

## Automating Clinical Data Cleaning and Analysis Using SAS Macros

Arvind Uttiramerur

Programmer Analyst at Thermofisher Scientific, USA

### \*Corresponding author

Arvind Uttiramerur, Programmer Analyst at Thermofisher Scientific, USA.

Received: May 12, 2023; Accepted: May 18, 2023, Published: May 25, 2023

### Introduction

Clinical trials data complexity involves integrating various data elements like demographic, laboratory, clinical, medication, and medical history, which are valuable but often lack completeness and cleanliness. Maintaining data integrity and cleanliness in clinical trials data requires strategic planning and execution of data edit checks, cleaning, and documentation using SAS procedures like PROC UNIVARIATE and PROC FREQ.

SAS provides real-time documentation of data cleaning procedures and results, ensuring transparency and traceability in the process. Data screening and cleaning processes are crucial in clinical trials to identify outliers, missing values, duplicates, and incorrect data, followed by executing SAS procedures to ensure data quality. Utilizing SAS macros, arrays, and data steps can efficiently flag data for cleaning, identify errors, and create uniform clinical analytic datasets for further analysis.

### Efficiency Improvements in Data Cleaning Using SAS Macros

Creation of a Data Screening and Cleaning Plan is essential for success in working with clinical trials data. Outline anticipated data ranges, necessary variables, primary outcome variables, and qualitative variables for screening. Utilize SAS procedures like PROC UNIVARIATE and PROC FREQ to identify outliers, missing values, duplicates, and incorrect data. Use SAS macros, data steps, and arrays to efficiently flag data for cleaning and create uniform clinical analytic datasets. Maintain documentation throughout the data cleaning process to provide a record of the steps performed and results obtained.

Record information on duplicate observations, outlying values, and expected ranges in the cleaning process spreadsheet for reference. Automation tools like SAS macros streamline the identification of outliers, duplicates, and errors, enhancing the efficiency of data cleaning processes. Automated procedures ensure consistent application of data cleaning checks across datasets, reducing the risk of human error and ensuring data quality.

Automation facilitates real-time documentation of data cleaning steps and results, providing transparency and traceability in the process. Automation tools help in quickly flagging data discrepancies, allowing for prompt identification and resolution of errors in the dataset. Utilizing automation in data cleaning processes saves time by efficiently executing repetitive tasks,

allowing data managers to focus on more complex data quality issues.

SAS Macros are pieces of code that are created to automate repetitive tasks in SAS programming, enhancing efficiency and reducing manual effort. Macros in SAS are designed to perform specific actions or tasks, allowing for the automation of processes like data cleaning, quality control, and report generation. SAS Macros save time and effort by automating routine tasks, ensuring consistency in data processing and analysis. Macros are commonly used in tasks that involve repetitive coding patterns, such as data validation, error checking, and data cleaning processes.

SAS Macros provide flexibility in programming by allowing users to define and customize their own functions and commands, tailored to specific project requirements. Begin by clearly defining the task that the SAS macro needs to automate, ensuring a thorough understanding of the specific requirements. Plan the structure of the macro by breaking down the task into smaller, manageable steps, ensuring a logical flow of operations within the macro code.

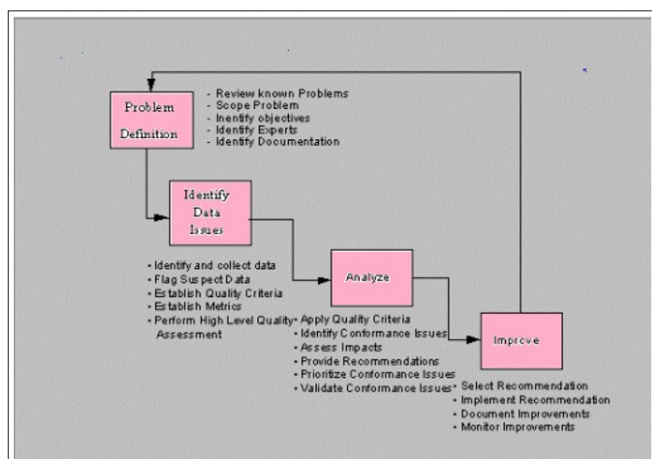
Parameterize inputs within the macro to make it adaptable to different datasets or scenarios, enhancing the flexibility and reusability of the macro. Implement robust error handling mechanisms within the macro to anticipate and address potential issues during execution, ensuring the reliability of the automation process. Thoroughly test the macro with various datasets to validate its functionality and accuracy, ensuring that it performs as intended across different scenarios.

Clearly define the data cleaning requirements and expected ranges of variables before initiating the cleaning process. Use SAS Macros to efficiently check continuous variables for outliers, providing statistics like minimum, maximum, mean, and median for each variable.

Employ a SAS Macro with PROC FREQ to identify outliers in categorical variables by comparing variable values with expected ranges, recording outlying values in the Cleaning Process spreadsheet. Utilize SAS Enterprise Guide to organize cleaning programs, assign project libraries, and document datasets undergoing the cleaning process for efficient data management. Use a SAS Macro in conjunction with PROC SORT to efficiently identify duplicate observations based on specified variables,

recording information on duplicate observations in the Cleaning Process spreadsheet.

Develop robust error handling mechanisms within SAS Macros to anticipate and address potential issues during data cleaning, ensuring the reliability of the cleaning process.



### Streamlining Data Analysis with SAS Macros

Data analysis is a crucial process that involves inspecting, cleansing, transforming, and modeling data to uncover meaningful insights, trends, and patterns that can drive informed decision-making. Before embarking on data analysis, it is essential to clearly define the objectives and goals of the analysis to ensure that the insights derived align with the intended outcomes. The initial steps of data analysis involve collecting relevant data from various sources and preparing it for analysis by cleaning, organizing, and structuring the data in a format suitable for analysis.

Exploratory Data Analysis (EDA) is a critical phase where analysts explore the data visually and statistically to understand its characteristics, identify patterns, outliers, and relationships between variables, providing valuable insights into the dataset. Statistical analysis techniques are then applied to the prepared data to uncover correlations, trends, and associations, followed by modeling processes such as regression, clustering, or machine learning to derive predictive insights from the data.

The final step involves interpreting the results of the analysis, drawing conclusions, and effectively communicating the insights to stakeholders through visualizations, reports, and presentations to support data-driven decision-making. SAS Macros can be leveraged to automate repetitive data analysis tasks, such as data cleaning, outlier detection, and report generation, streamlining the analysis process.

Automation using SAS can facilitate dynamic data monitoring by providing metrics and tools to track enrollment progression, site metrics reporting, and safety-based medical monitoring in real-time, enhancing the quality of research study data. The SAS Output Delivery System (ODS) can be utilized to generate custom reports in various formats like HTML, PDF, and Excel, automating the reporting process and ensuring easy accessibility to critical information.

Web-based Electronic Data Capture (EDC) systems can automate data collection processes, assist in clinical trial management, and provide backend tables for efficient trial monitoring and risk-based monitoring during clinical trials SAS Macros can efficiently

check continuous variables for outliers, ensuring data quality by identifying and addressing data entry errors that may impact the analysis phase, improving the overall integrity of the data.

Efficient Outlier Detection SAS Macros can efficiently check continuous variables for outliers using PROC UNIVARIATE, providing statistics like minimum, maximum, mean, median, and IQR for each variable, aiding in data quality assessment Automated Reporting with ODS.

The SAS Output Delivery System (ODS) enhances reporting capabilities by generating reports in multiple formats like HTML, RTF, PDF, and Excel files, automating the reporting process for trial monitoring and management.

### Quality Control Automation

SAS Macros are ideal for automating quality control tasks, detecting, reporting, and repairing data quality errors, saving time and promoting the importance of quality control in data analysis Dynamic Data Monitoring. SAS Macros can assist in dynamic clinical trial data monitoring by providing tools for real-time tracking of study activities, ensuring efficient data management and monitoring processes.

### Enhanced Data Cleaning

SAS Macros can be used to identify outliers in categorical variables using PROC FREQ, facilitating the identification of incorrect values and ensuring data integrity during the cleaning process.

### Automated Outlier Detection

SAS Macros can efficiently check continuous variables for outliers using PROC UNIVARIATE, providing essential statistics like minimum, maximum, mean, median, and IQR for each variable, aiding in data quality assessment.

### Enhanced Reporting Capabilities Utilizing

SAS Output Delivery System (ODS), SAS Macros can generate custom reports in various formats such as HTML, RTF, PDF, and Excel files, improving the efficiency of reporting processes for trial monitoring and management.

### Dynamic Data Monitoring in Clinical Trials

SAS Macros can automate tracking enrollment progression, real-time site metrics reporting, and safety-based medical monitoring, enhancing data quality and human subject protection in clinical trials.

### Quality Control Automation

SAS Macros are ideal for automating quality control tasks, detecting, reporting, and repairing data quality errors, saving time and promoting the importance of quality control in data analysis.

### Efficient Data Cleaning Process

SAS Macros can assist in organizing data cleaning processes efficiently, providing a structured platform to manage programs related to data cleaning, ensuring a systematic approach to data management.

### Enhanced Data Quality Automation

Using SAS Macros can improve data quality by facilitating dynamic data monitoring, tracking enrollment progression, and providing real-time site metrics reporting, ensuring clean and high-quality data for analysis.

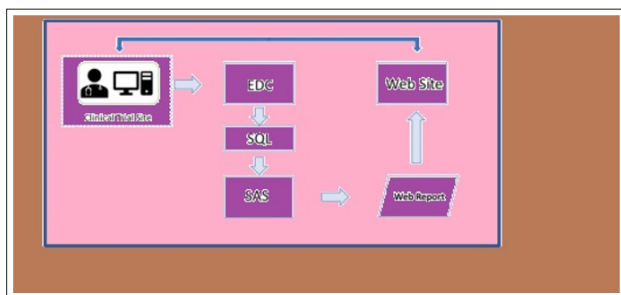
### Efficiency in Reporting Automation

Through SAS Macros and ODS can streamline reporting processes, generating reports in multiple formats like HTML, RTF, PDF, and Excel files, enhancing the efficiency of reporting for trial monitoring and management.

Time-Saving Data Cleaning Automation with SAS Macros allows for efficient data cleaning processes, identifying outliers, missing values, and incorrect data through simple SAS PROCs, DATA STEPS, MACROS, and ARRAYS, saving time and ensuring data integrity.

### Improved Trial Management Automation

In data analysis using SAS Macros can lead to improved trial management processes, ensuring timely capture of clinical data, monitoring forms, and data quality in multi-site trials, enhancing the overall study outcomes.



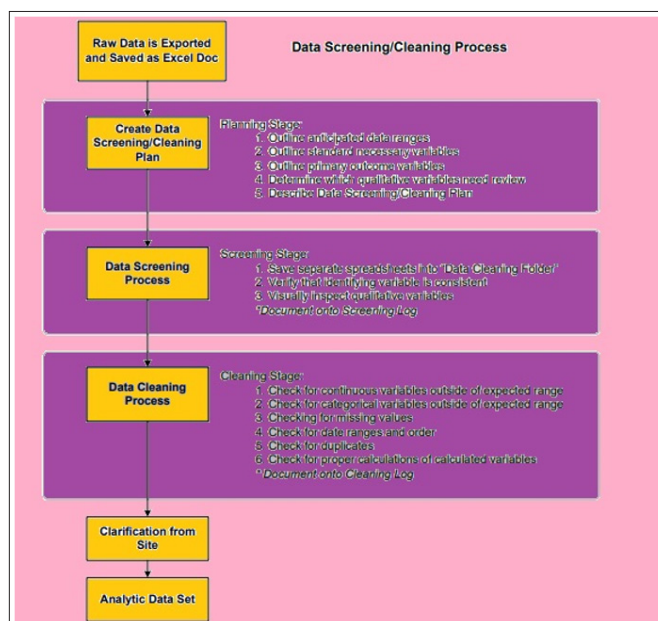
### Enhancing Data Integrity Through Automation

- **Importance of Data:** Quality data integrity refers to the accuracy, consistency, and reliability of data, ensuring that information is trustworthy and suitable for analysis.
- **Factors Affecting Data:** Integrity data integrity can be influenced by various factors such as accuracy, completeness, consistency, timeliness, uniqueness, and validity of the data.
- **Data Quality Characteristics:** Data is evaluated based on six quality characteristics: accuracy, completeness, consistency, timeliness, uniqueness, and validity, ensuring that data meets acceptance criteria and produces desired results.
- **Role of Data Cleaning Maintaining:** Data often involves data cleaning processes to identify and correct errors, inconsistencies, and missing values, ensuring data accuracy and reliability for analysis.
- **Automation for Data Integrity:** Automation tools like SAS Macros can play a crucial role in enhancing data integrity by automating quality control tasks, detecting errors, and improving data quality through efficient data cleaning processes.
- **Complex Data Sources Integrating:** Data from diverse sources with varying formats and structures can pose challenges in maintaining data integrity, requiring careful data cleaning and transformation processes.
- **Human Error:** Manual data entry and processing can lead to human errors, affecting data accuracy and integrity, highlighting the need for automated data validation and cleaning procedures.
- **Data Security Concerns:** Ensuring data security and privacy while maintaining data integrity can be challenging, especially in handling sensitive clinical trial data, requiring robust security measures and access controls.
- **Data Duplication:** Managing and identifying duplicate data entries can be a challenge in maintaining data integrity, as duplicate records can lead to inconsistencies and errors in analysis, necessitating deduplication processes.
- **Lack of Standardization:** Inconsistent data formats,

naming conventions, and quality standards across different data sources can hinder data integrity efforts, emphasizing the importance of standardization and data quality control measures.

- **Complex Data Sources:** Integrating data from diverse sources with varying formats and structures can pose challenges in maintaining data integrity, requiring careful data cleaning and transformation processes.

SAS Macros can be used to efficiently check continuous variables for outliers by running PROC UNIVARIATE and providing statistical summaries like minimum, maximum, mean, median, and IQR SAS Macros are ideal for automating quality control tasks, detecting, reporting, and repairing data quality errors, ensuring data integrity in projects SAS Macros offer functionalities like data cleansing and scrubbing routines to identify and remove outliers or incorrect values from datasets SAS Macros can assist in standardizing data formats, naming conventions, and documenting data cleaning processes to maintain data integrity SAS Macros can enhance reporting capabilities by generating reports in multiple formats using the Output Delivery System (ODS), aiding in monitoring data quality and integrity. Creation and organization of SAS EG projects to efficiently flag data for cleaning, identifying outliers, missing, incorrect, or duplicate data Implementation of SAS Macros to automate outlier detection, generate reports on data quality, and facilitate data corrections for improved data integrity Leveraging SAS Macros to standardize data formats, naming conventions, and document data cleaning procedures to ensure consistent and well-documented data integrity efforts. Integration of Advanced Machine Learning Algorithms Future trends involve integrating advanced machine learning algorithms into data cleaning processes to automate outlier detection, pattern recognition, and data quality assessment Enhanced Automation through AI and Natural Language Processing Automation through AI and natural language processing is expected to streamline data cleaning tasks by enabling automated data transformation, error detection, and data quality enhancement Utilization of Cloud Computing for Scalability The future trend includes leveraging cloud computing for scalable data cleaning and analysis processes, allowing for efficient handling of large datasets and complex data cleaning operations [1-3].



### **Conclusion**

Utilizing SAS Macros for Quality Control The papers emphasize the importance of utilizing SAS macros for automating quality control tasks, saving time and money while ensuring data integrity. Efficient Data Cleaning and Standardization The research highlights the efficiency of using SAS EG projects and macros for data cleaning, flagging outliers, and creating uniform clinical analytic datasets. Automation for Monitoring and Reporting The papers showcase how automation through SAS macros can enhance monitoring of site activities, safety-based medical monitoring, and tracking enrollment progression in clinical trials.

### **References**

1. G McQuown (2004) SAS® Macros are the Cure for Quality Control Pains. SUGI 29 093-29.
2. L Philpot, G Cantu (2012) Dirty Data? Clean it up with SAS. SCSUG [https://www.researchgate.net/publication/264492066\\_Dirty\\_Data\\_Clean\\_it\\_up\\_with\\_SAS](https://www.researchgate.net/publication/264492066_Dirty_Data_Clean_it_up_with_SAS).
3. W He (2022) Automation using SAS Makes it Easy to Monitor Dynamic Data in Clinical Trial. PharmaSUG 176.

**Copyright:** ©2023 Arvind Uttiramerur. This is an open-access article distributed under the terms of the Creative Commons Attribution License, which permits unrestricted use, distribution, and reproduction in any medium, provided the original author and source are credited.