

Logistic Regression, MLP and Random Forest on Baxter Dataset

Let's make sure we are using the Python version we want

```
import sys
print(sys.version)
```

[Output] 3.6.6 (v3.6.6:4cf1f54eb7, Jun 26 2018, 17:02:57)

[Output] [GCC 4.2.1 (Apple Inc. build 5666) (dot 3)]

Let's make sure we are in our project directory

```
import os
print(os.getcwd())
```

[Output] /Users/Begum/Documents/DeepLearning

Start by importing modules that will be necessary

```
##### IMPORT MODULES #####
# We need to use a backend for matplotlib
# https://stackoverflow.com/questions/21784641/installation-issue-with-matplotlib-python/21789908#21789908
import matplotlib as mpl
mpl.use('TkAgg')
import matplotlib.pyplot as plt
from sklearn.linear_model import LogisticRegression
from sklearn.model_selection import train_test_split
from sklearn.model_selection import KFold, cross_val_score, validation_curve
from sklearn import linear_model
import numpy as np # linear algebra
import pandas as pd # data processing, CSV file I/O (e.g. pd.read_csv)

from sympy import *
from sklearn.model_selection import GridSearchCV
from sklearn.neural_network import MLPClassifier
from scipy import interp
from itertools import cycle
from sklearn import svm, datasets
from sklearn.metrics import roc_curve, auc
from sklearn.model_selection import StratifiedKFold
# dependencies for statistic analysis
from scipy import stats
# importing our parameter tuning dependencies
from sklearn.neighbors import KNeighborsClassifier
from sklearn.model_selection import (cross_val_score, GridSearchCV, StratifiedKFold, ShuffleSplit )
# importing our dependencies for Feature Selection
from sklearn.feature_selection import (SelectKBest, chi2, RFE, RFECV)
from sklearn.linear_model import LogisticRegression, RandomizedLogisticRegression
from sklearn.ensemble import RandomForestClassifier
```

```

from sklearn.metrics import f1_score
# Importing our sklearn dependencies for the modeling
from sklearn.ensemble import RandomForestClassifier
from sklearn.neighbors import NearestNeighbors
from sklearn.preprocessing import StandardScaler
from sklearn.model_selection import train_test_split
from sklearn import metrics
from sklearn.metrics import (accuracy_score, confusion_matrix, classification_report, roc_curve, auc)
from sklearn.neural_network import MLPClassifier
from itertools import cycle
from scipy import interp
import warnings
from sklearn.model_selection import StratifiedKFold
warnings.filterwarnings('ignore')

```

Read and prepare data

```

## Read in the data
shared = pd.read_table("data/baxter.0.03.subsample.shared")
shared.head()
meta = pd.read_table("data/metadata.tsv")
## Check and visualize the data
print(meta.head())

```

```

[Output]      sample  fit_result      Site  ...  Diabetes_Med stage  Location
[Output] 0  2003650           0  U Michigan  ...           0      0      NaN
[Output] 1  2005650           0  U Michigan  ...           0      0      NaN
[Output] 2  2007660          26  U Michigan  ...           0      0      NaN
[Output] 3  2009650          10    Toronto  ...           0      0      NaN
[Output] 4  2013660           0  U Michigan  ...           0      0      NaN
[Output]
[Output] [5 rows x 27 columns]

```

```

print(shared.head())
## Remove unnecessary columns from meta

```

```

[Output]      label      Group  numOtu  ...  Otu11278  Otu11280  Otu11281
[Output] 0   0.03  2003650    6920  ...           0           0           0
[Output] 1   0.03  2005650    6920  ...           0           0           0
[Output] 2   0.03  2007660    6920  ...           0           0           0
[Output] 3   0.03  2009650    6920  ...           0           0           0
[Output] 4   0.03  2013660    6920  ...           0           0           0
[Output]
[Output] [5 rows x 6923 columns]

```

```

meta = meta[['sample', 'dx']]
## Rename the column name "Group" to match the "sample" in meta
shared = shared.rename(index=str, columns={"Group": "sample"})
## Merge the 2 datasets on sample
data = pd.merge(meta, shared, on=['sample'])
## Remove adenoma samples
data = data[data.dx.str.contains("adenoma") == False]
## Drop all except OTU columns for x
x = data.drop(["sample", "dx", "numOtu", "label"], axis=1)
## Cancer =1 Normal =0

```

```

diagnosis = { "cancer":1, "normal":0}
##Generate y which only has diagnosis as 0 and 1
y = data["dx"].replace(diagnosis)
# y = np.eye(2, dtype='uint8')[y]
## Drop if NA elements
y.dropna()
x.dropna()
print(x.head())

```

```

[Output]      Otu00001  Otu00002  Otu00003  ...      Otu11278  Otu11280  Otu11281
[Output] 0          350          268          213  ...          0          0          0
[Output] 1          568         1320          13  ...          0          0          0
[Output] 2          151          756          802  ...          0          0          0
[Output] 4         1409          174           0  ...          0          0          0
[Output] 5          167          712          213  ...          0          0          0
[Output]
[Output] [5 rows x 6920 columns]

```

```
print(y.head())
```

```

[Output] 0      0
[Output] 1      0
[Output] 2      0
[Output] 4      0
[Output] 5      0
[Output] Name: dx, dtype: int64

```

```
print(len(x))
```

```
[Output] 292
```

```
print(len(y))
```

```
[Output] 292
```

Split the data to generate training and testing set %80-20

Here we also want to shuffle the data and set a random state

```
x_train, x_test, y_train, y_test = train_test_split(x, y, test_size=0.2, random_state=82089, shuffle=True)
```

Now let's define a L2 regularized logistic regression model

```

## Define L2 regularized logistic classifier
logreg = linear_model.LogisticRegression(C=0.01)

```

Plot ROC curve for Logistic Regression training set

We split the training set to 10 to cross-validate

```

cv = StratifiedKFold(n_splits=10)
tprs = []
aucs = []
mean_fpr = np.linspace(0, 1, 100)
## Convert from pandas dataframe to numpy
X=x_train.values

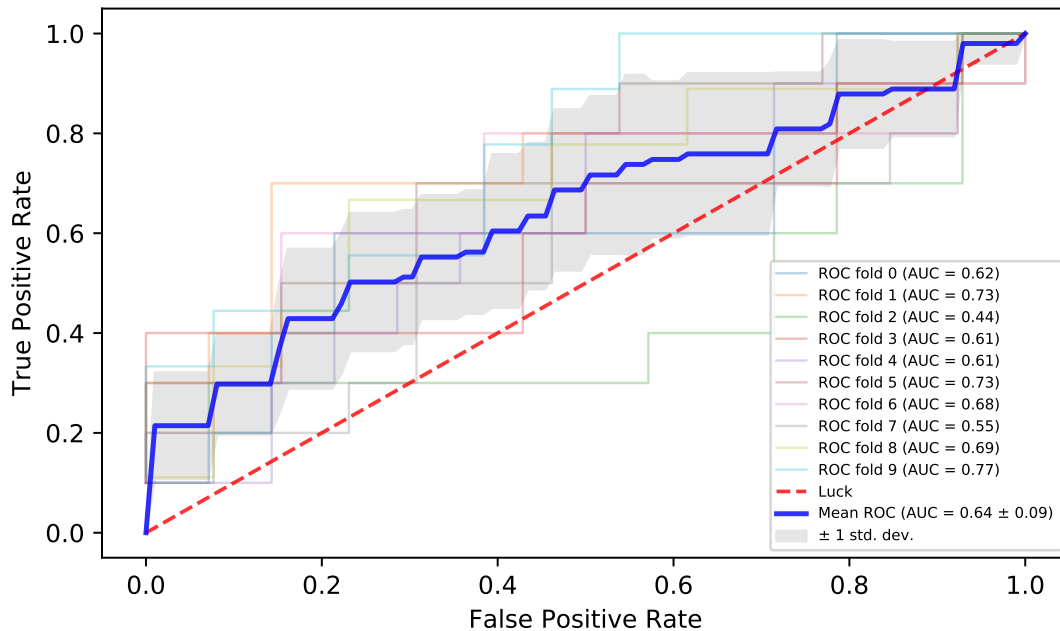
```

```

Y=y_train.values
# Plot the ROC curve over 10 iterations for each split
#logistic_plot = plt.figure()
i = 0
for train, test in cv.split(X,Y):
    probas_ = logreg.fit(X[train], Y[train]).predict_proba(X[test])
    # Compute ROC curve and area the curve
    fpr, tpr, thresholds = roc_curve(Y[test], probas_[ :, 1])
    tprs.append(interp(mean_fpr, fpr, tpr))
    tprs[-1][0] = 0.0
    roc_auc = auc(fpr, tpr)
    aucs.append(roc_auc)
    plt.plot(fpr, tpr, lw=1, alpha=0.3, label='ROC fold %d (AUC = %0.2f)' % (i, roc_auc))
    i += 1
plt.plot([0, 1], [0, 1], linestyle='--', color='r', label='Luck', alpha=.8)
mean_tpr = np.mean(tprs, axis=0)
mean_tpr[-1] = 1.0
mean_auc = auc(mean_fpr, mean_tpr)
std_auc = np.std(aucs)
plt.plot(mean_fpr, mean_tpr, color='b', label=r'Mean ROC (AUC = %0.2f  $\pm$  %0.2f)' % (mean_auc, std_auc))
std_tpr = np.std(tprs, axis=0)
tprs_upper = np.minimum(mean_tpr + std_tpr, 1)
tprs_lower = np.maximum(mean_tpr - std_tpr, 0)
plt.fill_between(mean_fpr, tprs_lower, tprs_upper, color='grey', alpha=.2, label=r' $\pm$  1 std. dev.')
plt.xlim([-0.05, 1.05])
plt.ylim([-0.05, 1.05])
plt.xlabel('False Positive Rate')
plt.ylabel('True Positive Rate')
plt.title('Logistic Regression ROC\n')
plt.legend(loc="lower right", fontsize=6)
plt.show()
#logistic_plot.savefig('results/figures/Logit_Baxter.png', dpi=1000)

```

Logistic Regression ROC



Predict using the Logistic Regression model on the test set

```
y_pred = logreg.predict(x_test)
print('Accuracy of logistic regression classifier on test set: {:.2f}'.format(logreg.score(x_test, y_test)))
```

[Output] Accuracy of logistic regression classifier on test set: 67.80 %

Now let's define a Multi-layer Perceptron Neural Network Model

```
clf = MLPClassifier(activation='logistic', alpha=0.001, batch_size='auto',
                    beta_1=0.9, beta_2=0.999, early_stopping=True, epsilon=1e-08,
                    hidden_layer_sizes=(100,), learning_rate='adaptive',
                    learning_rate_init=0.001, max_iter=200, momentum=0.9,
                    nesterovs_momentum=True, power_t=0.5, random_state=1, shuffle=True,
                    solver='sgd', tol=0.0001, validation_fraction=0.1, verbose=False,
                    warm_start=False)
```

Plot ROC curve for MLP on training set

We split the training set to 10 to cross-validate

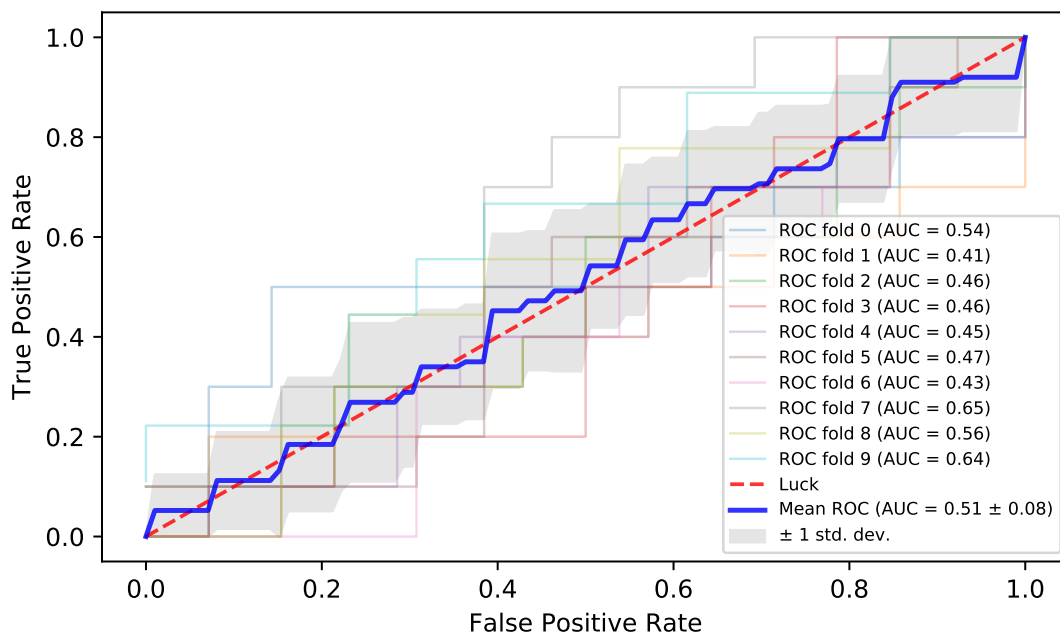
```
cv = StratifiedKFold(n_splits=10)
tprs = []
aucs = []
mean_fpr = np.linspace(0, 1, 100)
mlp_plot = plt.figure()
i = 0
for train, test in cv.split(X, Y):
    probas_ = clf.fit(X[train], Y[train]).predict_proba(X[test])
    # Compute ROC curve and area the curve
```

```

fpr, tpr, thresholds = roc_curve(Y[test], probas[:, 1])
tprs.append(interp(mean_fpr, fpr, tpr))
tprs[-1][0] = 0.0
roc_auc = auc(fpr, tpr)
aucs.append(roc_auc)
plt.plot(fpr, tpr, lw=1, alpha=0.3, label='ROC fold %d (AUC = %0.2f)' % (i, roc_auc))
i += 1
plt.plot([0, 1], [0, 1], linestyle='--', color='r', label='Luck', alpha=.8)
mean_tpr = np.mean(tprs, axis=0)
mean_tpr[-1] = 1.0
mean_auc = auc(mean_fpr, mean_tpr)
std_auc = np.std(aucs)
plt.plot(mean_fpr, mean_tpr, color='b', label=r'Mean ROC (AUC = %0.2f  $\pm$  %0.2f)' % (mean_auc, std_auc))
std_tpr = np.std(tprs, axis=0)
tprs_upper = np.minimum(mean_tpr + std_tpr, 1)
tprs_lower = np.maximum(mean_tpr - std_tpr, 0)
plt.fill_between(mean_fpr, tprs_lower, tprs_upper, color='grey', alpha=.2, label=r' $\pm$  1 std. dev.')
plt.xlim([-0.05, 1.05])
plt.ylim([-0.05, 1.05])
plt.xlabel('False Positive Rate')
plt.ylabel('True Positive Rate')
plt.title('Neural Network ROC\n')
plt.legend(loc="lower right", fontsize=7)
plt.show()
#mlp_plot.savefig('results/figures/MLP_Baxter.png', dpi=1000)

```

Neural Network ROC



Predict using the MLP model on the test set

```

y_pred = clf.predict(x_test)
## Print accuracy

```

```
print("Performance Accuracy on the Testing data:", round(clf.score(x_test, y_test) *100))
```

[Output] Performance Accuracy on the Testing data: 66.0

```
print("Number of correct classifiers:", round(accuracy_score(y_test, y_pred, normalize=False)))
```

[Output] Number of correct classifiers: 39

Now let's define a Random Forest model

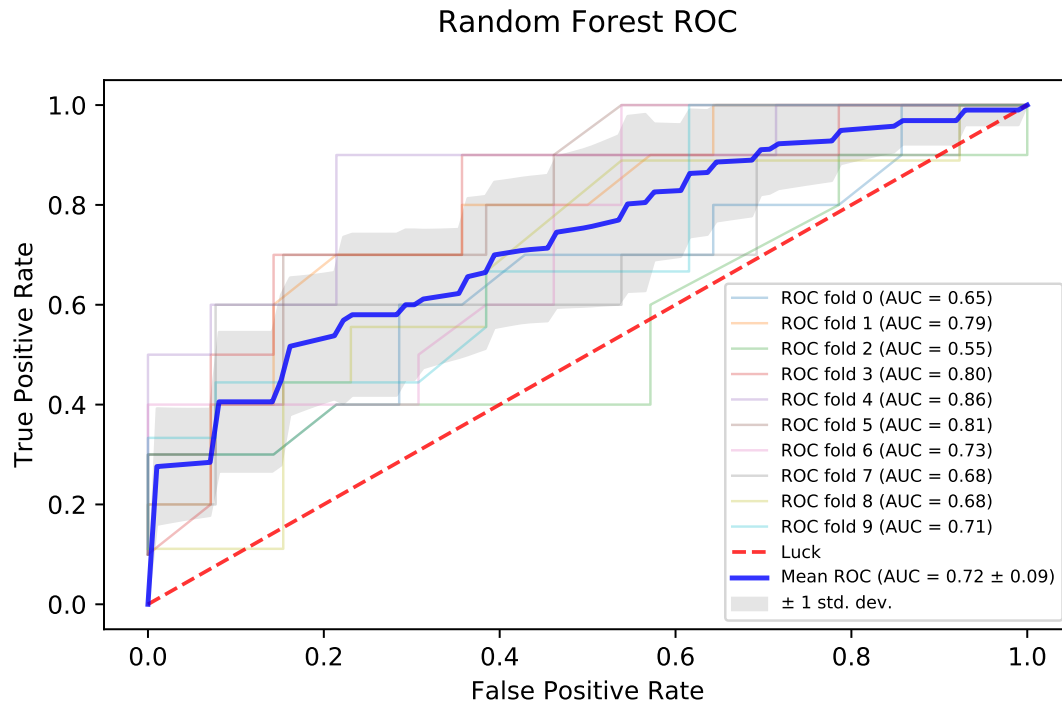
```
rfc = RandomForestClassifier(bootstrap=True, class_weight=None, criterion='gini',
                             max_depth=None, max_features='auto', max_leaf_nodes=None,
                             min_impurity_split=None, min_samples_leaf=1,
                             min_samples_split=2, min_weight_fraction_leaf=0.0,
                             n_estimators=150, n_jobs=1, oob_score=True, random_state=None,
                             verbose=0, warm_start=False)
```

Plot ROC curve for Random Forest on training set

We split the training set to 10 to cross-validate

```
RF_plot = plt.figure()
cv = StratifiedKFold(n_splits=10)
tprs = []
aucs = []
mean_fpr = np.linspace(0, 1, 100)
i = 0
for train, test in cv.split(X,Y):
    probas_ = rfc.fit(X[train], Y[train]).predict_proba(X[test])
    # Compute ROC curve and area the curve
    fpr, tpr, thresholds = roc_curve(Y[test], probas[:, 1])
    tprs.append(interp(mean_fpr, fpr, tpr))
    tprs[-1][0] = 0.0
    roc_auc = auc(fpr, tpr)
    aucs.append(roc_auc)
    plt.plot(fpr, tpr, lw=1, alpha=0.3, label='ROC fold %d (AUC = %0.2f)' % (i, roc_auc))
    i += 1
plt.plot([0, 1], [0, 1], linestyle='--', color='r', label='Luck', alpha=.8)
mean_tpr = np.mean(tprs, axis=0)
mean_tpr[-1] = 1.0
mean_auc = auc(mean_fpr, mean_tpr)
std_auc = np.std(aucs)
plt.plot(mean_fpr, mean_tpr, color='b', label=r'Mean ROC (AUC = %0.2f  $\pm$  %0.2f)' % (mean_auc, std_auc),
         lw=2, alpha=.8)
std_tpr = np.std(tprs, axis=0)
tprs_upper = np.minimum(mean_tpr + std_tpr, 1)
tprs_lower = np.maximum(mean_tpr - std_tpr, 0)
plt.fill_between(mean_fpr, tprs_lower, tprs_upper, color='grey', alpha=.2, label=r'$\pm$ 1 std. dev.')
plt.xlim([-0.05, 1.05])
plt.ylim([-0.05, 1.05])
plt.xlabel('False Positive Rate')
plt.ylabel('True Positive Rate')
plt.title('Random Forest ROC\n')
plt.legend(loc="lower right", fontsize=7)
```

```
plt.show()
#RF_plot.savefig('results/figures/Random_Forest_Baxter.png', dpi=1000)
```



Predict using the Random Forest model on the test set

```
y_pred = rfc.predict(x_test)
print("Performance Accuracy on the Testing data:", round(rfc.score(x_test, y_test) * 100))
```

[Output] Performance Accuracy on the Testing data: 71.0

```
print("Number of correct classifiers:", round(accuracy_score(y_test, y_pred, normalize=False)))
```

[Output] Number of correct classifiers: 42