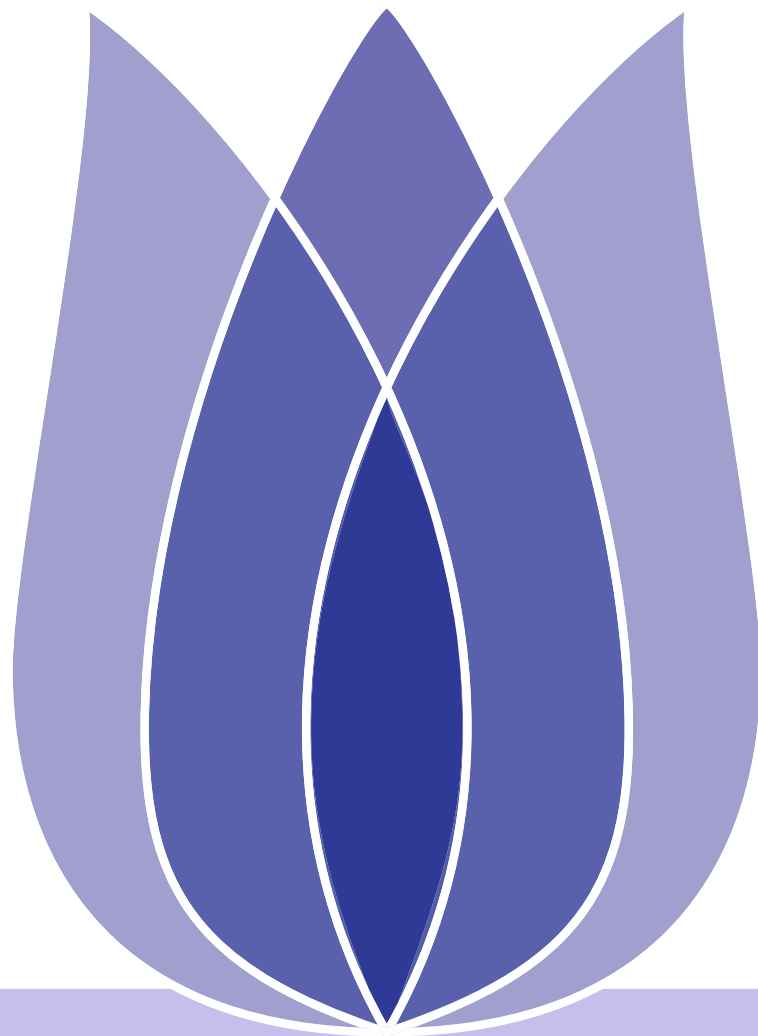


Dealing with Missing Values in Dataset

Mousa Tayseer Jafar

Deakin University - Australia

2022-08-07





Overview

[Problem Definition](#)

[Data Processing](#)

[Handle Missing Values](#)

[Build The Models](#)

[Conclusion](#)

Problem Definition

Definition

Why Do We Need To Care About Handling Missing Data?

Data Processing

Loading Data

Check for missing values in the dataset

Handle Missing Values

A-Deleting the Missing values

B-Replacing With Arbitrary Value (Zero)

C-Replacing with Previous Value – Forward Fill

D-Replacing with Next Value – Backward Fill

E-Interpolation Method

Build The Models

1- Linear Regression Model

2- Decision Tree Model

3- Random Forest Model

Conclusion



Problem Definition

Definition
Why Do We Need To Care About
Handling Missing Data?

Data Processing

Handle Missing Values

Build The Models

Conclusion

Problem Definition



- Problem Definition
- Definition
- Why Do We Need To Care About Handling Missing Data?
- Data Processing
- Handle Missing Values
- Build The Models
- Conclusion

Defn

What is a Missing Value?

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



Figure 1: Missing vlaues



Definition

- Problem Definition
- Definition**
 - Why Do We Need To Care About Handling Missing Data?
- Data Processing
- Handle Missing Values
- Build The Models
- Conclusion

How is Missing Value Represented In The Dataset? In Pandas, usually, missing values are represented by **NaN**.

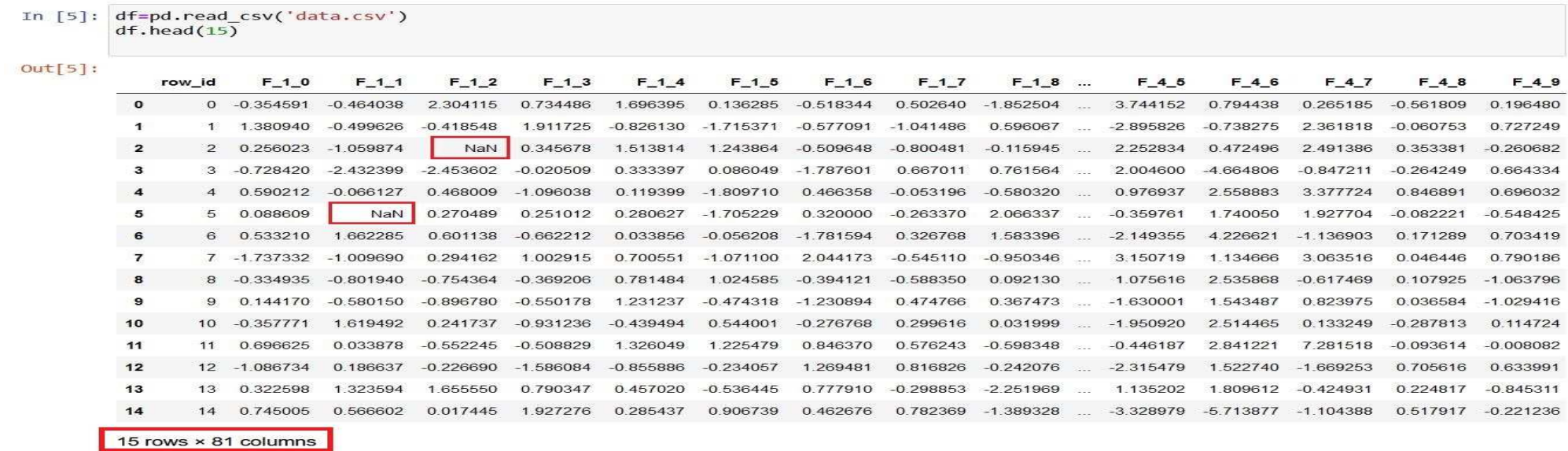


Figure 2: in each Row

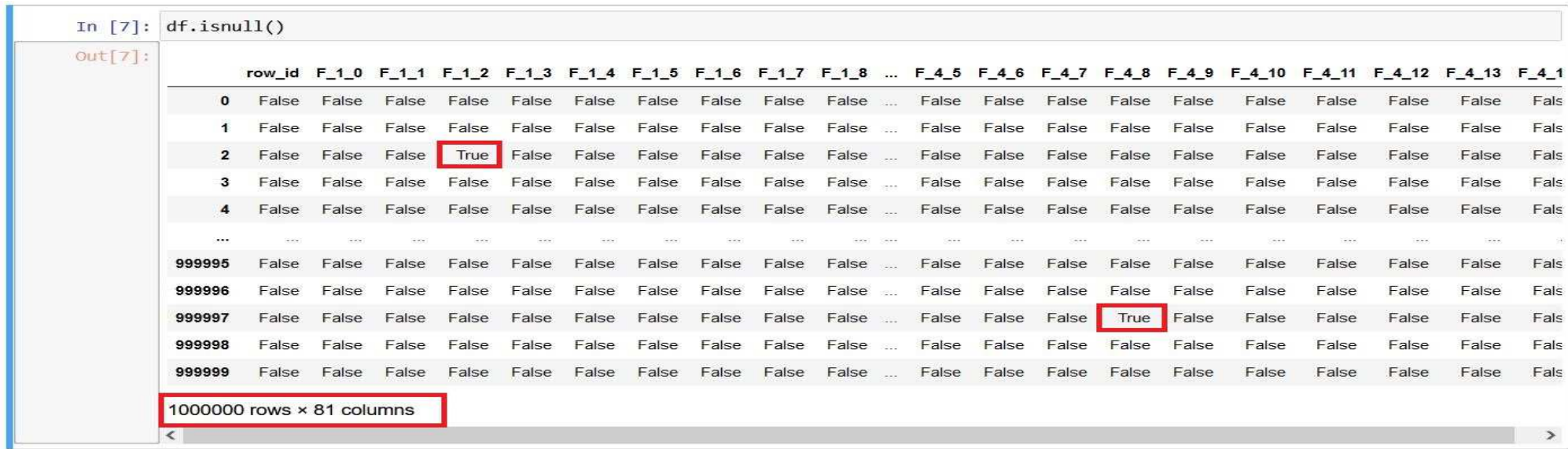


Figure 3: in all dataset

Why Do We Need To Care About Handling Missing Data?

Problem Definition

Definition

Why Do We Need To Care About Handling Missing Data?

Data Processing

Handle Missing Values

Build The Models

Conclusion

- 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.



Figure 4: Types Of Missing Values



TULIP

Team for Universal Learning and Intelligent Processing



[Problem Definition](#)

[Data Processing](#)

[Loading Data](#)

[Check for missing values in the dataset](#)

[Handle Missing Values](#)

[Build The Models](#)

[Conclusion](#)

Data Processing



Loading Data

- Problem Definition
- Data Processing
- Loading Data**
- Check for missing values in the dataset
- Handle Missing Values
- Build The Models
- Conclusion

```
In [1]: import pandas as pd

In [2]: df=pd.read_csv('data.csv')
df.head()

Out[2]:
```

	row_id	F_1_0	F_1_1	F_1_2	F_1_3	F_1_4	F_1_5	F_1_6	F_1_7	F_1_8	...	F_4_5	F_4_6	F_4_7	F_4_8
0	0	-0.354591	-0.464038	2.304115	0.734486	1.696395	0.136285	-0.518344	0.502640	-1.852504	...	3.744152	0.794438	0.265185	-0.561809
1	1	1.380940	-0.499626	-0.418548	1.911725	-0.826130	-1.715371	-0.577091	-1.041486	0.596067	...	-2.895826	-0.738275	2.361818	-0.060753
2	2	0.256023	-1.059874	NaN	0.345678	1.513814	1.243864	-0.509648	-0.800481	-0.115945	...	2.252834	0.472496	2.491386	0.353381
3	3	-0.728420	-2.432399	-2.453602	-0.020509	0.333397	0.086049	-1.787601	0.667011	0.761564	...	2.004600	-4.664806	-0.847211	-0.264249
4	4	0.590212	-0.066127	0.468009	-1.096038	0.119399	-1.809710	0.466358	-0.053196	-0.580320	...	0.976937	2.558883	3.377724	0.846891

5 rows × 81 columns

<>

```
In [3]: df.shape

Out[3]: (1000000, 81)
```

Figure 5: Load dataset

- Load the Tabular Playground Series - June 2022 dataset
- import the dataset from csv files

Check for missing values in the dataset

Problem Definition
Data Processing
Loading Data
Check for missing values in the dataset
Handle Missing Values
Build The Models
Conclusion

```
In [8]: df.isnull().sum()
Out[8]: row id      0
        F_1_0    18397
        F_1_1    18216
        F_1_2    18008
        F_1_3    18250
        ...
        F_4_10    18225
        F_4_11    18119
        F_4_12    18306
        F_4_13    17995
        F_4_14    18267
        Length: 81, dtype: int64
```

Figure 6: missing values in Each row

```
In [9]: df.isnull().sum().sum()
Out[9]: 1000000
```

Figure 7: missing values in the dataset

Visualization Missing Values

Problem Definition
Data Processing
Loading Data
Check for missing values in the dataset
Handle Missing Values
Build The Models
Conclusion

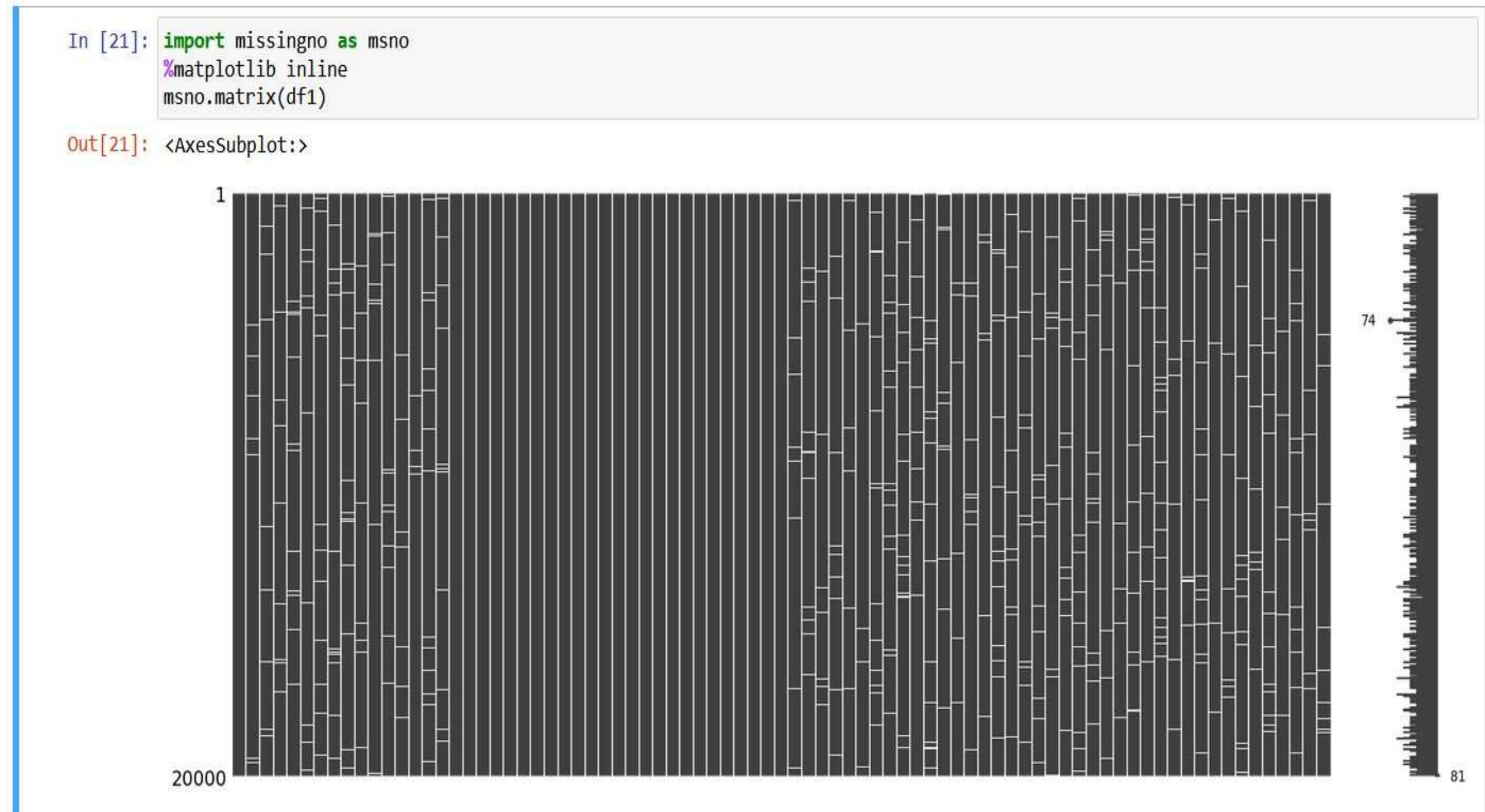


Figure 8: Missing Values

Visualization Missing Values

Problem Definition
Data Processing
Loading Data
Check for missing values in the dataset
Handle Missing Values
Build The Models
Conclusion



Figure 9: After Removing NAN vlaues





- [Problem Definition](#)
- [Data Processing](#)
- [Handle Missing Values](#)**
 - A-Deleting the Missing values
 - B-Replacing With Arbitrary Value (Zero)
 - C-Replacing with Previous Value – Forward Fill
 - D-Replacing with Next Value – Backward Fill
 - E-Interpolation Method
- [Build The Models](#)
- [Conclusion](#)

Handle Missing Values



How To Handle Missing Values?

Problem Definition

Data Processing

Handle Missing Values

A-Deleting the Missing values

B-Replacing With Arbitrary Value

(Zero)

C-Replacing with Previous Value –

Forward Fill

D-Replacing with Next Value –

Backward Fill

E-Interpolation Method

Build The Models

Conclusion

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:

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

Selecting Data

Problem Definition

Data Processing

Handle Missing Values

A-Deleting the Missing values

B-Replacing With Arbitrary Value

(Zero)

C-Replacing with Previous Value –

Forward Fill

D-Replacing with Next Value –

Backward Fill

E-Interpolation Method

Build The Models

Conclusion

- The dataset (Tabular Playground Series - June 2022) has huge data 1000000 Rows , and 81 columns.

```
In [6]: df.shape  
Out[6]: (1000000, 81)
```

Figure 10: Tabular Playground Series - June 2022 dataset

- We will select 5% of dataset.

```
In [7]: df1 = df.sample(frac = .05)  
df1.shape  
Out[7]: (50000, 81)
```

Figure 11: Selecting 5 % percent of data



TULIP

Team for Universal Learning and Intelligent Processing



Number of Missing Values in Selecting Data

- Problem Definition
- Data Processing
- Handle Missing Values
 - A-Deleting the Missing values
 - B-Replacing With Arbitrary Value (Zero)
 - C-Replacing with Previous Value – Forward Fill
 - D-Replacing with Next Value – Backward Fill
 - E-Interpolation Method
- Build The Models
- Conclusion

```
In [8]: df1.isnull().sum()|
Out[8]: row_id      0
        F_1_0      881
        F_1_1      873
        F_1_2      928
        F_1_3      932
        ...
        F_4_10     892
        F_4_11     839
        F_4_12     959
        F_4_13     873
        F_4_14     895
        Length: 81, dtype: int64
```

Figure 12: The Number of missing values for each row

```
In [9]: df1.isnull().sum().sum()|
Out[9]: 50064
```

Figure 13: The total of missing values in Selecting Data



A-Deleting the Missing values

- Problem Definition
- Data Processing
- Handle Missing Values
 - A-Deleting the Missing values**
 - B-Replacing With Arbitrary Value (Zero)
 - C-Replacing with Previous Value – Forward Fill
 - D-Replacing with Next Value – Backward Fill
 - E-Interpolation Method
- Build The Models
- Conclusion

Missing values can be handled by deleting rows or columns that have null values.

A- Handling Missing values by Drop Missing Values

```
In [10]: df2=df1.dropna()

In [11]: df2.isnull().sum()

Out[11]: row_id      0
         F_1_0      0
         F_1_1      0
         F_1_2      0
         F_1_3      0
         ..
         F_4_10     0
         F_4_11     0
         F_4_12     0
         F_4_13     0
         F_4_14     0
         Length: 81, dtype: int64

In [12]: df2.isnull().sum().sum()

Out[12]: 0
```

Figure 14: Drop Missing values



B-Replacing With Arbitrary Value (Zero)

- Problem Definition
- Data Processing
 - Handle Missing Values
 - A-Deleting the Missing values
 - B-Replacing With Arbitrary Value (Zero)**
 - C-Replacing with Previous Value – Forward Fill
 - D-Replacing with Next Value – Backward Fill
 - E-Interpolation Method
 - Build The Models
- Conclusion

We will replace the missing values with some arbitrary value using the following code.

B- Handling Missing values by Replacing with Zero

```
In [16]: df11=df1.fillna(0)

In [17]: df2.isnull().sum()

Out[17]: row_id      0
         F_1_0      0
         F_1_1      0
         F_1_2      0
         F_1_3      0
         ..
         F_4_10     0
         F_4_11     0
         F_4_12     0
         F_4_13     0
         F_4_14     0
         Length: 81, dtype: int64

In [18]: df11.isnull().sum().sum()

Out[18]: 0
```

Figure 15: Replace missing values with ‘0’.



C-Replacing with Previous Value – Forward Fill

- Problem Definition
- Data Processing
 - Handle Missing Values
 - A-Deleting the Missing values
 - B-Replacing With Arbitrary Value (Zero)
 - C-Replacing with Previous Value – Forward Fill**
 - D-Replacing with Next Value – Backward Fill
 - E-Interpolation Method
- Build The Models
- Conclusion

The missing value is imputed using the previous value and imputing the values with the previous value is more appropriate.

C-Handling Missing values by Replacing with Previous Value – Forward Fill

```
In [22]: df4=df1.fillna(method='pad')

In [23]: df4.isnull().sum()
Out[23]: row_id      0
         F_1_0      0
         F_1_1      0
         F_1_2      0
         F_1_3      0
         ..
         F_4_10     0
         F_4_11     0
         F_4_12     0
         F_4_13     0
         F_4_14     0
         Length: 81, dtype: int64

In [24]: df4.isnull().sum().sum()
Out[24]: 0
```

Figure 16: Replacing with previous value – Forward fill.



D-Replacing with Next Value – Backward Fill

- Problem Definition
- Data Processing
- Handle Missing Values
 - A-Deleting the Missing values
 - B-Replacing With Arbitrary Value (Zero)
 - C-Replacing with Previous Value – Forward Fill
 - D-Replacing with Next Value – Backward Fill**
 - E-Interpolation Method
- Build The Models
- Conclusion

In backward fill, the missing value is imputed using the next value.

D-Handling Missing values by Replacing with Next Value – Backward Fill

```
In [27]: df5=df1.fillna(method='bfill')

In [28]: df5.isnull().sum()

Out[28]: row_id      0
         F_1_0      0
         F_1_1      0
         F_1_2      0
         F_1_3      0
         ..
         F_4_10     0
         F_4_11     0
         F_4_12     0
         F_4_13     0
         F_4_14     0
         Length: 81, dtype: int64

In [29]: df5.isnull().sum().sum()

Out[29]: 0
```

Figure 17: We are replacing the missing values with next value.



E-Interpolation Method

- Problem Definition
- Data Processing
- Handle Missing Values
 - A-Deleting the Missing values
 - B-Replacing With Arbitrary Value (Zero)
 - C-Replacing with Previous Value – Forward Fill
 - D-Replacing with Next Value – Backward Fill
- E-Interpolation Method**
- Build The Models
- Conclusion

Missing values can also be imputed using interpolation. Pandas interpolate method can be used to replace the missing values with different interpolation methods like ‘polynomial’, ‘linear’, ‘quadratic’. Default method is ‘linear’.

E-Handling Missing values by Interpolate Method

```
In [19]: df3=df1.interpolate(method='linear')

In [20]: df3.isnull().sum()

Out[20]: row_id      0
F_1_0      0
F_1_1      0
F_1_2      0
F_1_3      0
..
F_4_10      0
F_4_11      0
F_4_12      0
F_4_13      0
F_4_14      0
Length: 81, dtype: int64

In [21]: df3.isnull().sum().sum()

Out[21]: 0
```

Figure 18: We are replacing the missing values with linear method.



[Problem Definition](#)

[Data Processing](#)

[Handle Missing Values](#)

[Build The Models](#)

1- Linear Regression Model

2- Decision Tree Model

3- Random Forest Model

[Conclusion](#)

Build The Models

1- Linear Regression Model

Problem Definition
Data Processing
Handle Missing Values
Build The Models
1- Linear Regression Model
2- Decision Tree Model
3- Random Forest Model
Conclusion

```
Linear Regression for Interpolate Method

In [66]: from sklearn.linear_model import LinearRegression
Lin=LinearRegression()
Lin.fit(X_train,y_train)

Out[66]: LinearRegression()

In [67]: print('The accuracy of LinearRegression: ',Lin.score(X_train,y_train))
The accuracy of LinearRegression: 0.6587766661021908
```

Figure 19: Interpolate Method

```
Linear Regression for Replacing with Previous Value – Forward Fill

In [96]: from sklearn.linear_model import LinearRegression
Lin=LinearRegression()
Lin.fit(X_train,y_train)

Out[96]: LinearRegression()

In [97]: print('The accuracy of LinearRegression: ',Lin.score(X_train,y_train))
The accuracy of LinearRegression: 0.6336116694335745
```

Figure 20: Replacing with Previous Value – Forward Fill

2- Decision Tree Model

Problem Definition
Data Processing
Handle Missing Values
Build The Models
1- Linear Regression Model
2- Decision Tree Model
3- Random Forest Model
Conclusion

Decision Tree for Interpolate Method

```
In [62]: from sklearn.tree import DecisionTreeClassifier
         tree=DecisionTreeClassifier(criterion='entropy',random_state=0)
         tree.fit(X_train,y_train)

Out[62]: DecisionTreeClassifier(criterion='entropy', random_state=0)

In [63]: print('The accuracy of Decision Tree: ',tree.score(X_train,y_train))
         The accuracy of Decision Tree:  1.0
```

Figure 21: Interpolate Method

Decision Tree for Replacing with Previous Value – Forward Fill

```
In [30]: from sklearn.tree import DecisionTreeClassifier
         tree=DecisionTreeClassifier(criterion='entropy',random_state=0)
         tree.fit(X_train,y_train)

Out[30]: DecisionTreeClassifier(criterion='entropy', random_state=0)

In [31]: print('The accuracy of Decision Tree: ',tree.score(X_train,y_train))
         The accuracy of Decision Tree:  1.0
```

Figure 22: Replacing with Previous Value – Forward Fill



3- Random Forest Model

Problem Definition

Data Processing

Handle Missing Values

Build The Models

1- Linear Regression Model

2- Decision Tree Model

3- Random Forest Model

Conclusion

Random Forest for Interpolate Method

```
In [64]: from sklearn.ensemble import RandomForestClassifier
forst=RandomForestClassifier(n_estimators=10,criterion='entropy',random_state=0)
forst.fit(X_train,y_train)

Out[64]: RandomForestClassifier(criterion='entropy', n_estimators=10, random_state=0)

In [65]: print('The accuracy of Random Forest : ',forst.score(X_train,y_train))
The accuracy of Random Forest :  0.9994375
```

Figure 23: Interpolate Method

Random Forest for Replacing with Previous Value – Forward Fill

```
In [34]: from sklearn.ensemble import RandomForestClassifier
forst=RandomForestClassifier(n_estimators=10,criterion='entropy',random_state=0)
forst.fit(X_train,y_train)

Out[34]: RandomForestClassifier(criterion='entropy', n_estimators=10, random_state=0)

In [35]: print('The accuracy of Random Forest : ',forst.score(X_train,y_train))
The accuracy of Random Forest :  0.9993125
```

Figure 24: Replacing with Previous Value – Forward Fill



TULIP

Team for Universal Learning and Intelligent Processing



- [Problem Definition](#)
- [Data Processing](#)
- [Handle Missing Values](#)
- [Build The Models](#)
- [Conclusion](#)**

Conclusion



Evaluating models

- [Problem Definition](#)
- [Data Processing](#)
- [Handle Missing Values](#)
- [Build The Models](#)
- [Conclusion](#)

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.

Table 1: Accuracy

Type of 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

The presence of missing values in a dataset can affect the performance of a classifier constructed using that dataset as a training sample. Several methods have been proposed to treat missing data.



The Results

- [Problem Definition](#)
- [Data Processing](#)
- [Handle Missing Values](#)
- [Build The Models](#)
- [Conclusion](#)

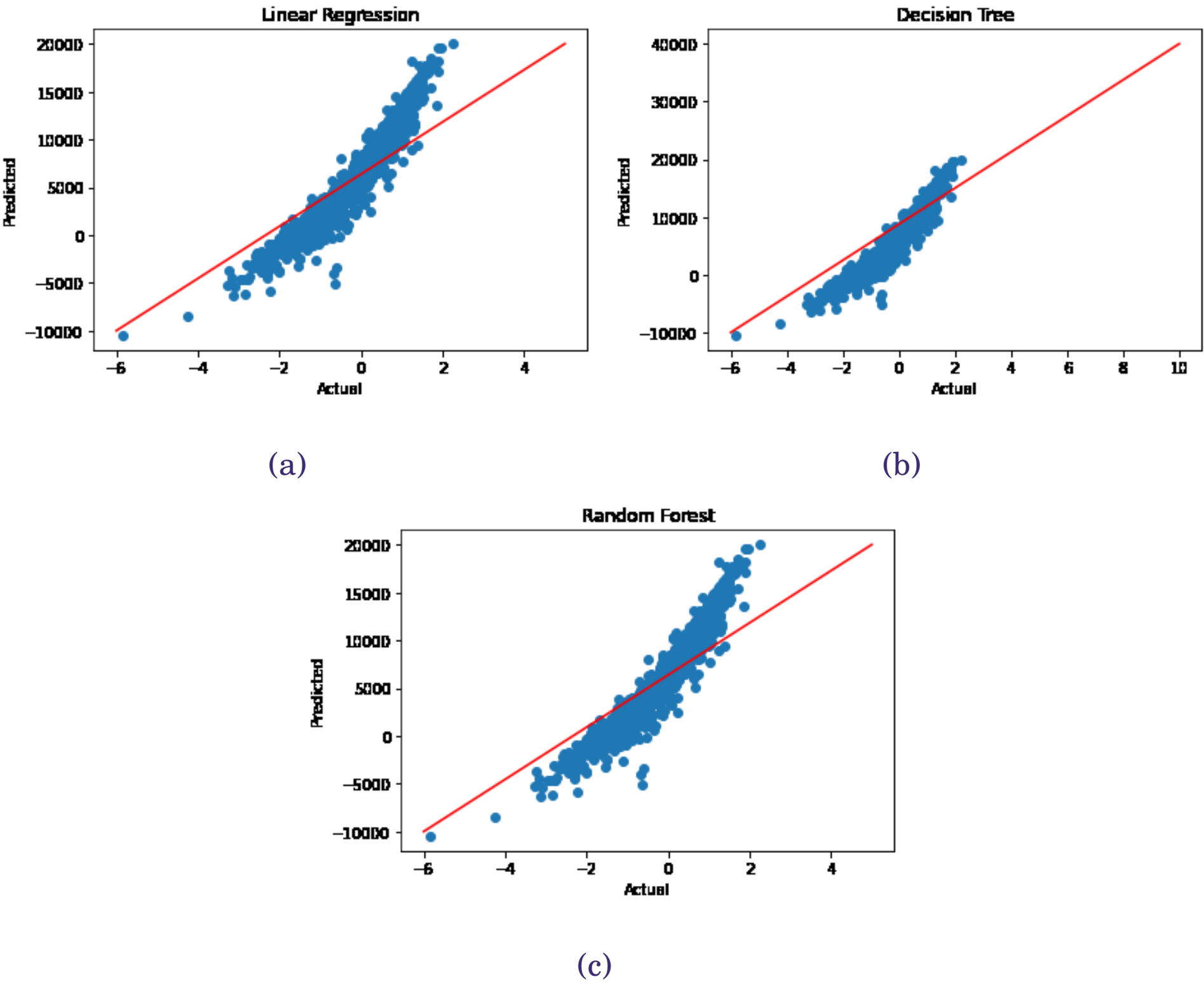


Figure 25: (Predicted for LR,DT,and RF)

Questions?

Problem Definition

Data Processing

Handle Missing Values

Build The Models

Conclusion



Thank you- Stay Safe



TULIP

Team for Universal Learning and Intelligent Processing



Contact Information

Mousa Tayseer Jafar
Deakin University, Australia



MUOSAJAFAR@GMAIL.COM



TEAM FOR UNIVERSAL LEARNING AND INTELLIGENT PROCESSING

