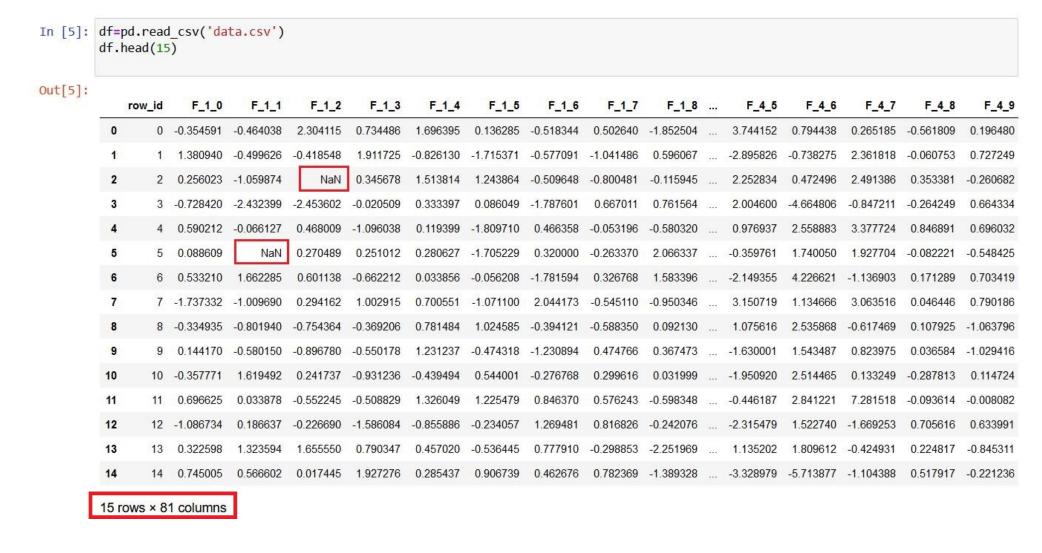# Dealing with Missing Values in Dataset

## Mousa Tayseer Jafar
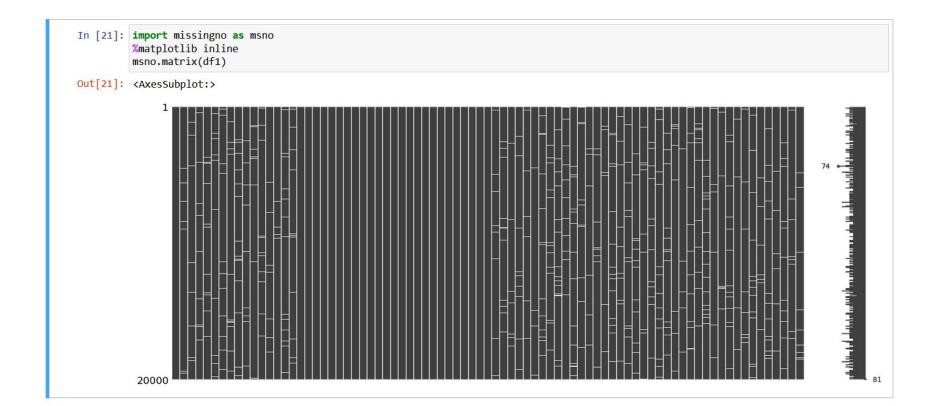
### Deakin University, Australia

## Introduction

**What is a Missing Value?** Missing data is defined as the values or data that is not stored (or not present) for some variable/s in the given dataset. Missing value can bias the results of the machine learning models and/or reduce the accuracy of the model. Below is a sample of the missing data from the Tabular Playground Series – June 2022 dataset.

## Visualization Missing Values

**How is Missing Value Represented In The Dataset?** In Pandas, usually, missing values are represented by NaN.
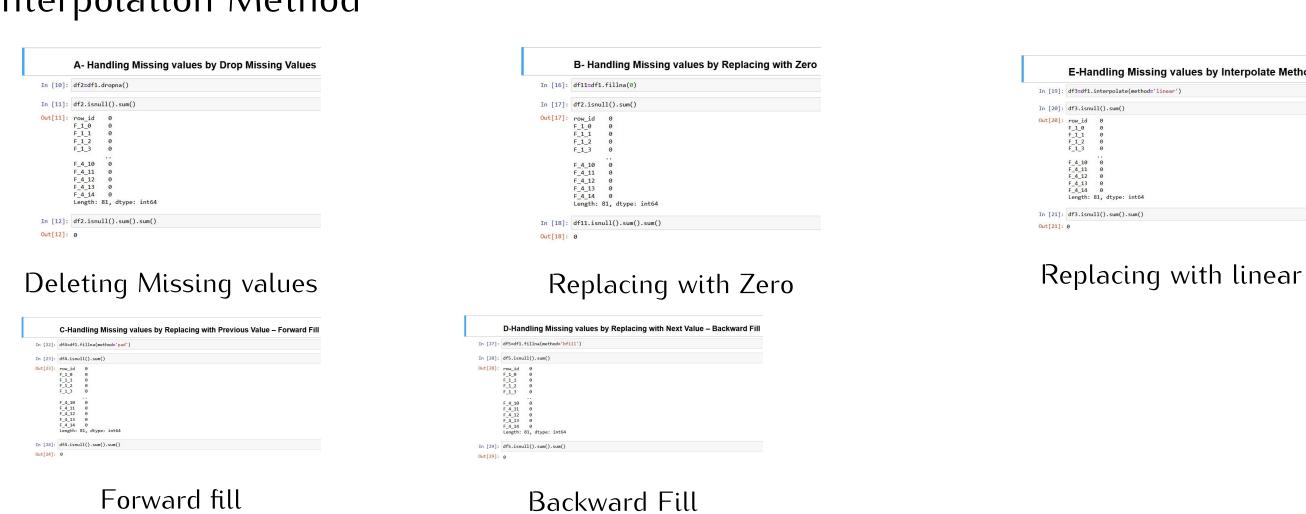
- The real-world data often has a lot of missing values.
- The cause of missing values can be data corruption or failure to record data.
- The handling of missing data is very important during the preprocessing of the dataset as many machine learning algorithms do not support missing values.
- Three problems associated with missing values:
  - Loss of efficiency,
  - Complications in handling and analyzing the data,
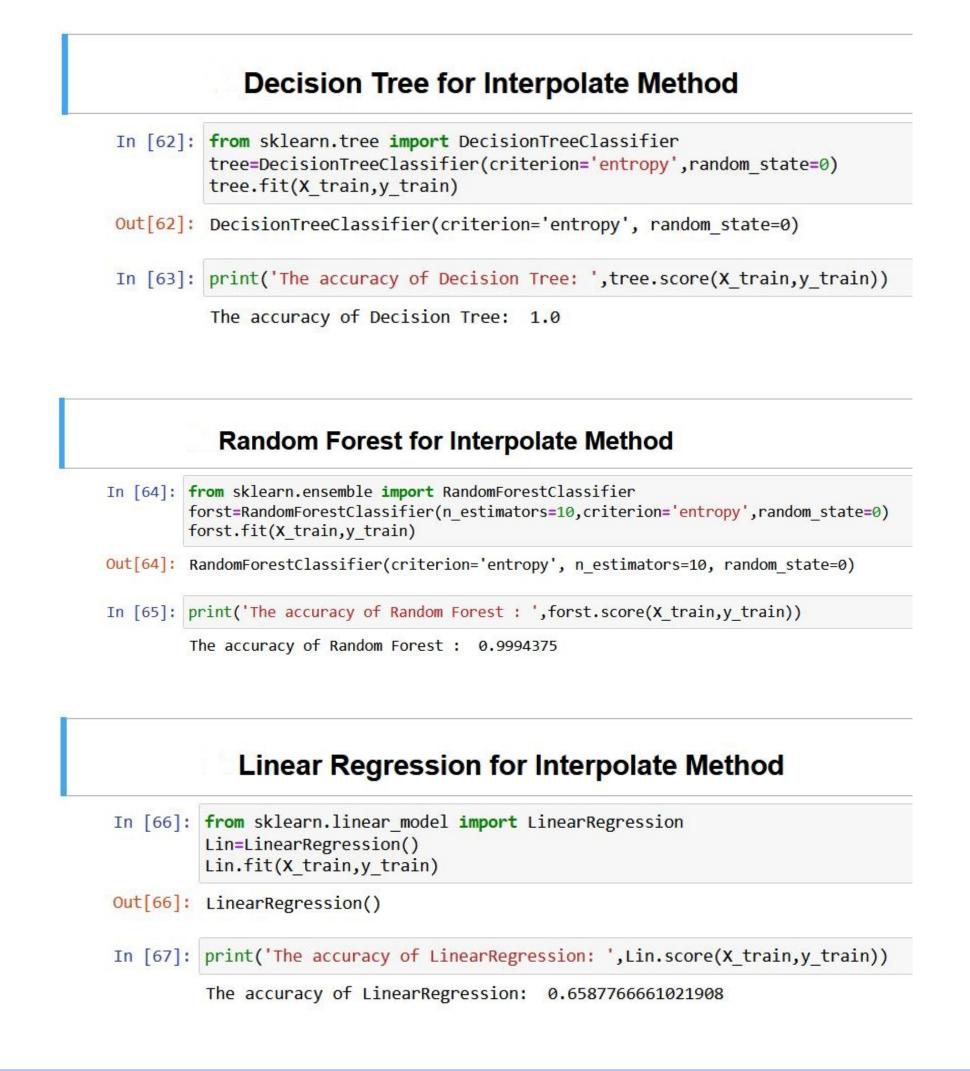  - Bias resulting from differences between missing and complete data.

## Data Preprocessing

**How To Handle Missing Values?** Analyze each column with missing values carefully to understand the reasons behind the missing values, as it is crucial to find out the strategy for handling the missing values. There are many ways of handling missing values We propose the below methods to handle the missing values:

- Deleting the Missing values
- Replacing With Arbitrary Value (Zero)
- Replacing with Previous Value  Forward Fill
- Replacing with Next Value  Backward Fill
- Interpolation Method

Deleting Missing values

Replacing with Zero

Replacing with linear

Forward fill

Backward Fill

## Applying Machine Learning Models

Missing data is a common problem in statistical analysis. In this work, we have been applying three Machine Learning Algorithms to measure accuracy after handling the missing values **Decision Tree, Random Forest, and Linear Regression.**

**Decision Tree for Interpolate Method**

**Random Forest for Interpolate Method**
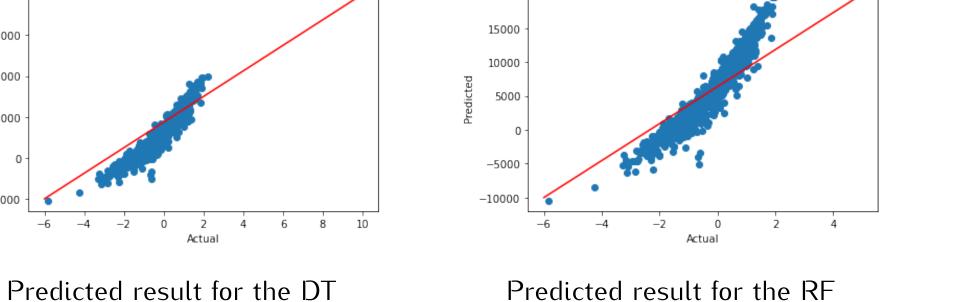
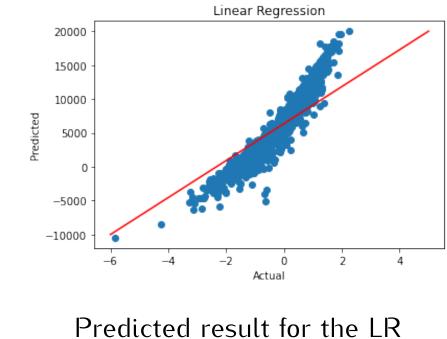**Linear Regression for Interpolate Method**

## Experiment

Missing data is a common problem in statistical analysis.
Random Forest, Decision Tree, and Linear Regression, were used as Machine Learning Methods to measure the Accuracy after handling the missing values. Accuracy is generated using the number of correct predictions on the test dataset to find the actual class label against the predicted class label for each category. The accuracy represents the total correct prediction overall the total prediction, as shown in equation 1.

$$Accuracy = (TP + TN)/(TP + TN + FP + FN) \qquad (1)$$

| Handling Missing Values | Decision Tree | Random Forest | Linear Regression |
|---|---|---|---|
| Drop Missing Values | 1.0 | 0.9994858 | 0.9047421 |
| Missing Values Zero | 1.0 | 0.9993125 | 0.7113801 |
| Interpolate Method | 1.0 | 0.9994375 | 0.6587766 |
| Forward Fill | 1.0 | 0.9993125 | 0.6336116 |
| Backward Fill | 1.0 | 0.9993125 | 0.6207112 |

Predicted result for the DT

Predicted result for the RF

Predicted result for the LR

## Conclusion

This work has introduced an approach for identifying and handling the missing values in the dataset. Several experiments have been conducted with three ML algorithms (DT, RF, and LR) to evaluate the efficiency and the performance of these approaches. All tests are experimented based on the Tabular Playground June 2022 dataset. Experiments have shown that the Decision be affected by handling missing values where the equals 100

Acknowledgement
- FLIP00 team