```python
import pandas as pd

import numpy as np

import seaborn as sns

import matplotlib.pyplot as plt

%matplotlib inline

df=pd.read_csv('diabetes.csv')

df.head()

#lets describe the data

df.describe()

#infromation of dataset

df.info()

#any null values

#not neccessary in above information we can see

df.isnull().values.any()

#histogram

df.hist(bins=10,figsize=(10,10))

plt.show()

#correlation

sns.heatmap(df.corr())

# we can see skin thickness,insulin,pregnencies and age are full independent to each other

#age and pregencies has negative correlation

#lets count total outcome in each target 0 1

#0 means no diabeted

#1 means patient with diabtes

sns.countplot(y=df['Outcome'],palette='Set1')

sns.set(style="ticks")

sns.pairplot(df, hue="Outcome")

#box plot for outlier visualization

sns.set(style="whitegrid")

df.boxplot(figsize=(15,6))

#box plot
```

```python
sns.set(style="whitegrid")

sns.set(rc={'figure.figsize':(4,2)})

sns.boxplot(x=df['Insulin'])

plt.show()

sns.boxplot(x=df['BloodPressure'])

plt.show()

sns.boxplot(x=df['DiabetesPedigreeFunction'])

plt.show()

#outlier remove


Q1=df.quantile(0.25)

Q3=df.quantile(0.75)

IQR=Q3-Q1


print("---Q1--- \n",Q1)

print("\n---Q3--- \n",Q3)

print("\n---IQR---\n",IQR)


#print((df < (Q1 - 1.5 * IQR))|(df > (Q3 + 1.5 * IQR)))

#outlier remove

df_out = df[~((df < (Q1 - 1.5 * IQR)) |(df > (Q3 + 1.5 * IQR))).any(axis=1)]

df.shape,df_out.shape

#more than 80 records deleted

#Scatter matrix after removing outlier

sns.set(style="ticks")

sns.pairplot(df_out, hue="Outcome")

plt.show()

#lets extract features and targets

X=df_out.drop(columns=['Outcome'])

y=df_out['Outcome']

#Splitting train test data 80 20 ratio
```

```python
from sklearn.model_selection import train_test_split

train_X,test_X,train_y,test_y=train_test_split(X,y,test_size=0.2)

train_X.shape,test_X.shape,train_y.shape,test_y.shape

from sklearn.metrics import confusion_matrix,accuracy_score,make_scorer

from sklearn.model_selection import cross_validate


def tn(y_true, y_pred): return confusion_matrix(y_true, y_pred)[0, 0]

def fp(y_true, y_pred): return confusion_matrix(y_true, y_pred)[0, 1]

def fn(y_true, y_pred): return confusion_matrix(y_true, y_pred)[1, 0]

def tp(y_true, y_pred): return confusion_matrix(y_true, y_pred)[1, 1]


#cross validation purpose

scoring = {'accuracy': make_scorer(accuracy_score),'prec': 'precision'}

scoring = {'tp': make_scorer(tp), 'tn': make_scorer(tn),

        'fp': make_scorer(fp), 'fn': make_scorer(fn)}


def display_result(result):

    print("TP: ",result['test_tp'])

    print("TN: ",result['test_tn'])

    print("FN: ",result['test_fn'])

    print("FP: ",result['test_fp'])

#Lets build the model


#Logistic Regression

from sklearn.linear_model import LogisticRegression

from sklearn.metrics import roc_auc_score


acc=[]

roc=[]


clf=LogisticRegression()
```

```python
clf.fit(train_X,train_y)

y_pred=clf.predict(test_X)

#find accuracy

ac=accuracy_score(test_y,y_pred)

acc.append(ac)


#find the ROC_AOC curve

rc=roc_auc_score(test_y,y_pred)

roc.append(rc)

print("\nAccuracy {0} ROC {1}".format(ac,rc))


#cross val score

result=cross_validate(clf,train_X,train_y,scoring=scoring,cv=10)

display_result(result)


#display predicted values uncomment below line

#pd.DataFrame(data={'Actual':test_y,'Predicted':y_pred}).head()

#Support Vector Machine

from sklearn.svm import SVC


clf=SVC(kernel='linear')

clf.fit(train_X,train_y)

y_pred=clf.predict(test_X)

#find accuracy

ac=accuracy_score(test_y,y_pred)

acc.append(ac)


#find the ROC_AOC curve

rc=roc_auc_score(test_y,y_pred)

roc.append(rc)

print("\nAccuracy {0} ROC {1}".format(ac,rc))
```

```python
#cross val score
result=cross_validate(clf,train_X,train_y,scoring=scoring,cv=10)
display_result(result)


#display predicted values uncomment below line
#pd.DataFrame(data={'Actual':test_y,'Predicted':y_pred}).head()
#KNN


from sklearn.neighbors import KNeighborsClassifier


clf=KNeighborsClassifier(n_neighbors=3)
clf.fit(train_X,train_y)
y_pred=clf.predict(test_X)
#find accuracy
ac=accuracy_score(test_y,y_pred)
acc.append(ac)


#find the ROC_AOC curve
rc=roc_auc_score(test_y,y_pred)
roc.append(rc)
print("\nAccuracy {0} ROC {1}".format(ac,rc))


#cross val score
result=cross_validate(clf,train_X,train_y,scoring=scoring,cv=10)
display_result(result)


#display predicted values uncomment below line
#pd.DataFrame(data={'Actual':test_y,'Predicted':y_pred}).head()
#Random forest
from sklearn.ensemble import RandomForestClassifier
```

```python
clf=RandomForestClassifier()

clf.fit(train_X,train_y)


y_pred=clf.predict(test_X)

#find accuracy

ac=accuracy_score(test_y,y_pred)

acc.append(ac)


#find the ROC_AOC curve

rc=roc_auc_score(test_y,y_pred)

roc.append(rc)

print("\nAccuracy {0} ROC {1}".format(ac,rc))


#cross val score

result=cross_validate(clf,train_X,train_y,scoring=scoring,cv=10)

display_result(result)


#display predicted values uncomment below line

#pd.DataFrame(data={'Actual':test_y,'Predicted':y_pred}).head()

#Naive Bayes Theorem

#import library

from sklearn.naive_bayes import GaussianNB


clf=GaussianNB()

clf.fit(train_X,train_y)

y_pred=clf.predict(test_X)

#find accuracy

ac=accuracy_score(test_y,y_pred)

acc.append(ac)
```

```python
#find the ROC_AOC curve
rc=roc_auc_score(test_y,y_pred)
roc.append(rc)
print("\nAccuracy {0} ROC {1}".format(ac,rc))


#cross val score
result=cross_validate(clf,train_X,train_y,scoring=scoring,cv=10)
display_result(result)


#display predicted values uncomment below line
#pd.DataFrame(data={'Actual':test_y,'Predicted':y_pred}).head()
#Gradient Boosting Classifier
from sklearn.ensemble import GradientBoostingClassifier
clf=GradientBoostingClassifier(n_estimators=50,learning_rate=0.2)
clf.fit(train_X,train_y)
y_pred=clf.predict(test_X)
#find accuracy
ac=accuracy_score(test_y,y_pred)
acc.append(ac)


#find the ROC_AOC curve
rc=roc_auc_score(test_y,y_pred)
roc.append(rc)
print("\nAccuracy {0} ROC {1}".format(ac,rc))


#cross val score
result=cross_validate(clf,train_X,train_y,scoring=scoring,cv=10)
display_result(result)


#display predicted values uncomment below line
#pd.DataFrame(data={'Actual':test_y,'Predicted':y_pred}).head()
```

```python
#lets plot the bar graph


ax=plt.figure(figsize=(9,4))

plt.bar(['Logistic Regression','SVM','KNN','Random Forest','Naivye Bayes','Gradient
Boosting'],acc,label='Accuracy')

plt.ylabel('Accuracy Score')

plt.xlabel('Algortihms')

plt.show()


ax=plt.figure(figsize=(9,4))

plt.bar(['Logistic Regression','SVM','KNN','Random Forest','Naivye Bayes','Gradient
Boosting'],roc,label='ROC AUC')

plt.ylabel('ROC AUC')

plt.xlabel('Algortihms')

plt.show()

#Great....

#Random forest has highest accuracy 98% and ROC_AUC curve 97%

#model can be improve more if we take same count of labels

#in our model 30% is diabetic and 70% no diabetic patient


#model can be improve with fine tunning
```