

[illegible]

```
from sklearn.metrics import classification_report

#Used for keeping time
import time

#0. Import the data

#Set the path for the CSV file
readPath = "J:\\DS\\Degree\\PennState\\DAAN_862\\Week 10\\Homework"

#Change the directory
os.chdir(readPath)

#Read the CSV file in
df = pd.read_csv("data2.csv")

In [3]: #1. Explore the Breast Cancer Data (10 points)

print("I will perform some exploration of the Breast Cancer data.")

#General Data Frame info
print("\n\nGeneral Formatting of the data frame:")
print(df.describe())
#Column info
print("\n\nColumn listing:")
print(df.columns)
#Looking for correlations
corr = df.corr()
print("Correlation Values:\n")
print(corr)
```

```
print("NA correlation plot:")

#plot a correlation matrix
plt.matshow(corr)
plt.title("Correlation Matrix")
plt.colorbar()
plt.xticks(range(10), list(df.columns))
plt.yticks(range(10), list(df.columns))

#Look for relationships in a pair plot
print("\nA pair plot of the data to explore for correlations")
sns.pairplot(df)
sns.pairplot(df, diag_kind='kde')

#split the variables and the classifiers
#Column variables
data = df.iloc[:, :9]
#Cancer classification
cancer = df.iloc[:,9]

I will perform some exploration of the Breast Cancer data.

General formatting of the data frame:
Age      BMI      Glucose      Insulin      HOMA      Leptin \
count    116.000000    116.000000    116.000000    116.000000    116.000000    116.000000
mean     59.301724    27.582111    97.793103    10.012086    2.694988    26.615080
std      16.112766    5.020336    22.552152    10.067768    3.642043    19.183294
min      24.000000    18.370000    60.000000    2.432000    0.467409    4.311000
25%      45.000000    22.973205    85.750000    4.359250    0.917966    12.313675
50%      56.000000    27.662416    92.000000    5.924500    1.380939    20.271000
75%      71.000000    31.241442    102.000000    11.189250    2.857787    37.378500
max      89.000000    38.578759    201.000000    58.460000    25.050342    90.280000

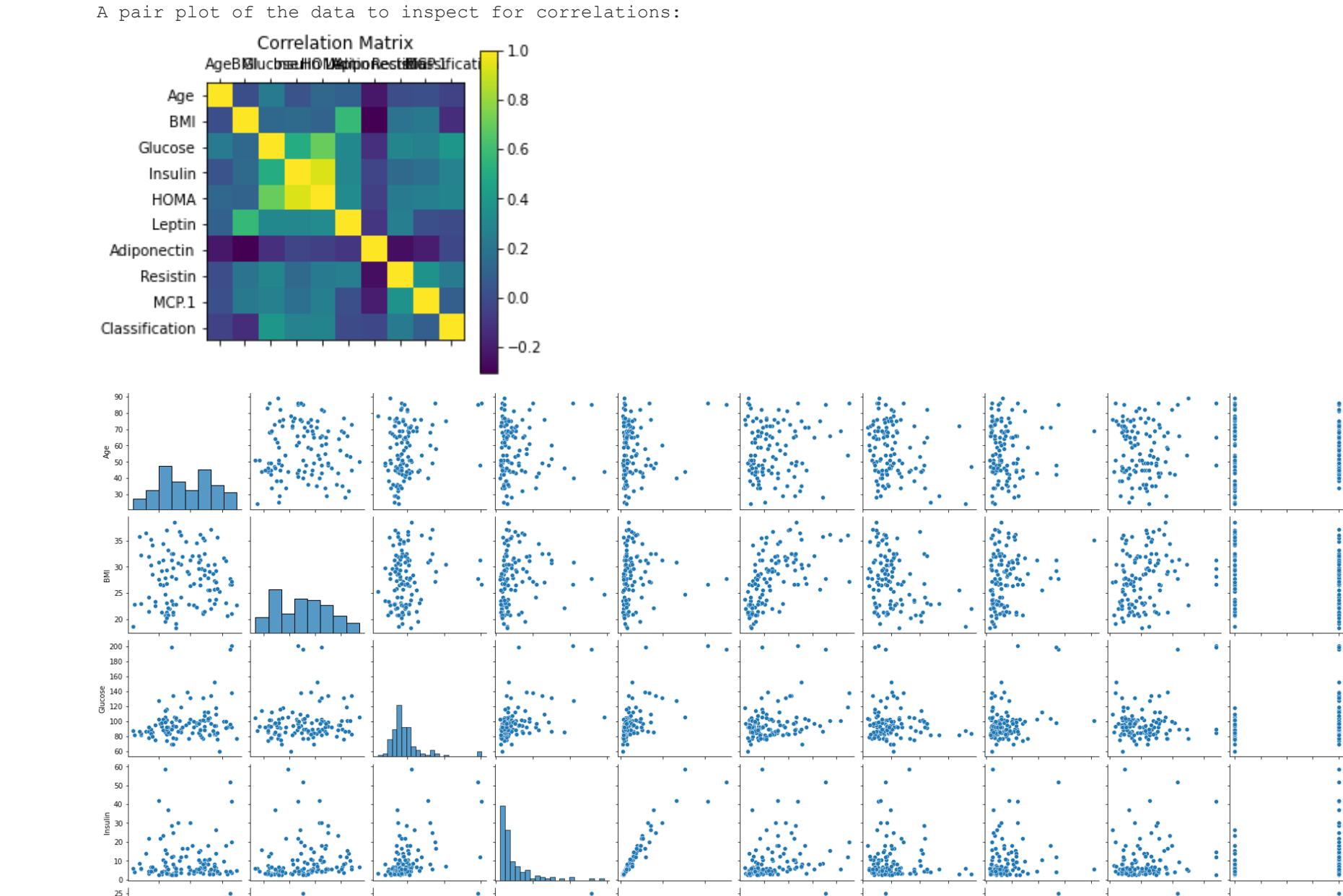
Adiponectin  Resistin      MCP.1  Classification
count    116.000000    116.000000    116.000000    116.000000
mean     10.180874    14.725966    534.647000    1.551724
std       6.843341    12.390646    345.912663    0.499475
1.656020    3.210000    45.843000    1.000000
25%       5.474283    6.881763    269.978250    1.000000
50%       8.352692    10.827740    471.322500    2.000000
75%      11.815970    17.755027    700.085000    2.000000
max      38.040000    82.100000    1698.440000    2.000000

Column listing :
Index('Age', 'BMI', 'Glucose', 'Insulin', 'HOMA', 'Leptin', 'Adiponectin',
      'Resistin', 'MCP.1', 'Classification'),
dtype='object')
Correlation Values:

Age      BMI      Glucose      Insulin      HOMA      Leptin \
Age      1.000000    0.008830    0.230106    0.032495    0.127033    0.102626
BMI      0.008830    1.000000    0.138845    0.145295    0.114480    0.569593
Glucose   0.230106    0.138845    1.000000    0.504653    0.696212    0.305080
Insulin   0.032495    0.145295    0.504653    1.000000    0.932198    0.301462
HOMA      0.127033    0.114480    0.696212    0.932198    1.000000    0.327210
Leptin    0.102626    0.569593    0.305080    0.301462    0.327210    1.000000
Adiponectin -0.219813    -0.302735    -0.122121    -0.031296    -0.056337    -0.095389
Resistin  0.002742    0.195350    0.291327    0.146731    0.231101    0.256234
MCP.1     0.013462    0.224038    0.264879    0.174356    0.239529    0.014099
Classification -0.043555    -0.126586    0.384315    0.276804    0.284012    -0.001078

Adiponectin  Resistin      MCP.1  Classification
Age      -0.219813    0.002742    0.013462    -0.043555
BMI       0.302735    0.195350    0.224038    -0.125286
Glucose   -0.122121    0.291327    0.264879    0.384315
Insulin   -0.031296    0.146731    0.174356    0.276804
HOMA      -0.056337    0.231101    0.239529    0.284012
Leptin    -0.095389    0.256234    0.014099    -0.001078
Adiponectin  1.000000    -0.252363    -0.200694    -0.019490
Resistin     0.252363    1.000000    0.366474    0.201710
MCP.1        0.200694    0.366474    1.000000    0.091381
Classification -0.019490    0.227310    0.091381    1.000000

A correlation plot:
```



```
print("time:" + str(end-start))

The metrics of the Linear SVC is:
fit_time      1.209549
score_time     0.002456
test_accuracy  0.763889
train_accuracy 0.839479
test_precision_macro 0.787500
train_precision_macro 0.836431
test_recall_macro 0.760833
train_recall_macro 0.840027
dtype: float64
time:12.149397899999713

The metrics of the Radial Basis Function SVC is:
fit_time      0.001892
score_time     0.002101
test_accuracy  0.568056
train_accuracy 0.567903
test_precision_macro 0.284028
train_precision_macro 0.283952
test_recall_macro 0.500000
train_recall_macro 0.500000
dtype: float64
time:0.07035910000013246

The Scores of the Polynomial SVC are:
precision      recall      f1-score      support
1      0.50      0.07      0.12      15
2      0.48      0.93      0.63      14
accuracy      0.49      0.50      0.48      29
macro avg      0.49      0.50      0.38      29
weighted avg    0.49      0.48      0.37      29

The polynomial used has the following estimators:
SVC(degree=4, kernel='poly')
{'degree': 4}

Results for all polynomials:
Degree  mean_train_score  mean_test_score
0      2      0.287366      0.285516
1      3      0.338106      0.285516
2      4      0.794140      0.392460

NOTE: The precision scores of the polynomial SVC came back ill-defined multiple times as all data came back positive for cancer.
This was the case at test_size=0.25. At lower values there was a split.

This can be seen in the confusion matrix:
[[ 1 14]
 [ 1 13]]

I was only able to get what looked like some sort of healthy split once and it was only when the poly degree was 4.
time:0.1668417999999292

The metrics of the Sigmoid Function SVC is:
fit_time      0.007700
score_time     0.002500
test_accuracy  0.565278
train_accuracy 0.559429
test_precision_macro 0.336111
train_precision_macro 0.360319
test_recall_macro 0.492500
train_recall_macro 0.491471
dtype: float64
NOTE: The Sigmoid SVC also came back ill-defined multiple times.
This is also due to the data all being classified as having cancer.

A check for classification also can be seen in the confusion matrix:
[[ 1 14]
 [ 1 13]]
time:0.07320069999968837
```

```

from sklearn.ensemble import RandomForestClassifier

#empty list to append scores into later

RFdf = pd.DataFrame()

#Use the train test from before as it is outside the loop
runs = range(0,6)

n_estim = range(2,201,2)

#Loop the number of runs
for j in runs:

    #reset accuracy to an empty list
    accuracy = []

    #Reset the train test split
    X_train, X_test, y_train, y_test = train_test_split(data,
                                                         cancer,
                                                         test_size=0.3)

    #loop over various values of n_estimators
    for i in n_estim:

        #Build the RFC, set a random state so that you can pin down n_estimators
        RF = RandomForestClassifier(n_estimators = i, random_state = 226)

        #Fit and get scores
        RF_scores = cross_val_score(RF, X_train, y_train, cv = 10)

        #append score values into accuracy list
        accuracy.append(RF_scores.mean())

    #put int the new data frame column
    runname = "run " + str(j)
    RFdf.insert(j, runname, accuracy)

#create some formatting for the plot
s = []
for i in range(2,201,2):
    s.append(i)
#set the index to count by 2's for the plot
RFdf = RFdf.set_index(s)

#plot the Ensemble Accuracy
plt.figure()
RFdf.plot()
plt.title("Random Forest Accuracy")
plt.ylabel('Accuracy')
plt.xlabel('N_Estimators in Model')
plt.legend(loc = 4, fontsize = 6)

#Get the average best value of n
#rather, get the n which hits the maximum at a moment from all of the runs
#and take the mean of those values

#Do not readjust the index since you've done this above
#This value is correct
avg_best_n = (int(RFdf.idxmax()).mean())

#Build the RF object
RF_best_est = RandomForestClassifier(n_estimators = avg_best_n)

#Fit the Random Forest
RF_best_est.fit(X_train, y_train)

#Run the prediction on the test data
RF_best_est_pred = RF_best_est.predict(X_test)

RF_best_report = metrics.classification_report(y_test, RF_best_est_pred)

```



```
ada_df = pd.DataFrame()

#make a smaller step size for adaboost
ada_n_estim = range(1,101,1)

#Loop for 10 runs
for j in runs:

    #empty list to append to for later
    ada_accuracy = []

    #Reset the train test split
    X_train, X_test, y_train, y_test = train_test_split(data,
                                                         cancer,
                                                         test_size=0.3)

    #loop over the count of n_estimators 0 - 200 by 2 from above
    for i in ada_n_estim:

        #Build teh adaboost object, modify learning rate, maybe use a grid?
        ada = AdaBoostClassifier(n_estimators=i, learning_rate=1,
                                 random_state=226)

        #run the CV on the AdaBoost classifier and get scores
        ada_score = cross_val_score(ada, X_train, y_train, cv=10)

        #put scores for each n_estimator into the ada_accuracy
        ada_accuracy.append(ada_score.mean())

    #put int the new data frame column
    runname = "run." + str(i)
    ada_df.insert(3, runname, ada_accuracy)
```

```
#create some formatting for the plot
s = []

for i in range(1,101,1):
    s.append(i)

#set the index to count by 2's for the plot
ada_df = ada_df.set_index(s))

#plot the Ensemble Accuracy
plt.figure()
ada_df.plot()
plt.title("AdaBoost Accuracy")
plt.ylabel("Accuracy")
plt.xlabel("N Estimators in Model")
plt.legend(loc = 1, fontsize = 6)

#don't reindex this
avg_best_ada_n = int(ada_df.idxmax().mean())

#run the adaBoost with the best average n_estimator
ada_best_est = AdaBoostClassifier(n_estimators=avg_best_ada_n,
                                  learning_rate = 0.005)

#fit
ada_best_est.fit(X_train, y_train)

#predict
ada_best_est_pred = ada_best_est.predict(X_test)

#metrics
ada_best_report = metrics.classification_report(y_test, ada_best_est_pred)

#Had to add one since the ada_n_estimators run from 1 - 100
print("\nMetrics for the AdaBoost at the Best n_estimator value of: " +
      str(avg_best_ada_n))
print(ada_best_report)
```

[illegible]



