



City Research Online

City, University of London Institutional Repository

Citation: Liu, S., Andrienko, G. ORCID: 0000-0002-8574-6295, Wu, Y., Cao, N., Jiang, L., Shi, C., Wang, Y. S. and Hong, S. (2018). Steering data quality with visual analytics: The complexity challenge. *Visual Informatics*, 2(4), pp. 191-197. doi: 10.1016/j.visinf.2018.12.001

This is the accepted version of the paper.

This version of the publication may differ from the final published version.

Permanent repository link: <http://openaccess.city.ac.uk/22366/>

Link to published version: <http://dx.doi.org/10.1016/j.visinf.2018.12.001>

Copyright and reuse: City Research Online aims to make research outputs of City, University of London available to a wider audience. Copyright and Moral Rights remain with the author(s) and/or copyright holders. URLs from City Research Online may be freely distributed and linked to.

City Research Online:

<http://openaccess.city.ac.uk/>

publications@city.ac.uk

Manuscript Details

Manuscript number	VISINF_2018_22
Title	Steering Data Quality with Visual Analytics: the Complexity Challenge
Article type	Review Article

Abstract

Data quality management, especially data cleansing, has been extensively studied for many years in the areas of data management and visual analytics. In the paper, we first review and explore the relevant work from the research areas of data management, visual analytics and human computer interaction. Then for different types of data such as multimedia data, textual data, trajectory data, and graph data, we summarize the common methods for improving data quality by leveraging data cleansing techniques at different analysis stages. Based on a thorough analysis, we propose a general visual analytics framework for interactively cleansing data. Finally, the challenges and opportunities are analyzed and discussed in the context of data and humans. Data quality management, especially data cleansing, has been extensively studied for many years in the areas of data management and visual analytics. In the paper, we first review and explore the relevant work from the research areas of data management, visual analytics and human computer interaction. Then for different types of data such as multimedia data, textual data, trajectory data, and graph data, we summarize the common methods for improving data quality by leveraging data cleansing techniques at different analysis stages. Based on a thorough analysis, we propose a general visual analytics framework for interactively cleansing data. Finally, the challenges and opportunities are analyzed and discussed in the context of data and humans.

Keywords	Data-quality-management, visual-analytics, data-cleansing,
Corresponding Author	Shixia Liu
Corresponding Author's Institution	Tsinghua University
Order of Authors	Shixia Liu, Gennady Andrienko, Yingcai Wu, Nan Cao, Liu Jiang, Conglei Shi, Yu-Shuen Wang, Seokhee Hong

Submission Files Included in this PDF

File Name [File Type]

first page.pdf [Title Page (with Author Details)]

Steering_Data_Quality_with_Interactive_Visualization (2).pdf [Manuscript File]

To view all the submission files, including those not included in the PDF, click on the manuscript title on your EVISE Homepage, then click 'Download zip file'.

Steering Data Quality with Visual Analytics: the Complexity Challenge

Shixia Liu^{a,*}, Gennady Andrienko^{b,c}, Yingcai Wu^d, Nan Cao^e, Liu Jiang^a,
Conglei Shi^f, Yu-Shuen Wang^g, Seokhee Hong^h

^a*Tsinghua University, Beijing, China*

^b*Fraunhofer Institute IAIS, Sankt-Augustin, Germany*

^c*City, University of London, London, UK*

^d*Zhejiang University, Zhejiang, China*

^e*Tongji University, Shanghai, China*

^f*Airbnb, San Francisco, CA, USA*

^g*National Chiao-Tung University, Hsinchu, Taiwan*

^h*University of Sydney, Sydney, Australia*

Abstract

Data quality management, especially data cleansing, has been extensively studied for many years in the areas of data management and visual analytics. In the paper, we first review and explore the relevant work from the research areas of data management, visual analytics and human computer interaction. Then for different types of data such as multimedia data, textual data, trajectory data, and graph data, we summarize the common methods for improving data quality by leveraging data cleansing techniques at different analysis stages. Based on a thorough analysis, we propose a general visual analytics framework for interactively cleansing data. Finally, the challenges and opportunities are analyzed and discussed in the context of data and humans.

Keywords: Data quality management, visual analytics, data cleansing

2010 MSC: 00-01, 99-00

*Corresponding author

Email address: shixia@tsinghua.edu.cn (Shixia Liu)

Steering Data Quality with Visual Analytics: the Complexity Challenge

Shixia Liu^{a,*}, Gennady Andrienko^{b,c}, Yingcai Wu^d, Nan Cao^e, Liu Jiang^a,
Conglei Shi^f, Yu-Shuen Wang^g, Seokhee Hong^h

^a*Tsinghua University, Beijing, China*

^b*Fraunhofer Institute IAIS, Sankt-Augustin, Germany*

^c*City, University of London, London, UK*

^d*Zhejiang University, Zhejiang, China*

^e*Tongji University, Shanghai, China*

^f*Airbnb, San Francisco, CA, USA*

^g*National Chiao-Tung University, Hsinchu, Taiwan*

^h*University of Sydney, Sydney, Australia*

Abstract

Data quality management, especially data cleansing, has been extensively studied for many years in the areas of data management and visual analytics. In the paper, we first review and explore the relevant work from the research areas of data management, visual analytics and human computer interaction. Then for different types of data such as multimedia data, textual data, trajectory data, and graph data, we summarize the common methods for improving data quality by leveraging data cleansing techniques at different analysis stages. Based on a thorough analysis, we propose a general visual analytics framework for interactively cleansing data. Finally, the challenges and opportunities are analyzed and discussed in the context of data and humans.

Keywords: Data quality management, visual analytics, data cleansing

2010 MSC: 00-01, 99-00

*Corresponding author

Email address: shixia@tsinghua.edu.cn (Shixia Liu)

1. Introduction

With the increasing predominance of data-centric approaches to business, scientific, and engineering problems, data and its quality have become more and more important [1, 2, 3, 4]. However, during the data collection and processing
5 stage, some incomplete, inconsistent, duplicate, or inaccurate data can be fused, which usually affects the further usage of data (e.g. lower the accuracy of the learning models) and can lead to loss of consumer trust and revenues. As a result, in the era of big data, one key issue on preparing and processing data is to ensure the quality and usability of data (data quality management), including
10 detecting, removing, and correcting errors and inconsistencies in the data.

Data quality management has been studied for many years in the area of database and data management [5, 2]. Its major goal is to efficiently detect and correct errors in the data. As an important part of data-driven analysis, data quality management takes over 30%-80% of the time and resource [6].
15 Most mature works focus on tabular data, such as assessing the data quality [7], interactive data cleansing [8], and data wrangler [9]. However, due to the increasing complexity of the data (e.g., multimedia data, texts, graphs, sequences and trajectories etc.) collected through a variety of ways, it is more and more challenging to effectively and accurately improve data quality. In most of the
20 cases, domain knowledge of experts is important to guide for better performance of data quality management algorithms [10]. As a consequence, there has been a growing interest in recent years to study how to better combine user-guided methods with system-guided methods during the analysis, where information visualization and visual are the important parts to achieve this
25 goal [3, 4, 11, 12, 13].

Data cleansing is a widely used practice for effective data quality management. As a result, most existing data quality management efforts focus on data cleansing. In this paper, we first report the related work from different research areas, including data management, visual analytics, and human computer interaction. Then for different types of data, we summarize the common methods
30

for improving data quality by leveraging data cleansing techniques at different stages. In addition, a high level abstraction of a framework on designing a visual analytic system for data cleansing is needed as a general guideline for the research in this direction. Thus, inspired by pipeline proposed in [14], we develop
35 a visual analytics framework, focusing on iteratively and progressively improving data quality from the screening stage, to diagnosis stage and the correction stage. Finally, we explore the research challenges and opportunities and align them with the our visual analytics framework, which we hope can better guide the future visual analytics research on data cleansing.

40 2. Related Work

Researchers have been extensively studying a variety of data cleansing techniques to improve the data quality for the past twenty years. Most efforts are mainly from two research areas: data management and visual analytics.

In the area of data management, researchers have developed a number of
45 approaches to checking, repairing, and correcting inconsistencies and errors in the data. Existing efforts can be classified into three categories: rule-based detection methods for cleaning data based by a set of rules [15, 16, 17, 18], quantitative error detection methods for discovering and resolving outliers and glitches in the data [19, 20, 21, 22], and record linkage and de-duplication meth-
50 ods for detecting duplicate data items [23, 24]. Recently, Abedjan et al. [25] conducted a comprehensive evaluation to analyze the performance of existing algorithms on four common types of data errors, including outliers, duplicates, rule violations, and pattern violations. These error types are relatively general and can be applied beyond tabular data. However, these works do not pro-
55 vide an end-to-end data cleansing pipeline. In order to enable effective and efficient data cleansing practices, several frameworks have been developed. For example, Florescu and modeled the cleansing application as a graph-based data transformation, which can be applied to SQL-based database [26]. Gill and Lee developed a distributed data cleansing framework specific for data streams [27].

60 Due to the limited usage scope, these frameworks cannot be applied to a general data cleansing application. Broeck et. al., proposed a three-stage framework on data cleansing, either manually or automatically. The framework categorizes the whole process into three stages, screen stage, diagnosis stage, and correction stage [14]. For each stage, the key problems are identified. The major feature
65 of this framework is that it covers the whole analysis workflow well, from the raw data exploration to actual error correcting.

In most real-world cases, the data cleansing process cannot be done fully automatically due to the ambiguity of the errors and the need of human knowledge to verify the cleansing results. To effectively loop human into the data cleansing
70 ing process, visual analytics researchers have developed several works focusing on interactive data cleansing. These works in most cases aim to solve some specific tasks for certain types of data, mostly structured table representations. Krishnan et. al. [28] designed ActiveClean to interactively cleanse the data for statistical modeling. Profiler [7] was designed to interactively detect and vi-
75 sually summarize the outliers from data. von Zernichow et al. [29] presented a prototype system for visual data profiling, with a focus on discovering and correcting missing values and outliers in tabular data. Wrangler [30] targets at interactively creating data transformation scripts. Guo et al. [31] later extended Wrangler to a mixed-initiative system by integrating a proactive recommenda-
80 tion model, which suggests applicable data transforms to users for a more guided exploration of the transformation space. Both tools (Profiler and Wrangler) only support tabular data cleansing. Beyond these works, Kandel et. al. [9] further summarized the research direction on how visualizations and interaction techniques can help data wrangling. The aforementioned works greatly demonstrate
85 the usefulness and effectiveness of visual analytics techniques in helping improve the data quality, however, it is not easy to apply these techniques to other types of data or different cleansing tasks.

While tabular data is the predominant focus of data cleansing researchers, we also note some efforts on visually cleaning time-oriented data. Gschwandt-
90 ner et al. [13] derived a taxonomy of quality issues with time-oriented data, and

envision the need of visualization tools for analyzing the quality issues with human in the loop. Following this line, they later proposed TimeCleanser [32], an interactive approach specifically for cleansing time-oriented data. The approach includes several syntax checks that are common for tabular data, including time
95 checks (valid temporal range, consistent interval lengths, missing time points or intervals), time-oriented value checks (e.g. identifying values that don't change for a long time), and consistency about multiple data sets (same temporal range, resolution etc.) Though TimeCleanser demonstrates its effectiveness on correcting data, it is less flexible to support reasoning about underlying causes to the
100 quality issues detected. In light of this, Arbesser et al. [33] designed Visplause, an interactive visualization system to facilitate the inspection of the quality of multiple time series. In particular, the system mainly integrates the meta information of data to provide a hierarchical overview that aggregates the results of data quality checks at different levels of detail. This enables a flexible
105 semantic reasoning about the data quality. Gschwandtner and Erhart [34] presented "Know Your Enemy", a visual analytics approach to the identification of issues in time series data. Compared with prior work, they more adequately tailored the visualization design to the time-oriented characteristics of data (e.g., providing the temporal contexts). We also note some recent works that focus
110 on cleaning the time-oriented data that has specialized characteristics in some particular application domains. For example, Schulz et al. [35] proposed a visual cleansing tool tailored to the low-level eye-tracking data. Dixit et al. [36] develop an interactive approach specifically for correcting event orderings in process logs.

115 3. Data Types and Their Relationships

The majority of the existing data cleansing tools focus on structured table data. The word 'structured' is the key here, as the structure of the table suggests how to assess and improve the quality of the data. Table data representation assumes that values in each column have the similar formatting.

120 Respectively, data cleansing procedures check data types, formatting, and help
users to perform basic data cleaning operations such as fixing typos (e.g. by
checking spelling of text values or ensuring consistency in using decimal dots
and commas in numeric values), modifying representations of dates and times,
computing statistics for numeric values, identifying and interactively inspecting
125 outliers, resolving encoding rules (e.g. 999 for *no_data*) etc. Special methods
exist for detecting missing values and replacing them by plausible values taking
into account values in other table columns. Modern tools like Tableau sup-
port cleansing of multiple related tables. For example, values in the connecting
columns of the two tables can be checked for consistency and, if some values do
130 not match, they can be modified interactively.

Non-table data require special methods for cleansing that take into account
the specifics of the data type. For example, preparation of textual data requires
language detection, removal of too short documents, detection of duplicated
content, fixing miss-spelled words and incorrect punctuation, just to name a few
135 operations. It is important to have a possibility to perform fast computations for
subset of the data (e.g. calculating the size of the vocabulary used in different
document or assessing the overall sentiment) for looking at the data set from
different perspectives and thus selecting appropriate documents for analysis.

For more complex data it is necessary to exploit their structures for develop-
140 ing appropriate data wrangling tools. For example, spatial event data [37] may
be considered as a special type of table data that include columns with temporal
and spatial references. Each event is described by its type (*‘what happened’*),
location (*‘where happened’*), time (*‘when happened’*) and attributes (e.g. event
magnitude, actors, etc.) Event types and attributes can be treated as regu-
145 lar columns, applying traditional methods for validating their quality. Spatial
and temporal references of the events can be used for attaching additional data
sources (that describe locations and times) for checking if events are plausible
at given locations and times. For example, traffic events are expected to happen
on roads but not in lakes, and visiting restaurant events are highly unlikely at
150 night times.

Another representative example is the graph data. Generally, a graph can be represented using two connected tables describing graph vertices and edges. Both vertices and edges can be described by their attributes. Moreover, additional attributes of vertices and edges can be computed from the topology
155 of the graph. Graph centrality measures can be used for identifying outliers and disconnected subgraphs. Respectively, the data cleansing tools for graphs need to take into account the structure of the graph data and potentially useful additional information that can be derived from the graph data.

Understanding the specifics and structure of the data is essential for design-
160 ing appropriate data transformations. Let us consider an example of trajectories of moving objects [38]. Each position is a spatial event that can be described by a reference to the moving object *id*, time stamp *t*, its coordinates *x* and *y*, and possible attributes: *id*, *t*, *x*, *y*, *attributes*. Sequences of events for the same moving object can be integrated into a trajectory (Fig. 1). Such inte-
165 gration allows computation of derived attributes such as displacement, time difference, speed estimate etc. These derived attributes can be used for extracting secondary events from trajectories (e.g. stops) and dividing trajectories into smaller subsets (e.g. trips between stops). Both trajectories and events can be aggregated by areas and links between areas, creating spatial time series re-
170 ferring to locations and links between them, respectively. Further events (e.g. extreme values) can be extracted from such derived spatial time series. This example demonstrates how data structure defines possible and potentially useful transformations. Such transformations can be used for looking at the data from a different perspective, thus facilitating the data cleansing process, identifying
175 and eventually fixing data problems.

4. Examples of different data types

In this section we consider the specifics of steering data quality for different types of data beyond regular data tables.

Multimedia. Multimedia covers many different types of data, which had

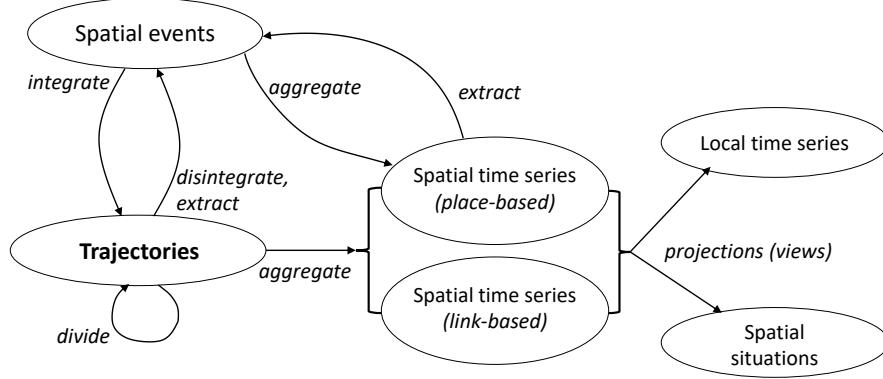


Figure 1: Potentially possible transformations of trajectory data.

180 been widely used for communication after cameras and recorders were invented.
 In this study, we focus on images, audios, and videos because they are the most
 popular. Typically, a pixel in images represent the color of a small area. While
 many pixels are arranged to form an image, which describe the appearance of a
 large area, they can be used to convey visual information. Compared to an im-
 185 age, a video has an additional time coordinate, in which each time span contains
 an image. Therefore, a video can show dynamic visual changes over time. In
 terms of audios, each sample in a time span describes the frequency and magni-
 tude of energy that would be heard by humans. Note that the most important
 characteristic of the above-mentioned data is that each sample/pixel is mean-
 190 ingless, whereas the integration of them is not. Therefore, there are no point
 anomalies, but contextual and collective anomalies in multimedia data. Cleans-
 ing abnormal multimedia data demands semantics and can only be achieved by
 complex algorithms and human guidance.

Textual data. Textual data consists of bag-of-words representation in the
 195 form of sentences, paragraphs, documents, and topics. Typical examples such
 as news articles, interviews, emails, field notes, as well as text descriptions
 from feeds and social media sources, are essential to communicating informa-

tion about people, events, and activities, sharing findings, knowledge, and best practices within a community, as well as to connecting people and propelling the organization forward [39].

In the era of big data, with the ever-increasing size of a document corpus, it is simply not possible in many cases for people to quickly locate key information or derive insights from such large amounts of textual data. Generally, text analysis tasks include information retrieval, cluster/topic analysis, natural language process, classification, outlier analysis, etc. [40].

To better perform the aforementioned tasks, the quality of textual data must be guaranteed. Examples of quality issues of textual data are duplicate documents that are intentionally rephrased, irrelevant content such as an advertisement in a news article, documents with mixed topics that are hard to disentangle for a specific application (e.g., information retrieval), multilingual documents, as well as inconsistent documents with varying document lengths, diverse writing styles, and different writing quality (e.g., news articles and tweets).

Trajectories of moving objects. Trajectory data is a sequence of events that correspond to recorded time-referenced positions of moving objects [38]. While each particular record has limited value, many applications require consideration of large collections of such position records. However, even simple tasks like detection of duplicates become difficult due to the complex structure of the data. For position records, a duplicate is a repeating combination of the moving object identity and temporal reference. If the positions and attributes are equal, this is a case of duplicate records that need to be eliminated. If positions or attributes differ, sophisticated conflict resolution procedures are needed.

Paper [41] analyses in detail the properties of trajectory data (properties of moving objects, space and time) and related properties of data collection procedures. These properties are used for identifying potential problems that may appear in trajectory data sets: missing data, accuracy problems and precision deficiency. The potential problems are considered in respect to all components of the data structure (moving objects, space and time). In addition, derived

attributes (e.g. speed or acceleration) and possible transformations (e.g. aggregating trajectories into occupancy indicators for visited areas and counts of moves between the areas) enable looking at data from different perspectives (see Fig. 1).

Graph data. A graph $G = (V, E)$ consists of a vertex set V and an edge set E , where a vertex represents an entity and an edge between two vertices represents a relationship between them. Both vertices and edges can have multiple *attributes*, such as numbers, texts, categorical data and images. An edge set defines a *topology* structure of the graph, and there can be more than one edge sets for a given vertex set, which represent multiple relationships between the entities.

For example, a vertex v can represent a student, where v has multiple attributes such as student id number, enrolled degree, year, home address, phone number, and photo id. Similarly, an edge e can have numbers such as weights, texts, time stamps and directions. For more complex data, an edge set E can be changing over time, i.e., change the values of their attributes as well as the topology of a graph. The student set can have multiple relationships, such as facebook friends, instagram friends, and tennis friends etc.

Therefore, the data cleansing for graph data need to take into account these attributes for vertices and edges, as well as the topology of the graph, defined by the edge set. For example, based on the types of attributes of vertices and edges (such as numbers, texts and images) and the topology of the graph, one need to use a variety of techniques for processing missing values, duplicates, uncertainty, imprecision and conflicts, as well as detection for outliers and specific patterns.

5. Analysis Pipeline

Based on the *Screening* \rightarrow *Diagnosis* \rightarrow *Correction* framework [14], we propose a visual analytics framework for analyzing and improving data quality. The goal of the proposed framework is to help a user (e.g., an analyzer) find the potential problems of the data to be analyzed/visualized and provide efficient

layer. Through screening, a user is able to choose to summarize and illustrate
275 the overview, statistic features, and data patterns (e.g., trend and cluster) via
intuitive visual representations. After that, the user can further make a di-
agnosis on the data to find out the potential problems (e.g., missing values,
duplications, pattern/constraint violations, inconsistencies) that effect the data
quality. Finally, a user can interactively correct the detected problems within
280 the data.

In the above analysis procedure, visualization plays an important role in sup-
porting data interpretation and decision making. In the framework, two types of
visualization designs are required: (1) the visualization designed for illustrating
and summarizing the data with the goal of showing overview, patterns, distri-
285 butions, and constrains of the data, thus handling the complexities associated
with data; (2) the visualization designed for data error correction with the goal
of identifying missing data, outliers, duplicates, pattern/constraint violations,
and data inconsistencies, which aims at tackling a variety of complexities related
to human.

290 6. Research Challenges and Opportunities

6.1. Data Complexity

Multimedia. Quality management for multimedia is challenging. Users have to
apply different methods to handle images, audio recordings, and videos because
they are of different nature. In addition, caused by different reasons, anomalies
295 of these media can be classified into low level and high level categories. On one
hand, low level anomalies are often caused by data transmission, compression, or
fault of a machine. These anomalies, such as blank images, white noise audios,
duplicated frames in a video, can be detected by simple rules. However, on
the other hand, high level anomalies are difficult to detect because of requiring
300 semantics. Specifically, each small segment of multimedia data, such as a pixel
in images or a sample in audios, is meaningless, but the integration of them is
not. To determine whether an image or a video is anomalous, users have to apply

image processing methods to assess visual quality, or object detection techniques to obtain semantics. Otherwise, anomalies such as overexposure, underexposure, and serious object occlusions cannot be detected. Similarly, speech recognition and emotion identification are often used to understand semantics in an audio, and signal processing techniques are used to identify low quality audios such as containing load but irrelevant background speech.

Textual data. Although textual data is widely used in many lines of work, data quality problems for such type of unstructured data remain largely unexplored. This is because, due to the unstructured nature of textual documents, quality management for textual data is challenging. First, textual data often contains several data fields and mixes the useful information with irrelevant information. As a result, one key challenge is how to retrieve interactively the useful content and remove the noisy information. For example, a web page is usually a mixture of many types of information, such as main textual content, advertising panels, navigation bars, copyright blocks, images, etc. In real-world applications, only part of information, typically the main textual content, is useful and the rest are treated as noise. Therefore, it is important to remove the irrelevant information, which is still a hot research topic in the area of information retrieval. Second, a text corpora may contain text strings of different distributions, such as different lengths and language usages. For example, news articles and formal publications are usually long and consist of sentences with grammatical rules, while tweets or micro-blogs are short, noisy, and with limited context information. This makes it impractical to use a unified text mining model to analyze them together. As a result, another challenge is how to effectively improve the quality of a text corpora with inconsistent data distributions.

Trajectory data. Quality management for trajectories is challenging. Users need to understand the nature of the data and of the phenomenon they represent for assessing the data quality properly and, subsequently, validating the data correctly. Concerning the phenomenon, it is necessary to distinguish different modes of movement (e.g. separate walking from biking or using public transportation, and then use different constraints concerning the physics

of movement: speed, acceleration, inertia) and take into account the context
335 of movement (e.g. multiple cars on a road must follow common direction).
Concerning the data, it is necessary to take into account the data collection
procedure (e.g. positions collected every minute, every 20m during straight
movement, by performing certain activities such as making a call).

It is necessary to understand the coverage properties of a data set under
340 inspection, ensuring that it corresponds to analysis tasks. Often data sets are
limited in spatial extent, causing complete trajectories or their parts being ab-
sent in data. Sometimes data collection is impossible in specific conditions (e.g.
positioning device does not work in tunnels or indoor). Proper temporal cov-
erage is essential, too. Another coverage aspect that needs to be considered is
345 population: for example, it is dangerous to make conclusions about mobility of
elderly people based on data derived from positions of social media activities
that are performed mostly by youngsters. Another example in vehicle traffic:
projecting mobility patterns of public buses onto individual cars is doubtful.

Graph data. Similarly to other data types, quality management for graph
350 data is challenging due to the complexity of the data. In addition to the chal-
lenges for various types of attributes, the topology structure of a graph adds
more challenges. Users need to apply a variety of techniques for missing val-
ues, duplicates, uncertainty, and outlier detection to handle different types of
attributes (such as numbers, texts and images) and the topology of the graph.
355 Furthermore, multiple relationships and relationships changer over time (i.e.,
dynamic graphs) add more challenges.

For example, users may need to compute statistics about the topology of
a graph, such as the density, diameter, clustering coefficient, connectivity and
average neighbour degree, as well as the property testing on the topology struc-
360 ture of a graph, such as testing whether a given graph is a tree (i.e., no cycle),
a planar graph (i.e., can be drawn in the plane without edge crossings), or a
Directed Acyclic Graph (DAG).

In addition to the attribute-based outlier detections, users need to perform
topology-based pattern detections. Examples include finding high frequency

365 patterns such as *motifs*, a small subgraph consists of three or four vertices, small cycle patterns such as *triads* (i.e., triangles), and other special subgraphs such as paths, trees, stars, and complete subgraphs. Algorithms for finding such special topological patterns are complex with high runtime complexity.

Users may need to analyze topology-based constraints. For example, a tree
370 has $n - 1$ edges, and a planar graph can have at most $3n - 6$ edges, where n is the number of vertices. Others include domain specific constraints; for example a family tree should not have a cycle, which represents a blood marriage.

Common challenges and opportunities. In addition to the aforementioned data-type-specific research challenges and opportunities, there are also some
375 common ones. First, existing visual data cleansing methods cannot be scalable to large-scale datasets. One potential solution to handle large-scale datasets is to sample only a small subset of the whole training set. The challenge here is how to develop effective sampling methods that can both keep the data density and preserve important data such as influential points, outliers, and exceptions.
380 Second, there is a lack of effective quality metrics to measure the quality of different types of data such as textual data, images, videos, graph data, and trajectory data. As a result, a potential research opportunity is to develop quality metrics from data content and evaluate them within specific usage contexts. In real-world applications, the analyst often needs to examine multiple
385 types of data and correct the errors among them. Accordingly, the third challenge is designing an integrated interface to visually illustrate the distributions of different types of data.

6.2. Human Complexity

Several challenges will arise when human intelligence is integrated into auto-
390 matic data cleaning pipeline. We classify these challenges into three categories and identify the associated opportunities as follows.

- **Lack of domain knowledge.** Better integration of human domain knowledge plays an important role in steering data quality. However,

395 sufficient knowledge or expertise regarding new types of data or new data
sets might not be always available. To overcome such a challenge, it is
important to explore how to tackle integration and calibration of insuffi-
cient or incomplete knowledge and expertise. One important direction is
to design progressive or collaborative visual interfaces that enable crowd-
sourcing, such that users without sufficient knowledge can gradually gain
400 more knowledge or seek support of other users with sufficient knowledge.
It would also be interesting to create a visual recommendation mechanism
for providing necessary automation in data cleaning, especially when users
lack domain knowledge.

- **Limitations of perception/cognition.** Prior psychological studies have
405 revealed limitations on visual perception and cognition, such as restricted
field of view in perception [42] and limited working memory in cogni-
tion [43]. These limitations can directly influence the way in which human
perceive and understand the world. For complicated types of data, it is
challenging to design a visual analytics system while keeping the complex-
410 ity of the system within the limitations. One worthy research direction
is to explore a mixed initiative mechanism which seamlessly integrates
system initiative guidance and user initiative guidance for better human
machine intelligence, such that perception or cognition limitations of users
can be largely addressed.

- **Difficulty in understanding uncertainty and its implications.** Un-
415 certainty might arise in any stage of a data cleaning process, and prop-
agate in subsequent stages [44]. Misunderstanding the uncertainty and
its implications would result in erroneous decisions and low-quality data.
However, understanding the uncertainty and its implications would be
generally difficult without a proper visual guidance. To circumvent the
420 problem, it is highly necessary to model and visualize the uncertainty in
data cleaning, such that users can make informed decisions during the
data cleaning process.

7. Conclusion

425 Data quality is of crucial importance to a wide variety of real-world applications. In this paper, we review and summarize research efforts on steering data quality, with a focus on data cleansing, a widely-used technique for effective data quality management. First, we summarize the relevant work from different research fields, including data management, visual analytics and human computer interaction. 430 Then, for different types of data, we discuss the common methods for improving data quality by leveraging data cleansing techniques. Building upon the existing analysis pipeline of data cleansing by Van den Broect et al.[14], we further propose a visual analytics framework for iteratively and progressively improving data quality from the screening, diagnosis 435 and correction stage. Finally, we analyze the research challenges and opportunities in the context of data and human complexities, which we believe are critical for future research on visual data cleansing.

References

- [1] W. Fan, F. Geerts, Foundations of data quality management, Morgan & 440 Claypool Publishers, 2012.
- [2] S. Liu, C. Chen, Y. Lu, F. Ouyang, B. Wang, An interactive method to improve crowdsourced annotations, IEEE transactions on visualization and computer graphics (accepted) 25 (1).
- [3] N. McCurdy, J. Gerdes, M. Meyer, A framework for externalizing implicit 445 error using visualization, IEEE Transactions on Visualization and Computer Graphics (accepted) 25 (1).
- [4] H. Song, D. A. Szafr, Where’s my data? evaluating visualizations with missing data, IEEE Transactions on Visualization and Computer Graphics (accepted) 25 (1).

- 450 [5] O. Kwon, N. Lee, B. Shin, Data quality management, data usage experience and acquisition intention of big data analytics, *International Journal of Information Management* 34 (3) (2014) 387–394.
- [6] B. Saha, D. Srivastava, Data quality: The other face of big data, in: *IEEE International Conference on Data Engineering*, 2014, pp. 1294–1297.
- 455 [7] S. Kandel, R. Parikh, A. Paepcke, J. M. Hellerstein, J. Heer, Profiler: Integrated statistical analysis and visualization for data quality assessment, in: *Proceedings of the International Working Conference on Advanced Visual Interfaces*, 2012, pp. 547–554.
- [8] V. Raman, J. M. Hellerstein, Potter’s wheel: An interactive data cleaning
460 system, in: *Very Large Database Conference*, Vol. 1, 2001, pp. 381–390.
- [9] S. Kandel, J. Heer, C. Plaisant, J. Kennedy, F. van Ham, N. H. Riche, C. Weaver, B. Lee, D. Brodbeck, P. Buono, Research directions in data wrangling: Visualizations and transformations for usable and credible data, *Information Visualization* 10 (4) (2011) 271–288.
- 465 [10] N. El Bekri, E. Peinsipp-Byma, Assuring data quality by placing the user in the loop, in: *International Conference on Computational Science and Computational Intelligence*, 2016, pp. 468–471.
- [11] S. Liu, X. Wang, M. Liu, J. Zhu, Towards better analysis of machine learning models: A visual analytics perspective, *Visual Informatics* 1 (1) (2017)
470 48–56.
- [12] J. Choo, S. Liu, Visual analytics for explainable deep learning, *IEEE Computer Graphics and Applications* 38 (4) (2018) 84–92.
- [13] T. Gschwandtner, J. Gärtner, W. Aigner, S. Miksch, A taxonomy of dirty time-oriented data, in: *International Conference on Availability, Reliability, and Security*, 2012, pp. 58–72.
475

- [14] J. Van den Broeck, S. A. Cunningham, R. Eeckels, K. Herbst, Data cleaning: detecting, diagnosing, and editing data abnormalities, *Public Library of Science, medicine* 2 (10) (2005) e267.
- 480 [15] Z. Abedjan, C. G. Akcora, M. Ouzzani, P. Papotti, M. Stonebraker, Temporal rules discovery for web data cleaning, *Proceedings of the Very Large Database Endowment* 9 (4) (2015) 336–347.
- [16] W. Fan, J. Li, S. Ma, N. Tang, W. Yu, Towards certain fixes with editing rules and master data, *The Very Large Database Journal* 21 (2) (2012) 213–238.
- 485 [17] F. Geerts, G. Mecca, P. Papotti, D. Santoro, Mapping and cleaning, in: *IEEE International Conference on Data Engineering*, 2014, pp. 232–243.
- [18] Z. Khayyat, I. F. Ilyas, A. Jindal, S. Madden, M. Ouzzani, P. Papotti, J.-A. Quiané-Ruiz, N. Tang, S. Yin, Bigdancing: A system for big data cleansing, in: *ACM Special Interest Group on Management of Data*, 2015, pp. 1215–1230.
- 490 [19] T. Dasu, J. M. Loh, Statistical distortion: Consequences of data cleaning, *Proceedings of the Very Large Database Endowment* 5 (11) (2012) 1674–1683.
- [20] N. Prokoshyna, J. Szlichta, F. Chiang, R. J. Miller, D. Srivastava, Combining quantitative and logical data cleaning, *Proceedings of the Very Large Database Endowment* 9 (4) (2015) 300–311.
- 495 [21] M. Vartak, S. Rahman, S. Madden, A. Parameswaran, N. Polyzotis, SeeDB: efficient data-driven visualization recommendations to support visual analytics, *Proceedings of the Very Large Database Endowment* 8 (13) (2015) 2182–2193.
- 500 [22] E. Wu, S. Madden, Scorpion: Explaining away outliers in aggregate queries, *Proceedings of the Very Large Database Endowment* 6 (8) (2013) 553–564.

- [23] A. K. Elmagarmid, P. G. Ipeirotis, V. S. Verykios, Duplicate record detection: A survey, *IEEE Transactions on knowledge and data engineering* 19 (1) (2007) 1–16.
- [24] M. Stonebraker, D. Bruckner, I. F. Ilyas, G. Beskales, M. Cherniack, S. B. Zdonik, A. Pagan, S. Xu, Data curation at scale: The data tamer system, in: *Conference on Innovative Data Systems Research*, 2013.
- [25] Z. Abedjan, X. Chu, D. Deng, R. C. Fernandez, I. F. Ilyas, M. Ouzzani, P. Papotti, M. Stonebraker, N. Tang, Detecting data errors: Where are we and what needs to be done?, *Proceedings of the Very Large Database Endowment* 9 (12) (2016) 993–1004.
- [26] D. Florescuand, An extensible framework for data cleaning, in: *IEEE International Conference on Data Engineering*, 2000, pp. 312–312.
- [27] S. Gill, B. Lee, A framework for distributed cleaning of data streams, *Procedia Computer Science* 52 (2015) 1186–1191.
- [28] S. Krishnan, J. Wang, E. Wu, M. J. Franklin, K. Goldberg, ActiveClean: interactive data cleaning for statistical modeling, *Proceedings of the Very Large Database Endowment* 9 (12) (2016) 948–959.
- [29] B. M. von Zernichow, D. Roman, Usability of visual data profiling in data cleaning and transformation, in: *OTM Confederated International Conferences” On the Move to Meaningful Internet Systems”*, 2017, pp. 480–496.
- [30] S. Kandel, A. Paepcke, J. Hellerstein, J. Heer, Wrangler: Interactive visual specification of data transformation scripts, in: *ACM Special Interest Group on ComputerHuman Interaction*, 2011, pp. 3363–3372.
- [31] P. J. Guo, S. Kandel, J. M. Hellerstein, J. Heer, Proactive wrangling: mixed-initiative end-user programming of data transformation scripts, in: *ACM Symposium on User Interface Software and Technology*, 2011, pp. 65–74.

- 530 [32] T. Gschwandtner, W. Aigner, S. Miksch, J. Gärtner, S. Kriglstein, M. Pohl, N. Suchy, TimeCleanser: A visual analytics approach for data cleansing of time-oriented data, in: International Conference on Knowledge Technologies and Data-driven Business, 2014, pp. 18:1–18:8.
- [33] C. Arbesser, F. Spechtenhauser, T. Mühlbacher, H. Piringer, Visplause: 535 Visual data quality assessment of many time series using plausibility checks, IEEE Transactions on Visualization and Computer Graphics 23 (1) (2017) 641–650.
- [34] T. Gschwandtner, O. Erhart, Know your enemy: Identifying quality problems of time series data, in: IEEE Pacific Visualization Symposium, 2018, 540 pp. 205–214.
- [35] C. Schulz, M. Burch, F. Beck, D. Weiskopf, Visual data cleansing of low-level eye-tracking data, in: Eye Tracking and Visualization, 2015, pp. 199–216.
- [36] P. M. Dixit, S. Suriadi, R. Andrews, M. T. Wynn, A. H. ter Hofstede, J. C. 545 Buijs, W. M. van der Aalst, Detection and interactive repair of event ordering imperfection in process logs, in: International Conference on Advanced Information Systems Engineering, Springer, 2018, pp. 274–290.
- [37] N. Andrienko, G. Andrienko, Exploratory Analysis of Spatial and Temporal Data: A Systematic Approach, Springer-Verlag, Berlin, Heidelberg, 2005.
- 550 [38] G. Andrienko, N. Andrienko, P. Bak, D. Keim, S. Wrobel, Visual Analytics of Movement, Springer Publishing Company, Incorporated, 2013.
- [39] S. Liu, M. X. Zhou, S. Pan, Y. Song, W. Qian, W. Cai, X. Lian, TIARA: Interactive, topic-based visual text summarization and analysis, ACM Transactions on Intelligent Systems and Technology 3 (2) (2012) 25.
- 555 [40] S. Liu, X. Wang, C. Collins, W. Dou, F. Ouyang, M. El-Assady, L. Jiang, D. Keim, Bridging text visualization and mining: A task-driven survey, IEEE Transactions on Visualization and Computer Graphics (accepted).

- [41] G. Andrienko, N. Andrienko, G. Fuchs, Understanding movement data quality, *Journal of Location Based Services* 10 (1) (2016) 31–46.
- 560 [42] S. H. Creem-Regehr, P. Willemsen, A. A. Gooch, W. B. Thompson, S. H. Creem-Regehr, P. Willemsen, A. A. Gooch, W. B. Thompson, The influence of restricted viewing conditions on egocentric distance perception: Implications for real and virtual indoor environments, *Perception* 34 (2) (2005) 191–204.
- 565 [43] A. Baddeley, Working memory: looking back and looking forward, *Nature Reviews Neuroscience* 4 (2003) 829–839.
- [44] Y. Wu, G.-X. Yuan, K.-L. Ma, Visualizing flow of uncertainty through analytical processes, *IEEE Transactions on Visualization and Computer Graphics* 18 (12) (2012) 2526–2535.