



Data Smells: Categories, Causes and Consequences, and Detection of Suspicious Data in AI-based Systems

Harald Foidl
University of Innsbruck
Austria
harald.foidl@uibk.ac.at

Michael Felderer
University of Innsbruck
Austria
Blekinge Institute of Technology
Sweden
michael.felderer@uibk.ac.at

Rudolf Ramler
Software Competence Center
Hagenberg GmbH
Austria
rudolf.ramler@scch.at

ABSTRACT

High data quality is fundamental for today's AI-based systems. However, although data quality has been an object of research for decades, there is a clear lack of research on *potential* data quality issues (e.g., ambiguous, extraneous values). These kinds of issues are latent in nature and thus often not obvious. Nevertheless, they can be associated with an increased risk of future problems in AI-based systems (e.g., technical debt, data-induced faults). As a counterpart to code smells in software engineering, we refer to such issues as *Data Smells*. This article conceptualizes data smells and elaborates on their causes, consequences, detection, and use in the context of AI-based systems. In addition, a catalogue of 36 data smells divided into three categories (i.e., Believability Smells, Understandability Smells, Consistency Smells) is presented. Moreover, the article outlines tool support for detecting data smells and presents the result of an initial smell detection on more than 240 real-world datasets.

ACM Reference Format:

Harald Foidl, Michael Felderer, and Rudolf Ramler. 2022. Data Smells: Categories, Causes and Consequences, and Detection of Suspicious Data in AI-based Systems. In *1st Conference on AI Engineering - Software Engineering for AI (CAIN'22)*, May 16–24, 2022, Pittsburgh, PA, USA. ACM, New York, NY, USA, 11 pages. <https://doi.org/10.1145/3522664.3528590>

1 INTRODUCTION

Applications based on artificial intelligence (AI) (e.g., automated driving, predictive maintenance) have grown in popularity over the past decade. However, the resulting AI-based systems pose several challenges [7, 38]. One of these challenges is their strong data dependency [52]. This dependency is caused by data-hungry machine learning (ML) algorithms, typically used in AI-based systems to make intelligent decisions automatically. As a result, poor quality data can lead to abnormal behaviour and false decisions in such systems, resulting in huge monetary losses or, in the worst case, even harming people [8].

Recent research (e.g., [28, 49]), however, suggests that data quality problems are pervasive in AI-based systems. Data quality issues

even became one of the main reasons why they suffer badly from technical debt [6].

To improve this situation and thus meet the demand for high data quality in the context of AI-based systems, research in the area of data validation has recently gained significant interest (e.g., [5, 9, 11, 37, 47]). To reliably detect data issues, data validation methods generally require some context-specific constraints, which are usually defined by schemes or rules [1]. However, this declarative and context-dependent nature of data validation combined with the constantly growing amount of data makes data validation a tedious and labor-intensive activity [17, 52].

To address this issue, we previously proposed using *potential* data issues as indicators of latent data quality problems to guide the data validation process in a risk-driven way [17]. These potential data problems are usually indicated by context-independent, suspicious data values, patterns, and representations, highlighting data that should be prioritised during the validation. We referred to these potential data issues as *data smells* by analogy with code smells.

In the field of software engineering, code smells have become established as indicators of bad design and programming practices that are not faults per se but increase the likelihood of introducing faults in the future [19, 60]. In that sense, code smells are *potential* faults or issues [50].

We assert that data smells share several characteristics with code smells. For example, they typically arise due to violated best practices in data handling (e.g., wrong sequence of operations) or data management (e.g., missing data catalogue). Further, they can lead to interpretation issues of software components and impede the evolution and maintenance of AI-based systems. Therefore, we claim that data smells contribute massively to the emergence of technical debt and data-induced bugs within AI-based systems and are thus an increasingly important area of research.

However, although research on data issues is quite mature, it only partly considered such potential data issues. In fact, previous studies (e.g., [31, 33, 35]) mainly focused on actual data errors, lacking a precise definition or categorization of such potential issues.

This paper aims to address this gap by providing a solid foundation of data smells. In detail, we describe the characteristics of data smells and outline their potential causes, consequences, and use in the context of AI-based systems. We further present a catalogue comprising 36 data smells divided into three categories (Believability, Understandability, and Consistency Data Smells). Moreover, the detection of data smells is discussed and corresponding tool support is presented. Although we consider data smells in the context of AI-based systems in this article, the concept presented is

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for components of this work owned by others than the author(s) must be honored. Abstracting with credit is permitted. To copy otherwise, or republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee. Request permissions from permissions@acm.org.

CAIN'22, May 16–24, 2022, Pittsburgh, PA, USA

© 2022 Copyright held by the owner/author(s). Publication rights licensed to ACM.

ACM ISBN 978-1-4503-9275-4/22/05...\$15.00

<https://doi.org/10.1145/3522664.3528590>

also partially applicable to other data-driven initiatives (e.g., data mining projects).

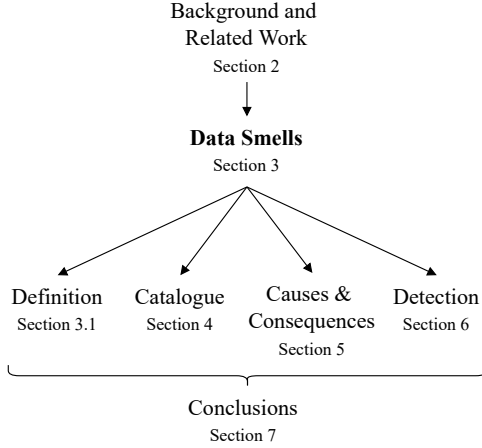


Figure 1: Structure of the article

The remaining article is structured as outlined in Figure 1. Section 2 provides background information on the term data smell and discusses related work. The concept of data smells is described in Section 3. Section 4 presents the data smell catalogue and Section 5 details the causes and consequences of data smells. Afterwards, Section 6 deals with detecting data smells, thereby presenting detection approaches, corresponding tool support including an experimental evaluation, and use cases. Finally, Section 7 concludes the paper.

2 BACKGROUND

This section first provides a brief overview of the origin and preliminary work on the term data smell in Section 2.1. Afterwards, Section 2.2 presents related work that deals with potential data issues and thus reflects our understanding of data smells.

2.1 Data Smells

The term "data smells" was first mentioned in the grey literature by Harris in 2014 [25]. In his article, Harris emphasized the importance of critically examining and questioning data before drawing results and conclusions from them. To this end and based on his software engineering background, he introduced the term data smells for the field of data analysis. Harris mentioned large standard deviations or double-counted records as concrete examples of data smells. However, Harris provided neither a precise definition nor an extensive list of data smells.

In the same year, Iubel [29] picked up the term and presented 13 smells for the domain of data journalism on GitHub. However, compared to our work, the presented smells lack a sound categorization and tangible definitions. Nevertheless, the smells Iubel proposed in the category of "unusable data" are generally applicable and not limited to the field of journalism.

Referring to the academic literature, there is one contribution besides our previous work [17], which used the term data smell. Sharma et al. [55] introduced the term data smell beside schema and

query smell as one kind of database smell. However, they provide a rather vague definition and only illustrate the concept with a single example of a data smell.

2.2 Related Work

There is a considerable amount of publications addressing issues in data. However, there are plenty of different terms used to refer to data issues, such as dirty data [33, 35], data error [1], data defect [31], data anomaly [18, 57] or data quality problem [41]. Following, we briefly discuss the most relevant related work that addresses potential data issues, albeit under different terms. See [4, 27] for an overview of the more general data quality and cleaning concepts.

The most relevant contribution regarding our work was published by Hynes et al. [26]. They presented a tool named data linter that aims to detect potential issues in ML training data features. The authors introduced the term data lint as a data- and model-specific inspection rule which identifies data feature representations that are suboptimal to specific ML models. Our work differs from their contribution in the following ways. First, unlike data lints, the data smells proposed in this paper describe a more comprehensive concept that also includes the causes and consequences of data quality issues. Second, data smells are generally applicable, while most of the proposed data lints are tailored to ML feature data. Third, we present 36 data smells in contrast to 13 lints proposed by Hynes et al. However, six of the 36 data smells presented are also detectable with the lint detection tool developed by Hynes et al.

Another field of related work worth to be mentioned has focused on detecting data issues in spreadsheets. Barowy et al. [3] presented a technique to investigate whether a data value is erroneous or not. The presented technique is based on the rationale that it is often impossible to know whether a data value has an issue or not. Thus, they propose investigating the impact of a data value on the computation result in spreadsheets. If there is an unusual effect on the result, the data may be erroneous and merit special attention. Similarly, Cunha et al. [30] proposed a catalogue of spreadsheet smells as indicators of possible issues in spreadsheet data. Although both articles share our understanding of data smells, many of their smells focus on spreadsheet characteristics (e.g., Reference to Empty Cells smell, Standard Deviation smell) and therefore only have limited applicability to other applications.

Other contributions mainly focused on taxonomies of data issues. Although most taxonomies propose proper categorizations of data issues, only a few separate potential data issues in an own category. One of these taxonomies was proposed by Kim et al. [33], who grouped such potential data issues as "not missing and not wrong but unusable". Within this category, they list issues such as abbreviations or homonyms. Also, Josko and Wöß [31] mention "less severe" data defects in their data defect ontology. According to them, such data defects do not threaten the functionality of a system directly but use them to spread to the rest of the system.

A further work to be mentioned is that of Kasunic [32]. In their work, they define data anomalies to might be erroneous but also as might represent correct data that is caused by unusual but actual circumstances. Although they do not provide further details, it is very close to our understanding of potential data issues.

Sambasivan et al. [49] recently published a study that, although not explicitly addressing potential data issues, is also related to our work. Simply put, the authors introduced the metaphor of data cascades to frame "compounding events causing negative, downstream effects from data issues" in the context of high-stakes AI systems. Although there are similarities in the origin of data cascades and data smells, both emerge through practices that undervalue data quality, there are significant differences. Data smells describe potential instead of general data issues. Further, data smells are grounded on an engineering perspective, whereas data cascades (e.g., conflicting reward systems) represent a much more high-level concept.

In sum, the present state of research shows that the idea of potential data issues is not entirely new. However, a sound concept and generally applicable list of data smells is still missing but highly needed. This is precisely where this paper ties in.

3 DATA SMELLS

This section first defines data smells and describes their main characteristics in Section 3.1. Section 3.2 then discusses how they relate to other types of data quality issues.

3.1 Definition and Characteristics

In its broadest sense, data smells describe *latent* data quality issues. By latent we mean issues that are present and can cause problems, but are not obvious now. In detail, we define data smells as *context-independent, data value-based indications of latent data quality issues caused by poor practices that may lead to problems in the future*. Following, we describe the main characteristics of data smells concerning their *suspicion, context, origin and consequences*.

Moderate degree of suspicion. Data smells are indicated by moderately suspicious data values, patterns and representations that make them challenging to identify. For example, "New York" can be considered smelly because it is semantically unclear whether it refers to the city or the federal state.

Context-Independence. Data smells indicate suspicions that are not tied to a specific context. Therefore, they are widely applicable and represent a potential threat for any data-driven system.

Caused by poor practices. The emergence of data smells is usually caused by the violation of recommended best practices in data management, software or data engineering (e.g., skip testing of data handling logic). Similarly, data smells are likely to occur due to poor data generation, acquisition or processing (e.g., glue code). Also, a poor intrinsic quality of the data sources from which the data originate (e.g., redundancies in database schema) causes them to arise.

May cause problems in the future. Data smells increase the likelihood of future problems in data processing (e.g., type conversion errors, data misclassifications) and in the evolution of data-driven systems (e.g., enhancement, maintenance). Important in this regard is that the consequences of data smells are usually uncertain and typically delayed in terms of time and location in a data-driven system.

3.2 Differentiation

To relate data smells to other types of data quality issues, we use the classification of data quality issues proposed by Ge and Helfert [22]. According to them, data quality issues can be divided into context-dependent and context-independent issues. Whereas the latter are independent of the context of use (e.g., missing values, synonyms), the former usually refer to issues that violate contextual rules (e.g., 2.8 as number of siblings). As context-specific knowledge is typically needed to decide if a data value is erroneous or not, we use the term *data error* to refer to context-dependent issues in the remaining paper.

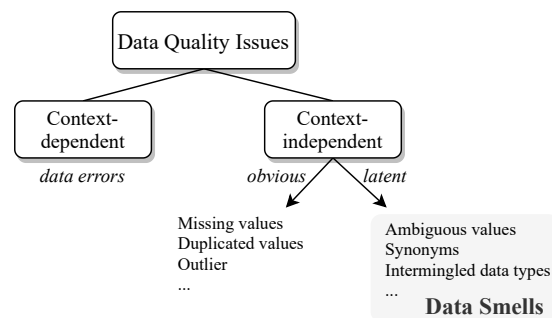


Figure 2: Classification of data smells in the context of data quality issues

To place data smells into this classification scheme, we propose further subdividing context-independent issues into *obvious* (e.g., missing values, duplicate values, obviously misspelled values) and *latent* issues (e.g., synonyms, ambiguous values) as depicted in Figure 2. In fact, we claim that data smells are a subgroup of context-independent issues representing latent data quality issues. Following, we describe how *obvious* data quality issues and *data errors* mainly differ from data smells.

Obvious data quality issues. By obvious issues, we mean issues that are indicated by highly suspicious data values, patterns or representations, which are usually recognizable at first glance (e.g., "New Yorc Zity"). Consequently, obvious issues have a higher degree of suspicion compared to data smells which are "only" indicated by moderate suspicions (e.g., "New York"). A further difference with data smells is that obvious issues are typically detectable through basic data profiling techniques (e.g., descriptive statistics). However, obvious data quality issues may have similar consequences as data smells but are much more likely to be discovered and dealt with.

Data errors. Data errors are strictly tied to a specific context as assured derivations of the ground truth. This contrasts to data smells which indicate suspicions regardless of the concrete context. In addition, data errors cause inevitable and usually timely consequences, whereas the effects of data smells are uncertain. A further difference to data smells, which typically emerge through poor practices, is that the causes of data errors are manifold (e.g., environmental effects on sensors, violation of domain constraints).

Note that both *obvious* and *latent* data quality issues (i.e., data smells) can constitute a data error when set into a concrete context.

For instance, the data smell "New York" can represent a data error if the concrete application context only allows country names (e.g., Hungary, Brazil). Thus, although data smells share many properties with code smells, they differ in this respect. Code smells do not represent a bug or an error. Instead, they are symptoms of poor design choices or implementation practices that can only lead to bugs and further issues in the future (e.g., costly maintenance).

4 CATALOGUE OF DATA SMELLS

In this section, we first outline the research method followed to develop the catalogue in Section 4.1. Subsequently, Section 4.2 presents the catalogue comprising 36 data smells divided into four categories.

Scope. We restrict the scope of data smells considered in the catalogue to data of tabular form as this is still the most widely used form of data organisation. In addition, we only consider data smells within one data attribute (i.e., data column). Furthermore, we limit the scope of data smells in this paper to three types of data: numerical data (i.e., integer and floating-point values), date/time data (i.e., timestamp values) and text data. These three data types were derived from the most common categories of data used in data-driven systems (i.e., categorical, numerical, time-series and text data).

4.1 Research Method

We conducted a multivocal literature review (MLR) [21] to develop a sound catalogue of data smells. A MLR was chosen because it allows a deeper insight into practice by additionally considering grey literature on data quality (e.g., blogs, books) and data quality tools (e.g., WinPure¹, QuerySurge²).

As a first step, we defined a set of keywords (e.g., "suspicious data", "potential data errors") and conducted a search in Google Scholar and the web search engine Google. We applied an effort bounded stop criterion to limit the search results to a manageable size. Therefore, we relied on the ranking algorithm of Google (Scholar), assuming that the most relevant hits usually appear on the first few pages. In detail, we only checked the first five pages of each keyword's results (i.e., 50 hits) and only continued if the last page contained a relevant hit.

In total, we checked more than 400 hits and excluded any sources that did not address (at least partly) latent data issues from the remainder of the review. After conducting forward and backward snowballing [62] on the identified sources, we extracted all latent data issues in a spreadsheet³. To ensure that we did not miss relevant issues due to different applied terminology, we additionally extracted common subtle data errors from established data error taxonomies (e.g., [33, 41]).

As a next step, we reviewed all extracted data issues regarding a moderate degree of suspicion and context-independence. To determine the degree of suspicion, we evaluated the syntactical (i.e., unusual use of characters, formats or data types) as well as the semantical (i.e., implausibility) suspicion of data issues. In addition,

we excluded apparent issues (i.e., obvious data quality issues) that are easily identifiable through data profiling techniques.

After excluding those that did not meet these criteria, we applied an inductive coding approach and assigned descriptive labels to all remaining issues. We then synthesised all labels and derived a list of 36 data smells. The naming of the smells was based on the idea to use the in our opinion most intuitive terms for practitioners. To ensure confirmability of our coding procedure, we conducted a second deductive coding cycle. In doing so, we labelled all extracted data issues again with the corresponding derived data smell(s).

4.2 Categories

We classified the 36 smells into different categories according to the data quality characteristic they primarily violate. This classification scheme is based on Ganesh et al. [20] who proposed the classification of software design smells according to the violation of object-oriented design principles. Because the terminology on data quality characteristics also differs between authors (e.g., Credibility versus Trustworthiness), we decided to apply the in our opinion most intuitive terms for practitioners. Following, we present the different categories and describe a corresponding smell example. The complete catalogue is depicted in Figure 3, while a definition and corresponding examples for each smell are available online³.

4.2.1 Believability Smells. Relate to semantically implausible data values. These smells may indicate a low believability of the data values but are usually understandable (i.e., readable) by humans and interpretable by software.

Dummy Value. This smell characterises a situation where a kind of substitute value may be used. As a concrete example, we claim that "999" represents such a smelly value because it is often used to represent missing values and thus is worth investigating. The value probably refers to a common emergency number if the application context is settled around a telephone directory. However, if the value should represent a person's age, it becomes clear that the value possibly indicates a missing entry.

4.2.2 Understandability Smells. Deal with the inappropriate, unusual or ambiguous use of characters, formats or data types. This may lead to problems in interpreting or reading the data values by humans or software, even though the data values are regarded as semantically true. These smells can further be subdivided into Encoding and Syntactic smells.

Encoding Smells. Represent the inappropriate, unusual or ambiguous use of data types that may lead to data de-/encoding issues.

Integer as String. This smell occurs when an integer is encoded as a string, usually indicated by double quotation marks (e.g., "5"). While some operations can be performed with this data value (e.g., concatenations), others will result in a fault (e.g., additions).

Syntactic Smells. Represent the inappropriate, unusual or ambiguous use of characters or formats that may lead to interpretation issues.

¹<https://winpure.com/>

²<https://www.querysurge.com/>

³<https://bit.ly/2WpPauB>

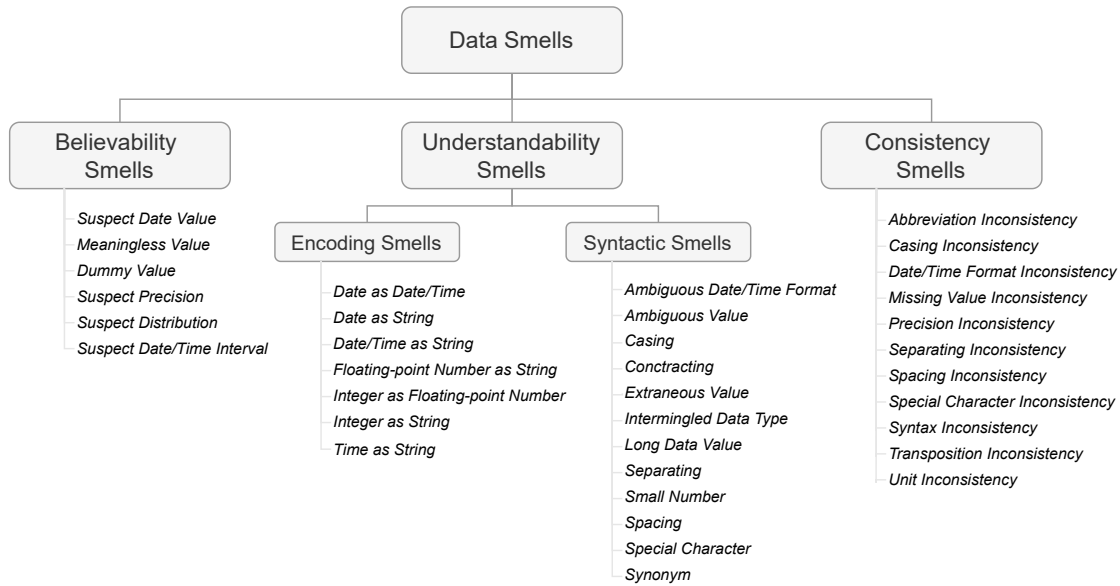


Figure 3: Data smell catalogue

Small Number. This smell occurs when data values represent numbers below 1. If several such values are multiplied in subsequent operations, the calculation result can be very small and has many decimal places. As decimal places are often limited manually (e.g., three decimal places) or by the system (i.e., arithmetic underflow), this can lead to incorrect further processing (e.g., $0.02 * 0.02 = 0.0004$).

4.2.3 *Consistency Smells.* Represent the use of inconsistent syntax with respect to data values in partitions of data.

Abbreviation Inconsistency. This smell arises when abbreviations, acronyms or contradictions are not used consistently. For example, the inconsistent use of titles (e.g., Doctor Hill versus Dr. Hill) can lead to malicious results in a natural language processing application.

5 CAUSES AND CONSEQUENCES OF DATA SMELLS

Based on our own experience with industrial data-driven projects and relevant literature (e.g., [9, 15, 46, 54]), we elaborate on the causes and consequences of data smells in this section. First, Section 5.1 describes the most common causes of data smells. Then, Section 5.2 illustrates how data smells may affect the correct functioning as well as the development and maintenance of AI-based systems.

5.1 Causes

5.1.1 *Data management.* Bad practices in data collection or preparation can cause the emergence of data smells (e.g., careless data entry, inconsistent data collection or transformation processes). Further, a lack of documentation (e.g., missing data lineage, no data dictionary) or poor communication between the different actors in the data life-cycle can introduce subtle issues in the data. For

example, incomplete metadata can lead to incorrect assumptions about the data by software engineers, resulting in incorrect implemented data processing logic, which causes smelly data. Recent studies outline that poor data practices [42, 49] and data-related mismatches [34] are pressing issues in the field of AI-based systems.

5.1.2 *Data handling.* According to a recent study [63], data handling code tends to be error-prone and often contains subtle issues. We claim that such poor data handling practices will likely cause data smells. For example, not being explicit when converting a date string (e.g., "2021-01-01") into a datetime object (e.g., `pandas.to_datetime`) causes the Date as Date/Time smell (i.e., "2021-01-01 00:00:00") to arise. As the conversion input represents just a date, developers may be unaware that without further declaration, a time suffix (i.e., "00:00:00") is added.

Such issues often go undetected because of the common programmatic style of method chaining when processing data [63]. By sequencing multiple data operations (i.e., method chaining), developers cannot see the intermediate processing results and thus identify problems introduced in the data [61].

5.1.3 *Data source quality.* Data smells can further arise through a poor intrinsic quality of the data sources from which the data originate. For example, a column name in a relational database is used in different tables but with different data types [46]. Accordingly, this can cause data encoding smells when developers retrieve the data as they are not expecting different data types. Often, however, data are already stored smelly in the data sources.

5.2 Consequences

5.2.1 *Defects and Failures.* The occurrence of data smells can cause interpretation problems of downstream software components in AI-based systems. For example, consider a smell (e.g., Integer as String, Intermingled Data Type) in one of ten thousand data instances of a

data attribute to be loaded and processed. When loading data, data processing software libraries typically apply a data type inference algorithm by default (e.g., `pandas.read_csv`). As this functionality selects the most flexible data type for a data attribute, a single data smell can lead to a wrong inferred data type [26]. Such interpretation problems can cause further incorrect data processing (e.g., type conversion errors) or lead to faults (e.g., arithmetical operations of strings, concatenations of integers) in AI-based systems.

In addition, data smells can also lead to incorrect knowledge generation in an AI-based system. As modern AI-based systems often pursue a continuous learning strategy (e.g., lifelong learning), they typically use serving data as training data. Thus, a well-trained ML model can degrade over time based on smelly serving data continuously used to update the model [15]. For example, consider a model that depends on countries as input data. Serving data that becomes smelly, for instance, "US" instead of "us" (i.e., Casing smell), could lead to a wrong result of the model [15, 43]. Due to their high ability to integrate new knowledge (i.e., plasticity), artificial neuronal networks are especially prone to such subtle data issues [2].

Moreover, the output of one AI-based system is often directly consumed by other systems or even influences its own training data (e.g., direct feedback loops) [9, 54]. In such cases, latent data issues can cascade to severe problems over time, causing a gradual regression of the performance of the ML models involved [11].

5.2.2 Development and Maintainability. Data smells often require additional data cleaning or preprocessing routines in AI-based systems. For instance, a further code fragment to lower-case all training data would be needed to address the casing smell mentioned above (i.e., "US" instead "us") [9]. As such additional data handling code tends to be error-prone [63] and difficult to understand for developers [61], it is supposed to affect the code comprehensibility of AI-based systems. In addition, data handling code is likely to introduce glue code or pipeline jungles in AI-based systems which results in the creation of technical debt over time [54]. Consequently, the occurrence of data smells and their corresponding treatment makes AI-based systems harder to understand and maintain and thus increases the risk of introducing further errors.

Further, data smells make it difficult for humans to understand the data concerned. Accordingly, this can result in wrong assumptions that lead to incorrect implemented functionality or faulty ML models. For example, consider timestamps (e.g., "08:00") represented in 12-hour clock format with missing a.m. or p.m. designators (i.e., Ambiguous Date/Time Format smell). Whereas such time values do not indicate an evident data quality issue, they can be misinterpreted by developers because it is not obvious whether the value refers to eight o'clock in the morning or evening. However, as long as the data describe a casual working day, this smell does not hinder the correct functionality of a ML model. Nevertheless, if shift work is introduced and the model needs to be updated, problems with the correct interpretation of the timestamp "08:00" may arise.

5.2.3 Real-life Examples. Below, we briefly describe how data smells negatively impacted a medical health system and caused problems in COVID-19 software projects.

Oncology Expert Advisor. The Oncology Expert Advisor (OEA) is an AI-based clinical decision support system developed by IBM [14, 59]. Based on IBM's cognitive computing system Watson, the main task of OEA is to provide evidence-based therapy recommendations to doctors. Therefore, the system ingests a huge amount of data (e.g., medical literature, practice guidelines, electronic health records). However, during a project with an American cancer center, OEA suffered heavily from latent data quality issues such as acronyms, shorthand phrases or different styles of writing [48]. Due to further issues, the project was cancelled in 2016 [36].

COVID-19 Software Projects. Consequences of data smells were also reported by an empirical study of COVID-19 software projects [45]. In this study, the authors investigated the occurrence of bugs in open source software projects that mine and aggregate COVID-19 data. As the second most common bug category, the study identified bugs that occur during processing data (i.e., data bugs). However, many of these data bugs were related to smelly data. For example, a web crawler suffered from a Suspect Date/Time Interval smell [12]. In another project, problems arose because the names of some districts in an Indian state differed slightly [13]. A different naming for some districts by official sources caused this Ambiguous Value smell to occur.

6 DETECTION OF DATA SMELLS

In this section, we first discuss detecting data smells in the realm of data validation in Section 6.1. Subsequently, we propose how data smell detection can be approached in Section 6.2. Section 6.3 then describes the developed tool support for detecting smells. In Section 6.4, we present the conducted experimental evaluation of the tool support. Lastly, Section 6.5 discusses the use of data smells in the context of AI-based systems.

6.1 Preliminaries

Data smell detection is closely related to the field of data validation. Validating data has become common practice in AI-based systems as algorithms in these systems strongly rely on the quality of the data fed to them [5]. Basically, data validation in AI-based systems is built upon context-specific schemes and constraints to detect issues in the input data [10]. Since AI-based systems in production often run continuously in real-time, immediate and straightforward human intervention in case of problems is crucial [11]. Thus, data validation systems are typically tailored to provide reliable, high-precision alerts about actual data errors [9, 15]. Notable solutions for validating data in AI-based systems are Amazon's Deequ [51, 52] or Google's TensorFlow Data Validation (TFDV) [11].

Although integrating data smell detection into data validation systems seems evident at first glance, we propose to decouple it for the following reasons.

First, detecting data smells will likely produce many alerts about potential issues. On the one hand, this vast amount of alerts is caused by the more generally specified detection methods needed to ensure context-independence detection. On the other hand, the stringent definition of detection thresholds needed because of the only moderate suspicion of smelly data also contributes to a large number of alerts. This large number of alerts is problematic because determining whether a detected smell really causes problems in

AI-based systems requires substantial human effort. In fact, a smell must be set into the concrete application context and requires a detailed examination of a system's downstream components and applications. Thus, detecting smelly data during validating data contrasts with the aim of data validation systems to provide precise, immediate and actionable notifications about actual errors in the data [44].

Second, applying data quality checks on continuously arriving data may affect the performance of AI-based systems (i.e., latency) [39]. Therefore, the amount and stringency of the detection methods used in data validation systems must be treated with caution. Accordingly, this also contradicts the integration of data smell detection into data validation systems.

In summary, smelly data are basically not data errors and are therefore not suitable to be detected in data validation systems. Since data smells typically arise through poor practices, the primary goal should not be to detect and clean them in real-time but to identify and resolve their root causes. This is closely related to the practice of refactoring to eliminate code smells in software engineering.

6.2 Detection Approach

As outlined in the previous section, detecting data smells as part of data validation modules in AI-based systems would come with several drawbacks. We thus propose to detect data smells in an offline manner without affecting systems running in production. To deal with the presumably large number of smell alerts and to vary the level of suspicion to be flagged, we further introduce two metrics.

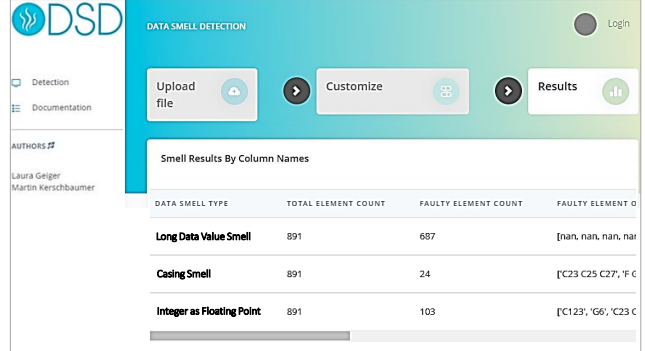
Data Smell Strength. This metric indicates the likelihood that a data value (i.e., data instance) or pattern (i.e., data partition) is treated as suspicious, and a smell is raised. Basically, this metric implies the concrete thresholds and/or hyperparameters of the individual detection methods. For example, the number of contiguous characters required to detect a Long Data Value smell.

Data Smell Density. This metric describes the relative number of detected smells of a data attribute (i.e., data column). Therefore, this metric can be used to focus on data attributes with a high density of smells. For example, a data attribute can only be considered smelly (i.e., *Smelly Data Attribute*) if at least 10 percent of its data instances represent a data smell.

6.3 Tool Support

To operationalize the data smell detection, we checked several data quality-related tools (e.g., OpenRefine⁴, WinPure¹) as well as several data validation tools and libraries, such as Deequ [52], TFDV [11], Data Sentinel [58], MobyDQ [40], Data Quality Toolkit [24], Data Quality Advisor [56] and Great Expectations [23] for their ability to detect smelly data. Most of the investigated tools have a particular focus (e.g., outlier detection), are not publicly available (e.g., [58]), or come with a limited, fixed set of built-in data checks (e.g., [24]). Thus, we did not find any tool or library that is able to detect most of the proposed data smells without adaptations. We

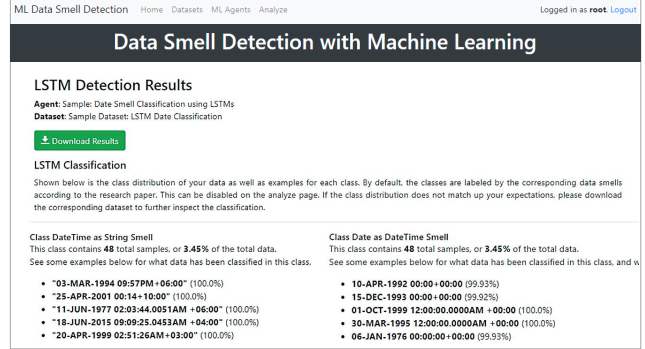
therefore decided to develop an own solution to enable automated data smell detection.



The screenshot shows the DSD web application interface. It has a sidebar with 'Detection' and 'Documentation' links. The main area has 'Upload file', 'Customize', and 'Results' buttons. Below, it displays 'Smell Results By Column Names' in a table.

DATA SMELL TYPE	TOTAL ELEMENT COUNT	FAULTY ELEMENT COUNT	FAULTY ELEMENT O
Long Data Value Smell	891	687	[nan, nan, nan, nar
Casing Smell	891	24	['C23 C25 C27', 'F C
Integer as Floating Point	891	103	['C123', 'G6', 'C23 C

(a) Rule-based detection



The screenshot shows the 'ML Data Smell Detection' web application interface. It has a sidebar with 'Home', 'Datasets', 'ML Agents', and 'Analyse' links. The main area is titled 'Data Smell Detection with Machine Learning' and shows 'LSTM Detection Results'.

LSTM Classification

Shown below is the class distribution of your data as well as examples for each class. By default, the classes are labeled by the corresponding data smells according to the research paper. This can be disabled on the analyze page. If the class distribution does not match up your expectations, please download the corresponding dataset to further inspect the classification.

Class DateTime as String Smell	Class Date as DateTime Smell
This class contains 48 total samples, or 3.45% of the total data. See some examples below for what data has been classified in this class.	This class contains 48 total samples, or 3.45% of the total data. See some examples below for what data has been classified in this class, and
<ul style="list-style-type: none"> • "03-MAR-1994 09:57PM+06:00" (100.0%) • "25-APR-2001 00:14+10:00" (100.0%) • "11-JUN-1977 02:03:44.0051AM +06:00" (100.0%) • "18-JUN-2015 09:09:25.0453AM +04:00" (100.0%) • "20-APR-1999 02:51:26AM+03:00" (100.0%) 	<ul style="list-style-type: none"> • 10-APR-1992 00:00+00:00 (99.93%) • 15-DEC-1993 00:00+00:00 (99.92%) • 01-OCT-1999 12:00:00.0000AM +00:00 (100.0%) • 30-MAR-1995 12:00:00.0000AM +00:00 (100.0%) • 06-JAN-1976 00:00+00:00 (99.93%)

(b) ML-based detection

Figure 4: Tool support

After reviewing each defined smell from the catalogue, we concluded that while some data smells (e.g., Long Data Value smell) are rather easy to spot using regular expressions, others (e.g., Synonym smell) require more advanced techniques to be reliably detected. Therefore, we propose to detect data smells using a *rule-based* and a *ML-based approach*. Traditional data validation techniques (e.g., range and data type checks) can be adapted and used to implement rule-based smell detection. For realizing a learning-based approach, techniques from natural language processing (e.g., recurrent neural networks) or anomaly detection (e.g., autoencoders) are worthwhile candidates.

Following, we present the two tools developed based on the outlined detection approaches. Several detection methods implemented in the tools were calibrated and trained with many publicly available datasets to ensure general applicability.

6.3.1 Rule-based Detection. The first tool⁵ is based on the open-source data validation tool *Great Expectations*⁶ and focuses on a rule-based smell detection. Great Expectations was chosen due to its ability to be easily extensible and thus suitable for implementing data smell-specific detections. By adjusting so called expectations

⁴<https://openrefine.org/>

⁵<https://github.com/mkerschbaumer/rb-data-smell-detection>

⁶<https://greatexpectations.io>

(i.e., assertions about data) provided by Great Expectations, the tool is able to detect smells such as Long Data Value, Casing, Integer as String, Floating Point Number as String or Integer as Floating Point Number. The tool provides a graphical user interface where CSV files can be uploaded. Further, the degree of suspicion to be flagged can be easily chosen based on predefined settings. However, it is also possible to define each parameter of a detection method individually.

6.3.2 Machine learning-based Detection. The second tool⁷ aims to detect data smells based on several ML algorithms. For example, we used the Word2Vec algorithms implemented in the Python library Gensim to realise the Synonym smell detection. Further, several inconsistency smells (e.g., Casing and Spacing Inconsistency) were implemented by using the Python library Dedupe. Neuronal networks (e.g., long short-term memory, autoencoders) for detecting smells such as Ambiguous Date/Time Format or Date/Time Format Inconsistency were realised with the library Keras. The developed tool also comes with a graphical user interface and accepts CSV files. Most of the implemented ML models were trained for general applicability but also allow to be retrained on user-specific data. A screenshot of detection results of both tools is depicted in Figure 4.

6.4 Experimental Evaluation

We applied several data smell detection methods on 246 Kaggle datasets to test and evaluate the tool support.

6.4.1 Setting. All datasets were randomly selected and are licensed under the Creative Commons CC0. To conduct the smell detection, the datasets were first downloaded and further processed with Python scripts to be analysable by the tool support. In total, the datasets resulted in more than 2,000 columns and more than 42 billion rows to be analysed. The majority of the columns were of type String (i.e., 1,130). Figure 5 visualizes the number of rows for each column grouped by its type. Although the number of rows varies widely (i.e., 2 to 8,405,079), 50 percent of each considered column type contained more than 41,500 rows.

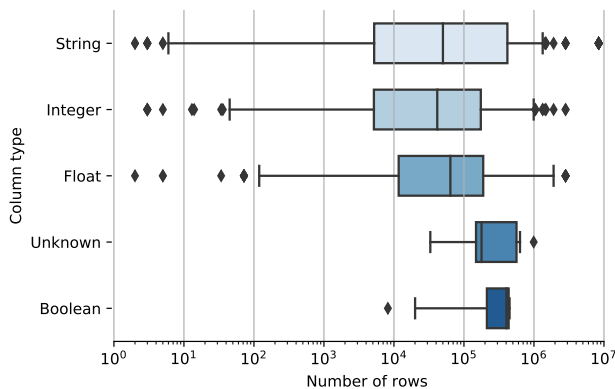


Figure 5: Boxplot of number of rows grouped by column type

⁷<https://github.com/georg-wenzel/ml-data-smell-detection>

6.4.2 Results. We analysed the number of Smelly Data Attributes for each dataset to get an initial impression on the occurrence of smells in real-world data. An attribute (i.e., column) was considered smelly if at least one smell was detected. Figure 6 shows the corresponding frequency distribution.

It is apparent from this figure that most of the datasets (i.e., 120) contained three to five smelly attributes. Only 45 datasets were detected with more than 10 smelly attributes. In contrast, three datasets had no data smells. These initial observations suggest that the number of detected smelly attributes seems to be at a manageable size for real-world data.

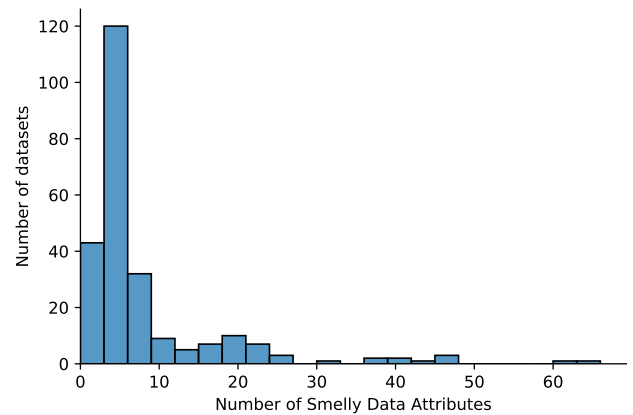


Figure 6: Histogram of smelly data attributes across all datasets

Exemplary Smells. Figure 7 shows an excerpt of detected smells in real-world data. On the left-hand side, Subfigure 7a shows an example of the Ambiguous Date/Time Format smell⁸. In fact, the date is represented in short format, which may lead to interpretation problems as it is unclear which date format is applied. Subfigure 7b also depicts an example of the Ambiguous Date/Time Format smell. The corresponding COVID-19 dataset⁹ represents case data but omits the year in the date column. Thus, without further information (e.g., metadata), it can lead to wrong assumptions about the corresponding year. In Subfigure 7c, an example of the Date as Date/Time smell¹⁰ is shown. The corresponding column represents the day of medical appointments with a time suffix (i.e., "00:00:00"). Thus, systems consuming these data may encounter problems in processing these data correctly.

6.5 Use

Detecting data smells is especially useful in the engineering phase of AI-based systems. A data validation system has typically not yet been implemented in this phase, and data quality assurance is thus often neglected. However, data smell detection can be applied with little effort due to the context-independent nature of smelly data. For instance, by checking training data for smells, potentially

⁸<https://www.kaggle.com/ravichaubey1506/covid19-india>

⁹<https://www.kaggle.com/joniarroba/noshowappointments>

¹⁰<https://www.kaggle.com/nasa/landslide-events>

date	country_n...	state/provi...	Date	# Daily Confi...	# Total Confi...	Appointment...	Appointment...	# Age
6/27/07	Ecuador	Pichincha	11-Feb	0	3	5638447	2016-04-29T00:00:00Z	21
7/1/07	United States	Texas	12-Feb	0	3	5629123	2016-04-29T00:00:00Z	19
7/4/07	Mexico	Veracruz-Llave	13-Feb	0	3	5638213	2016-04-29T00:00:00Z	30
7/8/07	Canada	Ontario	14-Feb	0	3	5628163	2016-04-29T00:00:00Z	29
7/13/07	Dominican Republic	Distrito Nacional	15-Feb	0	3	5634718	2016-04-29T00:00:00Z	22
7/24/07	United States	Texas	16-Feb	0	3	5636249	2016-04-29T00:00:00Z	28
8/9/07	Guatemala	Guatemala	17-Feb	0	3			
8/11/07	Jamaica	Portland	18-Feb	0	3			

(a) Ambiguous Date/Time Format smell

(b) Ambiguous Date/Time Format smell

(c) Date as Date/Time smell

Figure 7: Data smells in real-world data

problematic data sources can be identified before they are later used in production. Furthermore, data smells can guide data validation efforts in environments with many different data streams. For example, since validating thousands of data signals is impossible, data engineers can consider the smell density (i.e., Data Smell Density), besides other aspects (e.g., feature importance), when deciding which data fractions to validate. Additionally, data smells can uncover subtle issues in the data handling logic or software code that passed traditional quality assurance processes and thus would creep into the deployed system.

Further, checking data on the occurrence of smells is also useful on systems running in production. In fact, regularly detecting smells on archived serving data can identify data issues that were not caught by a system's data validation components. Thus, root causes of these smells can be identified and fixed before the smelly data may cause problems over time (e.g., model degradation, affecting other systems). If the root causes cannot be identified or resolved, data validation components can at least be adjusted (e.g., strengthen their constraints) to catch these smelly data.

In summary, detecting data smells effectively reduces technical debt and increases the quality of data in AI-based systems.

7 CONCLUSIONS

The main aim of this article was twofold. First, we have addressed the lack of research on latent data quality issues by conceptualizing them in the form of data smells. Therefore, we have presented a sound definition, the main characteristics, and a catalogue of data smells by analogy to the concept of code smells in software engineering.

Second, we have highlighted the importance of the often neglected yet pervasive, smelly data in AI-based systems. Therefore, we have explicitly elaborated on the causes and consequences of data smells in such systems. In detail, we have shown that bad practices (i.e., data management and handling) in engineering AI-based systems causes smelly data to arise. We further have demonstrated how data smells can negatively impact the correct functioning as well as the development and maintenance of AI-based systems. Two presented cases (i.e., IBM's Oncology Expert Advisor, COVID-19 software projects) illustrated the consequences and currentness of smelly data in real-world projects.

In addition, we have discussed the detection of data smells and proposed a rule-based as well as a ML-based detection approach decoupled from traditional data validation efforts in AI-based systems. We implemented both detection approaches as tool support and conducted an initial experimental evaluation on 246 Kaggle datasets. The evaluation indicated that the number of detected data smells seems to be at a manageable size for real-world data. Lastly, we have proposed application scenarios of detecting smells to provide a solid baseline for future research.

In fact, we intend to identify less and more influential smells regarding their impact on AI-based systems in the future. For this purpose, we plan to apply a framework proposed by Schelter et al. [53] (i.e., JENGA) which is able to study the impact of data issues on ML models. Moreover, future research should investigate whether using data smell detection as part of testing software components to spot subtle implementation errors is beneficial. On a wider level, research should also focus on determining data smells for specific domains (e.g., finance, industry, medicine). There are already contributions to build on, such as a publication by Ehrlinger et al. [16], which investigates patterns of missing industrial data.

ACKNOWLEDGMENTS

The research reported in this paper has been partly funded by the Federal Ministry for Climate Action, Environment, Energy, Mobility, Innovation and Technology (BMK) and the Federal Ministry for Digital and Economic Affairs (BMDW) as well as the State of Upper Austria in the frame of the COMET - Competence Centers for Excellent Technologies Programme (FFG-Nr. 865891) and the project ConTest (FFG-Nr. 888127) managed by Austrian Research Promotion Agency FFG. We also thank Martin Kerschbaumer, Georg Wenzel and Laura Geiger for their contribution in the implementation of the tool support.

REFERENCES

- [1] Ziawasch Abedjan, Xu Chu, Dong Deng, Raul Castro Fernandez, Ihab F. Ilyas, Mourad Ouzzani, Paolo Papotti, Michael Stonebraker, and Nan Tang. 2016. Detecting data errors: Where are we and what needs to be done? *Proceedings of the VLDB Endowment* 9, 12 (2016), 993–1004. <https://doi.org/10.14778/2994509.2994518>
- [2] Amina Adadi. 2021. A survey on data-efficient algorithms in big data era. *Journal of Big Data* 8, 1 (2021).

- [3] Daniel W. Barowy, Dimitar Gochev, and Emery D. Berger. 2014. CheckCell. In *Proceedings of the 2014 ACM International Conference on Object Oriented Programming Systems Languages & Applications*, Andrew Black and Todd Millstein (Eds.). ACM, 507–523.
- [4] Carlo Batini, Monica Scannapieco, et al. 2016. *Data and information quality*. Springer, Cham, Switzerland.
- [5] Felix Biessmann, Jacek Golebiowski, Tammo Rukat, Dustin Lange, and Philipp Schmidt. 2021. Automated Data Validation in Machine Learning Systems. *Bulletin of the IEEE Computer Society Technical Committee on Data Engineering* 44, 1 (2021), 51–65.
- [6] Justus Bogner, Roberto Verdecchia, and Ilias Gerostathopoulos. 2021. Characterizing Technical Debt and Antipatterns in AI-Based Systems: A Systematic Mapping Study. *arXiv preprint arXiv:2103.09783* (2021).
- [7] Jan Bosch, Helena Holmström Olsson, and Ivica Crnkovic. 2021. Engineering AI systems: A research agenda. In *Artificial Intelligence Paradigms for Smart Cyber-Physical Systems*. IGI Global, 1–19.
- [8] Houssein Ben Braiek and Foutse Khomh. 2020. On Testing Machine Learning Programs. *Journal of Systems and Software* 164 (2020), 110542.
- [9] Eric Breck, Neoklis Polyzotis, Sudip Roy, Steven Whang, and Martin Zinkevich. 2019. Data Validation for Machine Learning. In *Proceedings of the 2nd SysML Conference*.
- [10] Eric Breck, Neoklis Polyzotis, Sudip Roy, Steven Euijong Whang, and Martin Zinkevich. 2018. Data infrastructure for machine learning. In *SysML Conference*.
- [11] Emily Caveness, Paul Suganthan GC, Zhuo Peng, Neoklis Polyzotis, Sudip Roy, and Martin Zinkevich. 2020. Tensorflow data validation: Data analysis and validation in continuous ml pipelines. In *Proceedings of the 2020 ACM SIGMOD International Conference on Management of Data*. 2793–2796.
- [12] coronadatascraper. 2020. Data has a gap between 2020-3-11 and 2020-3-24 #375: GitHub. <https://github.com/covidatlas/coronadatascraper/issues/375>
- [13] covid19india react. 2020. Rajasthan District names are wrong #321: GitHub. <https://github.com/covid19india/covid19india-react/issues/321>
- [14] Thomas H. Davenport. 2015. Lessons-from-the-Cognitive-Front-Lines-Early-Adopters-of-IBMs-Watson. *The Wall Street Journal* 1 (2015). <https://www.tomdavenport.com/wp-content/uploads/2019/01/Lessons-from-the-Cognitive-Front-Lines-Early-Adopters-of-IBMs-Watson.pdf>
- [15] Mike Dreves, Gene Huang, Zhuo Peng, Neoklis Polyzotis, Evan Rosen, and Paul Suganthan GC. 2021. Validating Data and Models in Continuous ML pipelines. *Bulletin of the IEEE Computer Society Technical Committee on Data Engineering* 44, 1 (2021), 42–50.
- [16] Lisa Ehrlinger, Thomas Grubinger, Bence Varga, Mario Pichler, Thomas Natschlager, and Jürgen Zeindl. 2018. Treating missing data in industrial data analytics. In *Thirteenth International Conference on Digital Information Management (ICDIM)*. 148–155.
- [17] Harald Foidl and Michael Felderer. 2019. Risk-based data validation in machine learning-based software systems. In *Proceedings of the 3rd ACM SIGSOFT International Workshop on Machine Learning Techniques for Software Quality Evaluation*. 13–18.
- [18] Ralph Foorhuis. 2018. A Typology of Data Anomalies. In *Information Processing and Management of Uncertainty in Knowledge-Based Systems. Theory and Foundations*, Jesús Medina, Manuel Ojeda-Aciego, José Luis Verdegay, David A. Pelta, Inma P. Cabrera, Bernadette Bouchon-Meunier, and Ronald R. Yager (Eds.). Communications in Computer and Information Science, Vol. 854. Springer International Publishing, Cham, 26–38.
- [19] Martin Fowler, Kent Beck, J. Brant, W. Opdyke, and D. Roberts. 1999. Refactoring: improving the design of existing code, ser. In *Addison Wesley object technology series*. Addison-Wesley.
- [20] S. G. Ganesh, Tushar Sharma, and Girish Suryanarayana. 2013. Towards a Principle-based Classification of Structural Design Smells. *The Journal of Object Technology* 12, 2 (2013), 1–29.
- [21] Vahid Garousi, Michael Felderer, and Mika V. Mäntylä. 2019. Guidelines for including grey literature and conducting multivocal literature reviews in software engineering. *Information and Software Technology* 106 (2019), 101–121.
- [22] Mouzhi Ge and Markus Helfert. 2007. A review of Information Quality Research: Develop a Research Agenda. In *Proceedings of the 2007 MIT International Conference on Information Quality (ICIQ)*.
- [23] Great Expectations. 2022. <https://greatexpectations.io/>
- [24] Nitin Gupta, Hima Patel, Shazia Afzal, Naveen Panwar, Ruhi Sharma Mittal, Shanmukha Guttula, Abhinav Jain, Lokesh Nagalapatti, Sameep Mehta, Sandeep Hans, et al. 2021. Data Quality Toolkit: Automatic assessment of data quality and remediation for machine learning datasets. *arXiv preprint arXiv:2108.05935* (2021).
- [25] Jacob Harris. 2014. Distrust Your Data. <https://source.opennews.org/articles/distrust-your-data/>
- [26] Nick Hynes, D. Sculley, and Michael Terry. 2017. The data linter: Lightweight, automated sanity checking for ml data sets. In *31st Conference on Neural Information Processing Systems (NIPS): Workshop on ML Systems*.
- [27] Ihab F. Ilyas and Xu Chu. 2019. *Data cleaning*. Morgan & Claypool.
- [28] Md Johirul Islam, Giang Nguyen, Rangeet Pan, and Hridesh Rajan. 2019. A comprehensive study on deep learning bug characteristics. In *Proceedings of the 2019 27th ACM Joint Meeting on European Software Engineering Conference and Symposium on the Foundations of Software Engineering*. 510–520.
- [29] Nikolas Iubel. 2014. Ensuring Accuracy in Data Journalism. <https://github.com/nikeiubel/data-smells/wiki/Ensuring-Accuracy-in-Data-Journalism>
- [30] Jácome Cunha, João P. Fernandes, Hugo Ribeiro, and João Saraiva. 2012. Towards a Catalog of Spreadsheet Smells. *International Conference on Computational Science and Its Applications* (2012), 202–216.
- [31] Joao Marcelo Borovina Josko, Lisa Ehrlinger, and Wolfram Wöß. 2019. Towards a Knowledge Graph to Describe and Process Data Defects. In *DBKDA 2019: The Eleventh International Conference on Advances in Databases, Knowledge, and Data Applications*. 57–60.
- [32] Mark Kasunic, James McCurley, Dennis Goldenson, and David Zubrow. 2011. An Investigation of Techniques for Detecting Data Anomalies in Earned Value Management Data. <https://doi.org/10.21236/ADA591417>
- [33] Won Kim, Byoung-Ju Choi, Eui-Kyeong Hong, Soo-Kyung Kim, and Doheon Lee. 2003. A Taxonomy of Dirty Data. *Data Mining and Knowledge Discovery* 7, 1 (2003), 81–99.
- [34] Grace A. Lewis, Stephany Bellomo, and Ipek Ozkaya. 2021. Characterizing and Detecting Mismatch in Machine-Learning-Enabled Systems. *arXiv preprint arXiv:2103.14101* (2021).
- [35] Lin Li, Taoxin Peng, and Jessie Kennedy. 2011. A rule based taxonomy of dirty data. *GSTF Journal on Computing (JoC)* 1, 2 (2011).
- [36] Steve Lohr. 2021. What Ever Happened to IBM's Watson? <https://www.nytimes.com/2021/07/16/technology/what-happened-ibm-watson.html>
- [37] Lucy Ellen Lwakatare, Ellinor Ränge, Ivica Crnkovic, and Jan Bosch. 2021. On the experiences of adopting automated data validation in an industrial machine learning project. In *IEEE/ACM 43rd International Conference on Software Engineering: Software Engineering in Practice (ICSE-SEIP)*. 248–257.
- [38] Silverio Martínez-Fernández, Justus Bogner, Xavier Franch, Marc Oriol, Julien Siebert, Adam Trendowicz, Anna Maria Vollmer, and Stefan Wagner. 2021. Software Engineering for AI-Based Systems: A Survey. *arXiv preprint arXiv:2105.01984* (2021).
- [39] Jorge Merino, Xiang Xie, Ian Lewis, Ajith Parlikad, and Duncan McFarlane. 2020. Impact of Data Quality in Real-Time Big Data Systems. In *CEUR Workshop Proceedings*. Vol. 2716. 73–86.
- [40] MobyDQ. 2022. <https://ubisoft.github.io/mobydq/>
- [41] Paulo Oliveira, Fátima Rodrigues, Pedro Henriques, and Helena Galhardas. 2005. A taxonomy of data quality problems. In *2nd Int. Workshop on Data and Information Quality*. 219–233.
- [42] Amandalynne Paullada, Inioluwa Deborah Raji, Emily M. Bender, Emily Denton, and Alex Hanna. 2021. Data and its (dis)contents: A survey of dataset development and use in machine learning research. *Patterns* 2, 11 (2021), 100336.
- [43] Neoklis Polyzotis, Sudip Roy, Steven Euijong Whang, and Martin Zinkevich. 2018. Data lifecycle challenges in production machine learning: a survey. *ACM SIGMOD Record* 47, 2 (2018), 17–28.
- [44] Neoklis Polyzotis and Matei Zaharia. 2021. What can Data-Centric AI Learn from Data and ML Engineering? *arXiv preprint arXiv:2112.06439* (2021).
- [45] Akond Ashfaq Ur Rahman and Effat Farhana. 2021. An Empirical Study of Bugs in COVID-19 Software Projects. *Journal of Software Engineering Research and Development* 9 (2021).
- [46] redgate. 2017. SQL Code Smells. <http://assets.red-gate.com/community/books/sql-code-smells.pdf>
- [47] Sergey Redyuk, Zoi Kaoudi, Volker Markl, and Sebastian Schelter. 2021. Automating Data Quality Validation for Dynamic Data Ingestion. In *EDBT*. 61–72.
- [48] Casey Ross and Ike Swetlitz. 2017. IBM pitched its Watson supercomputer as a revolution in cancer care. It's nowhere close. (2017).
- [49] Nithya Sambasivan, Shivani Kapania, Hannah Highfill, Diana Akrong, Praveen Paritosh, and Lora M. Aroyo. 2021. "Everyone wants to do the model work, not the data work": Data Cascades in High-Stakes AI. In *Proceedings of the 2021 CHI*. 1–15.
- [50] José Amancio M. Santos, João B. Rocha-Junior, Luciana Carla Lins Prates, Rogeres Santos do Nascimento, Mydiã Falcão Freitas, and Manoel Gomes de Mendonça. 2018. A systematic review on the code smell effect. *Journal of Systems and Software* 144 (2018), 450–477.
- [51] Sebastian Schelter, Stefan Grafberger, Philipp Schmidt, Tammo Rukat, Mario Kiessling, Andrey Taptunov, Felix Biessmann, and Dustin Lange. 2018. Deequdata quality validation for machine learning pipelines. In *Machine Learning Systems workshop at the conference on Neural Information Processing Systems (NeurIPS)*.
- [52] Sebastian Schelter, Dustin Lange, Philipp Schmidt, Meltem Celikel, Felix Biessmann, and Andreas Grafberger. 2018. Automating large-scale data quality verification. *Proceedings of the VLDB Endowment* 11, 12 (2018), 1781–1794.
- [53] Sebastian Schelter, Tammo Rukat, and Felix Biessmann. 2021. JENGA - A Framework to Study the Impact of Data Errors on the Predictions of Machine Learning Models. In *24th International Conference on Extending Database Technology (EDBT), March 23-26, 2021*. 529–534.

- [54] David Sculley, Gary Holt, Daniel Golovin, Eugene Davydov, Todd Phillips, Dietmar Ebner, Vinay Chaudhary, Michael Young, Jean-Francois Crespo, and Dan Dennison. 2015. Hidden technical debt in machine learning systems. *Advances in neural information processing systems* 28 (2015), 2503–2511.
- [55] Tushar Sharma, Marios Fragkoulis, Stamatia Rizou, Magiel Bruntink, and Diomidis Spinellis. 2018. Smelly relations: measuring and understanding database schema quality. In *Proceedings of the 40th International Conference on Software Engineering: Software Engineering in Practice*. 55–64.
- [56] Shrey Shrivastava, Dhaval Patel, Anuradha Bhamidipaty, Wesley M. Gifford, Stuart A. Siegel, Venkata Sitaramagiridharganesh Ganapavarapu, and Jayant R. Kalagnanam. 2019. Dqa: Scalable, automated and interactive data quality advisor. In *2019 IEEE International Conference on Big Data*. 2913–2922.
- [57] Dina Sukhobok, Nikolay Nikolov, and Dumitru Roman. 2017. Tabular Data Anomaly Patterns. In *International Conference on Big Data Innovations and Applications (Innovate-Data)*, 21–23 August 2017. 25–34.
- [58] Arun Swami, Sriram Vasudevan, and Joojay Huyn. 2020. Data sentinel: A declarative production-scale data validation platform. In *IEEE 36th International Conference on Data Engineering (ICDE)*. 1579–1590.
- [59] Koichi Takahashi, Hagop M. Kantarjian, Guillermo Garcia-Manero, Rick J. Stevens, Courtney Denton Dinardo, Joshua Allen, Emily Hardeman, Scott Carrier, Cynthia Powers, Pat Keane, Sherry Pierce, Mark Routbort, Thai Nguyen, Brett Smith, Jeffery Frey, Keith Perry, John C. Frenzel, Rob High, Andrew Futreal, and Lynda Chin. 2014. MD Anderson's Oncology Expert Advisor powered by IBM Watson: A Web-based cognitive clinical decision support tool. *Journal of Clinical Oncology* 32, 15_suppl (2014), 6506.
- [60] Eva van Emden and Leon Moonen. 2012. Assuring software quality by code smell detection. In *19th Working Conference on Reverse Engineering*. IEEE, xix.
- [61] Zehao Wang. 2021. Understanding the Challenges and Assisting Developers with Developing Spark Applications. In *IEEE/ACM 43rd International Conference on Software Engineering: Companion Proceedings (ICSE-Companion)*. 132–134.
- [62] Claes Wohlin. 2014. Guidelines for snowballing in systematic literature studies and a replication in software engineering. In *Proceedings of the 18th International Conference on Evaluation and Assessment in Software Engineering - EASE '14*, Martin Shepperd, Tracy Hall, and Ingunn Myrtveit (Eds.). ACM Press, New York, New York, USA, 1–10.
- [63] Chenyang Yang, Shurui Zhou, Jin L. C. Guo, and Christian Kästner. 2021. Subtle bugs everywhere: Generating documentation for data wrangling code. In *Proceedings of the 36th IEEE/ACM International Conference on Automated Software Engineering (ASE)*. IEEE, Vol. 11.