

OASIS - EDA and ML Prediction

Data and Library Setup

```
In [1]: import pandas as pd
import seaborn as sns
import numpy as np
import matplotlib.pyplot as plt

In [2]: import os
print(os.listdir("F:/Dementia Prediction/data"))

['oasis_cross-sectional.csv', 'oasis_longitudinal.csv']

In [3]: df = pd.read_csv('F:/Dementia Prediction/data/oasis_longitudinal.csv')
```

Exploratory Data Analysis

```
In [4]: df.info()

<class 'pandas.core.frame.DataFrame'>
RangeIndex: 373 entries, 0 to 372
Data columns (total 15 columns):
#   Column      Non-Null Count  Dtype
---  ---
0   Subject ID  373 non-null    object
1   MRI ID      373 non-null    object
2   Group       373 non-null    object
3   Visit       373 non-null    int64
4   MR Delay    373 non-null    int64
5   M/F        373 non-null    object
6   Hand       373 non-null    object
7   Age        373 non-null    int64
8   EDUC       373 non-null    int64
9   SES        354 non-null    float64
10  MMSE       371 non-null    float64
11  CDR        373 non-null    float64
12  eTIV       373 non-null    int64
13  nWBV       373 non-null    float64
14  ASF        373 non-null    float64
dtypes: float64(5), int64(5), object(5)
memory usage: 43.8+ KB

In [5]: print("Tota Rows and Columns (Rows,Columns) : ",df.shape)
#print first five rows of the dataset
df.head(5)
```

Tota Rows and Columns (Rows,Columns) : (373, 15)

	Subject ID	MRI ID	Group	Visit	MR Delay	M/F	Hand	Age	EDUC	SES	MMSE	CDR	eTIV	nWBV	ASF
0	OAS2_0001	OAS2_0001_MR1	Nondemented	1	0	M	R	87	14	2.0	27.0	0.0	1987	0.696	0.883
1	OAS2_0001	OAS2_0001_MR2	Nondemented	2	457	M	R	88	14	2.0	30.0	0.0	2004	0.681	0.876
2	OAS2_0002	OAS2_0002_MR1	Demented	1	0	M	R	75	12	NaN	23.0	0.5	1678	0.736	1.046
3	OAS2_0002	OAS2_0002_MR2	Demented	2	560	M	R	76	12	NaN	28.0	0.5	1738	0.713	1.010
4	OAS2_0002	OAS2_0002_MR3	Demented	3	1895	M	R	80	12	NaN	22.0	0.5	1698	0.701	1.034

```
In [6]: df.describe()
```

	Visit	MR Delay	Age	EDUC	SES	MMSE	CDR	eTIV	nWBV	ASF
count	373.000000	373.000000	373.000000	373.000000	354.000000	371.000000	373.000000	373.000000	373.000000	373.000000
mean	1.882038	595.104558	77.013405	14.597855	2.460452	27.342318	0.290885	1488.128686	0.729568	1.195461
std	0.922843	635.485118	7.640957	2.876339	1.134005	3.683244	0.374557	176.139286	0.037135	0.138092
min	1.000000	0.000000	60.000000	6.000000	1.000000	4.000000	0.000000	1106.000000	0.644000	0.876000
25%	1.000000	0.000000	71.000000	12.000000	2.000000	27.000000	0.000000	1357.000000	0.700000	1.099000
50%	2.000000	552.000000	77.000000	15.000000	2.000000	29.000000	0.000000	1470.000000	0.729000	1.194000
75%	2.000000	873.000000	82.000000	16.000000	3.000000	30.000000	0.500000	1597.000000	0.756000	1.293000
max	5.000000	2639.000000	98.000000	23.000000	5.000000	30.000000	2.000000	2004.000000	0.837000	1.587000

```
In [7]: df.isna().sum()
```

```
Subject ID    0
MRI ID       0
Group         0
Visit         0
MR Delay      0
M/F           0
Hand          0
Age           0
EDUC          0
SES           19
MMSE          2
CDR           0
eTIV          0
nWBV          0
ASF           0
dtype: int64
```

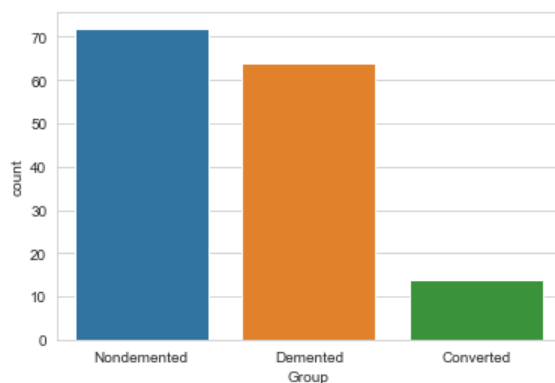
```
In [8]: sum(df.duplicated())
```

```
0
```

```
In [9]: df["SES"].fillna(df["SES"].median(), inplace=True)
df["MMSE"].fillna(df["MMSE"].mean(), inplace=True)
```

```
In [10]: sns.set_style("whitegrid")
ex_df = df.loc[df['Visit'] == 1]
sns.countplot(x='Group', data=ex_df)
```

```
<matplotlib.axes._subplots.AxesSubplot at 0x235af669790>
```



```
In [11]: ex_df['Group'] = ex_df['Group'].replace(['Converted'], ['Demented'])  
df['Group'] = df['Group'].replace(['Converted'], ['Demented'])  
sns.countplot(x='Group', data=ex_df)
```

<ipython-input-11-f052a051643f>:1: SettingWithCopyWarning:

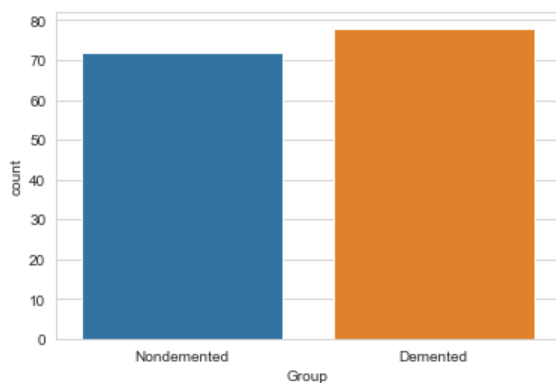
A value is trying to be set on a copy of a slice from a DataFrame.

Try using .loc[row_indexer,col_indexer] = value instead

See the caveats in the documentation: https://pandas.pydata.org/pandas-docs/stable/user_guide/indexing.html#returning-a-view-versus-a-copy

ex_df['Group'] = ex_df['Group'].replace(['Converted'], ['Demented'])

<matplotlib.axes._subplots.AxesSubplot at 0x235afdfd220>



```
In [12]: def bar_chart(feature):
    Demented = ex_df[ex_df['Group']=='Demented'][feature].value_counts()
    Nondemented = ex_df[ex_df['Group']=='Nondemented'][feature].value_counts()
    df_bar = pd.DataFrame([Demented,Nondemented])
    df_bar.index = ['Demented','Nondemented']
    df_bar.plot(kind='bar',stacked=True, figsize=(8,5))
    print(df_bar)
```

```
# Gender and Group ( Female=0, Male=1)
```

```
bar_chart('M/F')
```

```
plt.xlabel('Group',fontsize=13)
```

```
plt.xticks(rotation=0,fontsize=12)
```

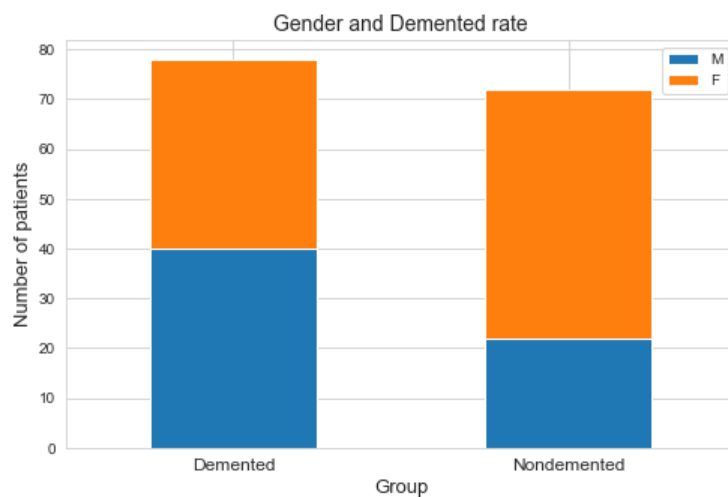
```
plt.ylabel('Number of patients',fontsize=13)
```

```
plt.legend()
```

```
plt.title('Gender and Demented rate',fontsize=14)
```

	M	F
Demented	40	38
Nondemented	22	50

```
Text(0.5, 1.0, 'Gender and Demented rate')
```

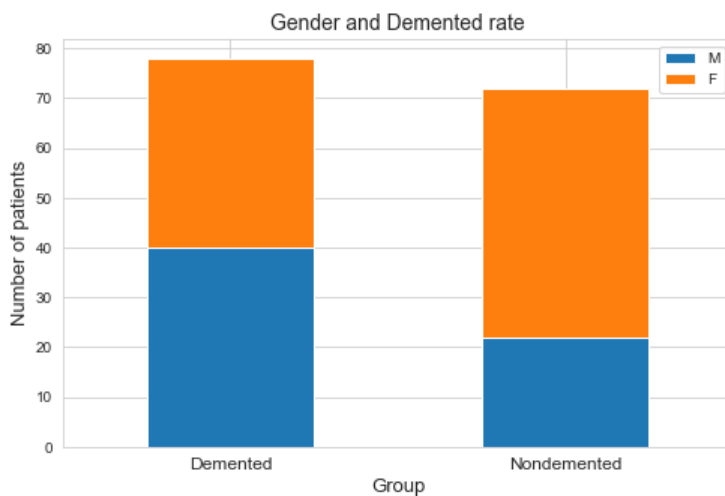


```
In [13]: def bar_chart(feature):
    Demented = ex_df[ex_df['Group']=='Demented'][feature].value_counts()
    Nondemented = ex_df[ex_df['Group']=='Nondemented'][feature].value_counts()
    df_bar = pd.DataFrame([Demented,Nondemented])
    df_bar.index = ['Demented','Nondemented']
    df_bar.plot(kind='bar',stacked=True, figsize=(8,5))
    print(df_bar)

# Gender and Group ( Female=0, Male=1)
bar_chart('M/F')
plt.xlabel('Group',fontsize=13)
plt.xticks(rotation=0,fontsize=12)
plt.ylabel('Number of patients',fontsize=13)
plt.legend()
plt.title('Gender and Demented rate',fontsize=14)
```

	M	F
Demented	40	38
Nondemented	22	50

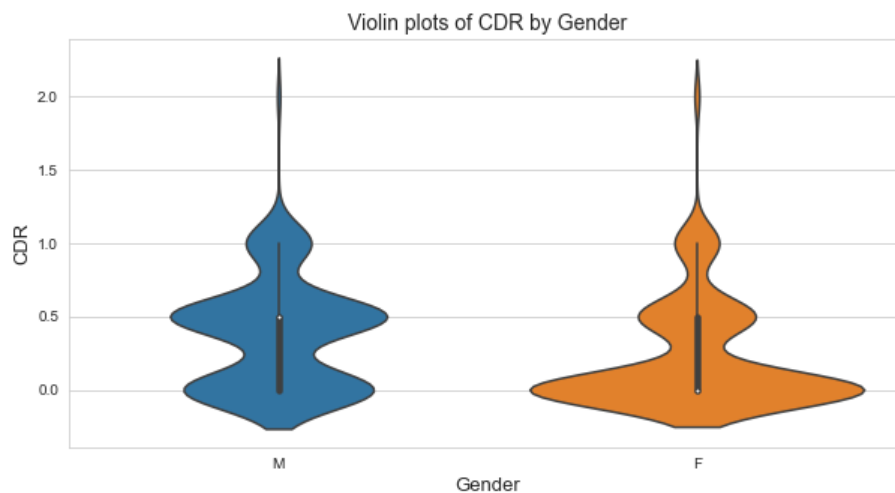
```
Text(0.5, 1.0, 'Gender and Demented rate')
```



CDR (Clinical Dementia Rating) : Ratings are assigned on a 0–5 point scale, (0 = absent; 0.5 = questionable; 1= present, but mild; 2 = moderate; 3 = severe; 4 = profound; 5 = terminal). A global summary score is obtained, leading to the use of the CDR for grouping patients on severity of dementia.

CDR By Gender

```
In [14]: plt.figure(figsize=(10,5))
sns.violinplot(x='M/F', y='CDR', data=df)
plt.title('Violin plots of CDR by Gender',fontsize=14)
plt.xlabel('Gender',fontsize=13)
plt.ylabel('CDR',fontsize=13)
plt.show()
```



CDR By Age

```
In [15]: plt.figure(figsize=(10,5))
sns.violinplot(x='CDR', y='Age', data=df)
plt.title('Violin plot of Age by CDR',fontsize=14)
plt.xlabel('CDR',fontsize=13)
plt.ylabel('Age',fontsize=13)
plt.show()
```



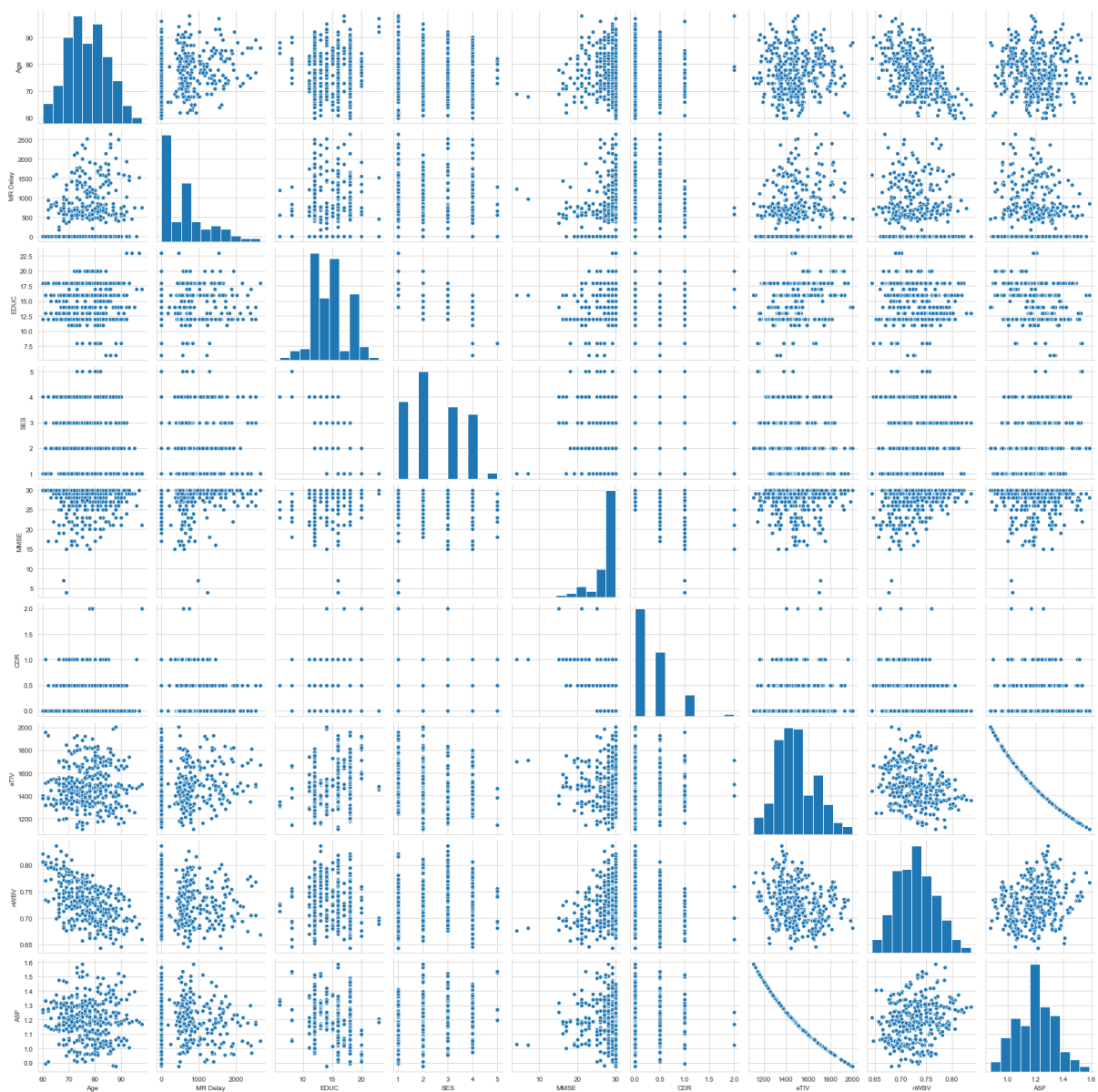
```
In [16]: def outliers_iqr(ys):
    quartile_1, quartile_3 = np.percentile(ys, [25, 75])
    iqr = quartile_3 - quartile_1
    lower_bound = quartile_1 - (iqr * 1.5)
    upper_bound = quartile_3 + (iqr * 1.5)
    return np.where((ys > upper_bound) | (ys < lower_bound))

list_attributes = ['MR Delay', 'EDUC', "SES", "MMSE", 'eTIV', "nWBV", "ASF"]
print("Outliers: \n")
for item in list_attributes:
    print(item, ': ', outliers_iqr(df[item]))
```

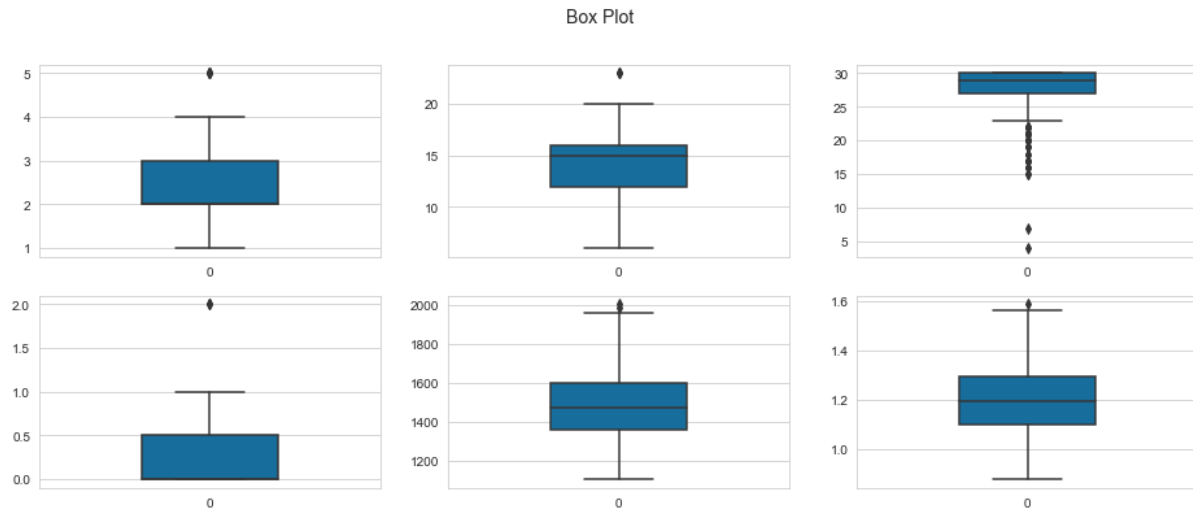
Outliers:

```
MR Delay : (array([ 32,  71,  75, 153, 159, 160, 265, 369], dtype=int64),)
EDUC : (array([107, 108, 109], dtype=int64),)
SES : (array([136, 137, 138, 161, 162, 179, 180], dtype=int64),)
MMSE : (array([ 4, 25, 26, 43, 44, 51, 52, 60, 88, 89, 90, 93, 94,
    97, 98, 99, 100, 101, 105, 106, 138, 162, 172, 173, 184, 185,
    186, 222, 225, 226, 231, 232, 234, 251, 299, 300, 316, 317, 328,
    332, 360, 366], dtype=int64),)
eTIV : (array([0, 1], dtype=int64),)
nWBV : (array([], dtype=int64),)
ASF : (array([282], dtype=int64),)
```

```
In [17]: from pylab import rcParams
rcParams['figure.figsize'] = 8,5
cols = ['Age', 'MR Delay', 'EDUC', 'SES', 'MMSE', 'CDR', 'eTIV', 'nWBV', 'ASF']
x=df.fillna('')
sns_plot = sns.pairplot(x[cols])
```




```
In [18]: fig, axes = plt.subplots(2,3,figsize = (16,6))
fig.suptitle("Box Plot",fontsize=14)
sns.set_style("whitegrid")
sns.boxplot(data=df['SES'], orient="v",width=0.4, palette="colorblind",ax = axes[0][0]);
sns.boxplot(data=df['EDUC'], orient="v",width=0.4, palette="colorblind",ax = axes[0][1]);
sns.boxplot(data=df['MMSE'], orient="v",width=0.4, palette="colorblind",ax = axes[0][2]);
sns.boxplot(data=df['CDR'], orient="v",width=0.4, palette="colorblind",ax = axes[1][0]);
sns.boxplot(data=df['eTIV'], orient="v",width=0.4, palette="colorblind",ax = axes[1][1]);
sns.boxplot(data=df['ASF'], orient="v",width=0.4, palette="colorblind",ax = axes[1][2]);
```

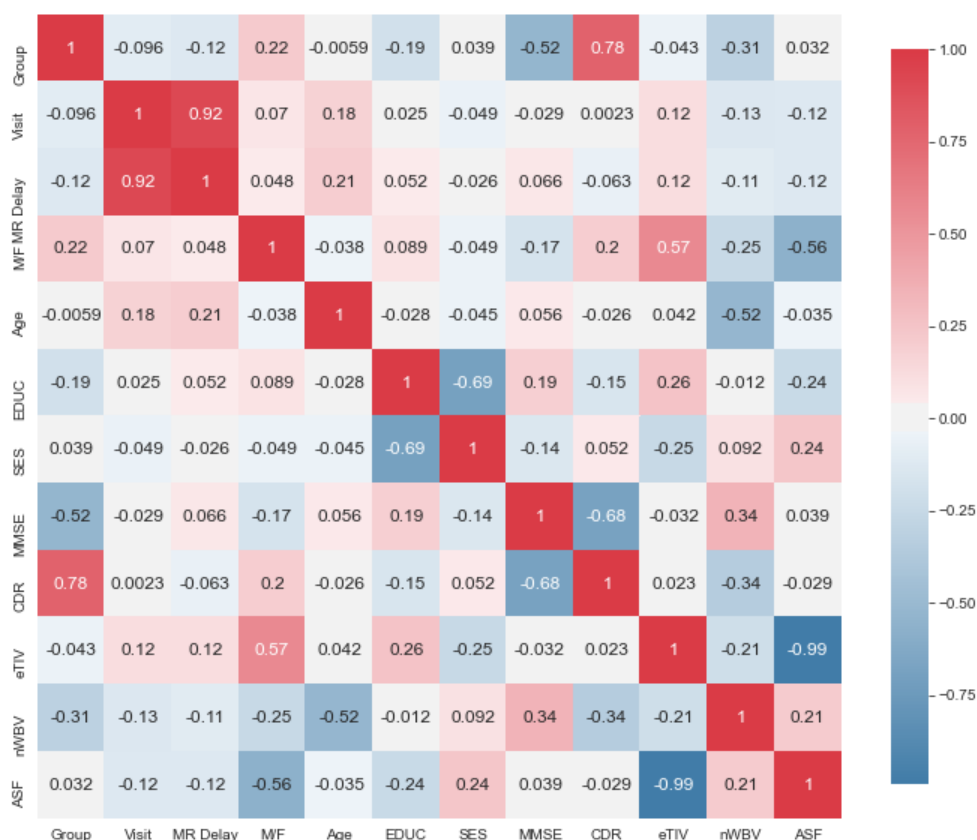


```
In [19]: #convet the charecter data into numeric
group_map = {"Demented": 1, "Nondemented": 0}

df['Group'] = df['Group'].map(group_map)
df['M/F'] = df['M/F'].replace(['F','M'], [0,1])
```

```
In [20]: def plot_correlation_map( df ):
    corr = df.corr()
    _, ax = plt.subplots( figsize=( 12 , 10 ) )
    cmap = sns.diverging_palette( 240 , 10 , as_cmap = True )
    _ = sns.heatmap(corr,cmap = cmap,square=True, cbar_kws={ 'shrink' : .9 }, ax=ax, annot = True,
    annot_kws = { 'fontsize' : 12 })
```

```
In [21]: plot_correlation_map(df)
```



Prediction task using Machine Learning (TBD)

```
In [22]: # Encode columns into numeric
from sklearn.preprocessing import LabelEncoder
for column in df.columns:
    le = LabelEncoder()
    df[column] = le.fit_transform(df[column])
```

```
In [23]: from sklearn.model_selection import train_test_split

feature_col_names = ["M/F", "Age", "EDUC", "SES", "MMSE", "eTIV", "nWBV", "ASF"]
predicted_class_names = ['Group']

X = df[feature_col_names].values
y = df[predicted_class_names].values

X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.30, random_state=42)
```

```

In [24]: from sklearn import metrics
def plot_confusion_matrix(y_test,model_test):
    cm = metrics.confusion_matrix(y_test, model_test)
    plt.figure(1)
    plt.clf()
    plt.imshow(cm, interpolation='nearest', cmap=plt.cm.Wistia)
    classNames = ['Nondemented', 'Demented']
    plt.title('Confusion Matrix')
    plt.ylabel('True label')
    plt.xlabel('Predicted label')
    tick_marks = np.arange(len(classNames))
    plt.xticks(tick_marks, classNames)
    plt.yticks(tick_marks, classNames)
    s = [['TN','FP'], ['FN', 'TP']]
    for i in range(2):
        for j in range(2):
            plt.text(j,i, str(s[i][j])+ " = " +str(cm[i][j]))
    plt.show()

In [25]: from sklearn.metrics import roc_curve, auc
def report_performance(model):

    model_test = model.predict(X_test)

    print("\n\nConfusion Matrix:")
    print("{0}".format(metrics.confusion_matrix(y_test, model_test)))
    print("\n\nClassification Report: ")
    print(metrics.classification_report(y_test, model_test))
    #cm = metrics.confusion_matrix(y_test, model_test)
    plot_confusion_matrix(y_test, model_test)

total_fpr = {}
total_tpr = {}
def roc_curves(model):
    predictions_test = model.predict(X_test)
    fpr, tpr, thresholds = roc_curve(predictions_test,y_test)
    roc_auc = auc(fpr, tpr)
    total_fpr[str((str(model).split('(')[0]))] = fpr
    total_tpr[str((str(model).split('(')[0]))] = tpr
    plt.figure()
    plt.plot(fpr, tpr, color='darkorange', lw=1, label='ROC curve (area = %0.2f)' % roc_auc)
    plt.plot([0, 1], [0, 1], color='navy', lw=2, linestyle='--')
    plt.xlim([0.0, 1.0])
    plt.ylim([0.0, 1.05])
    plt.xlabel('False Positive Rate')
    plt.ylabel('True Positive Rate')
    plt.title('Receiver operating characteristic')
    plt.legend(loc="lower right")
    plt.show()

In [26]: total_accuracy = {}
def accuracy(model):
    pred = model.predict(X_test)
    accu = metrics.accuracy_score(y_test,pred)
    print("\nAccuracy Of the Model: ",accu,"\n\n")
    total_accuracy[str((str(model).split('(')[0]))] = accu

```

Will be doing the ML using the following models:

- GAUSSIAN CLASSIFIER
- SVC
- Decision Tree
- XGBoost
- Gradient Boost
- Bagging
- Adaboost
- Random Forest
- LGBM
- Naive Bayes

```
In [29]: from sklearn.gaussian_process import GaussianProcessClassifier
from sklearn.svm import SVC
from sklearn.tree import DecisionTreeClassifier
from sklearn import metrics
from xgboost import XGBClassifier
from sklearn.ensemble import ExtraTreesClassifier
from sklearn.tree import ExtraTreeClassifier
from sklearn.tree import DecisionTreeClassifier
from sklearn.ensemble import GradientBoostingClassifier
from sklearn.ensemble import BaggingClassifier
from sklearn.ensemble import AdaBoostClassifier
from sklearn.ensemble import RandomForestClassifier
from lightgbm import LGBMClassifier
from sklearn.model_selection import GridSearchCV
from sklearn.naive_bayes import GaussianNB
```

```
In [30]: rfc=RandomForestClassifier(criterion='gini',max_depth=8,max_features='auto',n_estimators=200)

param_grid = {
    'n_estimators': [200],
    'max_features': ['auto'],
    'max_depth' : [4,5,6,7,8],
    'criterion' :['gini']
}

#CV_rfc = GridSearchCV(estimator=rfc, param_grid=param_grid, cv= 5,scoring = 'roc_auc')
rfc.fit(X_train, y_train.ravel())
#print("Best parameters set found on development set:")
#print(rfc.best_params_)
report_performance(rfc)
roc_curves(rfc)
accuracy(rfc)

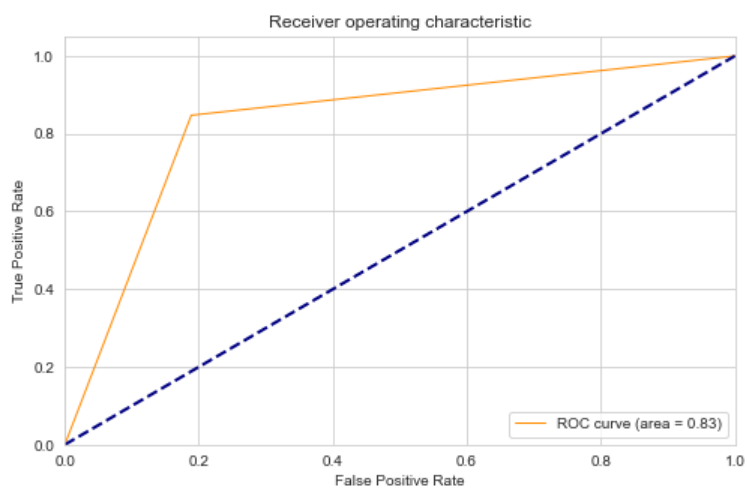
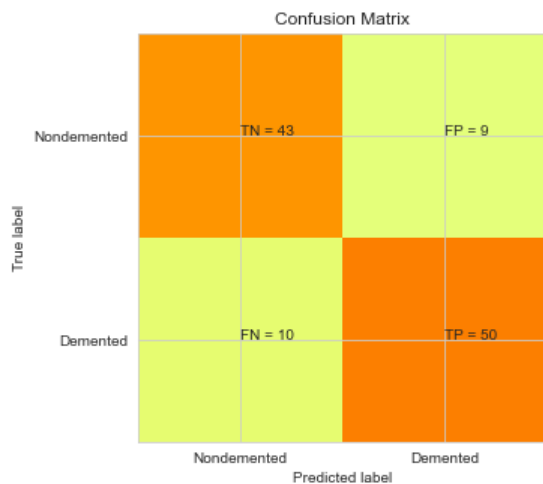
feat_importances = pd.Series(rfc.feature_importances_, index=feature_col_names)
feat_importances.nlargest(8).plot(kind='barh')
plt.title("Feature Importance:")
plt.show()
```

Confusion Matrix:

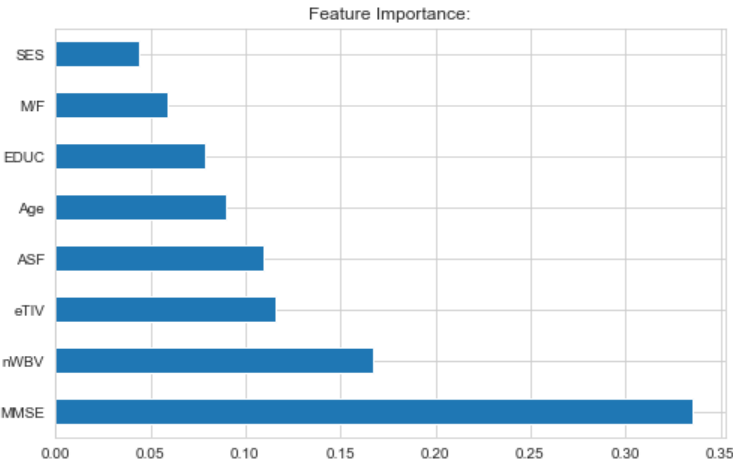
```
[[43  9]
 [10 50]]
```

Classification Report:

	precision	recall	f1-score	support
0	0.81	0.83	0.82	52
1	0.85	0.83	0.84	60
accuracy			0.83	112
macro avg	0.83	0.83	0.83	112
weighted avg	0.83	0.83	0.83	112



Accuracy Of the Model: 0.8303571428571429



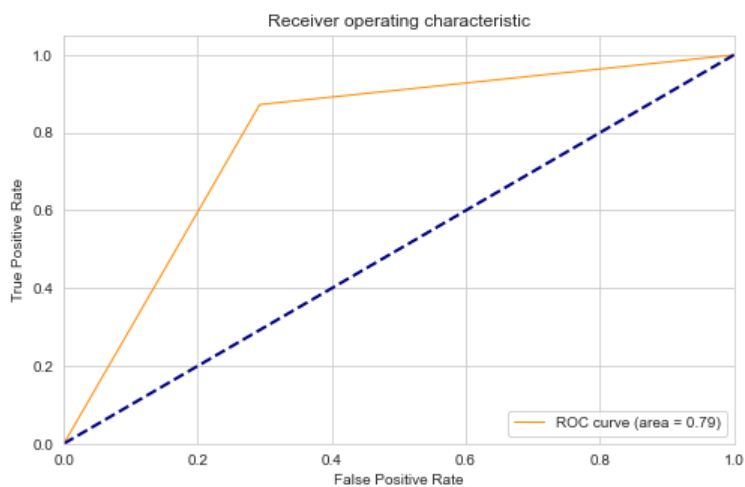
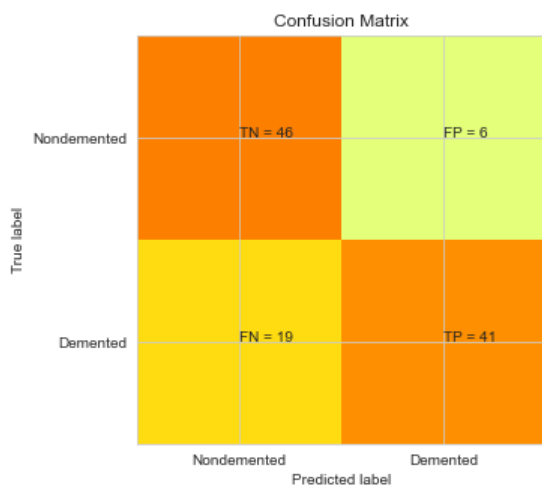
```
In [31]: svm = SVC(kernel="linear", C=0.1, random_state=0)
svm.fit(X_train, y_train.ravel())
report_performance(svm)
roc_curves(svm)
accuracy(svm)
```

Confusion Matrix:

```
[[46  6]
 [19 41]]
```

Classification Report:

	precision	recall	f1-score	support
0	0.71	0.88	0.79	52
1	0.87	0.68	0.77	60
accuracy			0.78	112
macro avg	0.79	0.78	0.78	112
weighted avg	0.80	0.78	0.78	112



Accuracy Of the Model: 0.7767857142857143

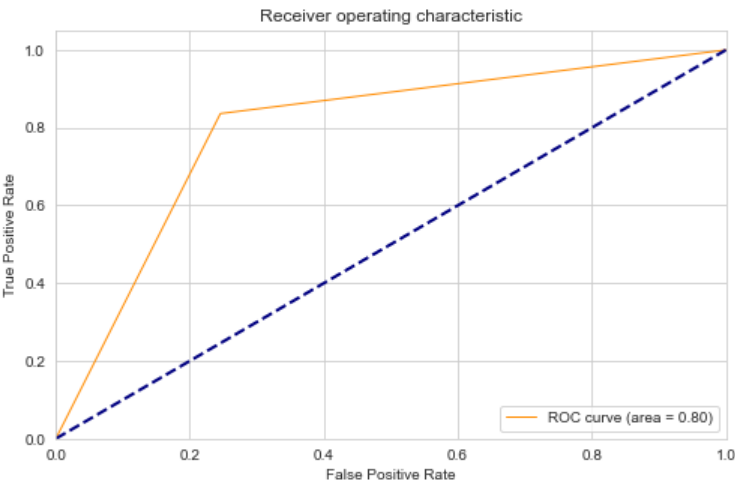
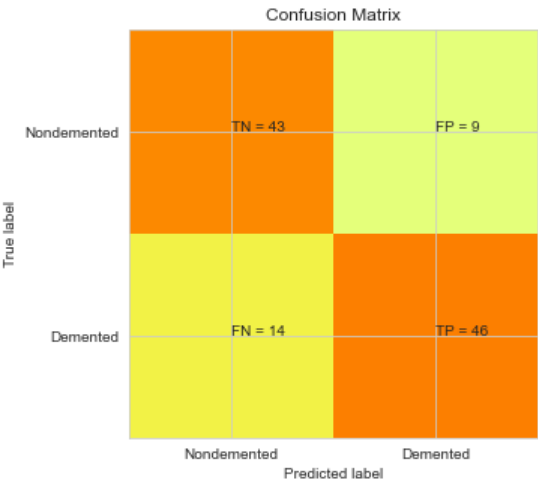

```
In [32]: clf_dtc = DecisionTreeClassifier(criterion='entropy',max_depth=5,random_state=0)
         clf_dtc.fit(X_train, y_train.ravel())
         report_performance(clf_dtc)
         roc_curves(clf_dtc)
         accuracy(clf_dtc)
         #importances = clf.feature_importances_

         feat_importances = pd.Series(clf_dtc.feature_importances_, index=feature_col_names)
         feat_importances.nlargest(8).plot(kind='barh')
         plt.title("Feature Importance:")
         plt.show()
```

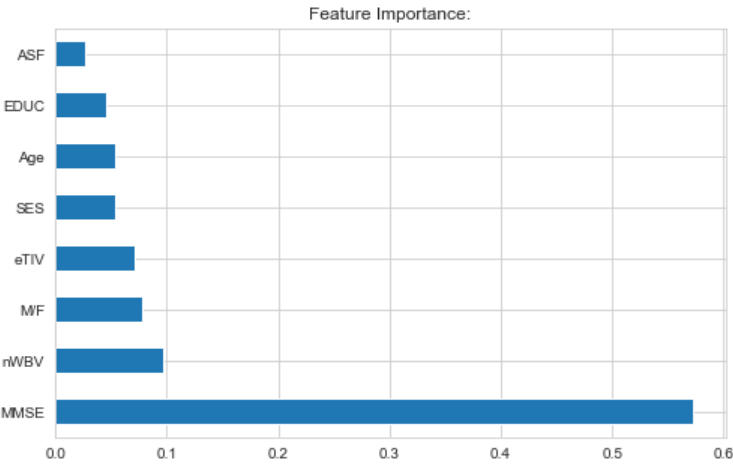
Confusion Matrix:
[[43 9]
[14 46]]

Classification Report:

	precision	recall	f1-score	support
0	0.75	0.83	0.79	52
1	0.84	0.77	0.80	60
accuracy			0.79	112
macro avg	0.80	0.80	0.79	112
weighted avg	0.80	0.79	0.79	112



Accuracy Of the Model: 0.7946428571428571



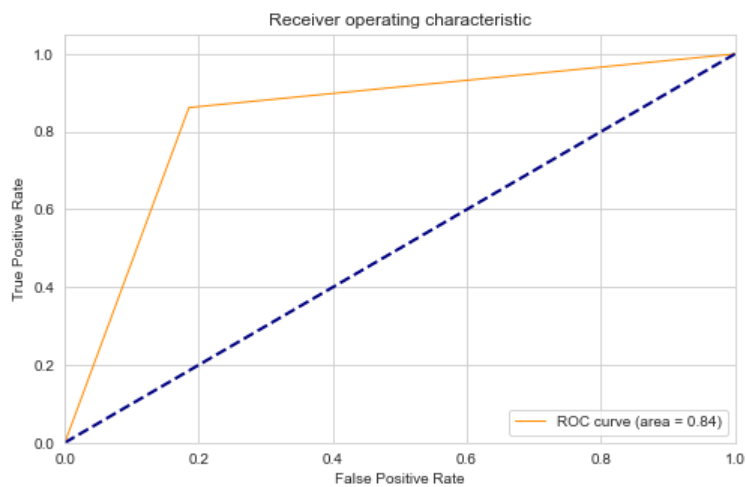
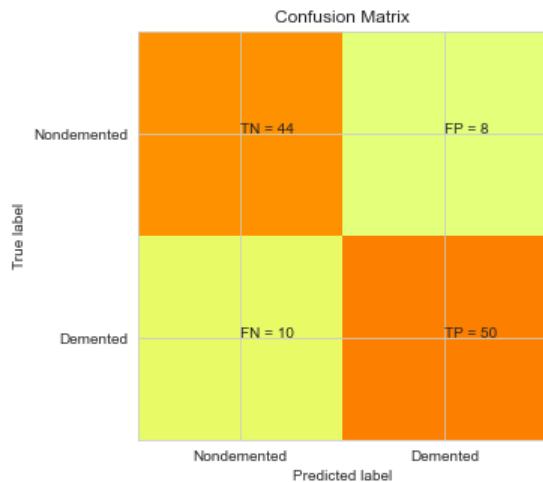
```
In [33]: params = {  
    'min_child_weight': [1, 5, 10],  
    'gamma': [0.5, 1, 1.5, 2, 5],  
    'subsample': [0.6, 0.8, 1.0],  
    'colsample_bytree': [0.6, 0.8, 1.0],  
    'max_depth': [1,2,3,4,5]  
}  
  
clf_xgb = XGBClassifier(random_state=0)  
clf_xgb.fit(X_train, y_train.ravel())  
report_performance(clf_xgb)  
roc_curves(clf_xgb)  
accuracy(clf_xgb)
```

Confusion Matrix:

```
[[44  8]
 [10 50]]
```

Classification Report:

	precision	recall	f1-score	support
0	0.81	0.85	0.83	52
1	0.86	0.83	0.85	60
accuracy			0.84	112
macro avg	0.84	0.84	0.84	112
weighted avg	0.84	0.84	0.84	112



Accuracy Of the Model: 0.8392857142857143

```
In [34]: from sklearn import svm, tree, linear_model, neighbors, naive_bayes, ensemble, discriminant_analys
is, gaussian_process
from xgboost import XGBClassifier
from sklearn.preprocessing import LabelEncoder
from sklearn import feature_selection
from sklearn import model_selection
```

```
In [35]: vote_est = [('etc', ensemble.ExtraTreesClassifier()),
                    ('gb', GradientBoostingClassifier()),
                    ('abc', AdaBoostClassifier()),
                    ('rfc', ensemble.RandomForestClassifier(criterion='gini', max_depth=8, max_features='auto', n_
estimators=200)),
                    #('svc', svm.SVC(probability=True)),
                    #('xgb', XGBClassifier()),
                    ('lbgm', LGBMClassifier())
                    ]

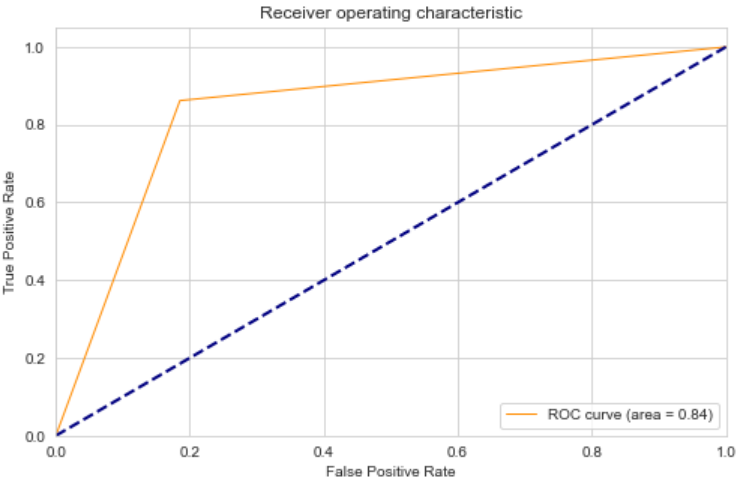
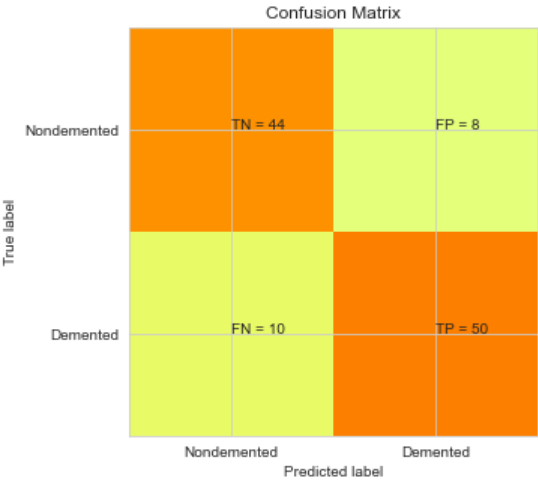
vote_hard = ensemble.VotingClassifier(estimators = vote_est , voting = 'hard')
vote_hard_cv = model_selection.cross_validate(vote_hard, X_train, y_train.ravel())
vote_hard.fit(X_train, y_train.ravel())
report_performance(vote_hard)
roc_curves(vote_hard)
accuracy(vote_hard)
#pred = vote_hard.predict(X_test)
#accu = metrics.accuracy_score(y_test, pred)
#print("\nAccuracy Of the Model: ", accu, "\n\n")

vote_soft = ensemble.VotingClassifier(estimators = vote_est , voting = 'soft')
vote_soft_cv = model_selection.cross_validate(vote_soft, X_train, y_train.ravel())
vote_soft.fit(X_train, y_train.ravel())
report_performance(vote_soft)
roc_curves(vote_soft)
accuracy(vote_soft)
#pred = vote_soft.predict(X_test)
#accu = metrics.accuracy_score(y_test, pred)
#print("\nAccuracy Of the Model: ", accu, "\n\n")
```

Confusion Matrix:
[[44 8]
[10 50]]

Classification Report:

	precision	recall	f1-score	support
0	0.81	0.85	0.83	52
1	0.86	0.83	0.85	60
accuracy			0.84	112
macro avg	0.84	0.84	0.84	112
weighted avg	0.84	0.84	0.84	112

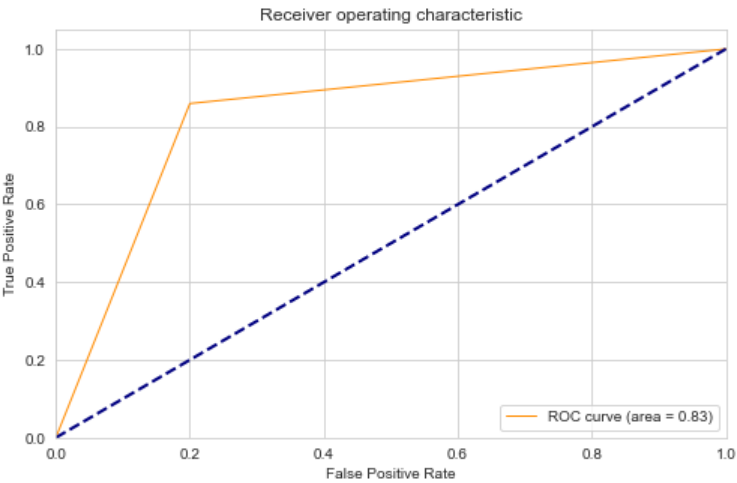
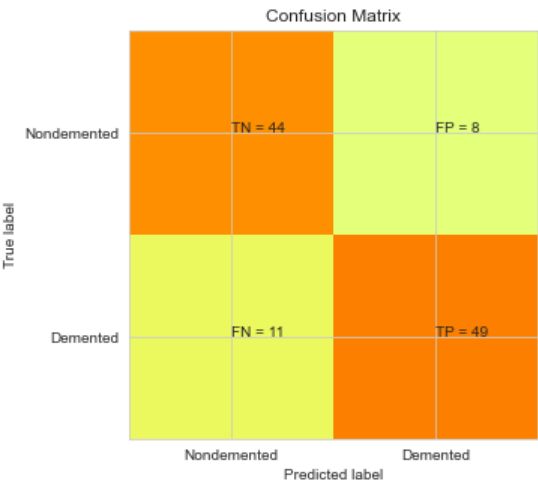


Accuracy Of the Model: 0.8392857142857143

Confusion Matrix:
[[44 8]
[11 49]]

Classification Report:

	precision	recall	f1-score	support
0	0.80	0.85	0.82	52
1	0.86	0.82	0.84	60
accuracy			0.83	112
macro avg	0.83	0.83	0.83	112
weighted avg	0.83	0.83	0.83	112



Accuracy Of the Model: 0.8303571428571429

```
In [36]: clfs =[ExtraTreesClassifier(),GradientBoostingClassifier(),AdaBoostClassifier()]
```



```
In [37]: for model in clfs:
          print(str(model).split('(')[0],": ")
          model.fit(X_train,y_train.ravel())
          X = pd.DataFrame(X_train)
          report_performance(model)
          roc_curves(model)
          accuracy(model)
```

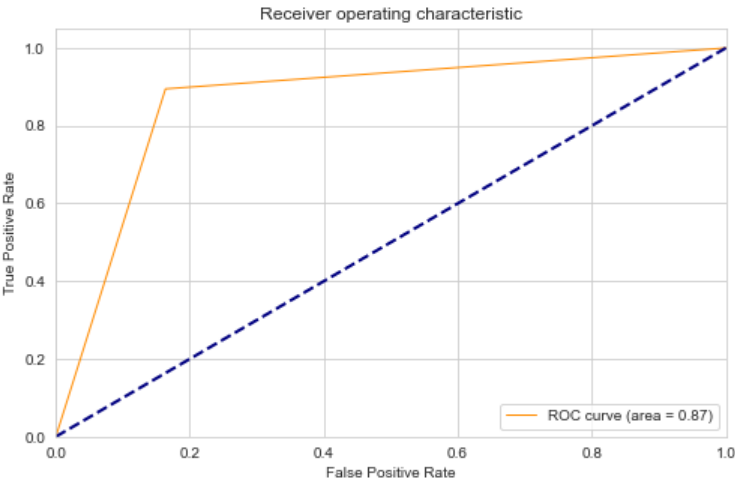
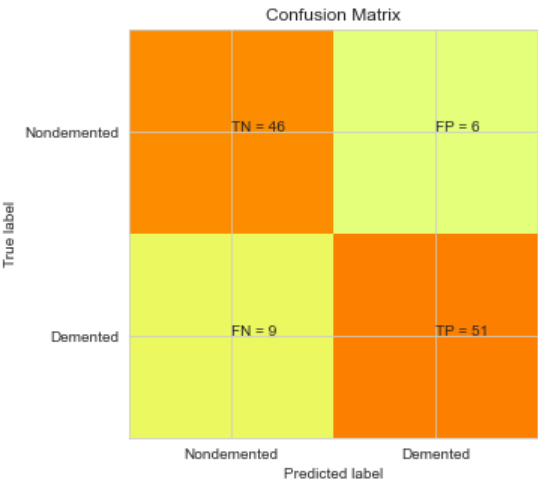
ExtraTreesClassifier :

Confusion Matrix:

```
[[46  6]
 [ 9 51]]
```

Classification Report:

	precision	recall	f1-score	support
0	0.84	0.88	0.86	52
1	0.89	0.85	0.87	60
accuracy			0.87	112
macro avg	0.87	0.87	0.87	112
weighted avg	0.87	0.87	0.87	112



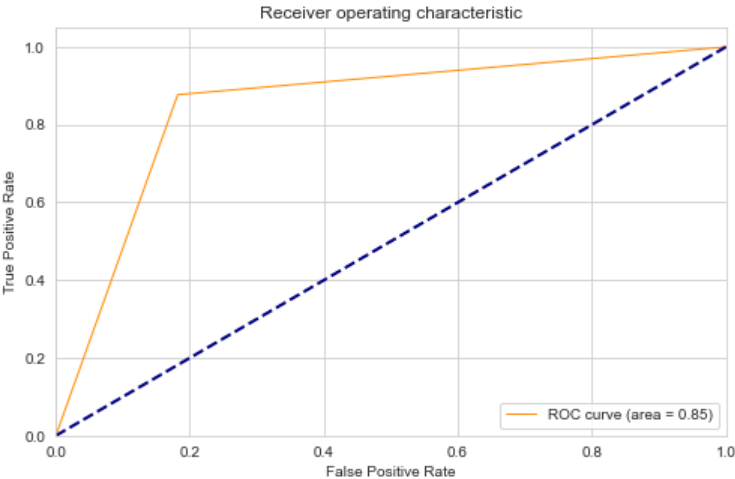
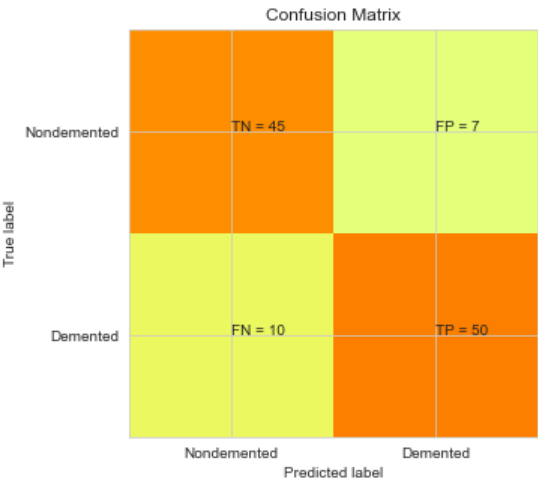
Accuracy Of the Model: 0.8660714285714286

GradientBoostingClassifier :

Confusion Matrix:
[[45 7]
[10 50]]

Classification Report:

	precision	recall	f1-score	support
0	0.82	0.87	0.84	52
1	0.88	0.83	0.85	60
accuracy			0.85	112
macro avg	0.85	0.85	0.85	112
weighted avg	0.85	0.85	0.85	112



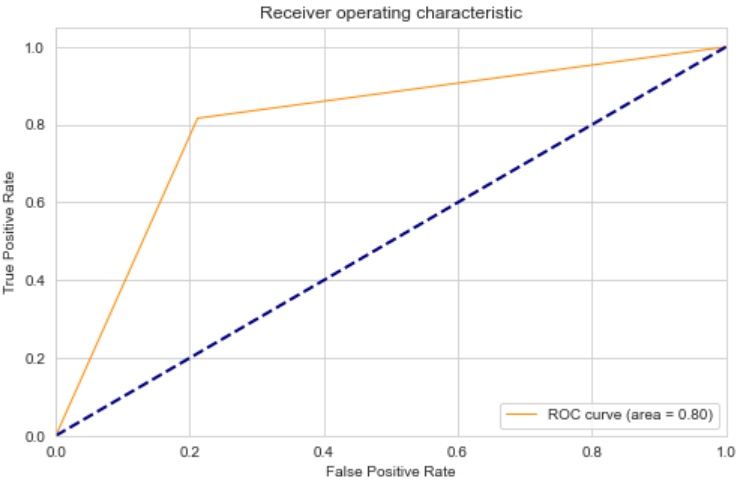
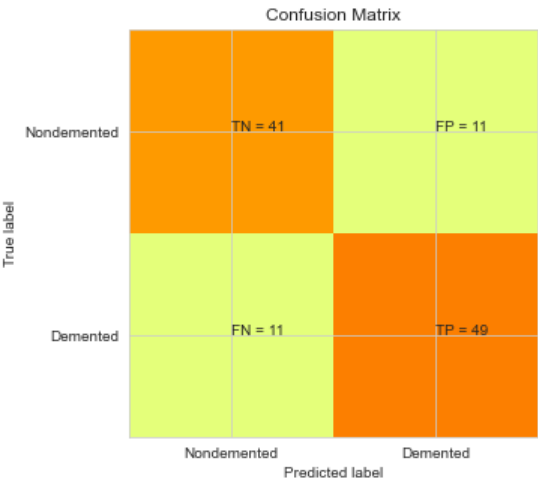
Accuracy Of the Model: 0.8482142857142857

AdaBoostClassifier :

Confusion Matrix:
[[41 11]
[11 49]]

Classification Report:

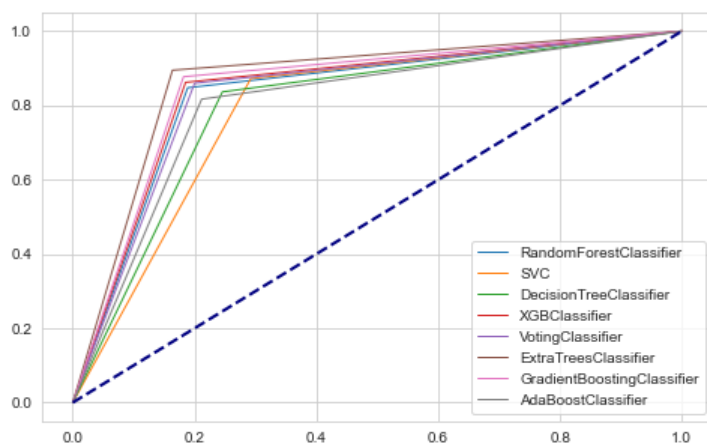
	precision	recall	f1-score	support
0	0.79	0.79	0.79	52
1	0.82	0.82	0.82	60
accuracy			0.80	112
macro avg	0.80	0.80	0.80	112
weighted avg	0.80	0.80	0.80	112



Accuracy Of the Model: 0.8035714285714286

```
In [40]: for i in total_fpr.keys():  
         plt.plot(total_fpr[i],total_tpr[i],lw=1, label=i)  
plt.plot([0, 1], [0, 1], color='navy', lw=2, linestyle='--')  
plt.legend()
```

<matplotlib.legend.Legend at 0x235b65eceb0>



```
In [ ]:
```