

DATE	18 OCT 2023
TEAM ID	344
PROJECT NAME	AI BASED DIABETES PREDICTION
NAME	MUTHU KUMAR.K

PROJECT NAME :AI BASED DIABETES PREDICTION ;

PHASE 3;

1. Import Required Libraries

```
import numpy as np # linear algebra
import pandas as pd # data processing, CSV file I/O (e.g. pd.read_csv)
import matplotlib.pyplot as plt #to plot charts
import seaborn as sns #used for data visualization
import warnings #avoid warning flash
warnings.filterwarnings('ignore')
```

```
import os
for dirname, _, filenames in os.walk('/kaggle/input'):
    for filename in filenames:
        print(os.path.join(dirname, filename))
```

2. Loading the dataset

```
df=pd.read_csv("../input/pima-indians-diabetes-database/diabetes.csv")
```

3. Exploratory Data Analysis

```
df.head() #get familier with dataset, display the top 5 data records
```

Out[3]:

	Pregnancies	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

```
df.shape #getting to know about rows and columns we're dealing with - 768 rows , 9 columns
```

output

(768, 9)

```
df.columns #learning about the columns
```

output

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

```
df.dtypes #knowledge of data type helps for computation
```

output

```
Pregnancies      int64  
Glucose           int64  
BloodPressure    int64  
SkinThickness    int64  
Insulin          int64  
BMI              float64  
DiabetesPedigreeFunction float64  
Age              int64  
Outcome          int64
```

```
dtype: object
```

```
df.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

df.describe()

Pregna ncies	Gluco se	BloodPr essure	SkinThi ckness	Insulin	BMI	DiabetesPedigr eeFunction	Age	Outco me	
count	768.00 0000	768.000 000	768.000 000	768.00 0000	768.00 0000	768.000000	768.00 0000	768.00 0000	768.00 0000
mean	3.8450 52	120.894 531	69.1054 69	20.536 458	79.799 479	31.992578	0.4718 76	33.240 885	0.3489 58
std	3.3695 78	31.9726 18	19.3558 07	15.952 218	115.24 4002	7.884160	0.3313 29	11.760 232	0.4769 51
min	0.0000 00	0.00000 0	0.00000 0	0.0000 00	0.0000 00	0.000000	0.0780 00	21.000 000	0.0000 00
25%	1.0000 00	99.0000 00	62.0000 00	0.0000 00	0.0000 00	27.300000	0.2437 50	24.000 000	0.0000 00
50%	3.0000 00	117.000 000	72.0000 00	23.000 000	30.500 000	32.000000	0.3725 00	29.000 000	0.0000 00
75%	6.0000 00	140.250 000	80.0000 00	32.000 000	127.25 0000	36.600000	0.6262 50	41.000 000	1.0000 00

max	17.000 000	199.000 000	122.000 000	99.000 000	846.00 0000	67.100000	2.4200 00	81.000 000	1.0000 00
-----	---------------	----------------	----------------	---------------	----------------	-----------	--------------	---------------	--------------

```
df=df.drop_duplicates()
```

```
df.isnull().sum()
```

output

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

```
print(df[df['BloodPressure']==0].shape[0])
```

```
print(df[df['Glucose']==0].shape[0])
```

```
print(df[df['SkinThickness']==0].shape[0])
```

```
print(df[df['Insulin']==0].shape[0])
```

```
print(df[df['BMI']==0].shape[0])
```

output

```
35
```

```
5
```

```
227
```

```
374
```

```
11
```