HW2 Report

R10946013 劉警瑄

- 步驟說明:
 - 1. 讀取資料表

讀取資料表

```
df = pd.read_csv("GDSC_PDX_Paclitaxel.csv")#,index_col=0)
df = df.rename(columns={'Unnamed: 0':'CELL_LINE_NAME'}) # 把cell_line_name欄位名稱統一(方便做合併)
df_info = pd.read_csv("GDSC_PDX_Paclitaxel_info.csv",index_col=0)
df_test = pd.read_csv("CCLE_PDX_Paclitaxel.csv")#,index_col=0)
df_test = df_test.rename(columns={'CCLE.Cell.Line.Name':'CELL_LINE_NAME'})

[2] 
$\square 23.7s$
```

2. 取欄位交集(兩資料表都有的欄位 features)

取欄位交集(兩個資料表都有的欄位)

```
train_features = set(df.columns)
test_features = set(df_test.columns)
# 把非交集的欄位從資料表中未除(減去差集)
(variable) df_test: DataFrame es.difference(test_features),axis=1)
df_test = df_test.drop(test_features.difference(train_features),axis=1)

✓ 0.2s
```

3. 計算 IC50 與 MAX_CONC_MICROMOLAR 欄位的關係

用IC50與MAX CONC MICROMOLAR欄位計算S或R 並將結果併入訓練資料表中

4. 資料標準化

資料標準化

```
X = df_train.drop("ans",1) #Feature Matrix
y = df_train["ans"] #目標變量
df_train = df_train.drop(['CELL_LINE_NAME'], axis=1)

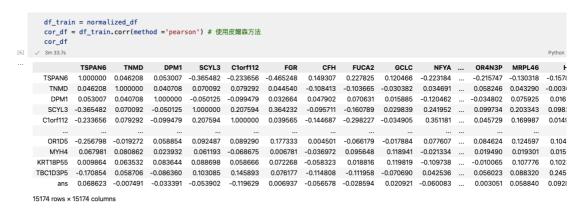
df_ = df_train.drop("ans",axis=1)

normalized_df=(df_ - df_.mean())/df_.std()
normalized_df['ans'] = y # 把標準化前刪掉的ans欄位併回來

/ 13.2s
```

5. 計算 correlation 以便做 Feature Selection (使用 Pearson 方法):

計算各欄位與欲預測欄位的correlation (pearson方法)



6. 設定閥值選取 features

(經測試後使用 cor 絕對值大於 0.17 的欄位組合模型表現最佳)

設定一個correlarion閥值來選取要用來訓練的feature(取絕對值後大於閥值)

```
#Correlation with output variable
cor_target = abs(cor_df["ans"])

#Selecting highly correlated features
relevant_features = cor_target[cor_target>0.17] #0.17 #設定閥值
print("用來訓練模型的feature數量:",len(relevant_features))

v 0.4s

... 用來訓練模型的feature數量: 108
```

- 7. 選定 Feature 後整理為建模資料
 - ※ 將feature 選定後整理為建模用資料

8. 調整參數與建模

我一共使用了三種模型進行測試,分別為: RandomForest Classifier、 XGBoost Classifier、KNN Classifier。

針對模型我也進行了調整參數,調整模型最主要、影響通常最大的參數:"n_neighbors",並觀察準確度的變化,以選取最佳的參數值。 此外,我也使用了 K-fold cross validation 來進行觀察,並取五次的平均 作為 Accuracy 的參考值。

結果如下:

a. Random Forest Classifier:

✓ Cross Validation 的結果(n_neighbors= range(10,290,20))

[0.7125	0.7125	0.725	0.8	0.72151899]
[0.75	0.7625	0.8	0.75	0.7721519]
[0.725	0.775	0.775	0.8125	0.75949367]
[0.7	0.725	0.825	0.7875	0.79746835]
[0.7375	0.7625	0.8125	0.7625	0.81012658]
[0.725	0.75	0.825	0.7625	0.75949367]
[0.725	0.775	0.8	0.8	0.78481013]
[0.7625	0.725	0.8125	0.8	0.79746835]
[0.7375	0.75	0.8	0.775	0.81012658]
[0.7625	0.7625	0.825	0.775	0.79746835]
[0.7375	0.7375	0.8	0.7875	0.79746835]
[0.725	0.7625	0.825	0.7625	0.78481013]
[0.7375	0.75	0.8125	0.8	0.81012658]
[0.7625	0.7625	0.8125	0.8	0.81012658]

✓ 取表現最佳的參數組合:

Accuracy = 0.7895253164556962 (N_neighbors = 270)

b. XGBoost Classifier:

✓ Cross Validation 的結果(n_neighbors= range(10,290,20))

[0.675	0.775	0.7875	0.725	0.81012658]
[0.6375	0.775	0.8	0.7375	0.75949367]
[0.65	0.775	0.7875	0.75	0.75949367]
[0.6625	0.8	0.8	0.7625	0.7721519]
[0.675	0.7875	0.8	0.75	0.7721519]
[0.675	0.7875	0.825	0.75	0.7721519]
[0.6875	0.7875	0.825	0.7375	0.78481013]
[0.7	0.8	0.8125	0.7375	0.79746835]
[0.725	0.7875	0.8125	0.7375	0.78481013]
[0.725	0.7875	0.8	0.7375	0.78481013]
[0.725	0.7875	0.775	0.7375	0.78481013]
[0.725	0.7875	0.775	0.7375	0.78481013]
[0.7125	0.7875	0.775	0.75	0.79746835]
[0.725	0.7875	0.775	0.75	0.78481013]

✓ 取表現最佳的參數組合:

Accuracy = 0.7694936708860759 (N_neighbors = 150)

c. KNN Classifier:

✓ Cross Validation 的結果(n_neighbors= range(10,290,20))

```
0.75
[0.7125
           0.725
                      0.725
                                           0.75949367]
[0.725
           0.725
                      0.7125
                                0.7125
                                           0.70886076]
[0.7
           0.7125
                      0.7125
                                0.7125
                                           0.70886076]
[0.7125
           0.7125
                      0.7125
                                0.7125
                                           0.70886076]
[0.7125
           0.7125
                     0.7125
                                0.7125
                                           0.70886076]
[0.7125
          0.7125
                    0.7125
                                0.7125
                                           0.70886076]
                                0.7125
[0.7125
          0.7125
                    0.7125
                                           0.70886076]
[0.7125
         0.7125
                    0.7125
                                0.7125
                                           0.70886076]
[0.7125
          0.7125
                    0.7125
                                0.7125
                                           0.70886076]
[0.7125
           0.7125
                     0.7125
                                0.7125
                                           0.70886076]
[0.7125
          0.7125
                     0.7125
                                0.7125
                                           0.70886076]
[0.7125
           0.7125
                     0.7125
                               0.7125
                                           0.70886076]
[0.7125
           0.7125
                     0.7125
                                0.7125
                                           0.70886076]
[0.7125
           0.7125
                      0.7125
                                0.7125
                                           0.70886076]
```

✓ 取表現最佳的參數組合:

Accuracy = 0.7343987341772152 (N_neighbors = 10)

● 結論與說明:

觀察上述三者模型,我們可以發現 Random Forest 在參數為 n_neighbors= 290 時的表現為三者中最佳。且在 Validation 時,我們也發現模型並未有 Overfitting 的狀態(如下圖)。因此我們選用此組合作為預測 CCLE 資料時使用的模型。

最好的模型與參數組合

print(accuracy)

[52] \(\sigma 0.1s \)

```
forest = ensemble.RandomForestClassifier(n_estimators = 270) #參數設為最好的
   forest_fit = forest.fit(train_X, train_y)
   # 預測
   test_y_predicted = forest.predict(test_X)
   # 績效
   scores = cross_val_score(forest,X,y,cv=5,scoring='accuracy')
   print(scores)
   for i in scores:
       sum+=i
   print(sum/5)
[0.725
           0.7625
                    0.8125
                              0.775
                                        0.7721519]
0.7694303797468354
  accuracy = metrics.accuracy_score(test_y, test_y_predicted)
```

Training accuracy 為 0.7694, Testing accuracy 為 0.75→未 overfitting

而針對欲預測的 CCLE 資料表,我們也做相同的處理步驟,需選取訓練模型時使用的特定 Feature 組合,並進行資料的標準化處理等,最後丟入模型預測,並輸出結果為 ans.csv 資料檔。