

objective : design a classifier using KNN to predict whether a patient suffers from diabetes or not

Roll no: 1562

```
In [1]: #Step 1: Import Libraries and Load the dataset
import pandas as pd
import numpy as np
```

```
In [4]: docs=pd.read_csv("C:/1562_AIML/diabetes.csv")
docs.head()
```

Out[4]:

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

```
In [10]: #Step 2: perform exploratory data analysis
docs.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
```

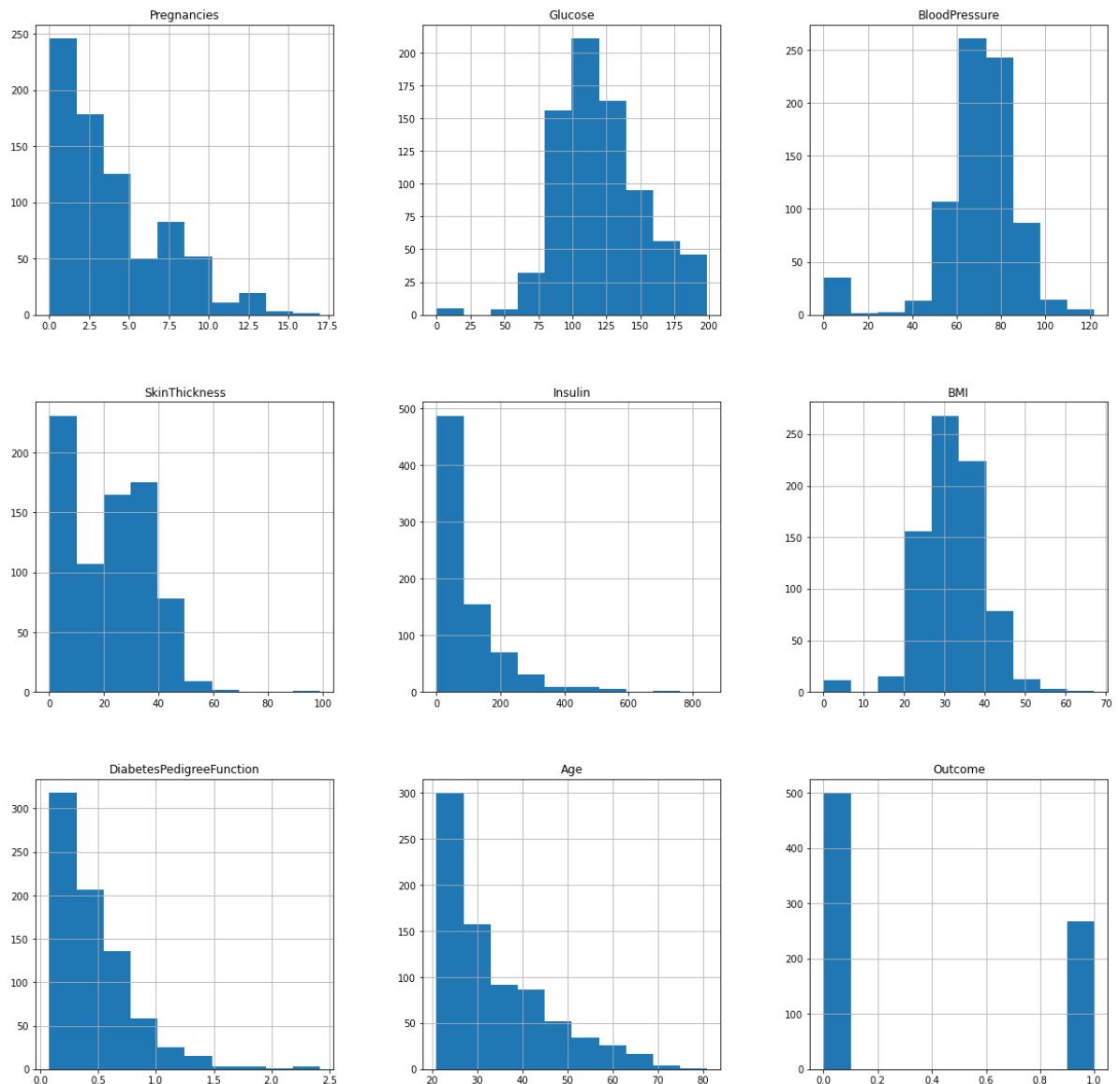
```
In [11]: # analyze the data  
docs.describe()
```

Out[11]:

	Pregnancies	Glucose	BloodPressure	SkinThickness	Insulin	BMI	Diabetes
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 [13]: # plot histogram
docs.hist(figsize=(20,20))
```

```
Out[13]: 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)
```

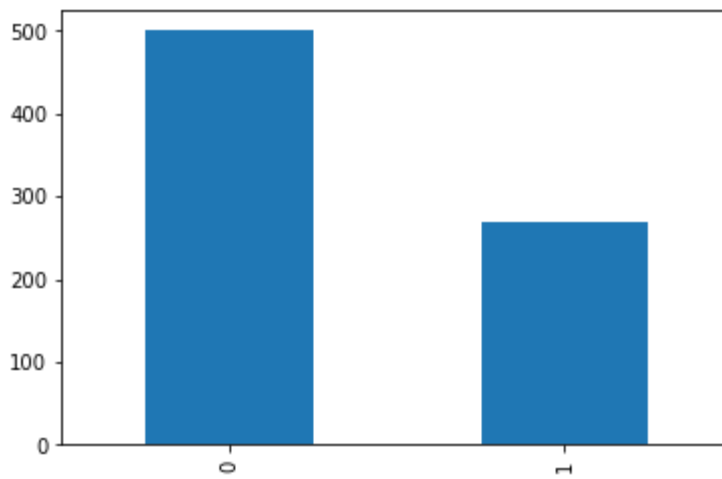


```
In [14]: #show distribution of class labels
print(docs.Outcome.value_counts())
```

```
0    500
1    268
Name: Outcome, dtype: int64
```

```
In [17]: docs.Outcome.value_counts().plot(kind='bar')
```

```
Out[17]: <AxesSubplot:>
```



```
In [19]: # identify the correlation  
import seaborn as sns
```

```
In [21]: sns.pairplot(docs,hue='Outcome')
```

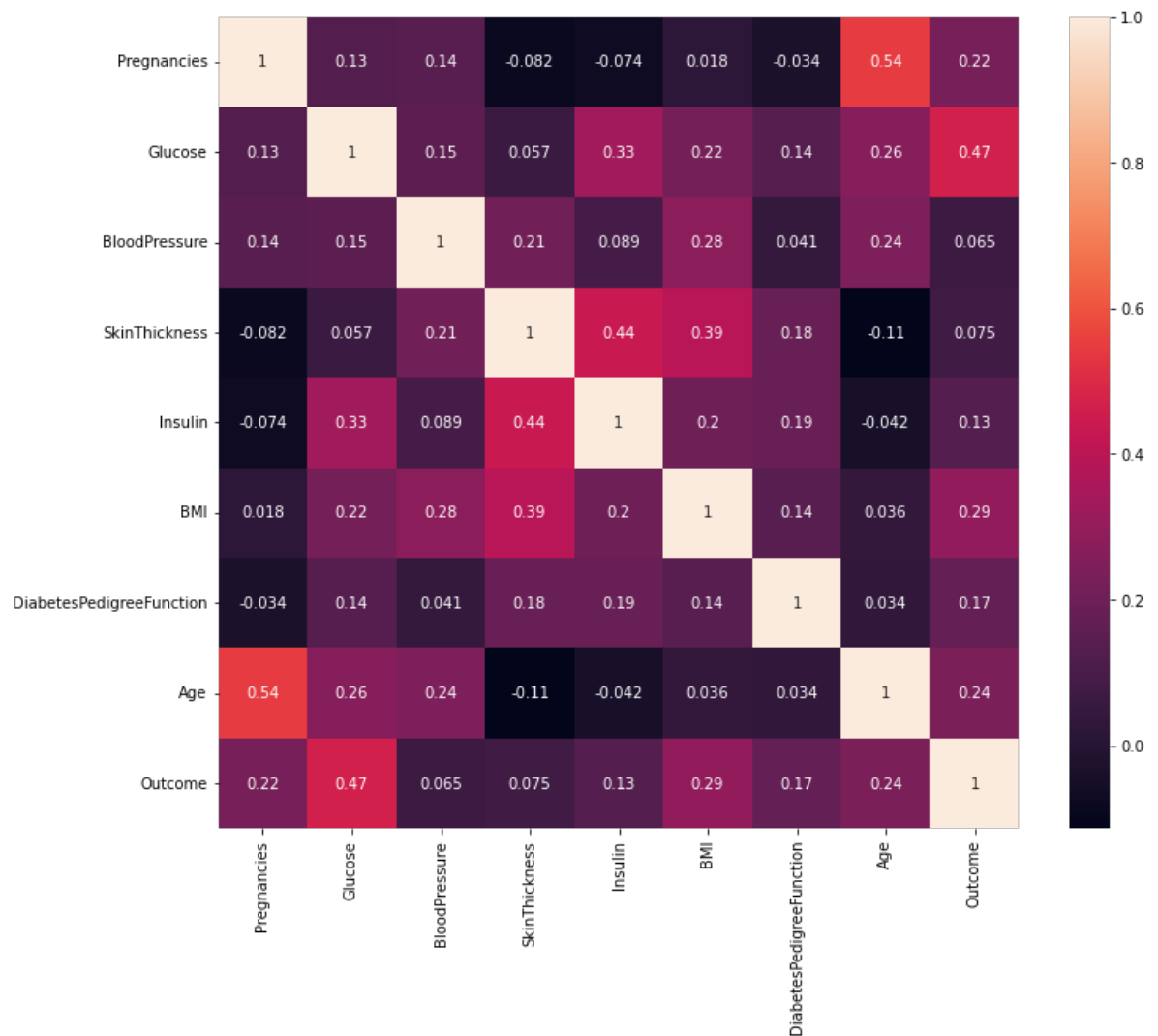
```
Out[21]: <seaborn.axisgrid.PairGrid at 0x2082cdfedf0>
```



```
In [23]: # use heat map
import matplotlib.pyplot as plt
```

```
In [25]: plt.figure(figsize=(12,10))  
sns.heatmap(docs.corr(),annot=True)
```

Out[25]: <AxesSubplot:>



```
In [26]: #implement z-score normalization  
# sepearate x and y  
y = docs.Outcome  
x= docs.drop('Outcome',axis=1)  
x.head()
```

Out[26]:

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

```
In [31]: from sklearn.preprocessing import StandardScaler
```

```
In [36]: sc=StandardScaler()
x_transform=pd.DataFrame(sc.fit_transform(x),columns=x.columns)
```

```
In [34]: x.head()
```

```
Out[34]:
```

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

```
In [37]: x_transform.describe()
```

```
Out[37]:
```

	Pregnancies	Glucose	BloodPressure	SkinThickness	Insulin	BMI
count	7.680000e+02	7.680000e+02	7.680000e+02	7.680000e+02	7.680000e+02	7.680000e+02
mean	2.544261e-17	3.614007e-18	-1.327244e-17	7.994184e-17	-3.556183e-17	2.295979e-16
std	1.000652e+00	1.000652e+00	1.000652e+00	1.000652e+00	1.000652e+00	1.000652e+00
min	-1.141852e+00	-3.783654e+00	-3.572597e+00	-1.288212e+00	-6.928906e-01	-4.060474e+00
25%	-8.448851e-01	-6.852363e-01	-3.673367e-01	-1.288212e+00	-6.928906e-01	-5.955785e-01
50%	-2.509521e-01	-1.218877e-01	1.496408e-01	1.545332e-01	-4.280622e-01	9.419788e-04
75%	6.399473e-01	6.057709e-01	5.632228e-01	7.190857e-01	4.120079e-01	5.847705e-01
max	3.906578e+00	2.444478e+00	2.734528e+00	4.921866e+00	6.652839e+00	4.455807e+00

```
In [39]: # split the dataset into training and test dataset
from sklearn.model_selection import train_test_split
x_train,x_test,y_train,y_test=train_test_split(x,y,train_size=.8,random_state=
x_train.shape
```

```
Out[39]: (614, 8)
```

```
In [ ]:
```