

# A Study on Data Cleaning Techniques for Large Datasets

Divya Kashid

*Sonopant Dandekar College, Palghar (W), Maharashtra, India;*

**Abstract**—Data cleaning is an essential phase in the data preparation process, especially when working with large datasets. These datasets include often missing value, duplicate records, noise and anomalies that must be addressed for reliable analysis and decision making. This research examines several types of data cleaning techniques, including missing value copying, deduction, external identification and generalization, and evaluating their application on a large-scale dataset. The paper also reviews modern equipment and outlines that facilitates automatic and scalable data cleaning.

identified issues like schema-level inconsistencies and duplicate detection. Other researchers have focused on the use of machine learning techniques for data cleaning, emphasizing their effectiveness in automating repetitive tasks and handling complex data structures. Popular open-source libraries such as Pandas, Dask, and data wrangling tools like OpenRefine and Trifacta have been recognized for their capabilities in large-scale data cleaning.

## 3. DATA CLEANING TECHNIQUES

### 1. INTRODUCTION

In recent years, the explosion of big data across industries has led to an increasing reliance on high-quality data for insights and decision-making. However, raw data collected from various sources is frequently incomplete, incorrect, or inconsistent. Data cleaning, also known as data scrubbing, plays a crucial role in ensuring the integrity of such data. Effective data cleaning processes help to minimize errors and improve the overall usability of datasets for data analysis and machine learning models. Large datasets, often in the order of terabytes or petabytes, present specific challenges in terms of scalability, performance, and accuracy of cleaning methods. This paper focuses on the importance of data cleaning and investigates various techniques that can be applied efficiently to large datasets.

### 2. LITERATURE REVIEW

The importance of data cleaning has been well-documented in various academic and industry research. Dasu and Johnson (2003) outlined key data quality dimensions and proposed exploratory techniques for assessing data quality. Rahm and Do (2000) discuss data cleaning architectures and

#### 3.1 Handling Missing Data

Missing data is one of the most common data quality problems. Techniques to handle missing data include deletion (listwise or pairwise), mean/median/mode imputation, forward or backward filling, and predictive modeling using algorithms such as K-Nearest Neighbors (KNN) or regression. The choice of method depends on the type of data, the percentage of missing values, and the analytical requirements.

#### 3.2 Deduplication

Duplicate records arise when data is collected from multiple sources or entered manually. Deduplication involves identifying and removing or merging duplicate records. Techniques include exact matching, fuzzy string-matching using algorithms like Levenshtein distance, and clustering-based deduplication approaches. Rule-based systems and ML-based entity resolution methods are increasingly used for large-scale deduplication.

#### 3.3 Outlier Detection

Outliers are values that deviate significantly from other observations and may indicate errors or rare events. Detection methods include statistical

techniques such as the Z-score and IQR, as well as machine learning approaches like isolation forests and DBSCAN clustering. Proper handling of outliers is essential for preserving the statistical integrity of the dataset.

### 3.4 Normalization and Transformation

Data normalization is used to scale numeric values to a common range, making the data more suitable for algorithms that are sensitive to magnitude. Common methods include min-max scaling, z-score normalization, and log transformation for skewed data. Standardizing the data can improve performance in machine learning tasks and visualization.

### 3.5 Automation Tools and Frameworks

Various tools are available for automating data cleaning. OpenRefine offers a graphical interface for cleaning messy data, while Trifacta provides intelligent suggestions and transformations. Python libraries like Pandas and NumPy support extensive data manipulation and cleaning functions. Dask and PySpark are preferred for distributed data cleaning on large datasets.

## 4. METHODOLOGY

To evaluate data cleaning techniques, we selected three datasets: - The Adult Income dataset from the UCI Machine Learning Repository - A customer sales dataset from Kaggle - A synthetic dataset with intentional noise and inconsistencies. Each dataset was cleaned using a combination of tools (Pandas, OpenRefine) and methods (imputation, deduplication, outlier removal). Metrics such as processing time, memory usage, and improvement in model accuracy were measured to compare the techniques.

## 5. RESULTS AND DISCUSSION

The experiments showed that Pandas is highly efficient for scripting data cleaning pipelines, especially when combined with Jupyter Notebooks. OpenRefine was easy to use but had performance limitations with large datasets (>1 million rows). Imputation improved dataset completeness, with

predictive methods performing better than statistical imputation in complex datasets. Deduplication improved classification model accuracy by reducing noise. Normalization was critical in improving convergence speed of learning algorithms. Overall, data quality improvements led to a 10–25% increase in model accuracy.

## 6. CONCLUSION

Data cleaning is a vital step in any data-driven project. As datasets grow in size and complexity, manual cleaning becomes impractical. This study highlights the effectiveness of various data cleaning techniques and tools in improving data quality. Future work should focus on developing intelligent, adaptive data cleaning systems that can learn and refine rules based on data context. Furthermore, integration of these tools into data pipelines and their application in real-time environments will enhance the scalability and reliability of data analytics systems.

## REFERENCES

- [1] Dasu, T., & Johnson, T. (2003). Exploratory Data Mining and Data Cleaning.
- [2] Rahm, E., & Do, H. H. (2000). Data cleaning: Problems and current approaches. IEEE Data Engineering Bulletin.
- [3] McKinney, W. (2010). Data Structures for Statistical Computing in Python.
- [4] UCI Machine Learning Repository. <https://archive.ics.uci.edu/ml>
- [5] Kaggle Datasets. <https://www.kaggle.com> - OpenRefine. <https://openrefine.org> - Trifacta Wrangler. <https://www.trifacta.com/products/wrangler>