# What is Missing Data and What Should You do About it?*

Gavin Crooks

March 4, 2024

This paper explores the challenges of missing data in research, emphasizing its inevitability and multifaceted origins. I advocate for a nuanced understanding of Missing Completely at Random (MCAR), Missing at Random (MAR), and Missing Not at Random (MNAR) patterns. Simulation is highlighted as a crucial tool for experimenting with missing data scenarios and evaluating imputation strategies. Complete Case Analysis (CCA) is discussed as a straightforward but limited approach, contrasted with more advanced techniques like Multiple Imputation and Data Augmentation. The importance of simulation studies in guiding decision-making is underscored, emphasizing the need for tailored strategies aligned with missing data characteristics.

## What is Missing Data

Firstly, we must understand what missing data is. Missing data, an inevitable aspect of any research endeavor, refers to the absence of observations or values in a dataset, signaling a critical void in our understanding of the phenomena under investigation. No matter the meticulousness of data collection processes, certain variables, or even entire observations, elude measurement. This phenomenon extends beyond mere oversight, encapsulating scenarios where respondents opt for non-response, leading to nuances like non-response bias. The absence of information, whether intentional or inadvertent, poses a challenge to the validity and reliability of analyses. Acknowledging and addressing the intricacies of missing data is fundamental to ensuring the fidelity of research outcomes and advancing methodologies for handling these gaps in information. In this context, a comprehensive exploration of missing data becomes a crucial precursor to navigating the complexities of statistical inference and drawing meaningful conclusions from empirical investigations.

---

*Paper available at this repo: https://github.com/Crooksyyy/MissingData

The occurrence of missing data within a dataset can be attributed to a multitude of factors, reflecting the inherent challenges of real-world data collection. Non-response, a significant contributor, manifests when individuals, for various reasons, decline to participate in surveys or selectively omit specific questions. This could be driven by a range of factors such as survey fatigue, privacy concerns, or a lack of interest in certain topics. Additionally, errors in the data acquisition process, including administrative oversights, technical glitches, or issues with survey instruments, can introduce unintentional gaps. Furthermore, the nature of the information sought may play a role; sensitive or stigmatized topics might elicit higher non-response rates. The complexity of human behavior and the diverse motivations behind survey participation contribute to the intricate landscape of missing data, necessitating a nuanced understanding of these factors to implement effective strategies for addressing and mitigating the impact of these gaps on analytical outcomes.

In an ideal scenario, missing data would follow the Missing Completely at Random (MCAR) pattern, ensuring inference reflective of the broader population. However, reality often aligns more closely with Missing at Random (MAR) or Missing Not at Random (MNAR) scenarios, necessitating a nuanced approach. Simulation plays a pivotal role in grappling with the intricacies of missing data, offering a methodological compass for researchers navigating the uncertainties inherent in various missing data scenarios. Aki Vehtari (2022) classification MCAR, MAR, and MNAR, provides a roadmap for understanding the nature of missingness. Simulation, in this context, serves as a virtual laboratory, enabling controlled experimentation with missing observations or values under different mechanisms. This proactive approach affords researchers the opportunity to comprehend the potential biases and ramifications introduced by missing data across diverse conditions. By strategically simulating missing data scenarios, researchers gain insights into the performance of various imputation strategies. The nuanced understanding derived from simulated experiments enhances the transparency and robustness of chosen strategies, equipping researchers with a well-informed toolkit for addressing the challenges associated with missing data.

## What Should you do About it

The first possible way to handle missing data is Complete Case Analysis (CCA), a common approach in handling missing data, involves excluding observations with any missing values from the analysis. In essence, this method retains only the complete cases with all required information for the variables of interest. While CCA is straightforward and easy to implement, its effectiveness is contingent on the assumption that the missing data is MCAR. However, this assumption is often challenging to verify in practice, and the application of CCA may lead to reduced sample sizes, potentially introducing biases if the missingness is related to specific patterns or variables within the dataset. Despite its limitations, CCA remains a practical option, especially when the missing data is minimal, and its randomness can be reasonably assumed.

In contrast to Complete Case Analysis (CCA), various other techniques offer more sophisticated approaches for handling missing data. **Mean/Median Imputation** involves replacing missing values with the mean or median of observed values for the respective variable, providing a simple yet limited strategy. **Multiple Imputation**, on the other hand, creates multiple datasets with imputed values based on statistical models, allowing for a more robust estimation of missing values and associated uncertainties. **Model-Based Imputation** leverages relationships between variables to estimate missing values, requiring a deeper understanding of data structures. The Bayesian method of **Data Augmentation** introduces latent variables to model the missing data distribution, offering a flexible approach within Bayesian frameworks. Each technique comes with its own set of assumptions and considerations, emphasizing the importance of selecting an approach that aligns with the specific characteristics and intricacies of the dataset at hand. Again, stressing the importance simulation studies play when assessing the performance of these techniques and guiding informed decisions in addressing missing data challenges.

In summary, effectively dealing with missing data requires a thorough understanding of its nature and employing suitable handling techniques. Aki Vehtari (2022) framework categorizes missing data into MCAR, MAR, and MNAR, guiding the choice of appropriate strategies. Simulations, as emphasized in the literature, prove essential for assessing the implications of these strategies. Complete Case Analysis (CCA) offers a straightforward but limited approach, whereas advanced techniques like Multiple Imputation, Model-Based Imputation, and Data Augmentation provide more sophisticated options. The choice of method should align with the missing data characteristics, and researchers are encouraged to use simulations for comprehensive evaluations and informed decision-making in handling missing data challenges.

It is important to note this paper draws heavily from Alexander (2023) as well as Aki Vehtari (2022). I Would also like to thank Denise Chang for editing this paper. Please note the use of ChatGPT (2024) in the LLMUsage file also found in the the repo

# References

Aki Vehtari, Jennifer Hill, Andrew Gelman. 2022. "ROS Examples." https://avehtari.github.io/ROS-Examples/.

Alexander, Rohan. 2023. "Telling Stories with Data." https://tellingstorieswithdata.com/.

ChatGPT. 2024. "ChatGPT Conversation." Online.