

```
In [1]: 1 import pandas as pd
        2 import numpy as np
        3 import seaborn as sns
        4 import matplotlib.pyplot as plt
        5 %matplotlib inline
```

```
In [2]: 1 df=pd.read_csv('diabetes.csv')
        2 df.head()
```

Out[2]:

	Pregnancies	Glucose	BloodPressure	SkinThickness	Insulin	BMI	DiabetesPedigreeFunction	Age
0	2	138	62	35	0	33.6	0.127	33
1	0	84	82	31	125	38.2	0.233	33
2	0	145	0	0	0	44.2	0.630	33
3	0	135	68	42	250	42.3	0.365	33
4	1	139	62	41	480	40.7	0.536	33

```
In [3]: 1 #lets describe the data
        2 df.describe()
```

Out[3]:

	Pregnancies	Glucose	BloodPressure	SkinThickness	Insulin	BMI	DiabetesPedigreeFunction	Age
count	2000.000000	2000.000000	2000.000000	2000.000000	2000.000000	2000.000000	2000.000000	2000.000000
mean	3.703500	121.182500	69.145500	20.935000	80.254000	32.193000	0.471000	33.430000
std	3.306063	32.068636	19.188315	16.103243	111.180534	8.149901	0.374605	4.167419
min	0.000000	0.000000	0.000000	0.000000	0.000000	0.000000	0.078000	21.000000
25%	1.000000	99.000000	63.500000	0.000000	0.000000	27.375000	0.243000	24.000000
50%	3.000000	117.000000	72.000000	23.000000	40.000000	32.300000	0.366000	33.000000
75%	6.000000	141.000000	80.000000	32.000000	130.000000	36.800000	0.625000	36.000000
max	17.000000	199.000000	122.000000	110.000000	744.000000	80.600000	0.812000	81.000000

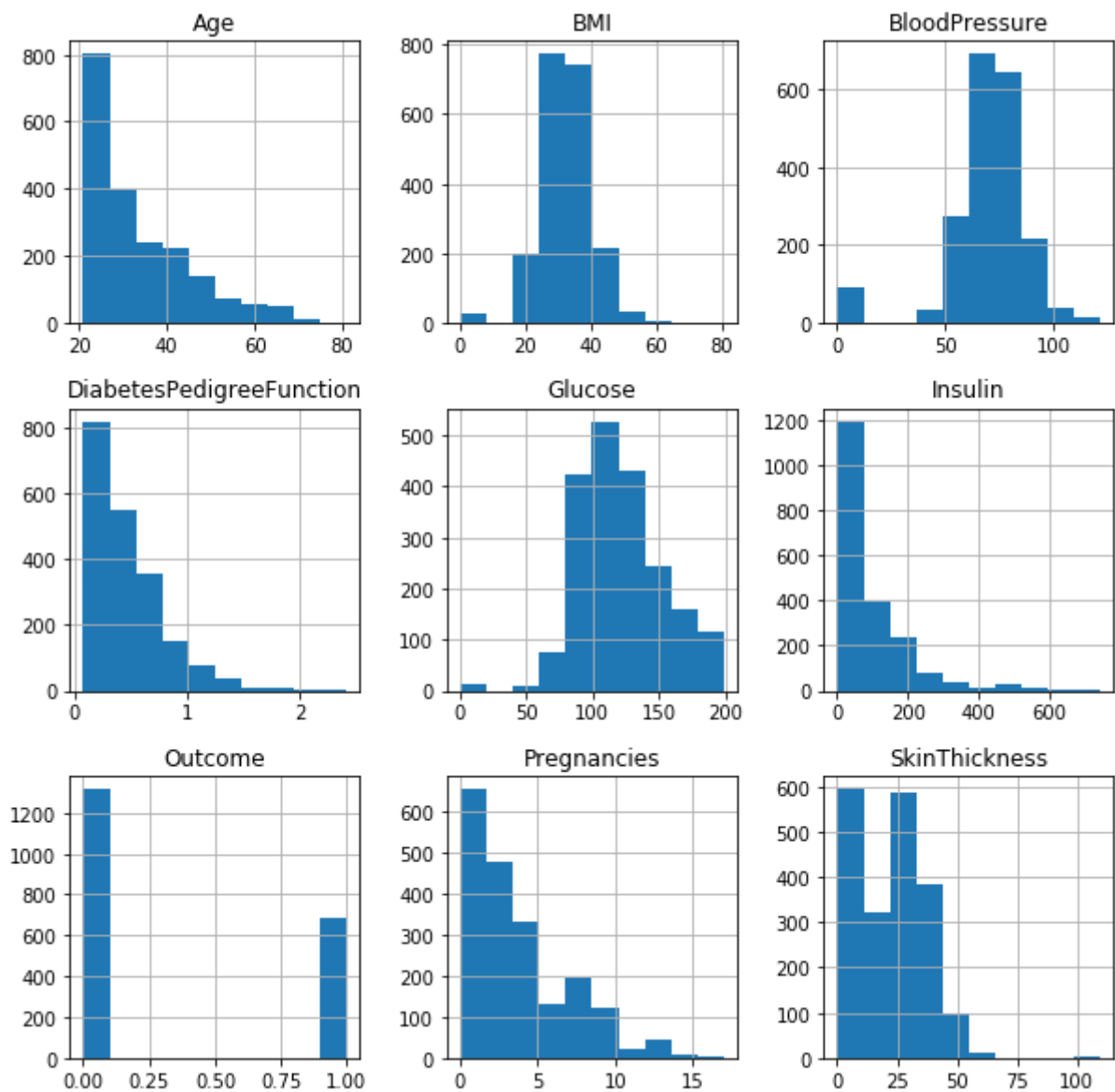
```
In [4]: 1 #information of dataset
        2 df.info()
```

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 2000 entries, 0 to 1999
Data columns (total 9 columns):
#   Column                                Non-Null Count  Dtype
---  -
0   Pregnancies                          2000 non-null   int64
1   Glucose                              2000 non-null   int64
2   BloodPressure                        2000 non-null   int64
3   SkinThickness                        2000 non-null   int64
4   Insulin                             2000 non-null   int64
5   BMI                                  2000 non-null   float64
6   DiabetesPedigreeFunction             2000 non-null   float64
7   Age                                  2000 non-null   int64
8   Outcome                              2000 non-null   int64
dtypes: float64(2), int64(7)
memory usage: 140.8 KB
```

```
In [5]: 1 #any null values
        2 #not necessary in above information we can see
        3 df.isnull().values.any()
```

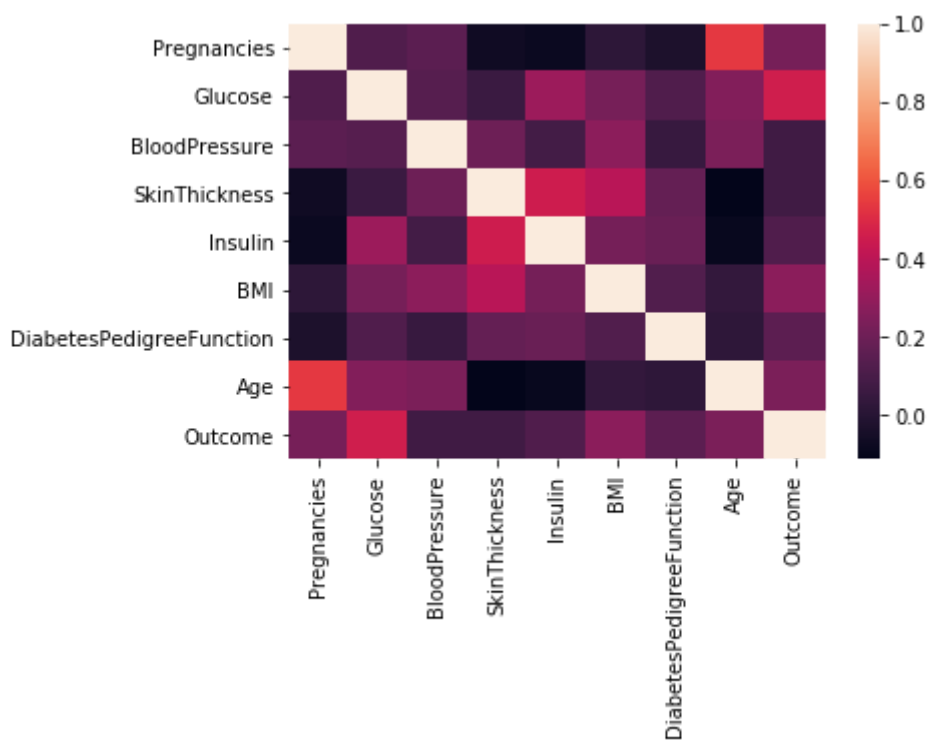
```
Out[5]: False
```

```
In [6]: 1 #histogram
2 df.hist(bins=10,figsize=(10,10))
3 plt.show()
```



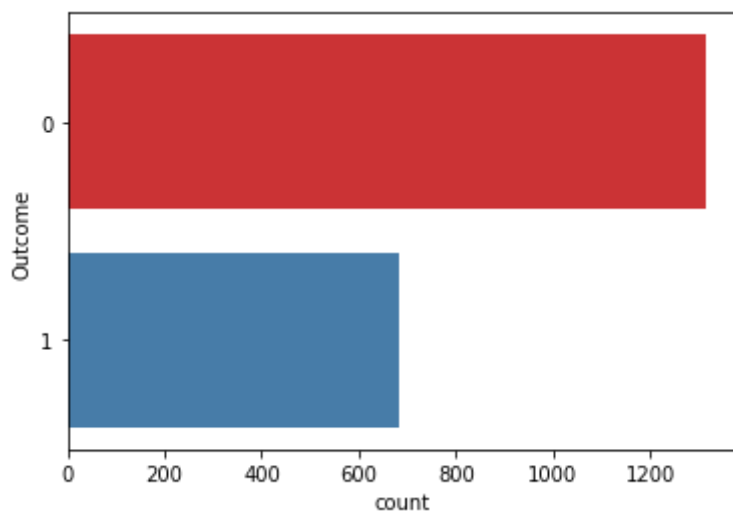
```
In [7]: 1 #correlation
        2
        3 sns.heatmap(df.corr())
        4 # we can see skin thickness,insulin,pregnencies and age are full indepe
        5 #age and pregencies has negative correlation
```

Out[7]: <matplotlib.axes.\_subplots.AxesSubplot at 0x7fd261b32e50>



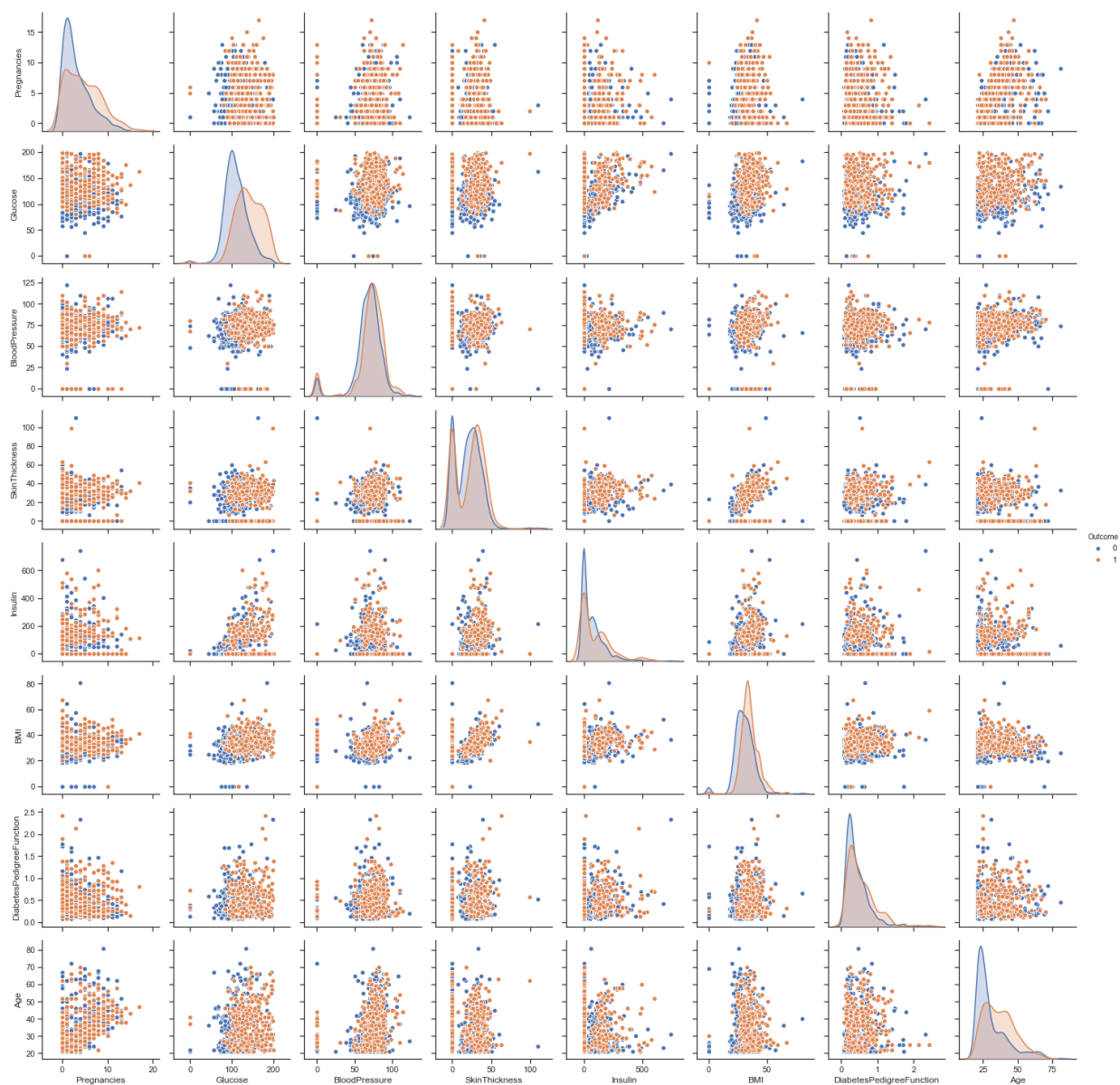
```
In [8]: 1 #lets count total outcome in each target 0 1  
2 #0 means no diabeted  
3 #1 means patient with diabtes  
4 sns.countplot(y=df['Outcome'],palette='Set1')
```

Out[8]: <matplotlib.axes.\_subplots.AxesSubplot at 0x7fd261bd7390>



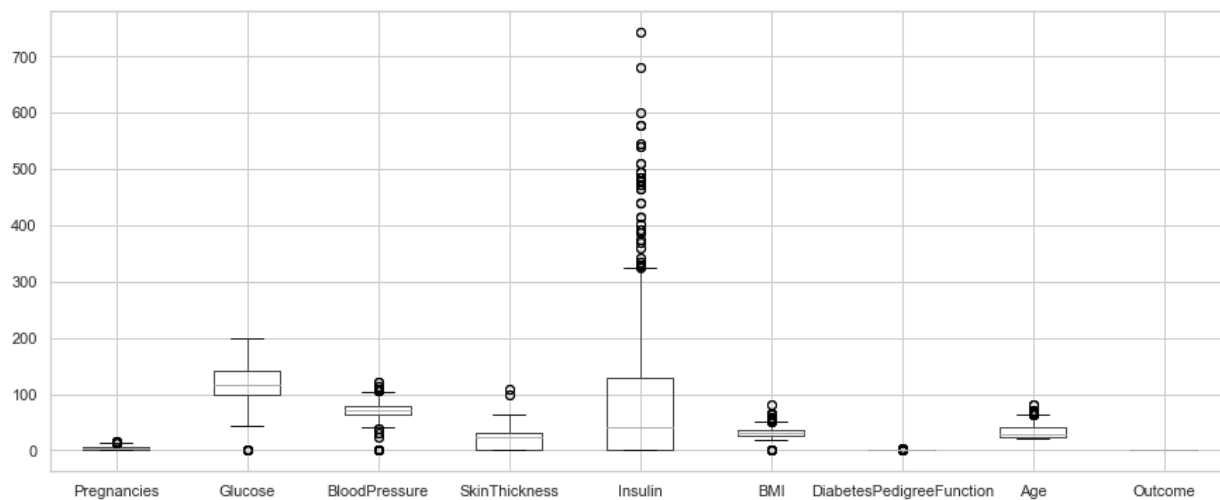
```
In [9]: 1 sns.set(style="ticks")  
2 sns.pairplot(df, hue="Outcome")
```

Out[9]: <seaborn.axisgrid.PairGrid at 0x7fd261c58c10>

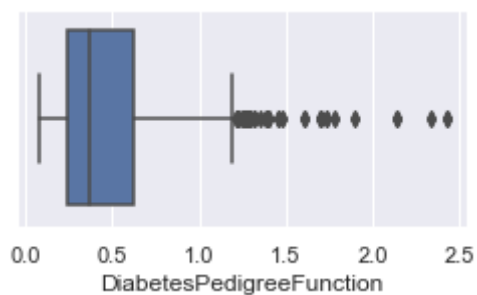
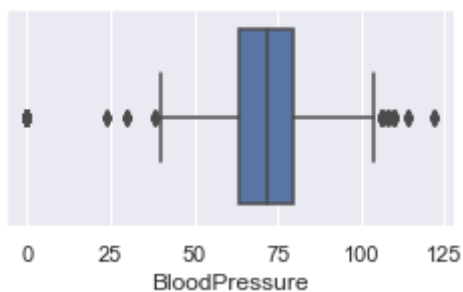
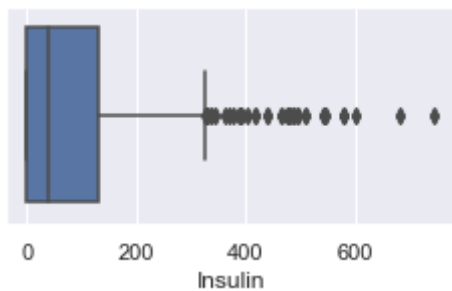


```
In [10]: 1 #box plot for outlier visualization  
2 sns.set(style="whitegrid")  
3 df.boxplot(figsize=(15,6))
```

Out[10]: <matplotlib.axes.\_subplots.AxesSubplot at 0x7fd2641fcf50>



```
In [11]: 1 #box plot
2 sns.set(style="whitegrid")
3
4 sns.set(rc={'figure.figsize':(4,2)})
5 sns.boxplot(x=df['Insulin'])
6 plt.show()
7 sns.boxplot(x=df['BloodPressure'])
8 plt.show()
9 sns.boxplot(x=df['DiabetesPedigreeFunction'])
10 plt.show()
```





```
In [12]: 1 #outlier remove
2
3 Q1=df.quantile(0.25)
4 Q3=df.quantile(0.75)
5 IQR=Q3-Q1
6
7 print("---Q1--- \n",Q1)
8 print("\n---Q3--- \n",Q3)
9 print("\n---IQR---\n",IQR)
10
11 #print((df < (Q1 - 1.5 * IQR))|(df > (Q3 + 1.5 * IQR)))
12
```

```
---Q1---
Pregnancies          1.000
Glucose              99.000
BloodPressure        63.500
SkinThickness         0.000
Insulin               0.000
BMI                  27.375
DiabetesPedigreeFunction 0.244
Age                  24.000
Outcome              0.000
Name: 0.25, dtype: float64
```

```
---Q3---
Pregnancies          6.000
Glucose             141.000
BloodPressure        80.000
SkinThickness        32.000
Insulin              130.000
BMI                  36.800
DiabetesPedigreeFunction 0.624
Age                  40.000
Outcome              1.000
Name: 0.75, dtype: float64
```

```
---IQR---
Pregnancies          5.000
Glucose              42.000
BloodPressure        16.500
SkinThickness        32.000
Insulin              130.000
BMI                   9.425
DiabetesPedigreeFunction 0.380
Age                  16.000
Outcome              1.000
dtype: float64
```

```
In [13]: 1 #outlier remove
2 df_out = df[-((df < (Q1 - 1.5 * IQR)) |(df > (Q3 + 1.5 * IQR))).any(axis=1)]
3 df.shape,df_out.shape
4 #more than 80 records deleted
```

```
Out[13]: ((2000, 9), (1652, 9))
```

```
In [14]: 1 #Scatter matrix after removing outlier
2 sns.set(style="ticks")
3 sns.pairplot(df_out, hue="Outcome")
4 plt.show()
```



```
In [15]: 1 #lets extract features and targets
2 X=df_out.drop(columns=[ 'Outcome' ])
3 y=df_out[ 'Outcome' ]
```

```
In [16]: 1 #Splitting train test data 80 20 ratio
2 from sklearn.model_selection import train_test_split
3 train_X,test_X,train_y,test_y=train_test_split(X,y,test_size=0.2)
```

```
In [17]: 1 train_X.shape,test_X.shape,train_y.shape,test_y.shape
```

```
Out[17]: ((1321, 8), (331, 8), (1321,), (331,))
```

```
In [18]: 1 from sklearn.metrics import confusion_matrix, accuracy_score, make_scorer
2 from sklearn.model_selection import cross_validate
3
4 def tn(y_true, y_pred): return confusion_matrix(y_true, y_pred)[0, 0]
5 def fp(y_true, y_pred): return confusion_matrix(y_true, y_pred)[0, 1]
6 def fn(y_true, y_pred): return confusion_matrix(y_true, y_pred)[1, 0]
7 def tp(y_true, y_pred): return confusion_matrix(y_true, y_pred)[1, 1]
8
9 #cross validation purpose
10 scoring = {'accuracy': make_scorer(accuracy_score), 'prec': 'precision'}
11 scoring = {'tp': make_scorer(tp), 'tn': make_scorer(tn),
12           'fp': make_scorer(fp), 'fn': make_scorer(fn)}
13
14 def display_result(result):
15     print("TP: ", result['test_tp'])
16     print("TN: ", result['test_tn'])
17     print("FN: ", result['test_fn'])
18     print("FP: ", result['test_fp'])
```

```

In [19]: 1  #Lets build the model
          2
          3  #Logistic Regression
          4  from sklearn.linear_model import LogisticRegression
          5  from sklearn.metrics import roc_auc_score
          6
          7  acc=[]
          8  roc=[]
          9
         10  clf=LogisticRegression()
         11  clf.fit(train_X,train_y)
         12  y_pred=clf.predict(test_X)
         13  #find accuracy
         14  ac=accuracy_score(test_y,y_pred)
         15  acc.append(ac)
         16
         17  #find the ROC_AOC curve
         18  rc=roc_auc_score(test_y,y_pred)
         19  roc.append(rc)
         20  print("\nAccuracy {0} ROC {1}".format(ac,rc))
         21
         22  #cross val score
         23  result=cross_validate(clf,train_X,train_y,scoring=scoring,cv=10)
         24  display_result(result)
         25
         26  #display predicted values uncomment below line
         27  #pd.DataFrame(data={'Actual':test_y,'Predicted':y_pred}).head()

```

Please also refer to the documentation for alternative solver options:

[https://scikit-learn.org/stable/modules/linear\\_model.html#logistic-regression](https://scikit-learn.org/stable/modules/linear_model.html#logistic-regression) ([https://scikit-learn.org/stable/modules/linear\\_model.html#logistic-regression](https://scikit-learn.org/stable/modules/linear_model.html#logistic-regression))

```

extra_warning_msg=_LOGISTIC_SOLVER_CONVERGENCE_MSG)
/Users/roshanremanan/opt/anaconda3/lib/python3.7/site-packages/sklearn/linear_model/_logistic.py:940: ConvergenceWarning: lbfgs failed to converge
(status=1):
STOP: TOTAL NO. of ITERATIONS REACHED LIMIT.

```

Increase the number of iterations (max\_iter) or scale the data as shown in:

<https://scikit-learn.org/stable/modules/preprocessing.html> (<https://scikit-learn.org/stable/modules/preprocessing.html>)

Please also refer to the documentation for alternative solver options:

[https://scikit-learn.org/stable/modules/linear\\_model.html#logistic-regression](https://scikit-learn.org/stable/modules/linear_model.html#logistic-regression) ([https://scikit-learn.org/stable/modules/linear\\_model.html#logistic-regression](https://scikit-learn.org/stable/modules/linear_model.html#logistic-regression))

```

extra_warning_msg=_LOGISTIC_SOLVER_CONVERGENCE_MSG)
/Users/roshanremanan/opt/anaconda3/lib/python3.7/site-packages/sklearn/linear_model/_logistic.py:940: ConvergenceWarning: lbfgs failed to converge
(status=1):
STOP: TOTAL NO. of ITERATIONS REACHED LIMIT.

```

```
In [20]: 1 #Support Vector Machine
2 from sklearn.svm import SVC
3
4 clf=SVC(kernel='linear')
5 clf.fit(train_X,train_y)
6 y_pred=clf.predict(test_X)
7 #find accuracy
8 ac=accuracy_score(test_y,y_pred)
9 acc.append(ac)
10
11 #find the ROC_AOC curve
12 rc=roc_auc_score(test_y,y_pred)
13 roc.append(rc)
14 print("Accuracy {0} ROC {1}".format(ac,rc))
15
16 #cross val score
17 result=cross_validate(clf,train_X,train_y,scoring=scoring,cv=10)
18 display_result(result)
19
20 #display predicted values uncomment below line
21 #pd.DataFrame(data={'Actual':test_y,'Predicted':y_pred}).head()
```

```
Accuracy 0.8187311178247734 ROC 0.7394548063127689
TP:  [26 25 21 22 18 20 20 23 24 26]
TN:  [86 80 80 79 77 76 81 81 81 87]
FN:  [18 18 22 21 25 23 23 20 19 17]
FP:  [ 3  9  9 10 12 13  8  8  8  2]
```

```
In [21]: 1 #KNN
2
3 from sklearn.neighbors import KNeighborsClassifier
4
5 clf=KNeighborsClassifier(n_neighbors=3)
6 clf.fit(train_X,train_y)
7 y_pred=clf.predict(test_X)
8 #find accuracy
9 ac=accuracy_score(test_y,y_pred)
10 acc.append(ac)
11
12 #find the ROC_AOC curve
13 rc=roc_auc_score(test_y,y_pred)
14 roc.append(rc)
15 print("\nAccuracy {0} ROC {1}".format(ac,rc))
16
17 #cross val score
18 result=cross_validate(clf,train_X,train_y,scoring=scoring,cv=10)
19 display_result(result)
20
21 #display predicted values uncomment below line
22 #pd.DataFrame(data={'Actual':test_y,'Predicted':y_pred}).head()
```

```
Accuracy 0.8972809667673716 ROC 0.9039454806312768
TP:  [34 33 34 37 29 29 30 29 30 36]
TN:  [87 81 79 76 79 80 78 81 85 82]
FN:  [10 10  9  6 14 14 13 14 13  7]
FP:  [ 2  8 10 13 10  9 11  8  4  7]
```

```
In [22]: 1 #Random forest
2 from sklearn.ensemble import RandomForestClassifier
3
4 clf=RandomForestClassifier()
5 clf.fit(train_X,train_y)
6
7 y_pred=clf.predict(test_X)
8 #find accuracy
9 ac=accuracy_score(test_y,y_pred)
10 acc.append(ac)
11
12 #find the ROC_AOC curve
13 rc=roc_auc_score(test_y,y_pred)
14 roc.append(rc)
15 print("\nAccuracy {0} ROC {1}".format(ac,rc))
16
17 #cross val score
18 result=cross_validate(clf,train_X,train_y,scoring=scoring,cv=10)
19 display_result(result)
20
21 #display predicted values uncomment below line
22 #pd.DataFrame(data={'Actual':test_y,'Predicted':y_pred}).head()
```

Accuracy 0.972809667673716 ROC 0.9817073170731707

TP: [43 42 42 40 42 41 43 41 38 41]

TN: [88 89 86 84 85 88 89 87 89 87]

FN: [1 1 1 3 1 2 0 2 5 2]

FP: [1 0 3 5 4 1 0 2 0 2]

```
In [23]: 1 #Naive Bayes Theorem
2 #import library
3 from sklearn.naive_bayes import GaussianNB
4
5 clf=GaussianNB()
6 clf.fit(train_X,train_y)
7 y_pred=clf.predict(test_X)
8 #find accuracy
9 ac=accuracy_score(test_y,y_pred)
10 acc.append(ac)
11
12 #find the ROC_AOC curve
13 rc=roc_auc_score(test_y,y_pred)
14 roc.append(rc)
15 print("\nAccuracy {0} ROC {1}".format(ac,rc))
16
17 #cross val score
18 result=cross_validate(clf,train_X,train_y,scoring=scoring,cv=10)
19 display_result(result)
20
21 #display predicted values uncomment below line
22 #pd.DataFrame(data={'Actual':test_y,'Predicted':y_pred}).head()
```

Accuracy 0.7764350453172205 ROC 0.7379483500717361

TP: [25 28 29 24 23 25 25 28 27 27]

TN: [78 73 75 76 71 75 76 77 76 79]

FN: [19 15 14 19 20 18 18 15 16 16]

FP: [11 16 14 13 18 14 13 12 13 10]



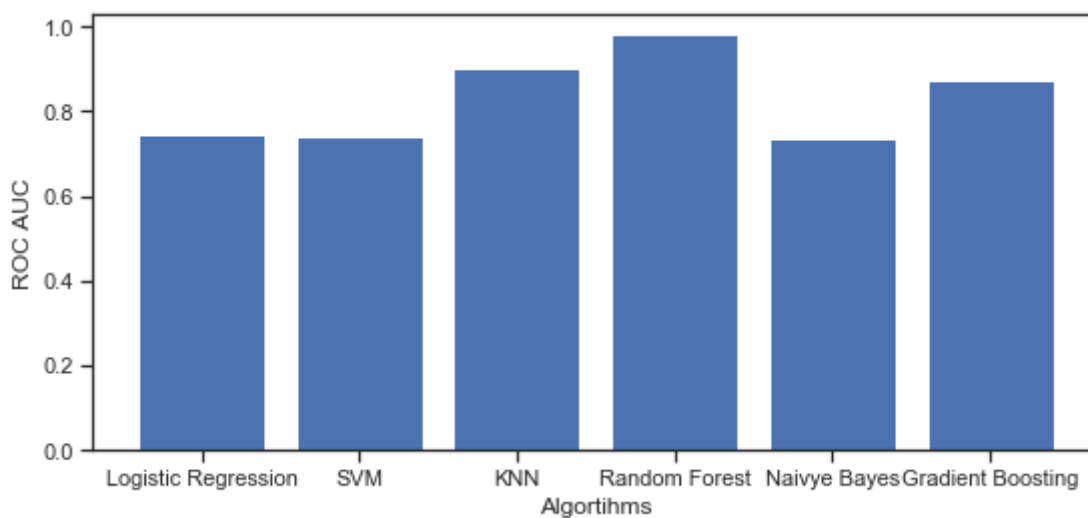
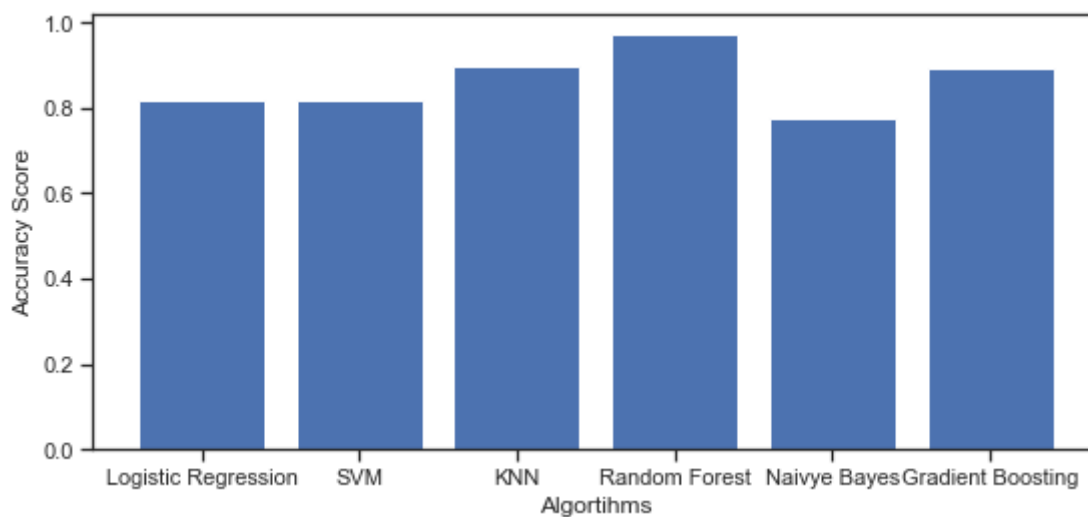
```
In [24]: 1 #Gradient Boosting Classifier
2 from sklearn.ensemble import GradientBoostingClassifier
3 clf=GradientBoostingClassifier(n_estimators=50,learning_rate=0.2)
4 clf.fit(train_X,train_y)
5 y_pred=clf.predict(test_X)
6 #find accuracy
7 ac=accuracy_score(test_y,y_pred)
8 acc.append(ac)
9
10 #find the ROC_AOC curve
11 rc=roc_auc_score(test_y,y_pred)
12 roc.append(rc)
13 print("\nAccuracy {0} ROC {1}".format(ac,rc))
14
15 #cross val score
16 result=cross_validate(clf,train_X,train_y,scoring=scoring,cv=10)
17 display_result(result)
18
19 #display predicted values uncomment below line
20 #pd.DataFrame(data={'Actual':test_y,'Predicted':y_pred}).head()
```

```
Accuracy 0.8912386706948641 ROC 0.8729316116690579
TP:  [35 36 30 33 37 29 31 33 36 33]
TN:  [86 89 83 81 83 83 81 87 84 82]
FN:  [ 9  7 13 10  6 14 12 10  7 10]
FP:  [3 0 6 8 6 6 8 2 5 7]
```

```

In [25]: 1 #lets plot the bar graph
2
3 ax=plt.figure(figsize=(9,4))
4 plt.bar(['Logistic Regression','SVM','KNN','Random Forest','Naivye Baye
5 plt.ylabel('Accuracy Score')
6 plt.xlabel('Algortihms')
7 plt.show()
8
9 ax=plt.figure(figsize=(9,4))
10 plt.bar(['Logistic Regression','SVM','KNN','Random Forest','Naivye Baye
11 plt.ylabel('ROC AUC')
12 plt.xlabel('Algortihms')
13 plt.show()

```



```

In [26]: 1 #Great....
2 #Random forest has highest accuracy 98% and ROC_AUC curve 97%
3 #model can be improve more if we take same count of labels
4 #in our model 30% is diabetic and 70% no diabetic patient
5
6 #model can be improve with fine tunning

```

```

In [ ]: 1

```

In [ ]:

1