

what is missing data and what should you do about it?*

Pengyu Sui

March 4, 2024

Abstract

In this study, I delve into the pervasive issue of missing data in statistical analysis, a challenge that significantly impacts the validity and reliability of research outcomes across various scientific domains. My exploration begins with an elucidation of the three primary types of missing data—Missing Completely At Random (MCAR), Missing At Random (MAR), and Missing Not At Random (MNAR)—and their respective implications on statistical inference. I engage in a thorough comparison of traditional and advanced methodologies for addressing missing data, including deletion methods, imputation techniques, and model-based approaches.

Introduction

Missing data is a ubiquitous challenge in statistical analysis and research, affecting nearly every domain where data collection and analysis occur. It can significantly impact the validity, reliability, and generalizability of research findings. Understanding the nature of missing data, its implications, and strategies for addressing it is crucial for producing accurate and meaningful analytical results. This essay delves into the concept of missing data, categorizes its types based on the underlying mechanisms, and explores approaches for handling it effectively.

Categories of Missing Data

The categorization of missing data into three main types—Missing Completely At Random (MCAR), Missing At Random (MAR), and Missing Not At Random (MNAR)—provides a framework for understanding its potential impacts and for selecting appropriate handling techniques:

1. **Missing Completely At Random (MCAR):** When the probability of missingness is the same for all units in the sample, it is termed MCAR. In this case, the missing data are independent of both observed and unobserved data. While this is the least problematic type of missingness, it is also the least common in practice. This very stringent condition is required in order for case deletion to be valid, and missing data is very rarely MCAR (Rubin, 1976).
2. **Missing At Random (MAR):** MAR occurs when the probability of missingness is related to observed data but not to the missing data itself. Though the term may be misleading, MAR allows for the relationship between missingness and observed data to be modeled, facilitating more robust handling techniques. MAR means is missing, but conditional on some other ‘X-variable’ observed in the data set, although not on the ‘Y-variable’ of interest (Schafer, 1997).
3. **Missing Not At Random (MNAR):** The most challenging type of missingness, MNAR, occurs when the probability of missingness is related to the missing data itself. This scenario requires careful consideration and sophisticated methods to address, as standard techniques may lead to biased results. MCAR and MAR are ignorable, for likelihood-based imputation methods, NMAR is not (Rubin, 1987).

Understanding these distinctions is crucial because the choice of handling technique may vary depending on the missing data mechanism. For instance, if data is MCAR, complete case analysis or simple imputation

*code and data are available at:<https://github.com/simon0202sui/what-is-missing-data-and-what-should-you-do-about-it-.git>
Acknowledge the review of Victor Chen

methods like mean or median substitution may be appropriate. However, if data is MAR or MNAR, more sophisticated techniques such as multiple imputation or maximum likelihood estimation may be necessary to mitigate bias and preserve the integrity of the analysis. By categorizing missing data into these types, researchers can better identify potential biases in their datasets and implement suitable strategies to handle missingness, thereby improving the validity and reliability of their analyses.

Handling Missing Data

With the understanding of why data is lost in the course of a study, the new goal is to deal with the lost data and minimize its impact on the results of the study. The approach to dealing with missing data depends on its type and the context of the research. However, common strategies include:

1. **Deletion Methods:** The simplest approach involves removing records with missing data from the analysis. While straightforward, this method can lead to significant data loss and potential bias, especially if the missing data are not MCAR. This can be either listwise (complete case only) or all value (Pairwise-available case), the cases are deleted which contain missing data, for the analysis being carried out (Schafer, 2002).
2. **Imputation Methods:** Imputation involves filling in missing values based on available data. Simple imputation techniques, such as mean or median imputation, can be useful but may underestimate variability. More sophisticated methods, like multiple imputation, generate several plausible datasets, analyze each, and then combine the results to account for the uncertainty introduced by the missing data.
3. **Model-Based Approaches:** Techniques such as maximum likelihood estimation (MLE) and Bayesian methods model the data and the missingness mechanism simultaneously, potentially offering a more rigorous solution to the missing data problem, particularly when data are not missing completely at random.

Choosing the Right Strategy

Selecting the appropriate strategy for handling missing data requires careful consideration of the type of missingness, the extent of the missing data, and the analysis goals. It is crucial to assess the potential impact of missing data on the research findings and to choose a method that minimizes bias and maximizes the validity of the results.

Conclusion

In this study, we tackled the pervasive issue of missing data in statistical analyses, highlighting its significant implications on the validity and reliability of research findings across various domains. By delineating the three primary types of missing data—MCAR, MAR, and MNAR—we provided insights into their distinct characteristics and the corresponding strategies required for effective management. Our comparative analysis of deletion methods, imputation techniques, and model-based approaches revealed the nuanced considerations necessary to choose the most appropriate method for handling missing data, depending on its underlying mechanism and the specific goals of the research.

Key findings underscore the importance of a tailored approach to missing data. For MCAR data, simple deletion might suffice without biasing the results significantly. However, for MAR and MNAR data, more sophisticated methods like multiple imputation or model-based approaches are essential to avoid misleading conclusions.

References

- Rubin, D.B., (1976) Inference and Missing Data. *Biometrika* 63 581-592
- Schafer, J.L., (1997) *The Analysis of Incomplete Multivariate Data*. Chapman & Hall
- Little, R.J.A., Rubin, D.B., (1987) *Statistical Analysis with Missing Data*. Wiley

Scheffer, J. (2002) Dealing with missing data, Research Letters in the Information and Mathematical Sciences, 3, 153-160