# Assignment #2

## EXPLORATORY DATA-ANALYSIS

## Group ID: 5

## Assignment Title: `missingpy` package and MCAR test

## Group Members:

Madhav Jivani
**202201285**

Ritwik Agarwal
**202411067**

Tarun Kumar
**202312126**

# 1   Introduction

Handling missing data is a critical step in data preprocessing and analysis. Missing values can skew results and lead to incorrect conclusions if not addressed properly. The `missingpy` library offers advanced imputation techniques that surpass the capabilities of basic methods provided by standard libraries. This document provides a comprehensive exploration of `missingpy`, its key functionalities, and how it interacts with statistical tests like Little's MCAR test to handle and understand missing data.

# 2   Overview of `missingpy` Library

`missingpy` is a specialized Python library designed to address the challenges associated with missing data in datasets. Unlike basic imputation techniques, `missingpy` employs advanced algorithms that are particularly suited for complex datasets where traditional methods may fall short.

## 2.1   Basic vs. Advanced Imputation Methods

Traditional imputation methods provided by libraries like `pandas` include:

- **Mean Imputation**: Replacing missing values with the mean of the observed values in a column.

- **Median Imputation**: Using the median value to fill missing entries, which is less sensitive to outliers compared to the mean.

- **Mode Imputation**: Filling missing values with the most frequent value in the dataset.

While these methods are straightforward, they do not account for the relationships between features and can introduce biases if the missing data is not randomly distributed. Advanced methods provided by `missingpy` address these limitations by considering the interactions between features.

## 2.2   Key Imputation Algorithms in `missingpy`

`missingpy` incorporates the following advanced imputation algorithms:

### 2.2.1   K-Nearest Neighbors Imputation (KNNImputer)

**Concept**: KNNImputer is based on the principle that similar records (neighbors) can provide meaningful estimates for missing values. The algorithm searches for K nearest neighbors to a record with missing data and uses their values to impute the missing entry.

**Steps in KNNImputer**:

1. **Distance Calculation**: Calculate the distance between records using metrics like Euclidean distance.

2. **Neighbor Selection**: Identify the K closest records to the one with missing data.

3. **Imputation**: Compute the mean (or weighted mean) of the K neighbors' values for the missing entry.

**Advantages**:

- Captures local patterns and relationships between features.

- Adapts to varying data distributions and correlations.

**Implementation Example**:

```
from missingpy import KNNImputer
import numpy as np

# Example dataset with missing values
data = np.array([[1, 2, np.nan], [3, np.nan, 1], [2, 5, 3]])

# Initialize KNNImputer
```

```
imputer = KNNImputer(n_neighbors=2)

# Fit and transform the data
imputed_data = imputer.fit_transform(data)
```

### 2.2.2 MissForest Imputation

**Concept**: MissForest is a non-parametric imputation method that uses an ensemble of decision trees (Random Forest) to predict missing values. It iteratively refines the imputation by using predictions from a Random Forest model built on the non-missing data.

    **Steps in MissForest**:

1. **Initial Imputation**: Start with simple imputation methods like mean imputation.

2. **Random Forest Modeling**: Build a Random Forest model using the non-missing data to predict missing values.

3. **Iterative Refinement**: Update the imputed values and re-train the Random Forest model until convergence.

    **Advantages**:

- Handles both numerical and categorical data.

- Captures complex relationships and interactions between features.

    **Implementation Example**:

```
from missingpy import MissForest

# Example dataset with missing values
data = np.array([[1, 2, np.nan], [3, np.nan, 1], [2, 5, 3]])

# Initialize MissForest
imputer = MissForest()

# Fit and transform the data
imputed_data = imputer.fit_transform(data)
```

# 3 Identifying MCAR Data with Little's MCAR Test

To effectively use imputation methods, it's crucial to understand the nature of missing data. Little's MCAR test is a statistical method used to determine whether the missing data is Missing Completely at Random (MCAR). If data is MCAR, it means that the missingness is unrelated to any observed or unobserved variables in the dataset, simplifying the imputation process.

## 3.1 Concept of MCAR

**MCAR (Missing Completely at Random)**: This type of missingness occurs when the probability of a data point being missing is entirely random and does not depend on any other observed or missing values in the dataset.

    **Implications of MCAR**:

- Simple imputation methods (e.g., mean or median) are generally appropriate.

- Analysis and conclusions drawn from the data are less likely to be biased by the missingness.

## 3.2 Performing Little's MCAR Test using `statsmodels`

`statsmodels` provides tools for conducting Little's MCAR test. The test assesses whether the pattern of missing data is consistent with the assumption that missingness is completely random.

**Steps to Perform the Test**:

1. **Install `statsmodels`**: Ensure that `statsmodels` is installed in your Python environment.

   ```
   pip install statsmodels
   ```

2. **Prepare the Dataset**: Convert your dataset into a format compatible with the MCAR test.

3. **Conduct the Test**: Use the function from `statsmodels` to perform Little's MCAR test.

   ```python
   import statsmodels.api as sm
   import pandas as pd

   # Load your dataset
   df = pd.DataFrame({
       'A': [1, 2, None, 4],
       'B': [None, 2, 3, 4],
       'C': [1, None, 3, 4]
   })

   # Perform Little's MCAR test
   test_results = sm.stats.mcar(df)

   # Print the results
   print(test_results)
   ```

   **Interpreting Results**:

   - **High p-value (typically greater than 0.05)**: Suggests that the missing data is MCAR.
   - **Low p-value (typically less than 0.05)**: Indicates that the missing data is not MCAR, suggesting the presence of a pattern.

# 4 Alternative Methods for Handling Missing Data

Besides KNNImputer and MissForest, several other techniques are available for handling missing data:

## 4.1 Multiple Imputation

**Concept**: Multiple imputation involves creating several different imputed datasets, analyzing each one separately, and then combining the results. This approach accounts for the uncertainty of the missing data.

**Advantages**:

- Provides a range of plausible values for missing data.
- Reduces bias and provides more robust statistical inferences.

## 4.2 Expectation-Maximization (EM) Algorithm

**Concept**: The EM algorithm estimates missing values through an iterative process that alternates between expectation (E) and maximization (M) steps.

**Advantages**:

- Handles complex datasets with missing values in a systematic manner.
- Provides maximum likelihood estimates for missing data.

## 4.3   Data Augmentation

**Concept**: Data augmentation involves generating synthetic data based on the observed data to replace missing values. Techniques like bootstrapping can be used to create additional datasets.

**Advantages**:

- Improves the robustness of the analysis by expanding the dataset.

- Provides additional information to improve imputation accuracy.

# 5   Conclusion

To effectively handle missing data and ensure accurate analysis, it is crucial to understand the nature of the missingness. The `missingpy` library offers advanced imputation techniques such as **KNNImputer** and **MissForest** that provide more sophisticated solutions than simple mean or median imputation methods.

While `missingpy` excels in imputing missing data, it does not directly test the nature of missingness. Combining it with tools like `statsmodels` for conducting Little's MCAR test allows for a comprehensive approach to missing data handling. Identifying whether data is MCAR helps in choosing the appropriate imputation method, ensuring that the analysis is both accurate and reliable.

By exploring various imputation methods and understanding the underlying missingness mechanisms, data scientists can improve the quality of their data analysis and draw more accurate conclusions from their datasets.