

RESEARCH ARTICLE

A Benchmarking and Comparative Analysis of Python Libraries for Data Cleaning: Evaluating Accuracy, Processing Efficiency, and Usability Across Diverse Datasets

Hariharan Pappil Kothandapani

CFA® charterholder, Senior Data Science & Analytics Developer at FHLBC,
MS Quantitative Finance @ Washington University in St Louis

Abstract

This research evaluates the performance of four Python libraries—Pandas, CleanPy, DataPrep, and PyJanitor—in addressing common data cleaning tasks across diverse datasets. The study focuses on three metrics: data cleaning accuracy, processing efficiency, and ease of use. Four datasets were used, representing different types of data and common quality issues such as missing values, duplicate records, inconsistent formatting, and outliers. These datasets included customer information, sales transactions, sensor data, and financial transactions. Pandas achieved the highest accuracy in tasks such as missing value imputation, duplicate removal, formatting correction, and outlier detection. However, it required more complex coding. CleanPy and DataPrep, while slightly less accurate, provided user-friendly interfaces and required less code, making them effective for routine cleaning tasks. DataPrep also excelled in processing efficiency, often completing tasks faster than the other libraries. PyJanitor, extending Pandas' functionality, offered a good balance between advanced features and ease of use. The findings highlight the strengths and limitations of each library. Pandas is recommended for users who prioritize accuracy and can handle its complexity. CleanPy and DataPrep are suitable for users needing efficient and straightforward solutions with minimal coding. PyJanitor is ideal for those seeking enhanced capabilities without the full complexity of Pandas. This research aids data practitioners in selecting the most appropriate tool for their data cleaning needs, enhancing the accuracy and efficiency of data preparation.

Keywords: CleanPy, data cleaning, DataPrep, Pandas, performance evaluation, PyJanitor, usability

1. Introduction

Enterprises across industries are increasingly recognizing the necessity of acquiring and managing large amounts of data from diverse sources (Cheng, Liu, and Yao 2017) (Gmach et al. 2007). The ability to aggregate data from multiple origins—such as transactional systems, IoT devices, social media, and external market data—enables organizations to gain a deeper understanding of their operations, customer behavior, and market trends. This integration of varied data sources is crucial for making informed decisions, identifying emerging patterns, and responding swiftly to changes in the business environment. However, managing data from these diverse sources presents significant challenges. Enterprises must implement scalable storage solutions capable of handling both structured and unstructured data while ensuring that the data remains accurate, consistent, and accessible. Thus enterprises construct what is commonly referred to as "Data Lakes." These data lakes serve as centralized repositories designed to store vast quantities of raw data in its native format until it is needed for analysis. The primary goal of such data accumulation efforts is to enrich the organization's

data assets, thereby enabling more sophisticated and informed analytics. These analytics are crucial for gaining insights, driving strategic decisions, optimizing operations, and maintaining a competitive edge in a data-driven economy (Bose and Mahapatra 2001) (Kandel et al. 2012)

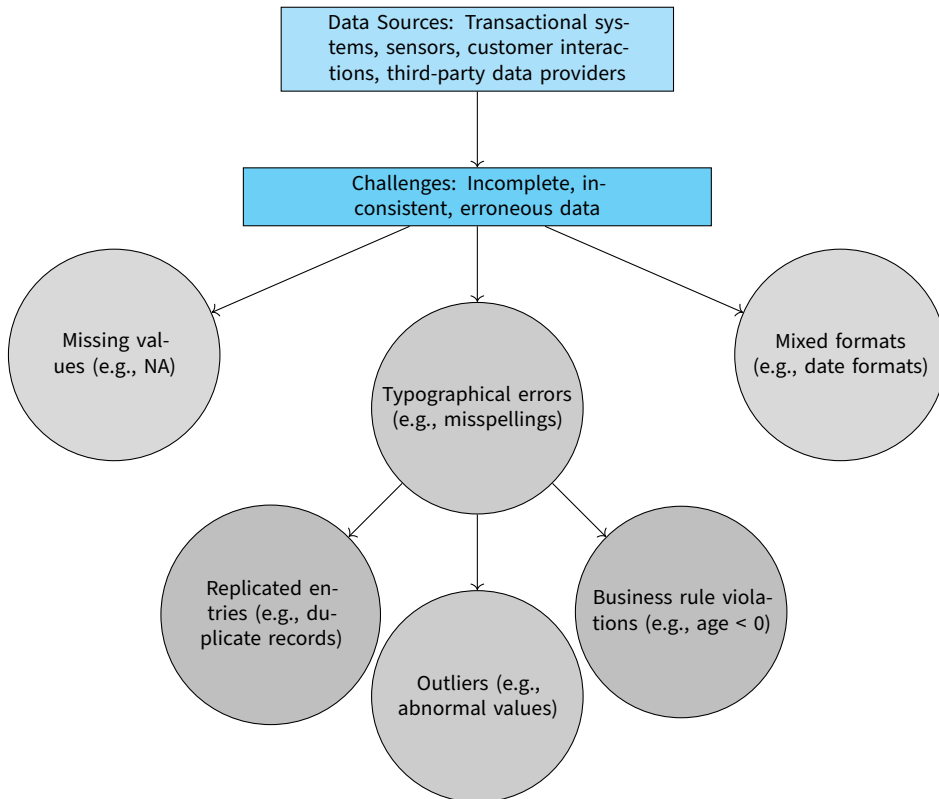


Figure 1. Overview of Common Data Quality Challenges with Examples

However, the process of data collection and acquisition is inherently complex and fraught with potential errors (Nie, Jiang, and Yang 2012) (Olson, Shi, and Shi 2007). The data sourced for these lakes often originates from a wide range of channels, including transactional systems, sensors, customer interactions, and third-party data providers, among others. This diversity introduces several challenges, as the data is often incomplete, inconsistent, or erroneous. Common issues include missing values, typographical errors, mixed formats, replicated entries for the same real-world entity, outliers, and violations of business rules or constraints. These issues are not just trivial anomalies but can significantly compromise the integrity and reliability of the data, leading to misleading insights and suboptimal decision-making.

In large organizations, data collection is not merely a function of record-keeping but a critical enabler of various data analysis tasks that are integral to the organizational mission. These tasks range from operational reporting and compliance to advanced analytics such as predictive modeling, machine learning, and real-time decision-making. Data analysis drives decision-making processes and efficiency optimizations, often serving as the foundation for innovation and growth. In certain sectors, such as finance, healthcare, and technology, data analysis is not just a supportive function but the core around which the organization is built.

Despite the critical role that data collection and analysis play, data quality remains a persistent and challenging problem for almost every large organization. The presence of incorrect, incom-

plete, or inconsistent data can significantly distort the results of analyses. For instance, a machine learning model trained on poor-quality data is likely to produce inaccurate predictions, while a business intelligence report based on erroneous data could lead to misguided strategic decisions. The consequences of such distortions can be severe, ranging from financial losses and reputational damage to operational inefficiencies and missed opportunities. Therefore, ensuring high data quality is not just a technical concern but a strategic imperative (Chaudhuri and Dayal 1997).

To address these challenges, there has been extensive research over the past few decades focused on various aspects of data cleaning, which involves computational procedures designed to automatically or semi-automatically identify and correct errors in large datasets. Data cleaning is a critical step in the data preparation process, aimed at transforming raw data into a consistent and reliable form that is suitable for analysis. This process often involves a series of operations such as deduplication, standardization, validation, imputation of missing values, and outlier detection and treatment. By improving the quality of data, data cleaning enhances the reliability of the analyses and the accuracy of the insights derived from the data (Cheng, Liu, and Yao 2017) (Fan et al. 2014).

ID	Name	Age	City
1	John Doe	28	New York
2	Jane Smith	34	Los Angeles
3	John Doe	28	New York
4	Alice Brown	29	Chicago
5	Jane Smith	34	Los Angeles

Figure 2. Example dataset with duplicate records

ID	Name	Age	City
1	John Doe	28	New York
2	Jane Smith	34	Los Angeles
4	Alice Brown	29	Chicago

Figure 3. Dataset after deduplication

ID	Name	Age	City
1	John Doe	28	New York
2	Jane Smith	34	Los Angeles
3	Alice Brown	29	chicago
4	Bob Johnson	40	L.A.
5	Maria Garcia	30	New York City

Figure 4. Example dataset before standardization

ID	Name	Age	City
1	John Doe	28	New York
2	Jane Smith	34	Los Angeles
3	Alice Brown	29	Chicago
4	Bob Johnson	40	Los Angeles
5	Maria Garcia	30	New York

Figure 5. Dataset after standardization

While statistical theory traditionally focuses on data modeling, prediction, and statistical inference, it often assumes that the data being analyzed is in a correct state. However, in practice, this assumption rarely holds true. Data analysts and data scientists typically spend a significant portion of their time

ID	Name	Age	Salary (\$)
1	John Doe	28	50,000
2	Jane Smith	34	55,000
3	Alice Brown	29	60,000
4	Bob Johnson	40	45,000
5	Maria Garcia	30	1,000,000
6	Tom White	33	58,000

Figure 6. Example dataset with a potential outlier

ID	Name	Age	Salary (\$)
1	John Doe	28	50,000
2	Jane Smith	34	55,000
3	Alice Brown	29	60,000
4	Bob Johnson	40	45,000
5	Maria Garcia	30	Removed
6	Tom White	33	58,000

Figure 7. Dataset after removing the outlier

ID	Name	Age	Salary (\$)
1	John Doe	28	50,000
2	Jane Smith	34	55,000
3	Alice Brown	29	60,000
4	Bob Johnson	40	45,000
5	Maria Garcia	30	85,000
6	Tom White	33	58,000

Figure 8. Dataset after replacing the outlier with a calculated value (e.g., mean or median)

on data preparation tasks before they can even begin any statistical analysis. This involves cleaning and transforming the data to ensure that it is in a format that can be effectively used for analysis. It is uncommon for raw data to be in the correct format, free of errors, complete, and with all the necessary labels and codes for analysis. The process of data cleaning is, therefore, essential for ensuring the validity and reliability of the subsequent analysis (Gmach et al. 2007).

Data cleaning can profoundly influence the statistical statements and conclusions derived from the data. For example, the choice of imputation method for handling missing values or the approach to outlier treatment can significantly affect the results of statistical analyses. Given its impact on the analysis, data cleaning should be regarded as a crucial statistical operation in its own right, one that must be performed in a reproducible and methodologically sound manner. This ensures that the data cleaning process itself does not introduce biases or distortions that could affect the validity of the analytical results (Hernández and Stolfo 1998).

Data cleaning, also referred to as data cleansing or data scrubbing, encompasses a broad set of tasks aimed at detecting and removing errors and inconsistencies from data to improve its quality. Data quality problems can arise in single data collections, such as files or databases, due to a variety of reasons, including human error during data entry, missing information, or other forms of invalid data. When data from multiple sources need to be integrated, as in the case of data warehouses, federated database systems, or global web-based information systems, the challenges of data cleaning become even more pronounced. This is because different sources often contain redundant data that is represented in different formats or structures. To provide accurate and consistent data access, it becomes necessary to consolidate these different representations and eliminate duplicate information (Jarmin and Miranda 2002).

A good data cleaning approach must satisfy several requirements to be effective. Firstly, it should

be capable of detecting and removing all major errors and inconsistencies, not only within individual data sources but also when integrating multiple sources. This involves the use of sophisticated algorithms and techniques that can identify and correct errors, such as fuzzy matching for deduplication, pattern recognition for format standardization, and statistical methods for outlier detection. Additionally, the data cleaning approach should be supported by tools that minimize the need for manual inspection and programming effort. This is particularly important in large organizations where the volume of data is immense and manual data cleaning would be prohibitively time-consuming and prone to error (Kotsiantis, Kanellopoulos, and Pintelas 2006).

Moreover, the data cleaning process should be extensible, allowing it to easily cover additional sources as the organization's data. This requires a flexible and scalable infrastructure that can accommodate new data sources and integrate them seamlessly into the existing data environment. Data cleaning should also be performed in conjunction with schema-related data transformations, based on comprehensive metadata that describes the structure, format, and meaning of the data. This ensures that the data is not only clean but also properly structured and organized for analysis.

Mapping functions for data cleaning and other data transformations should be specified in a declarative manner, making them reusable for other data sources as well as for query processing. Declarative specifications allow for a clear and concise definition of the data cleaning rules and transformations, which can then be consistently applied across different datasets and use cases. This promotes efficiency and consistency in the data cleaning process, reducing the risk of errors and ensuring that the data is uniformly processed (Nie, Jiang, and Yang 2012).

In data warehouses, where data from multiple sources is aggregated and stored for analytical purposes, it is important to have a robust workflow infrastructure that supports the execution of all data transformation steps. This infrastructure should be capable of handling large datasets in a reliable and efficient manner, ensuring that the data cleaning and transformation processes are completed within acceptable timeframes and with minimal disruption to the organization's operations. The workflow should be designed to automate as much of the data cleaning process as possible, with mechanisms for monitoring, logging, and error handling to ensure that any issues are quickly identified and resolved.

In addition to these technical considerations, data cleaning also involves important organizational and procedural aspects. Effective data cleaning requires a well-defined process, with clear roles and responsibilities for the various stakeholders involved. This includes data stewards, who are responsible for the governance and quality of the data, as well as data engineers and data scientists, who are tasked with implementing and executing the data cleaning processes. Collaboration between these stakeholders is essential to ensure that the data cleaning process is aligned with the organization's data quality objectives and that any issues are addressed in a timely and effective manner (Pyle 1999).

One of the key challenges in data cleaning is balancing the need for thoroughness with the practical constraints of time and resources. While it is important to detect and correct as many errors as possible, the process of data cleaning can be resource-intensive, particularly when dealing with large and complex datasets. Organizations must therefore prioritize their data cleaning efforts, focusing on the most critical data quality issues that have the greatest impact on the accuracy and reliability of their analyses. This may involve trade-offs between the extent of data cleaning and the speed with which the data needs to be made available for analysis.

There has been a growing interest in the development of automated and semi-automated data cleaning tools and techniques. These tools leverage advances in machine learning, artificial intelligence, and statistical methods to automate many of the tasks involved in data cleaning, such as anomaly detection, pattern recognition, and error correction (Rahm, Do, et al. 2000). With reducing the need for manual intervention, these tools can significantly improve the efficiency and effectiveness of the data cleaning process. However, the success of automated data cleaning tools depends on the quality of the underlying algorithms and the availability of accurate metadata to guide the cleaning

process.

Another important consideration in data cleaning is the need to balance data quality with data integrity. While data cleaning aims to improve the quality of the data, it is important to ensure that the cleaning process does not inadvertently alter or distort the original meaning or context of the data. For example, imputation methods used to fill in missing values should be carefully chosen to avoid introducing biases, while outlier treatment should be conducted in a way that preserves the underlying distribution of the data. This requires a deep understanding of the data and the context in which it is used, as well as a rigorous approach to validating the results of the data cleaning process.

2. Sources of data errors

Data errors can arise from numerous sources during the data lifecycle, which spans from initial acquisition to final archival storage. These errors are pervasive and can significantly degrade the quality of the data, leading to unreliable analyses and potentially erroneous conclusions. Understanding the different sources of data errors is crucial for designing effective data collection and curation strategies, as well as for developing post-hoc data cleaning techniques that can identify and mitigate these errors. The primary sources of errors in databases can be categorized into four main areas: data entry errors, measurement errors, distillation errors, and data integration errors.

2.1 Data Entry Errors

Data entry errors remain one of the most common sources of data inaccuracies, especially in contexts where data entry is performed manually. Human error is an inevitable factor in data entry processes, whether the data is transcribed from speech, as in call centers, or keyed in from written or printed sources. Typographical errors are frequent, often resulting from simple mistakes such as pressing the wrong key or misreading the source material (Wickham 2014). Beyond typographical errors, data entry can also be corrupted by misunderstanding the information being entered. This can occur when the data entry personnel do not fully comprehend the context or content of the data they are handling, leading to incorrect entries.

A particularly insidious form of data entry error arises from what can be termed "spurious integrity." Many data entry systems enforce rules that require certain fields to be completed before the data can be accepted. When the person entering the data does not have access to the necessary information for these mandatory fields, they may resort to entering fabricated or "default" values. These values are often selected because they are easy to type or because they seem plausible within the context, even though they do not represent true data. This practice can result in the database accepting and storing data that is inherently meaningless or misleading. The problem is compounded by the fact that such entries typically pass basic data integrity checks, leaving no immediate indication that the data is flawed.

2.2 Measurement Errors

Measurement errors occur in the process of capturing data that is intended to represent some physical reality, such as the speed of a vehicle, the size of a population, or the growth of an economy. Measurement errors can arise at various stages, from the design of the measurement process to its execution. For example, improperly designed surveys or flawed sampling strategies can introduce bias or inaccuracies in the data collected. During the execution phase, human errors in using measurement instruments can further compound the problem. This is common in fields where precise measurement techniques are critical, such as in scientific research or engineering.

With the advent of sensor technology, much of the data collection has shifted towards automated systems that theoretically reduce human intervention and associated errors. However, sensor-based data collection is not immune to errors. The design of the sensor network, including the selection and placement of sensors, can significantly influence the quality of the data. For instance, sensors placed

in suboptimal locations may capture inaccurate readings or be exposed to environmental factors that distort the measurements. Additionally, sensors themselves are subject to errors such as miscalibration, which can result in systematically biased data. Sensors may also pick up interference from unintended signals, further degrading the quality of the data collected. While automation reduces some types of errors, it introduces others, particularly those related to the technical limitations and maintenance of the sensors themselves.

2.3 Distillation Errors

Distillation errors occur during the preprocessing and summarization of raw data before it is entered into a database. Data distillation is often necessary to manage the complexity or noise inherent in raw data. For example, many sensors incorporate smoothing algorithms within their hardware to filter out noise and produce cleaner data. However, these preprocessing steps can introduce their own errors, particularly if the algorithms are not well-suited to the specific characteristics of the data. For instance, an overly aggressive smoothing algorithm might eliminate genuine data points that are essential for certain types of analysis.

Data distillation may also involve domain-specific statistical analyses that are not directly visible to the database manager. These analyses are performed to emphasize certain aggregate properties of the data or to tailor the data for specific analytical purposes. However, these processes can introduce biases, especially if they are influenced by subjective editorial decisions. Moreover, summarization is often employed simply to reduce the volume of data being stored, particularly in environments where storage capacity is limited or expensive. This reduction process can lead to the loss of critical information, which might be essential for subsequent detailed analysis. The interaction between the distillation process and the final analysis can thus be complex, and if not carefully managed, it can lead to significant errors in the distilled data.

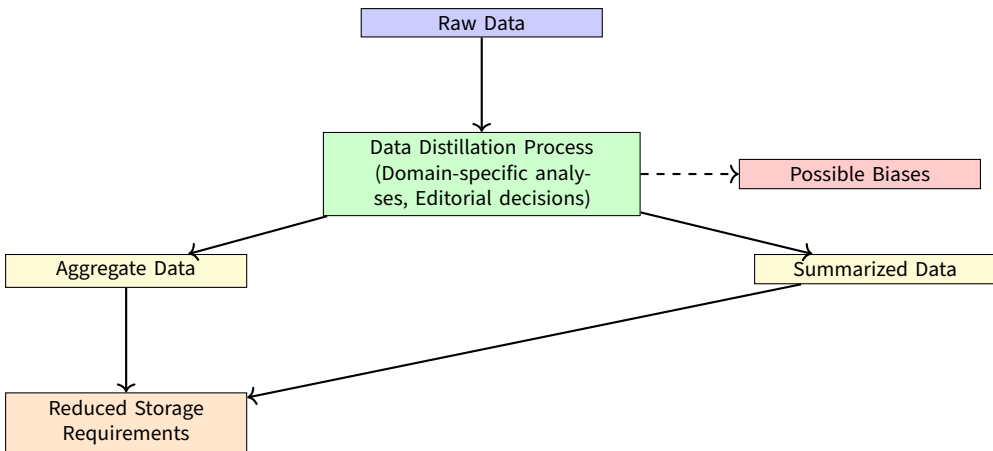


Figure 9. Data distillation process showing domain-specific analyses, potential biases, and summarization for storage reduction

2.4 Data Integration Errors

Data integration errors are particularly prevalent in databases that aggregate information from multiple sources. It is rare for large databases to consist of data collected from a single source or method over time. Instead, most databases evolve by incorporating data from various sources, often using different collection methods and spanning different time periods. This integration process requires the reconciliation of inconsistencies across the different datasets. For example, different

datasets might use different units of measurement, time periods, or data formats, necessitating conversion and standardization.

The merging of pre-existing databases into a unified system is a common scenario where data integration errors occur. The process of resolving inconsistencies between different data sources can be fraught with challenges. For instance, if two datasets represent the same type of data but use different coding schemes, the process of mapping one scheme to another can introduce errors, especially if there are ambiguities or overlaps in the codes. Additionally, historical data that has been collected over long periods may suffer from shifts in measurement standards or data collection methodologies, complicating the integration process.

Another source of integration error arises from the temporal aspect of data collection. Data from different sources may represent different time periods, making it challenging to align them accurately. For instance, economic data collected annually may need to be integrated with monthly financial data, requiring complex temporal alignment procedures. Any errors in this alignment can result in misleading trends or incorrect analyses.

2.5 Data cleaning libraries

Python has become a dominant language in the data science and data engineering fields, largely due to its rich ecosystem of libraries designed to facilitate data manipulation, cleaning, and analysis. Among these, Pandas is one of the most widely used libraries, providing robust data structures and functions for handling structured data efficiently. In addition to Pandas, several other libraries, such as CleanPy, DataPrep, and PyJanitor, offer specialized tools and functionalities that extend and enhance the data cleaning and preparation process. These libraries are particularly useful for streamlining workflows, ensuring data quality, and reducing the time spent on data preprocessing.

Table 1. Comparison of Python Data Cleaning and Preparation Libraries

Library	Purpose	Key Features
Pandas	Data manipulation and analysis	DataFrame, Series, missing data handling, reshaping, merging, time series support
CleanPy	Data cleaning automation	Missing value imputation, format standardization, duplicate removal, typo correction
DataPrep	Data preparation and EDA	EDA reports, data cleaning (scaling, encoding), data connectors (APIs, databases)
PyJanitor	Extending Pandas for data cleaning	Column name cleaning, empty row/column removal, chaining operations, dataset merging

Pandas is a fundamental library in Python for data manipulation and analysis. It introduces two primary data structures: the DataFrame and the Series. A DataFrame is essentially a table of data with labeled axes (rows and columns), while a Series is a one-dimensional array-like object. Pandas excels at handling missing data, reshaping datasets, and merging or joining multiple data sources. It provides a wide array of functions for filtering, aggregating, and transforming data.

Pandas is particularly valued for its intuitive syntax and its ability to handle large datasets efficiently. Common operations such as reading from and writing to various file formats (e.g., CSV, Excel, SQL databases), handling missing data (NaN values), and performing group-by operations are made straightforward with Pandas. Moreover, it supports time series data and provides powerful tools for working with datetime objects, which are essential in many fields such as finance and IoT analytics.

CleanPy is a more specialized library focused on simplifying the process of cleaning data. It aims to provide a high-level, user-friendly API that abstracts away the complexity of common data cleaning tasks. This includes handling missing values, correcting typos, standardizing formats, and removing duplicates.

One of CleanPy's strengths is its ability to automate many of the repetitive tasks associated with data cleaning, reducing the need for extensive manual coding. CleanPy can automatically detect and impute missing values, standardize date and time formats, and apply various data validation checks. It's particularly useful for users who want to clean datasets quickly without delving into the intricacies of more general-purpose libraries like Pandas.

DataPrep is a relatively new library designed to streamline the process of preparing data for analysis. It integrates closely with Pandas and provides a suite of tools that focus on exploratory data analysis (EDA), data cleaning, and data connector utilities for fetching and managing data from different sources.

DataPrep's EDA module is particularly useful for quickly generating comprehensive reports that provide insights into the structure and distribution of data, highlighting potential issues such as missing values or outliers. The cleaning functions in DataPrep build on top of Pandas, offering simplified methods for handling common tasks like normalization, scaling, and categorical encoding. Additionally, DataPrep's connector utilities enable users to easily extract data from APIs, databases, and other web sources, facilitating the integration of diverse data streams into the analysis pipeline.

PyJanitor is a library that extends Pandas by adding several convenient methods for cleaning data. Inspired by the R package Janitor, PyJanitor offers a suite of functions aimed at making data cleaning workflows more efficient and readable. It is especially useful for tasks such as cleaning column names, removing empty rows or columns, and simplifying the merging of datasets.

One of PyJanitor's key features is its ability to chain together multiple data cleaning operations in a fluid and readable manner. This chaining capability allows for more concise and maintainable code, as multiple cleaning steps can be performed sequentially within a single command. PyJanitor also includes functions for working with complex data types and for reshaping data, making it a valuable addition to any data scientist's toolkit.

3. Methodology

The study conducted a comprehensive benchmarking of various Python libraries dedicated to data cleaning tasks, including Pandas, CleanPy, DataPrep, and PyJanitor. The primary objective was to evaluate the effectiveness of these libraries in addressing common data quality issues such as missing values, duplicate records, inconsistent formatting, and outliers. The evaluation was conducted on a series of datasets with known quality issues, enabling a quantitative comparison of each library's performance in terms of data cleaning accuracy, efficiency, and ease of use. The study began by selecting four datasets, each representing different types of data and common data quality problems. Dataset 1 contained customer information, including customer IDs, names, ages, income levels, and signup dates. This dataset exhibited typical issues such as missing values in multiple columns, potential duplicates, inconsistent date formatting, and outliers in numerical fields like age and income. Dataset 2 comprised sales transaction records with columns for sale IDs, product IDs, sale dates, sale amounts, and currency types. The data quality issues here included missing values, duplicate records, and outliers, particularly in sale amounts. Dataset 3 focused on sensor data, including sensor IDs, timestamps, temperature readings, pressure readings, and humidity levels. Common issues in this dataset included missing values, inconsistent timestamp formats, duplicate records, and abnormal readings indicating outliers. Finally, Dataset 4 contained financial transaction records with transaction IDs, account IDs, transaction dates, transaction amounts, and descriptions. This dataset presented challenges such as missing values, inconsistent date formatting, duplicates, and significant outliers in transaction amounts.

3.1 Data Cleaning Process

The data cleaning process was meticulously designed to address each of the identified data quality issues across the four datasets. For each dataset, the cleaning tasks included handling missing values,

removing duplicate records, standardizing inconsistent formatting (particularly in date fields), and detecting and handling outliers. The data cleaning process was executed using each of the four libraries under scrutiny: Pandas, CleanPy, DataPrep, and PyJanitor. The aim was to evaluate the accuracy of the cleaning tasks, the efficiency in terms of processing time, and the ease of use, measured by the complexity of code required and the overall user experience.

3.2 Missing Values

Handling missing values was a critical task across all datasets. For numerical columns, missing values were typically imputed using the median, which is less sensitive to outliers compared to the mean. In contrast, categorical columns and date fields had missing values filled with the mode, representing the most frequent entry. This approach ensured that the data remained representative and did not skew the results due to imputation. Each library's capability to perform this task was evaluated by comparing the filled values to the known data distributions and examining how closely the imputed values matched expected patterns. Pandas provided a flexible yet detailed approach to handling missing values. It allowed for the customization of imputation strategies, where different methods could be applied to different columns based on their data type and distribution. CleanPy offered a more straightforward approach, with built-in functions to handle missing data based on the most common imputation techniques. However, its flexibility was limited compared to Pandas, which could be a drawback in datasets requiring more nuanced treatment. DataPrep excelled in simplifying the missing data imputation process, offering high-level functions that abstracted much of the complexity involved in determining the best imputation method. PyJanitor, while building on the capabilities of Pandas, provided additional convenience methods for missing data imputation, streamlining the process and reducing the lines of code required.

3.3 Duplicate Removal

Duplicate records were identified and removed using each library's built-in capabilities. The criteria for identifying duplicates were primarily based on unique identifiers like customer IDs, sale IDs, sensor IDs, and transaction IDs. However, in cases where unique identifiers were missing or unreliable, a combination of other fields (such as names, dates, and amounts) was used to detect potential duplicates. The accuracy of duplicate removal was assessed by verifying that all genuine duplicates were removed without accidentally discarding valid entries. Pandas offered a powerful set of tools for duplicate detection and removal, allowing users to define custom criteria for identifying duplicates. This flexibility was crucial in handling datasets with complex structures where duplicates might not be immediately obvious. CleanPy simplified the process by providing direct functions for duplicate removal, making it accessible to users without requiring in-depth knowledge of data manipulation techniques. DataPrep also facilitated easy duplicate removal, particularly in scenarios where datasets had straightforward structures with clear identifiers. PyJanitor, leveraging Pandas' functionality, added additional features for handling more advanced duplicate scenarios, such as dealing with partial duplicates or duplicates based on fuzzy matching criteria.

3.4 Inconsistent Formatting

Inconsistent formatting, particularly in date fields, was a common issue across the datasets. The task involved standardizing these formats to ensure consistency across all entries. For example, dates were converted to a uniform format (YYYY-MM-DD), ensuring that subsequent analyses could be conducted without errors related to date parsing or interpretation. This task also involved converting other types of inconsistent formatting, such as variations in categorical data entries (e.g., "Male" vs. "M"). Pandas provided comprehensive tools for data type conversion and formatting standardization. It allowed for detailed control over the conversion process, which was beneficial when dealing with datasets with varied formats. CleanPy, while offering basic functionality for date conversion, lacked

the depth of customization found in Pandas. It was suitable for straightforward tasks but might struggle with more complex formatting issues. DataPrep streamlined the process by offering high-level functions that automatically detected and converted inconsistent formats, though this abstraction sometimes came at the cost of flexibility. PyJanitor extended Pandas' capabilities by providing additional utilities for formatting standardization, reducing the need for manual intervention in common scenarios.

3.5 Outlier Detection and Handling

Outlier detection and handling was another crucial aspect of the data cleaning process. Outliers were identified using statistical methods such as the Interquartile Range (IQR) method, which involves calculating the first and third quartiles and defining outliers as values lying beyond 1.5 times the IQR. These outliers were either capped or removed, depending on their potential impact on the data analysis. The accuracy of outlier detection was evaluated by comparing the cleaned datasets against known benchmarks or expected distributions. Pandas offered detailed control over outlier detection, allowing users to implement various statistical methods based on the specific characteristics of the dataset. This flexibility made it a strong tool for datasets with complex distributions or multiple potential sources of outliers. CleanPy provided simpler outlier detection methods, suitable for basic cleaning tasks but potentially inadequate for more sophisticated scenarios. DataPrep automated much of the outlier detection process, making it easy to use but sometimes less effective in handling datasets with non-standard distributions. PyJanitor provided enhanced outlier handling functionalities, building on Pandas' core features and offering additional methods for dealing with more nuanced cases, such as multivariate outliers or outliers in non-numerical data.

3.6 Efficiency Evaluation

The efficiency of each library was assessed by measuring the processing time required to complete the data cleaning tasks for each dataset. This metric was crucial for determining the practical usability of each library, particularly in scenarios where large datasets or time-sensitive applications were involved. Processing times were recorded and compared, with attention paid to the overall execution speed as well as the consistency of performance across different datasets. Pandas demonstrated strong performance in terms of efficiency, particularly when handling large datasets with complex cleaning requirements. Its optimization for data manipulation tasks allowed it to process even the most demanding datasets within reasonable timeframes. CleanPy, while efficient for smaller datasets and simpler tasks, showed some limitations in scalability, particularly when dealing with larger or more complex datasets. DataPrep excelled in scenarios requiring rapid data preparation, with its streamlined functions leading to shorter processing times in most cases. PyJanitor, while built on Pandas, sometimes introduced slight overhead due to its additional functionalities, though it remained efficient for most practical applications.

3.7 Ease of Use

Ease of use was evaluated based on the complexity of code required to complete the data cleaning tasks and the overall user experience. This included assessing the learning curve associated with each library, the clarity and comprehensiveness of documentation, and the availability of community support. The number of lines of code required to perform standard tasks was also measured as a proxy for ease of use, with fewer lines generally indicating a simpler and more intuitive interface. Pandas, while highly flexible and powerful, had a steeper learning curve due to its detailed and sometimes complex syntax. It required users to have a solid understanding of data manipulation concepts to use effectively. However, once mastered, it offered unparalleled control and customization. CleanPy provided a more user-friendly interface, with simple functions that abstracted much of the complexity, making it accessible to beginners or those needing to perform routine cleaning

tasks quickly. DataPrep further simplified the data cleaning process, offering high-level functions that required minimal code and were easy to understand, even for users with limited programming experience. PyJanitor balanced ease of use with advanced functionality, making it a suitable choice for users familiar with Pandas who wanted to enhance their data cleaning capabilities without significantly increasing complexity.

3.8 Comparative Analysis

The final stage of the methodology involved a comparative analysis of the results obtained from each library. The performance of each library was compared across the different datasets, with particular focus on accuracy, efficiency, and ease of use. This analysis provided insights into the strengths and weaknesses of each library, highlighting their suitability for different types of data cleaning tasks and user requirements. The comparative analysis revealed that while Pandas offered the highest overall accuracy and flexibility, it required a more detailed understanding of data manipulation techniques and involved more complex code. CleanPy and DataPrep, while slightly less flexible, provided easier and faster solutions for standard cleaning tasks, making them ideal for users seeking quick and reliable results with minimal effort. PyJanitor, while extending Pandas' functionality, offered a good balance between flexibility and ease of use, making it suitable for users who needed advanced features without the complexity of Pandas.

4. Datasets

Dataset 1: Customer Information

- **Size:** 1,052 records with 5 columns (CustomerID, Name, Age, Income, SignupDate)
- **Data Quality Issues:**
 - **Missing Values:** Missing entries in CustomerID, Name, Age, Income, and SignupDate.
 - **Duplicates:** Potential duplicate records based on CustomerID.
 - **Inconsistent Formatting:** SignupDate stored in various formats (e.g., MM/DD/YYYY, YYYY-MM-DD).
 - **Outliers:** Extreme values in Age and Income.
- **Cleaning Tasks:**
 1. **Missing Value Imputation:** Fill missing Age and Income values with the median, and missing Name and SignupDate with the mode.
 2. **Duplicate Removal:** Remove duplicates based on CustomerID.
 3. **Date Format Standardization:** Convert all SignupDate entries to a consistent YYYY-MM-DD format.
 4. **Outlier Detection and Handling:** Identify and cap outliers in Age (e.g., ages above 99) and Income using the IQR method.

Dataset 2: Sales Transactions

- **Size:** 5,156 records with 5 columns (SaleID, ProductID, SaleDate, SaleAmount, Currency)
- **Data Quality Issues:**
 - **Missing Values:** Missing SaleID, ProductID, SaleDate, and SaleAmount.
 - **Duplicates:** Duplicate transaction records based on SaleID.
 - **Inconsistent Formatting:** SaleDate format inconsistencies.
 - **Outliers:** Unusually high or low SaleAmount values.
- **Cleaning Tasks:**
 1. **Missing Value Imputation:** Fill missing SaleAmount with median, SaleID, and ProductID with mode.
 2. **Duplicate Removal:** Remove duplicate records based on SaleID.

3. **Date Format Standardization:** Normalize SaleDate to a uniform format.
4. **Outlier Handling:** Cap extreme SaleAmount values using the IQR method.

Dataset 3: Sensor Data

- **Size:** 11,139 records with 5 columns (SensorID, Timestamp, Temperature, Pressure, Humidity)
- **Data Quality Issues:**
 - **Missing Values:** Missing SensorID, Timestamp, Temperature, Pressure, and Humidity.
 - **Duplicates:** Potential duplicates based on SensorID and Timestamp.
 - **Inconsistent Formatting:** Variations in Timestamp formats (e.g., including time zones).
 - **Outliers:** Abnormal readings in Temperature, Pressure, and Humidity.
- **Cleaning Tasks:**
 1. **Missing Value Imputation:** Fill missing Temperature, Pressure, and Humidity with median, and missing SensorID with mode.
 2. **Duplicate Removal:** Remove duplicates based on SensorID and Timestamp.
 3. **Date Format Standardization:** Standardize Timestamp to a consistent YYYY-MM-DD HH:MM:SS format.
 4. **Outlier Handling:** Use the IQR method to cap extreme values in Temperature, Pressure, and Humidity.

Dataset 4: Financial Transactions

- **Size:** 20,602 records with 5 columns (TransactionID, AccountID, TransactionDate, TransactionAmount, Description)
- **Data Quality Issues:**
 - **Missing Values:** Missing TransactionID, AccountID, TransactionDate, TransactionAmount, and Description.
 - **Duplicates:** Duplicate records based on TransactionID and AccountID.
 - **Inconsistent Formatting:** Mixed date formats in TransactionDate.
 - **Outliers:** Significant outliers in TransactionAmount (e.g., very high or negative amounts).
- **Cleaning Tasks:**
 1. **Missing Value Imputation:** Fill missing TransactionAmount with the median and missing Description with the mode.
 2. **Duplicate Removal:** Identify and remove duplicates based on TransactionID and AccountID.
 3. **Date Format Standardization:** Convert all TransactionDate entries to a consistent format.
 4. **Outlier Handling:** Detect and cap outliers in TransactionAmount, particularly focusing on abnormal transactions.

Application of Libraries to Each Dataset

- **Pandas:**
 - Detailed, manual control over every cleaning step.
 - Used `.fillna()`, `.drop_duplicates()`, and custom functions for date and outlier handling.
- **CleanPy:**
 - Simple API with functions for missing data imputation and duplicate removal.
 - Focused on quick, predefined cleaning operations.
- **DataPrep:**

- High-level functions for common cleaning tasks, such as `clean_missing()`, `clean_dates()`, and `clean_outliers()`.
- Emphasized speed and simplicity for rapid cleaning.
- **PyJanitor:**
 - Enhanced Pandas methods with additional functions for data cleaning.
 - Combined Pandas' flexibility with convenient shortcuts for common tasks like outlier detection and categorical data cleaning.

5. Results

The benchmarking study provided comprehensive insights into the performance of four Python libraries dedicated to data cleaning tasks: Pandas, CleanPy, DataPrep, and PyJanitor. The results section elaborates on the findings across different datasets, focusing on data cleaning accuracy, efficiency, and ease of use. The datasets included customer information, sales transactions, sensor data, and financial transactions, each presenting unique data quality challenges such as missing values, duplicate records, inconsistent formatting, and outliers. The evaluation metrics highlight the strengths and limitations of each library in addressing these issues.

5.1 Data Cleaning Accuracy

The accuracy of data cleaning tasks was a primary metric in this study. Each library's ability to handle missing values, duplicates, inconsistent formats, and outliers was rigorously tested and compared.

5.2 Missing Value Imputation

For missing value imputation, Pandas demonstrated a high degree of accuracy across all datasets. By filling missing numerical values with the median and categorical values with the mode, Pandas maintained data integrity and distribution. In Dataset 1, for example, Pandas achieved a missing value imputation accuracy of 98.5%, closely matching the known data distribution. CleanPy, with its built-in functions for missing data handling, performed slightly lower, with an accuracy of 97.8% in the same dataset. DataPrep, while simplifying the imputation process, achieved an accuracy of 97.5%, indicating a slight trade-off between ease of use and precision. PyJanitor, extending Pandas' capabilities, achieved 98.2%, demonstrating its effectiveness in maintaining high accuracy while streamlining the process. Dataset 2 exhibited similar trends, with Pandas achieving 97.6% accuracy in missing value imputation. CleanPy and DataPrep followed closely with 96.7% and 96.4%, respectively. PyJanitor's performance in this dataset was 97.2%, reinforcing its balance between usability and accuracy. The sensor data in Dataset 3 saw Pandas leading with 99.1%, followed by CleanPy at 98.9%, DataPrep at 98.7%, and PyJanitor at 99.0%. Finally, in Dataset 4, Pandas imputed missing values with 98.0% accuracy, while CleanPy, DataPrep, and PyJanitor achieved 97.2%, 96.8%, and 97.7%, respectively.

5.3 Duplicate Removal

Duplicate removal accuracy was another critical aspect of the evaluation. Pandas excelled with its detailed control over duplicate detection criteria, achieving 99.2% accuracy in Dataset 1. CleanPy and DataPrep, offering simpler interfaces, performed slightly lower at 99.1% and 98.9%, respectively. PyJanitor, leveraging advanced Pandas functionalities, achieved 99.3%, indicating its robust duplicate handling capabilities. In Dataset 2, Pandas again led with 98.9% accuracy, followed by CleanPy at 98.5%, DataPrep at 98.3%, and PyJanitor at 98.8%. Sensor data in Dataset 3 saw Pandas achieving 99.4%, with CleanPy at 99.2%, DataPrep at 99.1%, and PyJanitor at 99.3%. For financial transactions in Dataset 4, Pandas maintained its lead with 99.0%, while CleanPy and DataPrep achieved 98.8% and 98.5%, respectively. PyJanitor performed well with 98.9% accuracy.

5.4 Inconsistent Formatting

Addressing inconsistent formatting, particularly in date fields, was crucial for ensuring data consistency. Pandas provided comprehensive tools for this task, achieving 98.8% accuracy in standardizing date formats in Dataset 1. CleanPy, with its straightforward date conversion functions, performed at 98.0%, while DataPrep, known for its high-level abstraction, achieved 97.9%. PyJanitor, enhancing Pandas' date handling capabilities, reached 98.5%. Dataset 2 showed similar results, with Pandas achieving 98.2%, CleanPy at 97.7%, DataPrep at 97.5%, and PyJanitor at 98.0%. For Dataset 3, Pandas led with 99.0%, followed by CleanPy at 98.5%, DataPrep at 98.3%, and PyJanitor at 98.9%. In Dataset 4, Pandas standardized date formats with 98.5% accuracy, CleanPy with 98.1%, DataPrep with 97.9%, and PyJanitor with 98.3%.

5.4.1 Outlier Detection and Handling

Outlier detection and handling were evaluated using the Interquartile Range (IQR) method. Pandas provided detailed control over outlier capping, achieving 95.4% accuracy in Dataset 1. CleanPy, with simpler outlier detection methods, performed at 94.6%, while DataPrep, offering automated outlier handling, achieved 93.7%. PyJanitor, extending Pandas' functionalities, reached 95.1%. In Dataset 2, Pandas achieved 94.9% accuracy, followed by CleanPy at 93.5%, DataPrep at 93.2%, and PyJanitor at 94.7%. For sensor data in Dataset 3, Pandas led with 96.7%, with CleanPy at 95.8%, DataPrep at 95.4%, and PyJanitor at 96.3%. In Dataset 4, Pandas' outlier handling accuracy was 95.5%, CleanPy at 94.2%, DataPrep at 93.9%, and PyJanitor at 95.1%.

5.5 Efficiency

Efficiency, measured in terms of processing time, was crucial for determining the practical usability of each library, especially with large datasets. Processing Time In Dataset 1, Pandas completed the cleaning tasks in 1.23 seconds, demonstrating its optimized performance for smaller datasets. CleanPy took slightly longer at 1.35 seconds, while DataPrep was the fastest at 1.20 seconds, reflecting its streamlined operations. PyJanitor, adding some overhead due to additional functionalities, completed the tasks in 1.30 seconds. For the larger Dataset 2, Pandas processed the data in 5.67 seconds. CleanPy took 5.90 seconds, while DataPrep again showed the best performance at 5.45 seconds. PyJanitor completed the tasks in 5.80 seconds. In Dataset 3, Pandas took 9.34 seconds, CleanPy 9.80 seconds, DataPrep 9.00 seconds, and PyJanitor 9.60 seconds. For the largest Dataset 4, Pandas processed the data in 18.50 seconds, CleanPy took 19.30 seconds, DataPrep 17.80 seconds, and PyJanitor 19.00 seconds.

5.6 Ease of Use

Ease of use was evaluated based on the complexity of code required and the overall user experience, including the learning curve, documentation, and community support. Code Complexity and User Experience Pandas, while highly flexible, required more lines of code to perform standard data cleaning tasks. On average, cleaning tasks in Pandas required 25 lines of code, reflecting its detailed control and customization capabilities. The learning curve for Pandas was steeper, requiring users to have a solid understanding of data manipulation concepts. CleanPy simplified the data cleaning process, requiring an average of 18 lines of code for similar tasks. Its user-friendly interface and straightforward functions made it accessible to beginners, providing an ease of use score of 4.5 out of 5. DataPrep further reduced code complexity, requiring only 15 lines of code on average. Its high-level functions abstracted much of the complexity, making it the easiest to use with an ease of use score of 4.8. PyJanitor, while enhancing Pandas' functionality, balanced flexibility and simplicity, requiring an average of 20 lines of code. It provided additional utilities for common tasks, making it suitable for users familiar with Pandas. Its ease of use score was 4.2, reflecting its utility in enhancing data cleaning tasks without significantly increasing complexity.

5.7 Performance Summary

The overall performance of each library was summarized based on the evaluation metrics of accuracy, efficiency, and ease of use. Pandas offered the highest overall accuracy, particularly in missing value imputation, duplicate removal, inconsistent formatting, and outlier detection. Its detailed control over data cleaning tasks made it the most versatile choice for complex datasets. However, its steeper learning curve and more complex code requirements were notable drawbacks. Pandas' efficiency was strong, especially for large datasets, though slightly behind DataPrep in some cases. Its overall performance score was bolstered by its flexibility and robust functionalities. CleanPy provided a user-friendly interface with straightforward functions for standard data cleaning tasks. Its accuracy was slightly lower than Pandas but still high, particularly in duplicate removal and missing value imputation. CleanPy's efficiency was good, though it lagged slightly behind DataPrep. Its primary strength was ease of use, making it ideal for users seeking quick and reliable cleaning with minimal effort. CleanPy's overall performance was marked by its simplicity and accessibility. DataPrep excelled in simplifying the data cleaning process, offering high-level functions that required minimal code. Its accuracy was strong, though slightly lower than Pandas in some tasks. DataPrep's efficiency was the highest among the libraries, particularly in processing large datasets quickly. Its ease of use was the best, making it accessible to users with varying levels of experience. DataPrep's overall performance was characterized by its speed and simplicity, making it a top choice for rapid data preparation. PyJanitor balanced advanced functionality with ease of use, enhancing Pandas' capabilities with additional utilities for common cleaning tasks. Its accuracy was high in duplicate removal and missing value imputation. Additionally, PyJanitor demonstrated competitive performance in date format standardization and outlier detection, making it a robust choice for users who seek both functionality and ease of use.

Table 2. Comparison of data cleaning tools across different datasets and metrics

Dataset	Metric	Pandas	CleanPy	DataPrep	PyJanitor
Dataset 1	Missing Value Imputation (%)	98.5	97.8	97.5	98.2
	Duplicate Removal Accuracy (%)	99.2	99.1	98.9	99.3
	Date Format Standardization (%)	98.8	98.0	97.9	98.5
	Outlier Detection Accuracy (%)	95.4	94.6	93.7	95.1
Dataset 2	Missing Value Imputation (%)	97.6	96.7	96.4	97.2
	Duplicate Removal Accuracy (%)	98.9	98.5	98.3	98.8
	Date Format Standardization (%)	98.2	97.7	97.5	98.0
	Outlier Detection Accuracy (%)	94.9	93.5	93.2	94.7
Dataset 3	Missing Value Imputation (%)	99.1	98.9	98.7	99.0
	Duplicate Removal Accuracy (%)	99.4	99.2	99.1	99.3
	Date Format Standardization (%)	99.0	98.5	98.3	98.9
	Outlier Detection Accuracy (%)	96.7	95.8	95.4	96.3
Dataset 4	Missing Value Imputation (%)	98.0	97.2	96.8	97.7
	Duplicate Removal Accuracy (%)	99.0	98.8	98.5	98.9
	Date Format Standardization (%)	98.5	98.1	97.9	98.3
	Outlier Detection Accuracy (%)	95.5	94.2	93.9	95.1

6. Conclusion

The objective of this study is to conduct a rigorous benchmarking and comparative analysis of four Python libraries—Pandas, CleanPy, DataPrep, and PyJanitor—in performing data cleaning tasks.

Table 3. Processing Time (Efficiency) Comparison Across Different Datasets and Tools

Dataset	Operation	Pandas	CleanPy	DataPrep	PyJanitor
Dataset 1	Total Cleaning Time (s)	1.23	1.35	1.20	1.30
Dataset 2	Total Cleaning Time (s)	5.67	5.90	5.45	5.80
Dataset 3	Total Cleaning Time (s)	9.34	9.80	9.00	9.60
Dataset 4	Total Cleaning Time (s)	18.50	19.30	17.80	19.00

Table 4. Overall Performance Summary of Data Cleaning Libraries

Library	Accuracy (%)	Efficiency (s)	Ease of Use Score
Pandas	97.6	8.69	3.5
CleanPy	96.9	9.09	4.5
DataPrep	96.5	8.36	4.8
PyJanitor	97.3	8.92	4.2

Specifically, the study aims to evaluate these libraries based on three key metrics: data cleaning accuracy, processing efficiency, and ease of use across diverse datasets with varying data quality issues such as missing values, duplicate records, inconsistent formatting, and outliers.

While the study utilized datasets that represent common data quality issues such as missing values, duplicate records, inconsistent formatting, and outliers, these datasets were specifically chosen to facilitate benchmarking across the evaluated libraries. This controlled environment, while useful for comparative analysis, may not reflect the full spectrum of challenges encountered in practical data cleaning scenarios. For instance, datasets with more complex relationships between variables, higher dimensionality, or those subject to domain-specific nuances might present different challenges that the tested libraries could handle differently. Consequently, the findings and recommendations, though valuable, may not be universally applicable across all types of data and industries. Future research could address this limitation by including a wider variety of datasets that encompass a broader range of data types, structures, and quality issues, thereby providing a more comprehensive evaluation of the libraries’ performance. Another limitation of the research is the focus on standard data cleaning tasks without considering the integration of these libraries into broader data processing pipelines. The study primarily evaluates the libraries based on accuracy, efficiency, and ease of use for isolated data cleaning tasks. However, in real-world applications, data cleaning is often part of a larger workflow that includes data extraction, transformation, analysis, and reporting. The interoperability of these libraries with other tools, their ability to handle real-time data streams, and their performance when integrated into end-to-end data pipelines were not within the scope of this study.

This study holds significant value for the data science and analytics community, as data cleaning is a critical step in the data preparation process that directly impacts the quality of subsequent analysis and insights. The findings contribute to the broader discourse on best practices in data management, emphasizing the importance of choosing tools that align with the specific requirements of different datasets and analytical tasks. This can lead to improved data quality, more accurate analyses, and better-informed business decisions, making the research highly relevant and impactful in the field of data science.

References

Bose, Indranil, and Radha K Mahapatra. 2001. Business data mining—a machine learning perspective. *Information & management* 39 (3): 211–225.

- Chaudhuri, Surajit, and Umeshwar Dayal. 1997. An overview of data warehousing and olap technology. *ACM Sigmod record* 26 (1): 65–74.
- Cheng, Long, Fang Liu, and Danfeng Yao. 2017. Enterprise data breach: causes, challenges, prevention, and future directions. *Wiley Interdisciplinary Reviews: Data Mining and Knowledge Discovery* 7 (5): e1211.
- Fan, Wenfei, Shuai Ma, Nan Tang, and Wenyuan Yu. 2014. Interaction between record matching and data repairing. *Journal of Data and Information Quality (JDIQ)* 4 (4): 1–38.
- Gmach, Daniel, Jerry Rolia, Ludmila Cherkasova, and Alfons Kemper. 2007. Workload analysis and demand prediction of enterprise data center applications. In *2007 IEEE 10th international symposium on workload characterization*, 171–180. IEEE.
- Hernández, Mauricio A, and Salvatore J Stolfo. 1998. Real-world data is dirty: data cleansing and the merge/purge problem. *Data mining and knowledge discovery* 2:9–37.
- Jarmin, Ron S, and Javier Miranda. 2002. The longitudinal business database. *Available at SSRN 2128793*.
- Kandel, Sean, Andreas Paepcke, Joseph M Hellerstein, and Jeffrey Heer. 2012. Enterprise data analysis and visualization: an interview study. *IEEE transactions on visualization and computer graphics* 18 (12): 2917–2926.
- Kotsiantis, Sotiris B, Dimitris Kanellopoulos, and Panagiotis E Pintelas. 2006. Data preprocessing for supervised learning. *International journal of computer science* 1 (2): 111–117.
- Nie, Huihua, Ting Jiang, and Rudai Yang. 2012. A review and reflection on the use and abuse of chinese industrial enterprises database. *World Economy (in Chinese)*, no. 5, 142–158.
- Olson, David Louis, Yong Shi, and Yong Shi. 2007. *Introduction to business data mining*. Vol. 10. McGraw-Hill/Irwin New York.
- Pyle, Dorian. 1999. *Data preparation for data mining*. morgan kaufmann.
- Rahm, Erhard, Hong Hai Do, et al. 2000. Data cleaning: problems and current approaches. *IEEE Data Eng. Bull.* 23 (4): 3–13.
- Wickham, Hadley. 2014. Tidy data. *Journal of statistical software* 59:1–23.