```
Python 3
       View
                            Kernel Help
                                                                                                                Not Trusted
 Edit
              Insert
                      Cell
                      ► Run ■ C → Code
                                                    ~
 In [1]: import numpy as np
         import pandas as pd
         import matplotlib.pyplot as plt
         import seaborn as sns
         %matplotlib inline
         from IPython.display import display
         from sklearn.preprocessing import LabelEncoder, MinMaxScaler, StandardScaler
         import pickle
 In [2]: df = pd.read_csv('heart.csv')
         print('NO OF ROWS AND COLUMN IN DATASET ',df.shape,'\n')
         display(df.head(),"\n",df.dtypes)
         NO OF ROWS AND COLUMN IN DATASET (303, 14)
            age sex cp trestbps chol fbs restecg thalach exang oldpeak slope ca thal target
          0 63 1 3
                           145 233 1
                                                 150
                                                                     0 0 1
                                                               2.3
                                                         0
          1 37
                           130 250
                                                 187
                                                               3.5
                                                                     0 0
                1 2
                                    0
                           130 204 0
          2 41 0 1
                                                 172
                                                                     2 0 2
                           120 236
                                                 178
                                                                     2 0 2
                                   0
                 1 1
          4 57 0 0
                           120 354 0
                                                                     2 0 2
                                                 163
                                                               0.6
         '\n'
                      int64
         age
                      int64
         sex
                      int64
                      int64
         trestbps
                       int64
         chol
                      int64
         fbs
                      int64
         restecg
         thalach
                      int64
                      int64
         exang
         oldpeak
                     float64
         slope
                      int64
                      int64
         ca
         thal
                      int64
                      int64
         target
         dtype: object
 In [3]: print('COUNT NO OF NULL IN EACH COLUMN','\n',df.isnull().sum())
         COUNT NO OF NULL IN EACH COLUMN
          age
                     0
         sex
         trestbps
         chol
         fbs
         restecg
         thalach
         exang
         oldpeak
         slope
         ca
         thal
         target
         dtype: int64
 In [4]: df.describe().T
 Out[4]:
                                                    50% 75% max
                                          min 25%
                                     std
                 count
                           mean
                                                         61.0 77.0
                        54.366337
                                 9.082101
                                          29.0 47.5
             age 303.0
                        0.683168
                                 0.466011
                                          0.0
                 303.0
                                                0.0
                                                     1.0
                                                          1.0
                                                               1.0
              cp 303.0
                        0.966997
                                 1.032052
                                          0.0
                                               0.0
                                                   1.0 2.0
                 303.0 131.623762 17.538143
                                          94.0 120.0 130.0 140.0 200.0
                 303.0 246.264026 51.830751 126.0 211.0 240.0 274.5 564.0
                        0.148515 0.356198
             fbs 303.0
                                          0.0
                                                0.0
          restecg 303.0 0.528053 0.525860 0.0 0.0 1.0 1.0 2.0
                 303.0 149.646865 22.905161
                                          71.0 133.5 153.0 166.0 202.0
           exang 303.0
                                 0.469794
                        0.326733
                                           0.0
                                                0.0
                                                     0.0
                                                          1.0
                                                               1.0
                                 1.161075
          oldpeak 303.0
                         1.039604
                                           0.0
                                                0.0
                                                                6.2
            slope 303.0
                         1.399340
                                 0.616226
                                           0.0
                                                1.0
                                                     1.0
                                                          2.0
                                                                2.0
                                 1.022606
                        0.729373
                                          0.0
                                                0.0
                                                     0.0
              ca 303.0
                                                          1.0
                                                                4.0
             thal 303.0
                                 0.612277
                        2.313531
                                                2.0
            target 303.0
                        0.544554 0.498835
                                          0.0
                                                0.0
                                                     1.0
                                                          1.0
                                                               1.0
 In [5]: data = df.copy(deep = True)
         ##data[['age','sex','cp','trestbps','chol','fbs','restecg','thalch','exang','oldpeak','slope','ca','thal']] = data[['age','sex',
         ## showing the count of Nans
         print(data.isnull().sum())
                     0
         age
         sex
         trestbps
         chol
         fbs
         restecg
         thalach
         exang
         oldpeak
         slope
         ca
         thal
         target
         dtype: int64
 In [6]: corr=data.corr()
         corr.nlargest(15, 'target')['target']
 Out[6]: target
                    1.000000
                     0.433798
         thalach
                     0.421741
                     0.345877
         slope
                    0.137230
         restecg
         fbs
                    -0.028046
                    -0.085239
         chol
         trestbps
                   -0.144931
                    -0.225439
         age
                   -0.280937
         sex
         thal
                   -0.344029
                    -0.391724
         ca
                   -0.430696
         oldpeak
                   -0.436757
         exang
         Name: target, dtype: float64
         age,cp=cerebral palsy,chol=cholesterol,fbs=fasting bloodpressure,
         restecg= resting electrocardiographic,
         thalach= maximum heart rate achived, exang=exercise include angina, slope.
 In [7]: x = data[['age','cp','chol','fbs','restecg','thalach','exang','slope']]
         y = data['target']
         x.head()
 Out[7]:
            age cp chol fbs restecg thalach exang slope
          0 63 3 233 1
          1 37 2 250 0
                                     187
          2 41 1 204 0
                                0 172
          3 56 1 236 0
          4 57 0 354 0
                                     163
                             1
                                           1 2
 In [8]: from sklearn.model_selection import train_test_split
         x_train, x_test, y_train, y_test = train_test_split(x, y, test_size=0.2, random_state=0)
         print('X TRAIN DATA ', x_train.shape)
         print('X TEST DATA ', x_test.shape)
         print('Y TRAIN DATA ', y_train.shape)
         print('Y TEST DATA ', y_test.shape)
         X TRAIN DATA (242, 8)
         X TEST DATA (61, 8)
         Y TRAIN DATA (242,)
         Y TEST DATA (61,)
 In [9]: #from sklearn.preprocessing import StandardScaler
         #sc = StandardScaler()
         #xtrain = sc.fit_transform(x_train)
         \#xtest = sc.transform(x_test)
In [10]: from sklearn.naive_bayes import GaussianNB
         clf = GaussianNB()
         clf.fit(x_train, y_train)
Out[10]: GaussianNB(priors=None, var_smoothing=1e-09)
In [11]: filename = open('heartdiseasespredictmodel.pkl', 'wb')
         pickle.dump(clf,filename )
In [12]: filename.close()
In [13]: y_pred = clf.predict(x_test)
In [14]: from sklearn.metrics import confusion_matrix
         print(confusion_matrix(y_test,y_pred))
         [[19 8]
          [ 5 29]]
In [15]: from sklearn.metrics import classification_report
         print('CLASSIFICATION REPORT','\n',classification_report(y_test, y_pred))
         CLASSIFICATION REPORT
                                   recall f1-score support
                       precision
                           0.79
                                     0.70
                                              0.75
                                                          27
                           0.78
                                     0.85
                                              0.82
                                                          34
                                                          61
                                              0.79
             accuracy
                           0.79
                                              0.78
            macro avg
                                     0.78
                                                          61
         weighted avg
                           0.79
                                     0.79
                                                          61
                                              0.79
In [16]: #from sklearn.ensemble import RandomForestClassifier
         #clf = RandomForestClassifier(n_estimators=20)
         #clf.fit(xtrain, y_train)
In [17]: #y_pred = clf.predict(xtest)
In [18]: #from sklearn.metrics import confusion_matrix
         #print(confusion_matrix(y_test,y_pred))
In [19]: #from sklearn.metrics import classification_report
         #print('CLASSIFICATION REPORT','\n',classification_report(y_test, y_pred))
```

jupyter Heart_Diseases_Predict Last Checkpoint: 07/08/2020 (autosaved)

Logout