

## Practical No. 10

Aim: To perform and Data analysis with Co-relation Matrix

In [3]:

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

In [5]:

```
import os  
import pandas as pd
```

In [7]:

```
os.getcwd()
```

Out[7]:

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

In [9]:

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

In [11]:

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

In [13]:

```
data.head()
```

Out[13]:

|   | 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 [15]:

```
data.tail()
```

Out[15]:

|     | 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  |
| 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 [17]:

```
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 [19]:

```
data.info
```

Out[19]:

```
<bound method DataFrame.info of
Insulin    BMI    \
0          6    148    72    35    0  33.6
1          1    85    66    29    0  26.6
2          8   183    64    0    0  23.3
3          1    89    66    23   94  28.1
4          0   137    40    35  168  43.1
..          ...    ...    ...    ...    ...
763        10   101    76    48  180  32.9
764         2   122    70    27    0  36.8
765         5   121    72    23  112  26.2
766         1   126    60    0    0  30.1
767         1    93    70    31    0  30.4

    DiabetesPedigreeFunction  Age  Outcome
0                0.627    50         1
1                0.351    31         0
2                0.672    32         1
3                0.167    21         0
4                2.288    33         1
..                ...    ...        ...
763               0.171    63         0
764               0.340    27         0
765               0.245    30         0
766               0.349    47         1
767               0.315    23         0
```

```
[768 rows x 9 columns]>
```

In [21]:

```
data.describe()
```

Out[21]:

|       | Pregnancies | Glucose    | BloodPressure | SkinThickness | Insulin    | BMI        | DiabetesPedigr |
|-------|-------------|------------|---------------|---------------|------------|------------|----------------|
| count | 768.000000  | 768.000000 | 768.000000    | 768.000000    | 768.000000 | 768.000000 |                |

|      | Pregnancies | Glucose    | BloodPressure | SkinThickness | Insulin    | BMI       | DiabetesPedigreeFunction | Age | Outcome |
|------|-------------|------------|---------------|---------------|------------|-----------|--------------------------|-----|---------|
| 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 [23]:

```
data.shape
```

Out[23]:

```
(768, 9)
```

In [25]:

```
data.size
```

Out[25]:

```
6912
```

In [27]:

```
data.ndim
```

Out[27]:

```
2
```

In [29]:

```
data.columns
```

Out[29]:

```
Index(['Pregnancies', 'Glucose', 'BloodPressure', 'SkinThickness', 'Insulin',
       'BMI', 'DiabetesPedigreeFunction', 'Age', 'Outcome'],
      dtype='object')
```

Data pre-processing, data-cleaning, missing value treatment

In [32]:

```
data.isna()
```

Out[32]:

|     | 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   |

|     | Pregnancies | Glucose | BloodPressure | SkinThickness | Insulin | BMI   | DiabetesPedigreeFunction | Age   | Outcome |
|-----|-------------|---------|---------------|---------------|---------|-------|--------------------------|-------|---------|
| 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 [34]:

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

Out[34]:

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

In [36]:

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

Out[36]:

```
Pregnancies      0
Glucose           0
BloodPressure     0
SkinThickness     0
Insulin           0
BMI               0
DiabetesPedigreeFunction  0
Age               0
Outcome           0
dtype: int64
```

Co-relation Matrix

In [39]:

```
import seaborn as sns
import matplotlib.pyplot as plt
```

In [41]:

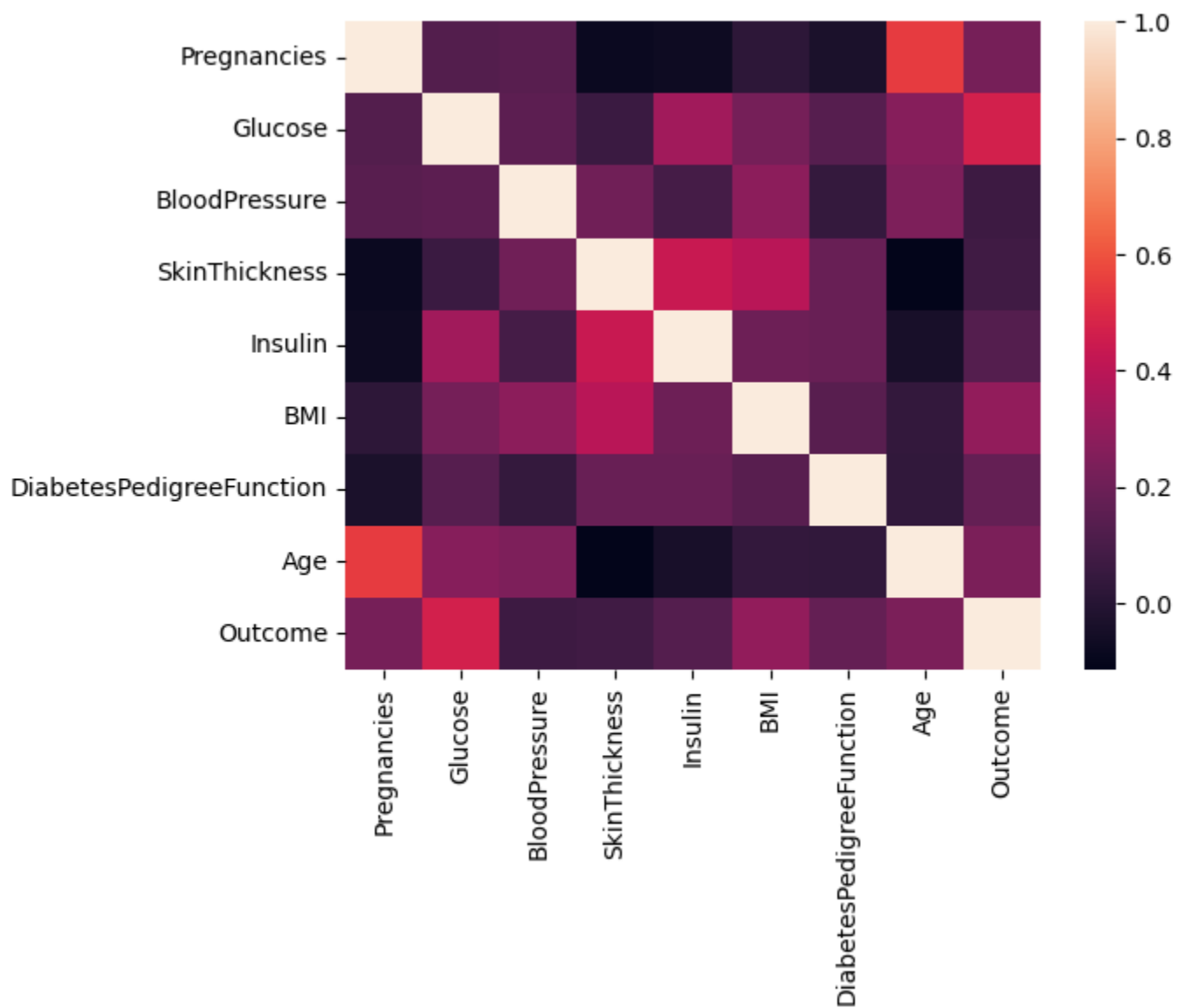
```
corr = data.corr()
```

In [43]:

```
sns.heatmap(data.corr())
```

Out[43]:

<Axes: >

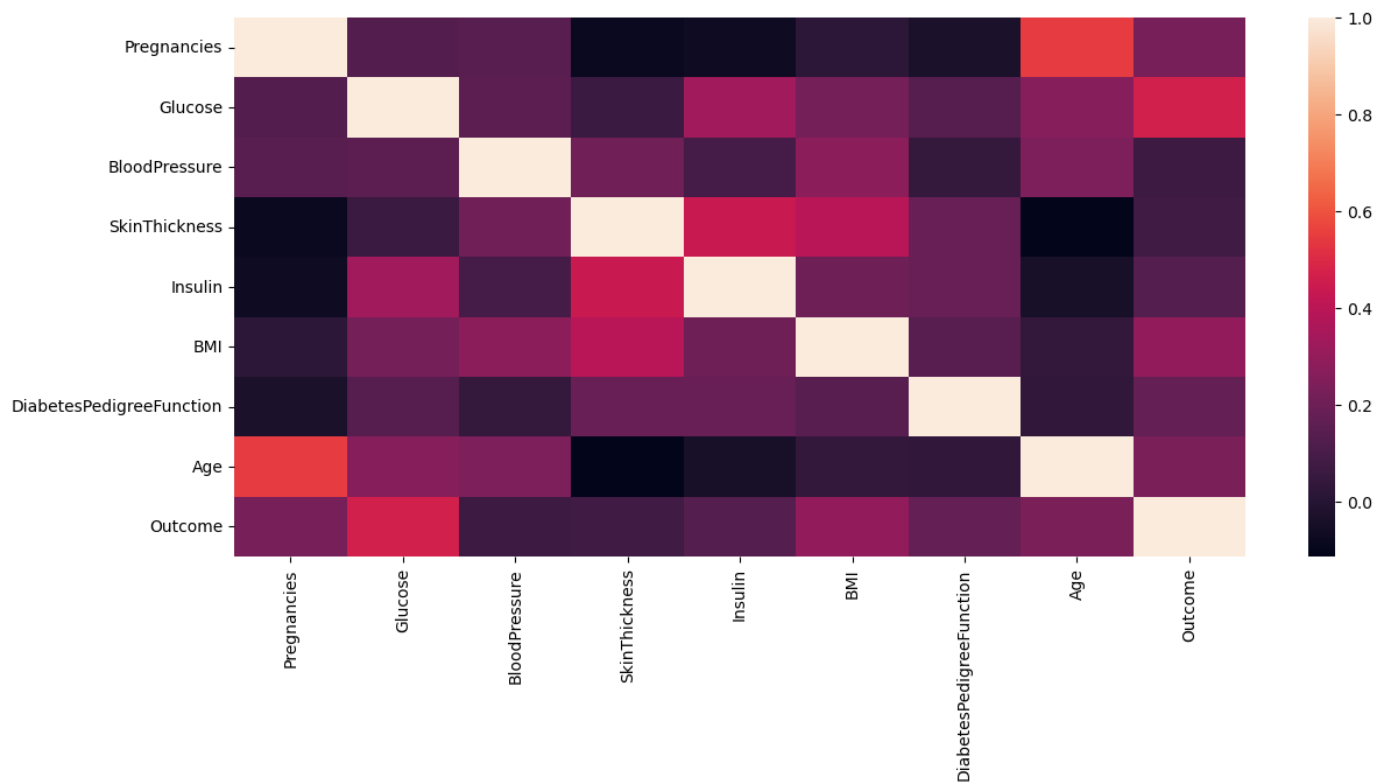


In [45]:

```
plt.figure(figsize=(14,6))  
sns.heatmap(data.corr())
```

Out[45]:

<Axes: >

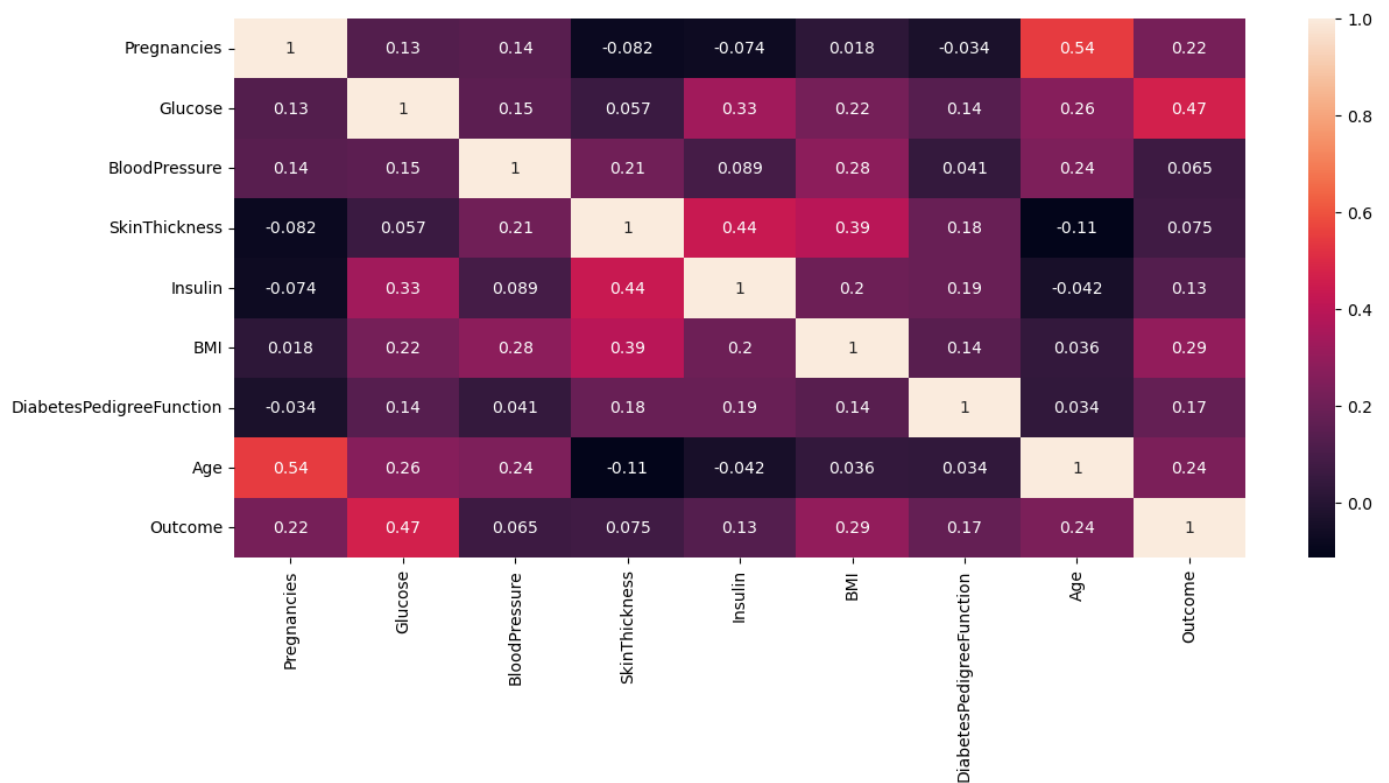


In [47]:

```
plt.figure(figsize=(14,6))
sns.heatmap(data.corr(),annot=True)
```

Out[47]:

<Axes: >



In [49]:

```
sns.distplot(data,bins=20)
plt.show()
```

C:\Users\USER\AppData\Local\Temp\ipykernel\_9668\1706651633.py:1: UserWarning:

`distplot` is a deprecated function and will be removed in seaborn v0.14.0.

Please adapt your code to use either `displot` (a figure-level function with similar flexibility) or `histplot` (an axes-level function for histograms).

For a guide to updating your code to use the new functions, please see <https://gist.github.com/mwaskom/de44147ed2974457ad6372750bbe5751>

```
sns.distplot(data,bins=20)
```

