

Microbiome Data Analysis with MicrobiomeAnalyst

User ID: guest4564085299521310481

November 27, 2019

1 Data Processing and Normalization

1.1 Reading and Processing the Raw Data

MicrobiomeAnalyst accepts count data in a variety of formats generated in microbiome studies including plain text table, biom format as well as output from mothur pipeline. User need to upload their data in one of three available formats and specify the taxonomic labels when uploading their data in order for MicrobiomeAnalyst to process the taxonomic information correctly. The hierarchical information for taxa can either be present within abundance data or uploaded as a separate taxonomy table file (.txt or .csv format). Also, uploading metadata file is required as plain text (.txt or.csv) with all three formats.

1.1.1 Reading abundance count data table

The abundance count data should be uploaded in tab-delimited text (.txt) or comma separated values (.csv) format. Samples are represented in columns, while rows contains the information about the features. Metadata file contains additional information about samples such as experimental factors or sample grouping.

A total of 12 samples and 437 features or taxa are present. The sample data contains a total of 12 samples and 2 sample variables. The OTUs are annotated as Others/Not_{specific}label.

1.1.2 Data Integrity Check

Before data analysis, a data integrity check is performed to make sure that all the necessary information has been collected. The sample variable should contain atleast two groups to perform most of the comparative analysis. *By default, sample variables which are found to be constant and continuous in nature will be removed from further analysis. Additionally, features just present in one sample will also be discarded from the data.* Figure 1 shows the library size for inspection of each sample.

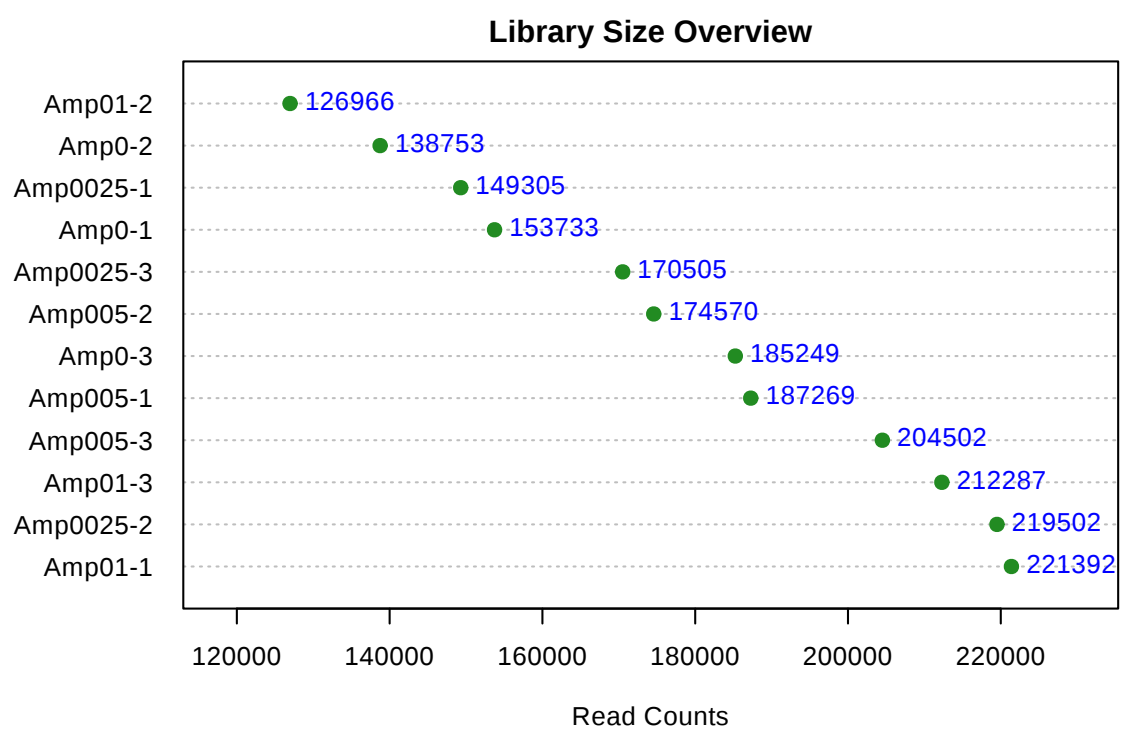


Figure 1: Library size Overview

1.1.3 Data Filtering

The purpose of the data filtering is to identify and remove features that are unlikely to be of use when modeling the data. No phenotype information are used in the filtering process, so the result can be used with any downstream analysis. This step can usually improves the results. Features having low count and variance can be removed during the filtration step. Features having very few counts are filtered based on their abundance levels (minimum counts) across samples (prevalence). Other than sample prevalence, such features can also be detected using minimum count cutoff based on their mean and median values. Features or taxa with constant or less variable abundances are invaluable for comparative analysis. Such features are filtered based on their inter-quantile ranges, standard deviations or coefficient of variations. *By default, features having zero counts across all the samples, or only appears in one sample will be removed from further analysis.*

A total of 36 low abundance features were removed based on prevalence. A total of 17 low variance features were removed based on iqr. The number of features remains after the data filtering step: 148

1.2 Data Normalization

The data is stored as a table with one sample per column and one variable (taxa or OTU) per row. The normalization procedures implemented below are grouped into three categories. Data rarefaction and scaling based methods deal with uneven sequencing depths by bringing samples to the same scale for comparison. While transformation based methods account for sparsity, compositionality, and large variations within the data. You can use one or combine all three to achieve better results. For more information about these methods, please refer to the paper by Weiss et al.¹ The normalization consists of the following options:

1. Data rarefying (with or without replacement)
2. Data scaling:
 - Total sum scaling (TSS)
 - Cumulative sum scaling (CSS)
 - Upper-quantile normalization (UQ)
3. Data transformation :
 - Relative log expression (RLE)
 - Trimmed mean of M-values (TMM)
 - Centered log ratio (CLR)

No data rarefaction was performed. Performed total sum normalization. No data transformation was performed.

2 Marker Gene Analysis

MicrobiomeAnalyst offers a variety of methods commonly used in microbiome data analysis. They include:

1. Visual exploration:
 - Stacked bar/area plot
 - Rarefaction curve

¹Weiss et al. *Normalization and microbial differential abundance strategies depend upon data characteristics*, Microbiome 2017

- Phylogenetic tree
 - Heat tree
2. Community profiling:
 - Alpha diversity analysis
 - Beta Diversity analysis
 - Core microbiome analysis
 3. Clustering analysis:
 - Heatmap
 - Dendrogram
 - Correlation analysis
 - Pattern Search
 4. Differential abundance analysis:
 - Univariate analysis
 - metagenomeSeq
 - RNAseq methods
 5. Biomarker analysis:
 - LEfSe
 - Random Forests
 6. Predictive functional profiling:
 - PICRUSt
 - Tax4Fun

2.1 Visual Exploration

These methods are used to visualize the taxonomic composition of community through direct quantitative comparison of abundances. MicrobiomeAnalyst provides an option to view this composition at various taxonomic levels (phylum, class, order) using either stacked bar/stacked area plot or piechart. Viewing composition at higher-levels (phylum) provides a better picture than lower-levels (species) when the number of species in a community is large and diversified. Additionally, such taxonomic abundance or composition can be viewed at community-level (all samples), sample-group level (based on experimental factor) or at individual sample-level. Taxa with very low read counts can also be collapsed into **Others** category using a count cutoff based on either sum or median of their counts across all samples or all groups. Merging such minor taxa will help in better visualization of significant taxonomic patterns in data. Figure 2 shows the taxonomic composition using Stacked bar/area plot.

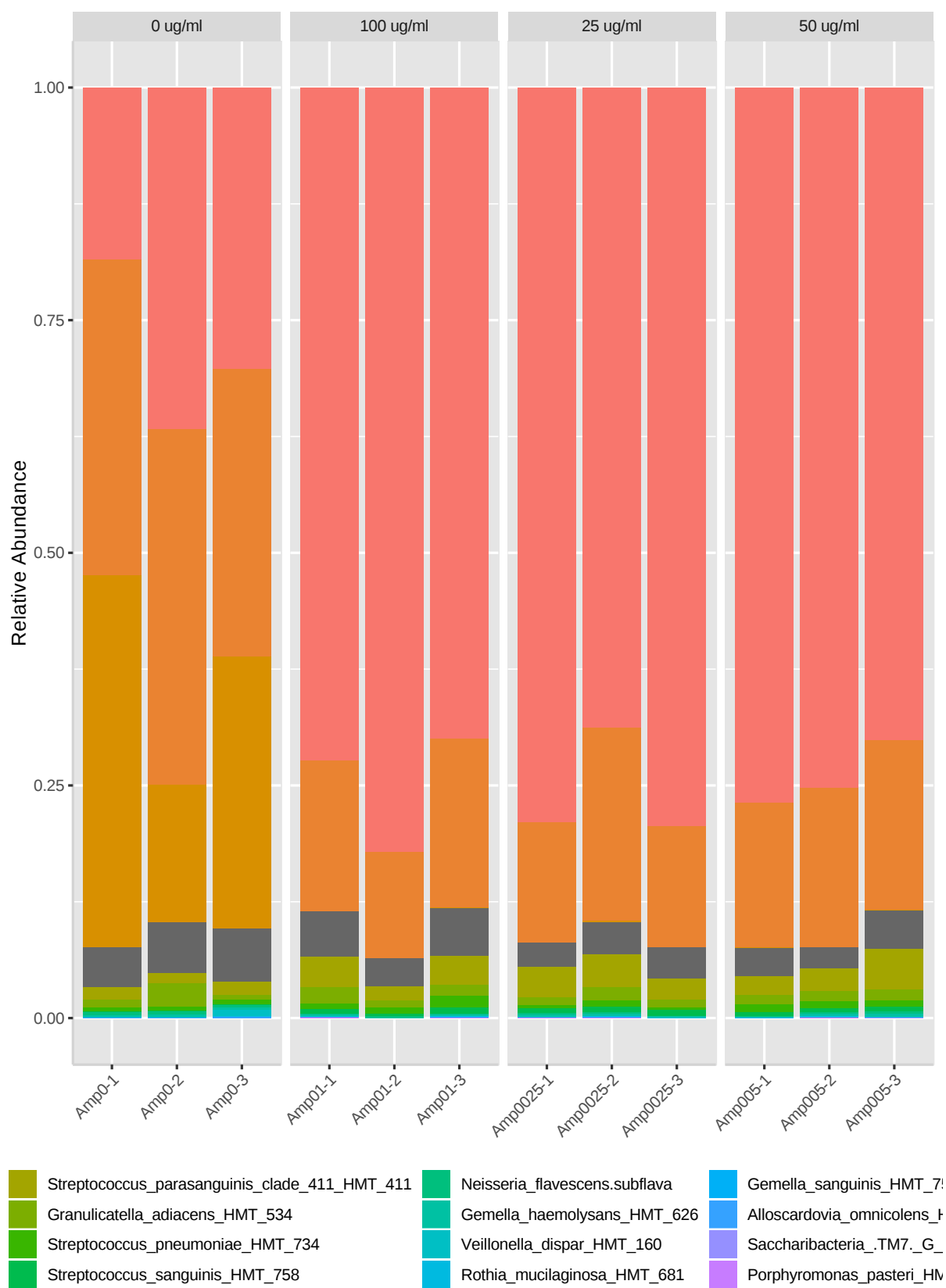


Figure 2: Taxonomic composition of community at Species level using Stacked Bar plot

2.2 Rarefaction curve analysis

This method is used to present relationship between number of OTUs and number of sequences. It can infer if the reads of a sample are enough to reach plateau, which means that with increasing of sequences, the gain of newly discovered OTUs is limited. If sequence depth of some samples are not enough, you may consider to resequence these samples or removed from downstream analysis. User can choose different metadata variables as group, line colors and line types. Rarefaction curve analysis is performed using the modified function `ggrare` originated from `ranacapa` package²

Figure 3 shows the rarefaction curve.

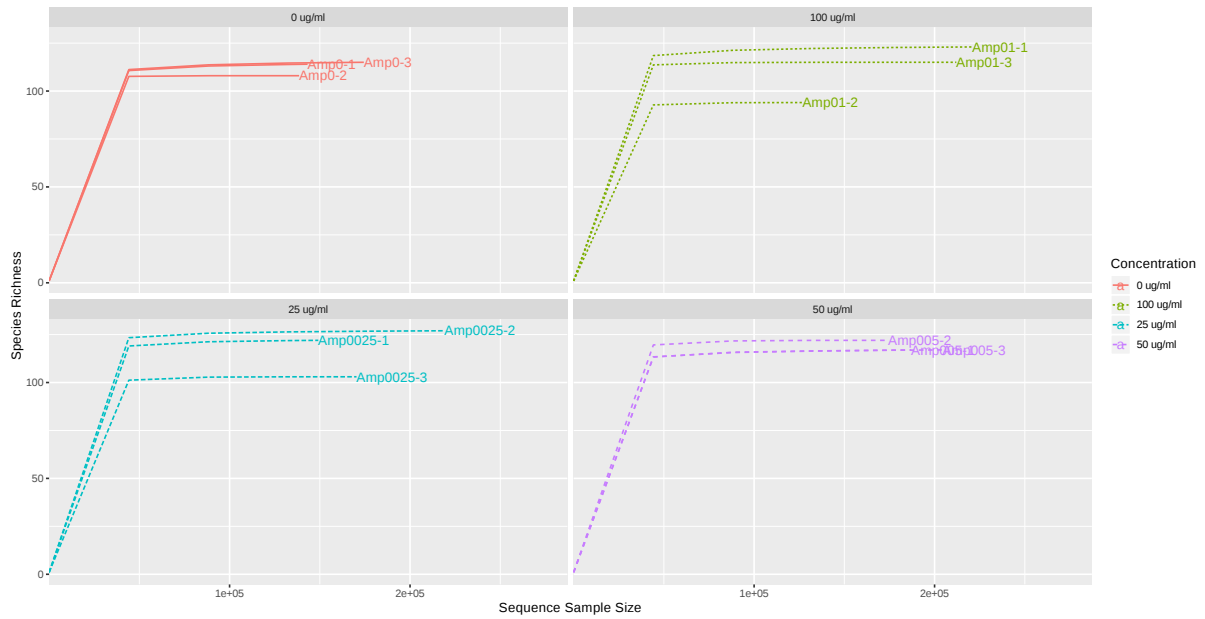


Figure 3: Rarefaction curve using `filtered` dataset

²Gaurav S. Kandlikar *ranacapa: An R package and Shiny web app to explore environmental DNA data with exploratory statistics and interactive visualizations.*, 2018.

2.3 Alpha diversity analysis

This method is used to measure the diversity present within a sample or community. Alpha diversity can be characterized via the total number of species (richness), the abundances of the species (evenness) or measures that considered both richness and evenness. How these measures estimates the diversity is need to be considered when performing alpha-diversity analysis. User can choose from richness based measure such as Observed index which calculates the actual number of unique taxa observed in each sample. While the Chao1 and ACE measures estimate the richness by inferring out the number of rare organisms that may have lost due to undersampling. Also,there are indices such as Shannon, Simpson and Fisher in which along with the number (richness), the abundance of organisms (evenness) is also measured to describe the actual diversity of a community.

Alpha diversity analysis is performed using the `phyloseq` package³. The results are plotted across samples and reviewed as box plots for each group or experimental factor. Further, the statistical significance of grouping based on experimental factor is also estimated using either parametric or non-parametric test. Figure 4 shows the alpha diversity measure across all the samples for given diversity index. Figure 5 shows the diversity distribution using box plot for a given group or experimental factor.

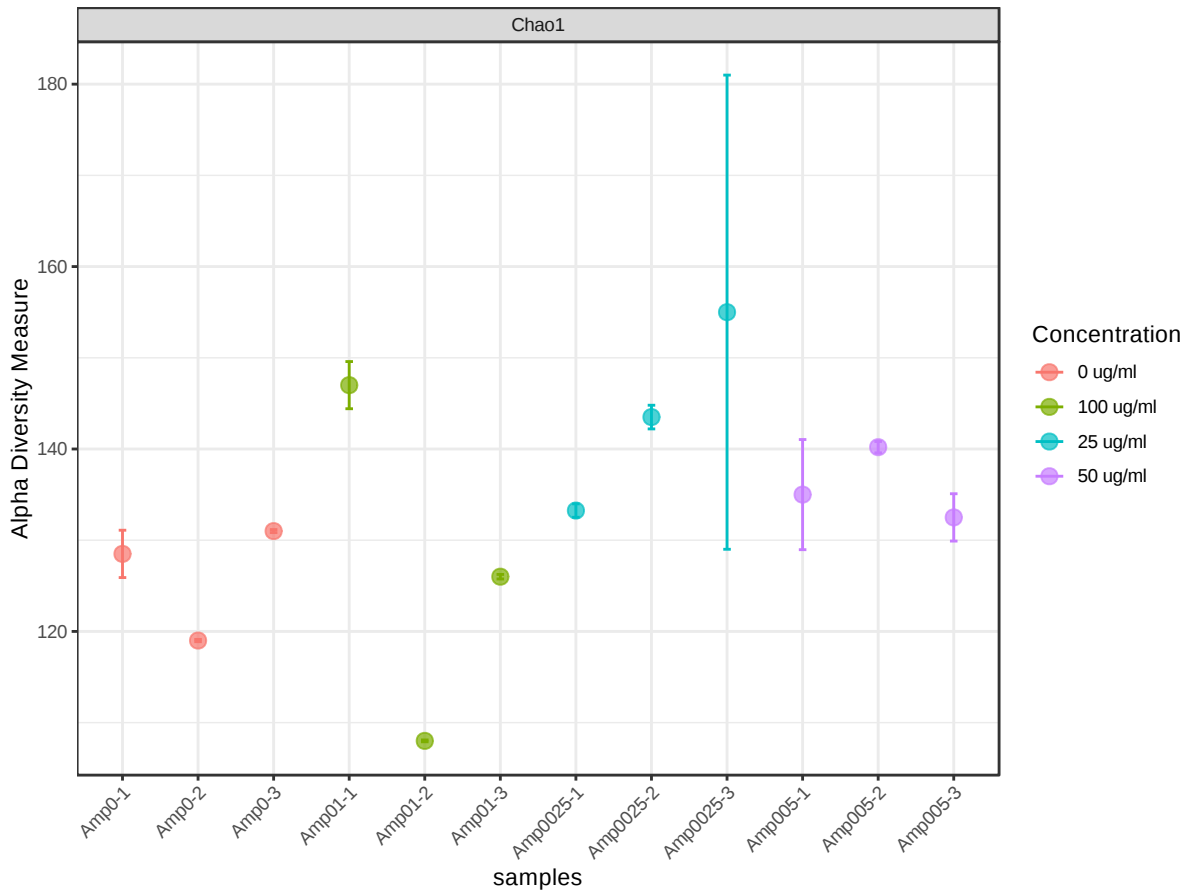


Figure 4: Alpha-diversity measure using `Chao1` at OTU level across all the samples. The samples are represented on X-axis and their estimated diversity on Y-axis. Each sample is colored based on `Concentration` class

³Paul J. McMurdie *phyloseq: An R package for reproducible interactive analysis and graphics of microbiome census data.*, 2013, R package version 1.19

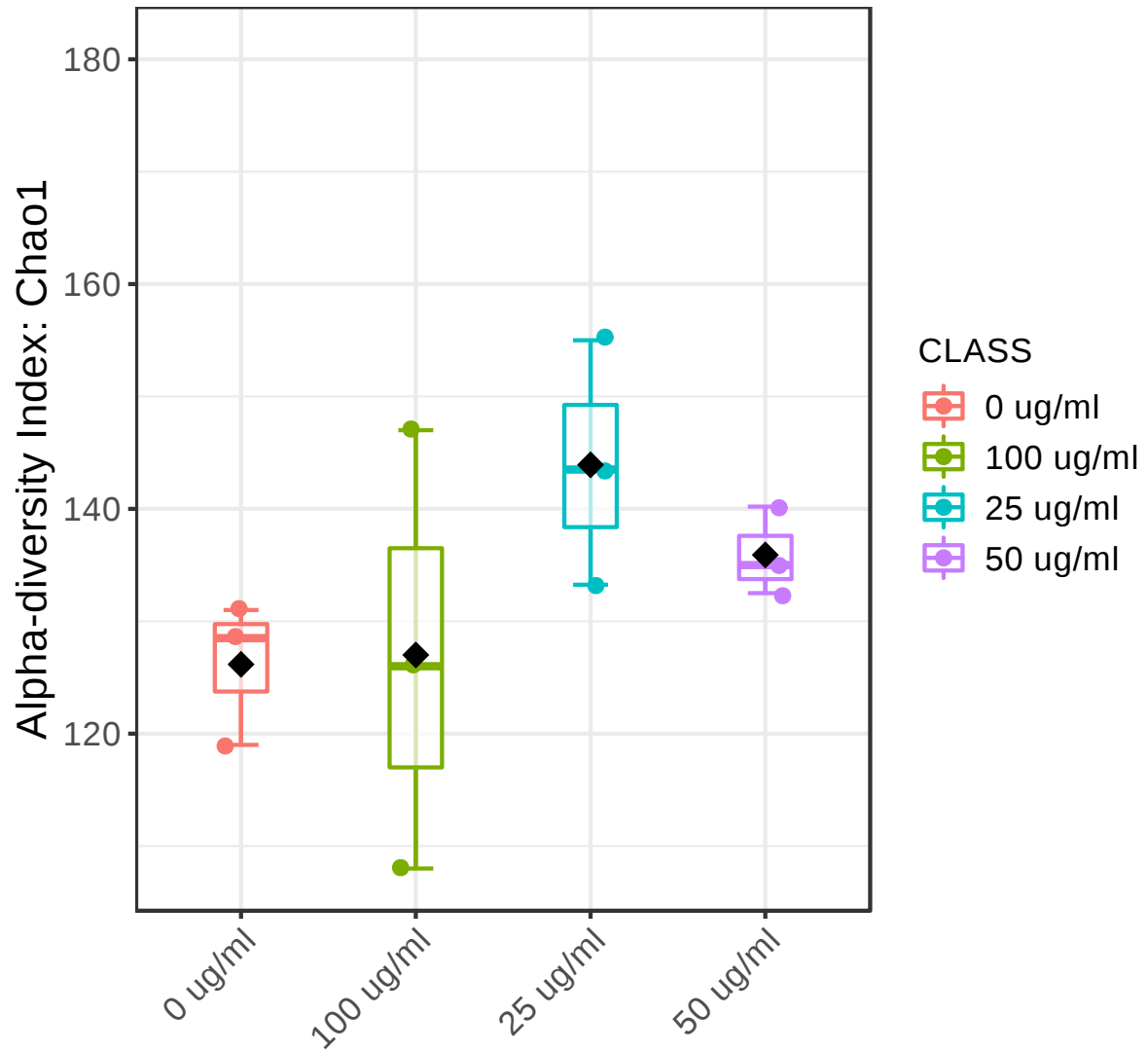


Figure 5: Alpha-diversity measure using Chao1 at OTU level represented as boxplot. Each boxplot represents the diversity distribution of a group present within Concentration class [Statistical significance: p-value: 0.1718; [Kruskal-Wallis] statistic: 5]

2.4 Beta diversity Analysis

This method provides a way to compare the diversity or composition between two samples or microbial communities. These methods compare the changes in the presence/absence or abundance of thousands of taxa present in a dataset and summarize these into how 'similar' or 'dissimilar' two samples. Each sample gets compared to every other sample generating a distance or dissimilarity matrix. Two parameters need to be considered when performing beta diversity analysis. The first one is how similarity or distance between sample is measured which includes non-phylogenetic (Bray-Curtis distance, Shannon index, Jaccard index) and phylogenetic-based (weighted and unweighted UniFrac) distances. The other parameter is how to visualize such dissimilarity matrix in lower dimensions. Ordination-based methods such as Principle Coordinate Analysis (PCoA) and non-metric multidimensional scaling (NMDS) are used to visualize these matrix in 2 or 3-D plot where each point represents the entire microbiome of a single sample. Each axis reflects the percent of the variation between the samples with the X-axis representing the highest dimension of variation and the Y-axis representing the second highest dimension of variation. Further, each point or sample displayed on PCoA or NMDS plots is colored based on either sample group, features alpha diversity measures, or the abundance levels of a specific feature.

Also, the statistical significance of the clustering pattern in ordination plots can be evaluated using anyone among Permutational ANOVA (PERMANOVA), Analysis of group Similarities (ANOSIM) and Homogeneity of Group Dispersions (PERMDISP).

Beta diversity analysis is performed using the `phyloseq` package⁴. Figure 6 shows the ordination plot represented in 2-D; Statistical significance is found out using [PERMANOVA] F-value: 1.9851; R-squared: 0.42674; p-value < 0.048 .

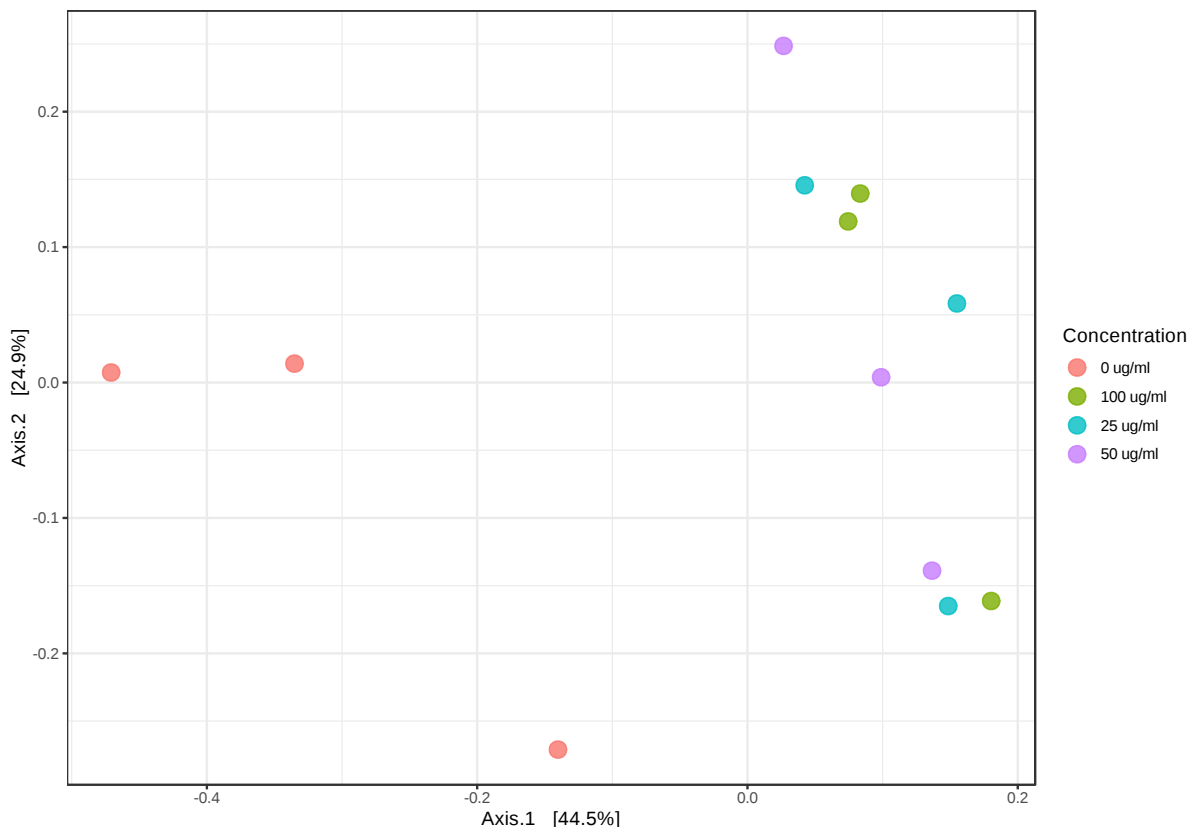


Figure 6: 2-D PCoA plot using `bray` distance. The explained variances are shown in brackets.

⁴Paul J. McMurdie *phyloseq: An R package for reproducible interactive analysis and graphics of microbiome census data.*, 2013, R package version 1.19

2.5 Hierarchical Clustering

In hierarchical cluster analysis, each sample begins as a separate cluster and the algorithm proceeds to combine them until all samples belong to one cluster. Two parameters need to be considered when performing hierarchical clustering. The first one is how similarity or distance between sample is measured which includes Bray-Curtis distance, Shannon index, Jaccard index, weighted and unweighted UniFrac. The other parameter is clustering algorithms, including average linkage (clustering uses the centroids of the observations), complete linkage (clustering uses the farthest pair of observations between the two groups), single linkage (clustering uses the closest pair of observations) and Ward's linkage (clustering to minimize the sum of squares of any two clusters). In MicrobiomeAnalyst, the result of clustering analysis are supported using Heatmap and dendrogram.

Hierarchical clustering is performed with the `hclust` function in package `stat`. Figure 7 shows the clustering result in the form of a heatmap.

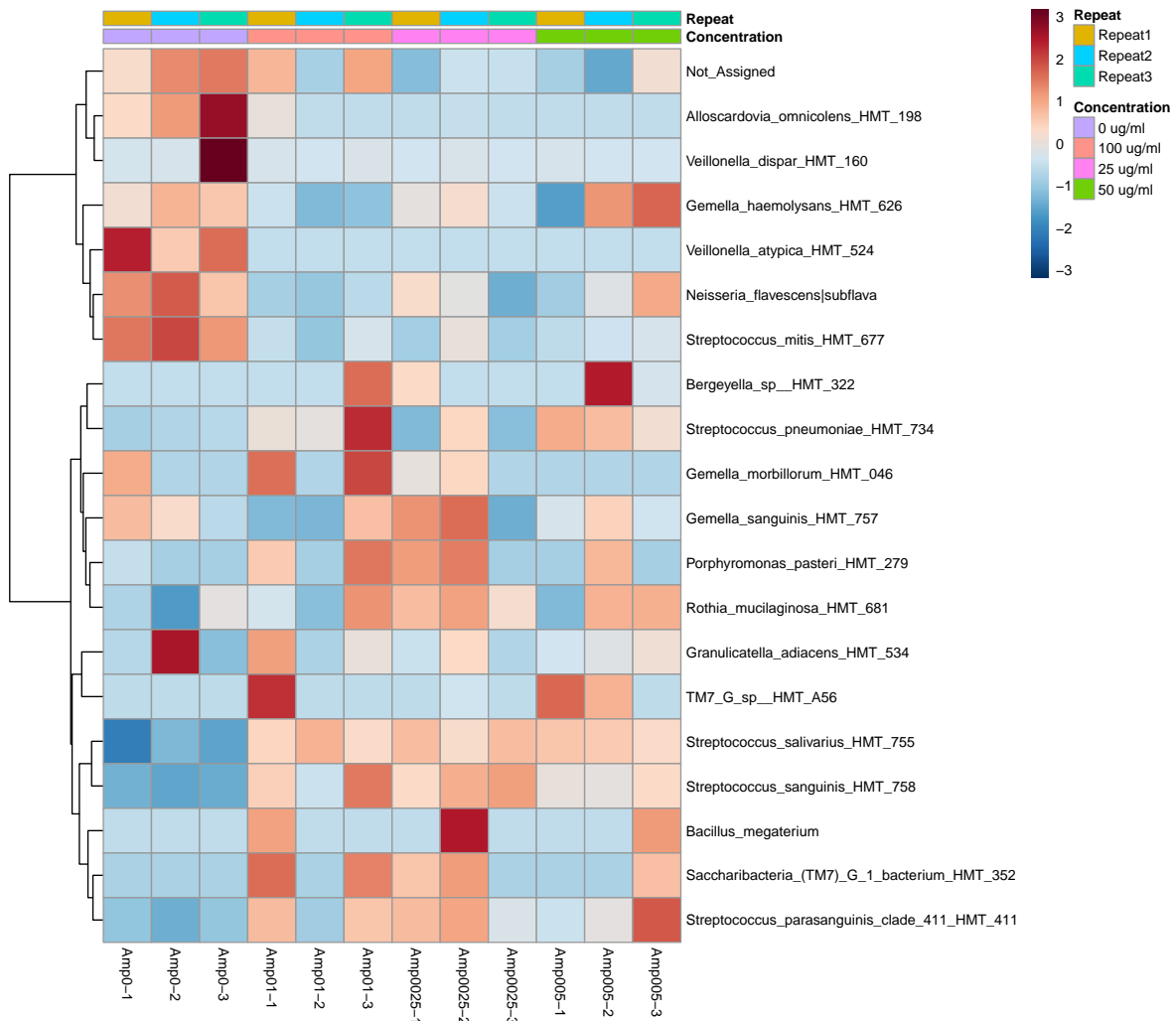


Figure 7: Clustering result shown as heatmap (distance measure using `euclidean` and clustering algorithm using `ward.D` at `Species` level)

2.6 metagenomeSeq

This method is specifically designed to evaluate differential abundance in sparse marker-gene survey data. This method combines Cumulative Sum Scaling (CSS) normalization with zero-inflated Log-Normal mixture model (fitFeature) or zero-inflated Gaussian (fitZIG) distribution mixture to account for undersampling and sparsity in OTU count data. In the fitZIG model, the count distribution is modeled as a mixture of two distributions. The zeros present in count data are modeled using point mass at zero, while remaining log transformed counts follows a normal distribution. On the other hand, fitFeature model shapes the count distribution using zero-inflated lognormal model. fitFeature model can be only used and recommended for two groups, whereas fitZIG model works on multiple groups for differential abundance testing. This method outperforms other methods in the detection of differentially abundant rare features. Features are considered to be significant based on their adjusted p-value. The default is `adj.p-value cutoff = 0.05`.

This method is available as a `metagenomeSeq` R package ⁵. Table 1 shows the important features identified by metagenomeSeq at `Species` level

Table 1: Important features identified by metagenomeSeq

	Features	Pvalues	FDR
1	<i>Veillonella atypica</i> _HMT_524	0.00	0.00
2	<i>Saccharibacteria</i> _(TM7)_G_1_bacterium_HMT_352	0.00	0.01
3	<i>Streptococcus sanguinis</i> _HMT_758	0.01	0.08
4	<i>Bacillus megaterium</i>	0.03	0.15
5	<i>Porphyromonas pasteri</i> _HMT_279	0.04	0.15
6	<i>Streptococcus salivarius</i> _HMT_755	0.06	0.20
7	<i>Neisseria flavescens</i> subflava	0.13	0.37
8	<i>Bergeyella sp.</i> _HMT_322	0.19	0.47
9	<i>Gemella morbillorum</i> _HMT_046	0.24	0.48
10	<i>Streptococcus parasanguinis</i> _clade_411_HMT_411	0.25	0.48
11	<i>Streptococcus pneumoniae</i> _HMT_734	0.27	0.48
12	<i>Streptococcus mitis</i> _HMT_677	0.29	0.48
13	<i>Gemella sanguinis</i> _HMT_757	0.40	0.61
14	<i>Veillonella dispar</i> _HMT_160	0.44	0.63
15	<i>Gemella haemolysans</i> _HMT_626	0.51	0.69
16	TM7_G_sp_HMT_A56	0.62	0.78
17	<i>Rothia mucilaginosa</i> _HMT_681	0.72	0.84
18	<i>Granulicatella adiacens</i> _HMT_534	0.76	0.85
19	<i>Alloscardovia omnicolens</i> _HMT_198	0.97	0.97
20	Not_Assigned	0.97	0.97

⁵JN Paulson *metagenomeSeq: Statistical analysis for sparse high-throughput sequencing*, 2013, R Bioconductor package version 1.20.1

2.7 Random Forest (RF)

Random Forest is a supervised learning algorithm suitable for high dimensional data analysis. It uses an ensemble of classification trees, each of which is grown by random feature selection from a bootstrap sample at each branch. Class prediction is based on the majority vote of the ensemble. RF also provides other useful information such as OOB (out-of-bag) error and variable importance measure. During tree construction, about one-third of the instances are left out of the bootstrap sample. This OOB data is then used as test sample to obtain an unbiased estimate of the classification error (OOB error). Variable importance is evaluated by measuring the increase of the OOB error when it is permuted. The outlier measures are based on the proximities during tree construction.

RF analysis is performed using the **randomForest** package⁶. Table 2 shows the confusion matrix of random forest. Figure 8 shows the cumulative error rates of random forest analysis for given parameters. Figure 9 shows the important features ranked by random forest. The OOB error is 0.75

Figure 8: Cumulative error rates by Random Forest classification. The overall error rate is shown as the black line; the red and green lines represent the error rates for each class.

	0 ug/ml	100 ug/ml	25 ug/ml	50 ug/ml	class.error
0 ug/ml	3.00	0.00	0.00	0.00	0.00
100 ug/ml	0.00	0.00	2.00	1.00	1.00
25 ug/ml	0.00	2.00	0.00	1.00	1.00
50 ug/ml	0.00	1.00	2.00	0.00	1.00

Table 2: Random Forest Classification Performance

Figure 9: Significant features identified by Random Forest. The features are ranked by the mean decrease in classification accuracy when they are permuted.

⁶Andy Liaw and Matthew Wiener. *Classification and Regression by randomForest*, 2002, R News

The report was generated on Wed Nov 27 13:10:56 2019 with R version 3.6.1 (2019-07-05).