

Handling Missing Values

```
In [ ]: import pandas as pd
import numpy as np
import matplotlib.pyplot as plt
```

```
In [ ]: df = pd.read_csv('./data/Churn_Modelling.csv')
```

```
In [ ]: df.info()
```

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 10000 entries, 0 to 9999
Data columns (total 14 columns):
#   Column                Non-Null Count  Dtype  
---  -
0   RowNumber              10000 non-null  int64  
1   CustomerId             10000 non-null  int64  
2   Surname                10000 non-null  object  
3   CreditScore             10000 non-null  int64  
4   Geography              10000 non-null  object  
5   Gender                 9946 non-null   object  
6   Age                    9700 non-null   float64 
7   Tenure                  10000 non-null  int64  
8   Balance                 10000 non-null  float64 
9   NumOfProducts          10000 non-null  int64  
10  HasCrCard               10000 non-null  int64  
11  IsActiveMember          10000 non-null  int64  
12  EstimatedSalary         10000 non-null  float64 
13  Exited                  10000 non-null  int64  
dtypes: float64(3), int64(8), object(3)
memory usage: 1.1+ MB
```

1. Gender has 54 missing values
2. Age has 300 missing values

The second way of finding whether we have null values in the data is by using the `isnull()` function.

```
In [ ]: print(df.isnull().sum())
```

```

RowNumber      0
CustomerId     0
Surname        0
CreditScore    0
Geography      0
Gender         54
Age           300
Tenure         0
Balance        0
NumOfProducts 0
HasCrCard      0
IsActiveMember 0
EstimatedSalary 0
Exited         0
dtype: int64

```

Handling Missing Values

1. Deleting the columns with missing data

```
In [ ]: updated_df = df.dropna(axis=1)
```

```
In [ ]: updated_df.info()
```

```

<class 'pandas.core.frame.DataFrame'>
RangeIndex: 10000 entries, 0 to 9999
Data columns (total 12 columns):
#   Column                Non-Null Count  Dtype
---  -
0   RowNumber             10000 non-null  int64
1   CustomerId            10000 non-null  int64
2   Surname               10000 non-null  object
3   CreditScore           10000 non-null  int64
4   Geography             10000 non-null  object
5   Tenure               10000 non-null  int64
6   Balance               10000 non-null  float64
7   NumOfProducts        10000 non-null  int64
8   HasCrCard            10000 non-null  int64
9   IsActiveMember       10000 non-null  int64
10  EstimatedSalary       10000 non-null  float64
11  Exited               10000 non-null  int64
dtypes: float64(2), int64(8), object(2)
memory usage: 937.6+ KB

```

The problem with this method is that we may lose valuable information on that feature, as we have deleted it completely due to some null values.

Should only be used if there are too many null values.

2. Deleting the rows with missing data

```
In [ ]: updated_df = df.dropna(axis=0)
```

```
In [ ]: updated_df.info()
```

```
<class 'pandas.core.frame.DataFrame'>
Int64Index: 9655 entries, 0 to 9999
Data columns (total 14 columns):
#   Column                Non-Null Count  Dtype
---  -
0   RowNumber              9655 non-null   int64
1   CustomerId             9655 non-null   int64
2   Surname                9655 non-null   object
3   CreditScore             9655 non-null   int64
4   Geography              9655 non-null   object
5   Gender                 9655 non-null   object
6   Age                    9655 non-null   float64
7   Tenure                  9655 non-null   int64
8   Balance                 9655 non-null   float64
9   NumOfProducts          9655 non-null   int64
10  HasCrCard               9655 non-null   int64
11  IsActiveMember          9655 non-null   int64
12  EstimatedSalary         9655 non-null   float64
13  Exited                  9655 non-null   int64
dtypes: float64(3), int64(8), object(3)
memory usage: 1.1+ MB
```

In this case, there are possibilities of getting better accuracy than before. This might be because the columns contains more valuable information than we expected.

3. Filling the Missing Values – Imputation

In this case, we will be filling the missing values with a certain number.

The possible ways to do this are:

- Filling the missing data with the mean or median value if it's a numerical variable.
- Filling the missing data with mode if it's a categorical value.
- Filling the numerical value with 0 or -999, or some other number that will not occur in the data. This can be done so that the machine can recognize that the data is not real or is different.
- Filling the categorical value with a new type for the missing values.

```
In [ ]: df['Age'].mean()
```

Out[]: 38.92432989690722

```
In [ ]: df['Age'].median()
```

Out[]: 37.0

```
In [ ]: #fillna: fills the null records
        #dropna: drops the null records

        updated_df = df
        updated_df['Age'] = updated_df['Age'].fillna(df['Age'].mean())
        updated_df.info()
```

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 10000 entries, 0 to 9999
Data columns (total 14 columns):
#   Column                Non-Null Count  Dtype
---  -
0   RowNumber              10000 non-null  int64
1   CustomerId             10000 non-null  int64
2   Surname                10000 non-null  object
3   CreditScore             10000 non-null  int64
4   Geography              10000 non-null  object
5   Gender                 9946 non-null   object
6   Age                    10000 non-null  float64
7   Tenure                 10000 non-null  int64
8   Balance                10000 non-null  float64
9   NumOfProducts          10000 non-null  int64
10  HasCrCard               10000 non-null  int64
11  IsActiveMember          10000 non-null  int64
12  EstimatedSalary         10000 non-null  float64
13  Exited                  10000 non-null  int64
dtypes: float64(3), int64(8), object(3)
memory usage: 1.1+ MB
```

```
In [ ]: updated_df1 = df
        updated_df1['Age'] = updated_df1['Age'].fillna(df['Age'].median())
        updated_df1.info()
```

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 10000 entries, 0 to 9999
Data columns (total 14 columns):
#   Column                Non-Null Count  Dtype
---  -
0   RowNumber             10000 non-null  int64
1   CustomerId            10000 non-null  int64
2   Surname               10000 non-null  object
3   CreditScore           10000 non-null  int64
4   Geography             10000 non-null  object
5   Gender                9946 non-null   object
6   Age                   10000 non-null  float64
7   Tenure                10000 non-null  int64
8   Balance               10000 non-null  float64
9   NumOfProducts         10000 non-null  int64
10  HasCrCard             10000 non-null  int64
11  IsActiveMember        10000 non-null  int64
12  EstimatedSalary       10000 non-null  float64
13  Exited                10000 non-null  int64
dtypes: float64(3), int64(8), object(3)
memory usage: 1.1+ MB
```

4. Forward & Backward Filling – Imputation

```
In [ ]: df1 = df
```

```
In [ ]: df1['Age'] = df1['Age'].bfill()
```

```
In [ ]: df1.info()
```

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 10000 entries, 0 to 9999
Data columns (total 14 columns):
#   Column                Non-Null Count  Dtype
---  -
0   RowNumber             10000 non-null  int64
1   CustomerId            10000 non-null  int64
2   Surname               10000 non-null  object
3   CreditScore           10000 non-null  int64
4   Geography             10000 non-null  object
5   Gender                9946 non-null   object
6   Age                   10000 non-null  float64
7   Tenure                10000 non-null  int64
8   Balance               10000 non-null  float64
9   NumOfProducts         10000 non-null  int64
10  HasCrCard             10000 non-null  int64
11  IsActiveMember        10000 non-null  int64
12  EstimatedSalary       10000 non-null  float64
13  Exited                10000 non-null  int64
dtypes: float64(3), int64(8), object(3)
memory usage: 1.1+ MB
```