

Practical No. 01

Aim : To find Unique and Duplicate value count in given data set

In [54]:

```
#Name : Taufiq Rafik Nagori  
#Roll no. : 77 (BDA-B77)  
#Section : B  
#Subject : PE-II
```

In [56]:

```
#Performing EDA {Exploratory Data Analysis} on given dataset with additional operations:  
#finding unique value using unique() and finding correlation matrix.
```

In [8]:

```
import os
```

In [10]:

```
import pandas as pd
```

In [12]:

```
os.getcwd()
```

Out[12]:

```
'C:\\Users\\USER'
```

In [14]:

```
os.chdir("C:\\Users\\USER\\Desktop")
```

In [16]:

```
data=pd.read_csv("diabetes.csv")
```

In [18]:

```
data.head()
```

Out[18]:

	Pregnancies	Glucose	BloodPressure	SkinThickness	Insulin	BMI	DiabetesPedigreeFunction	Age
0	6	148	72	35	0	33.6	0.627	50
1	1	85	66	29	0	26.6	0.351	31
2	8	183	64	0	0	23.3	0.672	32
3	1	89	66	23	94	28.1	0.167	21
4	0	137	40	35	168	43.1	2.288	33

In [20]:

```
data.tail()
```

Out[20]:

	Pregnancies	Glucose	BloodPressure	SkinThickness	Insulin	BMI	DiabetesPedigreeFunction	Ag
763	10	101	76	48	180	32.9	0.171	6
764	2	122	70	27	0	36.8	0.340	2

	Pregnancies	Glucose	BloodPressure	SkinThickness	Insulin	BMI	DiabetesPedigreeFunction	Ag
765	5	121	72	23	112	26.2	0.245	3
766	1	126	60	0	0	30.1	0.349	4
767	1	93	70	31	0	30.4	0.315	2

In [22]:

```
data.describe()
```

Out[22]:

	Pregnancies	Glucose	BloodPressure	SkinThickness	Insulin	BMI	DiabetesPedigr
count	768.000000	768.000000	768.000000	768.000000	768.000000	768.000000	
mean	3.845052	120.894531	69.105469	20.536458	79.799479	31.992578	
std	3.369578	31.972618	19.355807	15.952218	115.244002	7.884160	
min	0.000000	0.000000	0.000000	0.000000	0.000000	0.000000	
25%	1.000000	99.000000	62.000000	0.000000	0.000000	27.300000	
50%	3.000000	117.000000	72.000000	23.000000	30.500000	32.000000	
75%	6.000000	140.250000	80.000000	32.000000	127.250000	36.600000	
max	17.000000	199.000000	122.000000	99.000000	846.000000	67.100000	

In [24]:

```
data.info()
```

```
<class 'pandas.core.frame.DataFrame'>
```

```
RangeIndex: 768 entries, 0 to 767
```

```
Data columns (total 9 columns):
```

#	Column	Non-Null Count	Dtype
0	Pregnancies	768 non-null	int64
1	Glucose	768 non-null	int64
2	BloodPressure	768 non-null	int64
3	SkinThickness	768 non-null	int64
4	Insulin	768 non-null	int64
5	BMI	768 non-null	float64
6	DiabetesPedigreeFunction	768 non-null	float64
7	Age	768 non-null	int64
8	Outcome	768 non-null	int64

```
dtypes: float64(2), int64(7)
```

```
memory usage: 54.1 KB
```

In [32]:

```
data.shape
```

Out[32]:

```
(768, 9)
```

In [34]:

```
data.size
```

Out[34]:

```
6912
```

In [36]:

```
data.ndim
```

```
Out[36]:
```

```
2
```

```
In [38]:
```

```
data.isna()
```

```
Out[38]:
```

	Pregnancies	Glucose	BloodPressure	SkinThickness	Insulin	BMI	DiabetesPedigreeFunction	Age	Outcome
0	False	False	False	False	False	False	False	False	False
1	False	False	False	False	False	False	False	False	False
2	False	False	False	False	False	False	False	False	False
3	False	False	False	False	False	False	False	False	False
4	False	False	False	False	False	False	False	False	False
...
763	False	False	False	False	False	False	False	False	False
764	False	False	False	False	False	False	False	False	False
765	False	False	False	False	False	False	False	False	False
766	False	False	False	False	False	False	False	False	False
767	False	False	False	False	False	False	False	False	False

768 rows × 9 columns

```
In [40]:
```

```
data.isna().any()
```

```
Out[40]:
```

```
Pregnancies      False
Glucose           False
BloodPressure     False
SkinThickness     False
Insulin           False
BMI               False
DiabetesPedigreeFunction  False
Age               False
Outcome           False
dtype: bool
```

```
In [42]:
```

```
data.isna().sum()
```

```
Out[42]:
```

```
Pregnancies      0
Glucose           0
BloodPressure     0
SkinThickness     0
Insulin           0
BMI               0
DiabetesPedigreeFunction  0
Age               0
```

```
Outcome                                0
dtype: int64
```

```
In [44]:
```

```
unique_values = data['Age'].unique()
```

```
In [46]:
```

```
print(unique_values)
```

```
[50 31 32 21 33 30 26 29 53 54 34 57 59 51 27 41 43 22 38 60 28 45 35 46
 56 37 48 40 25 24 58 42 44 39 36 23 61 69 62 55 65 47 52 66 49 63 67 72
 81 64 70 68]
```

```
In [22]:
```

```
unique_count = len(data.stack().unique())
print(unique_count)
```

```
1005
```

```
In [26]:
```

```
data['Age'].duplicated().sum()
```

```
Out[26]:
```

```
716
```

```
In [28]:
```

```
data['Age'].duplicated()
```

```
Out[28]:
```

```
0      False
1      False
2      False
3      False
4      False
```

```
...
```

```
763     True
764     True
765     True
766     True
767     True
```

```
Name: Age, Length: 768, dtype: bool
```

```
In [58]:
```

```
data['Glucose'].duplicated()
```

```
Out[58]:
```

```
0      False
1      False
2      False
3      False
4      False
```

```
...
```

```
763     True
764     True
765     True
766     True
767     True
```

```
Name: Glucose, Length: 768, dtype: bool
```

```
In [ ]:
```