

```
# Importing all the necessary packages
import pandas as pd
import numpy as np
import matplotlib.pyplot as plt
import seaborn as sns

# Loading the data set
data=pd.read_csv('C:\Users\parak\Downloads\diabetes.csv')

data.head(10)

# Pregancies Glucose BloodPressure SkinThickness Insulin BMI DiabetesPedigreeFunction Age Outcome
0 6 148 72 35 0 33.6 0.627 50 1
1 1 85 66 29 0 26.6 0.351 31 0
2 8 183 64 0 0 23.3 0.672 32 1
3 1 89 66 23 94 28.1 0.167 21 0
4 0 137 40 35 168 43.1 2.288 33 1
5 5 116 74 0 0 25.6 0.201 30 0
6 3 78 50 32 88 31.0 0.248 28 1
7 10 115 0 0 0 35.3 0.134 29 0
8 2 197 70 45 543 30.5 0.158 53 1
9 8 125 96 0 0 0.0 0.232 54 1

data.tail()

# Pregancies Glucose BloodPressure SkinThickness Insulin BMI DiabetesPedigreeFunction Age Outcome
763 10 101 76 48 190 32.9 0.171 63 0
764 2 122 70 27 0 36.8 0.340 27 0
765 5 121 72 23 112 26.2 0.245 30 0
766 1 126 60 0 0 30.1 0.349 47 1
767 1 93 70 31 0 30.4 0.315 23 0

# Lets check the dimension of the data set
data.shape

(768, 9)

# Printing the first 5 records of the data set
data.head()

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

# Lets start by checking the count of records in each column of the data set.
# If the count of records is lesser than the total number of records i.e. 768, we can conclude that there
# are null records.
data.info()

<class 'pandas.core.frame.DataFrame'>
RangeIndex: 768 entries, 0 to 767
Data columns (total 9 columns):
# Column Non-Null Count Dtype
---
0 6 148 72 35 0 33.6 0.627 50 1
1 1 85 66 29 0 26.6 0.351 31 0
2 8 183 64 0 0 23.3 0.672 32 1
3 1 89 66 23 94 28.1 0.167 21 0
4 0 137 40 35 168 43.1 2.288 33 1
5 5 116 74 0 0 25.6 0.201 30 0
6 3 78 50 32 88 31.0 0.248 28 1
7 10 115 0 0 0 35.3 0.134 29 0
8 2 197 70 45 543 30.5 0.158 53 1
9 8 125 96 0 0 0.0 0.232 54 1
dtypes: float64(2), int64(7)
memory usage: 54.1 KB

# Since all the predictor columns are continuous in nature, there might be a chance that 0's
# in these columns indicate missing data.
# Lets check the above claim.
data.describe()

# Pregancies Glucose BloodPressure SkinThickness Insulin BMI DiabetesPedigreeFunction Age Outcome
count 768.000000 768.000000 768.000000 768.000000 768.000000 768.000000 768.000000 768.000000 768.000000
mean 3.845052 120.884531 69.105469 20.536458 79.799479 31.992578 0.471876 33.240895 0.346998
std 3.369576 31.972618 19.355807 15.952218 115.244002 7.884160 0.331376 11.760232 0.476951
min 0.000000 0.000000 0.000000 0.000000 0.000000 0.000000 0.078900 24.000000 0.000000
25% 1.000000 99.000000 62.000000 0.000000 0.000000 27.300000 0.243750 24.000000 0.000000
50% 3.000000 117.000000 72.000000 23.000000 30.500000 32.000000 0.372500 29.000000 0.000000
75% 6.000000 140.250000 80.000000 32.000000 127.250000 36.600000 0.626250 41.000000 1.000000
max 17.000000 199.000000 122.000000 99.000000 846.000000 67.100000 2.420000 81.000000 1.000000

# Replacing the 0's with NaN.
# 0 records that have 0's in columns Glucose, Blood Pressure, Skin Thickness, Insulin and BMI will be replaced with NaN
from numpy import nan
data['Glucose']=data['Glucose'].replace(0, np.nan)
data['BloodPressure']=data['BloodPressure'].replace(0, np.nan)
data['SkinThickness']=data['SkinThickness'].replace(0, np.nan)
data['Insulin']=data['Insulin'].replace(0, np.nan)
data['BMI']=data['BMI'].replace(0, np.nan)

# Lets check if the above code has worked
data.head()
# All the 0's have been replaced by NaN

# Pregancies Glucose BloodPressure SkinThickness Insulin BMI DiabetesPedigreeFunction Age Outcome
0 6 148.0 72.0 35.0 NaN 33.6 0.627 50 1
1 1 85.0 66.0 29.0 NaN 26.6 0.351 31 0
2 8 183.0 64.0 NaN NaN 23.3 0.672 32 1
3 1 89.0 66.0 23.0 94.0 28.1 0.167 21 0
4 0 137.0 40.0 35.0 168.0 43.1 2.288 33 1

data.isnull().any()

Pregancies False
Glucose True
BloodPressure True
SkinThickness True
Insulin True
BMI True
DiabetesPedigreeFunction False
Age False
Outcome False
dtype: bool

# Count of NaN values in the dataset
print(data.isnull().sum())

Pregancies 0
Glucose 5
BloodPressure 35
SkinThickness 227
Insulin 374
BMI 11
DiabetesPedigreeFunction 0
Age 0
Outcome 0
dtype: int64

data.median()

Pregancies 3.0000
Glucose 117.0000
BloodPressure 72.0000
SkinThickness 29.0000
Insulin 125.0000
BMI 32.0000
DiabetesPedigreeFunction 0.3725
Age 29.0000
Outcome 0.0000
dtype: float64

# Inputing missing values with their respective columns median
data.fillna(data.median(), inplace=True)

# Checking if the missing values have been inputed
print(data.isnull().sum())

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

# Outlier detection using boxplots
plt.figure(figsize=(20,15))
plt.subplot(4,4,1)
sns.boxplot(data['Pregancies'])

plt.subplot(4,4,2)
sns.boxplot(data['Glucose'])

plt.subplot(4,4,3)
sns.boxplot(data['BloodPressure'])

plt.subplot(4,4,4)
sns.boxplot(data['SkinThickness'])

plt.subplot(4,4,5)
sns.boxplot(data['Insulin'])

plt.subplot(4,4,6)
sns.boxplot(data['BMI'])

plt.subplot(4,4,7)
sns.boxplot(data['DiabetesPedigreeFunction'])

plt.subplot(4,4,8)
sns.boxplot(data['Age'])

C:\Users\parak\Anaconda3\lib\site-packages\seaborn\decorators.py:36: FutureWarning: Pass the following variable as a keyword arg: x. From version 0.12, the only valid positional argument will be 'data', and passing other arguments without an explicit keyword will result in an error or misinterpretation.
C:\Users\parak\Anaconda3\lib\site-packages\seaborn\decorators.py:36: FutureWarning: Pass the following variable as a keyword arg: x. From version 0.12, the only valid positional argument will be 'data', and passing other arguments without an explicit keyword will result in an error or misinterpretation.
C:\Users\parak\Anaconda3\lib\site-packages\seaborn\decorators.py:36: FutureWarning: Pass the following variable as a keyword arg: x. From version 0.12, the only valid positional argument will be 'data', and passing other arguments without an explicit keyword will result in an error or misinterpretation.
C:\Users\parak\Anaconda3\lib\site-packages\seaborn\decorators.py:36: FutureWarning: Pass the following variable as a keyword arg: x. From version 0.12, the only valid positional argument will be 'data', and passing other arguments without an explicit keyword will result in an error or misinterpretation.
C:\Users\parak\Anaconda3\lib\site-packages\seaborn\decorators.py:36: FutureWarning: Pass the following variable as a keyword arg: x. From version 0.12, the only valid positional argument will be 'data', and passing other arguments without an explicit keyword will result in an error or misinterpretation.
C:\Users\parak\Anaconda3\lib\site-packages\seaborn\decorators.py:36: FutureWarning: Pass the following variable as a keyword arg: x. From version 0.12, the only valid positional argument will be 'data', and passing other arguments without an explicit keyword will result in an error or misinterpretation.
C:\Users\parak\Anaconda3\lib\site-packages\seaborn\decorators.py:36: FutureWarning: Pass the following variable as a keyword arg: x. From version 0.12, the only valid positional argument will be 'data', and passing other arguments without an explicit keyword will result in an error or misinterpretation.
C:\Users\parak\Anaconda3\lib\site-packages\seaborn\decorators.py:36: FutureWarning: Pass the following variable as a keyword arg: x. From version 0.12, the only valid positional argument will be 'data', and passing other arguments without an explicit keyword will result in an error or misinterpretation.
C:\Users\parak\Anaconda3\lib\site-packages\seaborn\decorators.py:36: FutureWarning: Pass the following variable as a keyword arg: x. From version 0.12, the only valid positional argument will be 'data', and passing other arguments without an explicit keyword will result in an error or misinterpretation.
C:\Users\parak\Anaconda3\lib\site-packages\seaborn\decorators.py:36: FutureWarning: Pass the following variable as a keyword arg: x. From version 0.12, the only valid positional argument will be 'data', and passing other arguments without an explicit keyword will result in an error or misinterpretation.
C:\Users\parak\Anaconda3\lib\site-packages\seaborn\decorators.py:36: FutureWarning: Pass the following variable as a keyword arg: x. From version 0.12, the only valid positional argument will be 'data', and passing other arguments without an explicit keyword will result in an error or misinterpretation.
C:\Users\parak\Anaconda3\lib\site-packages\seaborn\decorators.py:36: FutureWarning: Pass the following variable as a keyword arg: x. From version 0.12, the only valid positional argument will be 'data', and passing other arguments without an explicit keyword will result in an error or misinterpretation.
C:\Users\parak\Anaconda3\lib\site-packages\seaborn\decorators.py:36: FutureWarning: Pass the following variable as a keyword arg: x. From version 0.12, the only valid positional argument will be 'data', and passing other arguments without an explicit keyword will result in an error or misinterpretation.
C:\Users\parak\Anaconda3\lib\site-packages\seaborn\decorators.py:36: FutureWarning: Pass the following variable as a keyword arg: x. From version 0.12, the only valid positional argument will be 'data', and passing other arguments without an explicit keyword will result in an error or misinterpretation.
C:\Users\parak\Anaconda3\lib\site-packages\seaborn\decorators.py:36: FutureWarning: Pass the following variable as a keyword arg: x. From version 0.12, the only valid positional argument will be 'data', and passing other arguments without an explicit keyword will
```