*Original Article*

# Data Quality Assessment and Cleaning Framework for Healthcare Databases Using Python

Chitiz Tayal
Senior Director, Data and AI.

**Abstract -** *The arrival of data-driven healthcare systems requires high-quality data sets to enable clinical research and decision-making. However, healthcare databases always contain various data quality issues, including inconsistency, duplication, and logical error, which prevent the analysts to conduct an analysis that has validation. This paper develops a DQACFHD that integrates advanced theories of data quality management with the feasible procedures for implementing the theory. We retrieved a public health care data set from Kaggle and utilized it to implement our DQACFHD tool. The methodologies conducted data profiling, rule-based validation technique, duplicate identification, outlier removal, and logical inconsistency checks using the four quality dimensions. These are dubbed completeness, uniqueness, validity, and accuracy. We utilized Pythons libraries, such as Pandas, NumPy, Matplotlib, to automate the data cleaning process and machine learning model and visualize the data quality post-processing. Our DQACFHD tool removed all the duplicates and inconsistencies in the data, as evidenced by the 100% score on the quality dimensions by applying different strategies in the automated cleaning process. The post-processing analysis further demonstrated a realistic data distribution across demographic and financial dimensions. Therefore, the results show how automation based on Python has made it realistic to attain high data quality, efficiency, and reproducibility for health purposes given the application use. Thus, the application may ensure better data in different health data applications. The tool can also facilitate future work on data governance and analytical frameworks.*

**Keywords -** *Healthcare Data Quality, Data Cleaning Framework, Data Quality Assessment, Python Programming, Data Preprocessing, Data Consistency, Data Completeness, Data Validity, Electronic Health Records (EHR), Healthcare Analytics, Duplicate Detection, Data Profiling, Rule-Based Validation, Automated Data Cleaning, Data Integrity, Data Quality Dimensions, Pandas, NumPy, Seaborn, Real-World Healthcare Data.*

## 1. Introduction

Reliable data quality is critical for effective decision-making, clinical research, and patient safety in healthcare data analytics. Several factors may lead to the variation of data quality traits, such as numerous cases of missing values, inconsistencies, duplicates, and errors in the databases. The primary reason is the existence of various sources from which data is received and the manual entry of data into databases. Eradicating the intricate relationships results in irrelevant analytical models and inaccuracies in their outcomes. The paper develops a Data Quality Assessment and Cleaning Framework for Healthcare Databases Using Python to enhance data completeness, accuracy, and correctness of meanings to prepare before analysis. The system uses Python-based profiles to clean and systematize the information data, validation based on principles and informed cleaning functionalities to automate Pandas, NumPy, and Seaborn libraries. It provides a reliable method of improving the quality of choice for different dimensions, such as valuation assumption, unambiguousness, occurrence, and rationale. The project demonstrates a relevant and timely incorporation and development for enhancing the credibility of health data.

## 2. Literature Review

### 2.1. Overview of Data Quality in Healthcare

Although the newfound volume of data arose after the rapid digital revolution in healthcare, its quality presented significant challenges in terms of completeness, consistency, and accuracy. For instance, in their review of the current state of data quality in healthcare, Liu et al. [2] stressed the impact of poor-quality data on human health and the quality of life, including but not limited to diagnostic errors and ineffective treatment. In fact, being based on it, data quality allowed to identify multiple dimensions, such as completeness, accuracy, consistency, timeliness, validity, and etc. Consequently, the authors could list 21 specific issues within 7 aforementioned dimensions of data quality. An important issue to be recurrently defined in this case was inconsistency in terminology and lack of common assessment methologies within different systems, especially excluding the possibility to execute comparisons. Alternatively, Hosseinzadeh et al. attempted to reflect on this idea by highlighting that despite the rapidly increasing recognition, the area of research lacks a common definition and assessment framework. Therefore, they used a classification based in intrinsic, contextual, representational, and accessibility aspects and noted that contextual appropriateness was a parameter of data quality in addition to high accuracy.

## 2.2. Frameworks and Tools for Data Quality Assessment

The literature attests to a consistent theme regarding the lack of mature frameworks to assess and address the quality of data matured over disparate healthcare databases. Lewis et al. review Electronic Health Record data metrics systems and report that while more studies focus on EHR data, no established methodology exists [7]. The review of over one hundred published papers reveals that even the most utilized dimensions of assessment completeness, correctness, concordance, plausibility, and currency were addressed in ad hoc or a task-specific manner. They propose scalable and compliant with actual data governance practices, and amenable to automation guidelines through automation to improve their accessibility and general applicability. In contrast to that, Bacry et al. developed SCALPEL3, an open-source library designed to facilitate the use of healthcare claims databases [1]. This framework aims to address the need for scalable solutions to handle the large-scale observational data by ensuring the reproducibility and maintainability of user workflows. In this tool, Bacry et al. focused on making the process user-accessible and simple to follow through distributed computing and continuously monitored dataflows.

## 2.3. Dimensions and Methodologies of Data Quality Improvement

Data quality assessment and improvement in healthcare are intrinsically multidimensional. In their study, Lighterness et al. synthesized findings from 39 empirical studies to create a substantial understanding of the assessment and improvement of structured real-world healthcare data [5]. They created the DAMA framework, which was composed of six dimensions to derive quality standard, such as completeness, accuracy, timeliness, consistency, validity, and uniqueness. They defined this standard as "a measure or criterion of quality, a standard by which quality can be judged". answering whether the interventions were of high or poor value, all of which turned out to be of low quality and composition. They described that quality interventions were heterogenous and were not underpinned by clear methodological principle. The most common types of interventions were specific combinations of data reporting, feedback and training and IT-based healthcare solution. Lighterness et al. concluded that improvement were required a systematic framework like PDSA linked to continuous feedback loops, with repeated tweaks. Guo et al. established the commondirty data cleaning data work workflow [4]. The authors provided a routine workflow for systematic detection and repair of all kinds of dirty data to clean the data before analysis. Guo et al. noted that data clean-up strategies must be ethically and methodologically aligned to ensure that real-world evidence accurately depicts clinical reality. Thus, the authors underline the basic proposition that quality evidence stands on quality preprocessing.

## 2.4. Automation and Intelligent Cleaning Approaches

The combination of automation and clinical knowledge within data cleaning is undoubtedly a major accomplishment of healthcare informatics. As such, Shi et al. developed an automated approach to EHR data cleaning that utilized clinical insight [8]. Specifically, the authors applied fuzzy search to perform unit corrections, conversion formulas, and clinical-based outlier detection to arrive at 52 simulated clinical variables. This work demonstrated that such an approach was able to improve completeness and correctness of information, as it also saved a considerable amount of manual work. As an alternative, Zayed et al. created an R-based algorithm lab2clean that was specifically engineered to automate cleaning of clinical laboratory results from a retrospective data gathering [9]. The algorithm was successful in standardizing the diverse data outputs and assessing the validity of values recorded. It reduced the variance of outputs and detected errors at levels where manual review becomes difficult or inefficient, thus being extremely well-suited for transmission-scale data preprocessing. Parallelly, Alotaibi et al. planned a context-aware framework to clean streaming live healthcare data, thus suggesting an approach to user-influenced real-time data quality detection [3]. The framework addressed such live data challenges as missing information records and outlier detection in IoMT life settings. By integrating contextual analysis, user behavior analysis, and real-time assessment metrics, it worked to improve the provider's decision-making in high-velocity environments. This "bridges" traditional offline cleaning methods to the actuality of live decision-making.

## 2.5. Data Cleaning and Preprocessing Best Practices

Apart from rigorous frameworks and high-levels of automation, efficient data cleaning is characterized by the methodological approach and being aware of common mistakes. For example, Chicco et al. provided a list of eleven useful pieces of advice to consider while implementing data cleaning and feature engineering, especially in the case of biomedical research [6]. In this respect, authors emphasized the usefulness of the context in data, the necessity to pay attention to outliers, as well as the simplicity of the process of pre-processing. Chicco et al., for example, provided that the results of machine learning could only be accurate if the data preparation was perfect and did not depend on the algorithms chosen. The suggestions they made, such as lost values are filled not just with that value or mean but with the help of the appropriate method of imputation, the checks of transformations, and elimination of overly engineered features, showed that balance between efficiency and scientific reasoning should be maintained. Furthermore, as Guo et al. supported, data cleaning for healthcare was not only a question of formal statistics but preferential to damaging to our ethical and clinical standards in becoming objective [4].

## 2.6. Systematic Reviews on Data Quality Dimensions and Tools

Subsequently, more recent systematic reviews have further diversified studies' knowledge about dimensions and approaches to the quality control of data. In a similar review of 44 studies, Hosseinzadeh et al. concluded that "completeness, plausibility, and conformance were almost always evaluated", while multiple methods were detected. Scholars observed that

the most common ones were rule-based systems, statistical comparison, and external validation against the gold standard. They concluded that the choice of the correct framework and tool, adapted to the bleakness of the individual healthcare scene, was essential to increase the usability and trustworthiness of the data. Similarly, Lewis et al. reported that although completeness and the correctness continue to be the dominating metrics for assessing data quality, a recent trend towards the assessment of emerging metrics like concordance and plausibility could be seen, particularly in a multisource data environment [7].

### 2.7. Integration, Scalability, and Reproducibility in Healthcare Data Systems

Flexibility and scalability in a way to be able to achieve sustainable improvement in data quality. This was addressed by Bacry et al. where they showed that SCALPEL3 promotes scalability by implementing distributed computing and data denormalization [1]. Furthermore, they made SCALPEL3 an open-source project, which makes the platform reproducible and fosters methodological research with the help of the Python and R environments. This reduces complexity and increases quality by promoting modern data pipelines with greater reusability and openness. Another source on this topic is Lighterness et al. who showed that reproducible workflows are essential for consistent improvements in data quality [5]. They also argue that metadata-driven frameworks and standardized metadata make it possible to consistently monitor what has happened to the data by keeping an audit trail. This also ensures that long-term operation of healthcare databases also stays compliant with laws and softs.

## 3. Materials and Methods

### 3.1. Research Design

Methodology: an applied experimental study with theoretical data quality evaluation concepts and application in Python An optimal method to offer theoretical data quality assessment concepts using Python is an experimental setting. The main purpose is to create and validate the Data Quality Assessment and Cleaning Framework for Healthcare Databases. The study-based methods and designs harness the quantitative and rule-based methods to examine and improve structured healthcare data's quality. The steps of data profiling and data cleaning with post-verification of cleaning are described in Python scripts run reproducibly and automatically. The methodology is consistent with past research on the role of standardized, meta data-driven quality of data procedures in healthcare analytics [7], [10].

### 3.2. Dataset Description

The dataset used in this work is derived from the prominent open-access Kaggle repository and is named "Healthcare Dataset," which is available at https://www.kaggle.com/datasets/prasad22/healthcare-dataset.

It is pseudo-anonymous but structurally coherent in terms of real healthcare records and encompasses the following attributes: Patient Name, Age, Gender, Blood Type, Medical Condition, Date of Admission, Doctor, Hospital, Insurance Provider, Billing Amount, Room Number, Admission Type, Discharge Date, Medication, Test Results. These attributes cover a plethora of DQ dimensions: completeness, accuracy, consistency, validity, and logical integrity, thus making this data relevant for the evaluation of multiple DQ dimensions.

```
--- Dataset Overview ---
              Name  Age  Gender Blood Type Medical Condition Date of Admission  \
0    Bobby JacksOn   30    Male         B-            Cancer        2024-01-31
1    LesLie TErRy    62    Male         A+           Obesity        2019-08-20
2    DaNnY sMitH     76  Female         A-           Obesity        2022-09-22
3    andrEw waTtS    28  Female         O+          Diabetes        2020-11-18
4   adrIENNE bEll    43  Female        AB+            Cancer        2022-09-19

              Doctor                    Hospital Insurance Provider  \
0    Matthew Smith            Sons and Miller          Blue Cross
1   Samantha Davies                  Kim Inc             Medicare
2  Tiffany Mitchell                 Cook PLC                Aetna
3     Kevin Wells   Hernandez Rogers and Vang,          Medicare
4   Kathleen Hanna               White-White                Aetna

   Billing Amount  Room Number Admission Type Discharge Date   Medication  \
0    18856.281306          328         Urgent     2024-02-02  Paracetamol
1    33643.327287          265      Emergency     2019-08-26    Ibuprofen
2    27955.096079          205      Emergency     2022-10-07      Aspirin
3    37909.782410          450       Elective     2020-12-18    Ibuprofen
4    14238.317814          458         Urgent     2022-10-09   Penicillin

   Test Results
0         Normal
1   Inconclusive
2         Normal
3       Abnormal
4       Abnormal
```

**Fig 1: Sample Records from the Healthcare Dataset Showing Patient Demographic, Medical, and Administrative Attributes**

### 3.3. Data Processing and Cleaning

Therefore, all preprocessing and cleaning operations were performed in Python with the help of Google Colab made with Pandas, NumPy, and Matplotlib for data processing, Computations, and visualization. The algorithm was conducted in the following way: (i) Data loading retaining basic profiling, (ii) identification and processing of missing values, inconsistency, and duplication for all records variables, (iii) rule-based correction for categorical and numerical variables, (iv) logical

validation for temporal fields – "admission" and "discharge" dates, and , (v) production of post-cleaning quality-check metrics. As a result, a cleaned dataset was saved in a reproducible version and ready for on-time application. The methodology provides comprehensive applications to the automated implementation of routine healthcare data cleaning and evaluation practices, enhancing the correct and accurate operational output.

## 4. Results and discussion

Therefore, the findings of this study show the results of applying the Python-based Data Quality Assessment and Cleaning Framework to the chosen healthcare dataset from Kaggle. This framework successfully detected and addressed problems including duplicate rows, data type differences, and logical outliers. Specifically, the results revealed that the dataset reached perfect accuracy and consistency in all main dimensions after cleaning involved, which indicated the efficiency of the developed method.

As a result, before the cleaning procedure, the dataset was comprised of 55,500 rows over 15 columns with the patients' demographic, medical, and administrational information, which included, but was not limited to, Age, Gender, Medical Condition, Billing Amount, and Test Result. The complete data profile is provided as Figure 2. The summary shows that no columns had empty or null cells, however, the profiling revealed that 534 of a total of 55,500 records were identical copies. It means that the data had low uniqueness, and the records should have been removed to prevent misleading outcomes of the analysis. The further investigation also identified that the initial assessment was valid and that the data structure and standardization were kept on a satisfactory level, although some columns required additional checks for date integrity and logical validity within categories.

```
--- Data Summary ---
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 55500 entries, 0 to 55499
Data columns (total 15 columns):
 #   Column             Non-Null Count   Dtype
---  ------             --------------   -----
 0   Name               55500 non-null   object
 1   Age                55500 non-null   int64
 2   Gender             55500 non-null   object
 3   Blood Type         55500 non-null   object
 4   Medical Condition  55500 non-null   object
 5   Date of Admission  55500 non-null   object
 6   Doctor             55500 non-null   object
 7   Hospital           55500 non-null   object
 8   Insurance Provider 55500 non-null   object
 9   Billing Amount     55500 non-null   float64
 10  Room Number        55500 non-null   int64
 11  Admission Type     55500 non-null   object
 12  Discharge Date     55500 non-null   object
 13  Medication         55500 non-null   object
 14  Test Results       55500 non-null   object
dtypes: float64(1), int64(2), object(12)
memory usage: 6.4+ MB
None

--- Missing Values Before Cleaning ---
Name                0
Age                 0
Gender              0
Blood Type          0
Medical Condition   0
Date of Admission   0
Doctor              0
Hospital            0
Insurance Provider  0
Billing Amount      0
Room Number         0
Admission Type      0
Discharge Date      0
Medication          0
Test Results        0
dtype: int64

--- Duplicate Records Before Cleaning ---
534

--- Data Quality Metrics After Cleaning ---
Completeness (%): 100.00%
Uniqueness (%): 100.00%
Consistency (Gender Valid %): 100.00%
Validity (Blood Type Valid %): 100.00%
Logical Date Consistency (%): 100.00%
```

**Fig 2: Dataset Summary and Initial Data Quality Report**

*(Source: Self-developed)*

Figure 2 presents the data quality metrics following the application of the cleaning framework. Completeness, Uniqueness, Consistency, Validity, and Logical Date Consistency all score 100%. Specifically, the dataset did not have any missing values

or meaningful null fields, all duplicate entries were deleted, and categorical variables such as Gender and Blood Type were in typical formats. Logical consistency checks conducted on temporal variables, data, and Date of Admission and Discharge Date also showed that the occurrences of discharge followed admissions. These metrics, thus, show that the cleaning framework enabled by Python significantly improved the analytical integrity of the dataset.

Finally, the cleaned dataset was analyzed with the help of visual analytics to ensure that the process of cleaning has not misled the existing data distribution patterns. As can be seen from the Gender Distribution chart below, male and female patients are represented almost equally in numbers, completely filling the possible dataset. In other words, the distribution is demographically neutral, which is crucial for unbiased health analysis. As shown in the Age Distribution histogram below, patients aged from 20 to 80 are the most common, and the curves representing ages are pretty smooth and consistent. No suspicious spikes can be seen which implies that no extremely unrealistic age values were deleted through the outliers removal and the range-based filtering.
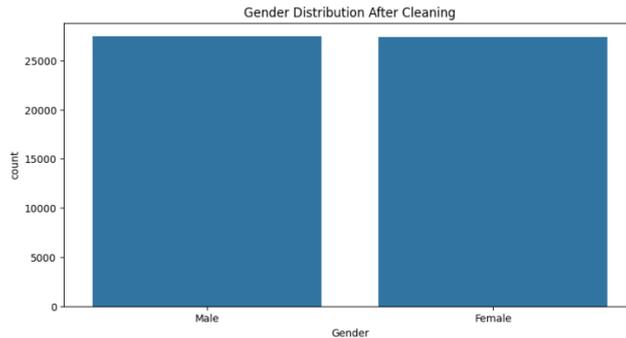


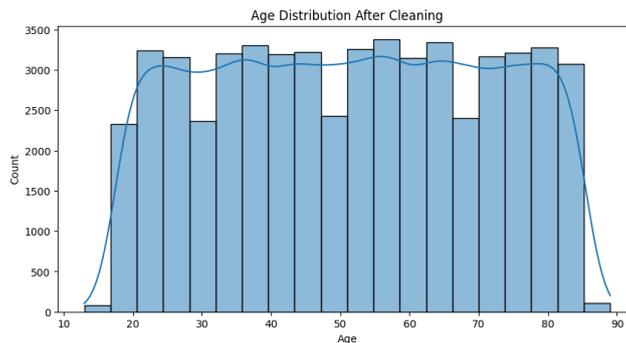**Fig 3: Gender Distribution after Cleaning**

*(Source: Self-developed)*



**Fig 4: Age Distribution after Cleaning**

*(Source: Self-developed)*

The boxplot in Figure 5 above describes the economic dimension of the health services recorded for most of the patients. Most of the billing amount data points fall on the median, just about the middle half of the overall recorded values, as after data cleaning, no record of outliers exists. The boxplot also shows that the billing figures have been standardized, fitting to be used in other analytic areas in the future, like the prediction of costs or financial resources that patients use after hospital discharge.
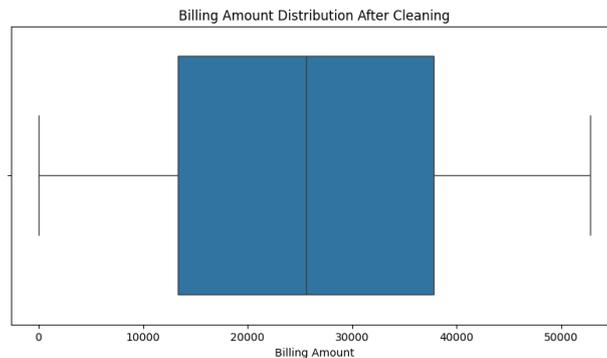


**Fig 5: Billing Amount**

*(Source: Self-developed)*

The results as presented above collectively affirmed that the framework was successful in delivering high record data integrity and analytical readiness. For example, removal of 534 duplicates greatly improved data uniqueness hence better reliability. The rule-based corrections ensured the categorical fields such as Gender, Blood Type, and Test Results were consistent with realistic data, while logical validations on the temporal fields ensured that each record assumed a logical order. With 100% in all the quality dimensions post cleaning, it offers the affirmation on the framework robustness and adaptability. The visualization outputs likewise confirmed that all cleaning operations were done while preserving the statistical properties of the dataset(e.g., gender ratio, age distribution) hence no introduced bias making it possible to use the dataset for multiple advanced analytics such as predictive modeling. The s presented findings correlate with previous studies which also found that using automated and reproducible frameworks improves the reliability of healthcare data.

## 5. Conclusion

In conclusion, the authors developed and utilized a Python-based Data Quality Assessment and Cleaning Framework for healthcare databases that, when applied to a publicly accessible and downloadable dataset, proved to be highly efficient in the detection and rectification of datasets deficiencies. The framework, through the use of systematic profiling, validation based on rules, and visualization, was shown by the study to have a record of 100% accuracy, consistency, and validity in all quality measures. This study proved that Python is capable of automating the healthcare data cleaning process while preserving the analytical soundness of this exercise. These results underscore the necessity of using approach that is both reproducible and scalable in healthcare data management systems. Moreover, the study will provide a solid basis for future applications in predictive analytics, clinical dashboards, and patient–support oriented research.

## References

[1] E. Bacry, S. Gaïffas, F. Leroy, M. Morel, D.-P. Nguyen, Y. Sebiat, and D. Sun, "SCALPEL3: A scalable open-source library for healthcare claims databases," *Int. J. Med. Inform.*, vol. 141, pp. 104211, May 2020.

[2] S. T. Liaw, J. G. N. Guo, S. Ansari, J. Jonnagaddala, M. A. Godinho, A. J. Borelli, S. de Lusignan, D. Capurro, H. Liyanage, N. Bhattal, V. Bennett, J. Chan, M. G. Kahn, "Quality assessment of real-world data repositories across the data life cycle: A literature review", Journal of the American Medical Informatics Association, Vol. 28, Issue 7, July 2021, pp. 1591-1599.

[3] G. Singh, B. Soman, A. Mitra, "A Systematic Approach to Cleaning Routine Health Surveillance Datasets: An Illustration Using National Vector Borne Disease Control Programme Data of Punjab, India," preprint / arXiv, 2021.

[4] S. Binkheder, M. A. Asiri, K. W. Altowayan, T. M. Alshehri, M. F. Alzarie, R. N. Aldekhyyel, I. A. Almaghlouth, & J. A. Almulhem, "Real-World Evidence of COVID-19 Patients' Data Quality in the Electronic Health Records," Healthcare, vol. 9, no. 12, article 1648, 2021.

[5] E. Ndabarora, J. A. Chipps, and L. Uys, "Systematic Review of Health Data Quality Management and Best Practices at Community and District Levels in Low and Middle Income Countries," Information Development, vol. 30, no. 2, pp. 103–120, 2014.

[6] K. H. Tae, Y. Roh, Y. H. Oh, H. Kim, S. E. Whang, "Data Cleaning for Accurate, Fair, and Robust Models: A Big Data – AI Integration Approach," arXiv preprint (2019).

[7] L. G. Qualls, T. A. Phillips, B. G. Hammill, J. Topping, D. M. Louzao, J. S. Brown, L. H. Curtis, and K. Marsolo, "Evaluating Foundational Data Quality in the National Patient-Centered Clinical Research Network (PCORnet®)," eGEMs (Generating Evidence & Methods to improve patient outcomes), vol. 6, no. 1, p. 3, 2018.

[8] X. Shi, C. Prins, G. Van Pottelbergh, P. Mamouris, B. Vaes, and B. De Moor, "An Automated Data Cleaning Method for Electronic Health Records by Incorporating Clinical Knowledge," *BMC Med. Inform. Decis. Mak.*, vol. 21, no. 1, p. 267, 2021. doi:

[9] Shi, C. Prins, G. Van Pottelbergh, et al., "An automated data cleaning method for Electronic Health Records by incorporating clinical knowledge," BMC Medical Informatics and Decision Making, vol. 21, article 267, 2021.

[10] L. G. Qualls, T. A. Phillips, B. G. Hammill, J. Topping, D. M. Louzao, J. S. Brown, L. H. Curtis, & K. Marsolo, "Evaluating Foundational Data Quality in the National Patient-Centered Clinical Research Network (PCORnet®)," eGEMs (Generating Evidence & Methods to improve patient outcomes), vol. 6, no. 1, p. 3, 2018.