# **Handle Missing Values**

Date	17 November 2022
Team ID	PNT2022TMID45545
Project name	MACHINE LEARNING BASED VEHICLE
	PERFORMANCE ANALAYZER

# **Handle Missing Values in Machine Learning**

Popular strategies to handle missing values in the dataset

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.

This article covers 7 ways to handle missing values in the dataset:

- 1. Deleting Rows with missing values
- 2. Impute missing values for continuous variable
- 3. Impute missing values for categorical variable

- 4. Other Imputation Methods
- 5. Using Algorithms that support missing values
- 6. Prediction of missing values
- 7. Imputation using Deep Learning Library Datawig

# Data used is <u>Titanic Dataset</u> from Kaggle

```
data = pd.read_csv("train.csv")
msno.matrix(data)
```

(Image by Author), Visualization of Missing Values: white lines denote the presence of missing value

### **Delete Rows with Missing Values:**

Missing values can be handled by deleting the rows or columns having null values. If columns have more than half of the rows as null then the entire column can be dropped. The rows which are having one or more columns values as null can also be dropped.

(Image by Author) Left: Data with Null values, Right: Data after removal of Null values

### **Pros:**

• A model trained with the removal of all missing values creates a robust model.

#### Cons:

- Loss of a lot of information.
- Works poorly if the percentage of missing values is excessive in comparison to the complete dataset.

### Impute missing values with Mean/Median:

Columns in the dataset which are having numeric continuous values can be replaced with the mean, median, or mode of remaining values in the column. This method can prevent the loss of data compared to the earlier method. Replacing the above two approximations (mean, median) is a statistical approach to handle the missing values.

(Image by Author) Left: Age column before Imputation, Right: Age column after imputation by the mean value

The missing values are replaced by the mean value in the above example, in the same way, it can be replaced by the median value.

- Prevent data loss which results in deletion of rows or columns
- Works well with a small dataset and is easy to implement.

### Cons:

- Works only with numerical continuous variables.
- Can cause data leakage
- Do not factor the covariance between features.

### Imputation method for categorical columns:

When missing values is from categorical columns (string or numerical) then the missing values can be replaced with the most frequent category. If the number of missing values is very large then it can be replaced with a new category.

(Image by Author) Left: Data before Imputation, Right: Cabin column after imputation by 'U'

- Prevent data loss which results in deletion of rows or columns
- Works well with a small dataset and is easy to implement.
- Negates the loss of data by adding a unique category

### Cons:

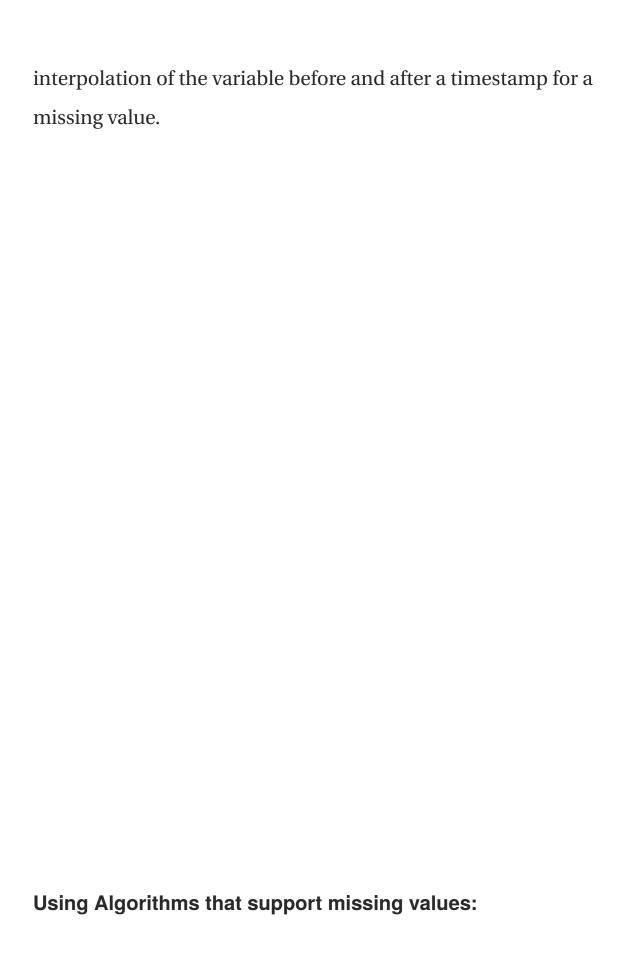
- Works only with categorical variables.
- Addition of new features to the model while encoding, which may result in poor performance

### **Other Imputation Methods:**

Depending on the nature of the data or data type, some other imputation methods may be more appropriate to impute missing values.

For example, for the data variable having longitudinal behavior, it might make sense to use the last valid observation to fill the missing





All the machine learning algorithms don't support missing values but some ML algorithms are robust to missing values in the dataset. The k-NN algorithm can ignore a column from a distance measure when a value is missing. Naive Bayes can also support missing values when making a prediction. These algorithms can be used when the dataset contains null or missing values.

The sklearn implementations of naive Bayes and k-Nearest Neighbors in Python do not support the presence of the missing values.

Another algorithm that can be used here is RandomForest that works well on non-linear and categorical data. It adapts to the data structure taking into consideration the high variance or the bias, producing better results on large datasets.

### **Pros:**

 No need to handle missing values in each column as ML algorithms will handle them efficiently.

### Cons:

• No implementation of these ML algorithms in the scikit-learn library.

### **Prediction of missing values:**

In the earlier methods to handle missing values, we do not use the correlation advantage of the variable containing the missing value and other variables. Using the other features which don't have nulls can be used to predict missing values.

The regression or classification model can be used for the prediction of missing values depending on the nature (categorical or continuous) of the feature having missing value.

```
Here 'Age' column contains missing values so for prediction of null values the spliting of data will be, y_train: rows from data["Age"] with non null values
```

y\_test: rows from data["Age"] with null values

**X\_train:** Dataset except data["Age"] features with non null values

**X\_test:** Dataset except data["Age"] features with null values



- Gives a better result than earlier methods
- Takes into account the covariance between the missing value column and other columns.

#### Cons:

• Considered only as a proxy for the true values

## Imputation using Deep Learning Library — <u>Datawig</u>

This method works very well with categorical, continuous, and nonnumerical features. Datawig is a library that learns ML models using Deep Neural Networks to impute missing values in the datagram.

```
Install datawig library,
pip3 install datawig
```

Datawig can take a data frame and fit an imputation model for each column with missing values, with all other columns as inputs.



- Quite accurate compared to other methods.
- It supports CPUs and GPUs.

#### Cons:

• Can be quite slow with large datasets.

### **Conclusion:**

Every dataset has missing values that need to be handled intelligently to create a robust model. In this article, I have discussed 7 ways to handle missing values that can handle missing values in every type of column. There is no thump rule to handle missing values in a particular manner, the method which gets a robust model with the best performance. One can use various methods on different features depending on how and what the data is about. Having domain knowledge about the dataset is important, which can give an insight into how to preprocess the data and handle missing values.

# **References:**

[1] Datawig: <a href="https://github.com/awslabs/datawig">https://github.com/awslabs/datawig</a>