```
In [4]:    import numpy as np
           import pandas as pd
           from matplotlib import pyplot as plt
           import seaborn as sns
```

```
In [5]:    diabetes = pd.read_csv('C:/Users/Administrator/Downloads/diabetes.csv')
```

```
In [6]:    diabetes
```

Out[6]:

| | 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 |
| ... | ... | ... | ... | ... | ... | ... | ... | ... |
| 763 | 10 | 101 | 76 | 48 | 180 | 32.9 | 0.171 | 63 |
| 764 | 2 | 122 | 70 | 27 | 0 | 36.8 | 0.340 | 27 |
| 765 | 5 | 121 | 72 | 23 | 112 | 26.2 | 0.245 | 30 |
| 766 | 1 | 126 | 60 | 0 | 0 | 30.1 | 0.349 | 47 |
| 767 | 1 | 93 | 70 | 31 | 0 | 30.4 | 0.315 | 23 |

768 rows × 9 columns

```
In [7]:    diabetes.head()
```

Out[7]:

| | Pregnancies | Glucose | BloodPressure | SkinThickness | Insulin | BMI | DiabetesPedigreeFunction | Age | O |
|---|---|---|---|---|---|---|---|---|---|
| 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 [8]:    diabetes.tail()
```

Out[8]:

| | Pregnancies | Glucose | BloodPressure | SkinThickness | Insulin | BMI | DiabetesPedigreeFunction | Age |
|---|---|---|---|---|---|---|---|---|

| | Pregnancies | Glucose | BloodPressure | SkinThickness | Insulin | BMI | DiabetesPedigreeFunction | Age |
|---|---|---|---|---|---|---|---|---|
| **763** | 10 | 101 | 76 | 48 | 180 | 32.9 | 0.171 | 63 |
| **764** | 2 | 122 | 70 | 27 | 0 | 36.8 | 0.340 | 27 |
| **765** | 5 | 121 | 72 | 23 | 112 | 26.2 | 0.245 | 30 |
| **766** | 1 | 126 | 60 | 0 | 0 | 30.1 | 0.349 | 47 |
| **767** | 1 | 93 | 70 | 31 | 0 | 30.4 | 0.315 | 23 |

In [9]:
```
diabetes.shape
```

Out[9]: (768, 9)

In [10]:
```
diabetes.info
```

Out[10]:
```
<bound method DataFrame.info of      Pregnancies  Glucose  BloodPressure  SkinThickness  \
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 [11]:
```
diabetes.describe()
```

Out[11]:

| | Pregnancies | Glucose | BloodPressure | SkinThickness | Insulin | BMI | DiabetesPedigreeF |
|---|---|---|---|---|---|---|---|
| **count** | 768.000000 | 768.000000 | 768.000000 | 768.000000 | 768.000000 | 768.000000 | 76 |
| **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 | |

|       | Pregnancies | Glucose    | BloodPressure | SkinThickness | Insulin     | BMI       | DiabetesPedigree |
|-------|-------------|------------|---------------|---------------|-------------|-----------|------------------|
| 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 [12]:
```python
diabetes.columns
```

Out[12]:
```
Index(['Pregnancies', 'Glucose', 'BloodPressure', 'SkinThickness', 'Insulin',
       'BMI', 'DiabetesPedigreeFunction', 'Age', 'Outcome'],
      dtype='object')
```

In [13]:
```python
diabetes.groupby('Outcome').mean()
```

Out[13]:

|         | Pregnancies | Glucose    | BloodPressure | SkinThickness | Insulin    | BMI       | DiabetesPedigr |
|---------|-------------|------------|---------------|---------------|------------|-----------|----------------|
| Outcome |             |            |               |               |            |           |                |
| 0       | 3.298000    | 109.980000 | 68.184000     | 19.664000     | 68.792000  | 30.304200 |                |
| 1       | 4.865672    | 141.257463 | 70.824627     | 22.164179     | 100.335821 | 35.142537 |                |

In [14]:
```python
#Check if any null value is present
diabetes.isnull().values.any()
```

Out[14]:
```
False
```

In [15]:
```python
diabetes.corr()
```

Out[15]:

|                          | Pregnancies | Glucose   | BloodPressure | SkinThickness | Insulin   | BMI       | Dia |
|--------------------------|-------------|-----------|---------------|---------------|-----------|-----------|-----|
| Pregnancies              | 1.000000    | 0.129459  | 0.141282      | -0.081672     | -0.073535 | 0.017683  |     |
| Glucose                  | 0.129459    | 1.000000  | 0.152590      | 0.057328      | 0.331357  | 0.221071  |     |
| BloodPressure            | 0.141282    | 0.152590  | 1.000000      | 0.207371      | 0.088933  | 0.281805  |     |
| SkinThickness            | -0.081672   | 0.057328  | 0.207371      | 1.000000      | 0.436783  | 0.392573  |     |
| Insulin                  | -0.073535   | 0.331357  | 0.088933      | 0.436783      | 1.000000  | 0.197859  |     |
| BMI                      | 0.017683    | 0.221071  | 0.281805      | 0.392573      | 0.197859  | 1.000000  |     |
| DiabetesPedigreeFunction | -0.033523   | 0.137337  | 0.041265      | 0.183928      | 0.185071  | 0.140647  |     |
| Age                      | 0.544341    | 0.263514  | 0.239528      | -0.113970     | -0.042163 | 0.036242  |     |
| Outcome                  | 0.221898    | 0.466581  | 0.065068      | 0.074752      | 0.130548  | 0.292695  |     |

In [16]:
```python
# separating  the data and labels
X = diabetes.drop(columns = 'Outcome',axis = 1)
Y = diabetes['Outcome']
```

In [17]:
```python
print(X)
```

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

[768 rows x 8 columns]
```

In [18]:
```python
print(Y)
```

```
0      1
1      0
2      1
3      0
4      1
      ..
763    0
764    0
765    0
766    1
767    0
Name: Outcome, Length: 768, dtype: int64
```

# Check the number of zeros value in dataset

In [19]:
```python
print('No. of zero value in Glucose ',diabetes[diabetes ['Glucose']==0].shape[0])
```

```
No. of zero value in Glucose  5
```

In [20]:
```python
print('No. of zero value in BloodPressure ',diabetes[diabetes ['BloodPressure']==0].sha
```

No. of zero value in BloodPressure  35

In [21]:
```python
print('No. of zero value in SkinThickness ',diabetes[diabetes ['SkinThickness']==0].sha
```

No. of zero value in SkinThickness  227

In [22]:
```python
print('No. of zero value in Insulin ',diabetes[diabetes ['Insulin']==0].shape[0])
```
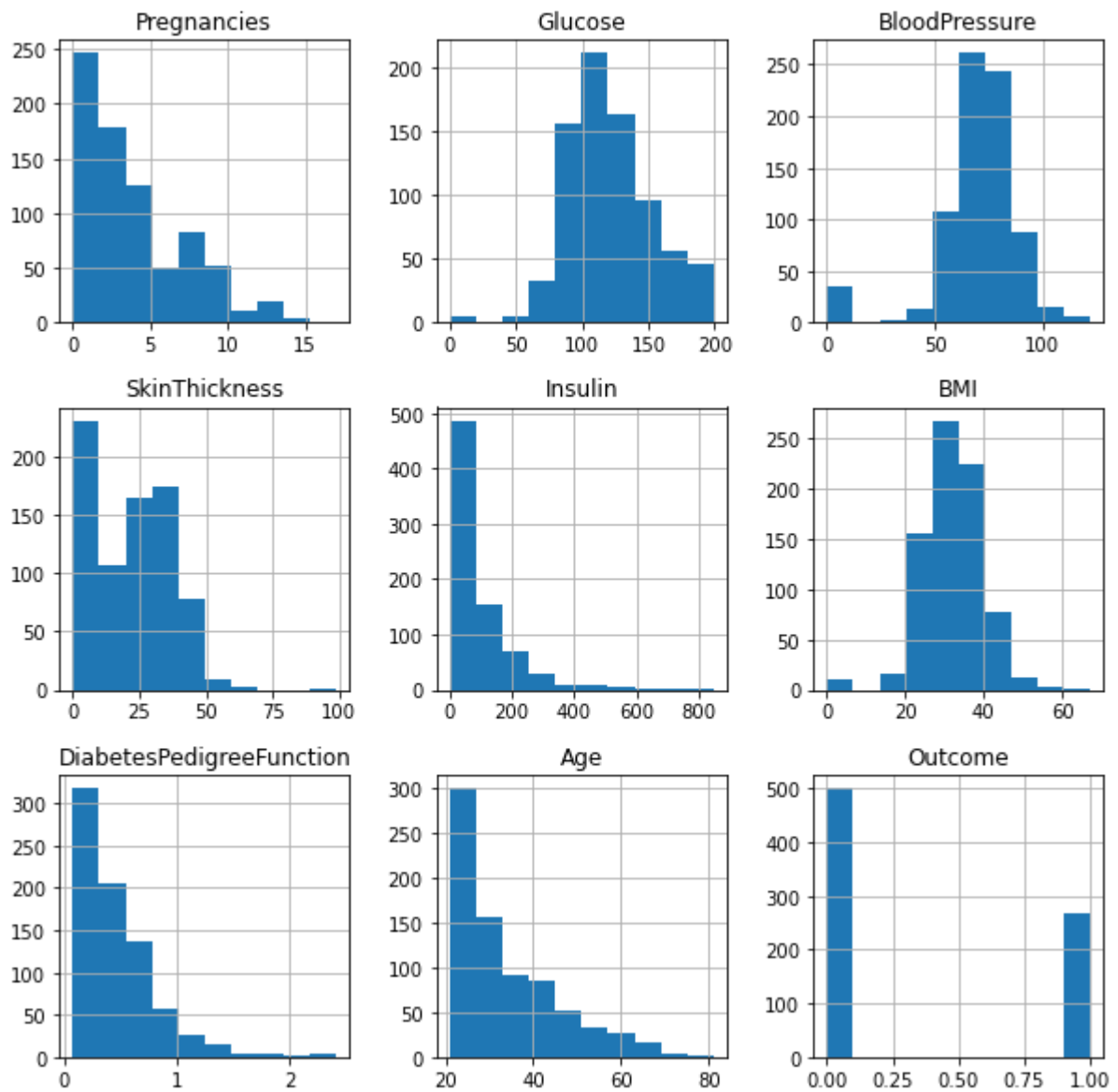
No. of zero value in Insulin  374

In [23]:
```python
print('No. of zero value in BMI ',diabetes[diabetes ['BMI']==0].shape[0])
```

No. of zero value in BMI  11

# Data visualization

In [24]:
```python
# histogram of each feature
diabetes.hist(bins=10,figsize=(10,10))
```

Out[24]:
```
array([[<AxesSubplot:title={'center':'Pregnancies'}>,
        <AxesSubplot:title={'center':'Glucose'}>,
        <AxesSubplot:title={'center':'BloodPressure'}>],
       [<AxesSubplot:title={'center':'SkinThickness'}>,
        <AxesSubplot:title={'center':'Insulin'}>,
        <AxesSubplot:title={'center':'BMI'}>],
       [<AxesSubplot:title={'center':'DiabetesPedigreeFunction'}>,
        <AxesSubplot:title={'center':'Age'}>,
        <AxesSubplot:title={'center':'Outcome'}>]], dtype=object)
```

In [29]:
```python
# scatter plot matrix
from pandas.plotting import scatter_matrix
scatter_matrix(diabetes,figsize = (20,20));
```

In [28]:
```python
# get correlation of each feature in dataset
corrmat= diabetes.corr()
top_corr_features = corrmat.index
plt.figure(figsize=(10,10))
#plot heat map
g=sns.heatmap(diabetes[top_corr_features].corr(),annot=True,cmap="RdYlGn")
```

```
In [30]:   target_name = 'Outcome'

           #Separated object for target feature
           y = diabetes[target_name]

           #separated object for input features
           x = diabetes.drop(target_name, axis=1)
```

```
In [31]:   x.head()
```

Out[31]:

| | 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 |

| | Pregnancies | Glucose | BloodPressure | SkinThickness | Insulin | BMI | DiabetesPedigreeFunction | Age |
|---|---|---|---|---|---|---|---|---|
| **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 [32]:
```python
y.head()
```

Out[32]:
```
0    1
1    0
2    1
3    0
4    1
Name: Outcome, dtype: int64
```

In [33]:
```python
# glucose for diabetes
fig = plt.figure(figsize = (16,6))
sns.distplot(diabetes['Glucose'][diabetes['Outcome']==1])
plt.ylabel('Glucose count')
plt.title('Glucose',fontsize = 20)
```
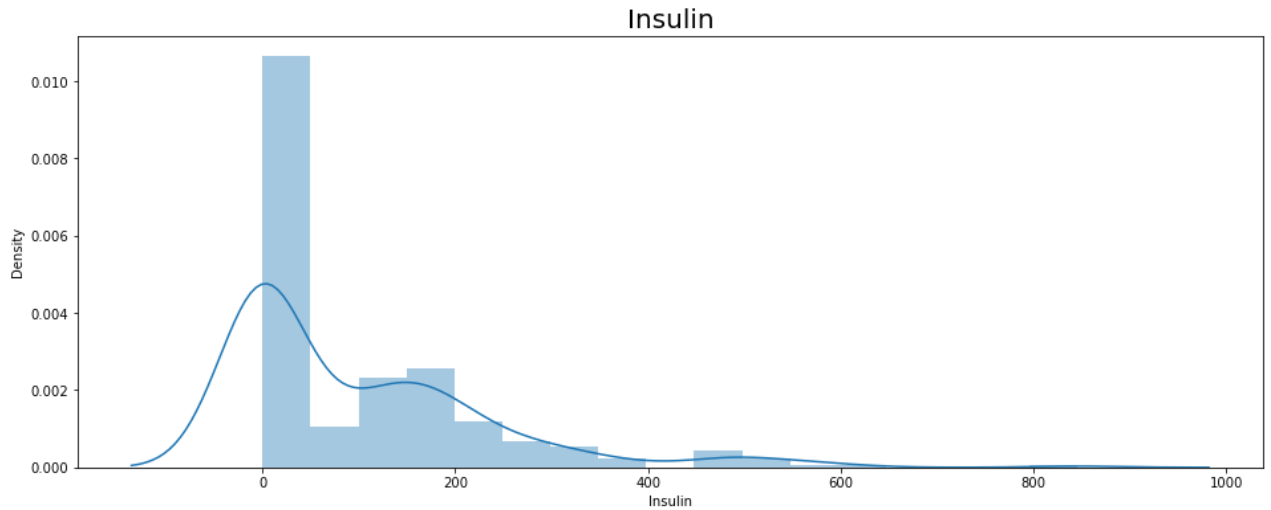
```
C:\ProgramData\Anaconda3\lib\site-packages\seaborn\distributions.py:2619: FutureWarning:
`distplot` is a deprecated function and will be removed in a future version. Please adap
t your code to use either `displot` (a figure-level function with similar flexibility) o
r `histplot` (an axes-level function for histograms).
  warnings.warn(msg, FutureWarning)
```
Out[33]:
```
Text(0.5, 1.0, 'Glucose')
```



In [35]:
```python
# Insuline for diabetes
fig = plt.figure(figsize = (16,6))
sns.distplot(diabetes["Insulin"][diabetes['Outcome']==1])
plt.xticks()
plt.title("Insulin",fontsize = 20)
```
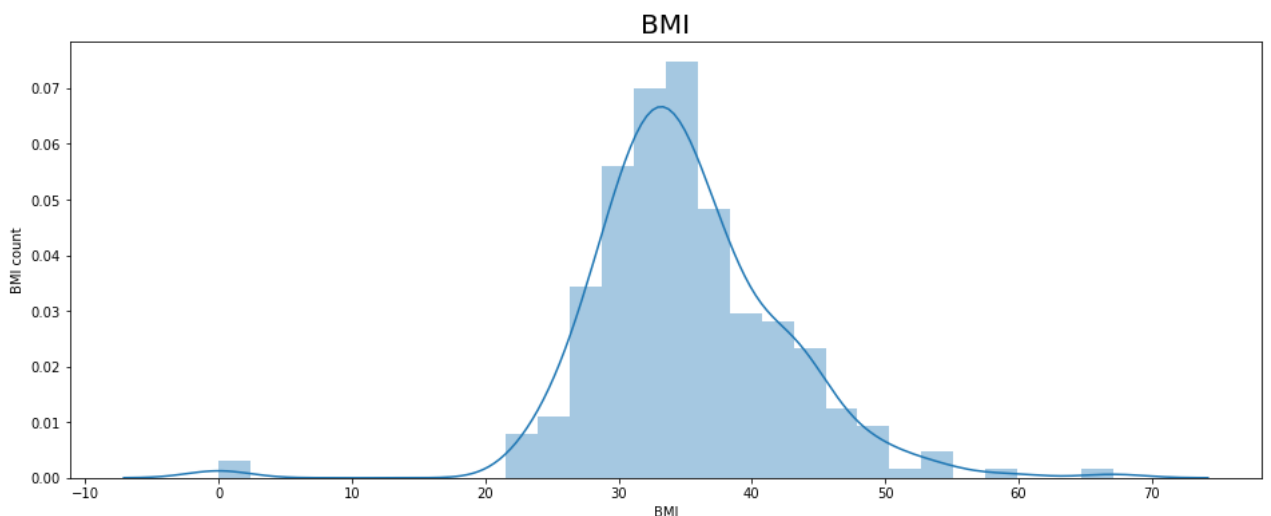
```
C:\ProgramData\Anaconda3\lib\site-packages\seaborn\distributions.py:2619: FutureWarning:
`distplot` is a deprecated function and will be removed in a future version. Please adap
t your code to use either `displot` (a figure-level function with similar flexibility) o
```

r `histplot` (an axes-level function for histograms).
  warnings.warn(msg, FutureWarning)

Out[35]: Text(0.5, 1.0, 'Insulin')



In [36]:
```python
# BMI for diabetes
fig = plt.figure(figsize = (16,6))
sns.distplot(diabetes['BMI'][diabetes['Outcome']==1])
plt.ylabel('BMI count')
plt.title('BMI',fontsize = 20)
```

C:\ProgramData\Anaconda3\lib\site-packages\seaborn\distributions.py:2619: FutureWarning:
`distplot` is a deprecated function and will be removed in a future version. Please adap
t your code to use either `displot` (a figure-level function with similar flexibility) o
r `histplot` (an axes-level function for histograms).
  warnings.warn(msg, FutureWarning)

Out[36]: Text(0.5, 1.0, 'BMI')



In [ ]: