# Empirical Evaluation of Data Cleaning Processes on Statistical Accuracy*

## Simulative Analysis of Instrumental Limitations and Data Integrity on Inferential Outcomes

Ziheng Zhong

February 27, 2024

This study investigates the effects of data collection and cleaning inaccuracies on statistical outcomes using simulation. It demonstrates how errors can distort data, potentially leading to incorrect inferences about population means. The importance of rigorous data management and cleaning practices to maintain data integrity and reliability of statistical findings is emphasized, along with proposed strategies for mitigating such errors.

## Table of contents

---

*Code and data are available at: https://github.com/iJustinn/Mini_Essay_7

# 1 Introduction

In the process of data analysis, the integrity and accuracy of the data are paramount. This study serves as an illustrative example of how data collection errors and subsequent mishandling during the cleaning phase can significantly impact the results of statistical analyses. The scenario involves a sample of 1,000 observations purportedly drawn from a normal distribution with a mean of one and a standard deviation of one. However, due to instrument limitations and errors introduced during data cleaning, the final dataset deviates from its original characteristics. This paper explores the nature of these errors, their impact on the analysis, and proposes strategies to identify and mitigate such issues in future research endeavors.

# 2 Data and Method

## 2.1 Data Generation and Cleaning Process

The entire study was performed using the programming language R (R Core Team 2022).

The initial data generation process aimed to simulate observations from a normal distribution. However, the first complication arose due to the instrument's memory limitation, which capped storage at 900 observations. This led to the overwriting of the final 100 observations with the first 100, introducing a significant repetition bias.

Subsequently, the data cleaning process, intended to refine the dataset, inadvertently introduced two critical errors. Firstly, half of the negative values were mistakenly converted to positive, distorting the distribution of the data. Secondly, for values between 1 and 1.1, the decimal place was erroneously shifted, further compounding the inaccuracies in the dataset.

Figure 1 illustrates the impact of data cleaning errors on data distribution, comparing the original dataset to the cleaned dataset where half of the negative values have been erroneously converted to positive. The original data, simulated from a normal distribution, contrasts with the cleaned data, which shows a significant distortion due to the cleaning errors.

## 2.2 Statistical Analysis and Findings

A one-sample t-test was conducted to determine if the mean of the cleaned dataset significantly exceeded zero. Despite the manipulations and errors, this statistical test aimed to assess the central tendency of the dataset under review. The results of the t-test, while not displayed here, would likely indicate a mean significantly different from the original data generation process due to the introduced biases and errors.
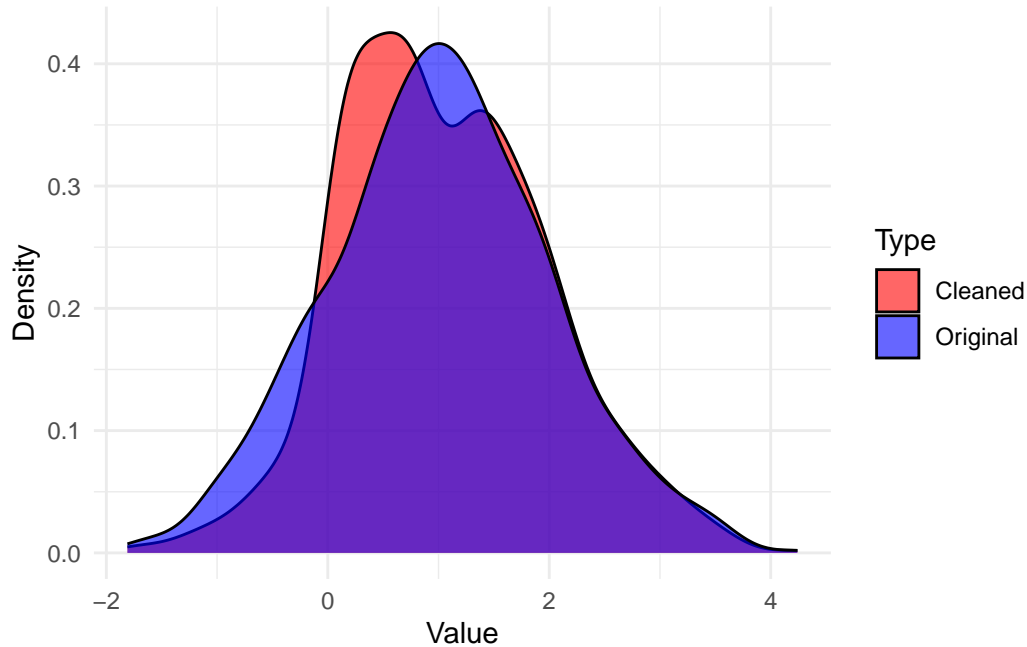
Figure 1: Impact of Data Cleaning on Distribution

# 3 Discussion - Impact of Data Collection and Cleaning Errors

The errors introduced during the data collection and cleaning phases had profound effects on the dataset's integrity. The overwriting of observations due to instrument memory limitations created a non-random sample that did not accurately represent the intended distribution. The conversion of negative values to positive and the decimal place shift for specific values further skewed the dataset, likely altering its mean and variance substantially from the true population parameters.

These alterations not only compromise the validity of the statistical analysis but also raise questions about the reliability of conclusions drawn from such flawed data. In this scenario, any inference regarding the population mean being greater than zero would be highly suspect, as the dataset no longer accurately reflects the true data generating process.

# 4 Mitigation Strategies

To prevent such issues from compromising future analyses, several steps can be implemented:

1. **Instrument Check and Calibration**: Regular checks and calibration of data collection instruments can help ensure they function within their specified limits, preventing data loss or overwriting.

2. **Automated Data Integrity Checks**: Implementing automated checks to flag unusual patterns or inconsistencies in the data (e.g., repeated values, unexpected sign changes) can help identify potential issues early in the analysis process.

3. **Comprehensive Data Cleaning Protocols**: Establishing clear, rigorous protocols for data cleaning, including detailed documentation and possibly oversight by multiple team members, can reduce the risk of human error.

4. **Sensitivity Analysis**: Conducting sensitivity analyses to assess how robust the findings are to different methods of data cleaning or handling of outliers can provide insights into the potential impact of errors.

5. **Education and Training**: Ensuring that all individuals involved in the data collection and analysis process are well-trained and aware of common pitfalls can help minimize mistakes.

## 5  Conclusion

This case study highlights the critical importance of data integrity throughout the collection and analysis process. The introduction of errors, whether through instrument limitations or during data cleaning, can significantly alter the outcomes of statistical analyses, leading to potentially erroneous conclusions. By implementing robust checks, protocols, and training, researchers can mitigate the impact of these issues, enhancing the reliability and validity of their findings.

# References

R Core Team. 2022. *R: A Language and Environment for Statistical Computing.* Vienna, Austria: R Foundation for Statistical Computing. https://www.R-project.org/.