

# A STATISTICAL PERSPECTIVE ON THE CHALLENGES IN MOLECULAR MICROBIAL BIOLOGY

BY PRATHEEPA JEGANATHAN<sup>1</sup> AND SUSAN P. HOLMES<sup>2,\*</sup>

<sup>1</sup>*Department of Mathematics and Statistics, McMaster University, [jeganp1@mcmaster.ca](mailto:jeganp1@mcmaster.ca)*

<sup>2</sup>*Department of Statistics, Stanford University, , [susan@stat.stanford.edu](mailto:susan@stat.stanford.edu)*

High throughput sequencing (HTS)-based technology enables identifying and quantifying non-culturable microbial organisms in all environments. Microbial sequences have enhanced our understanding of the human microbiome, the soil and plant environment, and the marine environment. All molecular microbial data pose statistical challenges due to contamination sequences from reagents, batch effects, unequal sampling, and undetected taxa. Technical biases and heteroscedasticity have the strongest effects, but different strains across subjects and environments also make direct differential abundance testing unwieldy. We provide an introduction to a few statistical tools that can overcome some of these difficulties and demonstrate those tools on an example. We show how standard statistical methods, such as simple hierarchical mixture and topic models, can facilitate inferences on latent microbial communities. We also review some nonparametric Bayesian approaches that combine visualization and uncertainty quantification. The intersection of molecular microbial biology and statistics is an exciting new venue. Finally, we list some of the important open problems that would benefit from more careful statistical method development

**1. Introduction.** High-throughput sequencing (HTS) enables the characterization of variation in microbial diversity in naturally changing or experimentally perturbed environments. Several technologies in molecular microbiology have revolutionized the resolution at which many environments can now be studied. The first is marker-gene sequencing (usually identified as microbiome studies); these use small regions (V4-V6) of a particular gene (the 16S rRNA gene, occasionally others) that serves as a “fingerprint” or signature for each species or strain of bacteria. More comprehensive profiling of all genes present in these microbial communities is available through complete shotgun sequencing. Shotgun sequence analysis, known as metagenomics, uses all the nucleic acids in a specimen to provide a comprehensive inventory of both the genes and taxa present, sometimes invoking what is known as the *pangenome* (Quince et al., 2017a). This can be done by Bayesian classifiers using short strings of nucleotide (k-mer) occurrences (Rosen, Reichenberger and Rosenfeld, 2011) or by genome assembly (Lu et al., 2017). Quince et al. (2017b) presents a review and comparison of these different approaches.

These two basic sequence-based methods have enabled major advances in biological, agricultural, and environmental research. For example, vaginal microbiome studies in pregnant women have shown that the reduced *Lactobacillus* species in the vagina is a risk factor for premature birth (Callahan et al., 2017; DiGiulio et al., 2015). The gut microbiome is associated with an increased risk of type 1 diabetes and inflammatory bowel diseases in children (Kostic et al., 2015; Gevers et al., 2014). Outside of human biology, microbial biota has been an important source for monitoring the marine environment (Thompson et al., 2017; Gilbert, Jansson and Knight, 2014). Recent work in soil science has also shown the

---

*Keywords and phrases:* Microbial ecology, Bayesian data analysis, hierarchical mixture models, latent Dirichlet allocation, Bayesian nonparametric ordination, sequencing data, quality control.

power of cross-sectional microbiome experimental designs to detect ecological perturbations (Delgado-Baquerizo et al., 2016; Ramirez et al., 2018) and enable climate change monitoring (Cavicchioli et al., 2019). A recent review of microbiome-based agro-management has shown how these can improve agricultural production, promote plant growth and health, maintain resistance against diseases, and quantify environmental stress (Compant et al., 2019).

The statistical analyses of abundance and diversity of microbial sequence data have many commonalities with standard ecological studies; this means that many of the downstream tools are already available in the statistical ecologists' toolbox. Spatial, multivariate, and longitudinal methods are central to the field and through the development of tools such as `phyloseq` (McMurdie and Holmes, 2013) we have tried to create bridges between the raw molecular genomic read data (as well as the phylogenetic relationships between taxa) and the data structures such as contingency tables already used in ecology and evolutionary biology.

The current review will concentrate on the statistical challenges inherent in the analyses of the sequencing reads organized as contingency tables. We will use the columns to represent the biological specimens, also called "specimen-samples," as in Table 1, the rows are labeled for the taxonomic strains known as Amplicon Sequence Variants (ASVs) (Callahan, McMurdie and Holmes, 2017). These rows are not predefined before the data become available and are inferred by denoising the raw sequences using the read frequencies and their quality scores, see Callahan et al. (2016a,b). Contrary to some recent statements in the literature (Quinn et al., 2018; Gloor et al., 2017; Silverman et al., 2017), the data themselves cannot be considered *compositional* since the number of rows (i.e., strains) and their definition is not known *a priori*, and there is always a substantial and variable proportion of reads that cannot be annotated. The total number of reads in each column corresponds to sequencing depths for each of the specimens and are often called the library sizes: we will show that these can be modeled as gamma-Poisson random variables. Strain level resolution is now also available for shotgun metagenomic data through Bayesian co-occurrence analyses (Quince et al. (2020)), and many new research problems arise when this higher resolution of analysis is used.

At high resolution, the unknown parameters of interest are the true prevalence of each of the microbial strains or ASVs and the differences between prevalences across different treatment groups or locations. The strain prevalence parameters sum to one within each specimen  $j$ , thus if we restrict ourselves to a finite number ( $m$ ) of possible taxa, the  $\mathbf{p}_j = (p_{i,j}, i = 1, \dots, m)$  do belong to the simplex. However, even if we postulate *a priori* knowledge on  $m$ , the estimation of the  $\mathbf{p}_j$  parameters poses several challenges. There is often between subject or location variability in the different strains detected in the specimens; thus, strains do not co-occur systematically across subjects or locations, and the count table is sparse with a large proportion of zeros (both technological and biological). There are also substantive experimental issues with DNA extraction and PCR amplification that preclude direct quantification of prevalences (McLaren, Willis and Callahan, 2019).

Many authors start their analysis by transforming the reads to relative abundances (transforming the counts by dividing by column sums) and then taking the log (Kurtz et al., 2015). However, all versions of centered log-ratio transformations ignore the underlying heteroscedasticities of the prevalence estimates due to library size variation and hinder valid downstream statistical inference. A statistical solution we recommend is variance stabilizing transforms similar to those applied to other genomic data such as RNA-seq, see (McMurdie and Holmes, 2014). If there is a small number of pre-specified taxa measured, it is possible to use weighted chi-square distances to account for compositionality, as proposed for ecological tables by Greenacre (2010a, 2011), this incorporates the column sums and thus does not ignore any of the information in the data. In cases where the rare taxa are the study focus, it can be beneficial to transform the data into a 0/1 table with just indicators of presence for

taxa that appear abundant above a certain number of reads (typically 2 to 4, depending on the library sizes). On the other hand, if comparisons are to be made between abundances in the core microbiome (a set of common taxa present in most of the specimens), the rare species will be filtered out. As is the case with all statistical transformations, the choice has to be an informed one, and there is no one method that fits all situations, as some environments are infinitely diverse, whereas others are very sparse. Starting by doing the robust rank-threshold transformation we describe in Section 5 is often a good first step.

Currently, many experiments that involve microbiota also include complementary assays that provide metabolic and gene expression profiles through the use of Mass Spectrometry (Grégory et al., 2018) and RNA-seq transcriptomics (Franzosa et al., 2014). These data also include clinical or chemical covariates measured on the same specimens facilitating a “holistic” understanding of the system under study. A recent review by Sankaran and Holmes (2019a) shows how many data integration approaches based on matrix decompositions and cross-table correlations can combine multiple data domain types effectively. We will not cover the multidomain-multimodal aspects here.

In Section 2, we introduce the form in which the data are collected, the format of the tables and contiguous information, and the statistical notation. Then, in Section 3, we describe an example problem and the data we will use as an illustration of the different tools available. Goodness of fit tests for each taxon enabled us to build generative models for these data that we can use for simulation or power studies such as the one illustrating the strain switching problem in Section 8.1. We also recommend Bayesian hierarchical models that we used to remove DNA contamination—a common additive source of measurement error. For the downstream analysis, it is first necessary to account for the library size - a source of multiplicative variation - discussed in Section 4. We then describe best practices for variance stabilization and robust rank-based transformation, as well as first-order rank-based dimensionality reductions in Section 5. We recommend starting all analyses with data exploration and visualizations of both the raw and transformed count data, including heatmaps and clustering of the specimens illustrated in Section 6. We discuss network analysis, many combinations of different distances, multivariate methods, and how we can enhance ordination using phylogenetic trees that account for evolutionary relationships in the taxa in Section 7. In each of the above sections, we show how these exploratory tools may not identify high-resolution variability and may be challenging to interpret.

An important goal in microbial ecology is the inference of differences in taxonomic abundances in different environments or treatment groups. In Section 8, we briefly review permutational analyses of distance matrices, generalized linear models for differential abundance methods and discuss the importance of identifying the bacterial communities and their differences across conditions. This motivates the use of latent Dirichlet allocation (LDA) topic models that we discuss and illustrate in Section 9. To help the biologists interpret the topics, we enhance this topic model analysis with a visualization that incorporates the phylogenetic tree, and we show it in the same section. We identify promising research direction is Bayesian nonparametric approaches that can directly account for the uncertainty in measured data, learning latent structure, and flexible enough through the use of hierarchies that can account for experimental design and random effects. For example, to account for the growing number of taxa and uncertainties in ordination analyses, we demonstrate a Bayesian nonparametric factor ordination method in Section 10. We conclude by discussing the example dataset results and open problems from a statistical perspective.

**2. Microbiome data.** In both marker-gene sequencing and shotgun metagenomics, the core of the data consists of a contingency table of read counts with specimens recorded in the columns and taxa (ASVs) identified in the rows as in Table 1. We will use marker-genes

throughout the text, but the methods and challenges discussed also apply to shotgun metagenomics. Associated with these counts is a matrix of specimen information with specimen identifiers on rows and column variables such as subject identity, sequencing batch, type of specimen (control or specimen types such as blood or gut or placenta, etc.), as shown in Table 2. It is often beneficial to annotate the Amplicon Sequence Variants (ASVs) (Callahan, McMurdie and Holmes, 2017) using a matrix of taxa identifiers at several selected taxonomy levels such as species name on rows and taxonomy levels (species, genus, phylum, family, order, class) on columns as shown in Table 3. Finally, evolutionarily relationship between taxa are formalized as a phylogenetic tree as in Figure 1.

TABLE 1

Count matrix  $\mathbf{K} \in \mathbb{R}^{m \times N}$  of  $n_1$  specimens,  $n_2$  controls and  $m$  taxa, where  $K_{ij}$  are the reads of taxon  $i$  in  $j$ -th specimen and  $K_{il}^0$  denote the number of reads of taxon  $i$  in  $l$ -th negative control.

	Specimen <sub>1</sub>	Specimen <sub>2</sub>	...	Specimen <sub><math>n_1</math></sub>	Control <sub><math>(n_1+1)</math></sub>	Control <sub><math>(n_1+2)</math></sub>	...	Control <sub><math>N</math></sub>
Taxa <sub>1</sub>	$K_{11}$	$K_{12}$	...	$K_{1n_1}$	$K_{1(n_1+1)}^0$	$K_{1(n_1+2)}^0$	...	$K_{1N}^0$
Taxa <sub>2</sub>	$K_{21}$	$K_{22}$	...	$K_{2n_1}$	$K_{2(n_1+1)}^0$	$K_{2(n_1+2)}^0$	...	$K_{2N}^0$
⋮	⋮					⋮		⋮
Taxa <sub><math>i</math></sub>	$K_{i1}$	$K_{i2}$	...	$K_{in_1}$	$K_{i(n_1+1)}^0$	$K_{i(n_1+2)}^0$	...	$K_{iN}^0$
⋮	⋮					⋮		⋮
Taxa <sub><math>m</math></sub>	$K_{m1}$	$K_{m2}$	...	$K_{mn_1}$	$K_{m(n_1+1)}^0$	$K_{m(n_1+2)}^0$	...	$K_{mN}^0$

TABLE 2

Sample data, a matrix of specimen information with specimen identifiers on rows and column variables.

	Specimen ID	Subject ID	Specimen type	Batch number
Specimen <sub>1</sub>	Specimen <sub>1</sub>	Subject <sub>1</sub>	Plasma	1
Specimen <sub>2</sub>	Specimen <sub>2</sub>	Subject <sub>2</sub>	Plasma	2
⋮	⋮	⋮	⋮	⋮
Specimen <sub><math>n_1</math></sub>	Specimen <sub><math>n_1</math></sub>	Subject <sub><math>n_1</math></sub>	Plasma	1
Control <sub><math>(n_1+1)</math></sub>	Control <sub><math>(n_1+1)</math></sub>	Reagent	Control	1
Control <sub><math>(n_1+2)</math></sub>	Control <sub><math>(n_1+2)</math></sub>	Library	Control	1
⋮	⋮	⋮	⋮	⋮
Control <sub><math>N</math></sub>	Control <sub><math>N</math></sub>	Reagent	Control	2

TABLE 3

Taxonomy table

	Kingdom	Phylum	Class	Order	Family	Genus
Taxa <sub>1</sub>	<i>Bacteria</i>	<i>Nitrospirae</i>	<i>Nitrospira</i>	<i>Nitrospirales</i>	<i>0319-6A21</i>	
Taxa <sub>2</sub>	<i>Bacteria</i>	<i>Acidobacteria</i>	<i>Blastocatellia</i>	<i>Blastocatellales</i>	<i>Blastocatellaceae_(SG_4)</i>	<i>DS-100</i>
⋮	⋮					
Taxa <sub><math>i</math></sub>	<i>Bacteria</i>	<i>Armatimonadetes</i>	<i>Armatimonadia</i>	<i>Armatimonadales</i>	<i>Armatimonadaceae</i>	<i>Armatimonas</i>
⋮	⋮					
Taxa <sub><math>m</math></sub>	<i>Bacteria</i>	<i>Chloroflexi</i>	<i>Chloroflexia</i>	<i>Herpetosiphonales</i>	<i>Herpetosiphonaceae</i>	<i>Herpetosiphon</i>

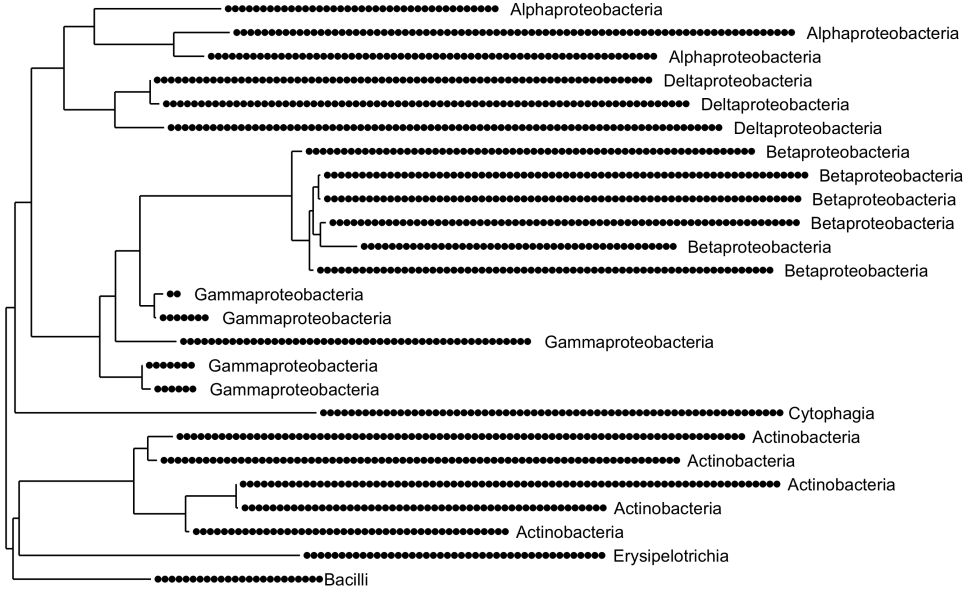


FIG 1. *Phylogenetic tree. Black points at each node corresponds to the specimens in which the corresponding taxa is present. Tip labels are class of the taxon.*

Table 1 shows the count table of  $n_1$  specimens,  $n_2$  controls and  $m$  taxa, where  $K_{ij}$  are the reads of taxon  $i$  in  $j$ -th specimen and  $K_{il}^0$  denote the number of reads of taxon  $i$  in  $l$ -th negative control - specimens consisting of molecular-grade water included in each extraction batch. We denote  $d_j$  and  $d_l^0$  the linear scaling factors for specimen  $j$  and control  $l$  that account for the library sizes — sums of the columns in the count Table 1 and can vary by orders of magnitude. For each specimen  $j$ , the read counts are in the vector  $\mathbf{K}_j = [K_{1j}, K_{2j}, \dots, K_{mj}]^T$ . Note that many bacteria are interdependent (in particular syntropic relations are frequent) precluding the use of a simple multinomial for these counts. We provide the details of a typical analysis for a real example in what follows.

**3. Example of 16S rRNA gene sequencing data.** We reanalyzed the 16S rRNA gene sequencing of the root endosphere specimens from [Fitzpatrick et al. \(2018\)](#). We retrieved this dataset from Google Dataset Search with the keyword of 16S rRNA gene sequencing and ASV. The authors described the specimen collection, library preparation and sequencing, and their microbial bioinformatics workflow. In plants, host organellar sequences such as mitochondrial and plastid share similar bacterial lineages and reduce the efficiency of quantifying microbial sequences. Universal peptide nucleic acids (PNA) have been used to limit the amplification of host-derived mitochondrial and plastid sequences. [Fitzpatrick et al. \(2018\)](#) designed their experiment to evaluate the efficacy of universal plastid peptide nucleic acids (O-pPNA) and Asteraceae-modified pPNA (M-pPNA) in limiting plastid host contaminants. Contamination varies across host plants, and [Fitzpatrick et al. \(2018\)](#) provided a validated framework for Asteraceae-modified pPNA so that there is no effect of pPNA type on bacterial detection. [Fitzpatrick et al. \(2018\)](#) identified less plastid contaminant sequences in Asteraceae with M-pPNA than O-pPNA. We used this data set as an example that shows typical analytical challenges posed by such data and provided some solutions inspired by statistical approaches.

We identified the host contaminant sequences from the taxonomy table. The data set had specimens from six Asteraceae and three non-Asteraceae plant types. The O-pPNA applied to three to ten replicates, whereas M-pPNA applied to three to five replicates. The root endosphere specimens from four out of six Asteraceae and all three non-Asteraceae plant types were sequenced with both pPNA (O and M) and had 18 paired specimens as in Supplementary Table 1. A common environment was chosen to grow all plant types from sterilized seeds, except *Centaurea solstitialis*, which were sampled in a field across France, Spain, and the USA.

We followed the filtering in Fitzpatrick et al. (2018). From 57,116 ASVs, we removed ASVs that lacked a kingdom assignment or were assigned to Archaea or Eukaryota, ASVs that lacked a phylum assignment, ASVs classified as mitochondria, ASVs classified as plastids. We retained plant types that had both endosphere or rhizosphere specimens and negative controls, removed specimens with less than 800 reads, and removed the fifth sequencing run that was bad. In the preprocessed data, there were 86 root endosphere and 25 control specimens for further analysis. The 6929 ASVs in specimens and controls were used in the Bayesian hierarchical model as implemented in the BARBI package to remove contaminants (Cheng et al., 2019).

**4. Models.** Although many analyses begin with exploratory analysis of transformed count data that identify outliers, biological variability and the interrelation between four different components of microbiome data discussed in Section 2, they cannot account for heterogeneity in the data. It can be beneficial to start by a simple generative model for the individual taxonomic counts. We start by describing a negative binomial (gamma-Poisson) goodness of fit test for each taxon. The test results justify this model and we will complement them with a Bayesian hierarchical model that removes DNA contamination, a supplementary additive error.

4.1. *Goodness of fit for taxon counts.* If we knew the true prevalence of different taxa  $p_{ij}$  and then sequence the same amplified DNA in technical replicates, we would expect to see a simple Poisson( $\lambda_{ij}$ ) variation in the  $K_{ij}$ . For instance, Grumaz et al. (2016) uses a Poisson model with healthy controls providing baseline proportions of each taxon to identify bacteria in blood infected patients after removing the well-known contaminant species *Xanthomonas*. A Poisson model with intensity estimated using negative and positive controls was used in Hong et al. (2018) to choose taxa for downstream analysis. However, this simple model needs to be enriched to accommodate other sources of variation, such as contaminant bacteria introduced during the sample preparation, sequencing run batch effects, and library size differences as well as the essential biological variation of interest.

We show that a negative binomial distribution (or equivalently gamma-Poisson) fits our example data set for the ASV counts well. We test the null hypothesis,  $H_0$ , that the ASV counts have a negative binomial distribution using a chi-square test statistic. We draw 1000 simulations from the negative binomial with the parameters estimated from the data (see the supplementary Rmd and html files in github at <https://pratheepaj.github.io/diffTop/>). We can then compute the p-value for the observed ASVs. The Supplementary Figure 1 shows how uniformly distributed the p-values are under the null hypothesis. After controlling for multiple testing, no ASVs reject the negative binomial distribution. For some ASVs (2.3% of total ASVs), the presence of zeros are larger than expected under negative binomial distribution. Some researchers have preferred zero-inflated negative binomial distribution (see (Holmes and Huber, 2018, Chapter 4) for a definition and formula) for such data Xu et al. (2015). Comparing the two models can be done using the Akaike Information Criterion (AIC) as in Romero et al. (2014a).

In the case of shotgun metagenomics data, it is also necessary to add taxonomic “bias” factors (McLaren, Willis and Callahan, 2019). Several of these unknown parameters have multiplicative effects, whereas the sequencing count data are also influenced by additive noise, such as contaminating DNA from reagents and the environment. Finally, we note that the sum of independent negative binomial random variables with the same parameters also follows a negative binomial distribution, so the library sizes are also expected to be negative binomial.

*4.2. Additive and Multiplicative Errors.* Specimens that are sequenced at much smaller library sizes will show many more zeros. Thus, zero-inflation is correlated with specimen library size and can be included in the relevant mixture models. In the zero-inflation case, the data can be split into the core taxa that present few zeros and whose abundances are used and a presence-absence table that encodes presence at a minimum number of reads and can take the rarer taxa into account.

Contamination of sequence-based data is modeled as an additive error (Davis et al., 2018; Salter et al., 2014). For instance, Salter et al. (2014) show reagents used in DNA extraction kits are heavily contaminated with microbial DNA, resulting in background noise that critically impacts results. Davis et al. (2018) proposed a simple statistical method called `decontam`, that identifies DNA contaminants in microbial studies using prevalences in designed experiments where negative controls have been included or by using frequencies if DNA concentrations of each specimen are included in the sample data Table 2. However, `decontam` may not identify specimen specific DNA contamination or rare taxa.

For contaminant removal, in the presence of negative controls, it has been shown that a statistical mixture model can estimate each taxon’s true intensity using reference priors and enable the removal of contaminant taxa (Cheng et al., 2019). The Bayesian hierarchical model for inferring DNA contamination in each specimen is as follows.

$$(4.1) \quad \begin{aligned} K_{ij} | \lambda_{ij}^{(r)}, \lambda_{ij}^{(c)}, d_j &\sim \text{Poisson} \left( \left( \lambda_{ij}^{(r)} + \lambda_{ij}^{(c)} \right) d_j \right), \\ \lambda_{ij}^{(r)} &\sim p \left( \lambda_{ij}^{(r)} \right) = \frac{|I(\lambda_{ij}^{(r)})|^{1/2}}{|I(0)|^{1/2}}, \quad \lambda_{ij}^{(c)} \sim \text{gamma} \left( \alpha_{ij}^{(c)}, \beta_{ij}^{(c)} \right), \end{aligned}$$

where  $\lambda_{ij}^{(r)}$  is the true intensity parameter,  $\lambda_{ij}^{(c)}$  is the contaminant intensity parameter,  $p(\lambda_{ij}^{(r)})$  is a marginal reference prior for the true intensities. We define it as a function of the Fisher information obtained through the marginal probability densities of  $\lambda_{ij}^{(r)}$ ,  $I(\cdot) = -\mathbb{E} \left[ \left( \frac{\partial^2}{\partial (\lambda_{ij}^{(r)})^2} \log p(k_{ij} | \lambda_{ij}^{(r)}, d_j) \right) \right]$ . We estimate hyper parameters  $\alpha_{ij}^{(c)}$  and  $\beta_{ij}^{(c)}$  using negative controls and find the reference for the library size using the median of ratios method (Anders and Huber, 2010) (see Section 4.3). We construct 95% highest posterior density (HPD) interval for the true intensity  $(L_{ij}^{(r)}, U_{ij}^{(r)})$  and the contaminant intensity  $(L_{ij}^{(c)}, U_{ij}^{(c)})$  for each taxon  $i$  in a specimen  $j$ . We declare a taxon to be contaminant taxon if the lower limit  $L_{ij}^{(r)}$  is smaller than the upper limit  $U_{ij}^{(c)}$ . For the 6,929 ASVs in specimens and controls, we applied BARBI (Cheng et al., 2019) and detected and removed 1,121 contaminants ASVs. We used the remaining 5,808 ASVs in 86 specimens for our downstream analysis.

*4.3. Library size scaling factor.* All downstream analyses need to account for library sizes (the column sums of the ASV contingency table); they are random multiplicative factors. In the context of RNA-seq, Anders and Huber (2010) propose a median-of-ratios algorithm that works well to estimate a library size scaling factor for each specimen. This method

divides each taxon count in Table 1 by the row’s geometric mean, and the library size scaling factor  $d_j$  is the median of the ratios for each specimen  $j$ . After the library size adjustment by  $d_j$ , an appropriate variance stabilization is applied and illustrated in Section 5.

Here we review some of the transformations we have found useful for the microbial count table. These transformations reduce the systematic variation in the microbiome count table and make the data approximately homoscedastic. Then we use dimensionality reduction methods to explore possible hierarchy factors or any batch effects.

## 5. Transformations.

5.1. *Variance stabilization.* Several parametric and nonparametric transformations were proposed for microbial count contingency tables in [McMurdie and Holmes \(2014\)](#). The variance of variance-stabilized transformed value is approximately independent of its mean value. Nonparametric regression methods are often used to characterize the mean-variance dependence of the library size normalized data. To do this, we can compute the mean and variance for each taxon. We then use a nonparametric regression method such as LOESS to model the relationship between mean and variance: the weights of each observation are adjusted using this fit, and an inverse transformation is applied that stabilizes the variance.

The transformations provided in `VOOM` ([Law et al., 2014](#)) scale the count table to count per million (CPM), then use the nonparametric fit for mean-variance dependence to compute the weight for each observation, take the log transformation of weighted observation to obtain variance-stabilized values. [Fukuyama et al. \(2017\)](#) shows that CPM transformations tend to produce false positives at the differential abundance analysis step. In our case, where the counts follow a negative binomial (NB) distribution, the [Anscombe \(1948\)](#) transformation provides an analytical solution. Given  $\mathbf{K}_j = [K_{1j}, K_{2j}, \dots, K_{mj}]^T$ , we let  $K_{ij}$  be a draw from a NB distribution with mean  $\mu_i$  and dispersion/exponent  $k_i$ . [Anscombe \(1948\)](#) shows that the optimal transformation for NB distribution is  $K_{ij}^*$ , where  $K_{ij}^* = \sinh^{-1} \left( \sqrt{\frac{K_{ij} + c}{k_i - 2c}} \right)$ . For  $k > 2$  and  $\mu_i$  large,  $c = \frac{3}{8}$  and  $\text{var} \left( K_{ij}^* \right) = \frac{1}{4} \psi' (k_i)$ , where  $\psi' (k_i) = \frac{1}{k_i - 1/2}$  for large  $k_i$ . All these types of transformations are now available in DESeq2 ([Love, Huber and Anders, 2014](#)). We propose [Anscombe \(1948\)](#) transformation to stabilize the variance (see Figure 2 in the Supplement, as compared to Figures 3, 4, and 5).

5.2. *Rank-based methods.* A robust approach to testing for treatment effects is to rank the taxa within each specimen from the most frequent to the least frequent. Most specimens do not contain representatives from all taxa, and noise level read counts could create large jumps in ranks at the lower end. Thus a threshold-rank approach is preferable, where the lower ranks are all assigned the same tied value of one. For instance, suppose there are 1,000 taxa present in the full data, and about one third occur at noise level in several specimens. The rank of noisy taxa could jump from 1 to 330 just by chance. Thus we assign scores from 670 for the most frequent in a specimen, down to a tied score of one for the last 330 taxa.

To choose the threshold 330 above, we performed a preliminary study on presence-absence patterns and choose the threshold for the number of reads as indicators that a taxon is present. Typically, three reads are sufficient. We then created a new table with a presence denoted by one and absence zero,  $B_{ij} = 1$ , if  $K_{ij} \geq 2$ . We can also represent the binary indicator  $B$  as a bipartite graph and rearrange the rows to reveal eventual block structures indicating communities stochastic block models can be convenient tools for such an approach ([Snijders and Nowicki, 1997](#)).



Principal component analysis (PCA) on truncated-rank transformed abundances produces a robust dimension reduction in the presence of a heavy-tailed distribution of count table. For the example data set, we choose ASVs in at least two specimens with at least 25 reads. This filtering results in 1418 ASVs in 86 specimens. Figure 2 shows a biplot resulting from the PCA after the truncated-rank transformation. *Actinobacteria* and *Betaproteobacteria* dominate in outlier *Centaurea solstitialis* in the positive direction of the first axis. The second axis explains the plant microbiome variability in non-Asteraceae specimens. The paired-specimens in Supplementary Table 1 are relatively close to each other.

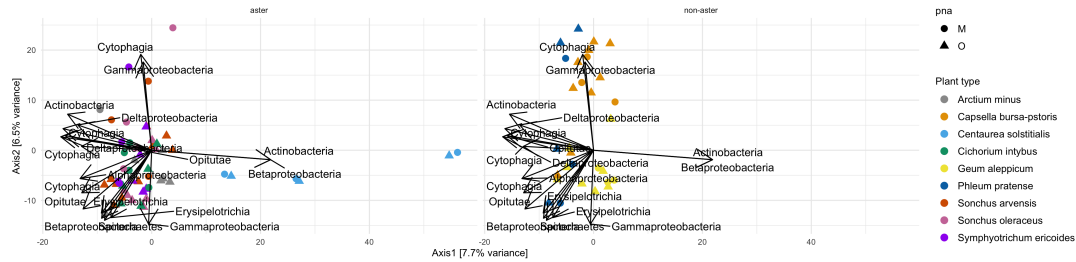


FIG 2. The biplot from the PCA after the truncated-rank transformation. The shape denotes universal (O) and Asteraceae-modified (M) pPNA types, respectively. Facet denotes Asteraceae or non-Asteraceae plants. *Actinobacteria* and *Betaproteobacteria* dominate in outliers *Centaurea solstitialis* in the positive direction of the first axis. The second axis explains the plant microbiome variability in non-Asteraceae specimens.

**6. Visualizations for heterogeneous data.** Microbiome data are high-dimensional and our example data have the four different components (count matrix, sample data, taxonomic table, phylogenetic tree) as elaborated in Section 2. Visualization methods identify taxonomic patterns or differences in various phenotypes or temporal variations in longitudinal experiments. These methods are also useful in getting insights from statistical inferences, such as differential abundance analysis. The `phyloseq` (McMurdie and Holmes, 2013) package incorporated wrappers for making interactive layered plots for microbiome data using the `ggplot2` (Wickham, 2016) package.

Traditionally, scientists prefer to use bar plots of the estimates of  $\mathbf{p}_j = (p_{i,j}, i = 1 \dots m)$ . More recently, heatmaps and interactive visualization of the phylogenetic tree jointly with taxonomic frequencies along longitudinal axes such as that provided by `treelapse` (Sankaran and Holmes, 2018, 2017) have become popular (see Kuntal and Mande (2019) for a review.) Networks and trees are useful aids to interpretation for microbial communities. Taxa co-occurrence graphs can serve as the basis for nonparametric tests inspired by the Friedman-Rafsky approach. One implementation is the `phyloseqGraphTest` package (Fukuyama (2020) illustrated in (Callahan et al., 2016a, Section)).

Simple statistical summaries show that in the example data set with 1418 ASVs in 86 samples, *Proteobacteria* is the most prevalent Phylum. Next largest prevalent Phyla are *Actinobacteria*, *Bacteroidetes*, and *Firmicutes*. We computed the relative abundance of ASVs in each specimen. Then, we removed Class with less than five ASVs. Supplementary Figures 6 and 7 show the distribution of the relative abundance of each Phylum in specimens faceted by Class. We can use the bar plots to compare the difference in scale and distribution of Phylum in both pPNA types. These figures show unimodal abundance profiles in four different

Classes of *Proteobacteria*. There are few replicates sequenced with M-pPNA than O-pPNA in non-Asteraceae.

The heatmap is constructed using transformed data to maximize the contrasts and understand how the specimens cluster and batch effects. In the example data set, Figure 3 shows that the ten most prevalent Class taxonomy is similar in Asteraceae and non-Asteraceae plants, except in *Centaurea solstitialis* plants, which have the most abundant ASVs of the Class *Actinobacteria*, *Alphaproteobacteria*, *Sphingobacteria*, and *Gammaproteobacteria*. The most prevalent Class *Actinobacteria* is more abundant in Asteraceae than it is non-Asteraceae plants.

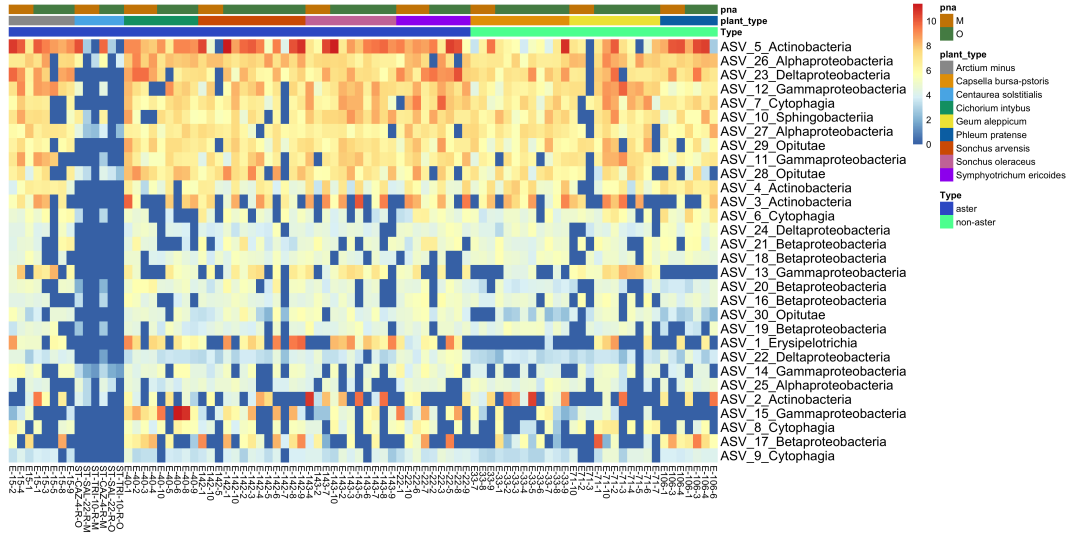


FIG 3. Thirty most abundant ASVs were selected in all specimens. Taxa are labeled by Class on rows, and specimens are on the columns of the heatmap. Some specimens from *Centaurea solstitialis* plant have the most abundant ASVs of the Class *Actinobacteria*, *Alphaproteobacteria*, *Sphingobacteria*, and *Gammaproteobacteria*.

## 7. Multivariate and Network Analyses.

7.1. *Ordination.* To identify outliers, clusters, and relative position or gradients of specimens in low dimensions, ordination methods can be used on any of the many distances available for measuring similarities using abundances or presence-absence. Historically, ecologists have been leaders in their careful choice of distances between specimens. The most popular include the Bray-Curtis, chi-square, Wasserstein, or Jaccard distances. Jaccard, Bray-Curtis, or unifracs distances are all popular choices in microbiome studies. The distance matrices are then used to create sample maps through ordination methods such as multidimensional scaling (MDS) (also known as principal coordinate analysis (PCoA)), double principal coordinate analysis (DPCoA), or nonmetric multidimensional scaling (NMDS). The resulting two or three dimensional maps can indicate clusters of samples when there is an underlying latent categorical variable (McMurdie and Holmes, 2013) or gradients — continuous latent variables such as water depth seen in the TARA ocean data (Nguyen and Holmes, 2017). Clusters can lead to a simplification of the data by assigning specimens a state type, as is done for the vaginal microbiome studies (Romero et al., 2014b; DiGiulio et al., 2015). However, sometimes clusters appear through artifacts, such as low-density sampling of the high dimensional space (Gorvitovskaia, Holmes and Huse, 2016). Nguyen and Holmes (2019)

provide a set of guidelines to use and interpret dimensionality reduction methods for different data types.

Supplementary Figure 8 shows an MDS plot built using weighted unifracs distances between all our specimens. Paired specimens are labeled, and we can detect some of them form clusters, except paired-specimens (E-143-7, E143-7), (E-142-5, E142-5), (E-71-2, E71-2), and (E-71-3, E71-3). *Centaurea solstitialis* specimens are outliers among Asteraceae plants in the positive direction of Axis 2, which were sampled in a field across three countries. Axis 1 explains the microbial variability in plant types. Supplementary Figure 9 shows biplots built using the double principal coordinate method (DPCoA) that also incorporates the phylogenetic information into the ordination. The labels in Supplementary Figure 9 (A) depict paired specimens with both pPNA types; these form clusters, except paired-specimens (E-143-7, E143-7), (E-142-5, E142-5), and (E-71-3, E71-3) that are in positive and negative axes. Axis 1 explains the microbial variability in all paired specimens and specimens from *Sonchus oleraceus* and *Sonchus arvensis* plants with highly abundant *Erysipelotrichia*. We scale the axis of ordination plots according to the eigenvalues to represent relative distances between specimens as faithfully as possible. It seems that the endosphere microbial variability comes from differences in plant types and specimens, much less from the amplification methods. Supplementary Figure 9 (B) suggests that specimens amplified with both types of pPNAs are composed of ASVs from different Classes. To enhance the visualizations of microbial count data, we can incorporate some of the phylogenetic information into the ordination summaries as described in the next section.

*7.2. Integrating the phylogenetic tree into the analyses.* When considering the abundance of the different bacteria associated with the denoised ASVs identified through pipelines such as DADA2 (Callahan et al., 2016b), strain variants are identified using standardized bacterial taxonomic databases (Cole et al., 2008; Pruesse et al., 2007) and phylogenies such as greengenes (DeSantis et al., 2006).

These phylogenetic relationships create a family tree, of which the tips are the different taxa or strains (ASVs are also called Operational Taxonomic Units, OTUs). Thus, the count table row identifiers in Table 1 are not evolutionarily independent, and there can be some benefit to taking this into account. Using the phylogenetic tree to inform distances or kernels between the abundance vectors was developed in Purdom (2011), extending the idea of doing a double principal coordinate analyses for ecological data as presented in Pavoine, Dufour and Chessel (2004).

One difficulty in tree-based analyses is that the phylogenetic relationships between taxa may only explain a small percentage of the abundance differences between specimens (or beta-diversity, as it is called). More recent work has shown that a modulated penalized-tree approach allows a more nuanced use of the phylogenetic tree to interpret sample differences (Fukuyama, 2019; Fukuyama et al., 2017). Testing procedures can help delineate what the tree can explain and how the residuals from the tree-based variability relate to the other specimen covariates. Figure 4 shows the results from adaptive generalized PCA (adaptive gPCA), available in `adaptiveGPCA` package. It reveals *Centaurea solstitialis* specimens are outliers among Asteraceae plants on the left of Axis 1. Axis 2 explains the microbial variability in plant types. Like truncated-rank PCA, the paired-specimens in Supplementary Table 1 are relatively close to each other.

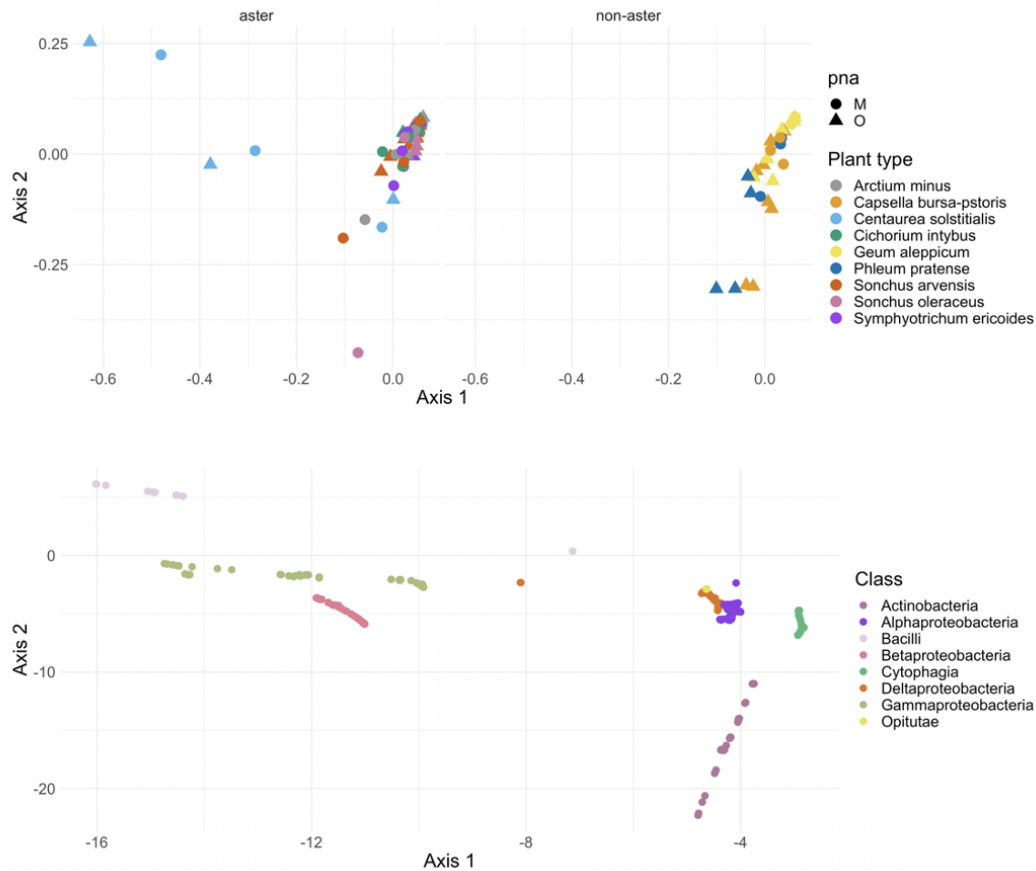


FIG 4. The results from adaptive gPCA reveal *Centaurea solstitialis* specimens are outliers among Asteraceae plants in the left of Axis 1. Axis 2 explains the microbial variability in plant types.

When testing for differential abundances between taxa, the unifrac score—a phylogenetic-based distance (Lozupone and Knight, 2005) or modifications thereof enable tree-based tests (for a review and comparisons between different tree-based distances see Fukuyama et al. (2012)). There have been several follow-up methods that use the development of phylogeny-based kernels to understand microbial diversity (Zhao et al., 2015; Washburne et al., 2019). There are still large areas that require additional research when it comes to propagating the uncertainties with which the phylogenies are known and the uncertainties with which the taxa are identified or the number of reads measured.

Several attempts have been made to leverage hierarchical Dirichlet processes to link the evolutionary processes at work with the ecological context. In particular, Harris et al. (2015) consider the result of Hubbell showing that under the neutrality assumption, the abundances within the neutral guild fluctuates and that the number of species at a single site is a balance between the immigration of new species and local extinctions.

**7.3. Correspondence analysis.** Correspondence analysis (CA) is a weighted bilinear method that provides a representation of a contingency table in a low-dimensional space (Greenacre, 2010b). CA can be understood as a generalized singular value decomposition that provides factor scores for columns and rows of the contingency table that are then used to represent the association between rows (taxa) and columns (specimens) (Holmes, 2008).

CA uses chi-square distances, which are sensitive to outliers. The scree plot shows that the dimensionality of the underlying variation is two dimensional, and we see the first two dimensions explain about 21.7% of the inertia (proportional to the sum of the chi-square distances). Figure 5 shows that CA is not robust to the outlier *Centaurea solstitialis*. We observe five clusters of specimens apparent in the data. *Bacilli* contributes more to the cluster of one paired *Centaurea solstitialis* in the first axis, *Sphingobacteriia*, *Gammaproteobacteria*, *Betaproteobacteria*, and *Actinobacteria* dominate in a cluster of another paired *Centaurea solstitialis*, and other ASVs comprise the other three clusters of all other specimens.

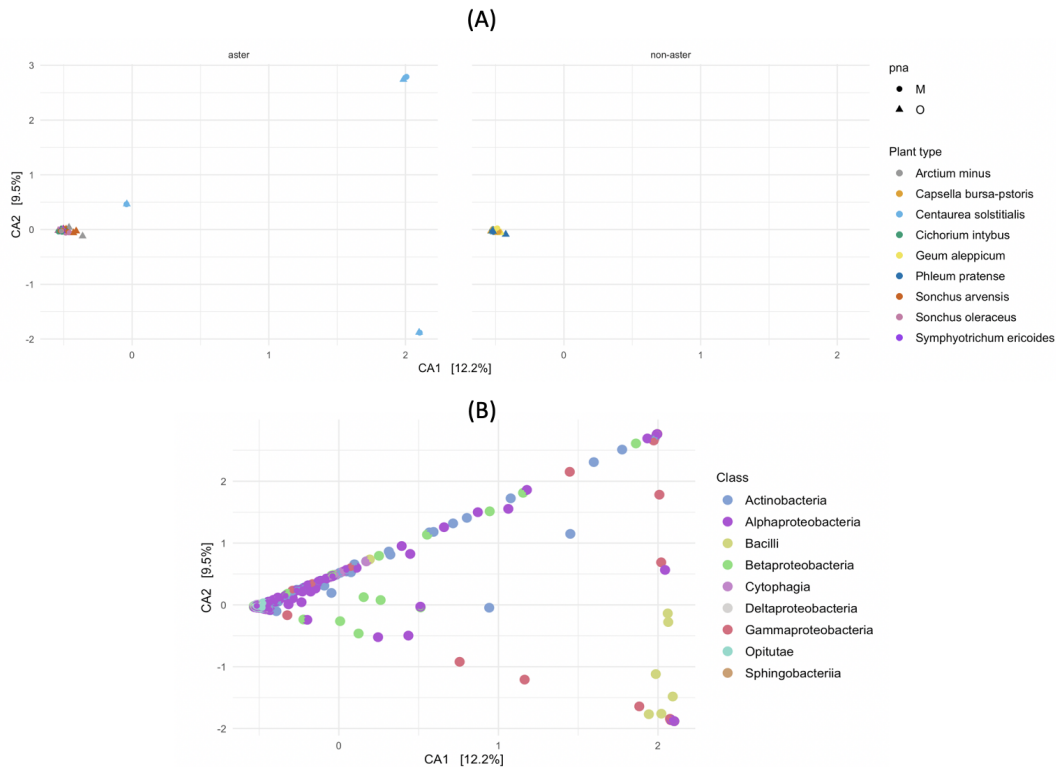


FIG 5. Correspondence analysis of all plant specimens. (A) plot shows the points of specimens and (B) plot shows the points of ASVs and the corresponding Class labels.

We see that several combinations of distances and multivariate methods provide similar conclusions at “larger levels.”

**7.4. Network analysis.** We subset the 18 paired-specimens to perform a graph-based test. Supplementary Figure 10 shows the network based on fixing a maximum threshold for the Jaccard dissimilarity matrix. All paired-specimens are connected, except (E-143-7, E143-7), (E-142-5, E142-5), (E-71-2, E71-2), and (E-71-3, E71-3). *Centaurea solstitialis* paired-specimens, and all other paired-specimens make subgraphs. We performed a test using the minimum spanning tree (MST) built with the Jaccard dissimilarity (without thresholding). The null hypothesis is that the two pPNA types (universal and Asteraceae-modified) have the same microbial distribution. This test has a larger p-value (0.872); thus, we do not reject the null hypothesis. This result mirrors the observations from the PCA after the truncated-rank transformation, adaptive gPCA, MDS, and DPCoA that paired-specimens in Supplementary Table 1 have similar microbial variability.

**8. Differential Abundance Analysis.** An important goal in microbiome research is often to find taxonomic differences across environments or groups. Although strain switching may obscure the interpretation of these differences, if the strains stay the same, an adaptation of differential abundance techniques used in RNA-seq has proved useful. After visualizing the data, several possible downstream analyses are commonly used depending on the experimental design and questions of interest. For instance, consider a study that presents several arms, such as treatments and control. The most straightforward approach is to rank taxa according to how differentially abundant or differentially prevalent they are in the two conditions. Then, we follow a standard decomposition of variability according to random effects or fixed effects. In human studies, the random effects could correspond to subjects and fixed effects to treatments.

Differential abundance analysis can be done on count data (McMurdie and Holmes, 2014; Love, Huber and Anders, 2014; Robinson, McCarthy and Smyth, 2010) or transformed data (Smyth, 2005). In the latter case, it can be necessary to add a pseudo count for zeros. It is preferable to do goodness of fit test to model counts for each taxon directly.

8.1. *Permutation tests using distance matrices.* A popular confirmatory analysis starts with the computation of dissimilarities between the samples using either the Jaccard, Bray Curtis, unweighted or weighted unifrac<sup>1</sup> distances. Then, under the null assumption that the sample abundance vectors are independent and identically distributed across groups, a null permutation distribution of a statistic dependent on the distances can be compared to the statistic calculated on the observed data. Examples of such procedures include the permutational multivariate analysis of variance "PERMANOVA", introduced in Anderson (2001) and available as the `adonis` function in the `vegan` package for instance Oksanen et al. (2020). However the method is dependent on the assumption that the samples are (conditionally) independent and permutations of the full set give false positives if the samples have nested a priori factors, in which case modified permutation procedures such as those suggested in Excoffier, Smouse and Quattro (1992) are more appropriate.

Some asymptotic theory is available for PERMANOVA in the most restrictive cases of weighted Euclidean distances between independent samples Anderson and Robinson (2003). These tests are sensitive to latent groupings and correlations between ASVs leading to common occurrences of false positives and over-interpretation of the significance of the differences between groups. Distance based tests agglomerate the different ASVs into one dissimilarity index so they do not provide indications as to which ASVs differ. One difficulty with current 16S rRNA data are that the distance based tests are very sensitive to the *choice* of distance and the presence of *strain switching* can substantially decrease the power of the test. As an example, in the article — from which we drew the data Fitzpatrick et al. (2018)— the authors use PERMANOVA and report that they do not detect a difference between specimens grouped according to pPNA types (universal and Asteraceae-modified pPNA types). In this case, one explanation is that these types of permutation tests are particularly underpowered when one set of samples have one strain (ASV) and another a slightly different strain, registered as a different ASV. In the supplementary material, we show that for the exemplary data this is in fact the case as strains ASV 153 is switched with ASVs 12, 354 and 345. To illustrate the decrease in power through a simulation, we generated negative binomial count data with parameters similar to these ASVs and show that when a species switches ASV, ie is present as one ASV in one set of samples and a close, distinct strain appears in the other set of samples, the power to

---

<sup>1</sup>the unifrac distance is a modification of the Wasserstein distance computed along the phylogenetic tree Fukuyama et al. (2012); Lozupone and Knight (2005); Evans and Matsen (2012).

detect a difference using a Bray Curtis distance and PERMANOVA is considerably diminished (see the code at [https://pratheepaj.github.io/diffTop/articles/appendix/08\\_differential\\_abundance\\_analysis.html](https://pratheepaj.github.io/diffTop/articles/appendix/08_differential_abundance_analysis.html) for the illustrative power calculations).

### 8.2. *Differential abundance through generalized linear modeling and transformations.*

The bacterial abundances vary between subjects and environments both because the underlying prevalences are different and because the sampling depths vary. The sampling depth variation causes the count data to have unequal variances (heteroscedasticity). Mixture/hierarchical models are useful for this type of data. Using a Poisson-gamma hierarchy achieves a first model that results in negative binomial [McMurdie and Holmes \(2014\)](#) count data for which the generalized linear models for differential abundance analysis of microbiome data ([Love, Huber and Anders, 2014](#); [Robinson, McCarthy and Smyth, 2010](#)) is well adapted. After an appropriate transformation on the count as a response, these generalized linear models can detect the differences in bacterial abundances ([Smyth, 2005](#)). Testing can be done at different levels in the phylogenetic tree and adjustment for nested testing is available through packages such as `structSSI` that implements a multiple testing procedure that accounts for hierarchical dependence in the taxonomy table or phylogenetic tree ([Sankaran and Holmes, 2014](#)).

Environmental differences, weather change, animal or human interaction can create substantial heterogeneity in bacterial communities within and between subjects or locations. Therefore, longitudinal experiments are often preferred because they account for the within and between-subject/location variability. We have developed a moving block bootstrap method for differential abundance analysis in longitudinal studies ([Jeganathan et al., 2018](#)) that accounts for the added dependences. This method resamples overlapping blocks of specimens within each subject to approximate the test statistic's distribution and tailors the pivoting procedure and block sizes to the data.

Some microbiome studies incorporate spatial dependence into the differential abundance analysis along spatial gradients ([Proctor and Relman, 2017](#); [Proctor et al., 2018](#)). [Singh et al. \(2019\)](#) proposed a nonparametric test to identify the effect of environmental factors that shape microbiome variability. This test can account for spatial dependence and inter-dependence among taxa.

For designed experiments, [Grantham et al. \(2020\)](#) developed a Bayesian mixed-effects model for testing the effect of environmental and treatment factors on microbiome variability. This method models the taxonomic counts as a multinomial and accounts for the interdependence between taxa by applying a hierarchical mixed-effects model.

Unfortunately, simple two-sample testing is marred by several technological difficulties. There are batch effects, technical biases, and heteroscedasticity in the prevalence estimates due to differences in the library sizes across specimens as well as strain switching (where ASV strains are replaced as we change locations or subjects). The taxonomic strains are often not pre-specified, requiring a nonparametric infinite-dimensional model and precluding the use of methods that assume a small well-defined set of categories — the case of compositional data methods — for instance. In general, we recommend hierarchical models that can account for undetected taxa and heterogeneity of microbial distribution in specimens.

### 8.3. *Communities instead of individual taxa.*

As noted above, high-resolution acquisition of data at the taxonomic strain level resulted in different subjects or environments presenting slightly different taxa that are “functionally synonymous”; thus we have to deal with these *strain switches*. The individual taxa may not be as important as their combination. This makes the problem similar to how synonyms occur in textual analyses. The co-occurrence

of bacteria and the departures from a simple multinomial model make the analogy between textual and microbiome data analysis useful. [Sankaran and Holmes \(2019b\)](#) demonstrated the utility of the analogy between textual analysis and molecular microbial ecology. In textual analyses, documents are of unequal lengths, and topics enable the simplification of document contents. Documents can have several topics; thus, mixed membership is the relevant model. The same is true in microbial ecology, where several communities can be present in a specimen. A probability distribution over bacteria characterizes each community. This has the advantage over the simpler Dirichlet multinomial mixture models [Holmes, Harris and Quince \(2012\)](#) that bacteria can be very interdependent. Syntrophy occurs when one bacterium cannot survive without another and can be modeled by having a simple two bacteria topic, whereas syntrophy is impossible to model with a simple multinomial. In the next section, we describe topic models and illustrate their use on our example data set. In particular, topic models provide useful aggregates that can be used for differential abundance analysis based on topics rather than individual strains.

**9. Latent Dirichlet Allocation.** We use the latent Dirichlet allocation (LDA) model as described by [Sankaran and Holmes \(2019b\)](#). We suppose that  $\mathbf{K}_{\cdot j}$  denotes the  $m \times 1$  vector of taxa abundance in specimen  $j$ , where  $j = 1, 2, \dots, n_1$  and  $T$  is a prespecified number of topics. The LDA generative process for each specimen  $\mathbf{K}_{\cdot j}$  of size  $S_j$  follows the following steps: Each specimen is associated to a set of topics drawn from a probability distribution  $\theta_j \stackrel{\text{iid}}{\sim} \text{Dirichlet}_T(\alpha)$ , over a mixture of latent topics (bacterial communities). Each topic is associated with a probability distribution  $\beta_t \stackrel{\text{iid}}{\sim} \text{Dirichlet}_m(\gamma)$  over taxa. We draw  $\mathbf{K}_{\cdot j}$  from  $\mathbf{K}_{\cdot j} | S_j, \theta_j, \mathbf{B} \stackrel{\text{iid}}{\sim} \text{Multinomial}(S_j, \mathbf{B}\theta_j)$ , where  $j = 1, \dots, n_1$ ,  $\mathbf{B} = [\beta_1, \dots, \beta_T]^T$ , and  $S_j$  is the library size of specimen  $j$ .

We set the hyper-parameters  $\alpha$  and  $\gamma$  less than one to generate sparse mixtures that are different from each other and avoid generating unrealistic topics. We estimated the model parameters using Hamiltonian Monte Carlo and the No-U-Turn (HMC-NUTS) method implemented in Stan ([Carpenter et al., 2017](#)). We denote the Bayesian posterior estimates  $\{\hat{\theta}_j, \hat{\beta}_t\}$ ,  $j = 1, 2, \dots, n_1$  and  $t = 1, 2, \dots, T$ .

For our example, we chose the number of topics  $T = 11$  based on an ordination analysis. For the data set in consideration, there were 1418 distinct ASVs across 86 specimens after selecting ASVs with at least 25 reads in at least two specimens. Ordination plots show that bacterial signatures vary between Asteraceae and non-Asteraceae plants, *Centaurea solstitialis* specimens are outliers among Asteraceae plants, and the paired specimens are relatively similar in microbial variation. For eight different plants and *Centaurea solstitialis* from three other countries, we chose the number of topics  $T = 11$ . We set  $\alpha$  to be 0.8 across all 86 specimens and  $\gamma$  to be 0.5 across all ASVs. We estimated the parameters using the HMC-NUTS sampler with four chains and 2000 iterations. Out of these 2000 iterations, 1000 iterations were used as warmup samples and discarded. Label switching across chains makes it difficult to directly compute log predictive density, split- $\hat{R}$  ([Vehtari et al., 4 July 2020](#); [Gelman and Rubin, 1992](#)), effective sample size ([Kruschke, 2014](#)) for model assessment, and evaluate convergence and mixing of chains. To address this issue, we fixed the order of topics in chain one and then found the permutation that best aligned the topics across all four chains. For each chain two to four, we identified the estimated topics pair with the highest correlation, then found the next highest pair among the remaining, and so forth.

We present the predictive model check for  $T = 11$ , with simulated and observed data using a statistic  $G(K_{ij}) = \max\{K_{ij}\}$ . Supplementary Figure 13 shows the histograms of  $G(K_{ij})$  of each ASV in simulated data from the fitted model and the horizontal line of the observed  $G(K_{ij})$ . According to the histograms, the LDA model with eleven topics makes a



realistic prediction. Supplementary Figures 14 and 15 show the effective sample size (ESS) and split  $\hat{R}$  with eleven topics. These diagnostics provide some evidence of good mixing and convergence of the chains to the stationary distribution.

The host contamination plastid is significantly reduced in Asteraceae plants with M-pPNA type (Fitzpatrick et al., 2018). But only four plants *Sonchus arvensis*, *Arctium minus*, *Sonchus oleraceus*, *Centaurea solstitialis* of Asteraceae have paired-specimens. Now, we infer whether the pPNA types affect microbial distributions in different plants. In endosphere specimens, bacteria that have similar functionalities in each plant type can co-occur as latent communities, and topic modeling provides insight into microbial communities across these types. With the goodness-of-fit of the LDA model, we can draw an informative summary of bacterial communities in each plant type with O and M-pPNA types. Hence after the topic analysis, we choose to study one plant type at a time. The analogous figures for the other plant types are available as Supplementary Figures 11 and 12.

Figure 6 shows the topic distribution in each specimen from the *Sonchus Arvensis* plant. Among the eleven topics for all plants, seven topics were predominately present in specimens from *Sonchus Arvensis*, an Asteraceae plant with 13 replicates specimens. Ten specimens were amplified with O and three with M-pPNA types, respectively. Among these specimens, E-142-1, E-142-5, and E-142-10 are paired because both pPNAs were used for the same DNA specimens. Paired-specimens (E-142-1, E142-1) and (E-142-5, E142-5) have a different mixture of topics. Figure 7 shows the ASVs distribution for each topic. Paired-specimen E-142-1 has largest proportion of Topic 6, which has distribution over ASVs from Classes *Acitinobacteria*, *Betaproteobacteria*, and *Eryscpelotrichia* whereas other paired-specimen E142-1 has largest proportion of Topic 1, which has similar ASV distribution, except of Class *Eryscpelotrichia*. Similarly, some paired-specimens from *Arctium minus* and *Sonchus oleraceus* plants have different mixtures of topics (Supplementary Figure 12). The replicates of *Centaurea solstitiali* that sampled across three countries show variation in microbial distribution, but the paired-specimens have similar distributions (Supplementary Figure 11).

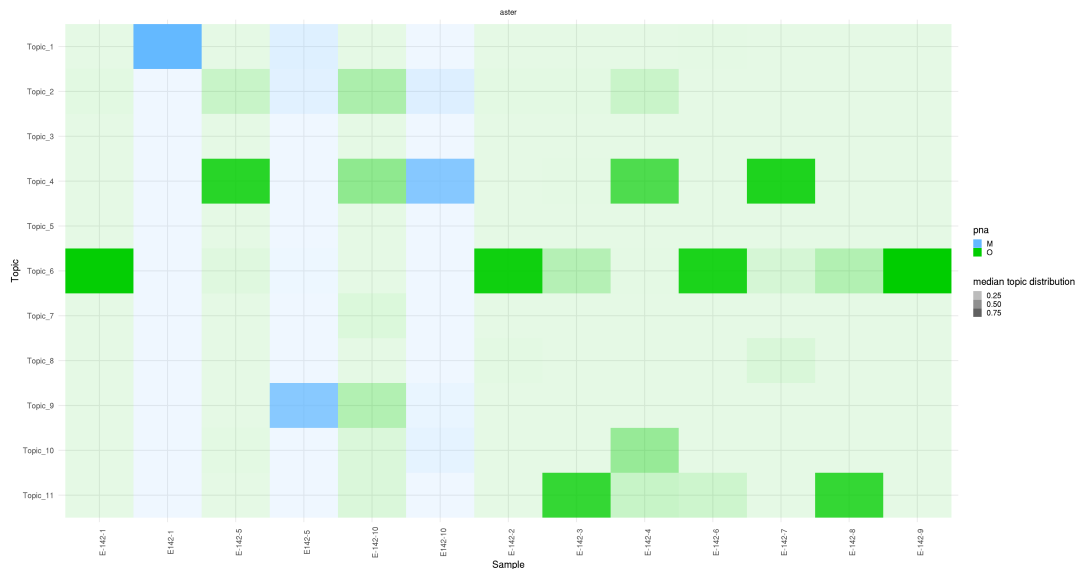


FIG 6. Topic distribution in specimens from *Sonchus Arvensis* plant. The color gradient represents the median topic distribution.

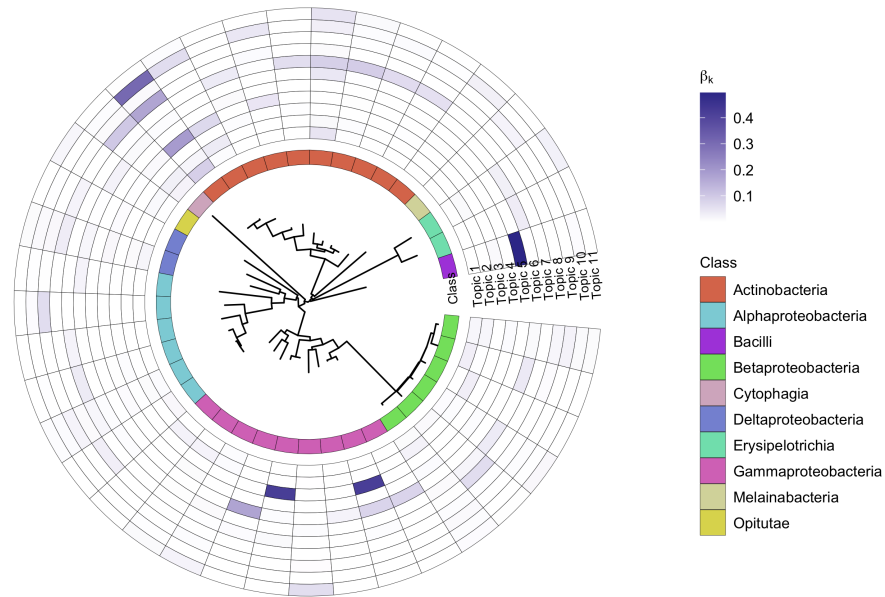


FIG 7. ASV distribution over topics in all specimens.

We use the microbial community-topics in each specimen to do a differential topic analysis to test whether topic memberships differ across conditions. First, we compute the median Bayesian posterior of topic proportions in each specimen. Then, we multiply the proportions by the library size and round to an integer. This will give an abundance of the topic in each specimen. Then we applied `DESeq2` to identify differentially abundant topics across O and M pPNA types. Supplementary Table 2 shows that Topics 4, 5 and 10 are differentially abundant in O and M pPNA types. We also consider testing on only paired-specimens from Asteraceae and non-Asteraceae plants. Supplementary Figure 12 shows the large proportion of Topics 1, 9, and 11 in Asteraceae paired-specimens. Supplementary Table 3 shows that these microbial communities are different in pPNA types. Topics 1, 2, 8, 9, 10 dominates in non-Asteraceae paired-specimens, and Supplementary Table 4 shows that we do not reject the hypothesis that these topics are differentially abundant. We conclude that microbial communities are significantly different in some Asteraceae plants.

Finally, we note that the choice of the number of topics is possible through a Bayesian non-parametric approach that incorporates an infinite number of topics as described by [Blei, Carin and Dunson \(2010\)](#). Sources of biological variability can be incorporated into the hierarchical model; subject or location variation can be modeled as random effects. Other covariates, experimental arm variability (for instance, one group of subjects/locations may be treated) can be easily added. In the study of the human microbiome, it is common to include a subject's age or the level of urbanization in the food supply at a given location ([Yatsunencko et al., 2012](#)).

An extension to the topic models described above is provided by a Bayesian non-parametric factor models that accomodates changing distributions of unbounded numbers of taxa in specimens. These models decompose the biological variation in bacterial communities'

composition in a low-dimensional space defined along *continuous* latent factors rather than through *discrete* topics (Ren et al., 2017). This method enables the propagation of uncertainty from the microbiome data to their multivariate ordination, and we illustrate this on our example data set below.

**10. Uncertainty quantification for ordination analysis.** Ren et al. (2017)’s nonparametric Bayesian approach provides one way of modeling infinitely many possible taxa and specimen-specific microbial distributions. Their approach assumes an underlying finite-dimensional factor model that represents dependencies and uses a kernel decomposition of the underlying communities. If we let unknown specimen-specific microbial distributions be  $P^j\{Z_i\}, j = 1, \dots, n_1$  and  $Z_i$  be the  $i$ -th taxon of unbounded set, then Gram matrix  $\phi(j_1, j_2)_{j_1, j_2 \in \{1, \dots, n_1\}}$  defines the similarity between microbial distributions in specimens  $j_1$  and  $j_2$ . Ren et al. (2017) represent in the prior the similarity between  $P^j$  through latent factor model in low dimensional space  $\mathbf{Y}^j$  and use the posterior samples of normalized Gram matrix for the ordination analyses. This Bayesian approach has the advantage of providing posterior probability credible sets.

We reused this technique on our example data set, where we choose ASVs that occur in at least eight specimens with at least 100 reads. This filtering resulted in 152 ASVs in 86 specimens. We start with 86 latent variables that assume one factor for each specimen. Then, we used Gibbs sampling with 50,000 iterations and a thinning size ten. Figure 8 shows Bayesian nonparametric ordination for all the specimens. The contours represent the variability in the position of each specimen in the consensus space. Since we filtered the data down to very consistent ASVs, the credible regions for all specimens are relatively small. The two paired specimens from *Centaurea Solstitialis* are outliers that corroborate the exploratory analysis results with credible regions.

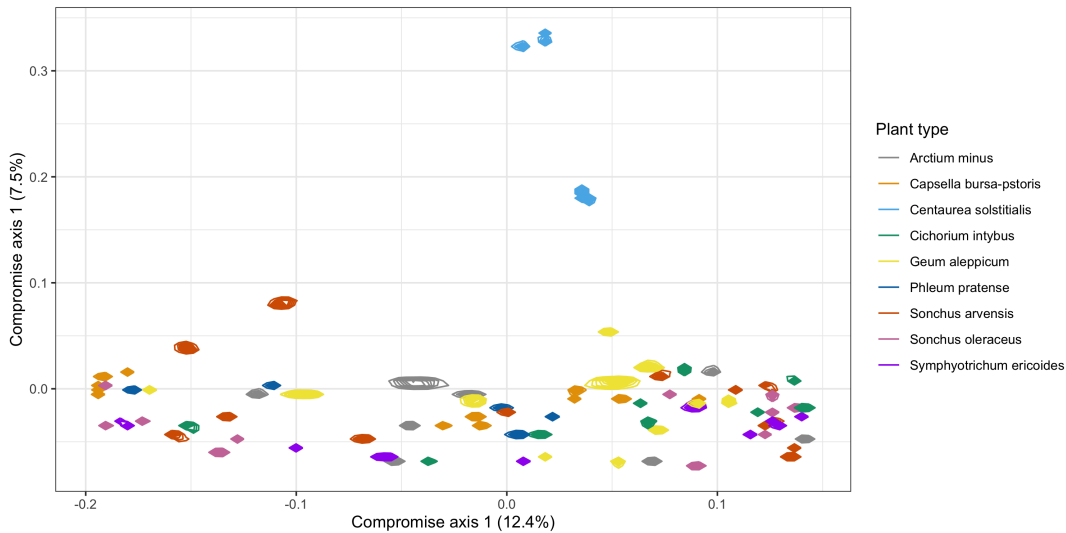


FIG 8. Ordination plot of specimens and 95% posterior credible regions. We plot the first two consensus axes. The percentages on the two axes are the ratios of the corresponding eigenvalues and the trace of the matrix. Color indicate the plant type.

**11. Discussion.** We have reused the 16S rRNA sequencing data from Fitzpatrick et al. (2018) to review and demonstrate goodness of fit tests, multiplicative and additive models, variance-stabilizing and truncated-rank transformations, multivariate methods, network

analysis, differential topic analysis, and uncertainty quantification for ordination analysis of sequencing reads in microbial ecology. The study used the 16S rRNA gene sequencing with universal (O) and Asteraceae-modified (M) plastid peptide nucleic acids (pPNA) types that limit the amplification of host-derived DNA, such as plastids, which share existing similar bacterial lineages. Fitzpatrick et al. (2018) designed the experiment to test the efficacy of pPNA types in limiting the plastid contamination and seeing if there was a difference in microbial variation between pPNA types. By applying a linear mixed-effects model for each taxon, Fitzpatrick et al. (2018) concluded that the pPNA types do not affect the identification of individual bacterial taxa and within specimen  $\alpha$ -diversity for different metrics. The authors found little evidence of a difference between specimen  $\beta$ -diversity using permutational multivariate analysis of variance (PERMANOVA) (Anderson, 2005). The original paper used a paired-designed experiment for some replicates of each plant. In this case, it would have been preferable to tailor the standard permutation method to the experimental design and incorporate ASVs that seemed to have switched labels. In general, strain switching precludes simple differential abundance analysis at the taxonomic level and we have introduced a more meaningful differential topic analysis that identifies some topics in the Asteraceae paired-specimens that are significantly different.

Bar plots and simple multinomial models have limitations in providing a complete picture of microbial variability, especially when syntropic relations create dependence between bacterial occurrences. The phylogenetic relationships between taxa can be important for biological interpretation, and interactive tools can enhance our understanding of the distribution of different taxonomies (McMurdie and Holmes, 2015; Sankaran and Holmes, 2018). If we use sorted prevalence to visualize each taxonomy, we may lose specific rare bacteria present in the data, so it is important to do both analyses with presence absence data and that of the core microbiome with abundances. Heatmaps can identify clusters of specimens and the contribution of ASVs for each cluster, although a large number of ASVs can make this difficult to do in one step.

Our review shows that ordination based on different distances can incorporate phylogenetic information, or incorporate sample depth weights as in correspondence analysis. For the example at hand, multidimensional scaling (MDS, PCoA) and correspondence analysis (CA) show that only a few paired-specimens from the Asteraceae plants and non-Asteraceae form clusters.

Among dimensional reduction methods, tree based metric ordination such as double principal coordinate analysis (DPCoA) is more interpretable and here it uncovers the dominant taxa in each cluster. These ordination plots depict only the most abundant or prevalent ASVs in clusters of specimens. As an alternative, topic models are shown to be the more interpretable as they unmask rare and synonymous taxa and can enhance our understanding of the assembly of microbial communities as topics which can be projected onto the phylogenetic tree. The ASV distributions over interpretable topics demonstrate that differential topic analysis can enhance our understanding of the differences in complex microbiome communities (here across different plants and pPNA types).

In our reanalysis of the data, there were six Asteraceae plants and three non-Asteraceae plants with 18 paired-specimens, as in Section 3. Compared to O-pPNA, the M-pPNA significantly reduced the host-derived DNA contamination in only Asteraceae plants. The topics we uncovered in the latent Dirichlet allocation (LDA) explained the differences in the mixture of topics in some paired-specimens of Asteraceae plants. Nevertheless, there were not enough paired-specimens to test the difference in microbial distribution. In addition, we found that eleven topics were distributed over mixture of ASVs from Classes *Acetivibacteria*, *Alphaproteobacteria*, *Betaproteobacteria*, *Deltaproteobacteria*, *Gammaproteobacteria*, *Bacilli*, *Cytophagia*, *Erysipelotricha*, and *Opitutae*. *Acetivibacteria* was not in Topics 3, 5, and 9, and

two paired-replicates of *Centaurea solstitialis* has the largest proportion of Topics 3 and 5, which would be the reason for identifying *Centaurea solstitialis* specimens as outliers in ordination analyses. In contrast to results on Asteraceae plants, [Fitzpatrick et al. \(2018\)](#) concluded that M-pPNA type did not limit deriving host-contamination in non-Asteraceae plants. We found fewer of topics in non-Asteraceae than Asteraceae plants as in Supplementary Figure 12. Supplementary Table 4 shows that there is no evidence to have different mixtures of topics in O and M-pPNA types in non-Asteraceae plants.

Latent microbial communities (topics) can share ASVs that are correlated. We could have enhanced our analyses by using correlated topic models ([Blei and Lafferty, 2006](#)) that can be used if some topics are expected to be exclusive or negatively correlated. LDA does not provide the number of topics or account for covariates; for this, we recommend using hierarchical Dirichlet process topic modeling with a stick-breaking process for the base measure that can account for covariates and enable the inference of the number of topics from the data.

Finally, we show how credible regions for ordination can help in the confirmatory phase of evaluating the relative similarity of specimens. Here, the Bayesian nonparametric ordination corroborated the LDA results on specimens showing dominant topics in most replicate specimens.

**12. Open problems and statistical challenges.** Microbiome data from designed experiments pose numerous statistical challenges because they are high-dimensional, heteroscedastic, and sparse. These count data are not normally distributed, and linear models that assume normality are not appropriate. Some transformations such as those in `voom-limma` ([Smyth, 2005](#)) have been proposed; these lead to the use of raw counts to count per million (CPM). Multiplying by large factors is tempting; however, the risk of invalidating the downstream statistical inferences is high. For instance, in the linear discriminant analysis effect size (LEfSe) method ([Segata et al., 2011](#)) multiply relative abundances by a million; this creates artificially large sample sizes and bloats the power of the experiment resulting in a large false-positive rate.

Our exemplary data were cross-sectional with an ideal paired-sample design. Complete randomized design, randomized complete block design, split-plot design, longitudinal experimental designs, spatial and factorial designed experiments require much more intricate hierarchical mixture models. Because of the strong between location/subject variability, longitudinal designs are currently the most used in microbiome studies. The transformation and visualization methods that we discussed in this paper can be used for any exploratory analysis. For the differential abundance analyses, hierarchical generalized linear mixed models or moving block bootstrap ([Jeganathan et al., 2018](#)) have proved successful in the cases where the strains stay consistent across samples, but many more extensions need to be developed to incorporate more complex structured designs.

The tools we have shown here for analyzing marker-gene counts use standard statistical techniques such as nonparametric or parametric variance stabilization through Anscombe's transformation, generalized log transformation, and smoothing techniques. However, research is still needed on how to invert these transformations after the latent factors in the data have been discovered. Open questions include propagation and quantifying uncertainty in the phylogenetic relationships, assessment of the effect of technical biases, and count estimations on the downstream analytics, even in well defined generative models.

Normalization and denoising methods for shotgun metagenomic tables are less developed than those for the marker-gene microbiome data; extra caution is necessary in accounting for the difference in gene sizes and the existence of pseudogenes and duplicate genes in the same genomes ([Quince et al., 2017b](#)). Some recent efforts leverage probabilistic methods

and use a few anchor genes, a similar procedure to the marker gene, see, for instance, [Quince et al. \(2017a\)](#), others use Bayesian hierarchical approaches to model subsets of short reads (called k-mers, see [Lu et al. \(2017\)](#)). Some methods do not explicitly align, assemble, or label reads, but simply embed the k-mers in a continuous space following modern unsupervised methods used in textual analysis ([Menegaux and Vert, 2019](#)), however inferential properties of these algorithms under realistic conditions are not understood. Recently, [Quince et al. \(2020\)](#) devised a bioinformatics tool for high-resolution shotgun metagenomics (STRONG). It may be useful to apply the statistical tools reviewed in this paper to metagenomic data output from that pipeline.

The analytics we have illustrated here used R, Bioconductor packages, and Stan for reproducible microbiome data analyses ([R Core Team, 2013](#); [Gentleman et al., 2004](#); [Carpenter et al., 2017](#)). No such high-level statistical toolbox exists as yet for metagenomic data.

Evaluation of the robustness of inferences to the multiple choices of distances and filtering parameters has only been done empirically up to now. As an example, the choice in the number of topics used in the Latent Dirichlet allocation models could be done by using another level in the Bayesian hierarchical model, but appropriate prior distributions for topics for microbiome data depend strongly on exactly what type of environment is under study, and more calibration experiments are needed to guide the user in their choices. The choice of the prior on the number of topics is important in determining how bacterial diversity changes when the number of specimens increases, as is shown in our example dataset, where the paired-specimen design improved the results of the analyses.

This review has only scratched the surface of the potential for the use of statistics in microbial ecology. [McMurdie and Holmes \(2013\)](#) built the `phyloseq` project as a bridge between bioinformatics pipelines and the statistical toolset enabling users to account for multivariate, spatiotemporal structure in the data. Opportunities abound for refining the basic analyses presented here: many problems have data with dependent sampling designs that benefit from the use of geostatistical methods, time series monitoring under perturbations, and ecological statistics with community assembly and network analyses.

Bayesian approaches have practical advantages; future research directions would focus on Bayesian nonparametric models. To use these methods, it would be worthwhile for the community to publish enhanced cases studies using Bayesian inference in Stan ([Carpenter et al., 2017](#)) and using R/Bioconductor S4 structures and packages. Microbiome research relates tightly to spatiotemporal and single-cell 'omics data, and we hope this review enables statisticians to extend ideas for more sophisticated, reproducible, interpretable microbiome data analyses.

Future work will certainly extend the existing statistical methods used in environmental and ecological studies to provide more complete experimental designs, hierarchical Bayesian models, and the incorporation of geostatistics and spatiotemporal structures. We see a bright future in applying Bayesian and hierarchical methods to analyze these data and look forward to seeing more statisticians involved in this area.

**Acknowledgements.** We are grateful for the thoughtful reading and suggestions made by David Relman and his group, Brian Reich and the referees that helped improve the manuscript. This work was funded by a VMRC grant from the Gates foundation and a grant R01AI112401 from the NIH. We are happy to acknowledge to the R and Bioconductor Core Teams and authors of the packages `BARBI`, `dada2`, `DESeq2`, `phyloseq`, `decontam`, `ggplot2`, `rstan`, `adaptiveGPCA`, `tidyverse` which were used for constructing figures and running the analyses in this paper.

## REFERENCES

- ANDERS, S. and HUBER, W. (2010). Differential expression analysis for sequence count data. *Nature Precedings* 1–1.
- ANDERSON, M. J. (2001). A new method for non-parametric multivariate analysis of variance. *Austral ecology* **26** 32–46.
- ANDERSON, M. J. (2005). Permutational multivariate analysis of variance. *Department of Statistics, University of Auckland, Auckland* **26** 32–46.
- ANDERSON, M. J. and ROBINSON, J. (2003). Generalized discriminant analysis based on distances. *Australian & New Zealand Journal of Statistics* **45** 301–318.
- ANSCOMBE, F. J. (1948). The transformation of Poisson, binomial and negative-binomial data. *Biometrika* **35** 246–254.
- BLEI, D., CARIN, L. and DUNSON, D. (2010). Probabilistic topic models. *IEEE Signal Processing Magazine* **27** 55–65.
- BLEI, D. and LAFFERTY, J. (2006). Correlated topic models. *Advances in Neural Information Processing Systems* **18** 147.
- CALLAHAN, B. J., MCMURDIE, P. J. and HOLMES, S. P. (2017). Exact sequence variants should replace operational taxonomic units in marker-gene data analysis. *The ISME Journal* **11** 2639.
- CALLAHAN, B. J., SANKARAN, K., FUKUYAMA, J. A., MCMURDIE, P. J. and HOLMES, S. P. (2016a). Bio-conductor workflow for microbiome data analysis: from raw reads to community analyses. *F1000Research* **5**.
- CALLAHAN, B. J., MCMURDIE, P. J., ROSEN, M. J., HAN, A. W., JOHNSON, A. J. A. and HOLMES, S. P. (2016b). DADA2: high-resolution sample inference from Illumina amplicon data. *Nature Methods* **13** 581.
- CALLAHAN, B. J., DIGIULIO, D. B., GOLTSMAN, D. S. A., SUN, C. L., COSTELLO, E. K., JEGANATHAN, P., BIGGIO, J. R., WONG, R. J., DRUZIN, M. L. and SHAW, G. M. (2017). Replication and refinement of a vaginal microbial signature of preterm birth in two racially distinct cohorts of US women. *Proceedings of the National Academy of Sciences* **114** 9966–9971.
- CARPENTER, B., GELMAN, A., HOFFMAN, M. D., LEE, D., GOODRICH, B., BETANCOURT, M., BRUBAKER, M., GUO, J., LI, P. and RIDDELL, A. (2017). Stan: A probabilistic programming language. *Journal of Statistical Software* **76**.
- CAVICCHIOLI, R., RIPPLE, W. J., TIMMIS, K. N., AZAM, F., BAKKEN, L. R., BAYLIS, M., BEHRENFELD, M. J., BOETIUS, A., BOYD, P. W. and CLASSEN, A. T. (2019). Scientists’ warning to humanity: microorganisms and climate change. *Nature Reviews Microbiology* **17** 569–586.
- CHENG, H. K., TAN, S. K., SWEENEY, T. E., JEGANATHAN, P., BRIESE, T., KHADKA, V., STROUTS, F., THAIR, S., DALAI, S. and HITCHCOCK, M. (2019). Combined use of metagenomic sequencing and host response profiling for the diagnosis of suspected sepsis. *BioRxiv* 854182.
- COLE, J. R., WANG, Q., CARDENAS, E., FISH, J., CHAI, B., FARRIS, R. J., KULAM-SYED-MOHIDEEN, A., MCGARRELL, D. M., MARSH, T., GARRITY, G. M. and TIEDJE, J. (2008). The Ribosomal Database Project: improved alignments and new tools for rRNA analysis. *Nucleic Acids Research* **37** D141–D145.
- COMPANT, S., SAMAD, A., FAIST, H. and SESSITSCH, A. (2019). A review on the plant microbiome: ecology, functions, and emerging trends in microbial application. *Journal of Advanced Research* **19** 29–37.
- DAVIS, N. M., PROCTOR, D. M., HOLMES, S. P., RELMAN, D. A. and CALLAHAN, B. J. (2018). Simple statistical identification and removal of contaminant sequences in marker-gene and metagenomics data. *Microbiome* **6** 226.
- DELGADO-BAQUERIZO, M., MAESTRE, F. T., REICH, P. B., JEFFRIES, T. C., GAITAN, J. J., ENCINAR, D., BERDUGO, M., CAMPBELL, C. D. and SINGH, B. K. (2016). Microbial diversity drives multifunctionality in terrestrial ecosystems. *Nature Communications* **7** 10541.
- DESANTIS, T. Z., HUGENHOLTZ, P., LARSEN, N., ROJAS, M., BRODIE, E. L., KELLER, K., HUBER, T., DALEVI, D., HU, P. and ANDERSEN, G. L. (2006). Greengenes, a Chimera-Checked 16S rRNA Gene Database and Workbench Compatible with ARB. *Applied and Environmental Microbiology* **72** 5069–5072.
- DIGIULIO, D. B., CALLAHAN, B. J., MCMURDIE, P. J., COSTELLO, E. K., LYELL, D. J., ROBACZEWSKA, A., SUN, C. L., GOLTSMAN, D. S., WONG, R. J. and SHAW, G. (2015). Temporal and spatial variation of the human microbiota during pregnancy. *Proceedings of the National Academy of Sciences* **112** 11060–11065.
- EVANS, S. N. and MATSEN, F. A. (2012). The phylogenetic Kantorovich–Rubinstein metric for environmental sequence samples. *Journal of the Royal Statistical Society: Series B (Statistical Methodology)* **74** 569–592.
- EXCOFFIER, L., SMOUSE, P. E. and QUATTRO, J. M. (1992). Analysis of molecular variance inferred from metric distances among DNA haplotypes: application to human mitochondrial DNA restriction data. *Genetics* **131** 479–491.

- FITZPATRICK, C. R., LU-IRVING, P., COPELAND, J., GUTTMAN, D. S., WANG, P. W., BALTRUS, D. A., DLUGOSCH, K. M. and JOHNSON, M. T. (2018). Chloroplast sequence variation and the efficacy of peptide nucleic acids for blocking host amplification in plant microbiome studies. *Microbiome* **6** 144.
- FRANZOSA, E. A., MORGAN, X. C., SEGATA, N., WALDRON, L., REYES, J., EARL, A. M., GIANNOUKOS, G., BOYLAN, M. R., CIULLA, D., GEVERS, D., IZARD, J., GARRETT, W. S., CHAN, A. T. and HUTTENHOWER, C. (2014). Relating the metatranscriptome and metagenome of the human gut. *Proceedings of the National Academy of Sciences* **111** E2329–E2338.
- FUKUYAMA, J. (2019). Adaptive gPCA: A method for structured dimensionality reduction with applications to microbiome data. *Annals of Applied Statistics* **13** 1043–1067.
- FUKUYAMA, J. (2020). phyloseqGraphTest: Graph-Based Permutation Tests for Microbiome Data R package version 0.1.0.
- FUKUYAMA, J., MCMURDIE, P. J., DETHLEFSEN, L., RELMAN, D. A. and HOLMES, S. (2012). Comparisons of distance methods for combining covariates and abundances in microbiome studies. *Pacific Symposium on Biocomputing* 213–224.
- FUKUYAMA, J., RUMKER, L., SANKARAN, K., JEGANATHAN, P., DETHLEFSEN, L., RELMAN, D. A. and HOLMES, S. P. (2017). Multidomain analyses of a longitudinal human microbiome intestinal cleanout perturbation experiment. *PLoS Computational Biology* **13** e1005706.
- GELMAN, A. and RUBIN, D. B. (1992). Inference from iterative simulation using multiple sequences. *Statistical Science* **7** 457–472.
- GENTLEMAN, R. C., CAREY, V. J., BATES, D. M., BOLSTAD, B., DETTLING, M., DUDOIT, S., ELLIS, B., GAUTIER, L., GE, Y. and GENTRY, J. (2004). Bioconductor: Open software development for computational biology and bioinformatics. *Genome Biology* **5** R80.
- GEVERS, D., KUGATHASAN, S., DENSON, L. A., VÁZQUEZ-BAEZA, Y., VAN TREUREN, W., REN, B., SCHWAGER, E., KNIGHTS, D., SONG, S. J. and YASSOUR, M. (2014). The treatment-naïve microbiome in new-onset Crohn’s disease. *Cell Host & Microbe* **15** 382–392.
- GILBERT, J. A., JANSSON, J. K. and KNIGHT, R. (2014). The Earth Microbiome project: successes and aspirations. *BMC Biology* **12** 69.
- GLOOR, G. B., MACKLAIM, J. M., PAWLOWSKY-GLAHN, V. and EGOZCUE, J. J. (2017). Microbiome datasets are compositional: and this is not optional. *Frontiers in Microbiology* **8** 2224.
- GORVITOVSKAIA, A., HOLMES, S. P. and HUSE, S. M. (2016). Interpreting Prevotella and Bacteroides as biomarkers of diet and lifestyle. *Microbiome* **4** 15.
- GRANTHAM, N. S., GUAN, Y., REICH, B. J., BORER, E. T. and GROSS, K. (2020). Mimix: A Bayesian mixed-effects model for microbiome data from designed experiments. *Journal of the American Statistical Association* **115** 599–609.
- GREENACRE, M. (2010a). Log-ratio analysis is a limiting case of correspondence analysis. *Mathematical Geosciences* **42** 129.
- GREENACRE, M. (2010b). Correspondence analysis of raw data. *Ecology* **91** 958–963.
- GREENACRE, M. (2011). Compositional data and correspondence analysis. *Compositional Data Analysis* 103–113.
- GRÉGORY, D., CHAUDET, H., LAGIER, J.-C. and RAOULT, D. (2018). How mass spectrometric approaches applied to bacterial identification have revolutionized the study of human gut microbiota. *Expert Review of Proteomics* **15** 217–229.
- GRUMAZ, S., STEVENS, P., GRUMAZ, C., DECKER, S. O., WEIGAND, M. A., HOFER, S., BRENNER, T., VON HAESELER, A. and SOHN, K. (2016). Next-generation sequencing diagnostics of bacteremia in septic patients. *Genome Medicine* **8** 73.
- HARRIS, K., PARSONS, T. L., IJAZ, U. Z., LAHTI, L., HOLMES, I. and QUINCE, C. (2015). Linking statistical and ecological theory: Hubbell’s unified neutral theory of biodiversity as a hierarchical Dirichlet process. *Proceedings of the IEEE* **105** 516–529.
- HOLMES, S. (2008). Multivariate data analysis: the French way. In *Probability and Statistics: Essays in honor of David A. Freedman* 219–233. Institute of Mathematical Statistics.
- HOLMES, I., HARRIS, K. and QUINCE, C. (2012). Dirichlet multinomial mixtures: generative models for microbial metagenomics. *PLoS One* **7** e30126.
- HOLMES, S. and HUBER, W. (2018). *Modern statistics for modern biology*. Cambridge University Press.
- HONG, D. K., BLAUWKAMP, T. A., KERTESZ, M., BERCOVICI, S., TRUONG, C. and BANAEI, N. (2018). Liquid biopsy for infectious diseases: sequencing of cell-free plasma to detect pathogen DNA in patients with invasive fungal disease. *Diagnostic Microbiology and Infectious Disease* **92** 210–213.
- JEGANATHAN, P., CALLAHAN, B. J., PROCTOR, D. M., RELMAN, D. A. and HOLMES, S. P. (2018). The Block Bootstrap Method for Longitudinal Microbiome Data. *ArXiv preprint ArXiv:1809.01832*.

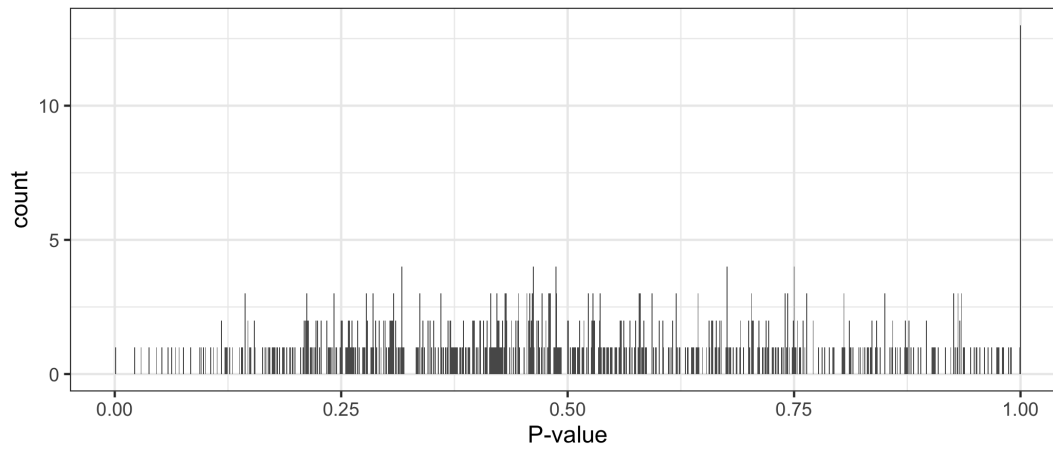


- KOSTIC, A. D., GEVERS, D., SILJANDER, H., VATANEN, T., HYÖTYLÄINEN, T., HÄMÄLÄINEN, A.-M., PEET, A., TILLMANN, V., PÖHÖ, P. and MATTILA, I. (2015). The dynamics of the human infant gut microbiome in development and in progression toward type 1 diabetes. *Cell Host & Microbe* **17** 260–273.
- KRUSCHKE, J. (2014). *Doing Bayesian Data Analysis: A tutorial with R, JAGS, and Stan*. Academic Press.
- KUNTAL, B. K. and MANDE, S. S. (2019). Visual exploration of microbiome data. *Journal of Biosciences* **44** 119.
- KURTZ, Z. D., MÜLLER, C. L., MIRALDI, E. R., LITTMAN, D. R., BLASER, M. J. and BONNEAU, R. A. (2015). Sparse and compositionally robust inference of microbial ecological networks. *PLoS Computational Biology* **11** e1004226.
- LAW, C. W., CHEN, Y., SHI, W. and SMYTH, G. K. (2014). voom: Precision weights unlock linear model analysis tools for RNA-seq read counts. *Genome Biology* **15** R29.
- LOVE, M. I., HUBER, W. and ANDERS, S. (2014). Moderated estimation of fold change and dispersion for RNA-seq data with DESeq2. *Genome Biology* **15** 550.
- LOZUPONE, C. and KNIGHT, R. (2005). UniFrac: A new phylogenetic method for comparing microbial communities. *Applied and Environmental Microbiology* **71** 8228.
- LU, J., BREITWIESER, F. P., THIELEN, P. and SALZBERG, S. L. (2017). Bracken: estimating species abundance in metagenomics data. *PeerJ Computer Science* **3** e104.
- MCLAREN, M. R., WILLIS, A. D. and CALLAHAN, B. J. (2019). Consistent and correctable bias in metagenomic sequencing experiments. *Elife* **8** e46923.
- MCMURDIE, P. J. and HOLMES, S. (2013). phyloseq: An R package for reproducible interactive analysis and graphics of microbiome census data. *PLoS One* **8** e61217.
- MCMURDIE, P. J. and HOLMES, S. (2014). Waste not, want not: why rarefying microbiome data is inadmissible. *PLoS Computational Biology* **10** e1003531.
- MCMURDIE, P. J. and HOLMES, S. (2015). Shiny-phyloseq: Web application for interactive microbiome analysis with provenance tracking. *Bioinformatics* **31** 282–283.
- MENEGAUX, R. and VERT, J.-P. (2019). Continuous embeddings of DNA sequencing reads and application to metagenomics. *Journal of Computational Biology* **26** 509–518.
- NGUYEN, L. H. and HOLMES, S. (2017). Bayesian unidimensional scaling for visualizing uncertainty in high dimensional datasets with latent ordering of observations. *BMC Bioinformatics* **18** 394.
- NGUYEN, L. H. and HOLMES, S. (2019). Ten quick tips for effective dimensionality reduction. *PLoS Computational Biology* **15**.
- OKSANEN, J., BLANCHET, F. G., FRIENDLY, M., KINDT, R., LEGENDRE, P., MCGLINN, D., MINCHIN, P. R., O'HARA, R. B., SIMPSON, G. L., SOLYMOS, P., STEVENS, M. H. H., SZOECS, E. and WAGNER, H. (2020). vegan: Community Ecology Package R package version 2.5-7.
- PAVOINE, S., DUFOUR, A.-B. and CHESSEL, D. (2004). From dissimilarities among species to dissimilarities among communities: a double principal coordinate analysis. *Journal of Theoretical Biology* **228** 523–537.
- PROCTOR, D. M. and RELMAN, D. A. (2017). The landscape ecology and microbiota of the human nose, mouth, and throat. *Cell Host & Microbe* **21** 421–432.
- PROCTOR, D. M., FUKUYAMA, J. A., LOOMER, P. M., ARMITAGE, G. C., LEE, S. A., DAVIS, N. M., RYDER, M. I., HOLMES, S. P. and RELMAN, D. A. (2018). A spatial gradient of bacterial diversity in the human oral cavity shaped by salivary flow. *Nature Communications* **9** 1–10.
- PRUESSE, E., QUAST, C., KNITTEL, K., FUCHS, B. M., LUDWIG, W., PEPLIES, J. and GLÖCKNER, F. O. (2007). SILVA: A comprehensive online resource for quality checked and aligned ribosomal RNA sequence data compatible with ARB. *Nucleic Acids Research* **35** 7188–7196.
- PURDOM, E. (2011). Analysis of a data matrix and a graph: Metagenomic data and the phylogenetic tree. *The Annals of Applied Statistics* **5** 2326–2358.
- QUINCE, C., DELMONT, T., RAGUIDEAU, S., ALNEBERG, J., DARLING, A., COLLINS, G. and EREN, M. (2017a). DESMAN: A new tool for de novo extraction of strains from metagenomes. *Genome Biology* **18** 1–22.
- QUINCE, C., WALKER, A., SIMPSON, J., LOMAN, N. and SEGATA, N. (2017b). Shotgun metagenomics, from sampling to analysis. *Nature Biotechnology* **35** 833–844.
- QUINCE, C., NURK, S., RAGUIDEAU, S., JAMES, R. S., SOYER, O. S., SUMMERS, J. K., LIMASSET, A., EREN, A. M., CHIKHI, R. and DARLING, A. E. (2020). Metagenomics Strain Resolution on Assembly Graphs. *BioRxiv*.
- QUINN, T. P., ERB, I., RICHARDSON, M. F. and CROWLEY, T. M. (2018). Understanding sequencing data as compositions: an outlook and review. *Bioinformatics* **34** 2870–2878.
- RAMIREZ, K. S., KNIGHT, C. G., DE HOLLANDER, M., BREARLEY, F. Q., CONSTANTINIDES, B., COTTON, A., CREER, S., CROWTHER, T. W., DAVISON, J., DELGADO-BAQUERIZO, M. and DORREPAAL, E. (2018). Detecting macroecological patterns in bacterial communities across independent studies of global soils. *Nature Microbiology* **3** 189.

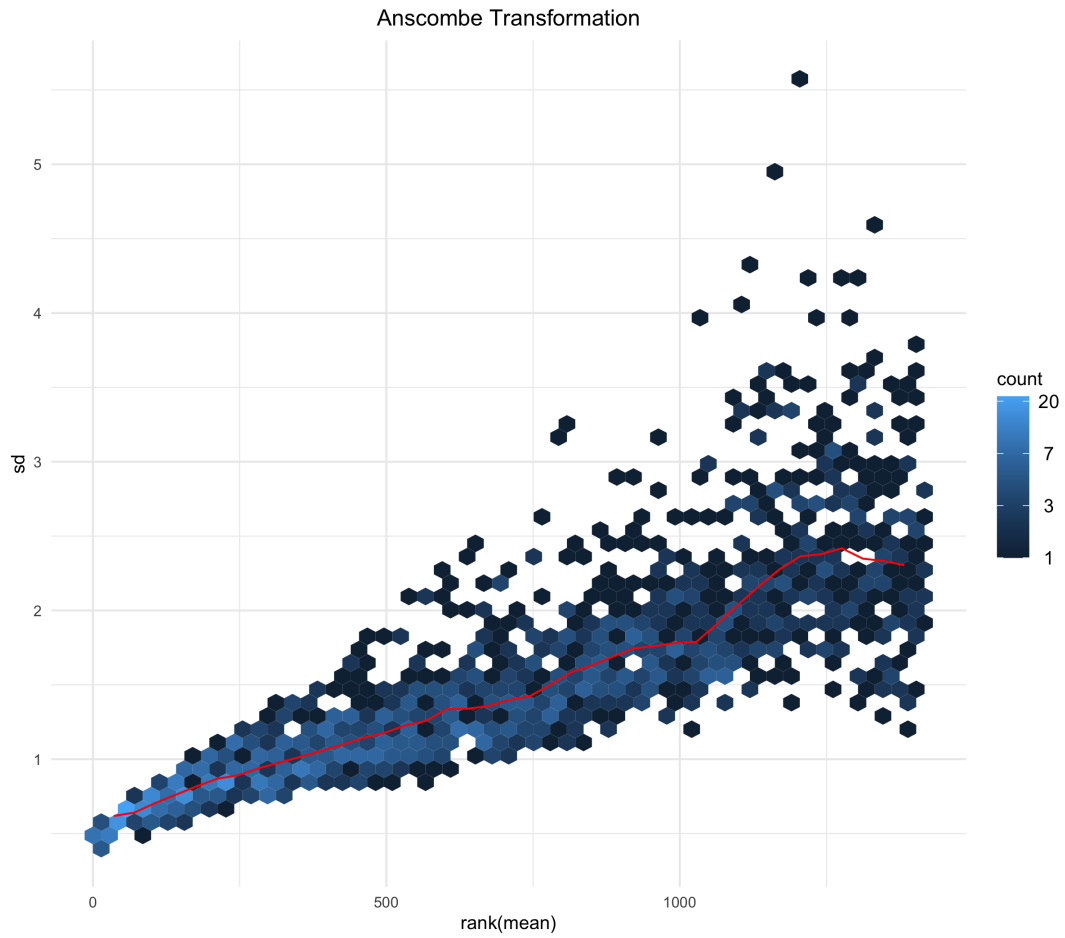
- REN, B., BACALLADO, S., FAVARO, S., HOLMES, S. and TRIPPA, L. (2017). Bayesian nonparametric ordination for the analysis of microbial communities. *Journal of the American Statistical Association* **112** 1430–1442.
- ROBINSON, M. D., MCCARTHY, D. J. and SMYTH, G. K. (2010). edgeR: A Bioconductor package for differential expression analysis of digital gene expression data. *Bioinformatics* **26** 139–140.
- ROMERO, R., HASSAN, S. S., GAJER, P., TARCA, A. L., FADROSH, D. W., BIEDA, J., CHAEMSAITHONG, P., MIRANDA, J., CHAIWORAPONGSA, T. and RAVEL, J. (2014a). The vaginal microbiota of pregnant women who subsequently have spontaneous preterm labor and delivery and those with a normal delivery at term. *Microbiome* **2** 18.
- ROMERO, R., HASSAN, S. S., GAJER, P., TARCA, A. L., FADROSH, D. W., NIKITA, L., GALUPPI, M., LAMONT, R. F., CHAEMSAITHONG, P., MIRANDA, J., CHAIWORAPONGSA, T. and RAVEL, J. (2014b). The composition and stability of the vaginal microbiota of normal pregnant women is different from that of non-pregnant women. *Microbiome* **2** 4.
- ROSEN, G. L., REICHENBERGER, E. R. and ROSENFELD, A. M. (2011). NBC: the Naive Bayes Classification tool webserver for taxonomic classification of metagenomic reads. *Bioinformatics* **27** 127–129.
- SALTER, S. J., COX, M. J., TUREK, E. M., CALUS, S. T., COOKSON, W. O., MOFFATT, M. F., TURNER, P., PARKHILL, J., LOMAN, N. J. and WALKER, A. W. (2014). Reagent and laboratory contamination can critically impact sequence-based microbiome analyses. *BMC Biology* **12** 87.
- SANKARAN, K. and HOLMES, S. (2014). structSSI: Simultaneous and Selective Inference for Grouped or Hierarchically Structured Data. *Journal of Statistical Software* **59** 1–21.
- SANKARAN, K. and HOLMES, S. P. (2017). treelapse: Visualization of hierarchically structured data.
- SANKARAN, K. and HOLMES, S. (2018). Interactive Visualization of Hierarchically Structured Data. *Journal of Computational and Graphical Statistics* **27** 553–563.
- SANKARAN, K. and HOLMES, S. P. (2019a). Multitable methods for microbiome data integration. *Frontiers in Genetics* **10** 627.
- SANKARAN, K. and HOLMES, S. P. (2019b). Latent variable modeling for the microbiome. *Biostatistics* **20** 599–614.
- SEGATA, N., IZARD, J., WALDRON, L., GEVERS, D., MIROPOLSKY, L., GARRETT, W. S. and HUTTENHOWER, C. (2011). Metagenomic biomarker discovery and explanation. *Genome Biology* **12** 1–18.
- SILVERMAN, J. D., WASHBURNE, A. D., MUKHERJEE, S. and DAVID, L. A. (2017). A phylogenetic transform enhances analysis of compositional microbiota data. *Elife* **6** e21887.
- SINGH, S. P., STAIU, A.-M., DUNN, R. R., FIERER, N. and REICH, B. J. (2019). A nonparametric spatial test to identify factors that shape a microbiome. *The Annals of Applied Statistics* **13** 2341–2362.
- SMYTH, G. K. (2005). Limma: linear models for microarray data. In *Bioinformatics and Computational Biology Solutions using R and Bioconductor* 397–420. Springer.
- SNIJEDERS, T. A. and NOWICKI, K. (1997). Estimation and prediction for stochastic blockmodels for graphs with latent block structure. *Journal of Classification* **14** 75–100.
- R CORE TEAM (2013). R: A language and environment for statistical computing.
- THOMPSON, L. R., SANDERS, J. G., MCDONALD, D., AMIR, A., LADAU, J., LOCEY, K. J., PRILL, R. J., TRIPATHI, A., GIBBONS, S. M. and ACKERMANN, G. (2017). A communal catalogue reveals Earth’s multiscale microbial diversity. *Nature* **551** 457.
- VEHTARI, A., GELMAN, A., SIMPSON, D., CARPENTER, B. and BÜRKNER, P.-C. (4 July 2020). Rank-normalization, folding, and localization: An improved  $\hat{R}$  for assessing convergence of MCMC. *Bayesian Analysis*, advance publication.
- WASHBURNE, A. D., SILVERMAN, J. D., MORTON, J. T., BECKER, D. J., CROWLEY, D., MUKHERJEE, S., DAVID, L. A. and PLOWRIGHT, R. K. (2019). Phylofactorization: A graph partitioning algorithm to identify phylogenetic scales of ecological data. *Ecological Monographs* **89** e01353.
- WICKHAM, H. (2016). *ggplot2: Elegant Graphics For Data Analysis*. Springer.
- XU, L., PATERSON, A. D., TURPIN, W. and XU, W. (2015). Assessment and selection of competing models for zero-inflated microbiome data. *PloS one* **10** e0129606.
- YATSUNENKO, T., REY, F. E., MANARY, M. J., TREHAN, I., DOMINGUEZ-BELLO, M. G., CONTRERAS, M., MAGRIS, M., HIDALGO, G., BALDASSANO, R. N. and ANOKHIN, A. P. (2012). Human gut microbiome viewed across age and geography. *Nature* **486** 222–227.
- ZHAO, N., CHEN, J., CARROLL, I. M., RINGEL-KULKA, T., EPSTEIN, M. P., ZHOU, H., ZHOU, J. J., RINGEL, Y., LI, H. and WU, M. C. (2015). Testing in microbiome-profiling studies with MiRKAT, the microbiome regression-based kernel association test. *The American Journal of Human Genetics* **96** 797–807.

## SUPPLEMENTARY MATERIAL

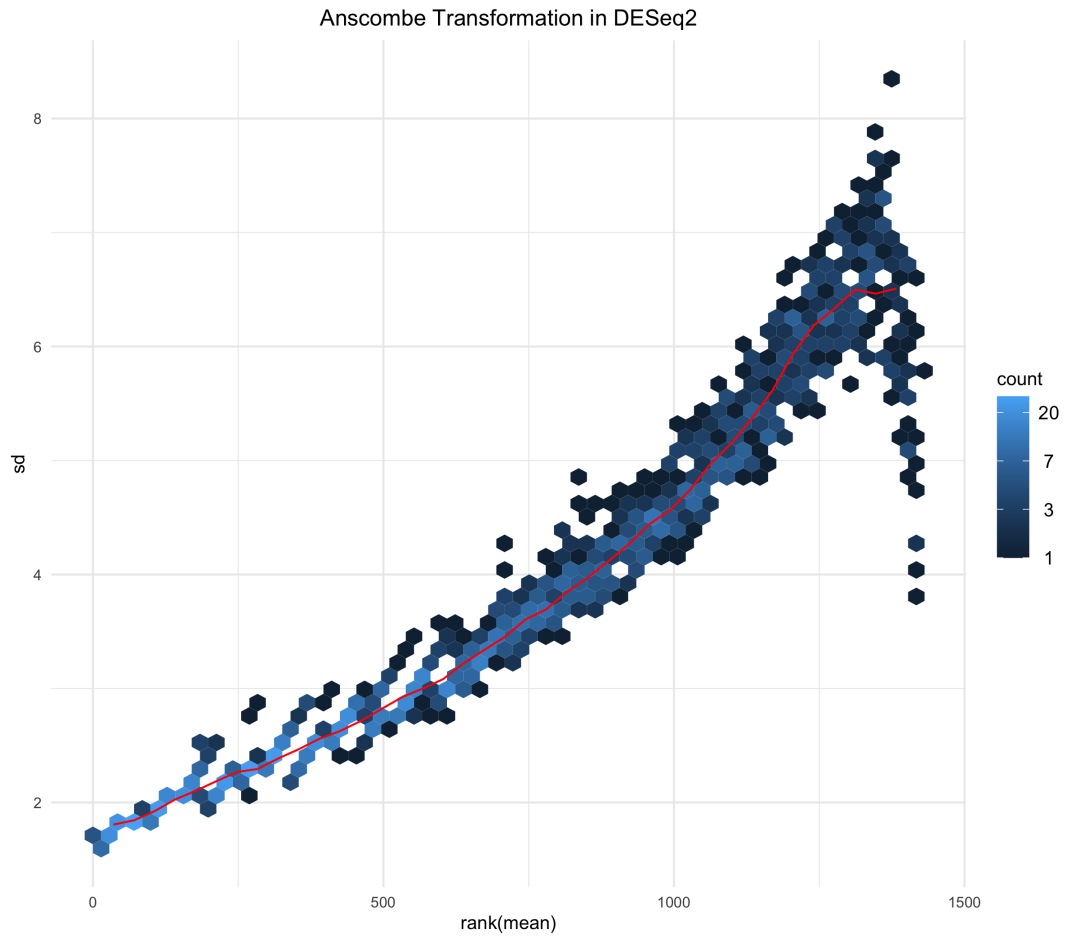
**Figures and tables produced in the analyses.**



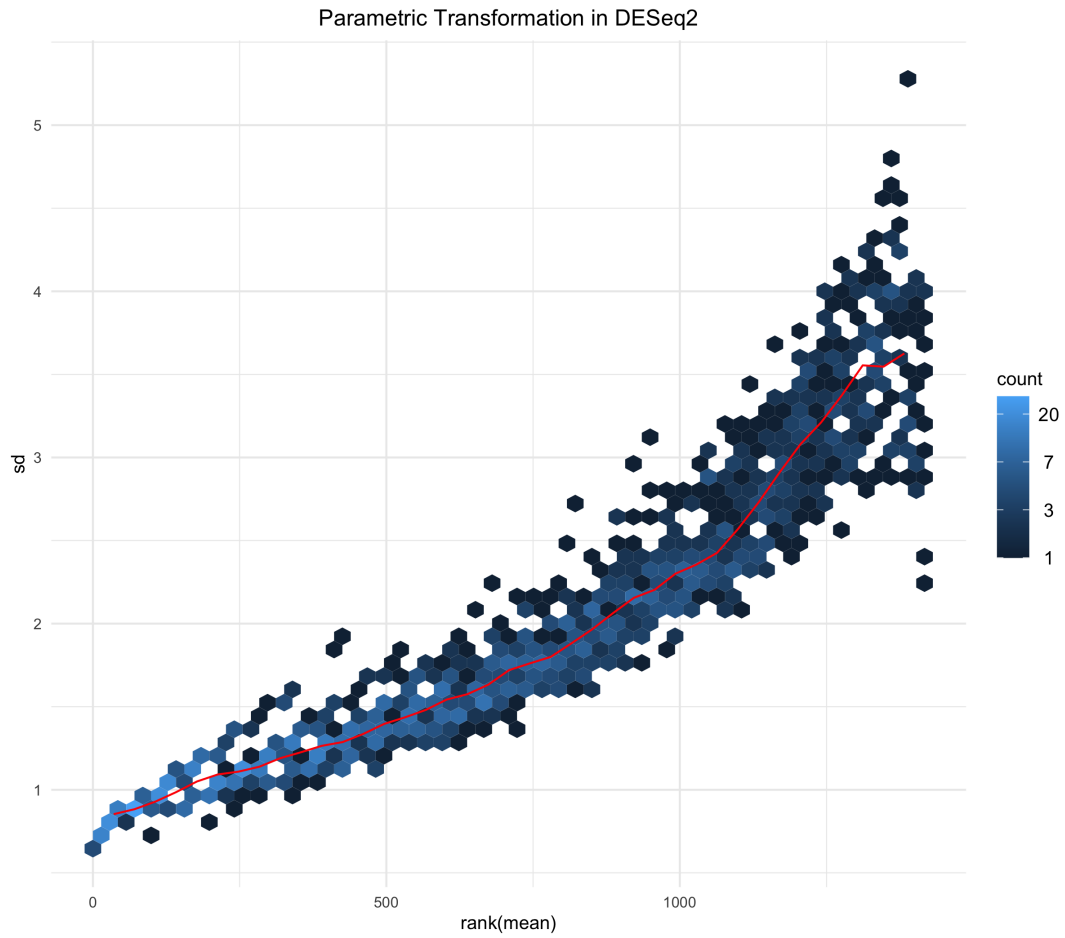
SUPPLEMENTARY FIGURE 1. *P* values for goodness of fit test of negative binomial distribution for each taxon.



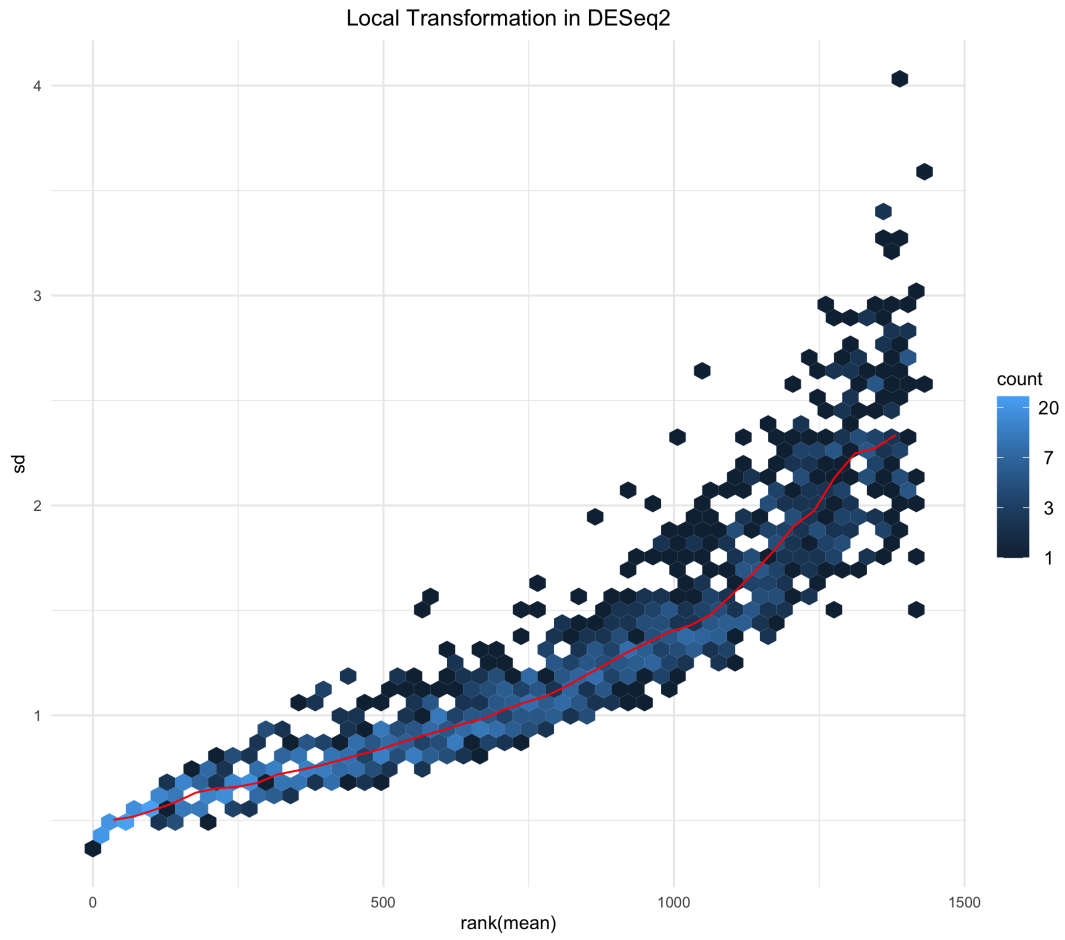
SUPPLEMENTARY FIGURE 2. *Anscombe's transformation of abundance data. Hexagonal binning avoids overplotting*



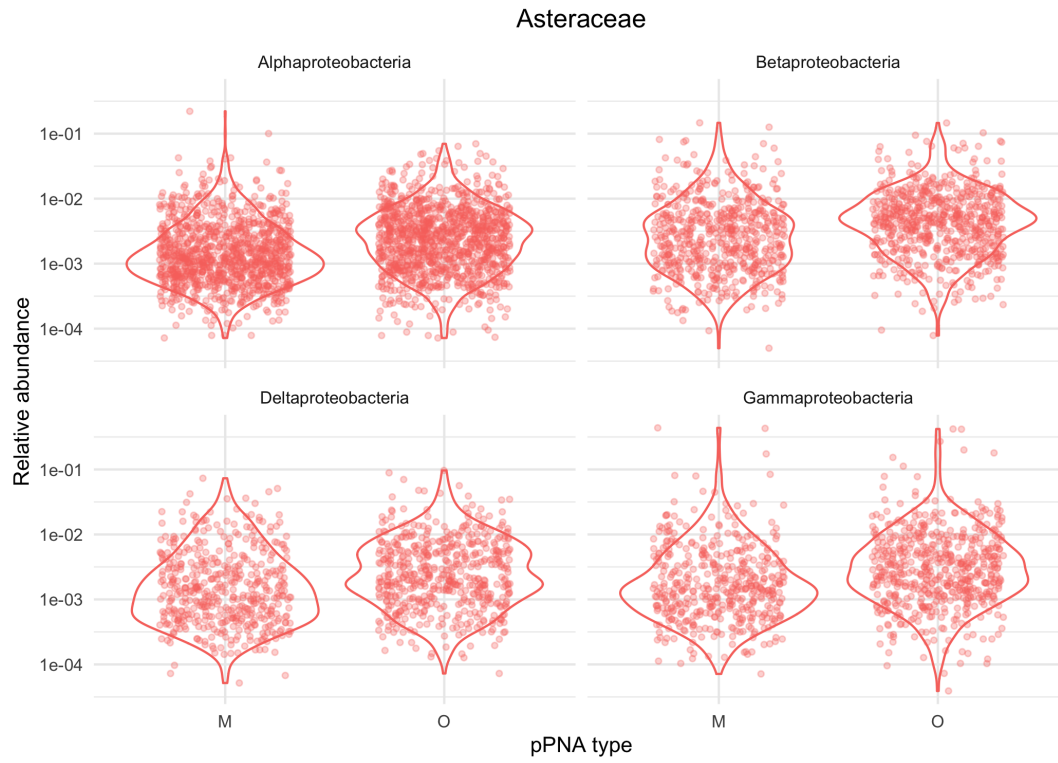
SUPPLEMENTARY FIGURE 3. *Anscombe's transformation implemented in DESeq2 package.*



SUPPLEMENTARY FIGURE 4. *Parametric transformation implemented in DESeq2 package.*



SUPPLEMENTARY FIGURE 5. *Nonparametric transformation implemented in DESeq2 package.*



SUPPLEMENTARY FIGURE 6. *Distribution of relative abundance of four different Classes of Proteobacteria in Asteraceae plants. O and M denote universal and Asteraceae-modified pPNA types, respectively.*

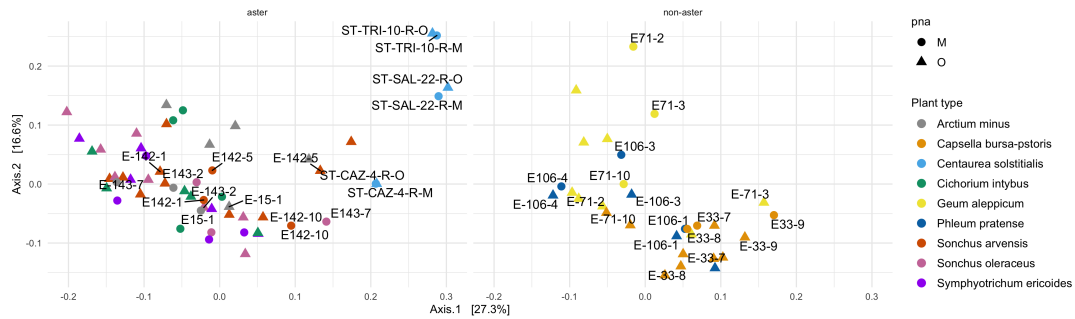




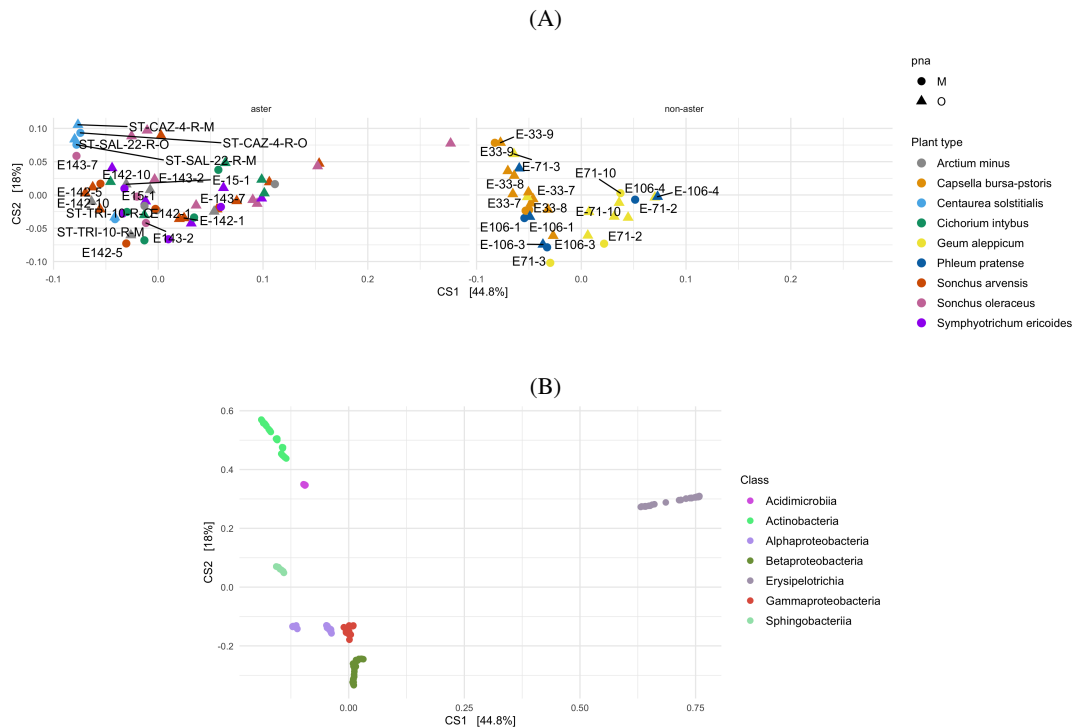
SUPPLEMENTARY FIGURE 7. Distribution of relative abundance of four different Classes of Proteobacteria in non-Asteraceae plants. O and M denote universal and Asteraceae-modified pPNA types, respectively.

SUPPLEMENTARY TABLE 1  
Paired specimens in all plant types.

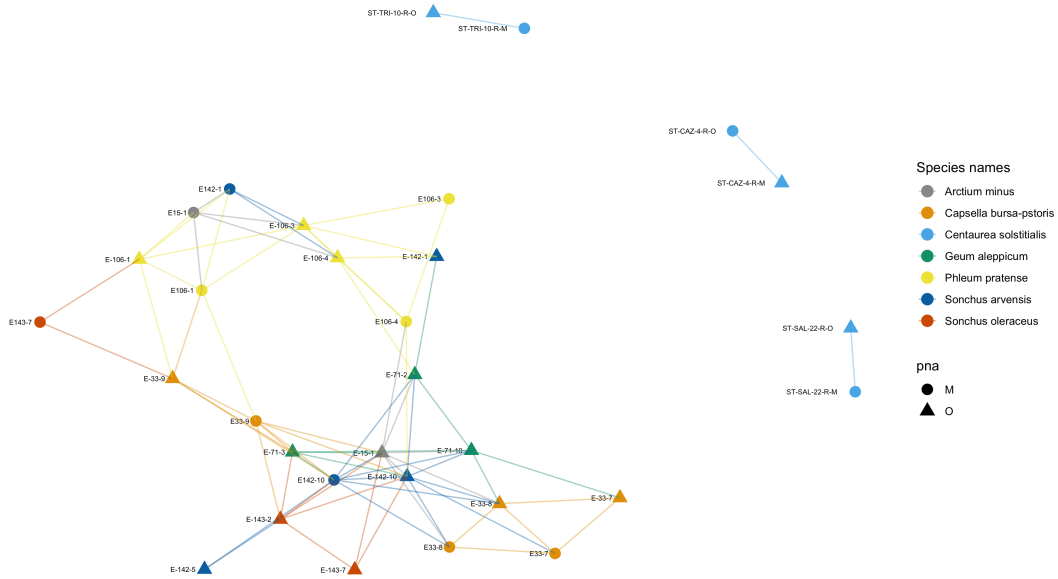
Plant	Type	universal pPNA	Asteraceae-modified pPNA
<i>Arctium minus</i>	Asteraceae	E-15-1	E15-1
<i>Sonchus arvensis</i>	Asteraceae	E-142-1	E142-1
<i>Sonchus arvensis</i>	Asteraceae	E-142-5	E142-5
<i>Sonchus arvensis</i>	Asteraceae	E-142-10	E142-10
<i>Sonchus oleraceus</i>	Asteraceae	E-143-2	E143-2
<i>Sonchus oleraceus</i>	Asteraceae	E-143-7	E143-7
<i>Centaurea solstitialis</i>	Asteraceae	ST-CAZ-4-R-O	ST-CAZ-4-R-M
<i>Centaurea solstitialis</i>	Asteraceae	ST-SAL-22-R-O	ST-SAL-22-R-M
<i>Centaurea solstitialis</i>	Asteraceae	ST-TRI-10-R-O	ST-TRI-10-R-M
<i>Capsella bursa-pstoris</i>	non-Asteraceae	E-33-7	E33-7
<i>Capsella bursa-pstoris</i>	non-Asteraceae	E-33-8	E33-8
<i>Capsella bursa-pstoris</i>	non-Asteraceae	E-33-9	E33-9
<i>Geum aleppicum</i>	non-Asteraceae	E-71-2	E71-2
<i>Geum aleppicum</i>	non-Asteraceae	E-71-3	E71-3
<i>Geum aleppicum</i>	non-Asteraceae	E-71-10	E71-10
<i>Phleum pratense</i>	non-Asteraceae	E-106-1	E106-1
<i>Phleum pratense</i>	non-Asteraceae	E-106-3	E106-3
<i>Phleum pratense</i>	non-Asteraceae	E-106-4	E106-4



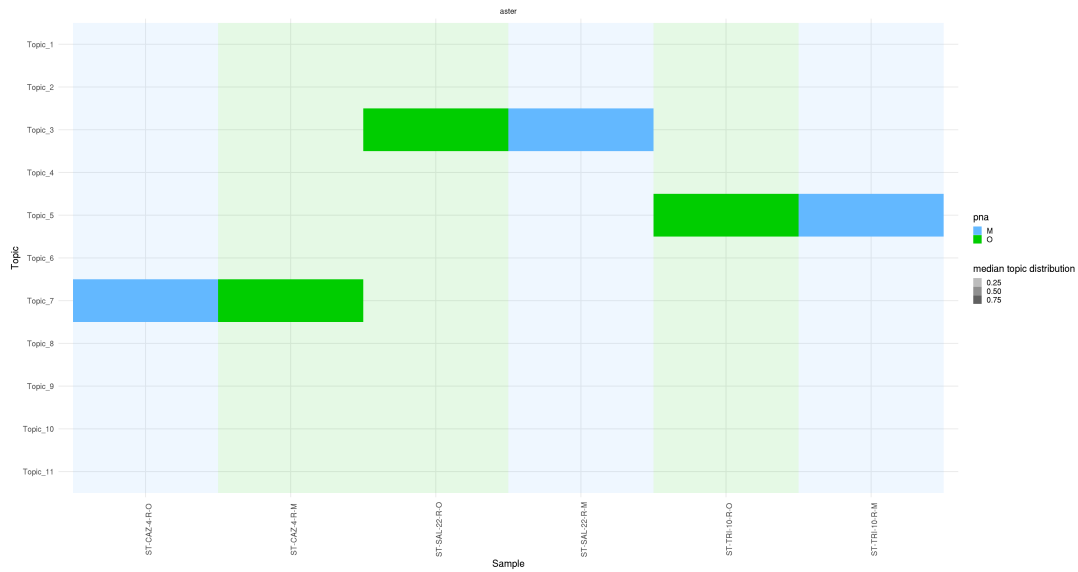
SUPPLEMENTARY FIGURE 8. Multidimensional scaling (MDS) with weighted unifrac distance on all specimens. The shape denotes universal (O) and Asteraceae-modified (M) pPNA types, respectively. Facet denotes Asteraceae or non-Asteraceae plants. Paired specimens are labeled and make clusters, except paired-specimens (E-143-7, E143-7), (E-142-5, E142-5), (E-71-2, E71-2), and (E-71-3, E71-3) that are in positive and negative axes. *Centaurea solstitialis* specimens are outliers among Asteraceae plants in the positive direction of Axis 2. Axis 1 explains the microbial variability in plant types.



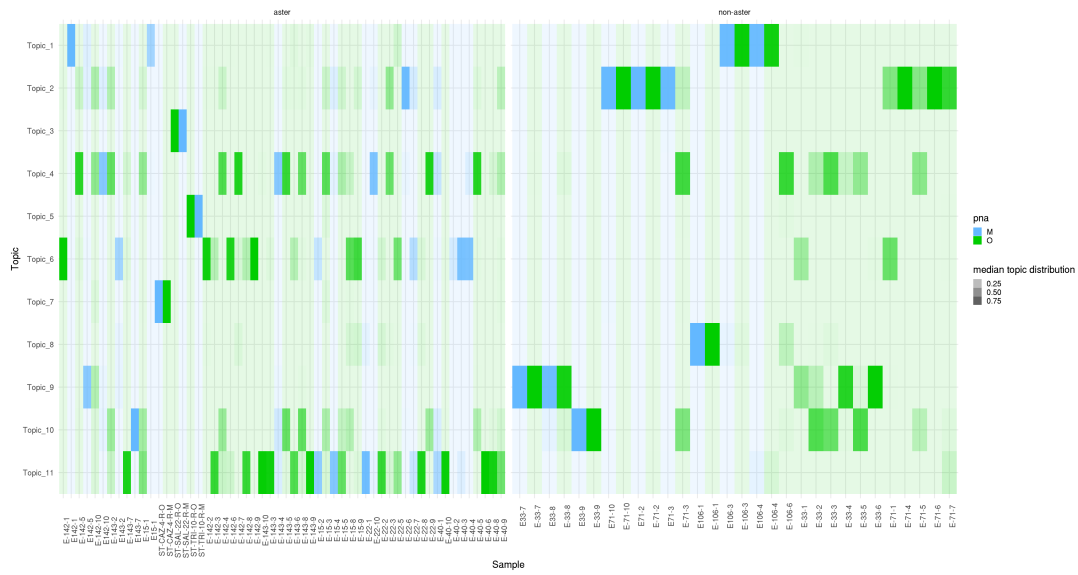
SUPPLEMENTARY FIGURE 9. (A) A DPCoA plot that incorporates phylogenetic information. The shape denotes universal (O) and Asteraceae-modified (M) pPNA types, respectively. Facet denotes Asteraceae or non-Asteraceae plants. Paired specimens are labeled and make clusters, except paired-specimens (E-143-7, E143-7), (E-142-5, E142-5), and (E-71-3, E71-3) that are in positive and negative axes. Axis 1 explains the microbial variability in all paired specimens and specimens from *Sonchus oleraceus* and *Sonchus arvensis* plants with highly abundant *Erysipelotrichia*. (B) The DPCoA specimen ordination that is interpreted with respect to the ASV coordinates.



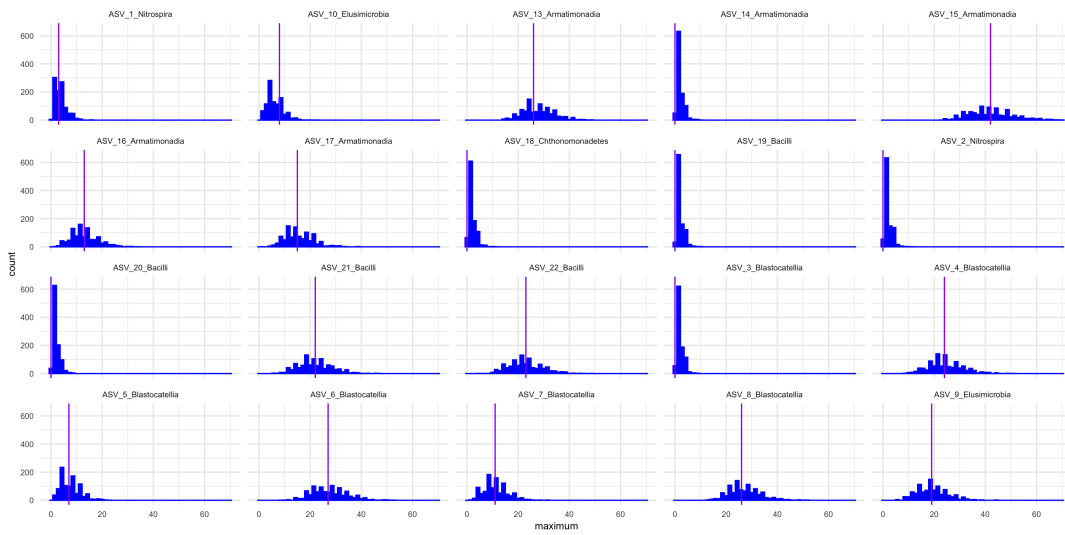
SUPPLEMENTARY FIGURE 10. A network created by thresholding the Jaccard dissimilarity matrix at 0.8. All paired-specimens are connected, except (E-143-7, E143-7), (E-142-5, E142-5), (E-71-2, E71-2), and (E-71-3, E71-3).



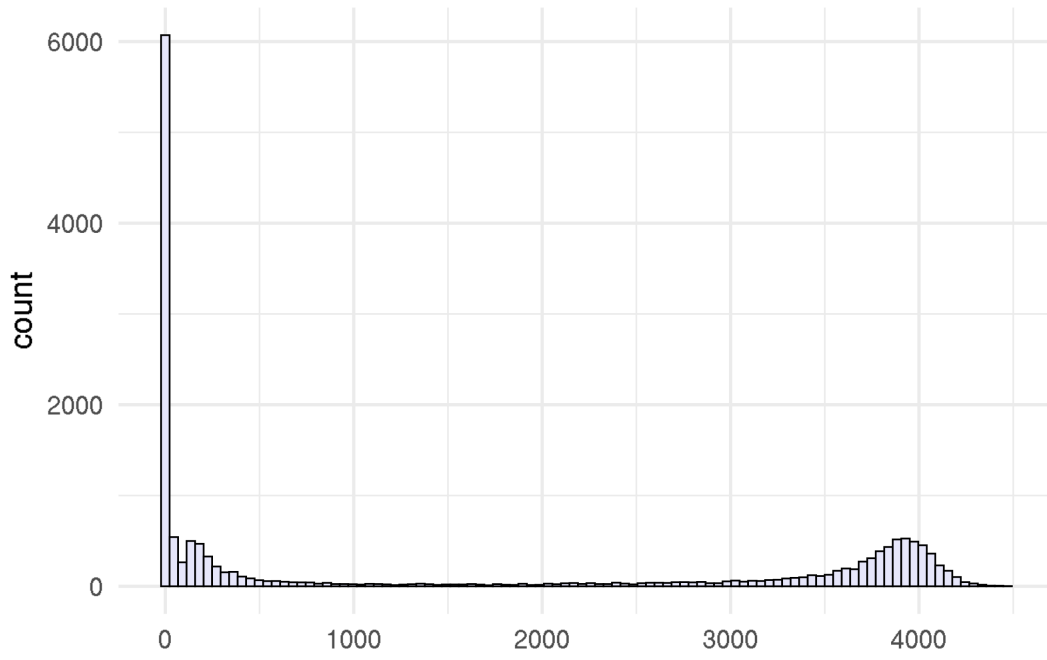
SUPPLEMENTARY FIGURE 11. Topic distribution in specimens from *Centaurea solstitialis* with eleven topics, which is from three different countries. This plant has three paired specimens sequenced with O and M-pNA types. The color gradient represents the median topic distribution.



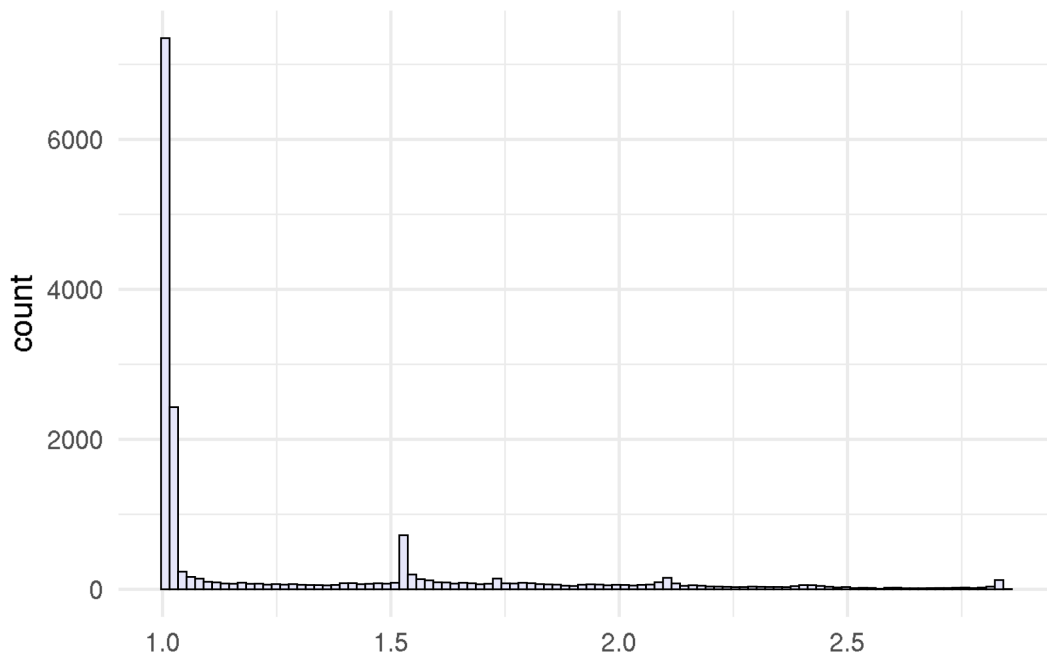
SUPPLEMENTARY FIGURE 12. *Topic distribution in all specimens with eleven topics. The color gradient represents the median topic distribution.*



SUPPLEMENTARY FIGURE 13. *Predictive model check with simulated data, observed data, and a statistic  $G(K_{ij}) = \max\{K_{ij}\}$ . Each facet shows the histogram of  $G(K_{ij})$  of each ASV in specimens from the posterior predictive distribution and the vertical line shows the value of  $G(K_{ij})$  of each ASV in observed data.*



SUPPLEMENTARY FIGURE 14. *Effective sample size (ESS) with eleven topics.*



SUPPLEMENTARY FIGURE 15. *Split  $\hat{R}$  with eleven topics.*

SUPPLEMENTARY TABLE 2

*Generalized linear model results on median of topic proportion and covariate pPNA type on all specimens.*

	Topic	lfc	lfcSE	WTS	pvalue	p.adj
1	Topic_1	-0.90	0.95	-0.95	0.34	0.6095
2	Topic_2	0.33	0.69	0.48	0.63	0.7196
3	Topic_3	-0.56	0.58	-0.96	0.34	0.6095
4	Topic_4	2.06	0.75	2.74	0.01	0.0224
5	Topic_5	2.69	0.67	4.05	0.00	<.0001
6	Topic_6	-0.29	0.80	-0.36	0.72	0.7196
7	Topic_7	-0.38	0.50	-0.77	0.44	0.6095
8	Topic_8	0.64	0.75	0.85	0.39	0.6095
9	Topic_9	0.78	0.78	1.00	0.32	0.6095
10	Topic_10	2.62	0.73	3.61	0.00	0.0017
11	Topic_11	0.33	0.78	0.42	0.67	0.7196

SUPPLEMENTARY TABLE 3

*Generalized linear model results on median of topic proportion and covariate pPNA type on paired-specimens from Asteraceae plants.*

	Topic	lfc	lfcSE	WTS	pvalue	p.adj
1	Topic_1	-10.43	1.51	-6.89	0.00	<.0001
2	Topic_2	1.29	1.36	0.95	0.34	0.6245
3	Topic_3	0.58	1.41	0.41	0.68	0.8343
4	Topic_4	0.97	1.54	0.63	0.53	0.8281
5	Topic_5	-0.15	1.28	-0.12	0.91	0.9482
6	Topic_6	0.09	1.33	0.07	0.95	0.9482
7	Topic_7	0.45	1.03	0.44	0.66	0.8343
8	Topic_8	-1.32	1.03	-1.28	0.20	0.4524
9	Topic_9	-7.81	1.33	-5.89	0.00	<.0001
10	Topic_10	-1.53	1.21	-1.27	0.21	0.4524
11	Topic_11	8.38	1.39	6.04	0.00	<.0001

SUPPLEMENTARY TABLE 4

*Generalized linear model results on median of topic proportion and covariate pPNA type on paired-specimens from non-Asteraceae plants.*

	Topic	lfc	lfcSE	WTS	pvalue	p.adj
1	Topic_1	-0.00	1.57	-0.00	1.00	0.9985
2	Topic_2	-0.24	1.49	-0.16	0.87	0.9985
3	Topic_3	-1.96	1.04	-1.89	0.06	0.1633
4	Topic_4	4.33	1.20	3.59	0.00	0.0036
5	Topic_5	2.24	1.18	1.89	0.06	0.1633
6	Topic_6	0.10	0.88	0.11	0.91	0.9985
7	Topic_7	-0.59	0.89	-0.67	0.51	0.7945
8	Topic_8	0.63	1.51	0.42	0.68	0.9283
9	Topic_9	-2.69	1.67	-1.62	0.11	0.1941
10	Topic_10	2.42	1.50	1.62	0.11	0.1941
11	Topic_11	3.26	1.16	2.82	0.00	0.0267

### Code References

Code to reproduce figures and simulations can be found at <https://pratheepaj.github.io/diffTop/>.