# Missing value treatment

In [3]: 
```
Exp : 6
```

In [5]: 
```python
#Name: Leena Rajeshwar Kale
#Roll No.:71
#Sec: C
#Subject:ET - 1
```

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

In [2]: 
```python
import matplotlib.pyplot as plt
import numpy as np
import seaborn as sns
from sklearn.model_selection import train_test_split
import warnings
warnings.filterwarnings('ignore')
```

In [3]: 
```python
import os
```

In [4]: 
```python
os.getcwd()
```

Out[4]: 
```
'C:\\Users\\dishi\\Downloads'
```

In [5]: 
```python
os.chdir("C:\\Users\\dishi\\Downloads")
```

In [22]: 
```python
df=pd.read_csv("framingham.csv")
```

In [23]: 
```python
df
```

Out[23]: 

| | male | age | education | currentSmoker | cigsPerDay | BPMeds | prevalentStroke | prevalentHyp |
|---|---|---|---|---|---|---|---|---|
| 0 | 1 | 39 | 4.0 | 0 | 0.0 | 0.0 | 0 | 0 |
| 1 | 0 | 46 | 2.0 | 0 | 0.0 | 0.0 | 0 | 0 |
| 2 | 1 | 48 | 1.0 | 1 | 20.0 | 0.0 | 0 | 0 |
| 3 | 0 | 61 | 3.0 | 1 | 30.0 | 0.0 | 0 | 1 |
| 4 | 0 | 46 | 3.0 | 1 | 23.0 | 0.0 | 0 | 0 |
| ... | ... | ... | ... | ... | ... | ... | ... | ... |
| 4233 | 1 | 50 | 1.0 | 1 | 1.0 | 0.0 | 0 | 1 |
| 4234 | 1 | 51 | 3.0 | 1 | 43.0 | 0.0 | 0 | 0 |
| 4235 | 0 | 48 | 2.0 | 1 | 20.0 | NaN | 0 | 0 |
| 4236 | 0 | 44 | 1.0 | 1 | 15.0 | 0.0 | 0 | 0 |
| 4237 | 0 | 52 | 2.0 | 0 | 0.0 | 0.0 | 0 | 0 |

4238 rows × 16 columns

In [24]: 
```python
df.info()
```

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 4238 entries, 0 to 4237
Data columns (total 16 columns):
 #   Column           Non-Null Count  Dtype
---  ------           --------------  -----
 0   male             4238 non-null   int64
 1   age              4238 non-null   int64
 2   education        4133 non-null   float64
 3   currentSmoker    4238 non-null   int64
 4   cigsPerDay       4209 non-null   float64
 5   BPMeds           4185 non-null   float64
 6   prevalentStroke  4238 non-null   int64
 7   prevalentHyp     4238 non-null   int64
 8   diabetes         4238 non-null   int64
 9   totChol          4188 non-null   float64
 10  sysBP            4238 non-null   float64
 11  diaBP            4238 non-null   float64
 12  BMI              4219 non-null   float64
 13  heartRate        4237 non-null   float64
 14  glucose          3850 non-null   float64
 15  TenYearCHD       4238 non-null   int64
dtypes: float64(9), int64(7)
memory usage: 529.9 KB
```

In [25]: `df.isna()`

Out[25]:

|      | male  | age   | education | currentSmoker | cigsPerDay | BPMeds | prevalentStroke | prevalentHyp |
|------|-------|-------|-----------|---------------|------------|--------|-----------------|--------------|
| 0    | False | False | False     | False         | False      | False  | False           | False        |
| 1    | False | False | False     | False         | False      | False  | False           | False        |
| 2    | False | False | False     | False         | False      | False  | False           | False        |
| 3    | False | False | False     | False         | False      | False  | False           | False        |
| 4    | False | False | False     | False         | False      | False  | False           | False        |
| ...  | ...   | ...   | ...       | ...           | ...        | ...    | ...             | ...          |
| 4233 | False | False | False     | False         | False      | False  | False           | False        |
| 4234 | False | False | False     | False         | False      | False  | False           | False        |
| 4235 | False | False | False     | False         | False      | True   | False           | False        |
| 4236 | False | False | False     | False         | False      | False  | False           | False        |
| 4237 | False | False | False     | False         | False      | False  | False           | False        |

4238 rows × 16 columns

In [26]: `df.isnull()`

Out[26]:

| | male | age | education | currentSmoker | cigsPerDay | BPMeds | prevalentStroke | prevalentHyp |
|---|---|---|---|---|---|---|---|---|
| 0 | False | False | False | False | False | False | False | False |
| 1 | False | False | False | False | False | False | False | False |
| 2 | False | False | False | False | False | False | False | False |
| 3 | False | False | False | False | False | False | False | False |
| 4 | False | False | False | False | False | False | False | False |
| ... | ... | ... | ... | ... | ... | ... | ... | ... |
| 4233 | False | False | False | False | False | False | False | False |
| 4234 | False | False | False | False | False | False | False | False |
| 4235 | False | False | False | False | False | True | False | False |
| 4236 | False | False | False | False | False | False | False | False |
| 4237 | False | False | False | False | False | False | False | False |

4238 rows × 16 columns

```python
In [27]: df['glucose'].fillna(value = df['glucose'].mean(),inplace=True)
```

```python
In [28]: print(df['BPMeds'].fillna(value = df['BPMeds'].mean(),inplace=True))
```
```
None
```

```python
In [29]: print(df['cigsPerDay'].fillna(value = df['cigsPerDay'].mean(),inplace=True))
```
```
None
```

```python
In [30]: df.isna().sum()
```

Out[30]:
```
male                 0
age                  0
education          105
currentSmoker        0
cigsPerDay           0
BPMeds               0
prevalentStroke      0
prevalentHyp         0
diabetes             0
totChol             50
sysBP                0
diaBP                0
BMI                 19
heartRate            1
glucose              0
TenYearCHD           0
dtype: int64
```

```python
In [31]: #Splitting the dependent and independent variables.
         x = df.drop("TenYearCHD",axis=1)
         y = df['TenYearCHD']
```

```python
In [32]: x
```

| | male | age | education | currentSmoker | cigsPerDay | BPMeds | prevalentStroke | prevalentHyp |
|---|---|---|---|---|---|---|---|---|
| 0 | 1 | 39 | 4.0 | 0 | 0.0 | 0.00000 | 0 | 0 |
| 1 | 0 | 46 | 2.0 | 0 | 0.0 | 0.00000 | 0 | 0 |
| 2 | 1 | 48 | 1.0 | 1 | 20.0 | 0.00000 | 0 | 0 |
| 3 | 0 | 61 | 3.0 | 1 | 30.0 | 0.00000 | 0 | 1 |
| 4 | 0 | 46 | 3.0 | 1 | 23.0 | 0.00000 | 0 | 0 |
| ... | ... | ... | ... | ... | ... | ... | ... | ... |
| 4233 | 1 | 50 | 1.0 | 1 | 1.0 | 0.00000 | 0 | 1 |
| 4234 | 1 | 51 | 3.0 | 1 | 43.0 | 0.00000 | 0 | 0 |
| 4235 | 0 | 48 | 2.0 | 1 | 20.0 | 0.02963 | 0 | 0 |
| 4236 | 0 | 44 | 1.0 | 1 | 15.0 | 0.00000 | 0 | 0 |
| 4237 | 0 | 52 | 2.0 | 0 | 0.0 | 0.00000 | 0 | 0 |

4238 rows × 15 columns