# What is Missing Data and What Should You Do about It?*

## Hannah Yu

## February 22, 2024

## 1 Understanding Missing Data and Strategies for Handling It

Missing data is a common issue encountered in data analysis across various fields, ranging from social sciences to healthcare and finance. It refers to the absence of values or information in a dataset for certain variables or observations. The presence of missing data can significantly impact the accuracy and reliability of statistical analyses and machine learning models. Therefore, it is essential for data analysts and researchers to understand the nature of missing data and employ appropriate strategies to handle it effectively.

## 2 Types of Missing Data:

Missing data can occur for various reasons, and understanding these reasons is crucial for determining the appropriate handling strategy. Generally, missing data can be classified into three main types:

Missing Completely At Random (MCAR): In MCAR, the missingness of data is completely random and unrelated to any other variables in the dataset, observed or unobserved. This means that the probability of data being missing is the same for all observations. For example, if survey responses are lost due to a technical error in data collection, it can be considered MCAR.

Missing at Random (MAR): Missing at Random (MAR) occurs when the probability of data being missing depends on other observed variables in the dataset but not on the missing data itself. In other words, the missingness is related to the observed data but not to the missing values. For instance, if respondents with higher income are less likely to disclose their income

---

*Code and data are available at: LINK.

in a survey, the missingness of income data is related to income but not to the missing values of income itself.

Missing Not At Random (MNAR): Missing Not At Random (MNAR) happens when the probability of data being missing is related to the missing values themselves or to unobserved variables. In MNAR, the missingness is systematically related to the missing values, which can bias the analysis if not handled properly. For example, if respondents with higher income are less likely to disclose their income, and their income values are missing from the dataset, it constitutes MNAR.

# 3 Strategies for Handling Missing Data:

Dealing with missing data requires careful consideration and appropriate strategies to mitigate potential biases and inaccuracies in the analysis. Some common strategies for handling missing data include:

Complete Case Analysis (CCA): Complete Case Analysis involves excluding observations with missing values from the analysis. While this approach is straightforward, it can lead to loss of valuable information and potential bias, especially if the missingness is not completely random.

Imputation Methods: Imputation methods involve replacing missing values with estimated values based on the available data. One simple imputation method is mean imputation, where missing values are replaced with the mean of the observed values for that variable. However, mean imputation can underestimate the variability and lead to biased estimates, especially if the data are not MCAR. Other imputation methods include median imputation, mode imputation, regression imputation, and multiple imputation.

Multiple Imputation: Multiple imputation is a sophisticated imputation technique that involves generating multiple plausible values for each missing value and then combining the results to obtain final estimates. Multiple imputation accounts for the uncertainty associated with imputed values and provides more accurate estimates compared to single imputation methods.

Model-Based Methods: Model-based methods involve using statistical models to predict missing values based on the observed data. These methods can provide more accurate estimates, especially when the missingness mechanism is complex and the data are not MCAR. Examples of model-based methods include linear regression, logistic regression, and machine learning algorithms.

# 4 Conclusion:

Missing data is a common challenge in data analysis, and addressing it effectively is essential for producing accurate and reliable results. Understanding the types of missing data and their implications is crucial for selecting appropriate handling strategies. While no single approach is universally applicable, a combination of techniques, such as complete case analysis, imputation methods, and model-based methods, can help mitigate biases and uncertainties associated with missing data. Additionally, transparency and documentation of the chosen approach are essential for ensuring the reproducibility and validity of the analysis results. By carefully handling missing data, researchers can enhance the quality and reliability of their findings, thereby making informed decisions and driving meaningful insights from the data.

# 5 References