Samuel L. Peoples
BBIO 393: Computational Biology
Project 1

# Beta Diversity of Human Gut Microbiome Between Lean and Obese Twins

## 1. Introduction

This project was completed based on Jeff Gordon's *A Core Gut Microbiome in Obese and Lean Twins* which was used to determine if there exists a noticeable dissimilarity when comparing the Bray-Curtis-Faith Beta Diversity between Lean-vs-Lean, Obese-vs-Obese, and Lean-vs-Obese Twins.

This metric is defined as "a non-phylogenetic β diversity metric" where $BC_{jk} = \frac{\Sigma_i |X_{ij} - X_{ik}|}{\Sigma_i (X_{ij} - X_{ik})}$, and $i$ represents OTUs, and $j,k$ are the samples being compared. This is completed using the *beta-diversity.py* script from within a *Qiime2* virtual machine. The annotated script is as-follows.

## 2. Import the appropriate packages.

The *Biom* package is used to load the table, which may be used further in Project 2, but a majority of the preprocessing is completed with separate scripts, so this may be removed if it is ultimately not used. *Matplotlib, Scipi, Pylab, Numpy,* and *Random* are used to facilitate any math and visualization that needs to be done, while *OS* and *CSV* are used to save and load different files throughout the process.

The main function is displayed first, so to provide a blueprint for the progression of the script. The *filepath* is defined so to keep track of the absolute path to the project folder location, which contains all the files necessary for this project. The original OTU table is passed into a rarefaction function, and the associated tables are created, ultimately to be visualized as distributions and boxplots for visual comparison.

```python
1.  from biom import load_table
2.  import matplotlib.pyplot as plt
3.  import scipy.stats as stats
4.  import matplotlib as mpl
5.  import pandas as pd
6.  import pylab as pl
7.  import numpy as np
8.  import random
9.  import csv
10. import os
11.
12. def main():
13.     """
14.     Define the path to the script, which should be in
15.     the same directory as study_77_011618-113533
16.     """
17.     filepath = '/home/qiime/Desktop/Project1'
18.
19.     # Rarefaction of biom table, assigns path as variable
20.     rarefaction = single_rarefaction(filepath)
21.
22.     # Parses the tables and labels
23.     table = parse_table(filepath,rarefaction)
24.     betas = betaDiversity(filepath)
25.     lean, obese = label_tables(filepath, betas)
26.
27.     # Lean-vs-Lean, Obese-vs-Obese, and Lean-vs-Obese tables
28.     LL_stats, OO_stats, LO_stats = break_table(filepath,betas, lean, obese)
29.
30.     # Create the boxplots
31.     boxplots(filepath, LL_stats,OO_stats,LO_stats)
32.
```

```
33.     # Plot the Distributions
34.     distribution(filepath,LL_stats,OO_stats,LO_stats)
```

# 3. Define the helper functions.

The *single_rarefaction* function is used to access the original *.biom* file, and calls the *single_rarefaction.py* script with a rarefaction sample-size of 1000. This makes sure that each column sums to the same value, and will ensure that our results are more accurate. The filepath of the new *biom* file is returned.

```
1.  def single_rarefaction(filepath):
2.      """
3.      Runs single_rarefaction.py and returns output file location
4.      Sample size set to 1000
5.      """
6.      input = '/study_77_011618-113533/processed_data/219_otu_table.biom'
7.      output = '/study_77_011618-113533/processed_data/219_otu_table_even1000.biom'
8.      sample_size = '1000'
9.      os.system('single_rarefaction.py'+
10.               ' -i '+filepath+input+
11.               ' -o '+filepath+output+
12.               ' -d '+sample_size)
13.     return output
```

This function is used to ensure that the counts of the columns in the new table sum to 1000, and saves the new *.biom* file to a variable, which is left unused during this project, but may be accessed during the second project.

```
1.  def parse_table(filepath,rarefaction):
2.      """
3.      Passes the rarified table to be parsed
4.      Returns parsed table
5.      """
6.      # Load the table
7.      biom = load_table(filepath+rarefaction)
8.
9.      # Define the observations and samples
10.     otus = biom.ids('observation')
11.     samples = biom.ids('sample')
12.     m = biom.matrix_data
13.     data = [pd.SparseSeries(m[i].toarray().ravel()) for i in np.arange(m.shape[0])]
14.     table = pd.SparseDataFrame(data, index=otus, columns = samples)
15.     counts(table)
16.     return table
17.
18. def counts(table):
19.     """
20.     Counts the sum of the columns to verify that the rarefaction
21.     had been completed properly.
22.     """
23.     for sum in table.sum(axis=0):
24.         if sum != 1000:
25.             print("single_rarefaction.py failure.")
26.     print("Rarefaction check completed.")
```

The *betaDiversity* function is used to calculate the dissimilarity matrix, which is saved to the *stats* directory. The Bray-Curtis-Faith metric is defined here, where a Uni-Frac metric will be implemented in Project 2 to improve accuracy with the introduction of a phylogenetic tree. The dissimilarity matrix is loaded from the output csv file, and returned as a DataFrame.

```
1.  def betaDiversity(filepath):
2.      """
3.      Runs beta_diversity.py with a bray_curtis_faith metric,
4.      and returns the beta diversity matrix as a DataFrame
```

```
5.         """
6.         input = '/study_77_011618-113533/processed_data/219_otu_table_even1000.biom'
7.         output_dir = '/stats'
8.         output = '/bray_curtis_faith_219_otu_table_even1000.txt'
9.         metric = 'bray_curtis_faith'
10.        os.system('beta_diversity.py'
11.                     +' -i '+filepath+input
12.                     +' -o '+filepath+output_dir
13.                     +' -m '+metric)
14.        return pd.DataFrame.from_csv(filepath+output_dir+output, sep ='\t')
```

The *label_tables* function is used to associate the sample IDs with their obesity category from the mapping file; a simpler mapping was used after removing excess feature columns. Two lists are returned containing sample IDs that are verified to be in the dissimilarity matrix for each of the respective categories.

```
1.  def label_tables(filepath, betas):
2.          """
3.          Returns label lists of sample IDs which
4.          correspond to their values in the mapping file
5.          """
6.          ez_path = '/study_77_011618-113533/mapping_files/easy_mapping.csv'
7.          lean=[]
8.          obese=[]
9.          lines = []
10.         with open(filepath+ez_path,'r') as mapping:
11.             for line in mapping:
12.                 lines.append(line.split(","))
13.             mapping.close()
14.         for line in lines:
15.             if str(line[2]).__contains__("Lean"):
16.                 if line[0] in betas.keys():
17.                     lean.append(line[0])
18.             elif str(line[2]).__contains__("Obese"):
19.                 if line[0] in betas.keys():
20.                     obese.append(line[0])
21.         return lean,obese
```

The *break_table* function is used to separate the comparisons into three categories, Lean-vs-Lean, Obese-vs-Obese, and Lean-vs-Obese. The stats are populated into separate lists, verifying that the distance is not equal to zero, which prevents self-comparisons from being considered in the distributions. Summary statistics printed to the console and saved to a text file, while all distances are also saved to their respective *CSV* files. Separating the values in this manner allows us to visualize the box-plots and distributions, making comparisons about the dissimilarities when comparing the different microbial communities.

```
1.  def break_table(filepath, betas, lean, obese):
2.          """
3.          Breaks the beta diversity matrix into Lean-vs-Lean,
4.          Obese-vs-Obese, and Lean-vs-Obese. The function will
5.          not allow comparison of self (no distances of zero),
6.          and will return lists of distances for each value.
7.          Each list is saved to its own csv file.
8.          """
9.          LL_stats = []
10.         OO_stats = []
11.         LO_stats = []
12.
13.         # Separate distances to comparisons of
14.         # Lean-Lean, Obese-Obese, Lean-Obese
15.         for i in range(len(lean)):
16.             for j in range(len(betas[lean[i]].keys())):
17.                 if betas[lean[i]].keys()[j] in lean:
18.                     if betas[lean[i]][j] != 0:
19.                         LL_stats.append(betas[lean[i]][j])
```

```
20.            elif betas[lean[i]].keys()[j] in obese:
21.                LO_stats.append(betas[lean[i]][j])
22.      for i in range(len(obese)):
23.          for j in range(len(betas[obese[i]].keys())):
24.              if betas[obese[i]].keys()[j] in obese:
25.                  if betas[obese[i]][j]!= 0:
26.                      OO_stats.append(betas[obese[i]][j])
27.
28.      # Print the stats of the three categories
29.      print("Lean-vs-Lean:\n "+
30.            "Low: "+str(min(LL_stats))+       # .345
31.            "\n Mean: "+str(np.mean(LL_stats))+  # .845
32.            "\n Std: "+str(np.std(LL_stats))+    # .072
33.            "\n High: "+str(max(LL_stats))+      # .982
34.            "\n Count: "+str(len(LL_stats)))     # 3422
35.      print("")
36.      print("Obese-vs-Obese:\n "+
37.            "Low: "+str(min(OO_stats))+       # .337
38.            "\n Mean: "+str(np.mean(OO_stats))+  # .851
39.            "\n Std: "+str(np.std(OO_stats))+    # .071
40.            "\n High: "+str(max(OO_stats))+      # .998
41.            "\n Count: "+str(len(OO_stats)))     # 34782
42.      print("")
43.      print("Lean-vs-Obese:\n "+
44.            "Low: "+str(min(LO_stats))+       # .582
45.            "\n Mean: "+str(np.mean(LO_stats))+  # .857
46.            "\n Std: "+str(np.std(LO_stats))+    # .064
47.            "\n High: "+str(max(LO_stats))+      # .994
48.            "\n Count: "+str(len(LO_stats)))     # 11033
49.
50.      # Save the Bray-Curtis-Faith distances to output csvs,
51.      with open(filepath+"/stats/LL_stats.csv", "wb") as f:
52.          wr = csv.writer(f,quoting=csv.QUOTE_ALL)
53.          wr.writerow(LL_stats)
54.          f.close()
55.      with open(filepath+"/stats/OO_stats.csv", "wb") as f:
56.          wr = csv.writer(f,quoting=csv.QUOTE_ALL)
57.          wr.writerow(OO_stats)
58.          f.close()
59.      with open(filepath+"/stats/LO_stats.csv", "wb") as f:
60.          wr = csv.writer(f,quoting=csv.QUOTE_ALL)
61.          wr.writerow(LO_stats)
62.          f.close()
63.
64.      return LL_stats, OO_stats, LO_stats
```

## 4. Visualize the output.

The boxplot figure is created from the respective stats lists, allowing the visualization of the different comparison categories. The data is defined as list of samples from the respective categories, which is a simple random sample of 1000 from each. This visualization is discussed in detail in the conclusion.

```
1.  def boxplots(filepath, LL_stats, OO_stats, LO_stats):
2.      """
3.      Creates three boxplots based on the LL, OO, and LO
4.      categories, and saves the boxplots to boxplot.png
5.      """
6.      mpl.use('agg')
7.      data = sample(LL_stats, OO_stats, LO_stats)
8.      fig = plt.figure(1, figsize=(9,6))
9.      ax = fig.add_subplot(111)
10.     bp = ax.boxplot(data)
11.     ax.set_title('Comparison of Beta Diversity Dissimilarity\n Sample of 1000')
12.     ax.set_xticklabels(['Lean vs Lean','Obese vs Obese', 'Lean vs Obese'])
13.     [flier.set(marker='o',color='#e7298a', alpha=.3) for flier in bp['fliers']]
14.     fig.savefig(filepath+"/figures/boxplot.png", bbox_inches='tight')
```

```
15.
16. def sample(LL_stats, OO_stats, LO_stats):
17.     """
18.         Returns a random sample of the three categories
19.     """
20.     return [random.sample(LL_stats, 1000),
21.             random.sample(OO_stats, 1000),
22.             random.sample(LO_stats, 1000)]
```

Individual distributions are plotted to better distinguish the differences in dissimilarity. Each visualization is discussed in detail in the conclusion.

```
1.  def distribution(filepath,LL_stats,OO_stats,LO_stats):
2.      """
3.          Visualize the Distributions of the
4.          Lean-vs-Lean, Obese-vs-Obese, and Lean-vs-Obese
5.          categories. Saves figures to output file
6.      """
7.      LL = sorted(LL_stats)
8.      OO = sorted(OO_stats)
9.      LO = sorted(LO_stats)
10.     LL_fit = stats.norm.pdf(LL,np.mean(LL),np.std(LL))
11.     OO_fit = stats.norm.pdf(OO,np.mean(OO),np.std(OO))
12.     LO_fit = stats.norm.pdf(LO,np.mean(LO),np.std(LO))
13.
14.     # Clear any figures that may be present.
15.     pl.clf()
16.
17.     # Plot the Lean vs Lean
18.     pl.plot(LL,LL_fit,'-o')
19.     pl.hist(LL,normed=True)
20.     pl.title("Bray-Curtis-Faith Dissimilarity of Lean vs Lean Twins")
21.     pl.savefig(filepath+"/figures/LL_Distribution.png")
22.     pl.clf()
23.
24.     # Plot the Obese vs Obese
25.     pl.plot(OO,OO_fit,'-o')
26.     pl.hist(OO,normed=True)
27.     pl.title("Bray-Curtis-Faith Dissimilarity of Obese vs Obese Twins")
28.     pl.savefig(filepath+"/figures/OO_Distribution.png")
29.     pl.clf()
30.
31.     #Plot the Lean vs Obese
32.     pl.plot(LO,LO_fit,'-o')
33.     pl.hist(LO,normed=True)
34.     pl.title("Bray-Curtis-Faith Dissimilarity of Lean vs Obese Twins")
35.     pl.savefig(filepath+"/figures/LO_Distribution.png")
36.     pl.clf()
```
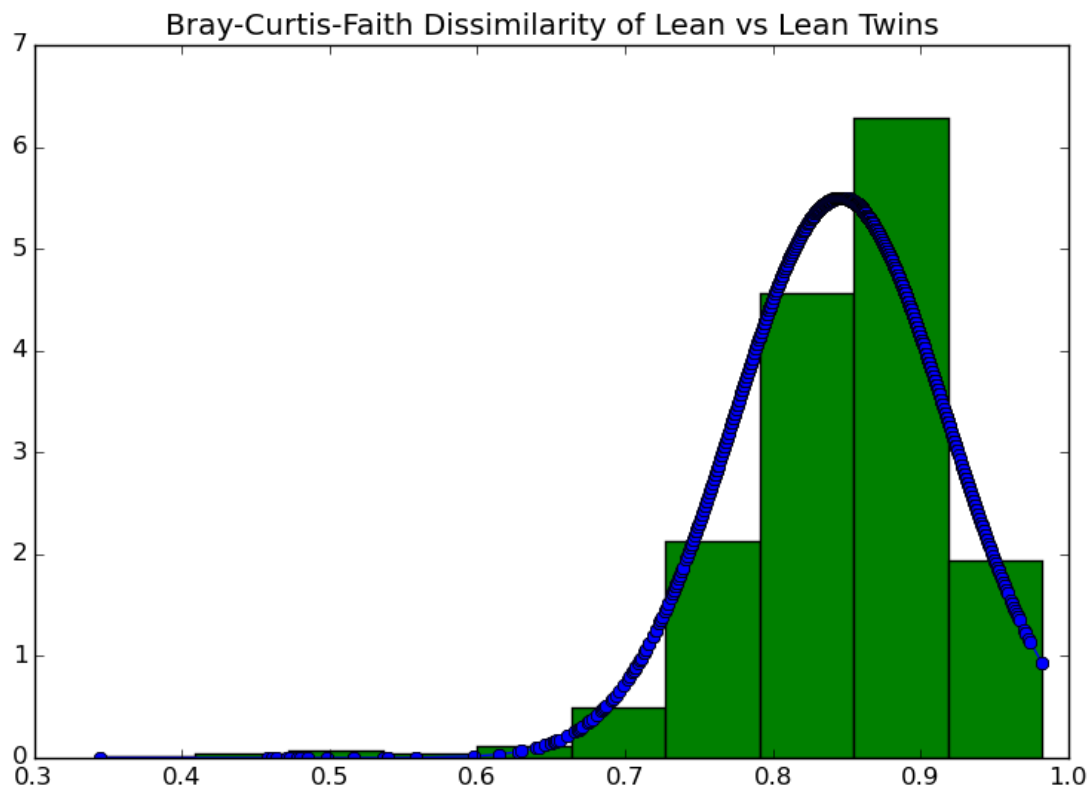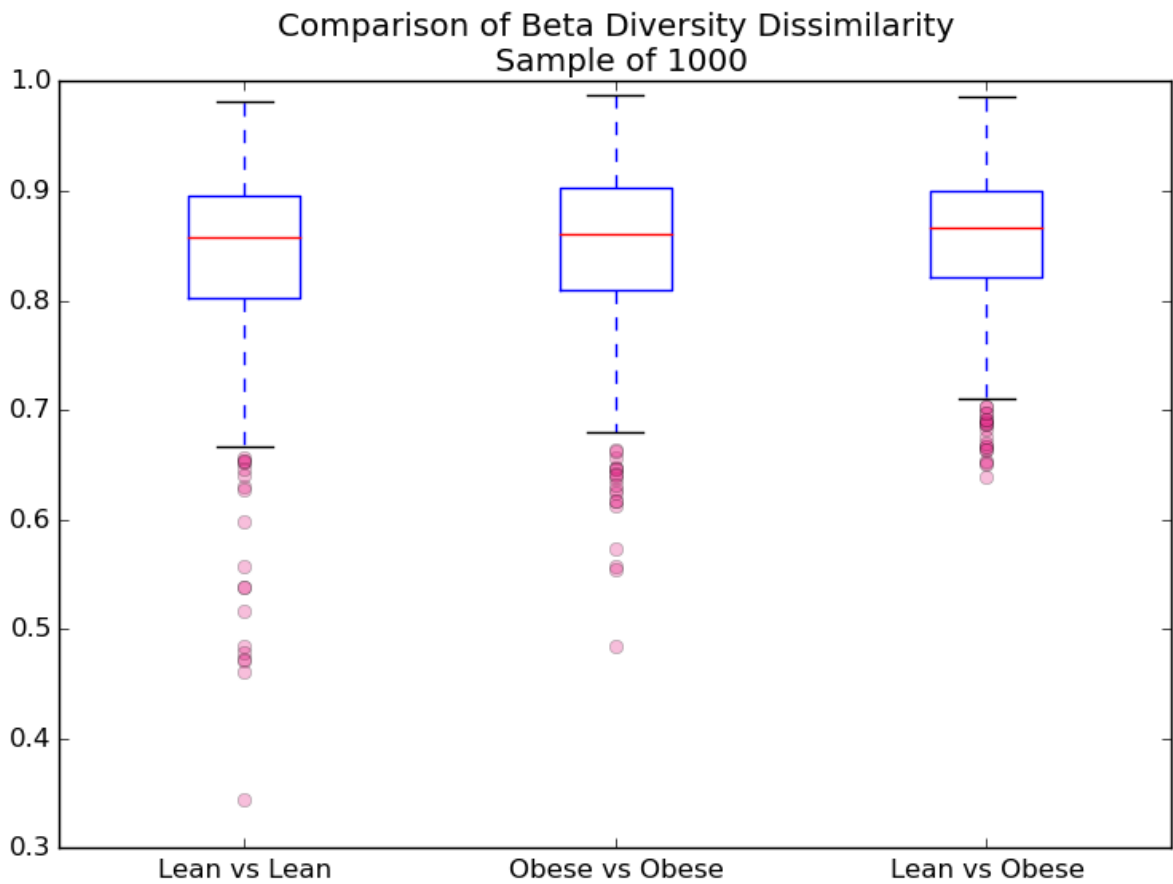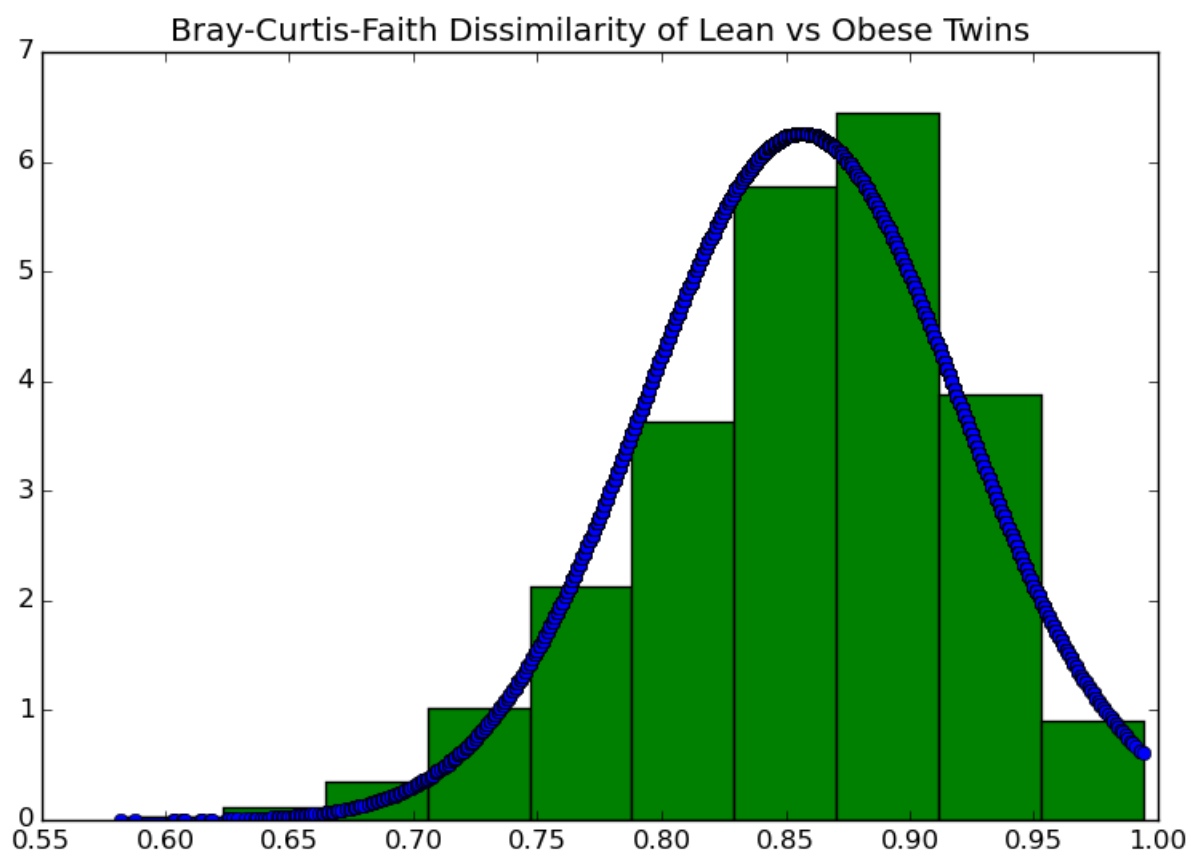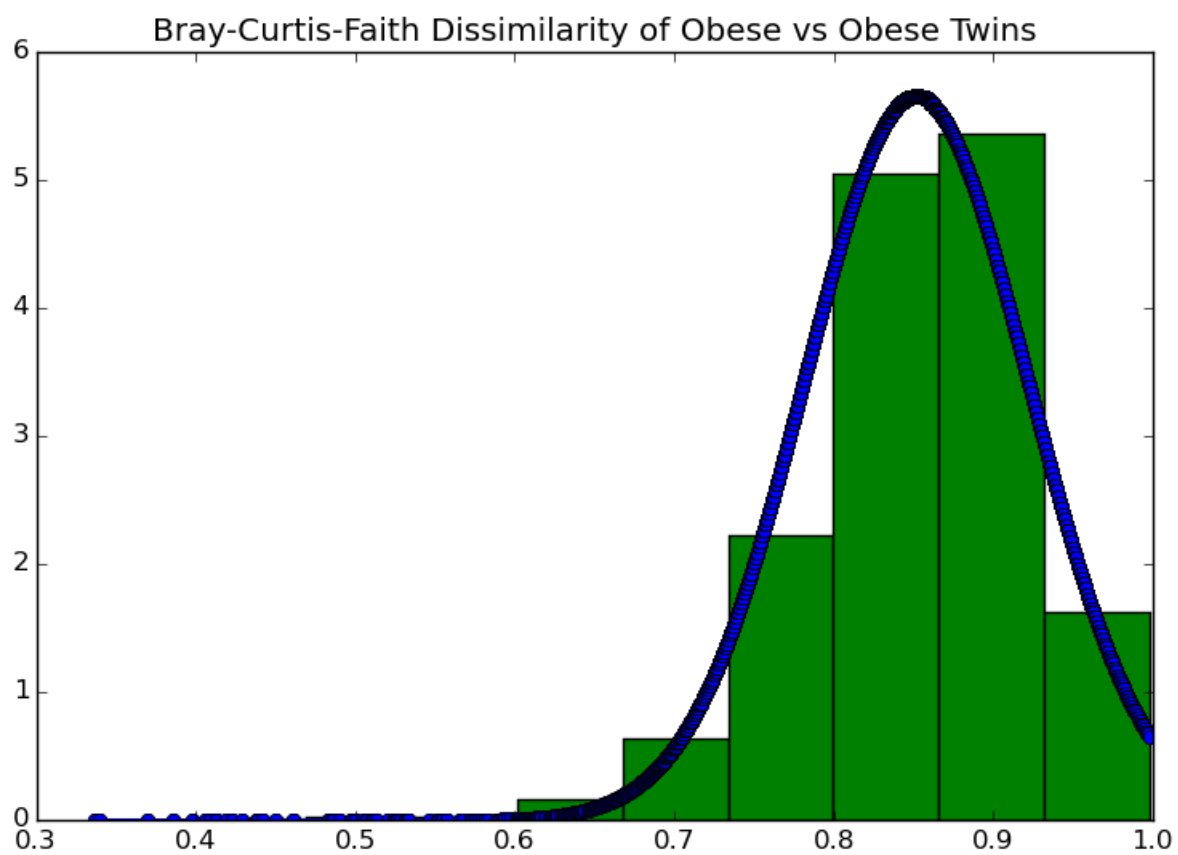
# 5. Conclusion

When viewing the boxplots and distributions, we can see that the distribution is reduced when comparing dislike obesity categories, where Lean-vs-Obese has no samples below roughly 0.60, and Lean-vs-Lean, Obese-vs-Obese each have samples with dissimilarities as low as roughly 0.30.

This difference opens up many questions which can be answered in Project 2, which include the statistical significance of this difference, and comparison to the same analysis under different diversity metrics. The increasing trend of higher dissimilarity in the boxplots suggests that there may be a relationship between the diversity of the gut microbiome and obesity, and could be answered with more investigation.

There is also room for improvement, where accuracy of the analysis can be increased with the use of a phylogeny-based beta diversity metric (Uni-Frac), which would affect the distances in the dissimilarity matrix. This evaluation can also be improved by implementing a Monte-Carlo (or similar) simulation with many samples of 1000 comparisons. The script could also be generalized following Project 2, allowing for more parameterized functionality, which would be applicable to more datasets.

Being able to visually recognize the dissimilarity reveals the potential differences that exist within the microbial communities of lean vs obese individuals, providing inspiration for further investigation and insights.

**Comparison of Beta Diversity Dissimilarity**
**Sample of 1000**



**Bray-Curtis-Faith Dissimilarity of Lean vs Lean Twins**

Bray-Curtis-Faith Dissimilarity of Obese vs Obese Twins

Bray-Curtis-Faith Dissimilarity of Lean vs Obese Twins

# 6. References

Qiime: http://qiime.org/
Qiita: https://qiita.ucsd.edu/
Bray-Curtis-Faith: http://readiab.org/book/0.1.3/3/1#4.1.1
Gut Microbiome Dataset: https://qiita.ucsd.edu/study/description/77
Biom-Format: http://biom-format.org/documentation/biom_format.html
Lozupone/ Knight's UniFrac: http://aem.asm.org/content/71/12/8228.full
Qualitative/ Quantitative: https://www.ncbi.nlm.nih.gov/pubmed/17220268
Weighted UniFrac: https://liorpachter.wordpress.com/2013/09/18/unifrac-revealed/