

all_data

February 26, 2020

```
[1]: import os
from IPython.display import display, Image
import pandas as pd
import numpy as np
import seaborn as sns
%matplotlib inline
import matplotlib.pyplot as plt
from matplotlib import colors
from matplotlib.ticker import PercentFormatter
from scipy.stats import linregress
import math
from functools import reduce
import matplotlib
import argparse
from Bio import SeqIO, Entrez, pairwise2
Entrez.email = 'hongyingsun1101@gmail.com'
from Bio.SeqRecord import SeqRecord
import re, time
import os, sys, glob
import random
import uuid
# from skbio.tree import TreeNode
# from skbio import read
# from skbio.stats.distance import DistanceMatrix
# from skbio.stats.distance import DissimilarityMatrix

from scipy import stats
from ast import literal_eval
import sqlite3
# roc curve and auc score
from sklearn.datasets import make_classification
from sklearn.neighbors import KNeighborsClassifier
from sklearn.ensemble import RandomForestClassifier
from sklearn.model_selection import train_test_split
from sklearn.metrics import roc_curve
from sklearn.metrics import roc_auc_score
import warnings
```

```

warnings.filterwarnings("ignore")

[2]: df = pd.read_csv("all_data.csv", index_col=0)
score= pd.read_csv("score_merged.csv", index_col=0)

[3]: reference ={'A':"RDP_10398", 'B':"RDP_5224", 'C':"RDP_1017", 'D':"RDP_92", 'E':
→'RDP_12'}

[4]: #pplacer_ref_list = ['A','B','C',␣
→'D','E'],pplacer_stats_list=['_adcl_log','_edpl',␣
→'_prichness'],community_list=['0','1','2','3','4'],cutoff_list=['mean','min','25%','50%', '7

[5]: def is_float(string):
    try:
        return float(string) and '.' in string # True if string is a number␣
→contains a dot
    except ValueError: # String is not a number
        return False

[6]: reference

[6]: {'A': 'RDP_10398',
      'B': 'RDP_5224',
      'C': 'RDP_1017',
      'D': 'RDP_92',
      'E': 'RDP_12'}

[7]: columnList=list(df.columns)

[8]: communityList = df.community

[9]: # df.describe()

[10]: def plot_pplacer(variable):
    fig, axes = plt.subplots(nrows=3, ncols=2)
    ax0, ax1, ax2, ax3, ax4, ax5 = axes.flatten()

    ax0.hist(df['A'+variable])
    ax0.set_title('A'+variable)

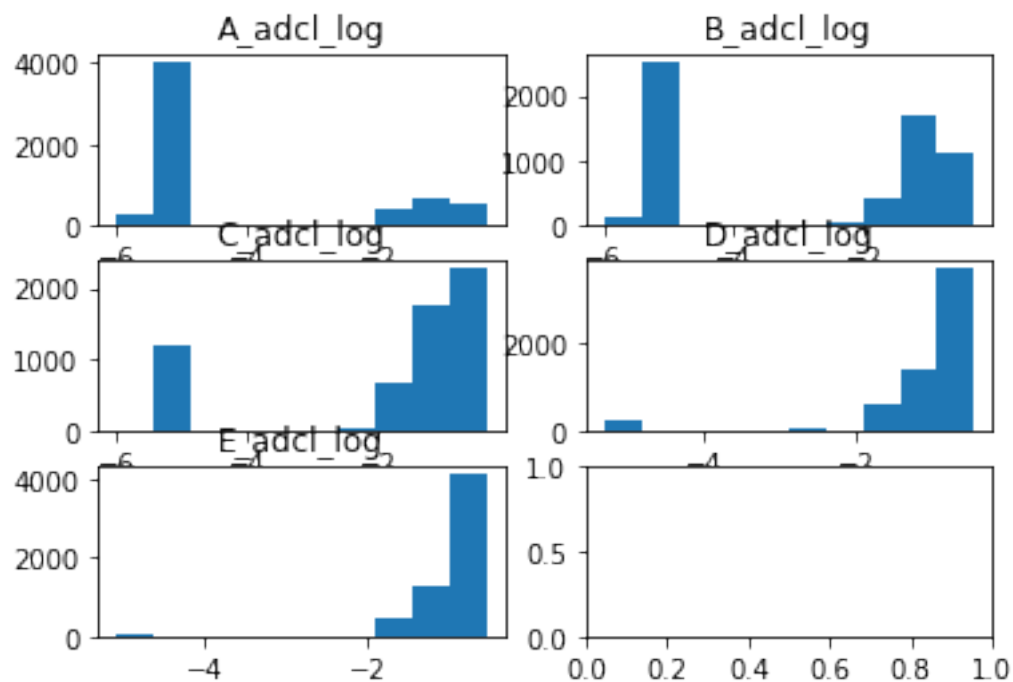
    ax1.hist(df['B'+variable])
    ax1.set_title('B'+variable)

    ax2.hist(df['C'+variable])
    ax2.set_title('C'+variable)

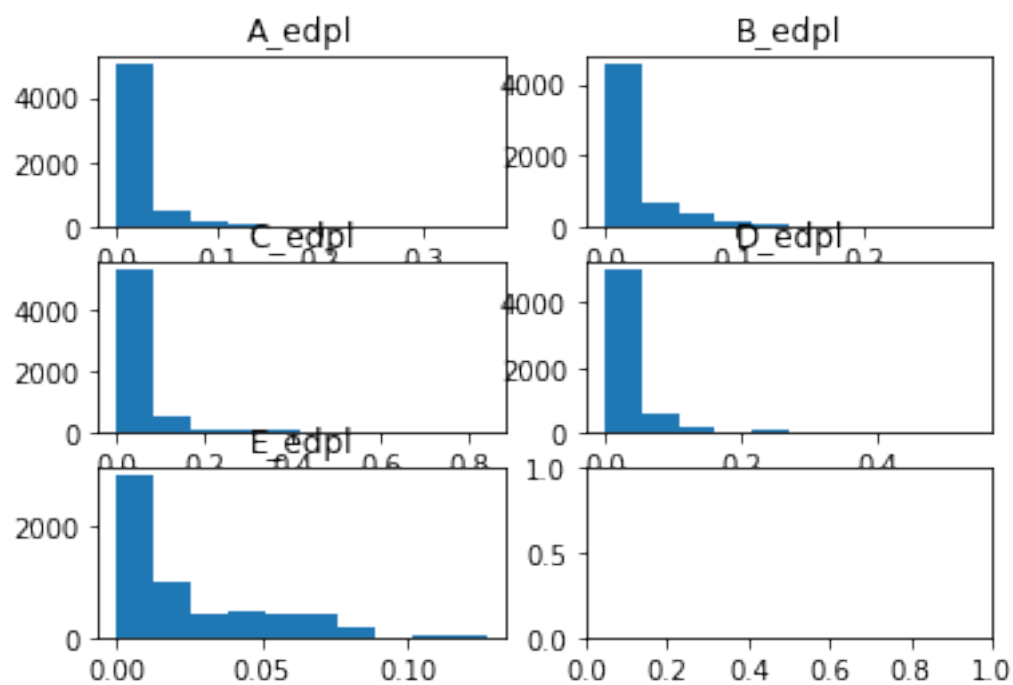
    ax3.hist(df['D'+variable])
    ax3.set_title('D'+variable)
    ax4.hist(df['E'+variable])
    ax4.set_title('E'+variable)

[11]: plot_pplacer('_adcl_log')

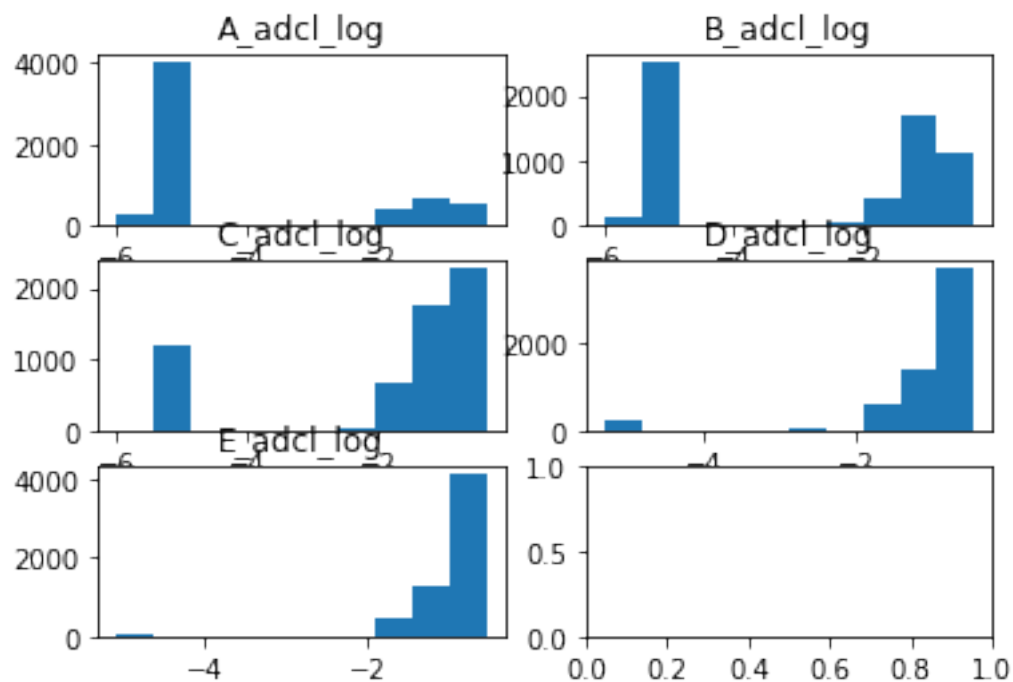
```



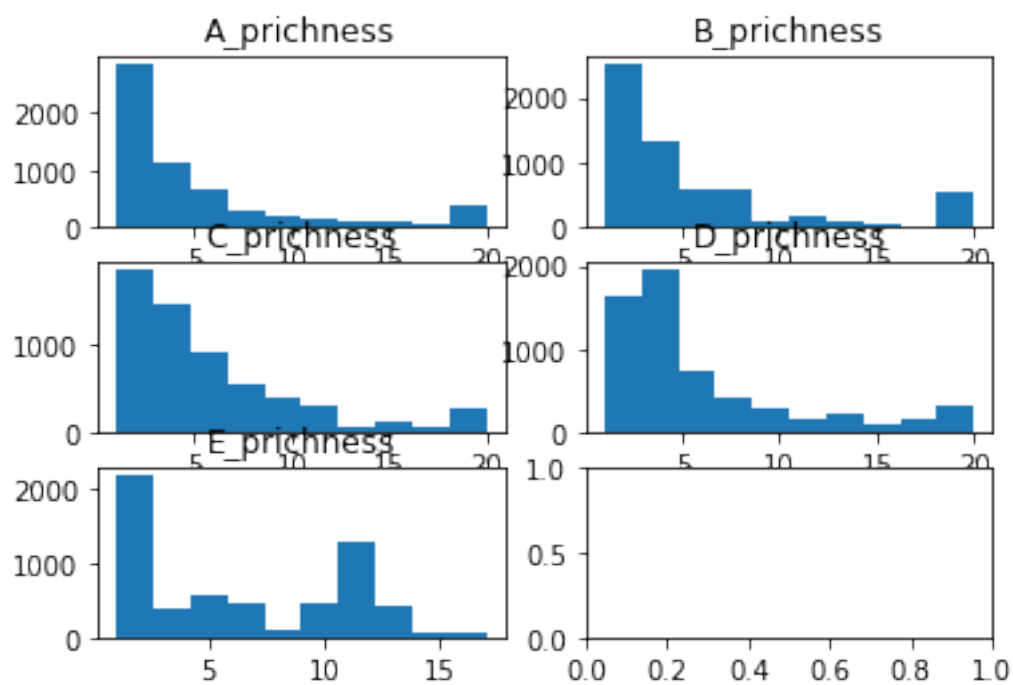
```
[12]: plot_pplacer('_edpl')
```



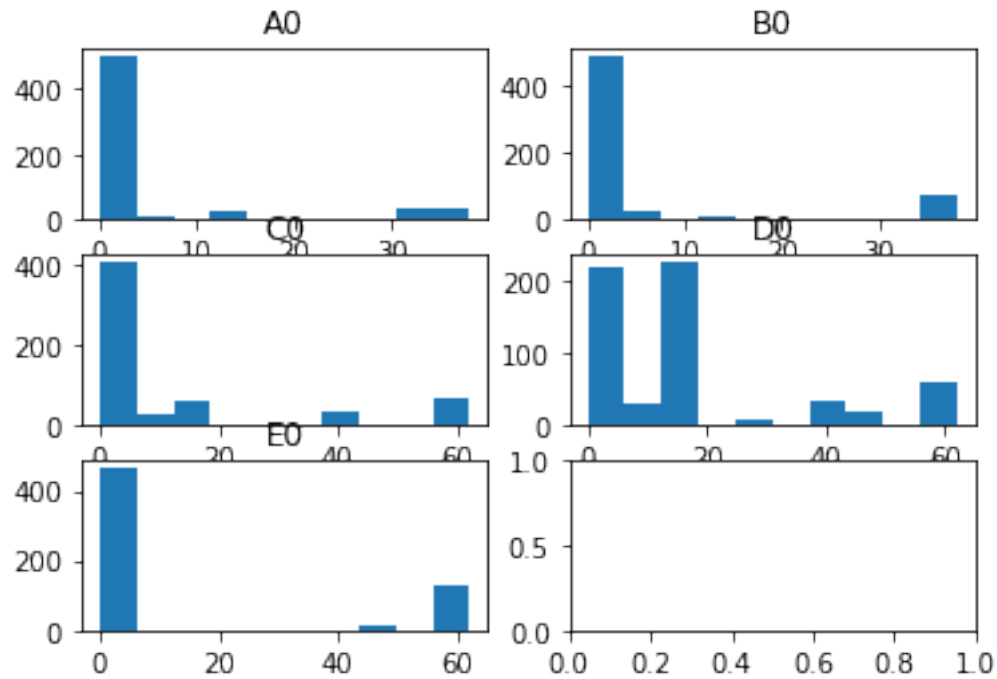
```
[13]: plot_pplacer('_adcl_log')
```



```
[14]: plot_pplacer('_prichness')
```



```
[15]: plot_pplacer('0');
```



```
[16]: def plotScatter(reference,community):
fig, axes = plt.subplots(nrows=2, ncols=2)
ax0, ax1, ax2, ax3 = axes.flatten()

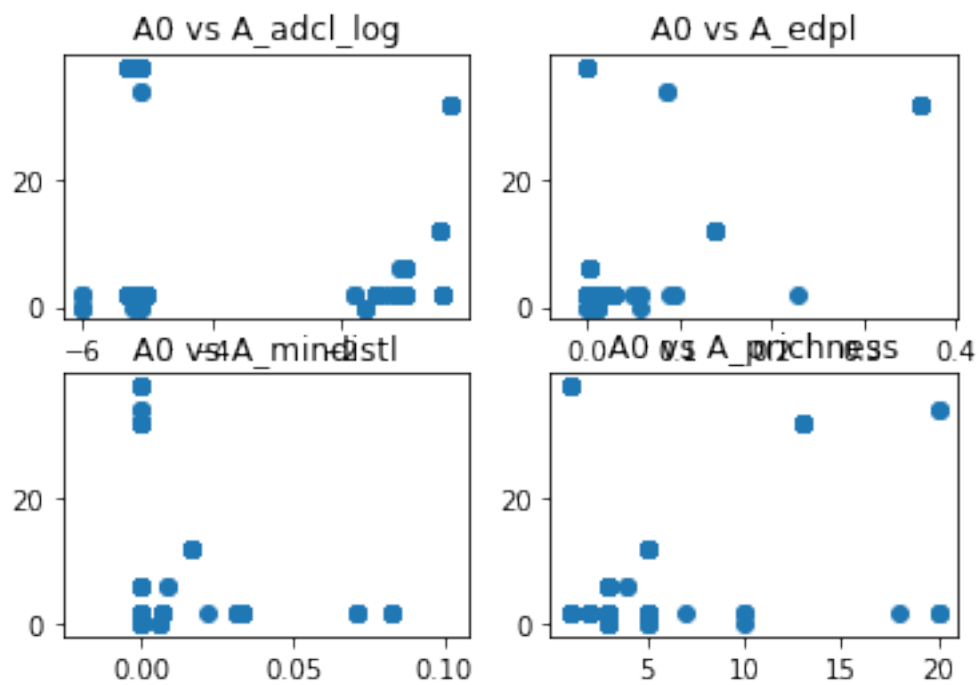
ax0.scatter(df[reference+'_adcl_log'], df[reference+community])
ax0.set_title(reference+community+' vs '+ reference + '_adcl_log')

ax1.scatter(df[reference+'_edpl'], df[reference+community])
ax1.set_title(reference+community+' vs '+ reference + '_edpl')

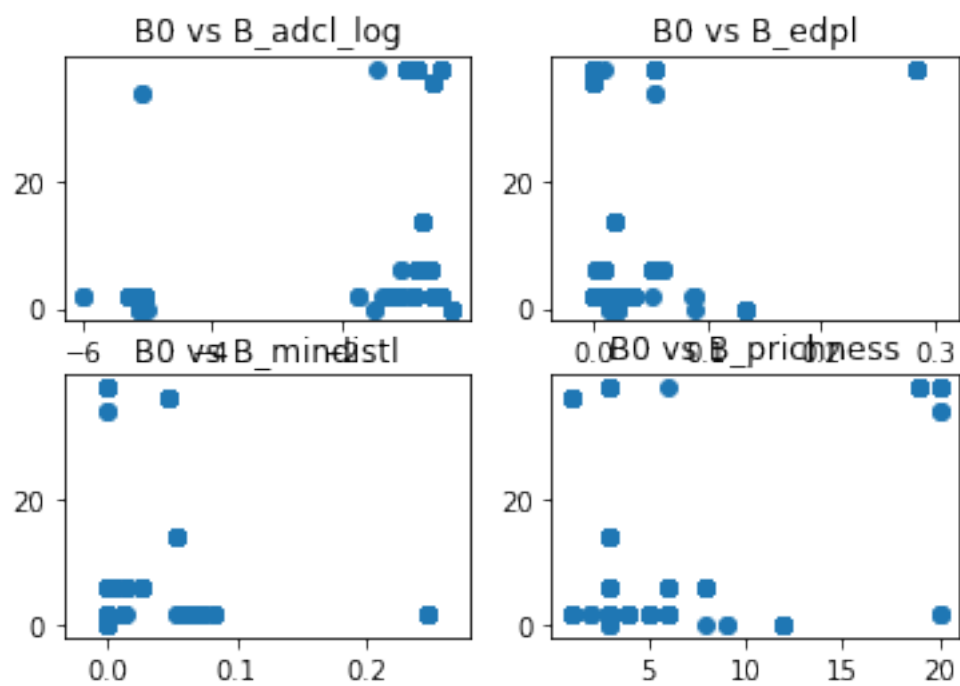
ax2.scatter(df[reference+'_mindistl'], df[reference+community])
ax2.set_title(reference+community+' vs '+ reference + '_mindistl')

ax3.scatter(df[reference+'_prichness'], df[reference+community])
ax3.set_title(reference+community+' vs '+ reference + '_prichness')
```

```
[17]: plotScatter('A','0')
```



[18]: `plotScatter('B','O')`



```
[19]: def plotScatterRef(variable,community):
fig, axes = plt.subplots(nrows=3, ncols=2)
ax0, ax1, ax2, ax3, ax4, ax5 = axes.flatten()

ax0.scatter(df['A'+variable], df['A'+community])
ax0.set_title('A' + community + ' vs A' + variable)
ax1.scatter(df['B'+variable], df['B'+community])
ax1.set_title('B' + community + ' vs B' + variable)

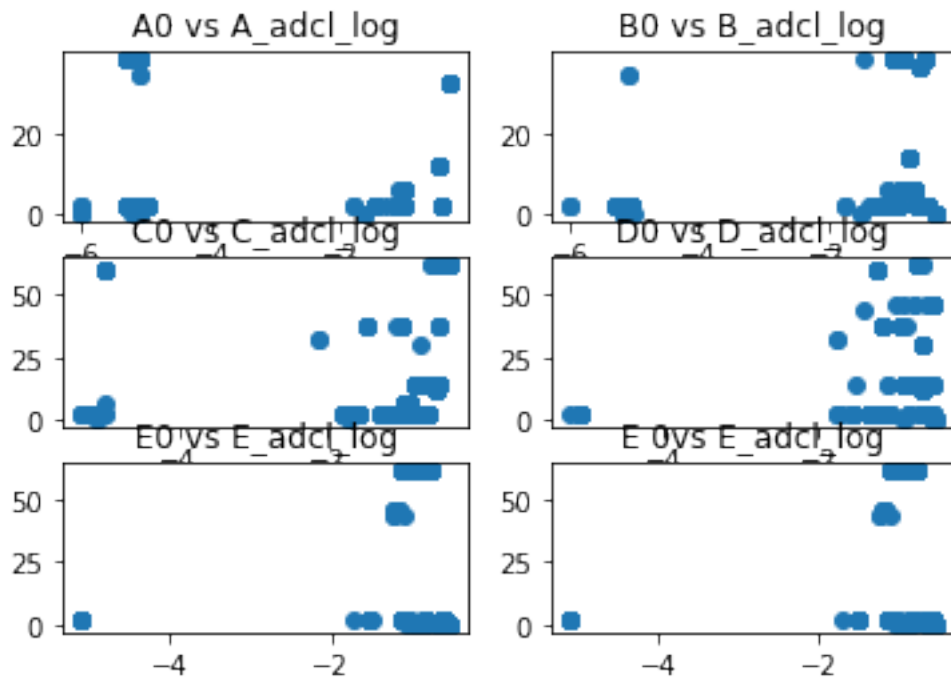
ax2.scatter(df['C'+variable], df['C'+community])
ax2.set_title('C' + community + ' vs C' + variable)

ax3.scatter(df['D'+variable], df['D'+community])
ax3.set_title('D' + community + ' vs D' + variable)

ax4.scatter(df['E'+variable], df['E'+community])
ax4.set_title('E' + community + ' vs E' + variable)

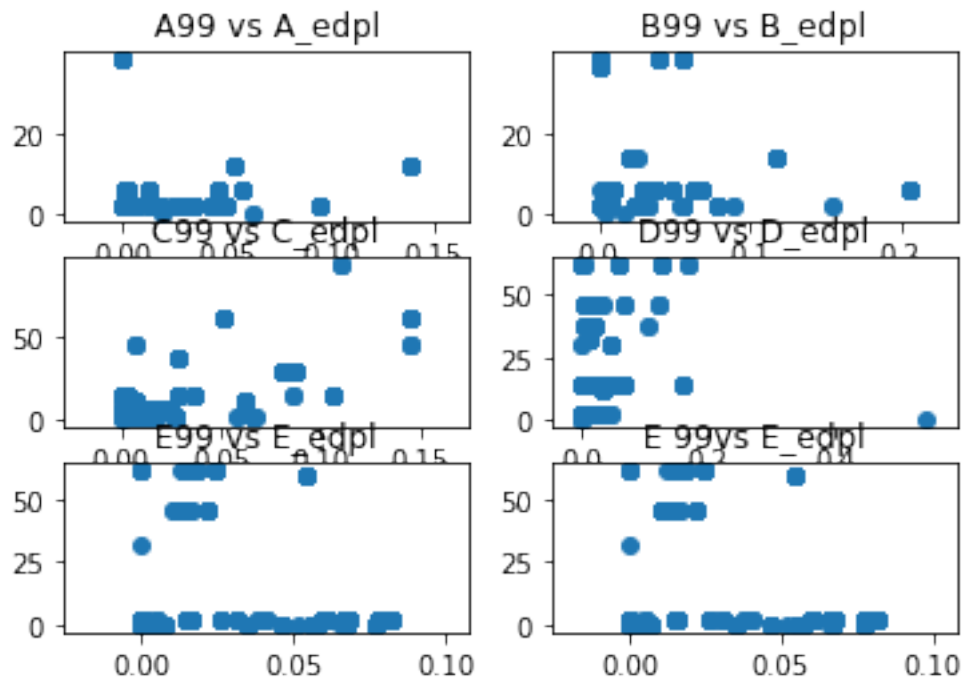
ax5.scatter(df['E'+variable], df['E'+community])
ax5.set_title('E ' + community + 'vs E' + variable)

plotScatterRef('_adcl_log','0');
```



[20]:

```
plotScatterRef('_edpl','99');
```



[21]:

```
cols=df.columns.tolist()
# cols[:20]
```

[164]:

```
def plot_roc_curve(fpr, tpr):
    plt.plot(fpr, tpr, color='orange', label='ROC')
    plt.plot([0, 1], [0, 1], color='darkblue', linestyle='--')
    plt.xlabel('False Positive Rate')
    plt.ylabel('True Positive Rate')
    plt.title('Receiver Operating Characteristic (ROC) Curve')
    plt.legend()
    plt.show()
```

[169]:

```
def plot_roc(data_X, class_label):
    trainX, testX, trainy, testy = train_test_split(data_X, class_label,
    →test_size=0.3, random_state=1)
    model = RandomForestClassifier()
    model.fit(trainX, trainy)
    probs = model.predict_proba(testX)
    probs = probs[:, 1]
    auc = roc_auc_score(testy, probs)
    fpr, tpr, thresholds = roc_curve(testy, probs)
    optimal_idx = np.argmax(tpr - fpr)
```



```

    optimal_threshold = thresholds[optimal_idx]
    print('optimal_threshold: %.2f' % optimal_threshold)
    print('AUC: %.2f' % auc)
    print(thresholds)
#     print(thresholds)
#     print('Model: ')
#     print(model)
plot_roc_curve(fpr, tpr)

```

```

[170]: def makeTable(headerRow, columnizedData, columnSpacing=2):
        """Creates a technical paper style, left justified table"""
        from numpy import array, max, vectorize

        cols = array(columnizedData, dtype=str)
        colSizes = [max(vectorize(len)(col)) for col in cols]

        header = ''
        rows = ['' for i in cols[0]]

        for i in range(0, len(headerRow)):
            if len(headerRow[i]) > colSizes[i]: colSizes[i] = len(headerRow[i])
            headerRow[i] += ' ' * (colSizes[i] - len(headerRow[i]))
            header += headerRow[i]
            if not i == len(headerRow) - 1: header += ' ' * columnSpacing

            for j in range(0, len(cols[i])):
                if len(cols[i][j]) < colSizes[i]:
                    cols[i][j] += ' ' * (colSizes[i] - len(cols[i][j]) + columnSpacing)
                rows[j] += cols[i][j]
                if not i == len(headerRow) - 1: rows[j] += ' ' * columnSpacing

        line = '-' * len(header)
        print(line)
        print(header)
        print(line)
        for row in rows: print(row)
        print(line)
    header = ['AUROC', 'Categoroy']
    cutoffs = ['0.9-1.0', '0.8-0.9', '0.7-0.8', '0.6-0.7', '0.5-0.6']
    evaluation = ['Very good', 'Good', 'Fair', 'Poor', 'Fail']
    makeTable(header, [cutoffs, evaluation])

```

```

-----
AUROC    Categoroy
-----
0.9-1.0  Very good
0.8-0.9  Good
0.7-0.8  Fair

```

0.6-0.7 Poor
0.5-0.6 Fail

```
[171]: def plot_roc_microbiome(data_X, class_label, x, y, data_test=False):
        if(not data_test):
            #         print("data_set is False")
            trainX, testX, trainy, testy = train_test_split(data_X, class_label,
→test_size=0.3, random_state=1)
            model = RandomForestClassifier()
            model.fit(trainX, trainy)
            probs = model.predict_proba(testX)
            probs = probs[:, 1]
            auc = roc_auc_score(testy, probs)
            fpr, tpr, thresholds = roc_curve(testy, probs)
            optimal_idx = np.argmax(tpr - fpr)
            optimal_threshold = thresholds[optimal_idx]
            print('optimal_threshold: %.2f' % optimal_threshold)
            print('AUC: %.2f' % auc)
            print('thresholds: ' + thresholds)

            plot_roc_curve(fpr, tpr)

        else:
            print("data_set is True")
            trainX, testX, trainy, testy = train_test_split(data_X, class_label,
→test_size=0.3, random_state=1)
            model = RandomForestClassifier()
            model.fit(trainX, trainy)
            probs1 = model.predict_proba(testX)
            probs1 = probs1[:, 1]
            auc1 = roc_auc_score(testy, probs1)
            fpr1, tpr1, thresholds1 = roc_curve(testy, probs1)
            optimal_idx1 = np.argmax(tpr1 - fpr1)
            optimal_threshold1 = thresholds1[optimal_idx1]
            print('AUC1: %.2f' % auc1)
            print('optimal_threshold1: %.2f' % optimal_threshold1)
            print(thresholds1)
            plot_roc_curve(fpr1, tpr1)

            probs2 = model.predict_proba(x)
            probs2 = probs2[:, 1]
            auc2 = roc_auc_score(y, probs2)
            fpr2, tpr2, thresholds2 = roc_curve(y, probs2)
            optimal_idx2 = np.argmax(tpr2 - fpr2)
            optimal_threshold2 = thresholds2[optimal_idx2]
            print('AUC2: %.2f' % auc2)
```

```

print('optimal_threshold2: %.2f' % optimal_threshold2)
print( thresholds2)
plot_roc_curve(fpr2, tpr2)

```

```

[172]: def plot_roc_curve_microbiome(pplacer_ref_list, pplacer_stats_list,
    community_list, cutoff_list, scoreOption=True):
    for (refIndex,pplacer_ref) in enumerate(pplacer_ref_list):
#         for refIndex in range(len(pplacer_ref_list)):
#             pplacer_ref = pplacer_ref_list[refIndex]
        for (statsIndex,pplacer_stats) in enumerate(pplacer_stats_list):
#             for statsIndex in range(len(pplacer_stats_list)):
#                 pplacer_stats = pplacer_stats_list[statsIndex]
            for (communityIndex,community) in enumerate(community_list):
#                 for communityIndex in range(len(community_list)):
#                     community = community_list[communityIndex]
                for (i, cutoff) in enumerate(cutoff_list):
#                     for i in range(len(cutoff_list)):
#                         cutoff=cutoff_list[i]
                    if(is_float(cutoff)):
                        cutoff_binary=float(cutoff)
                    else:
                        if(scoreOption):
                            cutoff_binary=float(df[pplacer_ref+community].
    community_list, cutoff_list, scoreOption=True):
    for (refIndex,pplacer_ref) in enumerate(pplacer_ref_list):
#         for refIndex in range(len(pplacer_ref_list)):
#             pplacer_ref = pplacer_ref_list[refIndex]
        for (statsIndex,pplacer_stats) in enumerate(pplacer_stats_list):
#             for statsIndex in range(len(pplacer_stats_list)):
#                 pplacer_stats = pplacer_stats_list[statsIndex]
            for (communityIndex,community) in enumerate(community_list):
#                 for communityIndex in range(len(community_list)):
#                     community = community_list[communityIndex]
                for (i, cutoff) in enumerate(cutoff_list):
#                     for i in range(len(cutoff_list)):
#                         cutoff=cutoff_list[i]
                    if(is_float(cutoff)):
                        cutoff_binary=float(cutoff)
                    else:
                        if(scoreOption):
                            cutoff_binary=float(df[pplacer_ref+community].
describe().loc[[cutoff]])

                        else:
                            cutoff_binary = float(df[pplacer_ref+pplacer_stats].
describe().loc[[cutoff]])

                        if(scoreOption):
                            mask = df[pplacer_ref+community] <= cutoff_binary
                            df.loc[mask, pplacer_ref+community+'_binary'] = 1
                            mask = df[pplacer_ref+community] >cutoff_binary
                            df.loc[mask, pplacer_ref+community+'_binary'] = 0
                            df_binary = df[[pplacer_ref+pplacer_stats,
pplacer_ref+community+'_binary']].dropna()
                            data_stats = df_binary[pplacer_ref+pplacer_stats].
to_numpy().reshape(-1,1)
                            binary_label = 
df_binary[pplacer_ref+community+'_binary'].to_numpy()
                            print(' The score cutoff ' + cutoff + ' for Reference ' +
pplacer_ref + ' community ' + community + ' with pplacer_stats '+
pplacer_stats[1:] + ': %.2f' % cutoff_binary )
                            plot_roc(data_stats,binary_label)
                        else:
                            mask = df[pplacer_ref+pplacer_stats] <= cutoff_binary
                            df.loc[mask, pplacer_ref+pplacer_stats+'_binary'] = 1
                            mask = df[pplacer_ref+pplacer_stats] >cutoff_binary

```

```

        df.loc[mask, pplacer_ref+pplacer_stats+'_binary'] = 0
        df_binary = df[[pplacer_ref+community,
→pplacer_ref+pplacer_stats+'_binary']].dropna()
        data_stats = df_binary[pplacer_ref+community].
→to_numpy().reshape(-1,1)
        binary_label = 0
→df_binary[pplacer_ref+pplacer_stats+'_binary'].to_numpy()
        print(' The pplacer_stats_cutoff ' + cutoff + ' for
→Reference ' + pplacer_ref + ' community ' + community + ' pplacer_stats ' +
→pplacer_stats[1:] + ' : %.2f' % cutoff_binary )
        plot_roc(data_stats,binary_label)

```

1 different reference same pplacer stats same community to test different cutoffs and different references for score

```

[173]: # plot_roc_curve_microbiome(pplacer_ref_list =
→['A', 'B', 'C', 'D', 'E'],pplacer_stats_list=['_adcl_log'],community_list=['A'],cutoff_list=['m
→'min', '25%', '50%', '75%'],scoreOption=False)

```

```

[174]: df['E0'].describe()

```

```

[174]: count      605.000000
mean         14.601653
std          25.354082
min           0.000000
25%           0.000000
50%           2.000000
75%           2.000000
max          62.000000
Name: E0, dtype: float64

```

2 Different reference same pplacer stats same community to test different cutoffs and different references for adcl_log

3 Fitting on large reference and test on small reference datasets

```

[175]: def plot_roc_curve_microbiome_test2(pplacer_ref_list, pplacer_stats_list,
→community_list, cutoff_list, test_data_list, scoreOption=True,
→testOption=False):
    for refIndex in range(len(pplacer_ref_list)):
        pplacer_ref = pplacer_ref_list[refIndex]
        for statsIndex in range(len(pplacer_stats_list)):
            pplacer_stats = pplacer_stats_list[statsIndex]

```

```

for communityIndex in range(len(community_list)):
    community = community_list[communityIndex]
    for i in range(len(cutoff_list)):
        cutoff=cutoff_list[i]
        if(is_float(cutoff)):
            cutoff_binary=float(cutoff)
        else:
            if(scoreOption):
                cutoff_binary=float(df[pplacer_ref+community].
→describe().loc[[cutoff]])

            else:
                cutoff_binary = float(df[pplacer_ref+pplacer_stats].
→describe().loc[[cutoff]])

            # no test situation, which is the default option
            if (not testOption):

                if(scoreOption):
                    mask = df[pplacer_ref+community] <= cutoff_binary
                    df.loc[mask, pplacer_ref+community+'_binary'] = 1
                    mask = df[pplacer_ref+community] >cutoff_binary
                    df.loc[mask, pplacer_ref+community+'_binary'] = 0
                    df_binary = df[[pplacer_ref+pplacer_stats,
→pplacer_ref+community+'_binary']].dropna()
                    data_stats = df_binary[pplacer_ref+pplacer_stats].
→to_numpy().reshape(-1,1)
                    binary_label = 
→df_binary[pplacer_ref+community+'_binary'].to_numpy()
                    print(' The score cutoff ' + cutoff + ' for Reference
→' + pplacer_ref + ' community ' + community + ' with pplacer_stats '+
→pplacer_stats[1:] + ': %.2f' % cutoff_binary )
                    # plot_roc(data_stats,binary_label)
                    
                    plot_roc_microbiome(data_stats,binary_label,x=None,y=None,data_test=False)
                else:
                    mask = df[pplacer_ref+pplacer_stats] <= 
→cutoff_binary
                    df.loc[mask, pplacer_ref+pplacer_stats+'_binary'] =
→1
                    mask = df[pplacer_ref+pplacer_stats] >cutoff_binary
                    df.loc[mask, pplacer_ref+pplacer_stats+'_binary'] =
→0
                    df_binary = df[[pplacer_ref+community,
→pplacer_ref+pplacer_stats+'_binary']].dropna()
                    data_stats = df_binary[pplacer_ref+community].
→to_numpy().reshape(-1,1)

```

```

        binary_label = □
→df_binary[pplacer_ref+pplacer_stats+'_binary'].to_numpy()
        print(' The pplacer_stats_cutoff ' + cutoff + ' for□
→Reference ' + pplacer_ref + ' community ' + community + ' pplacer_stats ' +□
→pplacer_stats[1:] + ' : %.2f' % cutoff_binary )
        # plot_roc(data_stats,binary_label)
□
→plot_roc_microbiome(data_stats,binary_label,x=None,y=None,data_test=False)

        # if there is test
    else:
        for j in range(len(test_data_list)):
            test=test_data_list[j]
            if(scoreOption):
                mask = df[pplacer_ref+community] <= □
→cutoff_binary
                df.loc[mask, pplacer_ref+community+'_binary'] =□
→1
                mask = df[pplacer_ref+community] >cutoff_binary
                df.loc[mask, pplacer_ref+community+'_binary'] =□
→0
                df_binary = df[[pplacer_ref+pplacer_stats,□
→pplacer_ref+community+'_binary']].dropna()
                data_stats =□
→df_binary[pplacer_ref+pplacer_stats].to_numpy().reshape(-1,1)
                binary_label = □
→df_binary[pplacer_ref+community+'_binary'].to_numpy()

                mask_test = df[test+community] <= cutoff_binary
                df.loc[mask_test, test+community+'_binary'] = 1
                mask_test = df[test+community] >cutoff_binary
                df.loc[mask_test, test+community+'_binary'] = 0
                df_binary = df[[test+pplacer_stats,□
→test+community+'_binary']].dropna()
                x = df_binary[test+pplacer_stats].to_numpy().
→reshape(-1,1)
                y = df_binary[test+community+'_binary'].
→to_numpy()

                print(' The score cutoff ' + cutoff + ' for□
→Reference ' + pplacer_ref + ' community ' + community + ' with□
→pplacer_stats ' + pplacer_stats[1:] + ' compared with test ' + test + ' : %.2f'□
→% cutoff_binary )
□
→plot_roc_microbiome(data_stats,binary_label,x,y,data_test=True)
        else:

```

```

mask = df[pplacer_ref+pplacer_stats] <= 0
→cutoff_binary
df.loc[mask,
pplacer_ref+pplacer_stats+'_binary'] = 1
mask = df[pplacer_ref+pplacer_stats]
→>cutoff_binary
df.loc[mask,
pplacer_ref+pplacer_stats+'_binary'] = 0
df_binary = df[[pplacer_ref+community,
pplacer_ref+pplacer_stats+'_binary']].dropna()
data_stats = df_binary[pplacer_ref+community].
→to_numpy().reshape(-1,1)
binary_label = 0
→df_binary[pplacer_ref+pplacer_stats+'_binary'].to_numpy()

mask_test = df[test+pplacer_stats] <= 0
→cutoff_binary
df.loc[mask_test, test+pplacer_stats+'_binary']
→= 1
mask_test = df[test+pplacer_stats]
→>cutoff_binary
df.loc[mask_test, test+pplacer_stats+'_binary']
→= 0
df_binary = df[[test+community,
test+pplacer_stats+'_binary']].dropna()
x = df_binary[test+community].to_numpy().
→reshape(-1,1)
y = df_binary[test+pplacer_stats+'_binary'].
→to_numpy()

print(' The pplacer_stats_cutoff ' + cutoff + '
→for Reference ' + pplacer_ref + ' community ' + community + ' pplacer_stats '
→ + pplacer_stats[1:] + ' compared with test ' + test + ': %.2f' %
→cutoff_binary )

→
→plot_roc_microbiome(data_stats,binary_label,x,y,data_test=True)

```

3.1 Model from larger reference sets to fit data used small reference set. Could be worse on both directions

```

[176]: plot_roc_curve_microbiome_test2(pplacer_ref_list =
→['A'],pplacer_stats_list=['_adcl_log','_edpl'],community_list=['0'],cutoff_list=['10.
→00'], test_data_list=['B','C','D','E'],testOption=True, scoreOption=True)

```

The score cutoff 10.00 for Reference A community 0 with pplacer_stats adcl_log

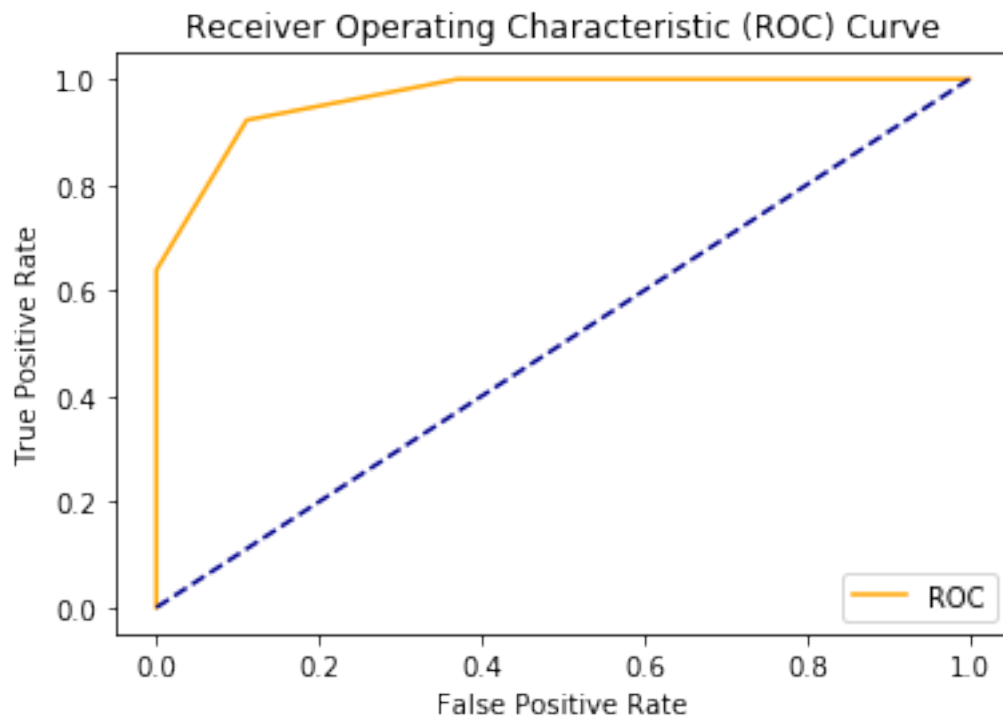
compared with test B: 10.00

data_set is True

AUC1: 0.97

optimal_threshold1: 0.84

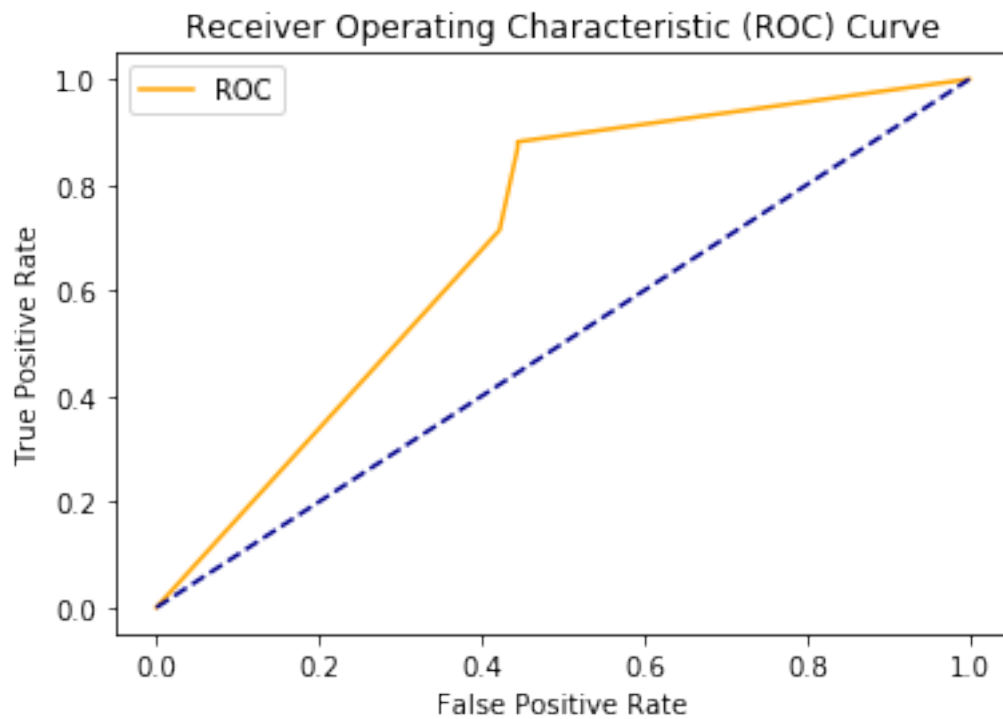
[2. 1. 0.84006066 0.59846994 0.]



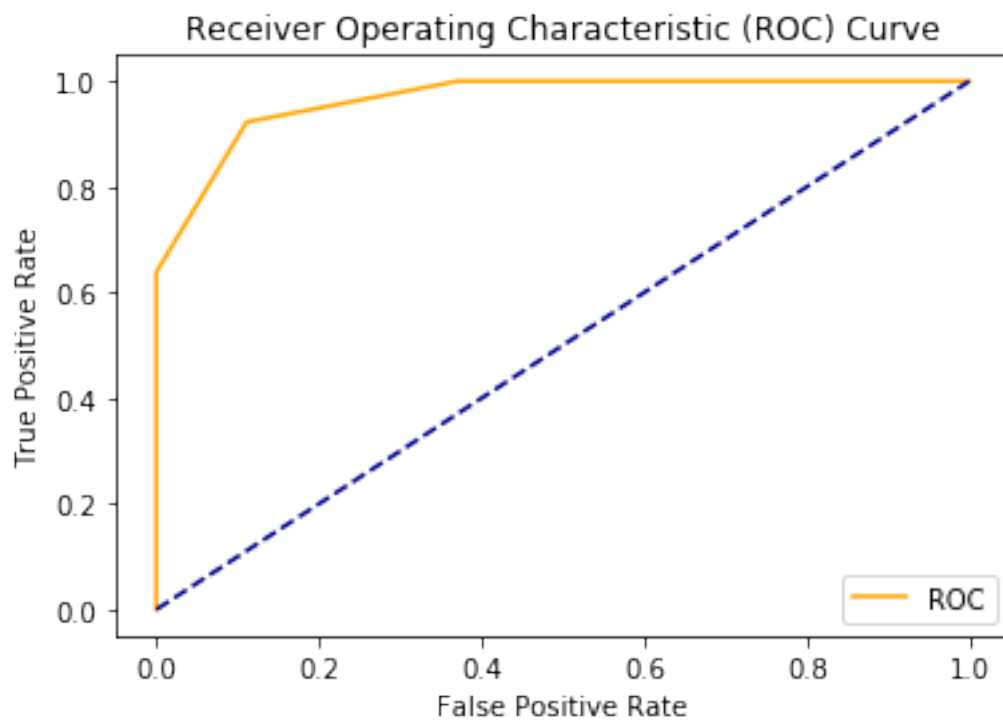
AUC2: 0.69

optimal_threshold2: 0.60

[2. 1. 0.84006066 0.59846994 0.]



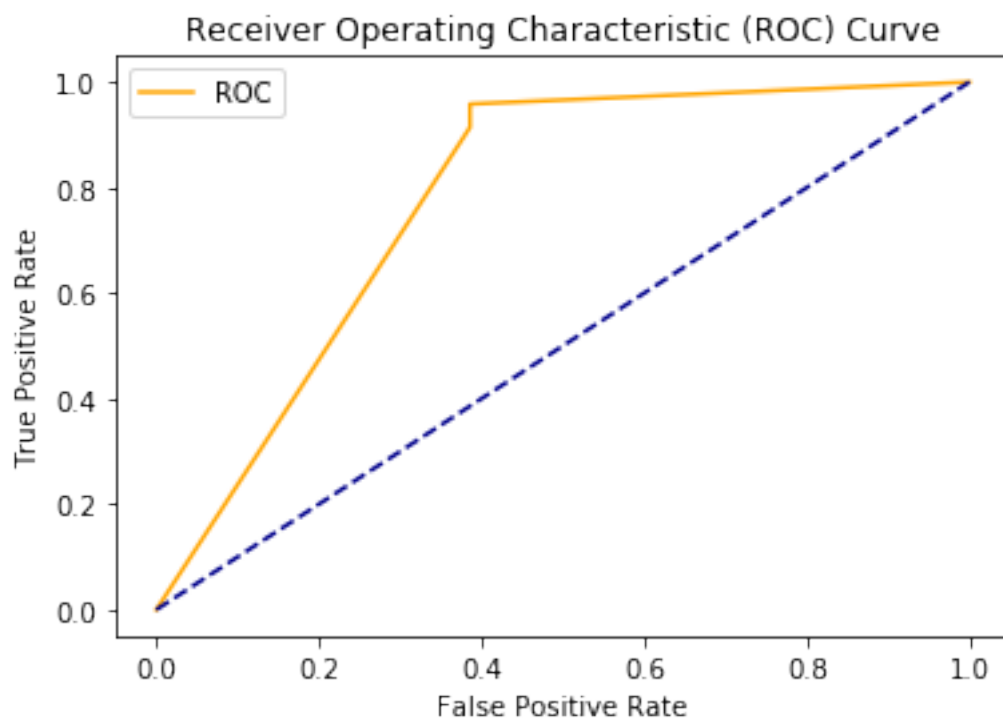
The score cutoff 10.00 for Reference A community 0 with pplacer_stats adcl_log
 compared with test C: 10.00
 data_set is True
 AUC1: 0.97
 optimal_threshold1: 0.84
 [2. 1. 0.83861265 0.6011308 0.]



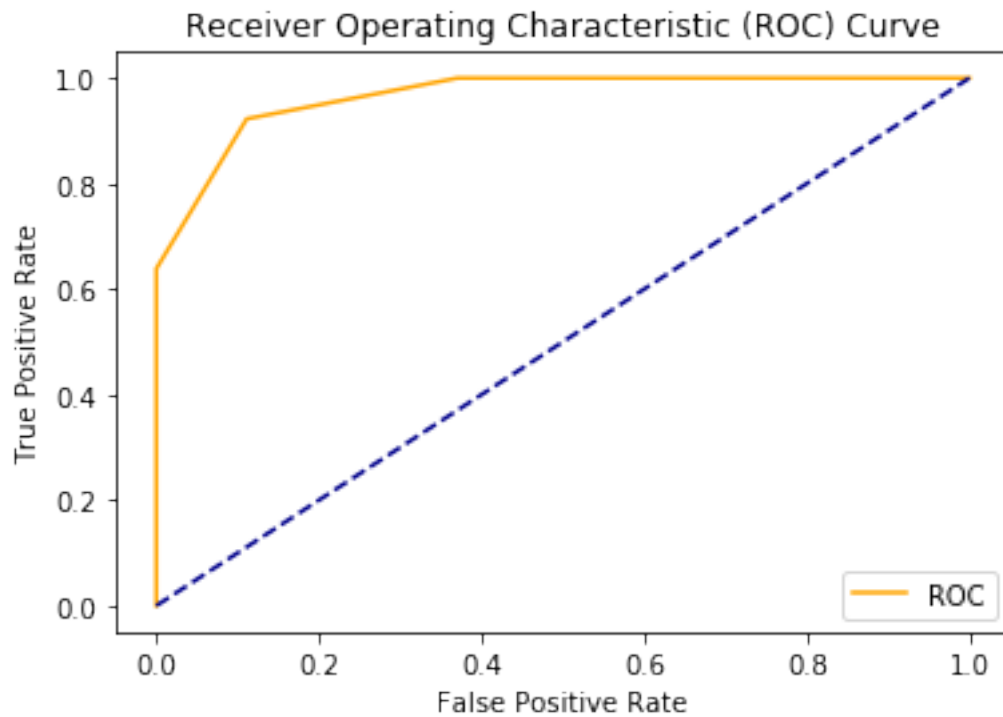
AUC2: 0.78

optimal_threshold2: 0.60

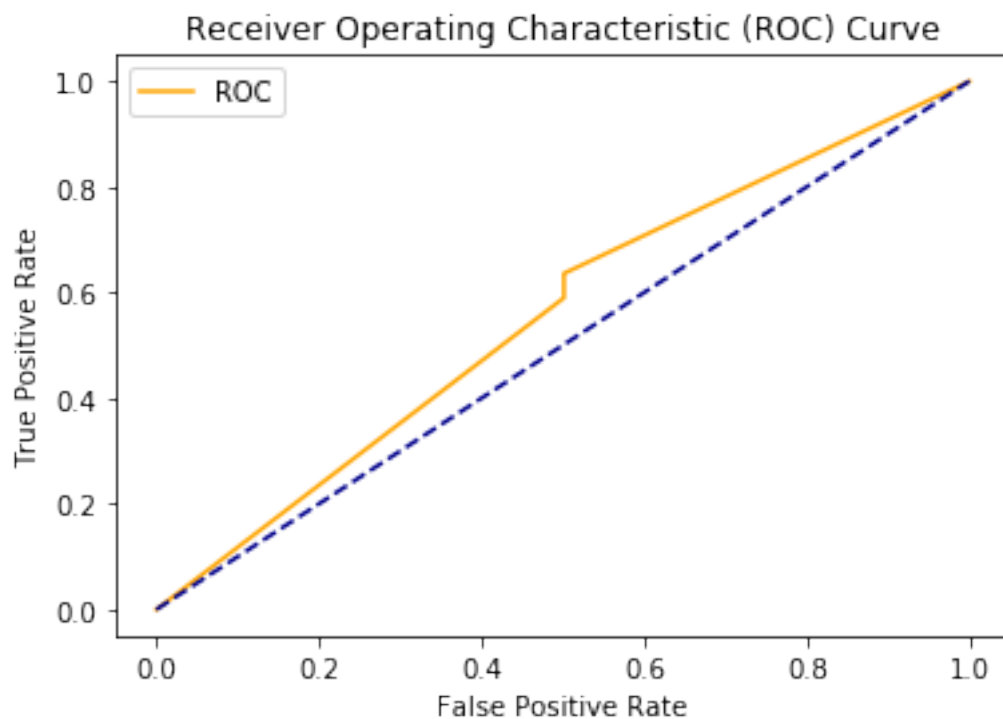
[2. 1. 0.83861265 0.6011308 0.]



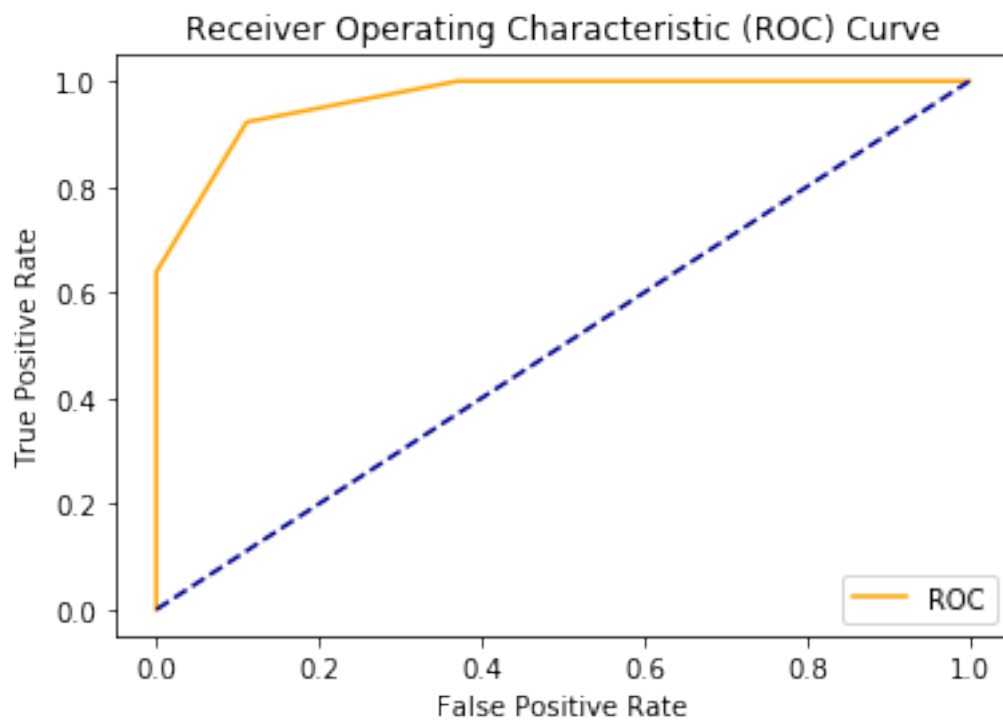
The score cutoff 10.00 for Reference A community 0 with pplacer_stats adcl_log
 compared with test D: 10.00
 data_set is True
 AUC1: 0.97
 optimal_threshold1: 0.84
 [2. 1. 0.83522539 0.58105212 0.]]



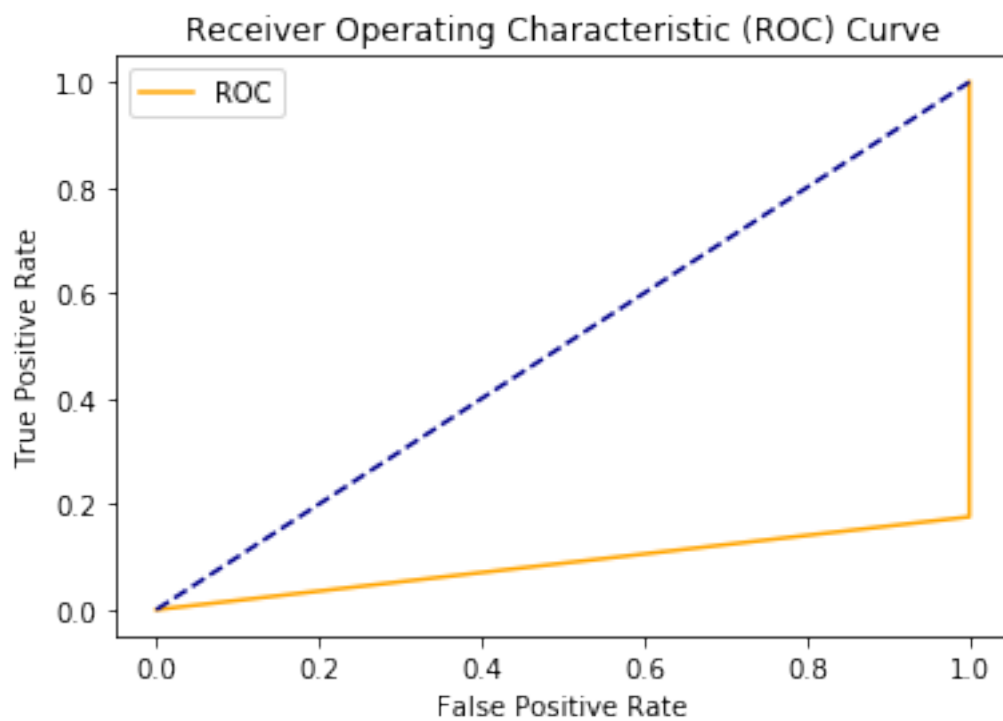
AUC2: 0.56
 optimal_threshold2: 0.84
 [2. 1. 0.83522539 0.]]



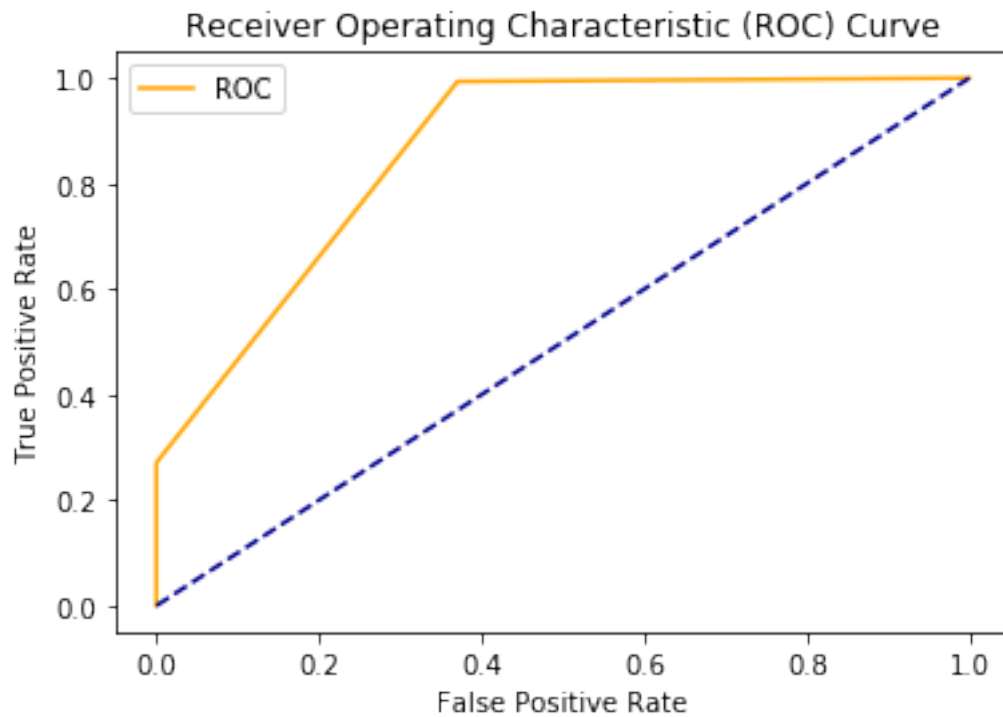
The score cutoff 10.00 for Reference A community 0 with pplacer_stats adcl_log
 compared with test E: 10.00
 data_set is True
 AUC1: 0.97
 optimal_threshold1: 0.86
 [2. 1. 0.85596766 0.60895695 0.]



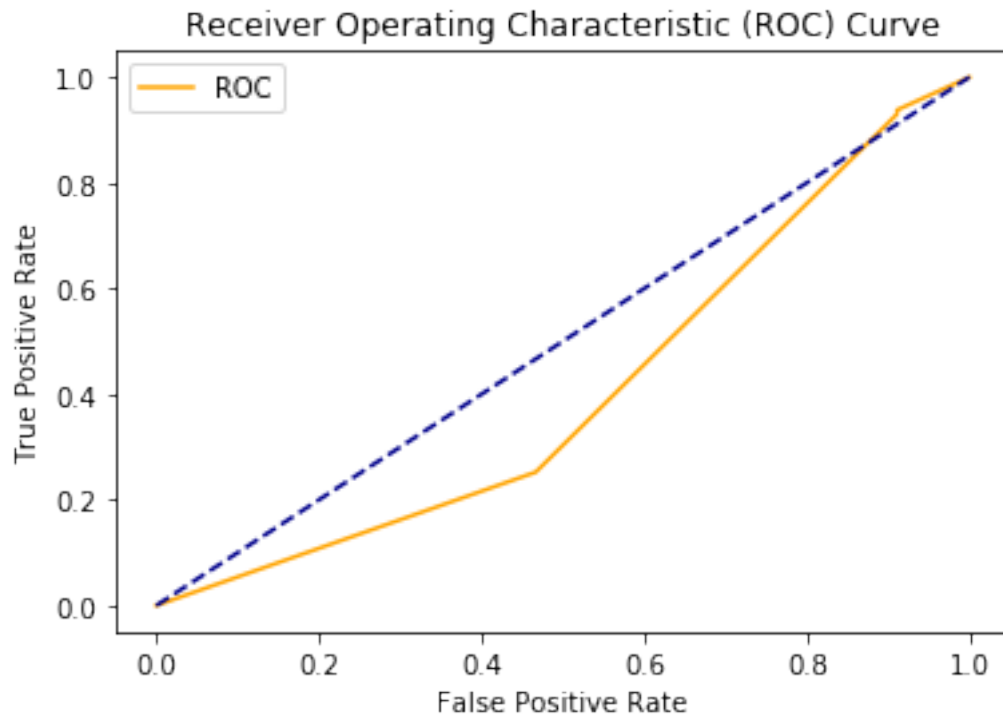
AUC2: 0.09
 optimal_threshold2: 2.00
 [2. 1. 0.85596766 0.]



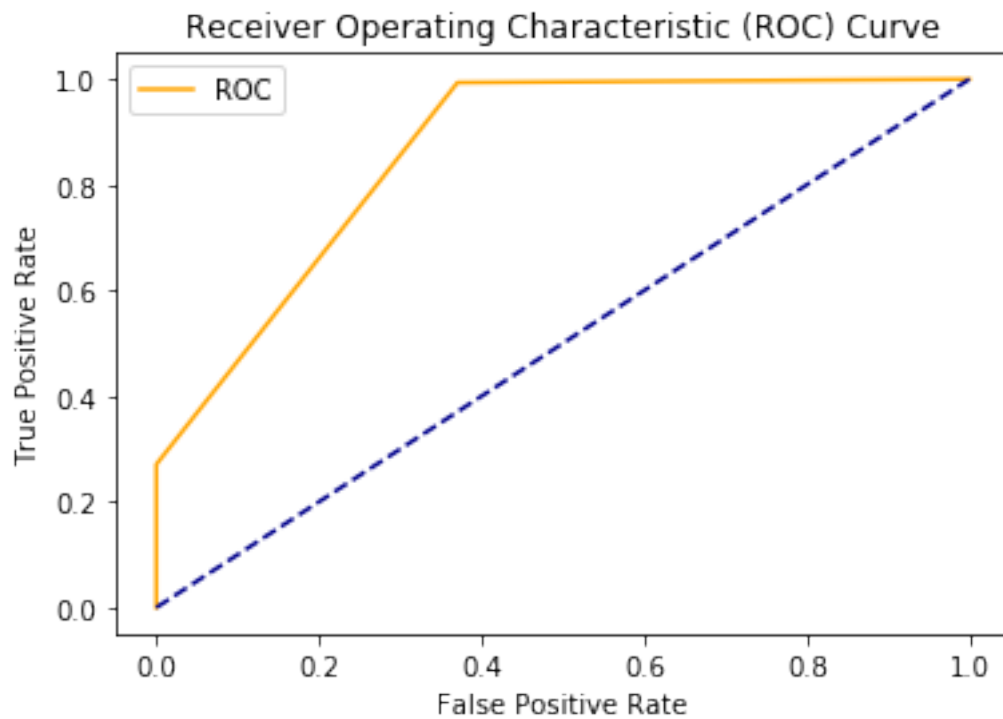
The score cutoff 10.00 for Reference A community 0 with pplacer_stats edpl
 compared with test B: 10.00
 data_set is True
 AUC1: 0.86
 optimal_threshold1: 0.90
 [2. 1. 0.9018491 0.]]



AUC2: 0.41
 optimal_threshold2: 0.20
 [2. 1. 0.9018491 0.2 0.]]



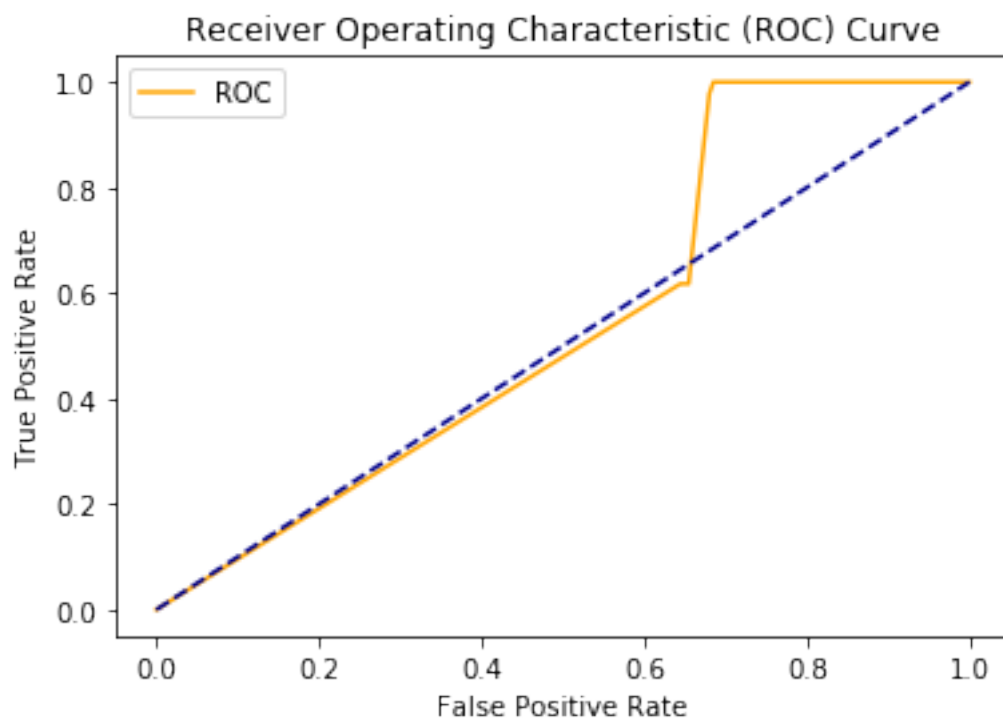
The score cutoff 10.00 for Reference A community 0 with pplacer_stats edpl
 compared with test C: 10.00
 data_set is True
 AUC1: 0.86
 optimal_threshold1: 0.88
 [2. 1. 0.88408218 0.]



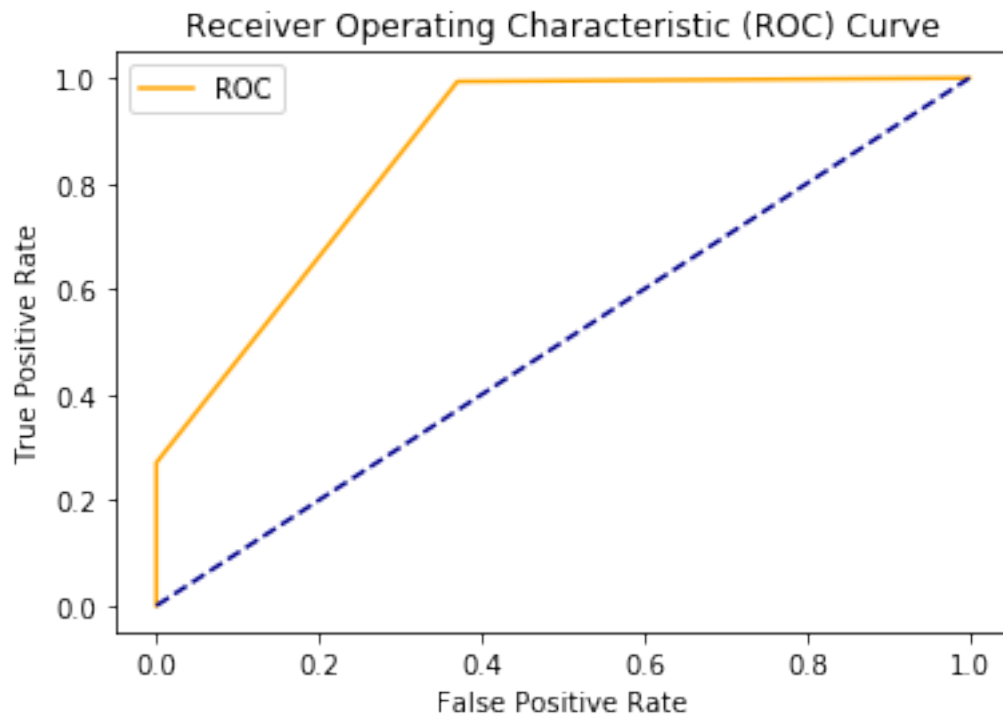
AUC2: 0.55

optimal_threshold2: 0.10

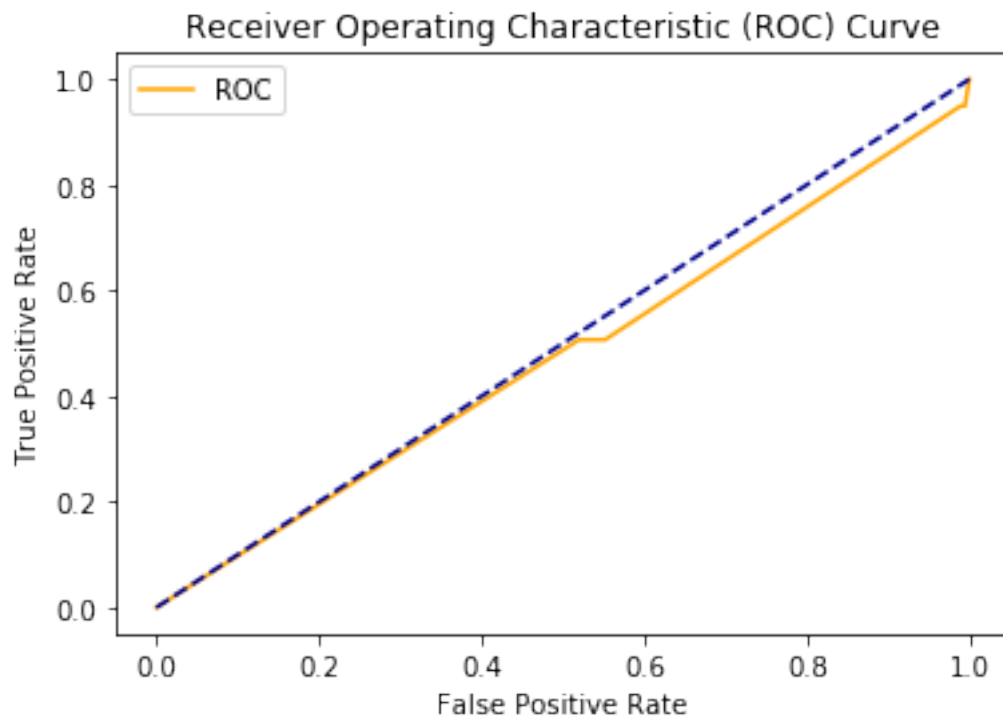
[2. 1. 0.9 0.88408218 0.1 0.]



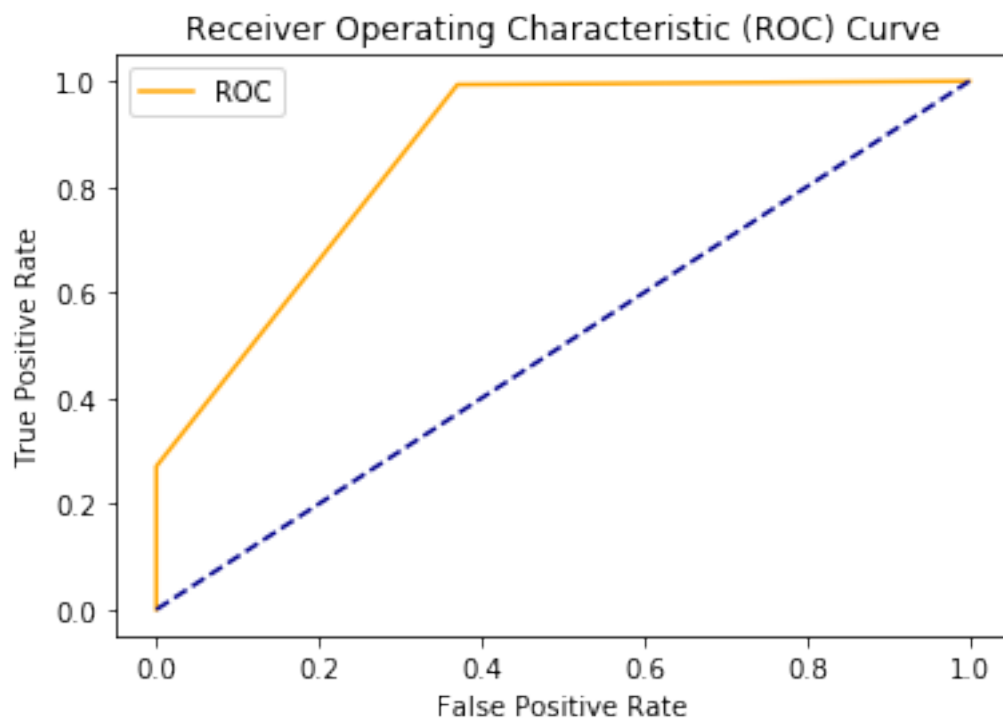
The score cutoff 10.00 for Reference A community 0 with pplacer_stats edpl
 compared with test D: 10.00
 data_set is True
 AUC1: 0.86
 optimal_threshold1: 0.87
 [2. 1. 0.87262189 0.]



AUC2: 0.48
 optimal_threshold2: 2.00
 [2. 1. 0.9 0.87262189 0.1 0.]



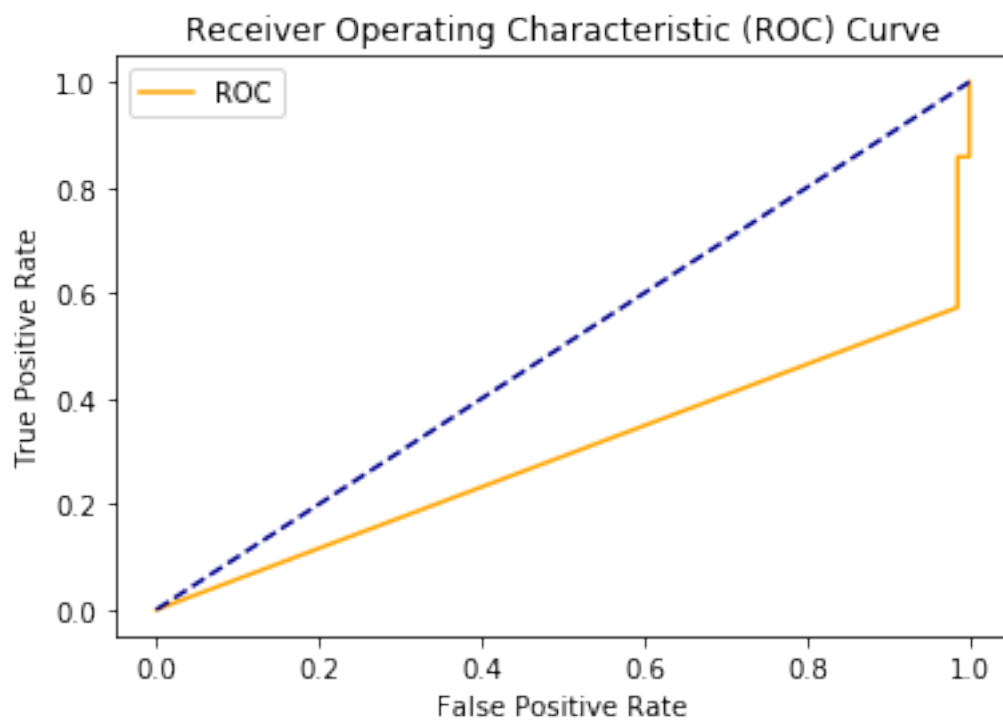
The score cutoff 10.00 for Reference A community 0 with pplacer_stats edpl
 compared with test E: 10.00
 data_set is True
 AUC1: 0.86
 optimal_threshold1: 0.89
 [2. 1. 0.98710937 0.89088363 0.]]



AUC2: 0.29

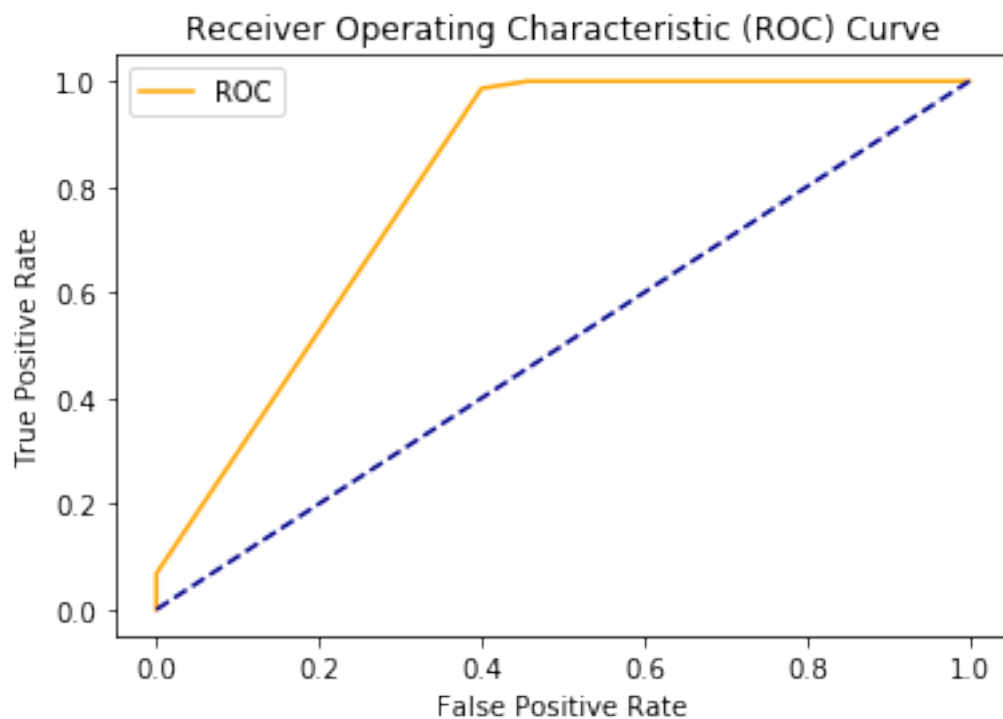
optimal_threshold2: 2.00

[2. 1. 0.89088363 0.7 0.1 0.]

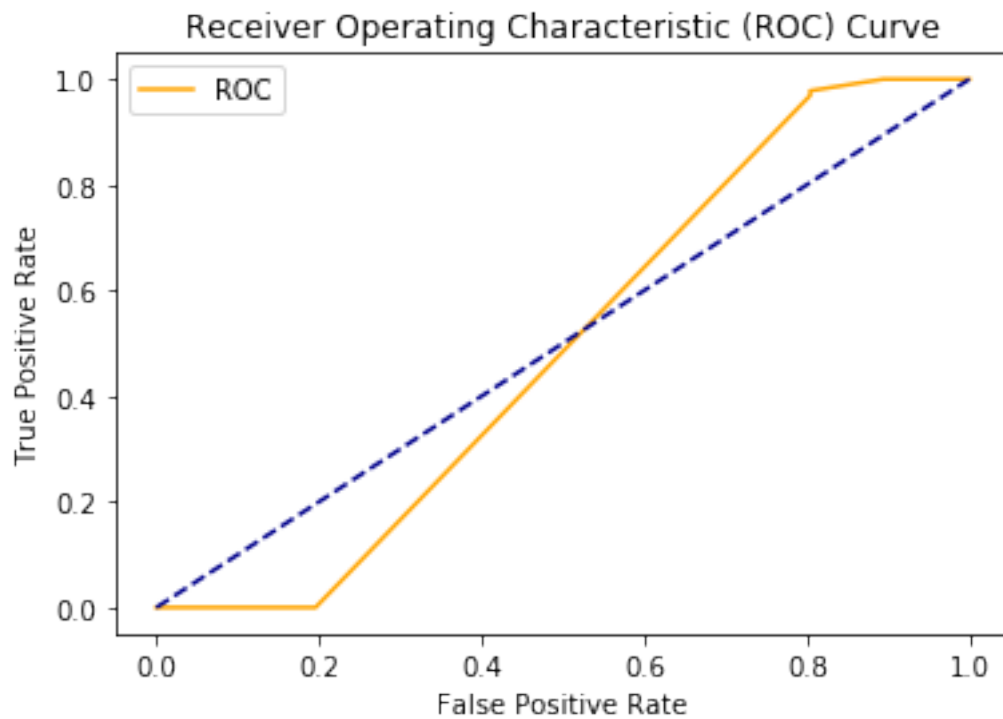


```
[177]: plot_roc_curve_microbiome_test2(pplacer_ref_list =
→ ['A'],pplacer_stats_list=['_adcl_log'],community_list=['0'],cutoff_list=['-4.
→ 00'], test_data_list=['B','C','D','E'],testOption=True, scoreOption=False)
```

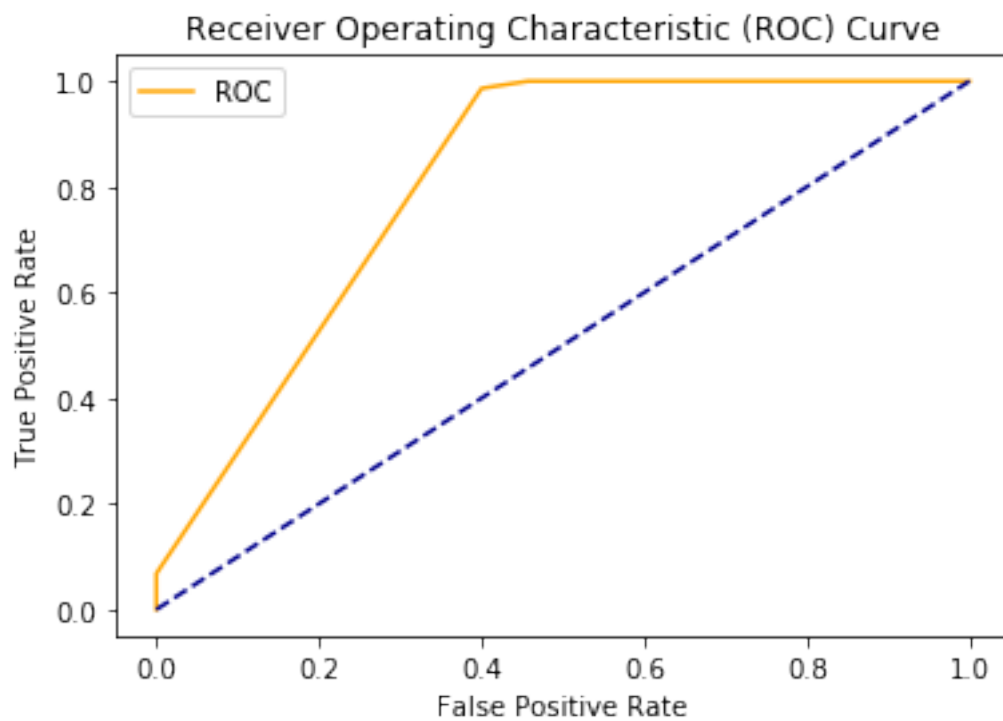
The pplacer_stats_cutoff -4.00 for Reference A community 0 pplacer_stats
adcl_log compared with test B: -4.00
data_set is True
AUC1: 0.81
optimal_threshold1: 0.92
[2. 1. 0.92196452 0.75654762 0.]



AUC2: 0.49
optimal_threshold2: 0.90
[2. 1. 0.92196452 0.9 0.75654762 0.]



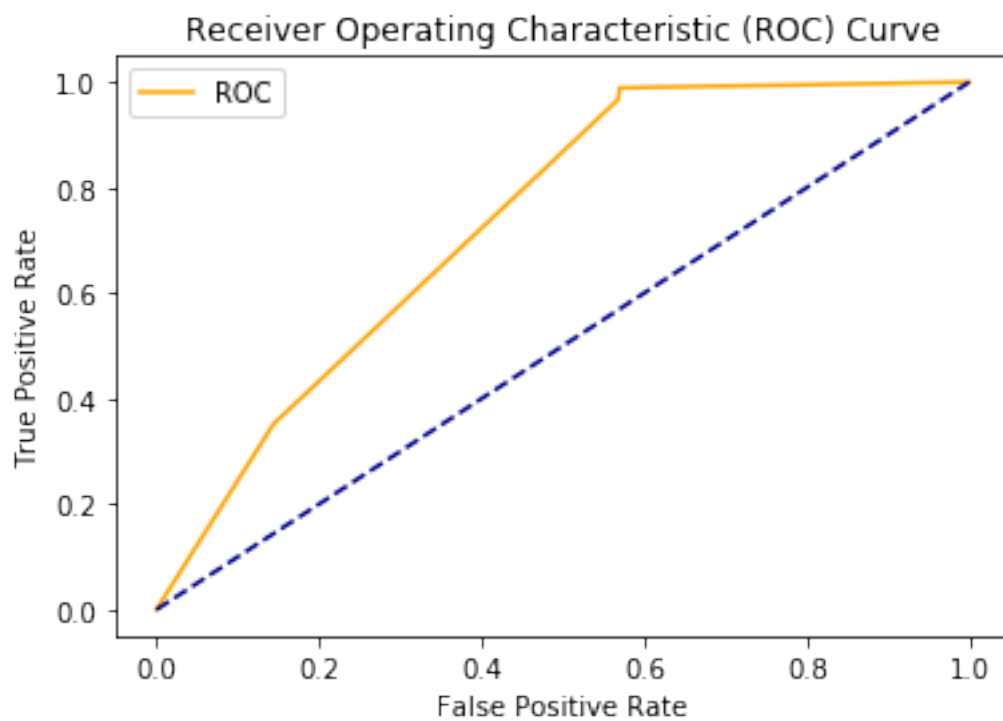
The pplacer_stats_cutoff -4.00 for Reference A community 0 pplacer_stats
 adcl_log compared with test C: -4.00
 data_set is True
 AUC1: 0.81
 optimal_threshold1: 0.92
 [2. 1. 0.92433893 0.6797619 0.]]



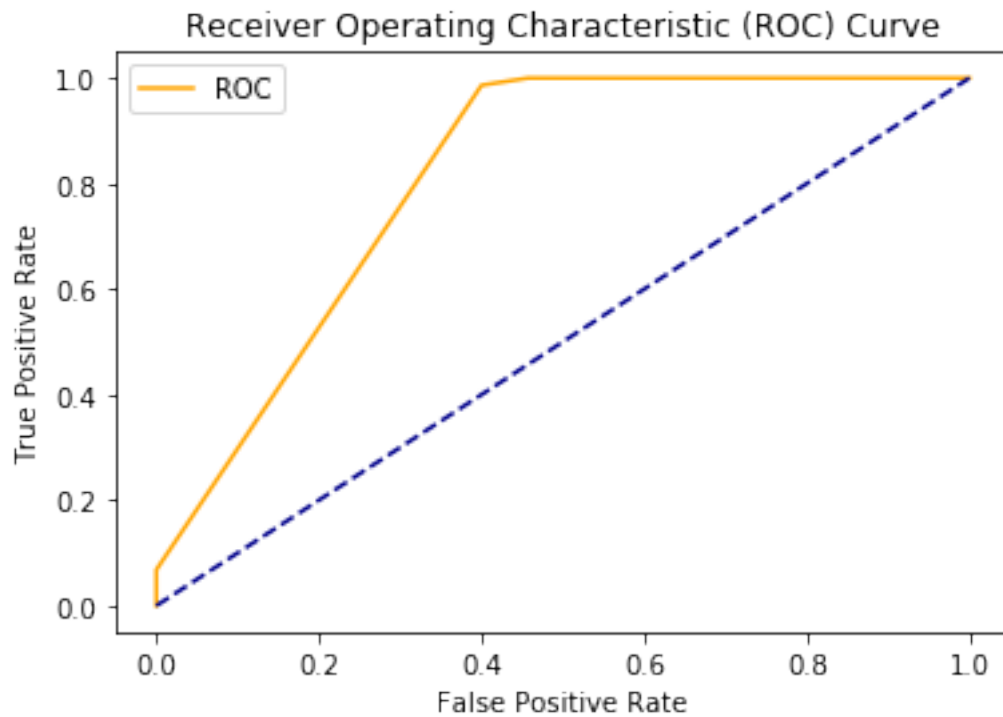
AUC2: 0.73

optimal_threshold2: 0.68

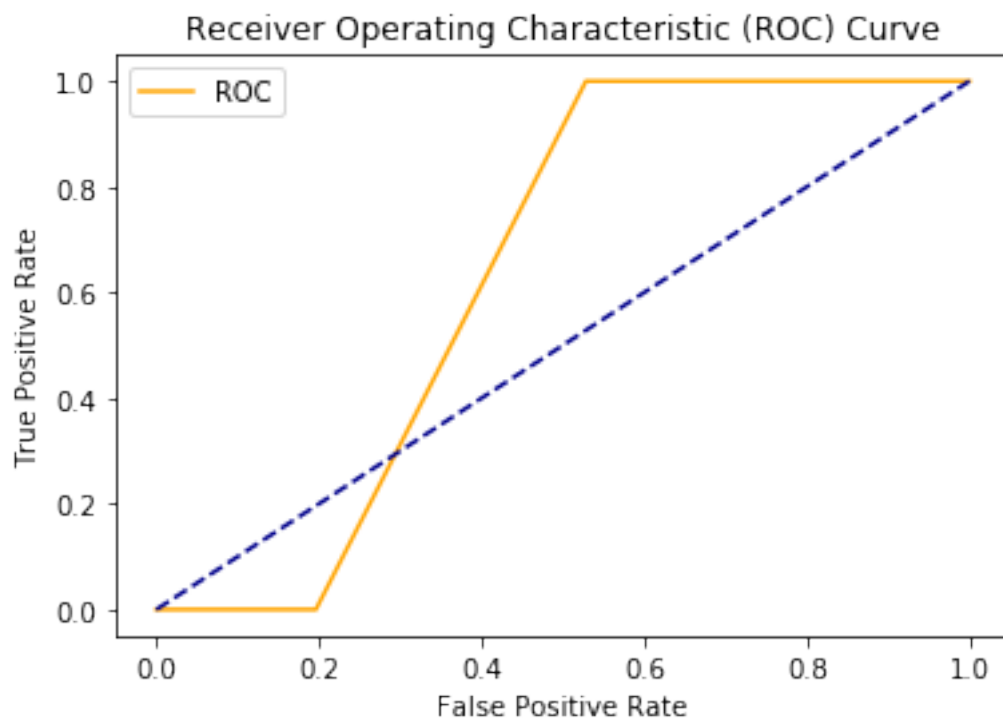
[2. 1. 0.92433893 0.6797619 0.]



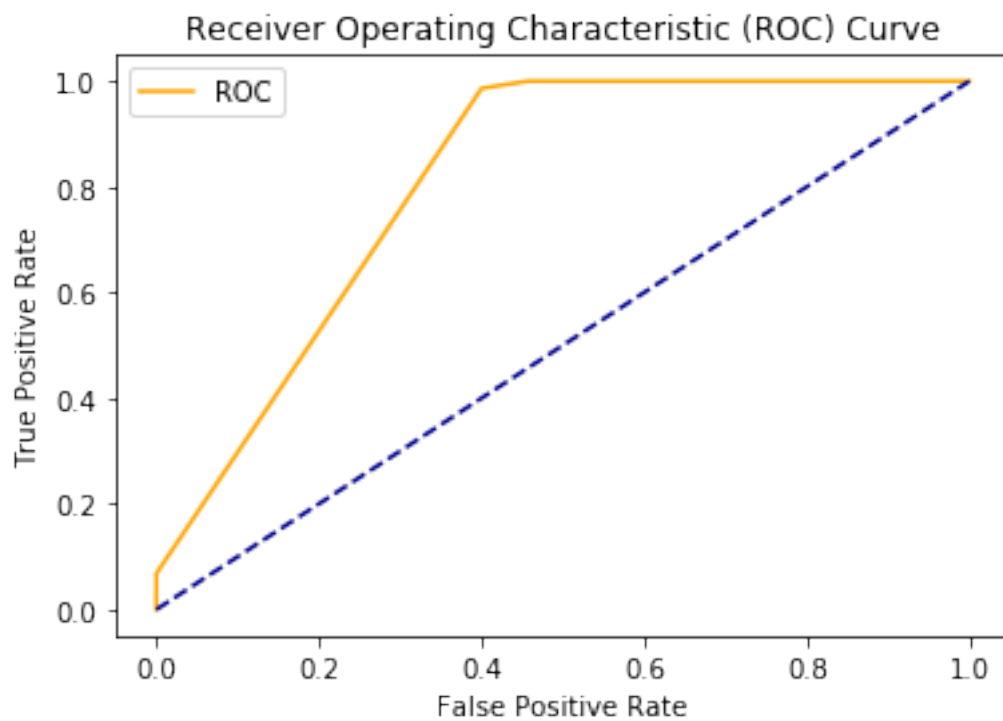
The pplacer_stats_cutoff -4.00 for Reference A community 0 pplacer_stats
adcl_log compared with test D: -4.00
data_set is True
AUC1: 0.81
optimal_threshold1: 0.92
[2. 1. 0.91710262 0.5984632 0.]]



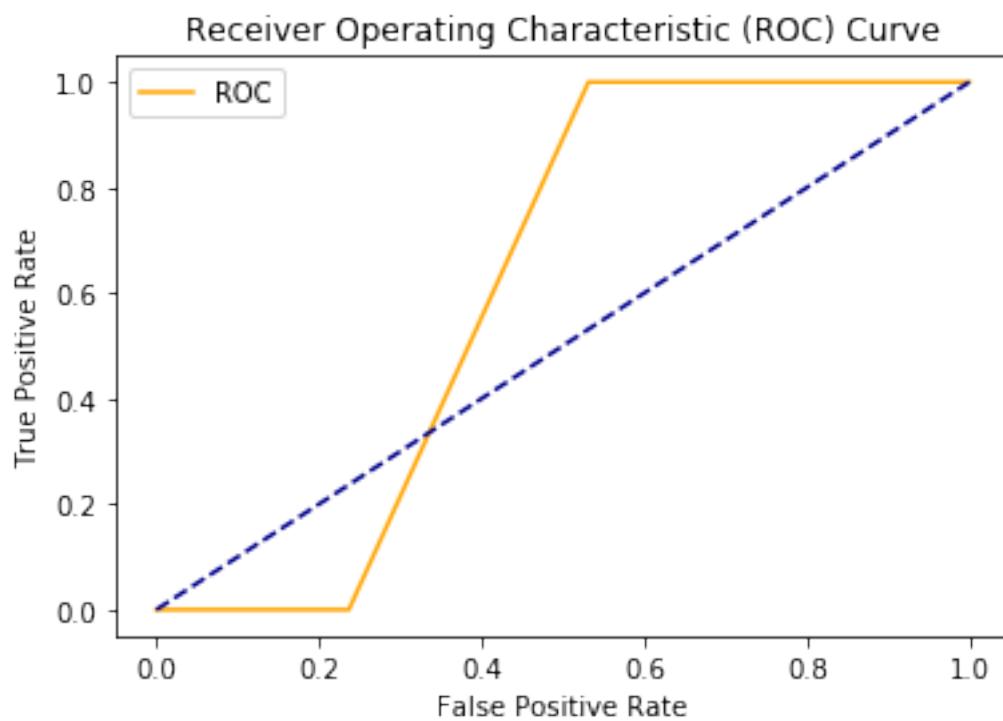
AUC2: 0.64
optimal_threshold2: 0.92
[2. 1. 0.91710262 0.5984632 0.]]



The pplacer_stats_cutoff -4.00 for Reference A community 0 pplacer_stats
 adcl_log compared with test E: -4.00
 data_set is True
 AUC1: 0.81
 optimal_threshold1: 0.93
 [2. 1. 0.93096139 0.69437908 0.]]

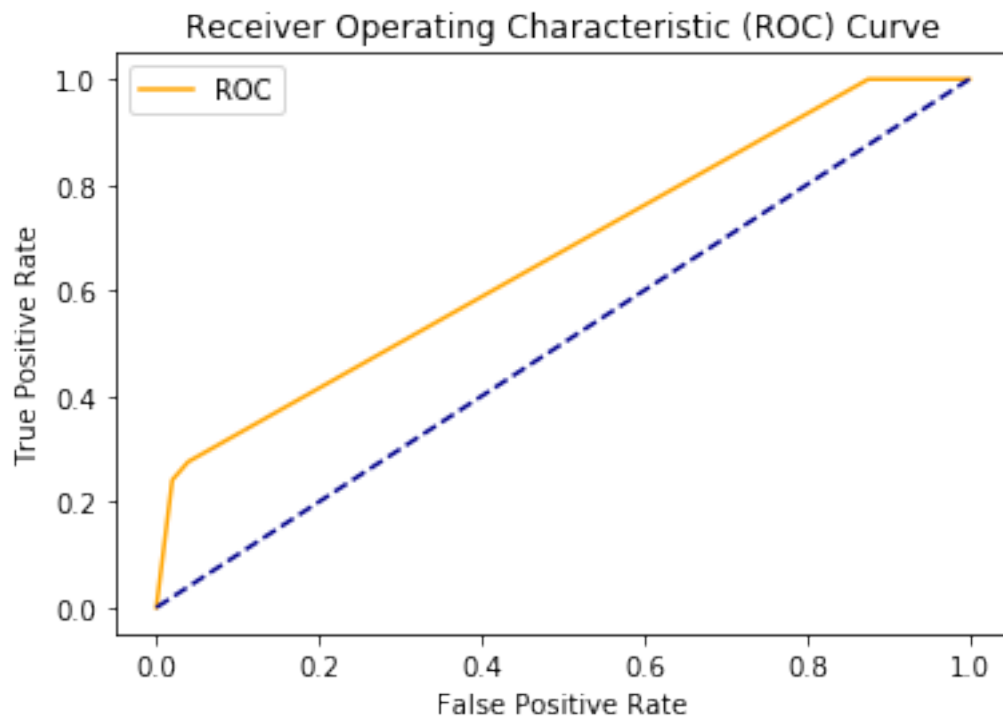


AUC2: 0.62
 optimal_threshold2: 0.93
 [2. 1. 0.93096139 0.69437908]

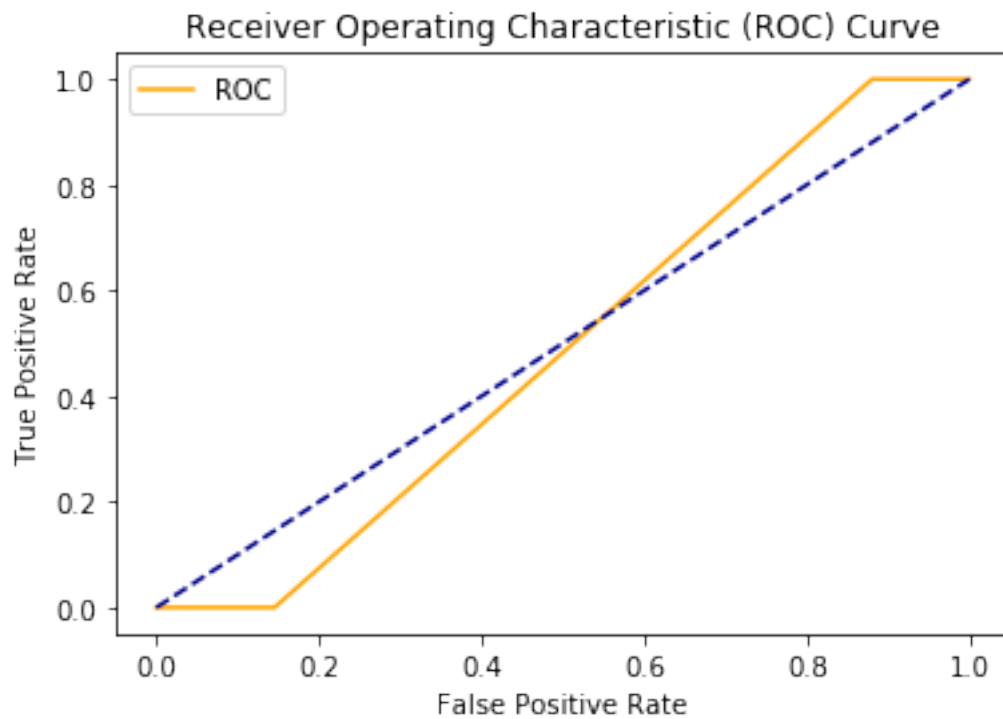


```
[136]: plot_roc_curve_microbiome_test2(pplacer_ref_list =  
→ ['A'],pplacer_stats_list=['_adcl_log'],community_list=['0'],cutoff_list=['25%'],  
→ test_data_list=['B','C'],testOption=True, scoreOption=False)
```

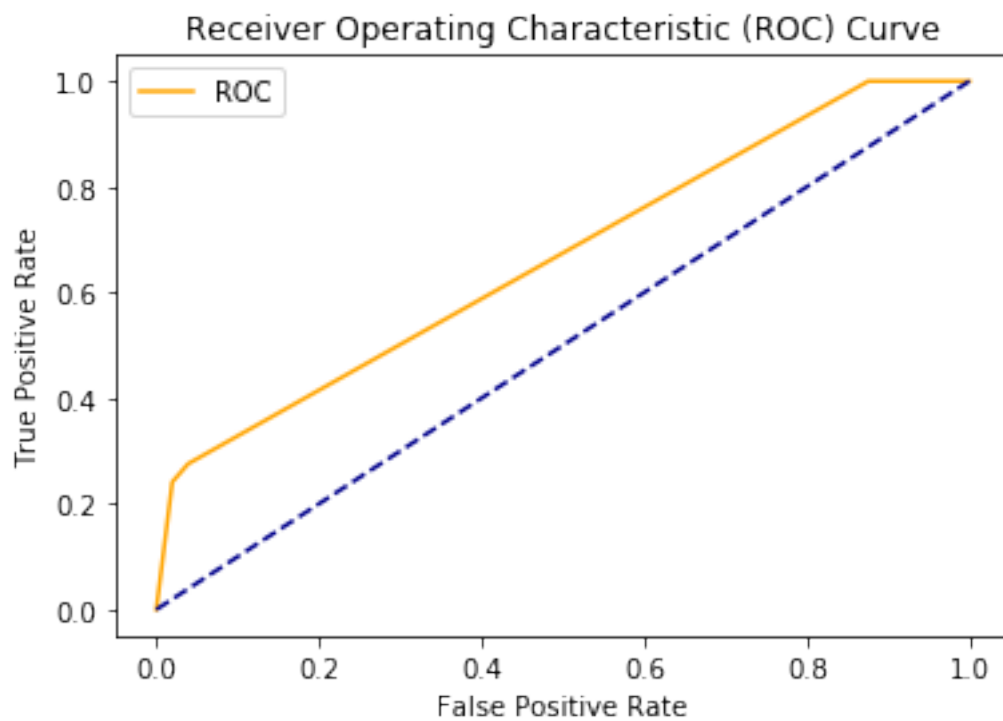
The pplacer_stats_cutoff 25% for Reference A community 0 pplacer_stats adcl_log compared with test B: -5.22
data_set is True
AUC1: 0.67
optimal_threshold1: 0.46



AUC2: 0.49
optimal_threshold2: 0.15

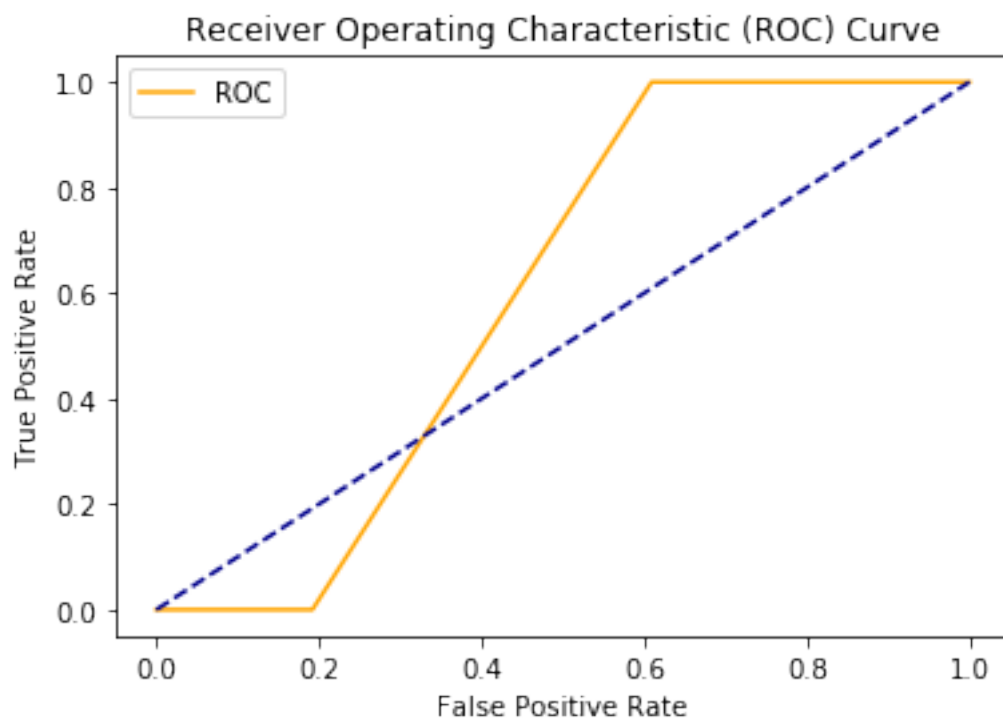


The pplacer_stats_cutoff 25% for Reference A community 0 pplacer_stats adcl_log
compared with test C: -5.22
data_set is True
AUC1: 0.67
optimal_threshold1: 0.42



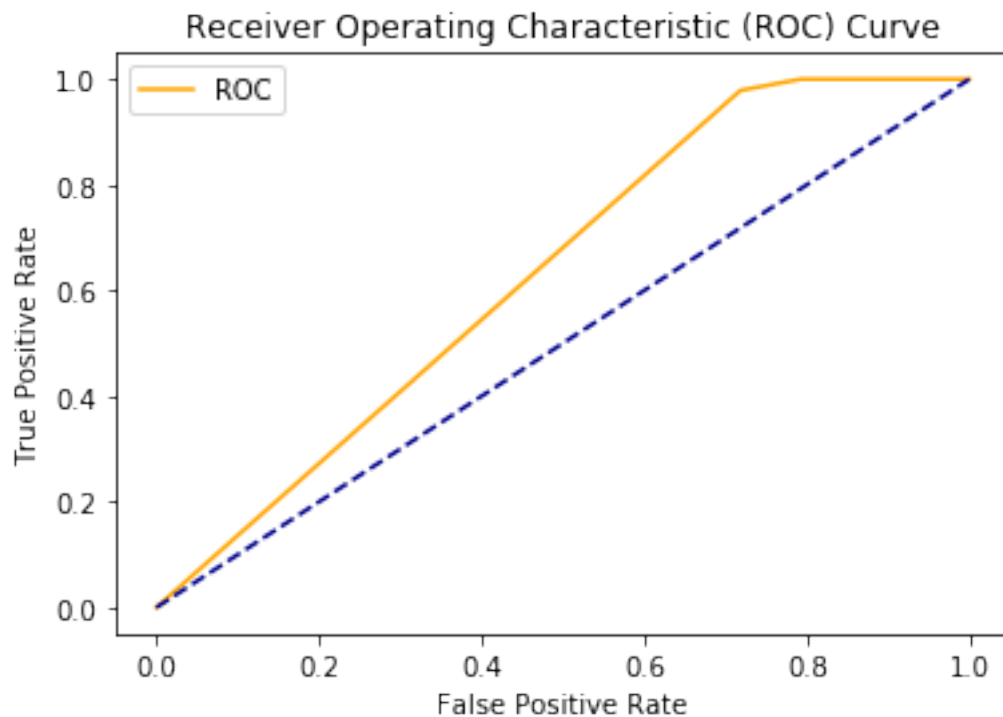
AUC2: 0.60

optimal_threshold2: 0.14

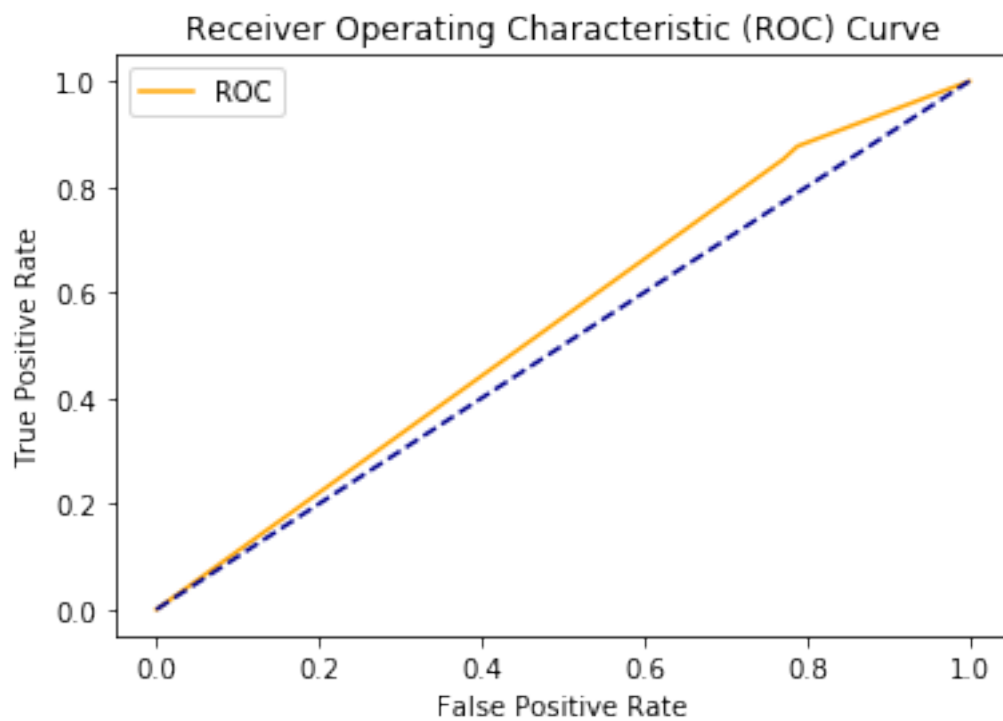


```
[128]: plot_roc_curve_microbiome_test2(pplacer_ref_list =  
→ ['B'], ppacer_stats_list=['_adcl_log'], community_list=['0'], cutoff_list=['25%'],  
→ test_data_list=['A', 'C'], testOption=True, scoreOption=False)
```

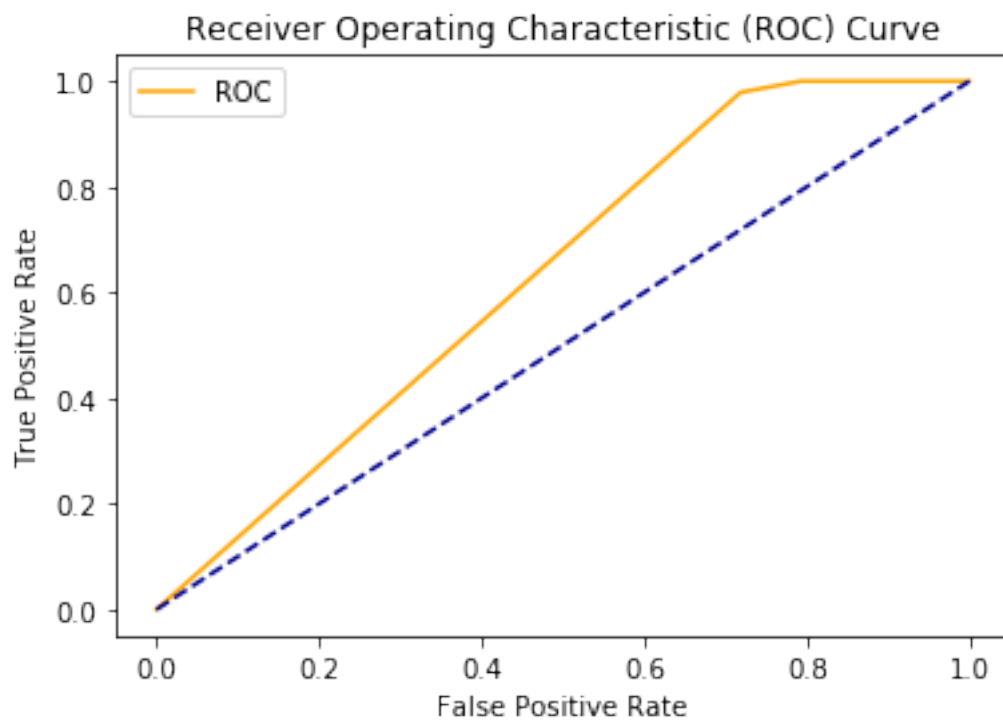
The ppacer_stats_cutoff 25% for Reference B community 0 ppacer_stats adcl_log compared with test A: -5.15
data_set is True
AUC1: 0.63
optimal_threshold1: 0.23



AUC2: 0.54
optimal_threshold2: 0.11

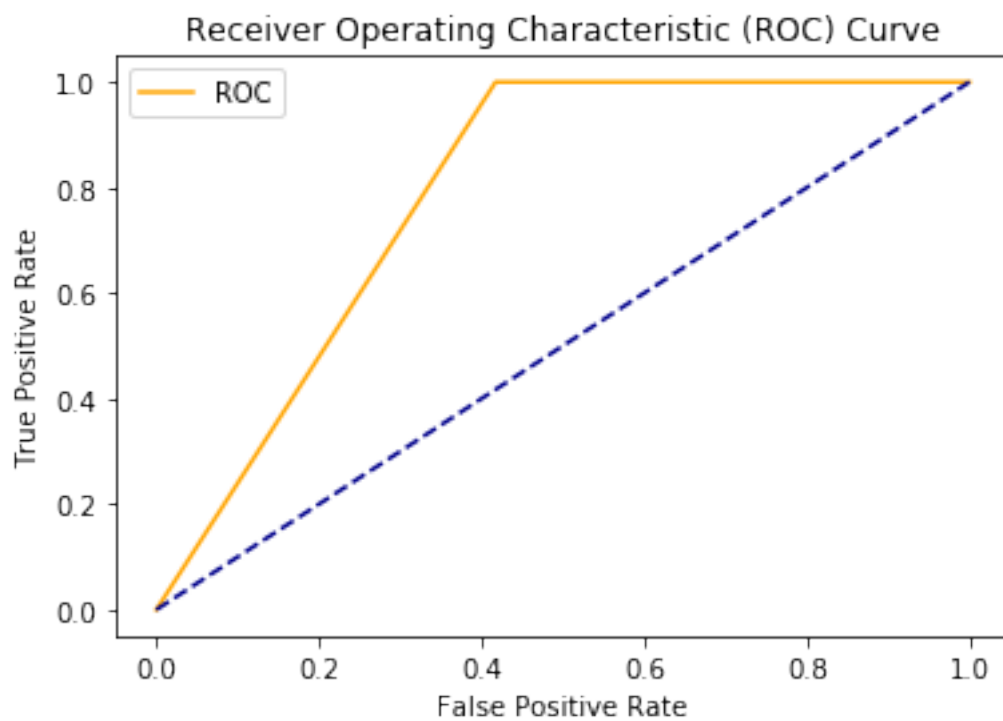


The pplacer_stats_cutoff 25% for Reference B community 0 pplacer_stats adcl_log
compared with test C: -5.15
data_set is True
AUC1: 0.63
optimal_threshold1: 0.23



AUC2: 0.79

optimal_threshold2: 0.23



```

[96]: # print("the head for df is {}".format(df.head)+ " the columns of the df is {}".
      ↪format(df.columns))
      #

[33]: # df['A0'].describe(), df['B0'].describe(), df['C0'].describe(), df['D0'].
      ↪describe(), df['E0'].describe()

[34]: # for community in ['A', 'B', 'C', 'D', 'E']:
      #     for i in range(10):
      #         print(df[community+str(i)].describe())

[35]: df_0 = df

[36]: # plot_pplacer('90')

[37]: # plotScatter('B', '0')

[38]: # plotScatterRef('_adcl_log', '0')

[39]: # plot_pplacer('_adcl_log')

[40]: df['A_adcl_log'].describe()

[40]: count      5974.000000
      mean       -4.083366
      std        1.837510
      min       -5.995679
      25%       -5.221126
      50%       -5.096367
      75%       -1.706947
      max       -0.344675
      Name: A_adcl_log, dtype: float64

[41]: # plot_pplacer('0')

[42]: df['A0'].describe()

[42]: count      605.000000
      mean       6.390083
      std       10.778008
      min       0.000000
      25%       2.000000
      50%       2.000000
      75%       2.000000
      max      38.000000
      Name: A0, dtype: float64

[43]: df1 = df[(df['A0']>10)]

[44]: df1['A0'].describe()

```



```
[44]: count    99.000000
      mean     28.626263
      std      10.791707
      min      12.000000
      25%      12.000000
      50%      32.000000
      75%      38.000000
      max      38.000000
      Name: A0, dtype: float64
```

```
[45]: 99/605
```

```
[45]: 0.16363636363636364
```

```
[46]: df1['B0'].describe()
```

```
[46]: count    99.000000
      mean     25.070707
      std      17.438670
      min       0.000000
      25%       0.000000
      50%      36.000000
      75%      38.000000
      max      38.000000
      Name: B0, dtype: float64
```

```
[47]: df2=df[['seqID', 'A0', 'B0', 'C0', 'D0', 'E0']].dropna()
```

```
[48]: df3 = df2[(df2['A0']>10) & (df2['B0']>10) & (df2['C0']>10) & (df2['D0']>10) &
      →(df2['E0']>10)]
```

```
[49]: # df2.describe(), df3.describe()
```

```
[50]: df3
```

```
[50]:
```

	seqID	A0	B0	C0	D0	E0
5313	CC11CM5SCR137ef78188b94db7b59504dc64363aa3	34.0	34.0	32.0	32.0	44.0
5314	CC11CM0SCR35529da454f0497fa16e04841e8e1639	34.0	34.0	32.0	32.0	44.0

```
[51]: 2/605
```

```
[51]: 0.003305785123966942
```

```
[52]: dfc90 = df[(df['A90']>10) & (df['B90']>10) & (df['C90']>10) & (df['D90']>10) &
      →& (df['E90']>10)]
```

```
[53]: dfc90['B0'].describe()
```

```
[53]: count    0.0
      mean    NaN
      std     NaN
      min     NaN
      25%     NaN
      50%     NaN
```

```

75%      NaN
max      NaN
Name: B0, dtype: float64

```

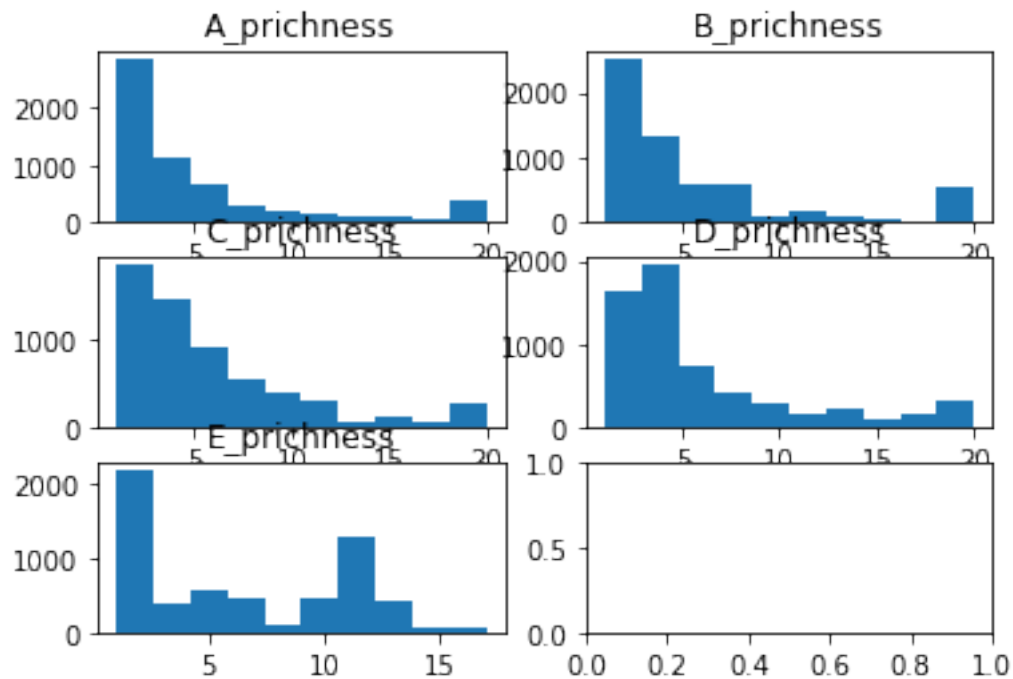
```
[54]: df[(df.community=='CC11CM0')]['C_adcl_log'].dropna().describe()
```

```

[54]: count      55.000000
      mean      -1.885119
      std       1.762885
      min      -5.300162
      25%      -1.773077
      50%      -1.040954
      75%      -0.682030
      max      -0.373058
      Name: C_adcl_log, dtype: float64

```

```
[55]: plot_pplacer('_prichness')
```



```
[56]: df['A_prichness'].describe()
```

```

[56]: count      5974.000000
      mean       4.727653
      std       5.465596
      min       1.000000
      25%       1.000000
      50%       3.000000
      75%       5.000000

```

```
max          20.000000
Name: A_prichness, dtype: float64
```

```
[57]: df[df.A0>10].A0.count()
```

```
[57]: 99
```

```
[58]: df[df.A0>10].A0.count()/df.A0.count()
```

```
[58]: 0.16363636363636364
```

```
[59]: # df.head()
```

```
[60]: df.A_adcl.count()
```

```
[60]: 5974
```

```
[61]: d={"a":1, "b":2}
```

```
[62]: d
```

```
[62]: {'a': 1, 'b': 2}
```

```
[63]: dd = pd.Series(d, name='score')
```

```
[64]: dd.index.name="community"
```

```
[65]: dd.reset_index()
```

```
[65]:   community  score
0          a       1
1          b       2
```

```
[66]: "CC11CM"+str(0)
```

```
[66]: 'CC11CM0'
```

```
[67]: c0=df['A0'][df['community']=='CC11CM0']
```

```
[68]: per=c0[c0>10].count()/c0.count()
```

```
[69]: def generateScore(stats, referenceID,scorecutoff,statscutoff):
      d1={}
      d2={}
      for i in range(100):
          values = df[referenceID+str(i)][df.community=='CC11CM'+str(i)]
          statsvalues = df[stats][df['community']=='CC11CM'+str(i)]
          d1['CC11CM'+str(i)] = values[values>scorecutoff].count()/values.count()
          d2['CC11CM'+str(i)] =statsvalues[statsvalues>statscutoff].count()/
      ↪statsvalues.count()
      d1=pd.Series(d1, name=referenceID)
      d1.index.name='community'
      d1=d1.reset_index()
      d2=pd.Series(d2, name=stats)
      d2.index.name='community'
      d2=d2.reset_index()
```

```

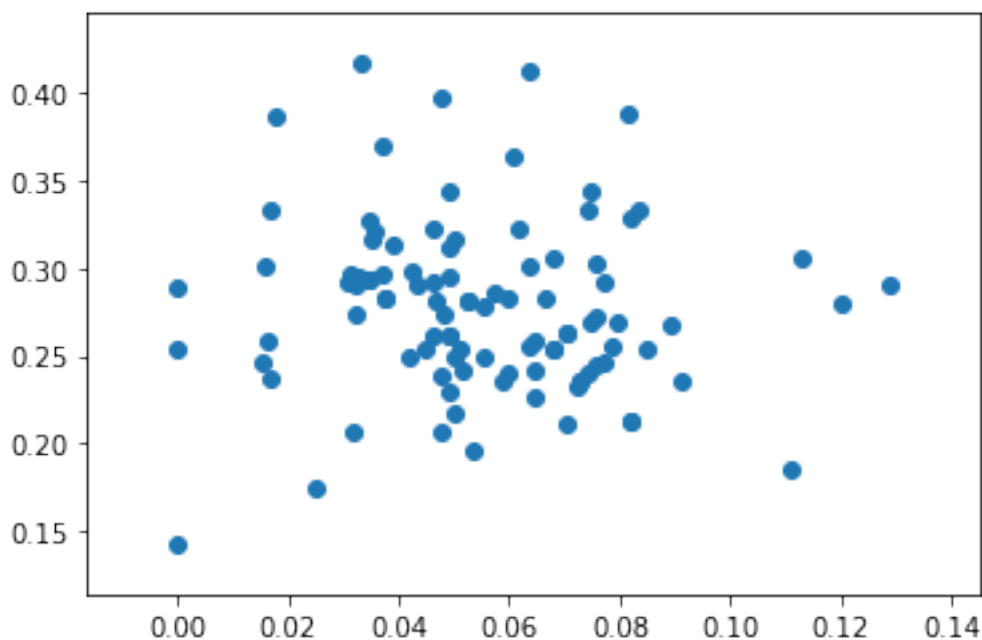
dt = pd.concat([d1,d2], axis=1)
# dt=dt.set_index('community')
return (dt)

```

```
[70]: dt=generateScore('A_adcl', 'A', 10, 0.001)
```

```
[71]: plt.scatter(dt.A, dt.A_adcl)
```

```
[71]: <matplotlib.collections.PathCollection at 0x1efb8cd0>
```



```

[72]: t=[]
for referenceID in ['A','B','C','D','E']:
    t.append(generateScore('A_adcl', referenceID, 10, 0.001))

```

```
[73]: t[0].head()
```

```

[73]:  community      A community  A_adcl
0  CC11CM0  0.090909  CC11CM0  0.236364
1  CC11CM1  0.082192  CC11CM1  0.328767
2  CC11CM2  0.025000  CC11CM2  0.175000
3  CC11CM3  0.050847  CC11CM3  0.254237
4  CC11CM4  0.000000  CC11CM4  0.288462

```

```
[74]: tt = pd.concat([t[0],t[1],t[2],t[3],t[4]], axis=1)
```

```

[75]: def generateScoreu(stats, referenceID,scorecutoff,statscutoff):
    d1={}
    d2={}
    for i in range(100):
        values = df[referenceID+str(i)][df.community=='CC11CM'+str(i)]

```

```

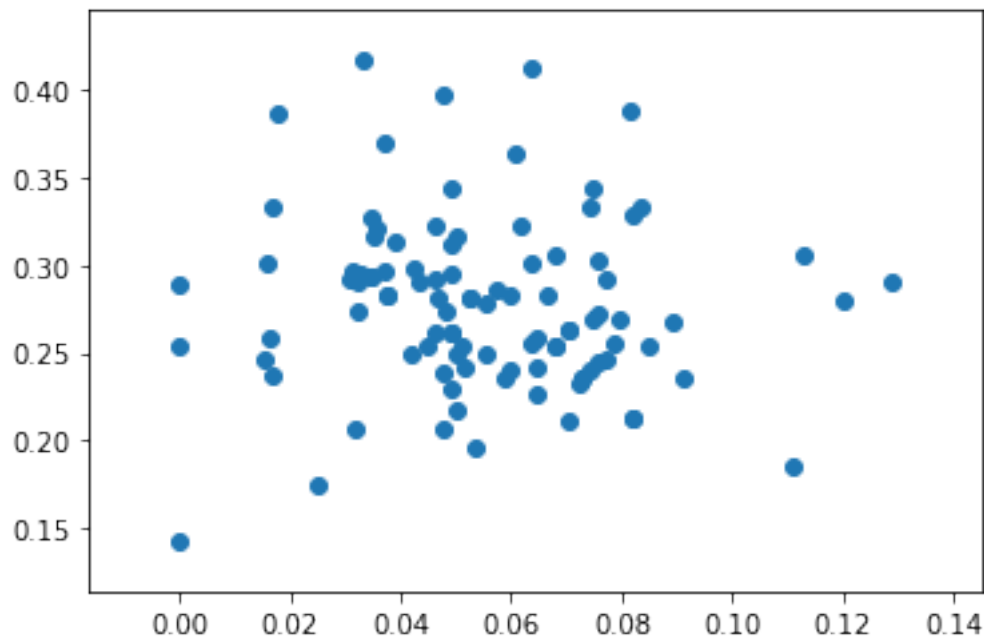
statsvalues = df[referenceID+stats][df['community']=='CC11CM'+str(i)]
d1['CC11CM'+str(i)] = values[values>scorecutoff].count()/values.count()
d2['CC11CM'+str(i)] =statsvalues[statsvalues>statscutoff].count()/
→statsvalues.count()
d1=pd.Series(d1, name=referenceID)
d1.index.name='community'
d1=d1.reset_index()
d2=pd.Series(d2, name=referenceID+stats)
d2.index.name='community'
d2=d2.reset_index()
dt = pd.concat([d1,d2], axis=1)
dt=dt.loc[:, ~dt.columns.duplicated()]
dt=dt.set_index('community')
return (dt)

```

```
[76]: dtu=generateScoreu('_adcl', 'A', 10, 0.001)
```

```
[77]: plt.scatter(dtu.A, dtu.A_adcl)
```

```
[77]: <matplotlib.collections.PathCollection at 0x1f294fd0>
```



```
[ ]:
```

```

[78]: t=[]
statsdir= {'_adcl':0.0001, '_edpl':0, '_prichness':10, '_mindistl':0.05}
for stats in statsdir.keys():

    for referenceID in ['A', 'B', 'C', 'D', 'E']:

```

```
t.append(generateScoreu(stats, referenceID, 10, statsdir[stats]))
```

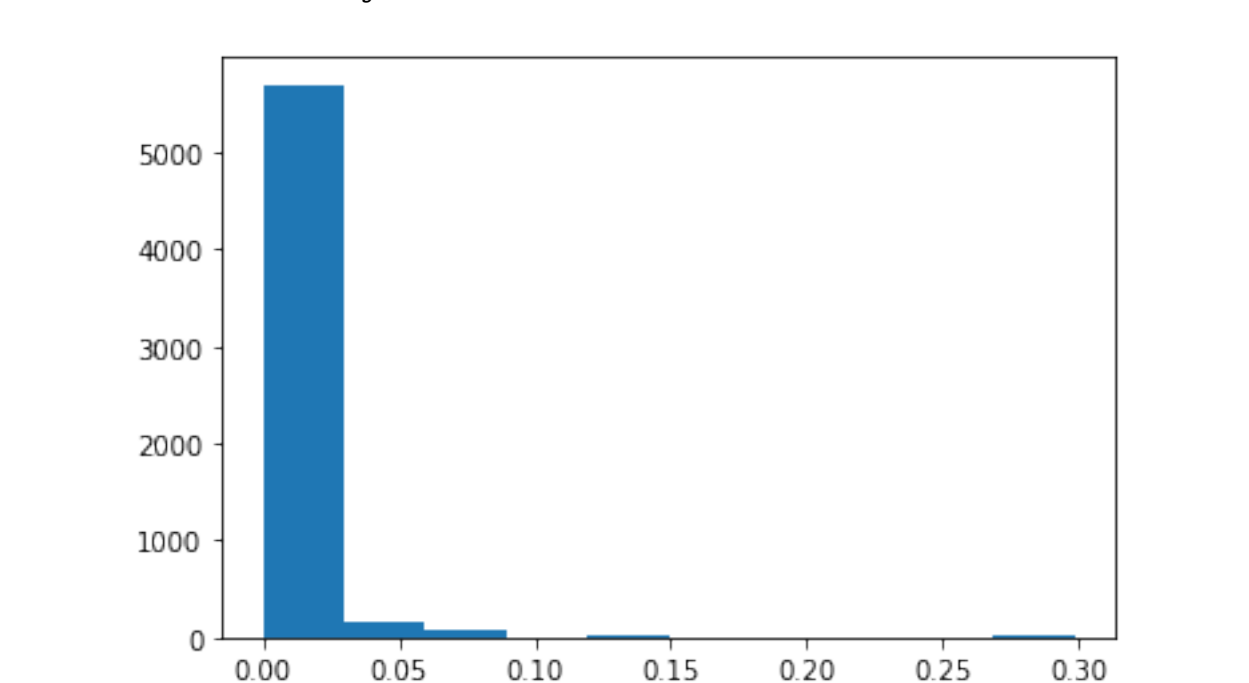
```
ttd = pd.  
    ↳ concat([t[0],t[1],t[2],t[3],t[4],t[5],t[6],t[7],t[8],t[9],t[10],t[11],t[12],t[13],t[14],t[15],t[16],t[17],t[18],t[19],t[20],t[21],t[22],t[23],t[24],t[25],t[26],t[27],t[28],t[29],t[30],t[31],t[32],t[33],t[34],t[35],t[36],t[37],t[38],t[39],t[40],t[41],t[42],t[43],t[44],t[45],t[46],t[47],t[48],t[49],t[50],t[51],t[52],t[53],t[54],t[55],t[56],t[57],t[58],t[59],t[60],t[61],t[62],t[63],t[64],t[65],t[66],t[67],t[68],t[69],t[70],t[71],t[72],t[73],t[74],t[75],t[76],t[77],t[78],t[79],t[80],t[81],t[82],t[83],t[84],t[85],t[86],t[87],t[88],t[89],t[90],t[91],t[92],t[93],t[94],t[95],t[96],t[97],t[98],t[99],t[100],t[101],t[102],t[103],t[104],t[105],t[106],t[107],t[108],t[109],t[110],t[111],t[112],t[113],t[114],t[115],t[116],t[117],t[118],t[119],t[120],t[121],t[122],t[123],t[124],t[125],t[126],t[127],t[128],t[129],t[130],t[131],t[132],t[133],t[134],t[135],t[136],t[137],t[138],t[139],t[140],t[141],t[142],t[143],t[144],t[145],t[146],t[147],t[148],t[149],t[150],t[151],t[152],t[153],t[154],t[155],t[156],t[157],t[158],t[159],t[160],t[161],t[162],t[163],t[164],t[165],t[166],t[167],t[168],t[169],t[170],t[171],t[172],t[173],t[174],t[175],t[176],t[177],t[178],t[179],t[180],t[181],t[182],t[183],t[184],t[185],t[186],t[187],t[188],t[189],t[190],t[191],t[192],t[193],t[194],t[195],t[196],t[197],t[198],t[199],t[200],t[201],t[202],t[203],t[204],t[205],t[206],t[207],t[208],t[209],t[210],t[211],t[212],t[213],t[214],t[215],t[216],t[217],t[218],t[219],t[220],t[221],t[222],t[223],t[224],t[225],t[226],t[227],t[228],t[229],t[230],t[231],t[232],t[233],t[234],t[235],t[236],t[237],t[238],t[239],t[240],t[241],t[242],t[243],t[244],t[245],t[246],t[247],t[248],t[249],t[250],t[251],t[252],t[253],t[254],t[255],t[256],t[257],t[258],t[259],t[260],t[261],t[262],t[263],t[264],t[265],t[266],t[267],t[268],t[269],t[270],t[271],t[272],t[273],t[274],t[275],t[276],t[277],t[278],t[279],t[280],t[281],t[282],t[283],t[284],t[285],t[286],t[287],t[288],t[289],t[290],t[291],t[292],t[293],t[294],t[295],t[296],t[297],t[298],t[299],t[300],t[301],t[302],t[303],t[304],t[305],t[306],t[307],t[308],t[309],t[310],t[311],t[312],t[313],t[314],t[315],t[316],t[317],t[318],t[319],t[320],t[321],t[322],t[323],t[324],t[325],t[326],t[327],t[328],t[329],t[330],t[331],t[332],t[333],t[334],t[335],t[336],t[337],t[338],t[339],t[340],t[341],t[342],t[343],t[344],t[345],t[346],t[347],t[348],t[349],t[350],t[351],t[352],t[353],t[354],t[355],t[356],t[357],t[358],t[359],t[360],t[361],t[362],t[363],t[364],t[365],t[366],t[367],t[368],t[369],t[370],t[371],t[372],t[373],t[374],t[375],t[376],t[377],t[378],t[379],t[380],t[381],t[382],t[383],t[384],t[385],t[386],t[387],t[388],t[389],t[390],t[391],t[392],t[393],t[394],t[395],t[396],t[397],t[398],t[399],t[400],t[401],t[402],t[403],t[404],t[405],t[406],t[407],t[408],t[409],t[410],t[411],t[412],t[413],t[414],t[415],t[416],t[417],t[418],t[419],t[420],t[421],t[422],t[423],t[424],t[425],t[426],t[427],t[428],t[429],t[430],t[431],t[432],t[433],t[434],t[435],t[436],t[437],t[438],t[439],t[440],t[441],t[442],t[443],t[444],t[445],t[446],t[447],t[448],t[449],t[450],t[451],t[452],t[453],t[454],t[455],t[456],t[457],t[458],t[459],t[460],t[461],t[462],t[463],t[464],t[465],t[466],t[467],t[468],t[469],t[470],t[471],t[472],t[473],t[474],t[475],t[476],t[477],t[478],t[479],t[480],t[481],t[482],t[483],t[484],t[485],t[486],t[487],t[488],t[489],t[490],t[491],t[492],t[493],t[494],t[495],t[496],t[497],t[498],t[499],t[500],t[501],t[502],t[503],t[504],t[505],t[506],t[507],t[508],t[509],t[510],t[511],t[512],t[513],t[514],t[515],t[516],t[517],t[518],t[519],t[520],t[521],t[522],t[523],t[524],t[525],t[526],t[527],t[528],t[529],t[530],t[531],t[532],t[533],t[534],t[535],t[536],t[537],t[538],t[539],t[540],t[541],t[542],t[543],t[544],t[545],t[546],t[547],t[548],t[549],t[550],t[551],t[552],t[553],t[554],t[555],t[556],t[557],t[558],t[559],t[560],t[561],t[562],t[563],t[564],t[565],t[566],t[567],t[568],t[569],t[570],t[571],t[572],t[573],t[574],t[575],t[576],t[577],t[578],t[579],t[580],t[581],t[582],t[583],t[584],t[585],t[586],t[587],t[588],t[589],t[590],t[591],t[592],t[593],t[594],t[595],t[596],t[597],t[598],t[599],t[600],t[601],t[602],t[603],t[604],t[605],t[606],t[607],t[608],t[609],t[610],t[611],t[612],t[613],t[614],t[615],t[616],t[617],t[618],t[619],t[620],t[621],t[622],t[623],t[624],t[625],t[626],t[627],t[628],t[629],t[630],t[631],t[632],t[633],t[634],t[635],t[636],t[637],t[638],t[639],t[640],t[641],t[642],t[643],t[644],t[645],t[646],t[647],t[648],t[649],t[650],t[651],t[652],t[653],t[654],t[655],t[656],t[657],t[658],t[659],t[660],t[661],t[662],t[663],t[664],t[665],t[666],t[667],t[668],t[669],t[670],t[671],t[672],t[673],t[674],t[675],t[676],t[677],t[678],t[679],t[680],t[681],t[682],t[683],t[684],t[685],t[686],t[687],t[688],t[689],t[690],t[691],t[692],t[693],t[694],t[695],t[696],t[697],t[698],
```

```
# ttt.describe()
```

```
ttt.to_csv("community-based.csv")
```

```
plt.hist(df.A_mindist1)
```

```
(array([5.696e+03, 1.570e+02, 6.800e+01, 0.000e+00, 2.200e+01, 0.000e+00,
        0.000e+00, 5.000e+00, 0.000e+00, 2.600e+01]),
 array([3.48920365e-07, 2.98496853e-02, 5.96990217e-02, 8.95483581e-02,
        1.19397695e-01, 1.49247031e-01, 1.79096367e-01, 2.08945704e-01,
        2.38795040e-01, 2.68644377e-01, 2.98493713e-01]),
 <a list of 10 Patch objects>)
```



```
dp =pd.read_csv("community-based.csv", index_col=0)
```

```
dp.describe()
```

	A	A_adcl	B	B_adcl	C	C_adcl	\
count	100.000000	100.000000	100.000000	100.000000	100.000000	100.000000	
mean	0.055034	0.277586	0.116043	0.557102	0.318386	0.799584	
std	0.024318	0.049060	0.032415	0.054352	0.048526	0.044362	

min	0.000000	0.142857	0.051724	0.406780	0.216667	0.666667
25%	0.037736	0.246154	0.095238	0.516532	0.285119	0.773585
50%	0.052178	0.275986	0.112007	0.563333	0.315789	0.800000
75%	0.072530	0.301587	0.136310	0.600000	0.346392	0.830769
max	0.129032	0.416667	0.209677	0.682540	0.448980	0.888889

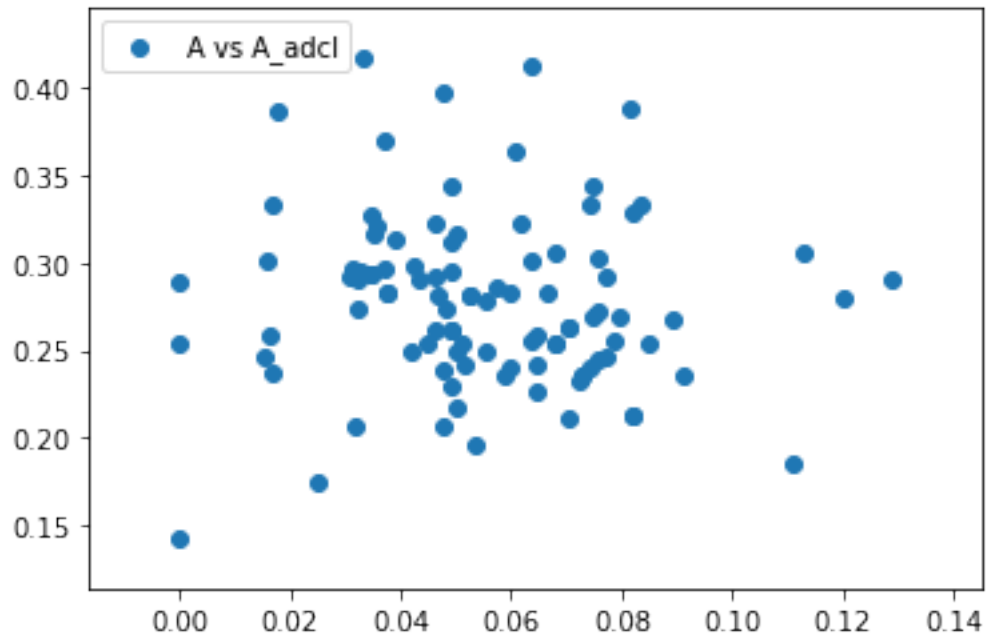
	D	D_adcl	E	E_adcl	...	A_prichness \
count	100.000000	100.000000	100.000000	100.000000	...	100.000000
mean	0.647292	0.962263	0.229600	0.989559	...	0.140234
std	0.051736	0.017520	0.045910	0.008343	...	0.035764
min	0.491803	0.916667	0.097222	0.975000	...	0.052632
25%	0.612455	0.949788	0.193768	0.983051	...	0.118395
50%	0.649561	0.966667	0.229508	0.984615	...	0.137147
75%	0.682738	0.978723	0.261943	1.000000	...	0.164801
max	0.786885	1.000000	0.339623	1.000000	...	0.235294

	B_prichness	C_prichness	D_prichness	E_prichness	A_mindistl \
count	100.000000	100.000000	100.000000	100.000000	100.000000
mean	0.147809	0.142149	0.160060	0.304304	0.025369
std	0.041895	0.037576	0.034324	0.042562	0.018591
min	0.042553	0.056604	0.080645	0.190476	0.000000
25%	0.123077	0.120690	0.144585	0.274194	0.015873
50%	0.146257	0.140351	0.155048	0.306452	0.020221
75%	0.179410	0.157540	0.180082	0.338524	0.034044
max	0.240741	0.254545	0.244898	0.393443	0.096774

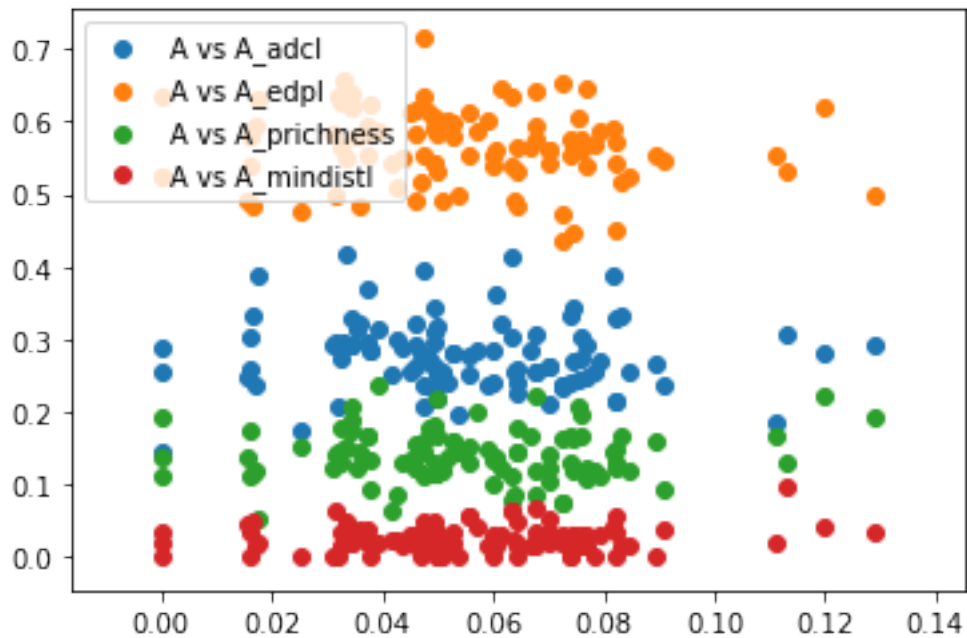
	B_mindistl	C_mindistl	D_mindistl	E_mindistl
count	100.000000	100.000000	100.000000	100.000000
mean	0.114361	0.097589	0.198869	0.290481
std	0.030794	0.030017	0.042808	0.047167
min	0.037037	0.033333	0.109091	0.156863
25%	0.095013	0.079132	0.172414	0.265789
50%	0.111111	0.096774	0.196400	0.298507
75%	0.134615	0.114754	0.229823	0.323077
max	0.187500	0.203390	0.333333	0.387755

[8 rows x 25 columns]

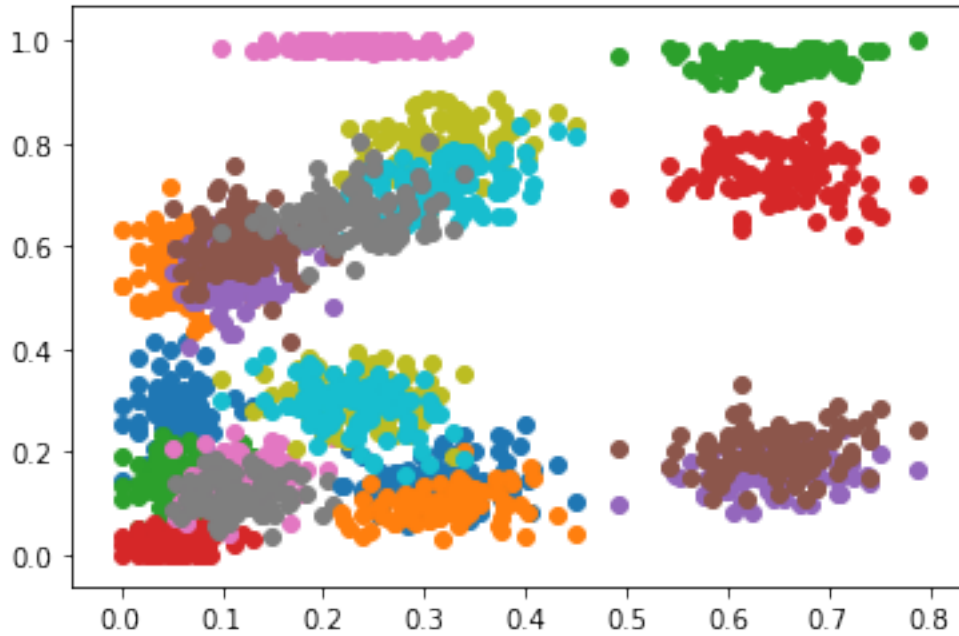
```
[85]: for score in ['A','B','C','D','E'][0:1]:
        for stats in ['_adcl', '_edpl', '_prichness', '_mindistl'][0:1]:
            plt.scatter(dp[score], dp[score+stats], label=score + ' vs ' +
→score+stats)
            plt.legend(loc='upper left')
            plt.show
```



```
[86]: for score in ['A', 'B', 'C', 'D', 'E'][0:1]:
      for stats in ['_adcl', '_edpl', '_prichness', '_mindistl'][0:5]:
          plt.scatter(dp[score], dp[score+stats], label=score + ' vs ' +
→score+stats)
      plt.legend(loc='upper left')
      plt.show
```




```
[87]: for score in ['A','B','C','D','E']:
        for stats in ['_adcl', '_edpl', '_prichness', '_mindistl']:
            plt.scatter(dp[score], dp[score+stats], label=score + ' vs ' +
→score+stats)
#         plt.legend(loc='upper left')
plt.show
```



```
[88]: dp.describe()
```

```
[88]:
```

	A	A_adcl	B	B_adcl	C	C_adcl	\
count	100.000000	100.000000	100.000000	100.000000	100.000000	100.000000	
mean	0.055034	0.277586	0.116043	0.557102	0.318386	0.799584	
std	0.024318	0.049060	0.032415	0.054352	0.048526	0.044362	
min	0.000000	0.142857	0.051724	0.406780	0.216667	0.666667	
25%	0.037736	0.246154	0.095238	0.516532	0.285119	0.773585	
50%	0.052178	0.275986	0.112007	0.563333	0.315789	0.800000	
75%	0.072530	0.301587	0.136310	0.600000	0.346392	0.830769	
max	0.129032	0.416667	0.209677	0.682540	0.448980	0.888889	

	D	D_adcl	E	E_adcl	...	A_prichness	\
count	100.000000	100.000000	100.000000	100.000000	...	100.000000	
mean	0.647292	0.962263	0.229600	0.989559	...	0.140234	
std	0.051736	0.017520	0.045910	0.008343	...	0.035764	
min	0.491803	0.916667	0.097222	0.975000	...	0.052632	
25%	0.612455	0.949788	0.193768	0.983051	...	0.118395	

50%	0.649561	0.966667	0.229508	0.984615	...	0.137147
75%	0.682738	0.978723	0.261943	1.000000	...	0.164801
max	0.786885	1.000000	0.339623	1.000000	...	0.235294

	B_prichness	C_prichness	D_prichness	E_prichness	A_mindist1	\
count	100.000000	100.000000	100.000000	100.000000	100.000000	
mean	0.147809	0.142149	0.160060	0.304304	0.025369	
std	0.041895	0.037576	0.034324	0.042562	0.018591	
min	0.042553	0.056604	0.080645	0.190476	0.000000	
25%	0.123077	0.120690	0.144585	0.274194	0.015873	
50%	0.146257	0.140351	0.155048	0.306452	0.020221	
75%	0.179410	0.157540	0.180082	0.338524	0.034044	
max	0.240741	0.254545	0.244898	0.393443	0.096774	

	B_mindist1	C_mindist1	D_mindist1	E_mindist1
count	100.000000	100.000000	100.000000	100.000000
mean	0.114361	0.097589	0.198869	0.290481
std	0.030794	0.030017	0.042808	0.047167
min	0.037037	0.033333	0.109091	0.156863
25%	0.095013	0.079132	0.172414	0.265789
50%	0.111111	0.096774	0.196400	0.298507
75%	0.134615	0.114754	0.229823	0.323077
max	0.187500	0.203390	0.333333	0.387755

[8 rows x 25 columns]

[]: