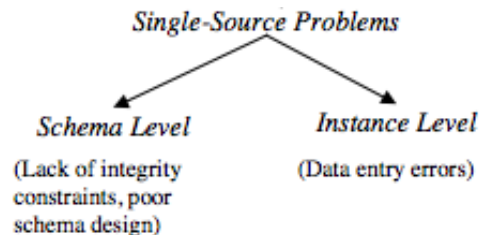


Figure 1: Single-source data quality problems [5]

Instance level problems in data arise as attributes of data *tuples* being out of their scope, or just wrong. These contain misspellings, and *outliers* by for example possible sensor errors, and processing errors before gathered data is actually stored. These problems can be detected via *quantitative techniques* related to data cleansing that focus on error detection and correction based on numerical attributes of data [2]. If for example some framework can reduce multi-dimensional data in two-dimensional space, it's fairly easy for an expert to detect outliers in that space, and thus possible errors in the data.

Schema level problems on the other hand arise as *violations of rules* such as *integrity constraints* (IC), and *functional dependencies* (FD) [3] (p. 289 – 291). As an example, functional dependency

$$ZIP \rightarrow STATE \quad (1)$$

defines a constraint between data attributes *ZIP* and *STATE* such that *ZIP* implies the *STATE*. In other words, if we had a dataset with two tuples having the same *ZIP* attribute, but different values of *STATE*, and we trust rule 1, at least one of them is erroneous.

Duplicate (or partially duplicate) data entries can also be considered both schema level, and instance level data quality problems [3] (p. 283). They usually appear as the same entries with some attributes having different values from each other, and thus should be classified as instance level problems. They can however be detected by similar rules that are used for other schema level problems. Consider us having a person register where each person has its name, phone number, address, zip code (*ZIP*), and state (*STATE*) defined. Now we could have two tuples with the same name, phone number, and address, but different zip codes. An error between these two tuples found be detected by rule 1, and could be solved either automatically, or via human interaction.

This seminar report focuses on highly scalable *qualitative* data cleansing techniques. They focus on detecting data errors in *big data* via rules such as FDs, *CFDs* (conditional functional dependencies), and *DCs* (denial constraints). Ilyas and Chu explain these in their paper [3] (p. 287 - 301).

In traditional ways, the data errors are detected by algorithms that compare data tuples trying to find violations based only on user-defined rules, or also such rules or *patterns* that the data cleansing framework itself generates by processing the data. The found errors can afterwards be fixed either via totally automatic, or (partially) human guided procedures.

The rest of the report is divided as follows: in Section 2, I will introduce recent findings related to big data cleansing, and where we stand on the field at the moment. There should also be a separate section for more accurate introductions to different kinds of state-of-art frameworks (especially HoloClean). Alongside that, I should also have a section for discussing weak points in the papers. So far, I haven't really found anything else than somewhat bad inner organisation of different topics - especially in [3]. And beside them, I would also of course add a section of more detailed analysis of the papers, and what there is left to be improved with data cleansing. As far as I understand, HoloClean is an awesome framework for big data cleansing tasks when all the data is located under a single source, or site. Problems however arise when applying data cleansing for multi-source data. This is especially because the shipment of huge amounts of data over the network is slow [3] (p. 379).

2. LITERATURE REVIEW

Data cleansing has been an interesting topic for decades, and there has been a large amount of research on the field in the recent years, Extensive summary from 2015 by Ilyas and Chu [3] discusses recent techniques on both error detection and repairing of relational data. The errors are generally detected in data as violations of rules such as FDs, CFDs and DCs.

Ilyas and Chu define violation respect to a rule r being the minimal subset of database cells such that at least one of the cells has to be modified to satisfy r . Here cell stands for an attribute value of a tuple, or a row. The same definition can of course be extended to non-relational, but structured data that NoSQL document databases can consist of.

The manual construction of data quality rules such as FDs, CFDs and DCs is time-consuming and requires lots of domain expertise. Thus having automatic tools for their generation is essential [1]. Previously implemented *TANE* and *FASTFD* that Ilyas and Chu discuss [3] can be used for finding FDs based on a relational data, and its schema, while *FASTDC* can be used for finding DCs.

The following part of the report is yet to be written out. But the plan was to especially discuss recent frameworks for error detection with relatively small datasets, and then introduce more recent, scalable ones like BigDancing (distributed rule-based error detection and correcting), SampleClean (probability based fixed data sampling while having dirty database on the background), DeDooP / Dis-Dedup (frameworks specifically designed for duplicate tuples detection), and HoloClean. HoloClean is by far the most interesting of the above-mentioned frameworks. This is because it's designed to be applicable for detecting errors that are untrackable by comparing tuples with each other by using

only a single rule.

Abedjan et al. show in their research paper from 2016 that even though frameworks such as BigDancing are highly scalable, and can be applied to cleansing all kinds of datasets, they can track only approximately 40% of all possible errors - even with fully optimised configuration parameters. If an expert is however used to cleanse data by applying several frameworks in the right order for specific data cleansing tasks, much more than 40% of the possible errors are trackable, but the cleansing accuracy still hardly reaches even 70% ??.

This specific research probably inspired partially the same research team in designing HoloClean of which purpose was to overcome the long lasted challenges with data cleansing tools not achieving good cleansing accuracy. According to the Rekatsinas et al., HoloClean achieves over twice the accuracy of existing state-of-art methods with average precision of 90% ???. What makes HoloClean so much better than the existing state-of-art methods is that accordingly to its name, it's a *holistic* data cleansing framework. Holistic means that it can do tuple-based comparisons by not just using one rule at the time, but the set of all defined rules.

Holistic in data cleansing was no way a new innovation. Chu et al. published a research paper on the topic in 2013 ???. For now, I haven't done enough research to find out, why HoloClean wasn't invented sooner though. Perhaps the real need for it wasn't really known before 2016 when Abedjan et al. published their work about cleansing accuracy going significantly up when carefully selecting, which data cleansing frameworks to use, for what tasks and in what order ??.

3. REFERENCES

- [1] X. Chu and I. F. Ilyas. Qualitative data cleaning. *Proc. VLDB Endow.*, 9(13):1605–1608, Sept. 2016.
- [2] J. M. Hellerstein. Quantitative data cleaning for large databases, 2008.
- [3] I. F. Ilyas and X. Chu. Trends in cleaning relational data: Consistency and deduplication. *Found. Trends databases*, 5(4):281–393, Oct. 2015.
- [4] Z. Khayyat, I. F. Ilyas, A. Jindal, S. Madden, M. Ouzzani, P. Papotti, J.-A. Quiané-Ruiz, N. Tang, and S. Yin. Bigdancing: A system for big data cleansing. In *Proceedings of the 2015 ACM SIGMOD International Conference on Management of Data*, SIGMOD '15, pages 1215–1230, New York, NY, USA, 2015. ACM.
- [5] E. Rahm and H. H. Do. Data cleaning: Problems and current approaches. *IEEE Data Engineering Bulletin*, 23:2000, 2000.