

```

#import the library
import pandas as pd
import numpy as np
import matplotlib.pyplot as plt
import sklearn
import matplotlib.pyplot as plt

#load the dataset
dataset = pd.read_csv('/Users/shraddhalipane/Downloads/Project
2/Healthcare - Diabetes/health care diabetes.csv')

#know dataset
data=pd.read_csv('/Users/shraddhalipane/Downloads/Project 2/Healthcare
- Diabetes/health care diabetes.csv')

# shape , info,describe , null values, neagative values , head, tail

```

```
data.shape
```

```
(768, 9)
```

```
data.describe()
```

	Pregnancies	Glucose	BloodPressure	SkinThickness
Insulin \				
count	768.000000	768.000000	768.000000	768.000000
mean	3.845052	120.894531	69.105469	20.536458
std	3.369578	31.972618	19.355807	15.952218
min	0.000000	0.000000	0.000000	0.000000
25%	1.000000	99.000000	62.000000	0.000000
50%	3.000000	117.000000	72.000000	23.000000
75%	6.000000	140.250000	80.000000	32.000000
max	17.000000	199.000000	122.000000	99.000000

	BMI	DiabetesPedigreeFunction	Age	Outcome
count	768.000000	768.000000	768.000000	768.000000
mean	31.992578	0.471876	33.240885	0.348958
std	7.884160	0.331329	11.760232	0.476951
min	0.000000	0.078000	21.000000	0.000000
25%	27.300000	0.243750	24.000000	0.000000
50%	32.000000	0.372500	29.000000	0.000000
75%	36.600000	0.626250	41.000000	1.000000
max	67.100000	2.420000	81.000000	1.000000

```
data.isnull().sum()
```

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

```
data.head()
```

	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

	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

```
data.info()
```

```
<class 'pandas.core.frame.DataFrame'>
```

```
RangeIndex: 768 entries, 0 to 767
```

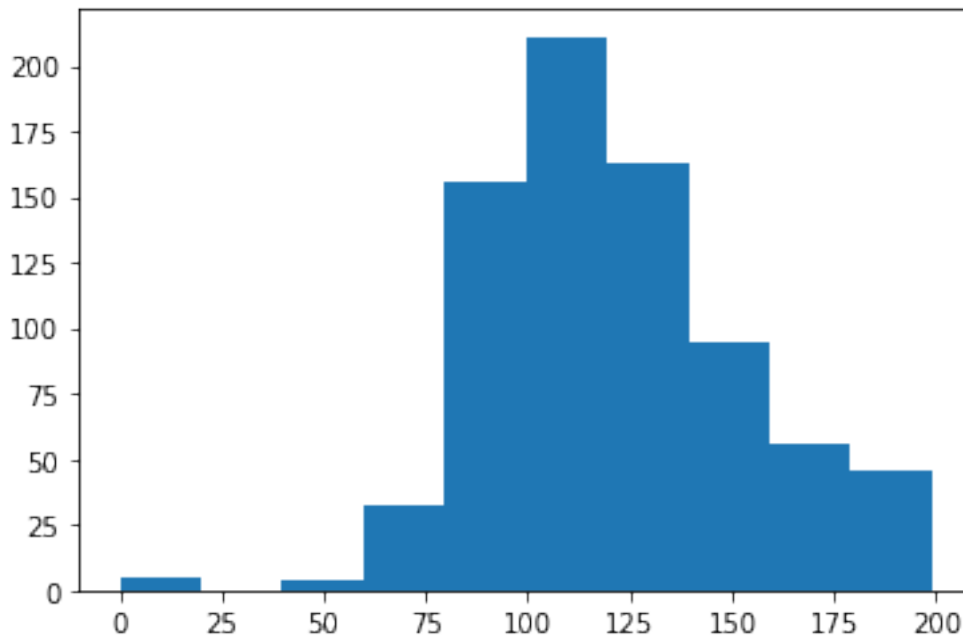
```
Data columns (total 9 columns):
```

#	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
```

```
plt.hist(data['Glucose'])
```

```
(array([ 5.,  0.,  4., 32., 156., 211., 163., 95., 56., 46.]),
 array([ 0., 19.9, 39.8, 59.7, 79.6, 99.5, 119.4, 139.3, 159.2,
        179.1, 199. ]),
 <BarContainer object of 10 artists>)
```



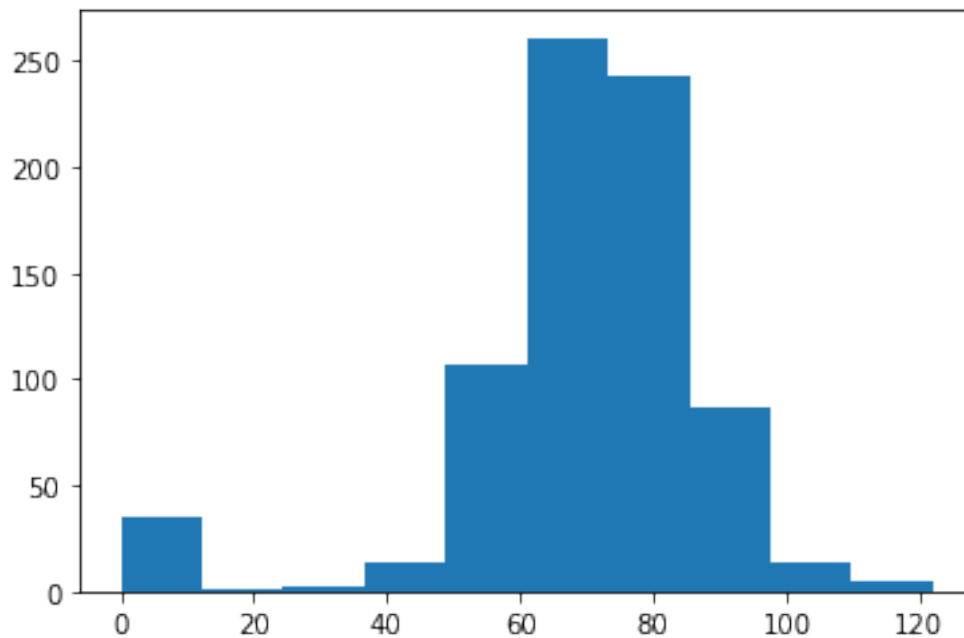
```
data['Glucose'].value_counts().head(7)
```

```
99      17
100     17
111     14
129     14
125     14
106     14
112     13
```

```
Name: Glucose, dtype: int64
```

```
plt.hist(data['BloodPressure'])
```

```
(array([ 35.,  1.,  2., 13., 107., 261., 243., 87., 14.,  5.]),
 array([ 0., 12.2, 24.4, 36.6, 48.8, 61., 73.2, 85.4, 97.6,
        109.8, 122. ]),
 <BarContainer object of 10 artists>)
```



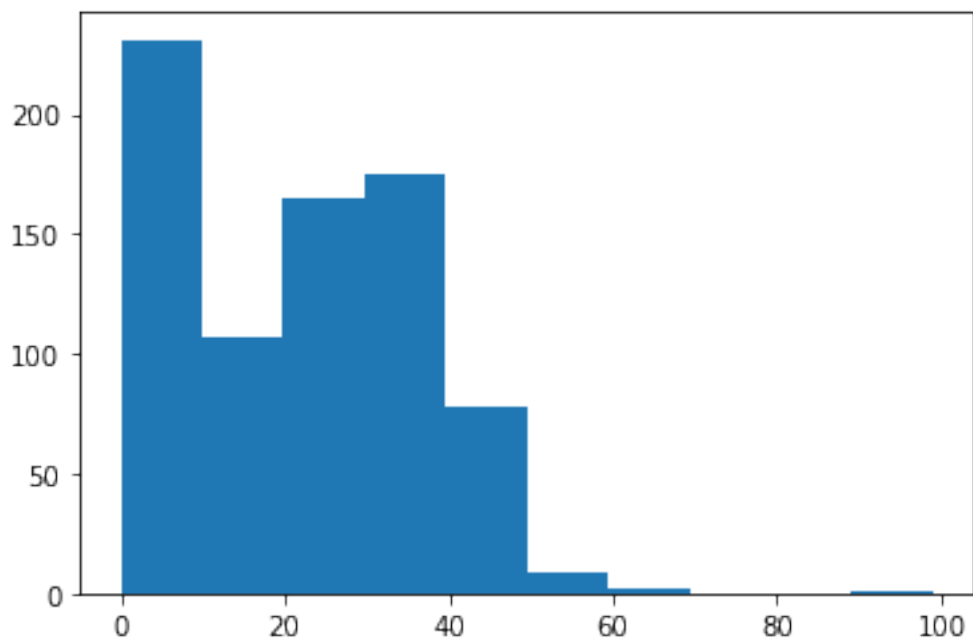
```
data['BloodPressure'].value_counts().head()
```

```
70    57
74    52
78    45
68    45
72    44
```

```
Name: BloodPressure, dtype: int64
```

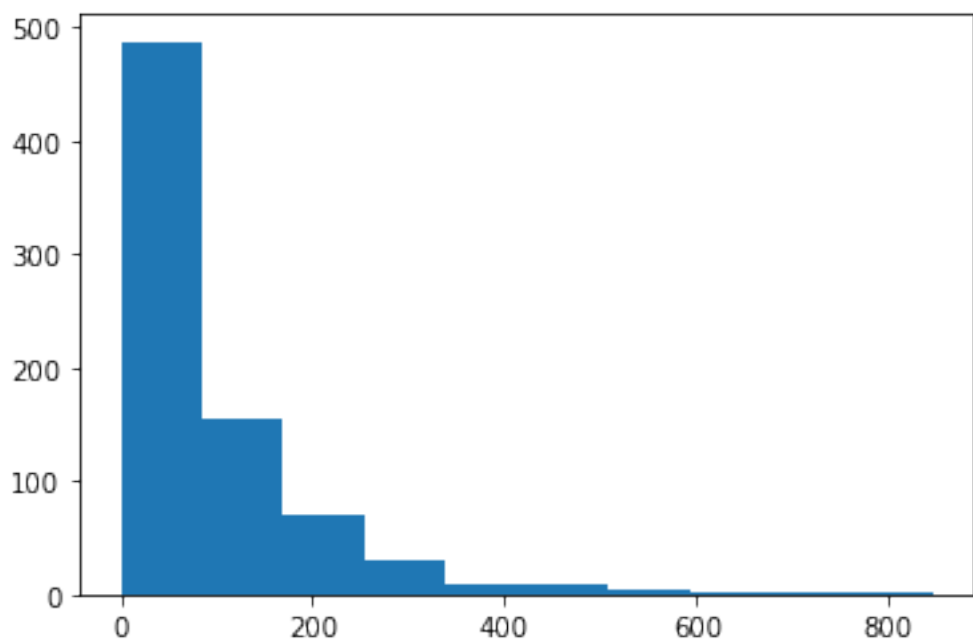
```
plt.hist(data['SkinThickness'])
```

```
(array([231., 107., 165., 175.,  78.,   9.,   2.,   0.,   0.,   1.]),
 array([ 0. ,  9.9, 19.8, 29.7, 39.6, 49.5, 59.4, 69.3, 79.2, 89.1,
 99. ]),
 <BarContainer object of 10 artists>)
```



```
plt.hist(data['Insulin'])
```

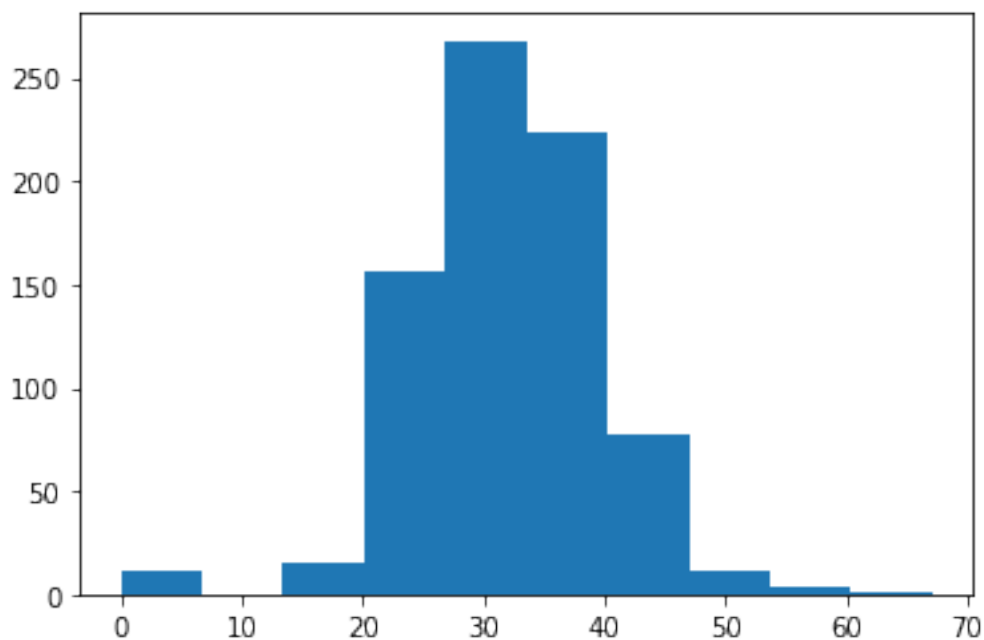
```
(array([487., 155., 70., 30., 8., 9., 5., 1., 2., 1.]),
 array([ 0., 84.6, 169.2, 253.8, 338.4, 423., 507.6, 592.2, 676.8,
        761.4, 846. ]),
 <BarContainer object of 10 artists>)
```



```
plt.hist(data['BMI'])
```

```
(array([ 11., 0., 15., 156., 268., 224., 78., 12., 3., 1.]),
 array([ 0., 6.71, 13.42, 20.13, 26.84, 33.55, 40.26, 46.97, 53.68,
```

```
60.39, 67.1 ]),
<BarContainer object of 10 artists>)
```



```
data['Glucose'].value_counts()
```

```
99      17
100     17
111     14
129     14
125     14
```

```
..
191      1
177      1
44        1
62        1
190       1
```

```
Name: Glucose, Length: 136, dtype: int64
```

```
data['BloodPressure'].value_counts()
```

```
70      57
74      52
78      45
68      45
72      44
64      43
80      40
76      39
60      37
0       35
62      34
```

66	30
82	30
88	25
84	23
90	22
86	21
58	21
50	13
56	12
52	11
54	11
75	8
92	8
65	7
85	6
94	6
48	5
96	4
44	4
100	3
106	3
98	3
110	3
55	2
108	2
104	2
46	2
30	2
122	1
95	1
102	1
61	1
24	1
38	1
40	1
114	1

Name: BloodPressure, dtype: int64

data['SkinThickness'].value_counts()

0	227
32	31
30	27
27	23
23	22
33	20
28	20
18	20
31	19
19	18
39	18

29	17
40	16
25	16
26	16
22	16
37	16
41	15
35	15
36	14
15	14
17	14
20	13
24	12
42	11
13	11
21	10
46	8
34	8
12	7
38	7
11	6
43	6
16	6
45	6
14	6
44	5
10	5
48	4
47	4
49	3
50	3
8	2
7	2
52	2
54	2
63	1
60	1
56	1
51	1
99	1

Name: SkinThickness, dtype: int64

data['Insulin'].value_counts()

0	374
105	11
130	9
140	9
120	8
...	
73	1


```
171      1
255      1
52       1
112      1
Name: Insulin, Length: 186, dtype: int64
```

```
data['BMI'].value_counts()
```

```
32.0      13
31.6      12
31.2      12
0.0       11
32.4      10
..
36.7       1
41.8       1
42.6       1
42.8       1
46.3       1
Name: BMI, Length: 248, dtype: int64
```

```
data.describe()
```

	Pregnancies	Glucose	BloodPressure	SkinThickness
Insulin \				
count	768.000000	768.000000	768.000000	768.000000
768.000000				
mean	3.845052	120.894531	69.105469	20.536458
79.799479				
std	3.369578	31.972618	19.355807	15.952218
115.244002				
min	0.000000	0.000000	0.000000	0.000000
0.000000				
25%	1.000000	99.000000	62.000000	0.000000
0.000000				
50%	3.000000	117.000000	72.000000	23.000000
30.500000				
75%	6.000000	140.250000	80.000000	32.000000
127.250000				
max	17.000000	199.000000	122.000000	99.000000
846.000000				

	BMI	DiabetesPedigreeFunction	Age	Outcome
count	768.000000	768.000000	768.000000	768.000000
mean	31.992578	0.471876	33.240885	0.348958
std	7.884160	0.331329	11.760232	0.476951
min	0.000000	0.078000	21.000000	0.000000
25%	27.300000	0.243750	24.000000	0.000000
50%	32.000000	0.372500	29.000000	0.000000
75%	36.600000	0.626250	41.000000	1.000000
max	67.100000	2.420000	81.000000	1.000000

```
data['Outcome'].value_counts().head(7)
```

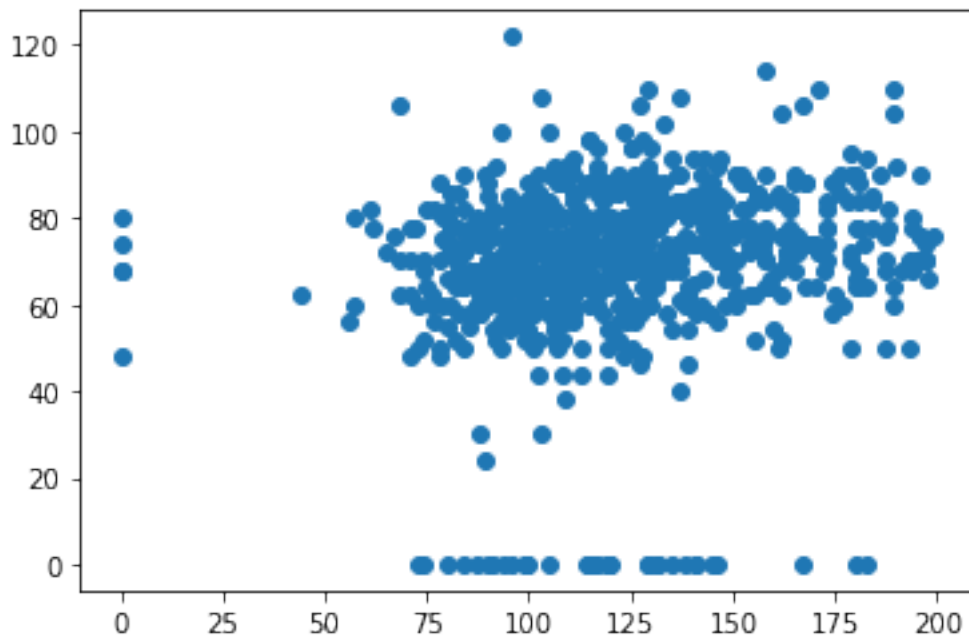
```
0    500
```

```
1    268
```

```
Name: Outcome, dtype: int64
```

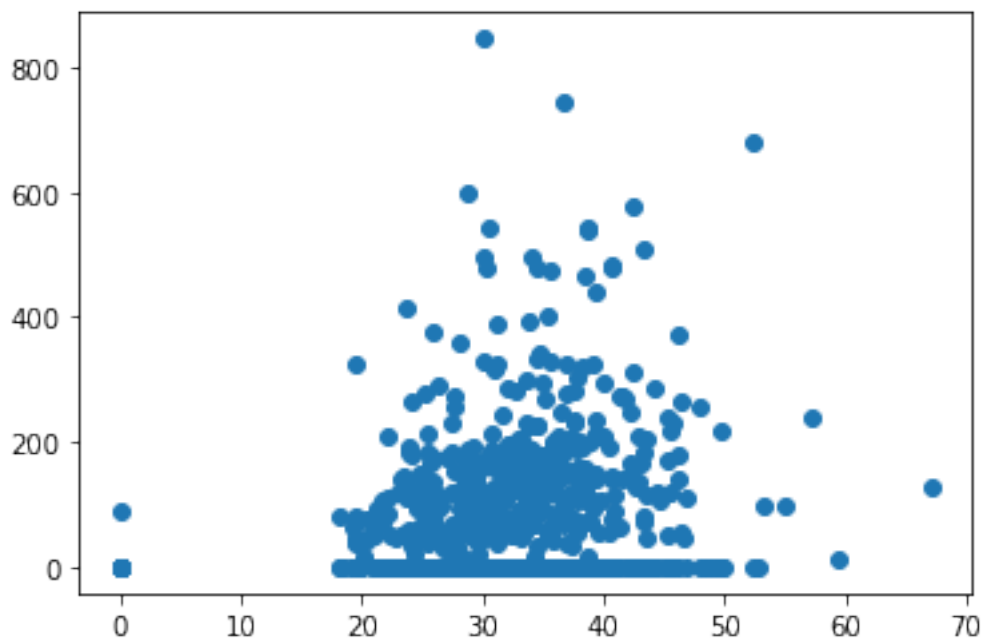
```
plt.scatter(data['Glucose'],data['BloodPressure'])
```

```
<matplotlib.collections.PathCollection at 0x7fc5dd60d310>
```



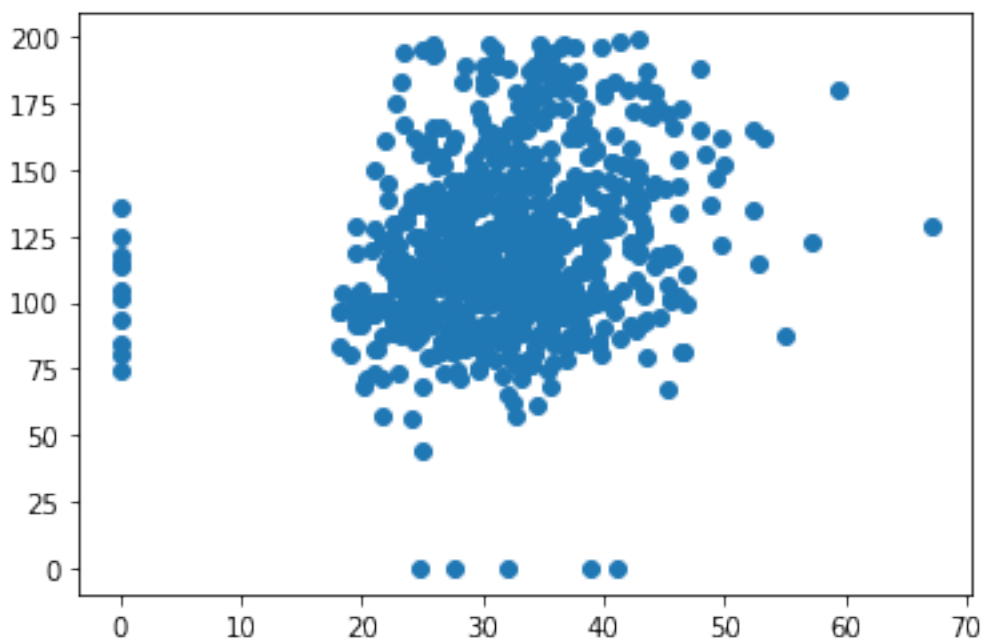
```
plt.scatter(data['BMI'],data['Insulin'])
```

```
<matplotlib.collections.PathCollection at 0x7fc5dd7e7e80>
```



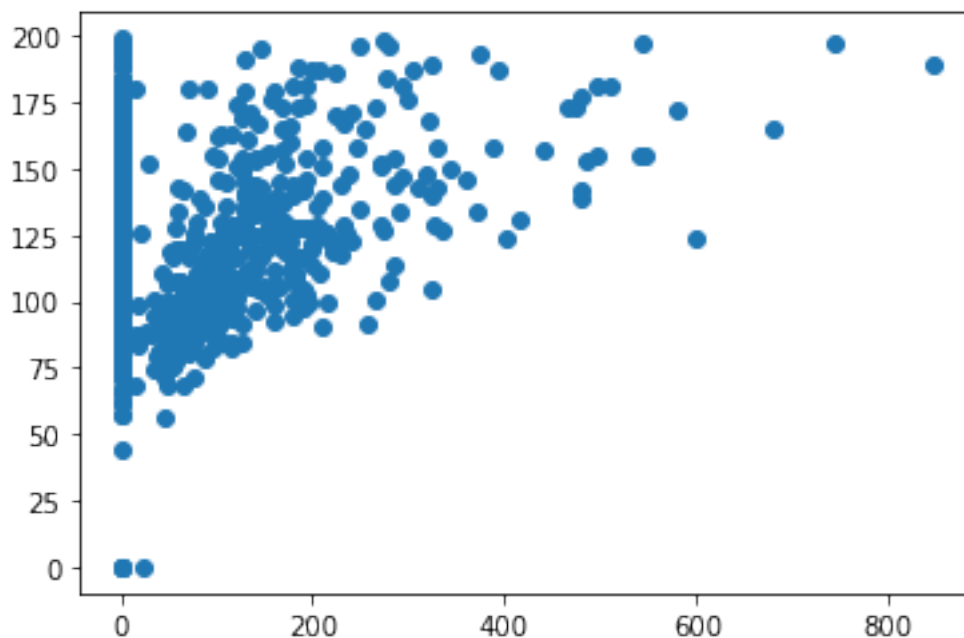
```
plt.scatter(data['BMI'],data['Glucose'])
```

```
<matplotlib.collections.PathCollection at 0x7fc5ca007910>
```

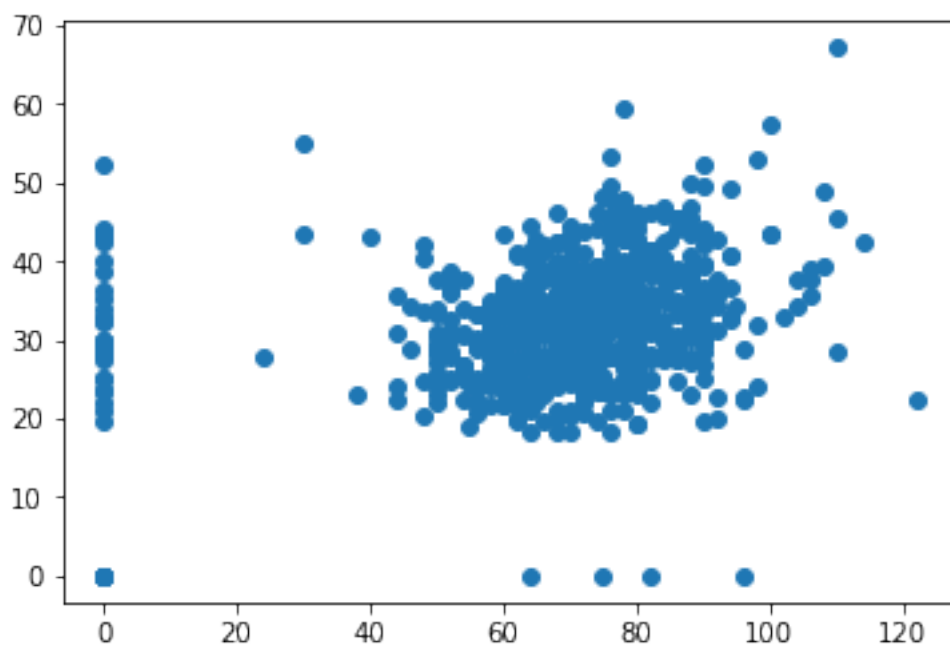


```
plt.scatter(data['Insulin'],data['Glucose'])
```

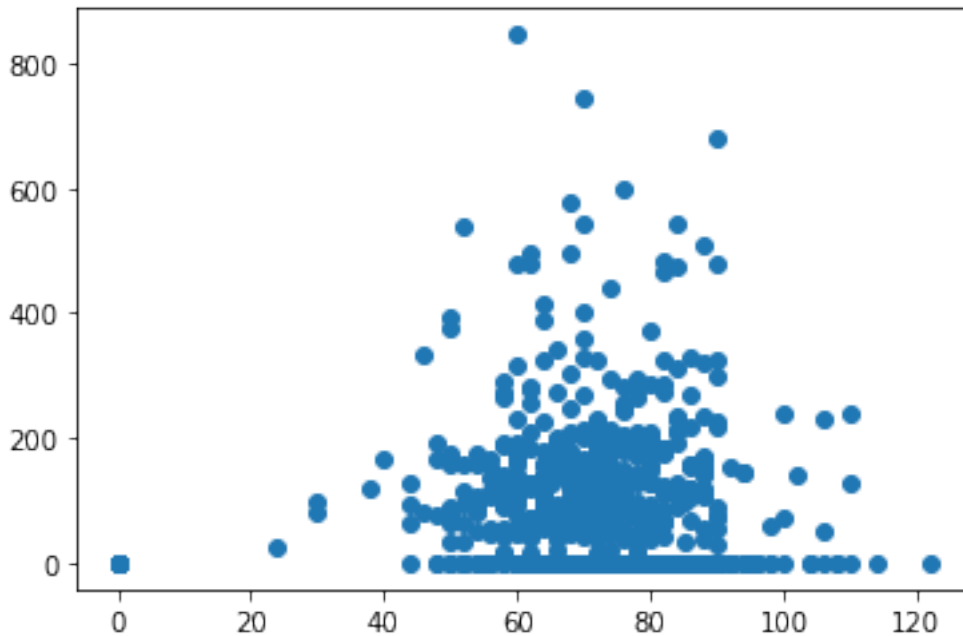
```
<matplotlib.collections.PathCollection at 0x7fc5b8215b20>
```



```
plt.scatter(data['BloodPressure'],data['BMI'])  
<matplotlib.collections.PathCollection at 0x7fc5c9f4faf0>
```



```
plt.scatter(data['BloodPressure'],data['Insulin'])  
<matplotlib.collections.PathCollection at 0x7fc5880b3e80>
```



```
#heatmap
data.corr()
```

	Pregnancies	Glucose	BloodPressure	
SkinThickness \				
Pregnancies	1.000000	0.129459	0.141282	-
0.081672				
Glucose	0.129459	1.000000	0.152590	
0.057328				
BloodPressure	0.141282	0.152590	1.000000	
0.207371				
SkinThickness	-0.081672	0.057328	0.207371	
1.000000				
Insulin	-0.073535	0.331357	0.088933	
0.436783				
BMI	0.017683	0.221071	0.281805	
0.392573				
DiabetesPedigreeFunction	-0.033523	0.137337	0.041265	
0.183928				
Age	0.544341	0.263514	0.239528	-
0.113970				
Outcome	0.221898	0.466581	0.065068	
0.074752				
	Insulin	BMI	DiabetesPedigreeFunction	
\				
Pregnancies	-0.073535	0.017683	-0.033523	
Glucose	0.331357	0.221071	0.137337	

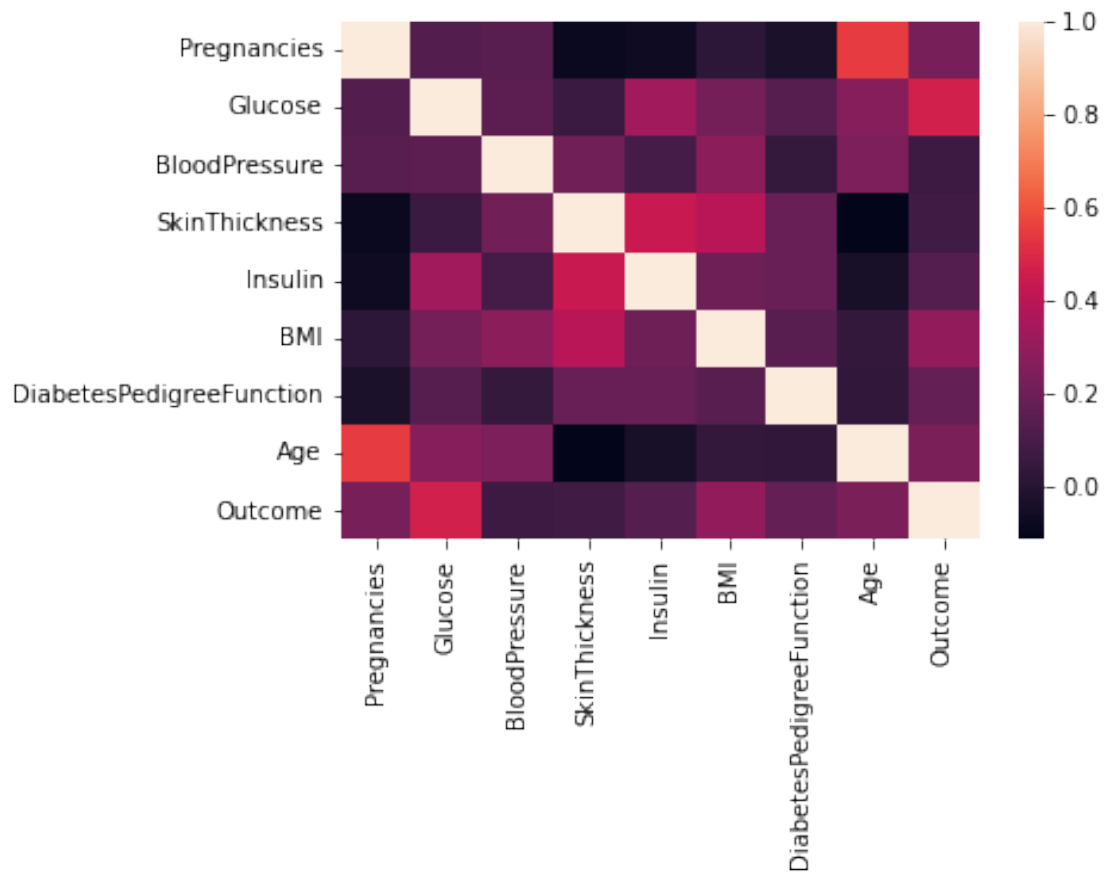
BloodPressure	0.088933	0.281805	0.041265
SkinThickness	0.436783	0.392573	0.183928
Insulin	1.000000	0.197859	0.185071
BMI	0.197859	1.000000	0.140647
DiabetesPedigreeFunction	0.185071	0.140647	1.000000
Age	-0.042163	0.036242	0.033561
Outcome	0.130548	0.292695	0.173844

	Age	Outcome
Pregnancies	0.544341	0.221898
Glucose	0.263514	0.466581
BloodPressure	0.239528	0.065068
SkinThickness	-0.113970	0.074752
Insulin	-0.042163	0.130548
BMI	0.036242	0.292695
DiabetesPedigreeFunction	0.033561	0.173844
Age	1.000000	0.238356
Outcome	0.238356	1.000000

```
import seaborn as sns
```

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

```
<AxesSubplot:>
```



```
data.head()
```

	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

	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

```
#model building
```

```
# Extract the features
```

```
features=data.iloc[:,0:8].values
```

```
label=data.iloc[:,8].values
```

```
features
```

```
array([[ 6.   , 148.   , 72.   , ..., 33.6 , 0.627, 50.   ],
       [ 1.   , 85.   , 66.   , ..., 26.6 , 0.351, 31.   ],
       [ 8.   , 183.   , 64.   , ..., 23.3 , 0.672, 32.   ],
       ...,
       [ 5.   , 121.   , 72.   , ..., 26.2 , 0.245, 30.   ],
       [ 1.   , 126.   , 60.   , ..., 30.1 , 0.349, 47.   ],
       [ 1.   , 93.   , 70.   , ..., 30.4 , 0.315, 23.   ]])
```

```
label
```

```
array([1, 0, 1, 0, 1, 0, 1, 0, 1, 1, 0, 1, 0, 1, 1, 1, 1, 1, 0, 1, 0,
0,
      1, 1, 1, 1, 1, 0, 0, 0, 0, 1, 0, 0, 0, 0, 0, 1, 1, 1, 0, 0, 0,
1,
      0, 1, 0, 0, 1, 0, 0, 0, 0, 1, 0, 0, 1, 0, 0, 0, 0, 1, 0, 0, 1,
0,
      1, 0, 0, 0, 1, 0, 1, 0, 0, 0, 0, 0, 1, 0, 0, 0, 0, 0, 1, 0, 0,
0,
      1, 0, 0, 0, 0, 1, 0, 0, 0, 0, 0, 0, 1, 1, 0, 0, 0, 0, 0, 0, 0,
1,
      1, 1, 0, 0, 1, 1, 1, 0, 0, 0, 1, 0, 0, 0, 1, 1, 0, 0, 1, 1, 1,
1,
      1, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 1, 0, 0, 0, 0, 0, 0, 0, 0, 1,
0,
      1, 1, 0, 0, 0, 1, 0, 0, 0, 0, 1, 1, 0, 0, 0, 0, 1, 1, 0, 0, 0,
1,
      0, 1, 0, 1, 0, 0, 0, 0, 0, 1, 1, 1, 1, 1, 0, 0, 1, 1, 0, 1, 0,
1,
      1, 1, 0, 0, 0, 0, 0, 0, 1, 1, 0, 1, 0, 0, 0, 1, 1, 1, 1, 0, 1,
1,
      1, 1, 0, 0, 0, 0, 0, 1, 0, 0, 1, 1, 0, 0, 0, 1, 1, 1, 1, 0, 0,
0,
      1, 1, 0, 1, 0, 0, 0, 0, 0, 0, 0, 0, 0, 1, 1, 0, 0, 0, 1, 0, 1,
0,
      1, 0, 1, 0, 0, 1, 1, 0, 0, 0, 0, 0, 1, 0, 0, 0, 1, 0, 0, 1, 1,
0,
      0, 1, 0, 0, 0, 1, 1, 1, 0, 0, 1, 0, 1, 0, 1, 1, 0, 1, 0, 0, 1,
0,
      1, 1, 0, 0, 1, 0, 1, 0, 0, 1, 0, 1, 0, 1, 1, 1, 0, 0, 1, 0, 1,
0,
      0, 0, 1, 0, 0, 0, 0, 1, 1, 1, 0, 0, 0, 0, 0, 0, 0, 0, 0, 1, 0,
0,
      0, 0, 0, 1, 1, 1, 0, 1, 1, 0, 0, 1, 0, 0, 1, 0, 0, 1, 1, 0, 0,
```



```

0,      0, 1, 0, 0, 1, 0, 0, 0, 0, 0, 0, 0, 1, 1, 1, 0, 0, 1, 0, 0, 1,
1,      0, 1, 0, 1, 1, 0, 1, 0, 1, 0, 1, 0, 1, 1, 0, 0, 0, 0, 1, 1, 0,
0,      0, 1, 0, 0, 0, 0, 1, 1, 0, 1, 0, 1, 0, 0, 0, 0, 0, 1, 0, 0, 0,
0,      1, 0, 0, 1, 1, 1, 0, 0, 1, 0, 0, 1, 0, 0, 0, 1, 0, 0, 1, 0, 0,
0,      0, 0, 0, 0, 0, 0, 1, 0, 0, 0, 0, 0, 0, 0, 0, 1, 0, 0, 0, 1, 0, 0,
0,      1, 1, 0, 0, 0, 0, 0, 0, 0, 1, 0, 0, 0, 0, 1, 0, 0, 0, 1, 0, 0,
0,      1, 0, 0, 0, 1, 0, 0, 0, 0, 1, 1, 0, 0, 0, 0, 0, 0, 1, 0, 0, 0,
0,      0, 0, 0, 0, 0, 0, 0, 1, 0, 0, 0, 1, 1, 1, 1, 0, 0, 1, 1, 0, 0,
0,      0, 0, 0, 0, 0, 0, 0, 0, 0, 1, 1, 0, 0, 0, 0, 0, 0, 0, 1, 0,
0,      0, 0, 0, 0, 0, 1, 0, 1, 1, 0, 0, 0, 1, 0, 1, 0, 1, 0, 1, 0, 1,
0,      0, 1, 0, 0, 1, 0, 0, 0, 0, 1, 1, 0, 1, 0, 0, 0, 0, 1, 1, 0, 1,
0,      0, 0, 1, 1, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 1, 0, 0, 0, 0, 1, 0,
1,      1, 0, 0, 0, 1, 0, 0, 0, 1, 1, 1, 0, 0, 0, 0, 0, 0, 1, 0, 0, 0,
1,      0, 1, 1, 1, 1, 0, 1, 1, 0, 0, 0, 0, 0, 0, 0, 1, 1, 0, 1, 0, 0,
0,      0, 1, 0, 0, 0, 0, 0, 1, 0, 1, 0, 1, 0, 1, 1, 0, 0, 0, 0, 1, 1,
0,      0, 0, 1, 0, 1, 1, 0, 0, 1, 0, 0, 1, 1, 0, 0, 1, 0, 0, 1, 0, 0,
0,      0, 0, 0, 0, 1, 1, 1, 0, 0, 0, 0, 0, 1, 1, 0, 0, 1, 0, 0, 1,
0,      1, 1, 1, 0, 0, 1, 1, 1, 0, 1, 0, 1, 0, 1, 0, 0, 0, 0, 1, 0])

```

```

from sklearn.model_selection import train_test_split

```

```

X_train,X_test,y_train,y_test = train_test_split(features,
                                                    label,
                                                    test_size=0.2,
                                                    random_state =10)

```

```

from sklearn.linear_model import LogisticRegression
classifier = LogisticRegression()
classifier.fit(X_train,y_train)

```

```

LogisticRegression()

```

```

print(classifier.score(X_train,y_train))
print(classifier.score(X_test,y_test))

```

```
0.7719869706840391
0.7662337662337663
```

```
from sklearn.metrics import confusion_matrix
cm = confusion_matrix(label, classifier.predict(features))
cm
```

```
array([[446,  54],
       [122, 146]])
```

```
from sklearn.metrics import classification_report
print(classification_report(label, classifier.predict(features)))
```

	precision	recall	f1-score	support
0	0.79	0.89	0.84	500
1	0.73	0.54	0.62	268
accuracy			0.77	768
macro avg	0.76	0.72	0.73	768
weighted avg	0.77	0.77	0.76	768

```
from sklearn.metrics import roc_curve
from sklearn.metrics import roc_auc_score      #ROC Curve

probability = classifier.predict_proba(features) # predict
probabilities
```

```
probability = probability[:, 1] # probability for positive outcomes
```

```
auc = roc_auc_score(label, probability) # calculating AUC
```

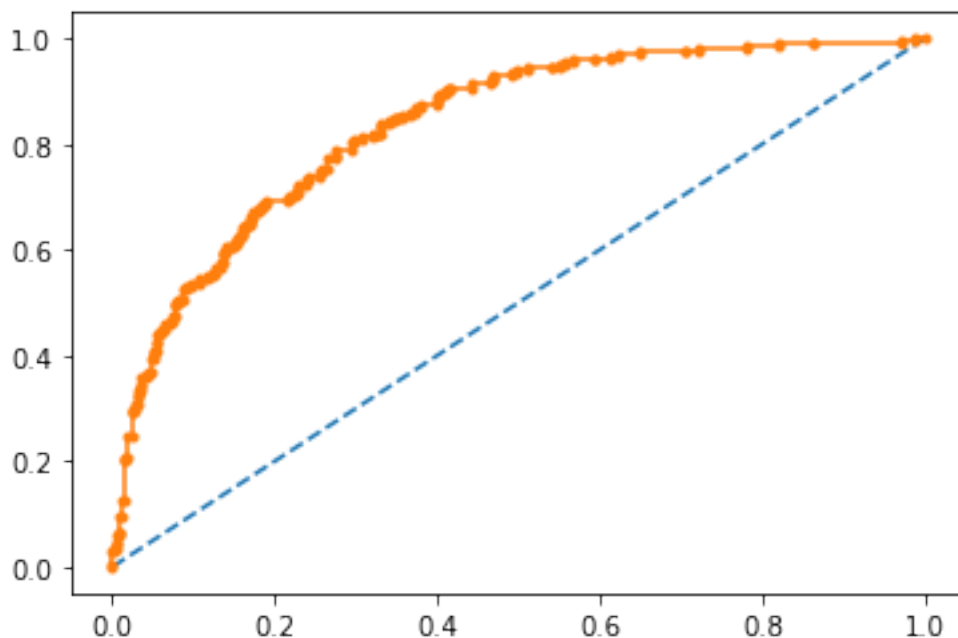
```
print('AUC: %.3f' % auc)
```

```
AUC: 0.837
```

```
fpr, tpr, thresholds = roc_curve(label, probability)      # ROC curve
calculation
```

```
plt.plot([0, 1], [0, 1], linestyle='--')
plt.plot(fpr, tpr, marker='.')
```

```
[<matplotlib.lines.Line2D at 0x7fc5dc5f20a0>]
```



```
from sklearn.metrics import classification_report
print(classification_report)
```

```
<function classification_report at 0x7fc5dc29f940>
```

```
print(classifier.score(X_train,y_train))
print(classifier.score(X_test,y_test))
```

```
0.7719869706840391
```

```
0.7662337662337663
```

```
classification_report
```

```
<function
sklearn.metrics._classification.classification_report(y_true, y_pred,
*, labels=None, target_names=None, sample_weight=None, digits=2,
output_dict=False, zero_division='warn')>
```

```
classifier.predict(X_test)
```

```
array([1, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 1, 0, 0, 1, 1, 0, 0, 1, 0,
0,
      0, 0, 0, 1, 0, 0, 1, 1, 0, 0, 0, 0, 0, 0, 0, 0, 1, 0, 0, 0, 0,
0,
      0, 0, 0, 0, 0, 0, 0, 1, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0,
0,
      1, 0, 0, 1, 1, 0, 0, 1, 0, 0, 1, 1, 0, 0, 0, 1, 0, 0, 0, 1, 0,
0,
      1, 0, 0, 0, 0, 0, 0, 0, 1, 0, 0, 0, 1, 0, 1, 0, 0, 0, 1, 1, 0,
1,
      0, 1, 0, 0, 0, 0, 0, 0, 0, 1, 0, 0, 0, 1, 0, 0, 0, 0, 0, 0, 1,
0,
      0,
```

```
1, 0, 0, 0, 0, 1, 1, 1, 0, 0, 0, 0, 1, 0, 0, 1, 1, 0, 0, 0, 1,  
0])
```

```
from sklearn.metrics import classification_report  
print(classification_report)
```

```
<function classification_report at 0x7fc5dc29f940>
```