# Assignment #1

## EXPLORATORY DATA-ANALYSIS

## Group ID: 5

**Assignment Title:** `missingno` package

## Group Members:

Madhav Jivani
**202201285**

Ritwik Agarwal
**202411067**

Tarun Kumar
**202312126**

# 1   Introduction to Missing Data

Missing data is a common occurrence in datasets, and it happens for a variety of reasons. This can include errors during data collection, accidental deletions, or even technical limitations when capturing measurements. Understanding missing data is important because it affects the quality of our analyses and machine learning models. In data science, there are three main types of missing data that we focus on:

- **Missing Completely at Random (MCAR):** The missing values occur completely by chance, and there's no specific pattern. For instance, if a sensor randomly fails in a small percentage of cases without any specific reason, that would be MCAR.

- **Missing at Random (MAR):** Here, the missingness is related to other observed data, but not the data that is missing. For example, in a survey, older participants might be less likely to answer a certain question, so the missingness is related to age (an observed factor).

- **Missing Not at Random (MNAR):** The missingness depends on the data that is missing itself. For instance, people with high incomes might choose not to reveal their income, meaning that the missingness is directly related to the data point that's missing.

Each type of missing data requires different handling techniques to avoid introducing bias or errors into our analyses.

# 2   Challenges of Missing Data in Data Science

Missing data presents several challenges:

- It can bias our results, meaning that our conclusions may be wrong if we don't handle the missing data correctly.

- Some machine learning algorithms require complete datasets, and missing data can reduce the performance or force us to discard useful data points.

- In some cases, the pattern of missing data itself can provide valuable information, helping us understand the data better.

Handling missing data properly is crucial to ensuring the accuracy and integrity of our analyses.

# 3   Overview of `missingno` Package

The `missingno` package is a Python library that helps us visualize missing data in a dataset. Instead of trying to look at numbers or rows directly, it allows us to see missing data as patterns and correlations, making it easier to decide how to handle it. The package integrates smoothly with the widely used `pandas` library.

# 4   Types of Visualizations in `missingno`

The `missingno` package provides several useful visualizations:

- **Bar Plot:** This plot shows the number of non-missing (or present) values in each column of a dataset. It helps to quickly assess how much data is missing in each column.

- **Matrix Plot:** A matrix visualization shows exactly where data is missing. It highlights the gaps row by row and column by column.

- **Heatmap:** The heatmap shows the correlation between missing values in different columns. A high correlation means if one column has missing data, the other one likely does too.

- **Dendrogram:** This is a clustering visualization that groups columns with similar patterns of missing data together, helping to find relationships between columns where missing data may be related.

# 5    Interpreting Missing Data Visualizations

Each of these visualizations tells us something different:

- The **bar plot** shows which columns have the most missing values, which is useful to decide if we should remove or impute those columns.

- The **matrix plot** helps identify where and how missing data is spread across the dataset.

- The **heatmap** helps identify columns with correlated missing values, which could mean we need to treat them together.

- The **dendrogram** shows which columns are related in terms of missingness, providing insights into underlying relationships.

# 6    Applications and Use Cases

The `missingno` package is valuable for many common data science tasks:

- **Data Cleaning:** Before applying machine learning algorithms, we can clean the data by understanding where missing values occur.

- **Exploratory Data Analysis (EDA):** During EDA, `missingno` can provide quick insights into missing data and help us understand its extent and impact.

- **Preprocessing for Machine Learning:** It helps decide on the best strategies for handling missing data before applying predictive models.

# 7    Types of Missing Data Targeted

As mentioned earlier, missing data falls into three main categories:

- **Missing Completely at Random (MCAR):** No specific pattern, and missingness is not related to any other variable in the dataset.

- **Missing at Random (MAR):** Missingness is related to other observed variables in the dataset, but not the missing value itself.

- **Missing Not at Random (MNAR):** Missingness is related to the missing value itself, making it hard to detect or handle.

The `missingno` package helps us visualize and identify which type of missing data we are dealing with by analyzing correlations and patterns in the missing values.

## 7.1 Correlation of Missing Data

The `missingno` package uses correlation to understand the relationships between columns with missing data.

- **High Positive Correlation:** This indicates that missing values tend to occur together in two or more columns. For example, if two columns have a high correlation, missing values in one column might imply that the other is also likely missing. This can be seen in MAR or MNAR data where some factors drive missingness.

- **Negative Correlation:** If one column is missing data, the other is less likely to have missing values. This might indicate complex dependencies between columns, possibly reflecting MNAR patterns.

Different correlation patterns can suggest the type of missing data:

- **MCAR:** In this case, there should be no or very low correlation between missing values in different columns. Since missingness occurs randomly, we do not expect to see any strong relationship in the heatmap.

- **MAR:** Missing values in one column may be associated with observed values in another. The heatmap might show moderate positive correlations between columns, indicating that missingness is related to some other observed variables.

- **MNAR:** While MNAR data is harder to detect, the heatmap could show positive correlations between missing values in different columns. Dendrograms may also cluster columns with similar missing data tendencies.

These visualizations help us determine which kind of strategy is appropriate for handling the missing data, whether it's imputation or removal.

# 8 Identifying Missing Data Patterns from Plots

Each type of visualization offers clues about the nature of the missing data:

- **Bar Plot:**
  - If the missing data is spread evenly across all columns, this might suggest MCAR, as missingness is random.
  - If one or two columns have much more missing data than others, this might suggest MAR or MNAR, as the missingness could be related to specific factors.

- **Matrix Plot:**
  - If missing values appear scattered randomly across the rows, this is a sign that the data could be MCAR.
  - If missing values appear in blocks or clusters, this might indicate that the missingness is related to specific patterns, which is typical in MAR or MNAR.

- **Heatmap:**
  - **Low Correlation:** If the heatmap shows very little correlation between columns, this suggests MCAR.
  - **Positive Correlation:** If the heatmap shows high positive correlations between missing values in some columns, this is likely MAR, where missingness is related to other observed variables.
  - **Negative Correlation:** Negative correlations between columns are rare, but they could indicate complex relationships, possibly MNAR.

- **Dendrogram:**

  - Strong clusters of columns with similar missing data patterns indicate a possible MAR or MNAR relationship.
  - If columns don't cluster well, this suggests that the missingness might be random (MCAR).

These visual patterns give us a better understanding of the structure and reason behind missing data in our dataset.

# 9    Conclusion

Handling missing data properly is critical to ensuring accurate analysis and reliable machine learning models. The `missingno` package offers a variety of tools to visualize missing data patterns and correlations, which can help guide the best strategies for handling missing data.

By using the bar plot, matrix plot, heatmap, and dendrogram, we can:

- Get a quick sense of how much data is missing and where.

- Understand how missing data is related between different columns.

- Decide if we can safely remove, impute, or apply more sophisticated techniques to deal with missing data.

With these insights, we can confidently approach missing data, ensuring that our analysis is more accurate and less prone to bias.