

RESEARCH ARTICLE

MicrobiomeCensus estimates human population sizes from wastewater samples based on inter-individual variability in gut microbiomes

Lin Zhang¹, Likai Chen², Xiaoqian (Annie) Yu³, Claire Duvallet^{4,5}, Siavash Isazadeh^{4,5}, Chengzhen Dai⁶, Shinkyu Park⁶, Katya Frois-Moniz^{4,5}, Fabio Duarte⁶, Carlo Ratti⁶, Eric J. Alm^{4,5,7}, Fangqiong Ling^{1,8,9,10*}

1 Department of Energy, Environmental and Chemical Engineering, Washington University in St. Louis, St. Louis, Missouri, United States of America, **2** Department of Mathematics, Washington University in St. Louis, St. Louis, Missouri, United States of America, **3** Department of Biology, Massachusetts Institute of Technology, Boston, Massachusetts, United States of America, **4** Department of Biological Engineering, Massachusetts Institute of Technology, Boston, Massachusetts, United States of America, **5** Center for Microbiome Informatics and Therapeutics, Massachusetts Institute of Technology, Boston, Massachusetts, United States of America, **6** SENSEable City Lab, Massachusetts Institute of Technology, Boston, Massachusetts, United States of America, **7** Eli and Edythe L. Broad Institute of MIT and Harvard, Boston, Massachusetts, United States of America, **8** Department of Computer Science and Engineering, Washington University in St. Louis, St. Louis, Missouri, United States of America, **9** Division of Biological and Biomedical Sciences, Washington University in St. Louis, St. Louis, Missouri, United States of America, **10** Division of Computational and Data Science, Washington University in St. Louis, St. Louis, Missouri, United States of America

* fangqiong@wustl.edu



OPEN ACCESS

Citation: Zhang L, Chen L, Yu X(Annie), Duvallet C, Isazadeh S, Dai C, et al. (2022) MicrobiomeCensus estimates human population sizes from wastewater samples based on inter-individual variability in gut microbiomes. *PLoS Comput Biol* 18(9): e1010472. <https://doi.org/10.1371/journal.pcbi.1010472>

Editor: Morgan Langille, DAL, CANADA

Received: January 27, 2022

Accepted: August 5, 2022

Published: September 23, 2022

Copyright: © 2022 Zhang et al. This is an open access article distributed under the terms of the [Creative Commons Attribution License](https://creativecommons.org/licenses/by/4.0/), which permits unrestricted use, distribution, and reproduction in any medium, provided the original author and source are credited.

Data Availability Statement: Source code is available through <https://github.com/linglab-washu/population-model>. Sewage metagenomic data is available at National Center for Biotechnology Information Short Read Archive at BioProject PRJNA683921.

Funding: FL acknowledges support from the National Science Foundation Directorate of Engineering Faculty Early Career Development Program (2047470), Alfred P. Sloan Foundation Microbiology of the Built Environment Postdoctoral Fellowship (2015-14164), as well as Washington

Abstract

The metagenome embedded in urban sewage is an attractive new data source to understand urban ecology and assess human health status at scales beyond a single host. Analyzing the viral fraction of wastewater in the ongoing COVID-19 pandemic has shown the potential of wastewater as aggregated samples for early detection, prevalence monitoring, and variant identification of human diseases in large populations. However, using census-based population size instead of real-time population estimates can mislead the interpretation of data acquired from sewage, hindering assessment of representativeness, inference of prevalence, or comparisons of taxa across sites. Here, we show that taxon abundance and sub-species diversity in gut-associated microbiomes are new feature space to utilize for human population estimation. Using a population-scale human gut microbiome sample of over 1,100 people, we found that taxon-abundance distributions of gut-associated multi-person microbiomes exhibited generalizable relationships with respect to human population size. Here and throughout this paper, the human population size is essentially the sample size from the wastewater sample. We present a new algorithm, MicrobiomeCensus, for estimating human population size from sewage samples. MicrobiomeCensus harnesses the inter-individual variability in human gut microbiomes and performs maximum likelihood estimation based on simultaneous deviation of multiple taxa's relative abundances from their population means. MicrobiomeCensus outperformed generic algorithms in data-driven

University in St. Louis McKelvey School of Engineering Faculty Startup Fund. EJA and CR acknowledge funding from the Kuwait Foundation for Advancement of Sciences (KFAS). The funders had no role in study design, data collection and analysis, decision to publish, or preparation of the manuscript.

Competing interests: I have read the journal's policy and the authors of this manuscript have the following competing interests: E.J.A has an equity stake in Biobot Analytics. C. Duvallet is employed by Biobot Analytics.

simulation benchmarks and detected population size differences in field data. New theorems are provided to justify our approach. This research provides a mathematical framework for inferring population sizes in real time from sewage samples, paving the way for more accurate ecological and public health studies utilizing the sewage metagenome.

Author summary

Wastewater-based epidemiology (WBE) is an emerging field that employs sewage as aggregated samples of human populations. This approach is particularly promising for tracking diseases that can spread asymptotically in large populations, such as the COVID-19. As a new type of biological data, sewage has its own unique challenges to utilize. While wastewater samples are usually assumed to represent large populations, the assumption is not guaranteed at locations closer to residences due to stochasticity in toilet flushes; thus, unlike epidemiological experiments collecting data from individuals, sample size, herein the human population size represented by a wastewater sample, is a fundamental yet difficult-to-characterize parameter for sewage samples. Researchers would need to aggregate data from large areas and week-long collection to stabilize data, during which, important spikes in small areas or short time scales may be lost. It also remains challenging to turn viral titers into case prevalences, evaluating representativeness, or comparing measurements across sites/studies.

This study provides a framework to estimate human population size from sewage utilizing human gut-associated microorganisms. Through analysis, we demonstrate that variance of taxon abundances and single-nucleotide polymorphism as two variables that change with population size. We provide a new tool MicrobiomeCensus that performs population size estimation from microbial taxon abundances. MicrobiomeCensus outperforms generic algorithms in terms of computational efficiency while at comparable or better accuracy. Using MicrobiomeCensus, we detected population size differences in sewage samples taken in Cambridge, MA, under two sampling approaches, i.e., “grab” or “composite” sampling. This study provides a framework to utilize individual-level microbiomes to learn from sewage, paving the way to prevalence estimation and improved spatio-temporal resolutions in WBE.

Introduction

The metagenome embedded in urban sewage is an attractive new data source to understand urban ecology and assess human health status at scales beyond a single host [1–3]. Sewage microbiomes are found to share a variety of taxa with human gut microbiomes, where the baseline communities are characterized by a dominance of human-associated commensal organisms from the *Bacteroidetes* and *Firmicutes* phyla [1, 3, 4]. Human viruses like SARS-CoV-2 and polioviruses were detected in sewage samples during the pandemic and silent spreads, respectively, and found to correlate to reported cases, suggesting that sewage samples could be useful for understanding the dynamics in the human-associated symbionts at a population level [5, 6]. Sewage has several advantages as samples of the population's collective symbionts. For instance, sewage samples are naturally aggregated, wastewater infrastructures are highly accessible, and data on human symbionts can be collected without visits to clinics, thus utilizing sewage samples can reduce costs and avoid biases associated with stigma and

accessibility [2, 7]. Consequently, SARS-CoV-2 surveillance utilizing sewage samples are underway globally and incorporated into the U.S. Centers for Disease Control and Prevention surveillance framework [8].

A pressing challenge in utilizing sewage for ecological and public health studies is the lack of methods to directly estimate human population size from sewage. Specifically, virus monitoring at finer spatial granularity, e.g., single university dorms and nursing homes, are informative for guiding contact tracing and protecting populations at higher risk, but real-time population size estimations at such fine granularity are not yet available. For a given area, the census population (*de jure* population) can be larger than the number of people who contributed feces to sewage at a given time (*de facto* population)[9]. Conversely, the *de jure* population can also be smaller than the *de facto* population due to the presence of undocumented individuals [10]. Population proxies that are currently used for monitoring at wastewater-treatment plants, such as the loading of pepper mild mottle viruses, likely have high error at the neighborhood level because of their large variability in human fecal viromes (10^6 - 10^9 virions per gram of dry weight fecal matter)[11]. Consequently, it is difficult to assess the representativeness of a sewage sample, infer the taxon abundance differences across time and space, or interpret errors. Lack of population size information could lead to false negatives in assessing virus eradication, because an absence of biomarkers might be caused by a sewage sample that under-represents the population size. Despite its importance, few studies have explicitly explored ways to estimate real-time human population size from sewage samples independent from census estimates [12].

Macroecological theories of biodiversity may offer clues to decipher and even enumerate the sources of a sewage microbiome. While we are only beginning to view sewage as samples of human symbionts beyond one person, generating multi-host microbiomes resembles a fundamental random additive process. Sizling et al. showed that lognormal species abundance distributions (SADs) can be generated solely from summing the abundances from multiple non-overlapping sub-assemblages to form new assemblages [13]. Likewise, adding multiple sub-assemblages can also give rise to common Species-Area Relationships [13]. For microbial ecosystems, Shoemaker et al. examined the abilities of widely known and successful models of SADs in predicting microbial SADs and found that Poisson Lognormal distributions outperformed other distributions across environmental, engineered, and host-associated microbial communities, highlighting the underpinning role of lognormal processes in shaping microbial diversity [14].

In this study, we conceptualize a sewage microbiome as a multi-person microbiome, where the number of human contributors can vary. We hypothesize that the species abundance distribution in the multi-person microbiome will vary as a function of the human population size, which would arise from summing taxon abundances from multiple hosts analogous to the Central Limit Theorem. We use human gut microbiome data comprising over a thousand human subjects and machine learning algorithms to explore these relationships. Upon discovering a generalizable relationship, we develop MicrobiomeCensus, a nonparametric model that utilizes relative taxon abundances in the microbiome to predict the number of people contributing to a sewage sample. MicrobiomeCensus utilizes a multivariate T statistic to capture the simultaneous deviation of multiple taxa's abundances from their means in a human population and performs maximum likelihood estimation. We provide proof on the validity of our approach. Next, we examine model performance through a simulation benchmark using human microbiome data. Last, we apply our model to data derived from real-world sewage. Our nonparametric method does not assume any underlying distributions of microbial abundances and can make inferences with just the computational power of a laptop computer.

Results

Species abundance distributions of multi-person microbiomes vary by population size

We consider the fraction of microorganisms observed in sewage that are human-associated anaerobes as an “average gut microbiome” sampled from residents of a catchment area. Hence, our task becomes to find the underlying relationship between the number of contributors and the observed microbiome profiles in sewage samples. We define an “ideal sewage mixture” scenario to illustrate our case, where the sewage sample consists only of gut-associated microorganisms and is an even mix of n different individuals’ feces (Fig 1). We denote the gut microbiome profile of an individual as $X_i = (X_{i,1}, X_{i,2}, \dots, X_{i,p})^T$, where each $X_{i,j}$ represents the

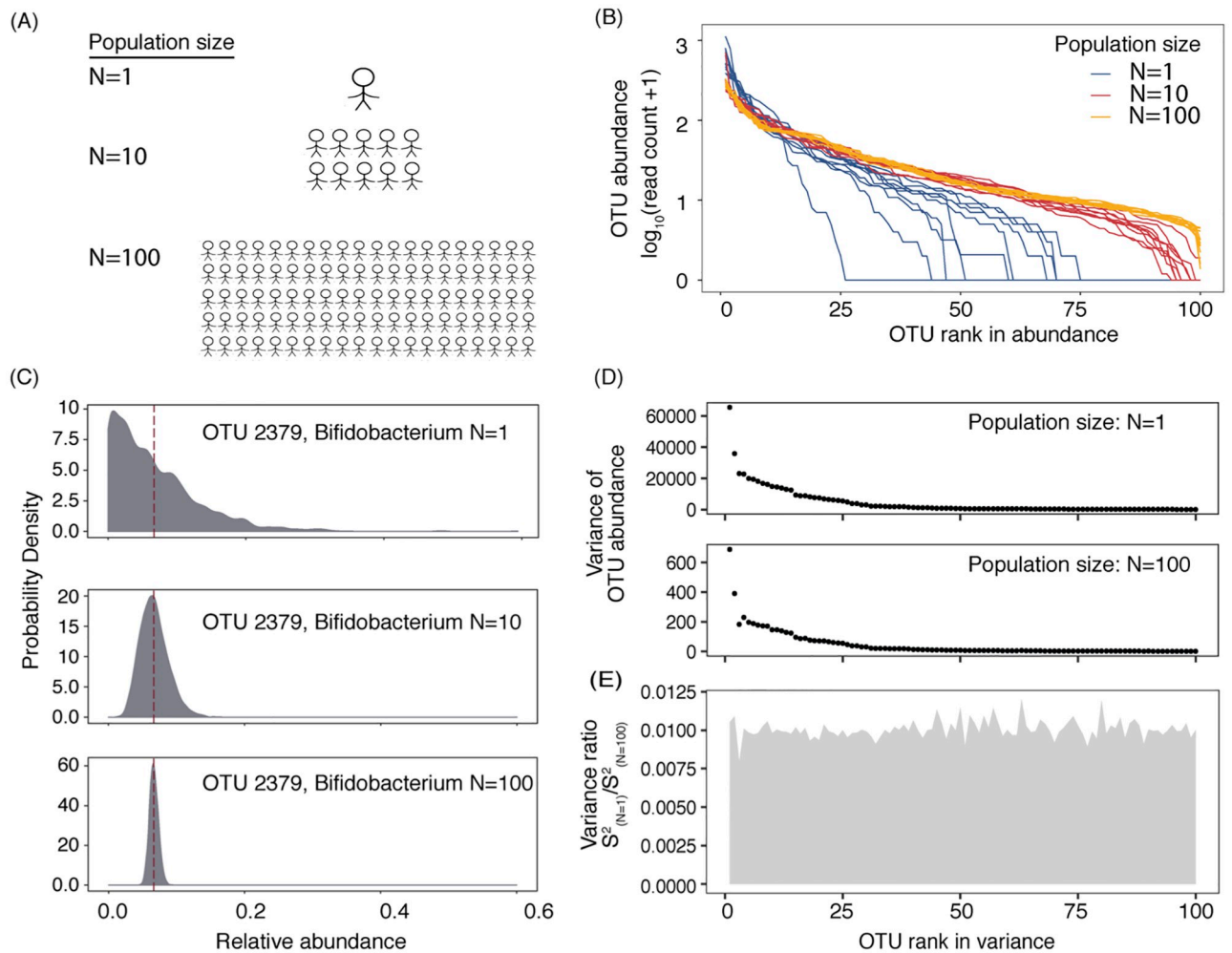


Fig 1. An ideal sewage mixture simulation shows the potential of microbiome taxon abundance profiles as population census information sources. (A) We generated an “ideal sewage mixture” consisting of gut microbiomes from different numbers of people. (B) Ranked abundance curves for gut microbiomes of one person and mixtures of multiple people exhibit different levels of dominance and diversity. Blue lines show the rank abundance curves in stool samples (one person), red lines show 10-person mixtures, and saffron lines show 100-person mixtures. In each scenario, ten examples are shown. All samples were rarefied to the same sequencing depths (4,000 seqs/sample). (C) The probability density function of the relative abundance of one taxon for different population sizes. OTU-2379, a *Bifidobacterium* taxon, was used as an example. Maroon dashed lines indicate the sample means. (D) Multiple taxa’s abundance variances in one-person samples and 100-person samples. The dominant taxa are shown (top100) and are sorted by their ranks in variance. (E) The ratios of the variances of one-person samples and 100-person samples across dominant gut microbial taxa.

<https://doi.org/10.1371/journal.pcbi.1010472.g001>

relative abundance of operational taxonomic unit (OTU) j from individual i . Hence our ideal sewage mixture can be represented as

$$\bar{X}_n = \sum_{i=1}^n X_i/n \quad (1)$$

where vectors $X_1, X_2, \dots, X_n \in \mathbb{R}^p$ are microbiome profiles from individuals $1, \dots, n$. Under the ideal sewage mixture scenario, if we can quantitatively capture the departure of the sewage microbiome profile from the population mean of the human gut microbiomes of people constituting the catchment area, we will be able to estimate the population size.

Using a dataset comprised of 1,100 individuals' gut microbiome taxonomic profiles [15], we created synthetic mixture samples of different numbers of contributors through bootstrapping (Fig 1A). First, examined from an ecological perspective, the shape of the ranked abundance curves of the gut microbiomes differed when the means of multiple individuals were examined: when the number of contributors increased, a normal distribution appeared (Fig 1B). For the single-person microbiomes, log-series and lognormal distributions explained 94% and 93% of the variations in the SADs, respectively, compared with 89% for Poisson log-normal, 87% for Zipf multinomial and 80% for the broken-stick model. Multi-person microbiomes were best predicted by log-series or lognormal models, but as the population increased to over a hundred, the multi-person SADs were best described by only lognormal SADs (S1 Table).

We explored the distributions of the relative abundances of gut bacteria as a function of population size. As expected, the distribution of a taxon's relative abundance changes with population size (Fig 1C). For instance, for OTU-2397, a *Bifidobacterium* taxon, the relative abundance distribution was approximately log-normal when the relative abundance in single-host samples was considered, yet converged to a Normal distribution when mixtures of multiple hosts were considered. Although the means of the distributions of the same taxon under different population sizes were close, the variation in the data changed. A smaller variance was observed when the number of contributors increased (Fig 1D). Notably, different taxa varied in the rates at which their variances decreased with population size (Fig 1E), suggesting that a model that considers multiple features would be useful in predicting the number of contributors.

Classifiers utilizing microbial taxon abundance features alone detects single-person and multi-person microbiomes

Inspired by the distinct shapes of SADs in multi-person gut-microbiomes from those of single-person microbiomes, we set up a classification task using the taxon relative abundances to separate synthetic communities constituting one, ten, and a hundred people. With algorithms of varying complexity, namely Logistic Regression (LR), Support Vector Machine (SVM), and Random Forest (RF) classifier, classification accuracies of 29.6%, 97.2%, and 100% were achieved (Fig 2). Between RF and SVM, RF showed higher sensitivity and specificity in classifying all population groups (S2 Table). This experiment suggests the usefulness of microbiome features in predicting human population counts from mixture samples.

MicrobiomeCensus is a statistical model that estimates population size from microbial taxon abundances

While the classification tasks described above demonstrated the usefulness of taxa's relative abundances in predicting the population size, a complex model like RF provided little

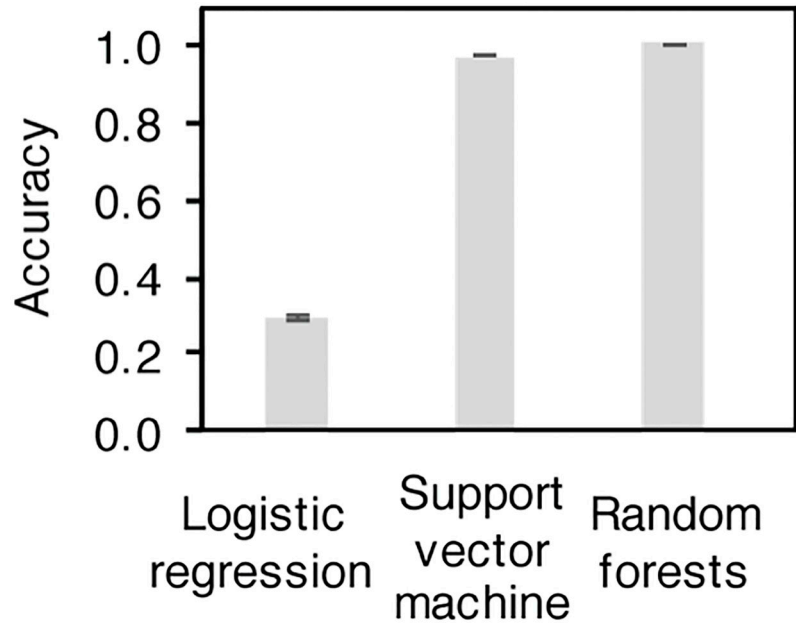


Fig 2. Classifier performance of models utilizing gut microbiome taxon abundances.

<https://doi.org/10.1371/journal.pcbi.1010472.g002>

explanatory power. We then ask, since the variance in the relative abundance of a given taxon decreases with population size, can we devise a statistic that captures the simultaneous deviation of several taxa’s abundances from their means, and estimate population size utilizing the statistic? Further, will this new method perform well despite inter-personal variation in gut microbiomes?

Our new method, MicrobiomeCensus, involves a statistic T_n to capture the simultaneous deviation of multiple taxa’s abundances from their means in relation to the variance of those taxa in the population (Fig 3A). We denote $\Sigma_0 = (\sigma_{ij})_{1 \leq i, j \leq p}$ as the covariance matrix for the individual microbiome profile and let Λ_0 be a diagonal matrix with $\Lambda_0 = \text{diag}(\sigma_{11}^{1/2}, \sigma_{22}^{1/2}, \dots, \sigma_{pp}^{1/2})$. Then the statistic takes form

$$T_n = \|\hat{\Lambda}_0^{-1}(\bar{X}_n - \hat{\mu})\|_2^2, \tag{2}$$

where $\bar{X}_n = \sum_{i=1}^n X_i/n$ denotes the observed microbiome profile in ideal sewage, μ represents the population mean for the catchment area and $\hat{\mu}$ is an estimator, $\hat{\Lambda}_0$ is an estimator of Λ_0 and $\|v\|_2 := (\sum_{i=1}^p v_i^2)^{1/2}$ for any vector $v \in \mathbb{R}^p$. This statistic is enlightened by the classical Hotelling T^2 statistic [16] $\tilde{T}_n = n(\bar{X}_n - \mu)\hat{\Sigma}_0^{-1}(\bar{X}_n - \mu)$, where $\hat{\Sigma}_0$ is the sample covariance matrix, an estimator of Σ_0 . Actually if we assume the covariance matrix is diagonal (no correlations between different taxa), then they are essentially the same statistic in view of $\tilde{T}_n = nT_n$. The reason we replace covariance matrix Σ_0 by its diagonal Λ_0 is because for high dimensional situations, it would be very difficult to estimate the covariance matrix. In cases when $p > n$, the sample covariance matrix is singular and thus \tilde{T}_n is not even well defined. Studies accommodating the Hotelling T^2 type statistic into the high-dimensional situation can be found, for example, in Bai and Saranadasa [17], Chen and Qi [18], Xu et al [19], etc. Our proposed statistic can handle the high dimensional cases as well, since the diagonal entities Λ_0 can be well

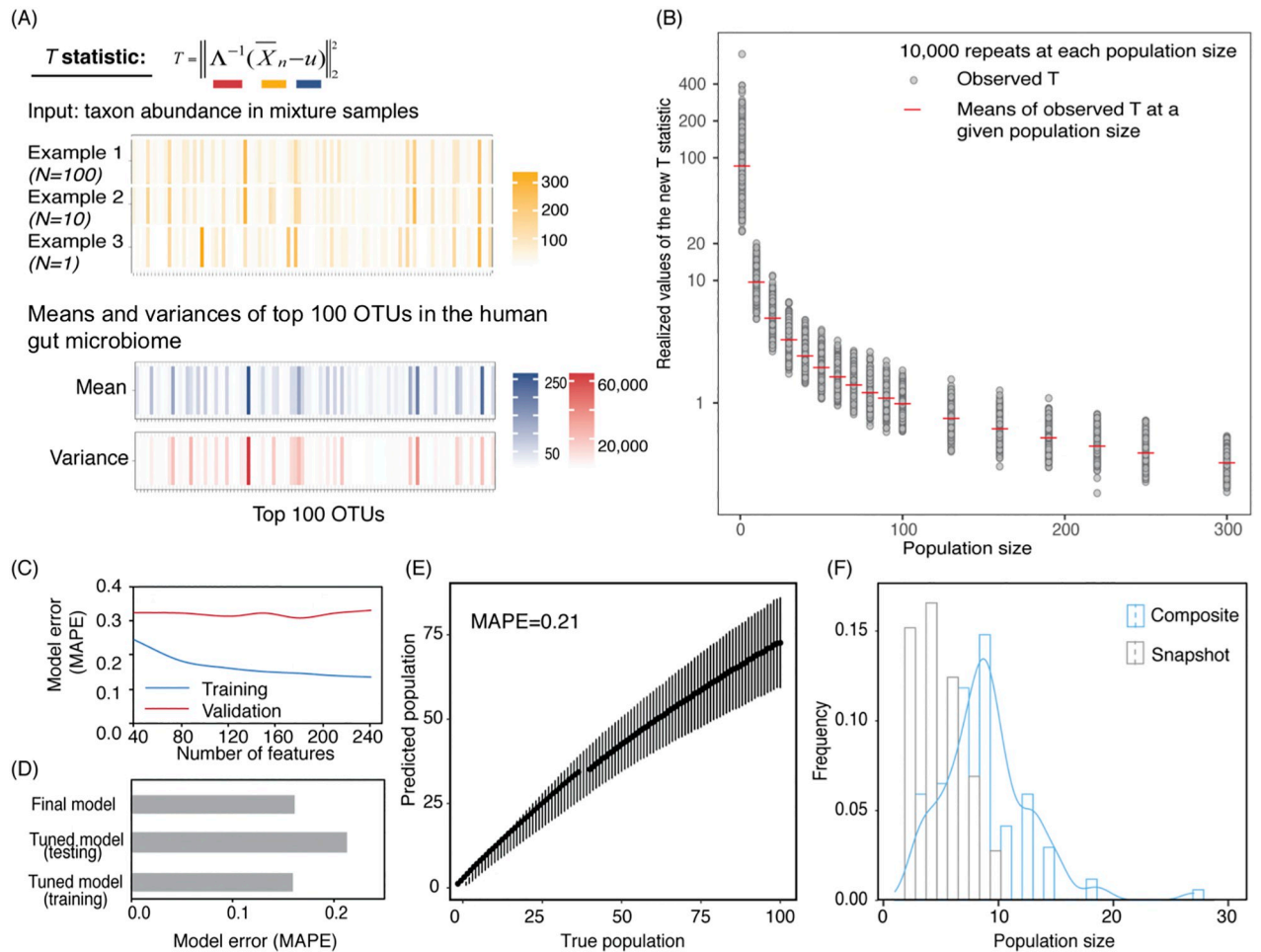


Fig 3. MicrobiomeCensus statistic definition, model training, validation, and application. (A) Example of computing the T statistic. (B) Simulation results for T with different population sizes. Grey points are simulation results. Red bars are means of 10,000 repeats performed for each population size. (C) Model training and tuning. We built the MicrobiomeCensus model using our T statistic and a maximum likelihood procedure. The training set consisted of 10,000 samples for population sizes ranging from 1–300, and 50% of the data were used to train and validate the model. Training and validation errors from different feature subsets are shown. Training errors are shown as blue lines, and validation errors are shown as red lines. (D) Model performance on simulation benchmark. After training and validation, the model utilized the top 120 abundant features. Model performance was tested on synthetic data generated from 550 different subjects not previously seen by the model. The training set consisted of 10,000 samples with population sizes from 1–300, and the testing set consisted of 10,000 repeats at the evaluated population sizes. The training error, testing error, and the error of the final model are shown. (E) Model performance evaluated using a testing set. Black solid dots indicate the means of the predicted values, and error bars indicate the standard deviations of the predicted values. (F) Application of the microbiome population model in sewage. Seventy-six composite samples (blue) were taken from three manholes on the MIT campus, and each sample was taken over 3 hours during the morning peak water usage hours. Twenty-five snapshot samples (grey) were taken using a peristaltic pump for 5 minutes at 1-hour intervals throughout a day.

<https://doi.org/10.1371/journal.pcbi.1010472.g003>

estimated even when p is large. And we extend its application beyond the problem of the significance of the multivariate means.

In developing this new method, we utilize the variance change by population, but without any *priori* assumption about the gut bacterium species taxon abundance distributions and the covariance between species. Our analysis showed that the statistic T_n changed monotonically with increasing population size, indicating the promise of a population estimation model (Fig 3B).

Leveraging our statistic T_n , we constructed an asymptotic maximum likelihood estimator to estimate size of the sample without the information of each individual, that is, we do not

observe X_1, \dots, X_n but only their mixture $\bar{X}_n = \sum_{i=1}^n X_i/n$. Here, the parameter of interest is the population size n , the test statistic is T_n , and a point estimate is made by maximizing the estimated likelihood of T_n with respect to n . We performed training and validation using 50% of the human microbiome data and held out the rest of the data for testing. Our model achieved a training error as low as 13% (mean absolute percentage error, MAPE) when up to 250 features are included. The model's training performance increased when more features were included, yet the validation error did not profoundly change with an increasing number of features (Fig 3C). Cross-validation on top 40, 80, 120, 150, 180, 210, and 240 abundant OTUs led to validation errors at 32.4%, 32.3%, 31.4%, 32.3%, 30.9%, 32.3%, and 33.1%, respectively. Upon training and validation, we chose the top 120 OTUs and tested the performance of the tuned model on a test set held out during training/validation. The model's MAPE was 21% (Fig 3D and 3E, testing errors at each population size evaluated are provided in S3 Table), indicating that our model generalized well across different hosts. We then used all data and tuned hyperparameter to acquire a final model. When applying the final model on the same testing data, our model achieved a testing error of 16.2% (Fig 3D).

It is worth noting that in this algorithm, for each size n , we only need to estimate the sampling distribution of the statistic T_n once. Hence it is not time-consuming regardless of the true population size. We also note that an RF regression model could not be trained in a reasonable time on the same dataset, even with high-performance computing (Methods). Our model performed remarkably better than a ten-fold cross-validated RF regression model utilizing a reduced dataset, which gave an MAPE of 32%, while the training time for our model was only a fraction of that of the RF regression model (S1 Fig).

MicrobiomeCensus detects human population size differences in sewage samples

With the newly developed population model, we set out to apply our model to sewage samples. Ideally, we would like to apply the model to samples generated from a fully controlled experiment with known human hosts contributing at a given time, yet such an experiment presents logistic challenges. Instead, we applied our model to sewage samples taken using one of two methods, either a snapshot (grab sample) sample taken from the sewage stream over 5 minutes, or an accumulative (composite sample) taken at a constant rate over 3 hours during morning peak human defecation [20] (S2 Fig). We hypothesized that the composite samples would represent more people than snapshot samples. Taking grab samples, we sampled at 1-hr intervals at one manhole ($n = 25$); using the accumulative method, we sampled at three campus buildings (classroom, dormitory, and family housing) multiple times over three months ($n = 76$). To remove sequences possibly contributed by the water, we applied a taxonomic filter to retain families associated with the gut microbiome and normalized the species abundance by the retained sequencing reads (Methods, S4 Table). We applied our final model to the sewage data set. Our model estimated 1–9 people's waste was captured by the snapshot samples (mean = 3, s.d.=3), and 3–27 people were represented by the composite samples (mean = 9, s.d.=7), where the composite samples represented significantly more people ($p < 0.0001$) (Fig 3F). The hypothesis that composite samples represent more people is well supported by our model results.

Sub-species diversity in sewage samples reflects adding microbiomes from multiple people

Independent from our MicrobiomeCensus model, we found that certain human gut-associated species were frequently detected in sewage samples by using shotgun metagenomics, e.g.,

Bacteroides vulgatus, *Prevotella copri*, and *Eubacterium rectale*. Further, their sub-species diversity, as indicated by nucleotide diversity and the number of polymorphic sites in house-keeping genes, was dramatically higher in sewage samples than in the gut microbiomes of individual human subjects (Fig 4A–4F and Text A in S1 Text).

To examine the effect of increasing population size on sub-species genetic variation in representative gut-associated microbial species, we simulated aggregate human gut samples using a sample without replacement procedure and computed the nucleotide diversity and numbers of polymorphic sites for the aggregate samples at different population sizes. This resulted in single nucleotide variant (SNV) profiles from 64 species. Our simulation showed increases in both nucleotide diversity and the number of polymorphic sites as more human gut samples were aggregated (Fig 4G and 4H). For instance, the nucleotide diversity and number of polymorphic sites in *Eubacterium rectale* increased from 0.029 (s.d. 0.026) to 0.149 (s.d. 0.002) and 64 (s.d. 54.33) to 1274 (s.d. 18.41), respectively, when the population size increased from 1 to 300. Further, the number of polymorphic sites strongly correlated with the population size (Pearson correlation coefficient >0.8) in 49 of the 64 species (S5 Table), suggesting the potential that the SNV profiles of a wide range of gut species could be developed into feature space for population size estimation. Our simulation further shows that the number of polymorphic sites increased with population size more slowly than nucleotide diversity, indicating its potential to reflect more subtle changes in population size (Fig 4G and 4H). Despite the need for further model developments, the analysis here shows the potential of the sub-species diversity of gut anaerobes as a feature space to be developed into a population size estimation model, independent from the taxon abundance-based model described here.

Discussion

The MicrobiomeCensus method we present here can, in theory, estimate the population size contributing to a sewage sample from the taxon abundance of multiple human gut microbiome taxa, using our T statistic and associated maximum likelihood estimation and application procedures. While the model is trained to perform accurate population estimation on a neighborhood scale, we expect the population range it can estimate to expand with increasing training gut microbiome data availability. We propose the MicrobiomeCensus model as a tool to drive further developments in quantitative sewage-based epidemiology. We have provided mathematical proof of the validity of our approach.

MicrobiomeCensus showed excellent performance in our simulation benchmark. In particular, the study subjects that we utilized in the training and testing sets are random samples out of 1,100 men and women across a wide range of age without any stratification, hence the model's testing performance indicates its generalizability. Our study is founded on the observations that healthy gut microbiomes are resilient, with inter-individual variability outweighing variability within individuals over time [21–23]. There are caveats to our approach; potentially, diets and regional effects on human microbiome composition could introduce noises to the prediction [24, 25]. In applications to sewage, future studies on water matrix effects should be performed to understand and further account for noises from the sewage collection network. It should be mentioned that while our model is trained on microbiome data, it is not limited to microbiome features because we did not impose any assumptions on the distribution of features. Other features, e.g., crAssphage titers, may also be incorporated once individual level data at a large population size become available to allow model validation.

Utilizing sewage to understand population-level dynamics of human symbionts presents a new scenario of sampling meta-communities. The gut microbiomes of humans can be viewed as local communities, and gut microbiomes of people living in a neighborhood

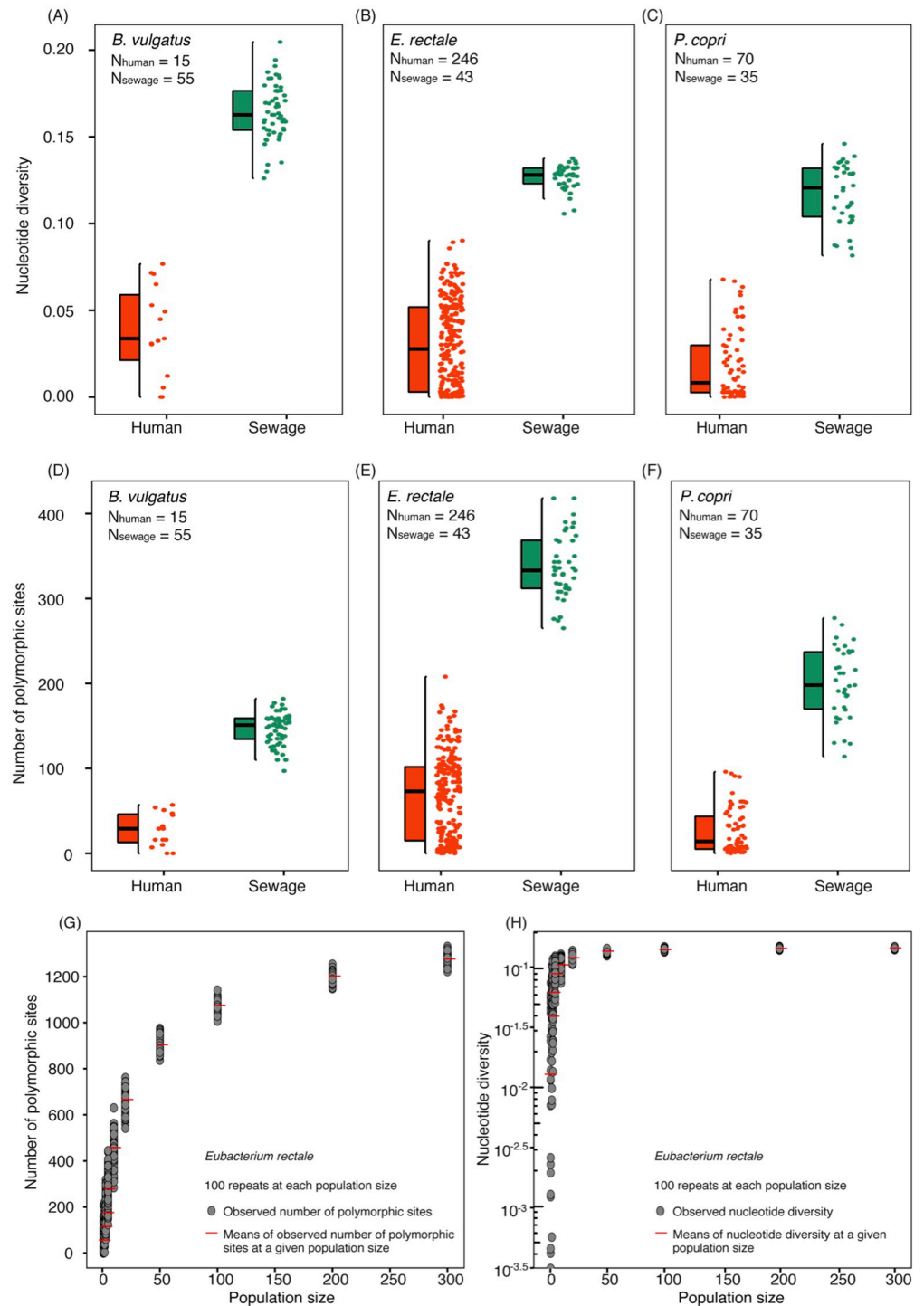


Fig 4. Sub-species diversity in gut-associated bacterial species as a potential marker for human population size. (A-F) Comparison of sub-species diversity of gut-associated bacteria in human gut microbiome samples (LifelinesDeep) and MIT sewage samples. Nucleotide diversity and numbers of polymorphic sites were computed from ten phylogenetic marker genes. (G) and (H) Simulation results showing intra-species diversity in response to increasing population size, as represented by the number of polymorphic sites (G) and nucleotide diversity (H).

<https://doi.org/10.1371/journal.pcbi.1010472.g004>

could be viewed as a kind of regional meta-communities, because these communities are linked by dispersal that can take place among people connected by social networks and through a shared built environment. The meta-community framework is considered to provide useful new conceptual tools to understand the largely unexplained inter-personal variability in gut microbiomes, with expansions of the theory to consider biotic interactions suggested by Miller, Svanbäck, and Bohannan [26]. In considering a sample of meta-communities, Leibold and Chase asked provocatively “what is a community?” and observed that the definition of a community is usually “user-defined and could be context-dependent” –“one community ecologist might explore the patterns of coexistence and species interactions among species within a delimited area, the other might ask the same question but define a community that encompasses more area and thus types of species, as well as different degrees of movements and heterogeneity patterns” [27]. The ambiguity between samples of meta-communities and local communities is particularly challenging for samples of microbial communities, because dispersal boundaries are difficult to delineate. Despite the conceptual importance, empirical methods that explicitly test whether a microbiome sample is a sample of a meta-community or a local community have not been available. MicrobiomeCensus directly distinguishes samples of meta-communities and local communities by enumerating the number of hosts contributing to a microbiome. While MicrobiomeCensus is trained on gut microbiome data, the procedure may have wide applications in other microbial ecosystems.

There are several limitations of this study. First, our approach requires a sizable set of individual gut microbiome data to generate the empirical distributions employed in maximum likelihood estimation. Second, accurate estimations of population sizes will depend on good estimations of the population means and variances of relative abundances in individual gut microbiomes. Future studies generating data of individual-level gut microbiome of residents in specific cities, as well as model training based on those data, will help bridge the applications to sewage samples at the corresponding areas. Third, our current application procedures require minimal decays of gut-associated bacteria, thus are suitable for applications at catchments near residential buildings (e.g., sentinel monitoring). Better understandings of the growth/death dynamics of human gut-associated microbial taxa in the sewer environment and feature selection leveraging that information will help bridge the applications to larger catchments.

In response to the COVID-19 pandemic now affecting the human population globally, sewage-based virus monitoring is underway [28]. Our analysis calls for attention to the denominator used in normalizing the biomarker measurements. While in practice, loading-based population proxies such as the copy numbers of pepper mild mottle viruses are used to normalize data generated from sewage, such proxies would likely have high error at the neighborhood level because of their variability in human fecal viromes (10^6 - 10^9 virions per gram of dry weight in fecal matters)[11], while they likely have reasonable performance when the population size is sufficiently large and the means of biomarker loadings converge under the Central Limit Theorem. Thus, the relationships between sewage measurements and true viral prevalence in small populations are hard to establish despite the need for sentinel population studies. Our model has immediate application in detecting false negatives, because it alerts us to the possibility that an absence of biomarkers might be caused by a sewage sample that underrepresents a population. With further developments incorporating local training data, the model can potentially generate a denominator that can help turn biomarker measurements into estimates of prevalence and enable the application of epidemiology models at finer spatio-temporal resolutions.

Methods

Rationales

If the distribution of X_i is known, then the distribution T_n is known and one can easily use maximum likelihood estimator (MLE) to estimate the human population size. Here the human population size is essentially the sample size n from the wastewater sample. However for generality, we do not want to impose any specific distribution assumptions on taxon abundance distributions, thus, we need to rely on asymptotic results to estimate the distribution of our statistic. Unlike the univariate case where the asymptotic distribution of the statistic T_n can be simply derived by central limit theorem, we are dealing with a much more challenging situation due to high dimensionality.

In the following, we shall firstly introduce some notations and assumptions that will be needed for the theorems. Then, in Theorem 1, we derive the Gaussian approximation for the test statistic T_n , which implies that we can use simulated Gaussian vectors to approximate the distribution of our statistic. To apply this approximation, one needs to further estimate the covariance matrix which is highly non-trivial due to high dimensionality. To get around this difficulty, we further propose a sub-sampling approach. Theorem 2 provides the theoretical foundation for this sub-sampling scheme.

Notation

Recall that vectors X_1, X_2, \dots, X_n are microbiome profiles from individuals $1, 2, \dots, n$. Assume that they are independent and identically distributed (i.i.d) random vectors in \mathbb{R}^p with mean $\mu \in \mathbb{R}^p$ and variance $\text{Var}(X_i) = \Sigma_0 \in \mathbb{R}^{p \times p}$. Denote $\Sigma_0 = (\sigma_{ij})_{1 \leq i, j \leq p}$ and $\sigma_i = \sigma_{ii}^{1/2}, 1 \leq i \leq p$. Let the diagonal matrix $\Lambda_0 = \text{diag}(\sigma_i, 1 \leq i \leq p)$.

To construct the Gaussian approximation, we shall firstly work with cases when both Λ_0 and μ are given, that is, consider the statistic

$$T_n^\circ = \|\Lambda_0^{-1}(\bar{X}_n - \mu)\|_2^2.$$

We will later extend all the results to cases when those parameters are unknown. For notation's simplicity, consider the normalized version of X_i :

$$Y_i = \Lambda_0^{-1}(X_i - \mu).$$

Then $T_n^\circ = \|\bar{Y}_n\|_2^2$, where $\bar{Y}_n = \sum_{i=1}^n Y_i/n$, and the covariance matrix Σ_Y of Y_i is the correlation matrix of X_i , with expression $\Sigma_Y = \Lambda_0^{-1}\Sigma_0\Lambda_0^{-1}$.

We need the following condition on Y_i for the main theorem.

ASSUMPTION 1 Let $s > 2$. Assume

$$K_s^s = \mathbb{E} \left| \frac{\|Y_1\|_2^2 - p}{\|\Sigma_Y\|_F} \right|^s < \infty, \quad \text{and} \quad D_s^s = \mathbb{E} \left| \frac{Y_1^T Y_2}{\|\Sigma_Y\|_F} \right|^s < \infty, \tag{3}$$

where $\|\cdot\|_F$ is the matrix Frobenius norm.

REMARK 1 Above conditions naturally hold if $Y_{1i}, 1 \leq i \leq p$, are independent and $\max_{1 \leq i \leq p} (\mathbb{E}|Y_{1i}|^s)^{1/s} < \infty$. Actually under this setting, $\Sigma_Y = I_p$ and thus $\|\Sigma_Y\|_F = p^{1/2}$. By Berkholder [29]'s inequality,

$$\mathbb{E} \left(\left| \|Y_1\|_2^2 - p \right|^s \right) \leq (s-1)^s \left(\sum_{i=1}^p (\mathbb{E}|Y_{1i}^2 - 1|^s)^{2/s} \right)^{s/2} \leq c_1 p^{s/2},$$

where $c_1 > 0$ is independent of p . This justifies K_s part in condition (3). Again by Berkholder [29]’s inequality,

$$\mathbb{E}(|Y_1^\top Y_2|^s) \leq (s - 1)^s \left(\sum_{i=1}^p (\mathbb{E}|Y_{1i} Y_{2i}|^s)^{2/s} \right)^{s/2} \leq c_2 p^{s/2},$$

where $c_2 > 0$ is independent of p . And thus D_s part in condition (3) holds.

Now we are ready to introduce the first asymptotic result. The following theorem essentially states that under certain regularity conditions, the distribution of our test statistic can be approximated by the distribution of some function of a Gaussian vector.

THEOREM 1 (THEOREM 2.2 IN XU ET AL. [19]) *Assume Assumption 1 holds with some $s > 2$, also assume*

$$K_2^2/n + K_s^s/n^{s-1} + D_s^s/n^{s/2-1} \rightarrow 0. \tag{4}$$

Then for $Z \sim N(0, \Sigma_Y)$, we have

$$\sup_{t \in \mathbb{R}} \left| \mathbb{P}(nT_n^\infty \leq t) - \mathbb{P}(\|Z\|_2^2 \leq t) \right| \rightarrow 0. \tag{5}$$

If K_s and D_s in (3) are bounded, then we can get a crude bound $O(n^{-s/(10 + 4s)})$ for (5). It is worth noticing that under settings in Remark 1, condition (4) holds. Based on the above theorem, if we can estimate the covariance matrix Σ_Y , then $\|Z\|_2^2$ can be generated and used for approximation of our statistic. However the estimation of Σ_Y is difficult due to high dimensionality unless some additional conditions are imposed on the covariance matrix. To get around this difficulty, we shall consider a bootstrap approach. The main idea is that for each n , we randomly generate n samples from the reference data $X_i, 1 \leq i \leq n_0$, and construct some generated statistic. Using the empirical distribution of the generated statistic to approximate the distribution of our statistic. In the following, we shall provide the theoretical justification for this procedure.

For some integer $J > 0$, let A_1, A_2, \dots, A_J be i.i.d uniformly sampled from the class $\mathcal{A}_n = \{A : A \subset \{1, 2, \dots, n_0\}, |A| = n\}$. Assume the sampling process are independent from our data $(X_i)_i$. Then for each $1 \leq k \leq J$, the set $\{X_i, i \in A_k\}$ is of size n and can be used to construct one test statistic $n\|\Lambda_0^{-1}(\bar{X}_{A_j} - \bar{X}_{n_0})\|_2^2$, where $\bar{X}_{A_j} = \sum_{i \in A_j} X_i / |A_j|$. After repeating this procedure J times, we would have J realizations of our test statistic which can be used to construct the empirical distribution $\hat{F}_n(t)$:

$$\hat{F}_n(t) = J^{-1} \sum_{j=1}^J \mathbf{1}_{\{n\|\Lambda_0^{-1}(\bar{X}_{A_j} - \bar{X}_{n_0})\|_2^2 \leq t(1-n/n_0)\}}.$$

We shall later show that this empirical distribution can be adopted for approximating the distribution of the target statistic T_n^∞ . Following result can be derived based on Theorem 3.5 in Xu et al.[19] and Theorem 1.

THEOREM 2 *Assume conditions in Theorem 1 hold, and moreover assume that $n \rightarrow \infty, n = o(n_0)$ and (4) holds. Then for $J \rightarrow \infty$, we have*

$$\sup_{t \in \mathbb{R}} |\hat{F}_n(t) - \mathbb{P}(nT_n^\infty \leq t)| \rightarrow 0. \tag{6}$$

Theorem 2 implies that we can use the empirical distribution generated from our sub-samples to estimate the distribution of our target. As mentioned previously, so far, we are assuming that we know the value of Λ_0 which is not realistic in applications. Therefore we need to

further estimate Λ_0 , that is, we need to estimate the standard deviation for each coordinate. This can be easily accomplished by considering the estimator

$$\hat{\Lambda}_{0,j}^2 = \sum_{i=1}^{n_0} (X_{i,j} - \bar{X}_{n_0,j})^2 / n_0,$$

where $X_{i,j}$ represents the j th coordinate of X_i , and $\Lambda_{0,j}$ is the j th diagonal entity of Λ_0 and n_0 is the size of the reference data. Also one can use the average of the reference data to replace μ . If $X_{i,j}$ has heavy tail, we can also consider a robust m-estimator for Σ_0 and μ , see for example, Catoni [30].

Bootstrap procedure

Below we describe the bootstrap procedure we use to approximate the distribution of T_n for different census counts. Recall that X_1, \dots, X_n represent arrays of taxon relative abundances in the gut microbiome of human subject $1, \dots, n$, and T_n is defined in (Eq 2).

- Step 1. Estimate the population mean $\hat{\mu}$, and the diagonal matrix $\hat{\Lambda}_0$, from a reference sample human gut microbiome data.
- Step 2. For each census count n , generate X_1^*, \dots, X_n^* from the reference data to compute T_1^* .
- Step 3. Repeat Step 2 B times (herein 10,000 times) to get T_1^*, \dots, T_B^* .
- Step 4. Obtain the density function \hat{f}_n of T_n based on T_1^*, \dots, T_B^* using a Gaussian kernel.
- Step 5. Repeat Steps 2–4 for all the census counts $n = 1, 2, \dots, N$ considered, herein $N = 300$. It should be noted that per Theorem 2 we require bootstrap sample size n much smaller than total reference sample size n_0 , thus up to 300-person samples were simulated here because the gut microbiome reference dataset we utilized consisted of a total of 1,100 people. The range can be expanded if a larger dataset is available.

Maximum likelihood estimation

We use a maximum likelihood estimation (MLE) procedure to achieve point estimates of the population size from a new mixture sample, W . The MLE procedure firstly computes T_0 by replacing the sample average by W , that is

$$T_0 = \|\hat{\Lambda}_0^{-1}(W - \hat{\mu})\|_2^2. \quad (7)$$

And then computes the likelihoods that T_0 was drawn from population sizes from 1 to N , respectively, using the estimated distributions generated from the bootstrap procedures described above. Next, the population size \hat{n} that yields the highest likelihood is chosen.

Confidence interval

Due to Theorem 1, the asymptotic distribution of nT_n is the same as $\|Z\|_2^2$ and is therefore independent of the parameter n . Hence for any confidence level $1 - \alpha$, we can firstly estimate the $1 - \alpha$ quantile of nT_n based on the simulated data described above, denoted as $\hat{q}_{1-\alpha}$. Then for any new mixture W and the corresponding T_0 as in (8), our $1 - \alpha$ confidence interval is $[1, \hat{q}_{1-\alpha}/T_0]$.

Model training, validation, and testing

We synthesized a mixed data set from a gut microbiome dataset of 1,135 healthy human hosts from the Lifelines Deep study [15], which was the largest single-center study of population-level human microbiome variations from a single sequencing center at the time of this study. The data set consisted of 661 women and 474 men. We considered OTUs defined by 99% similarity of partial ribosomal RNA gene sequences (Methods of OTU clustering are described in detail in Text B in [S1 Text](#)). After quality filtering, we retained 1,100 samples that had more than 4,000 sequencing reads/sample. We split the entire dataset approximately in half, using 550 subjects to generate the training/validation set and the other 550 subjects to generate the test set. We then used the aforementioned ideal sewage mixture approach to generate synthetic populations of up to 300 individuals, which is the relevant range for population estimation in upstream sewage. The training error was computed using the entire training data set. Five repeated holdout validations using a 50–50 split in the training set were performed to tune the hyperparameter for feature selection. The training and cross-validation errors were evaluated at integers from 1 to 100, using the error definition:

$$\delta = \left| \frac{N_{\text{predicted}} - N_{\text{actual}}}{N_{\text{actual}}} \right| \times 100\% \quad (8)$$

and the model's performance across all the population sizes was characterized by the mean absolute percentage error (MAPE):

$$\text{MAPE} = \frac{1}{n} \sum_{n=1}^{100} \left| \frac{N_{\text{predicted}} - N_{\text{actual}}}{N_{\text{actual}}} \right| \times 100\%. \quad (9)$$

We chose to use MAPE because for a problem of population estimation, the error relative to the true value is important to consider for performance evaluation. After training and validation, the hyperparameter (in this case, the top k abundant OTUs) that yielded the best performance in the validation step was used in the model. The tuned model was then tested on the test set. Our synthetic sewage microbiome approach captured the actual microbiome variation among individual hosts and demonstrated the model's generalizability.

Human gut microbiome 16S rRNA amplicon data source

The single-person and multi-person microbiome data were drawn from a gut microbiome dataset of 1,135 healthy human hosts from the Lifelines Deep study [15], which was the largest study of population-level human microbiome variations from a single sequencing center at the time of this study. The data set consisted of 661 women and 474 men. We considered operational taxonomic units defined by 99% similarity of partial ribosomal RNA gene sequences. After quality filtering, we retained 1,100 samples that had more than 4,000 sequencing reads/sample. The rarefaction depth was chosen to balance sample size and sequencing depth.

16S rRNA gene amplicon sequencing data analysis

Operational taxonomic units defined at 99% sequencing similarity were generated from the combined dataset by first denoising the samples with DADA2 [31], and then clustering the outputted exact sequence variants with the q2-vsearch plugin of QIIME2 [32]. Taxonomic assignments were performed using a multinomial naïve Bayes classifier against SILVA 132 [33, 34]. All 16S rRNA gene amplicon analyses were performed in the QIIME2 platform (QIIME 2019.10) [35].

Species abundance distribution

We examined the relationships between the performances of several widely used SAD models and the number of contributors (population size) to a multi-person microbiome. Multi-person microbiomes were generated by sampling N individuals from the quality-filtered gut microbiome 16S rRNA dataset and summing the abundances of the same taxa. At each population size, 10,000 repeats were performed. The repeats were chosen according to the constraints of computational efficiency. The SADs evaluated included the Lognormal, Poisson Lognormal, Broken-stick, Log series and the Zipf model, which were shown to have varied successes in predicting microbial SADs [14]. We examined the fit using a rank-by-rank approach as previously described by Shoemaker et al. [14]. First, maximum-likelihood coefficients for each of the SADs described above were estimated using the R package `sads` [36]. Next, SADs were predicted using each model, and tabulated as RADs. Then, we used a least-squares regression to assess the relationship between the performance of the predicted SADs against the observations and recorded the coefficient of determination (R-squared). Last, R-squared values from model fits of each SAD model were summarized as the means, and the models that resulted in the highest R-squared values for each simulated community were recorded.

Field data

We conducted a field sampling campaign, collecting sewage samples daily at manholes near three buildings (two dormitory buildings and one office building) on the campus of Massachusetts Institute of Technology. Seventy-six sewage samples were collected through a continuous peristaltic pump sampler operated at the morning peak (7–10 a.m. near the dormitory buildings and 8–11 a.m. near the office building) at 4 mL/min for 3 hours. This composite sampling approach was intended to capture the morning water usage peak. Wastewater was filtered through sterile 0.22- μ m mixed cellulose filters to collect microbial biomass. Environmental DNA was extracted with a Qiagen PowerSoil DNA extraction kit according to the manufacturer's protocol. The DNA was amplified for the V4 region of the 16S rRNA gene and sequenced in a Miseq paired-end format at the MIT BioMicro Center, according to a previously published protocol [37]. Included as a comparison are a set of snapshot sewage samples taken using a peristaltic pump sampler at 100 mL/min for 5 minutes over a day (10 a.m. on Wednesday April 8, 2015, to 9 a.m. on Thursday April 9, 2015). The sampling methods for snapshot samples are described in detail by Matus et al. [4]

Application to sewage data

The 16S rRNA gene amplicon sequencing data from the field sewage samples were trimmed to the same region, 16S V4 (534–786) with the LifeLines Deep data using Cutadapt 1.12 [38]. Forward reads were trimmed to 175bps, and reverse reads were first trimmed to 175bps and then further trimmed to 155bps during quality screening. We created a taxonomic filter based on the composition of the gut microbiome data set, which consisted of the abundant family-level taxa that accounted for 99% of the sequencing reads in the human gut microbiome data set, and excluded those that might have an ecological niche in tap water (*Enterobacteriaceae* and *Burkholderiaceae*). This exclusion resulted in 25 bacterial families and one archaeal family in our taxonomic filter, including *Lachnospiraceae*, *Ruminococcaceae*, *Bifidobacteriaceae*, *Erysipelotrichaceae*, *Bacteroidaceae*, and others (S4 Table). We applied our taxonomic filter to the sewage sequencing data, which retained 73.9% of the sequencing reads. This retention rate is consistent with our previous report of the human microbiome fraction in residential sewage samples [4]. We then normalized the relative abundance of taxa against the remaining sequencing reads in each sample. Welch's two-sample t-tests were performed to retain the

OTUs whose means did not differ significantly from the human microbiome data set ($p > 0.05$).

Deployment of generic machine learning models

Logistic regression, support vector machine, and random forest classifiers were employed to perform the classification task for population sizes of 1, 10, and 100. Model training, cross-validation, and testing were performed using the R Caret platform with the default setting [39]. For the support vector machine, the radial basis function kernel was employed. Ten-fold cross-validation and five repeats were performed for all the models considered. Model performance was evaluated using accuracy, sensitivity, and specificity. Based on the classifier performance, the RF regression model was used for comparison with our new model's performance. Initially, we trained the model using the same training data set used in training our maximum likelihood model, however, the computation was infeasible, even with a 36-thread, 3TB-memory computing cluster. We then introduced gaps in the population size range, using populations from the vector (1, 5, 10, 20, 30, 40, 50, 60, 70, 80, 90, 100, 110, 120, 150, 180, 240, 300)^T while maintaining the same sample size at each population size (10,000 samples). The training was performed in R Caret, using 10-fold cross-validation. Ten variables were randomly sampled as candidates at each split, $mtry = 10$. The performance was evaluated using the same testing set that was used to evaluate the maximum likelihood model.

Supporting information

S1 Fig. Performance of a cross-validated Random Forest regression model utilizing microbiome features to predict population size. Black solid dots indicate the means of the predicted values, and error bars indicate the standard deviations of the predicted values. (TIFF)

S2 Fig. Sewage sampling schematics. Sewage samples were collected from a manhole in upstream areas close to the sources. Snapshot samples were collected at a high pump speed for a short sampling time. Composite samples were collected at a low pump speed and for a long sampling time. Sewage flow is simplified to illustrate the sampling. Color bars indicate different sections of the flow, coming from different contributors. (TIFF)

S3 Fig. *Bacteroides vulgatus* sub-species diversity in human gut microbiome samples (LifelinesDeep) and MIT sewage samples, as revealed by ten phylogenetic marker genes. (TIFF)

S4 Fig. *Eubacterium rectale* sub-species diversity in human gut microbiome samples (LifelinesDeep) and MIT sewage samples, as revealed by ten phylogenetic marker genes. (TIFF)

S5 Fig. *Prevotella copri* sub-species diversity in human gut microbiome samples (LifelinesDeep) and MIT sewage samples, as revealed by ten phylogenetic marker genes. (TIFF)

S1 Text. Supplementary Texts. (DOCX)

S1 Table. Comparison of the performance of SAD models for single-person and multi-person gut microbial communities. (XLSX)

S2 Table. Classifier performance for population size prediction from microbiome features.
(XLSX)

S3 Table. MicrobiomeCensus model performance in simulation benchmark (1000 repeats at each population size).
(XLSX)

S4 Table. Family-level taxa for taxonomic filter in model application.
(XLSX)

S5 Table. Correlation of sub-species diversity with population size revealed by simulation.
(XLSX)

S6 Table. Sub-species diversity of gut-associated bacteria in a population-level human microbiome dataset (Lifeline Deep).
(XLSX)

S7 Table. Sub-species diversity of bacteria prevalent in sewage samples.
(XLSX)

Acknowledgments

We thank Cho C. Yiu, David E Hingston, and Joseph S. Monteiro from the MIT Facilities department for assistance with sewage sampling. We thank Noriko Endo, Sean Gibbons, Tami Lieberman, Xiaofang Jiang and Shijie Zhao for valuable discussions. We thank Mariana Matus and Newsha Ghaeli for acquisition and access of 24-hr sewage time series data.

Author Contributions

Conceptualization: Eric J. Alm, Fangqiong Ling.

Data curation: Lin Zhang, Claire Duvallet, Siavash Isazadeh, Chengzhen Dai, Shinkyu Park, Fangqiong Ling.

Formal analysis: Lin Zhang, Likai Chen, Xiaoqian (Annie) Yu, Fangqiong Ling.

Funding acquisition: Carlo Ratti, Eric J. Alm, Fangqiong Ling.

Investigation: Lin Zhang, Likai Chen, Xiaoqian (Annie) Yu, Siavash Isazadeh, Chengzhen Dai, Shinkyu Park, Fangqiong Ling.

Methodology: Likai Chen, Siavash Isazadeh, Shinkyu Park, Eric J. Alm, Fangqiong Ling.

Project administration: Katya Frois-Moniz, Fabio Duarte.

Resources: Carlo Ratti, Eric J. Alm, Fangqiong Ling.

Software: Lin Zhang, Likai Chen, Fangqiong Ling.

Validation: Lin Zhang, Likai Chen, Siavash Isazadeh, Chengzhen Dai, Shinkyu Park, Fangqiong Ling.

Visualization: Lin Zhang, Fangqiong Ling.

Writing – original draft: Lin Zhang, Likai Chen, Xiaoqian (Annie) Yu, Fangqiong Ling.

Writing – review & editing: Lin Zhang, Likai Chen, Xiaoqian (Annie) Yu, Claire Duvallet, Siavash Isazadeh, Chengzhen Dai, Shinkyu Park, Katya Frois-Moniz, Fabio Duarte, Carlo Ratti, Eric J. Alm, Fangqiong Ling.

References

1. Maritz JM, Ten Eyck TA, Elizabeth Alter S, Carlton JM. Patterns of protist diversity associated with raw sewage in New York City. *ISME J*. 2019; 13(11):2750–2763. <https://doi.org/10.1038/s41396-019-0467-z> PMID: 31289345
2. Berchenko Y, Manor Y, Freedman LS, Kaliner E, Grotto I, Mendelson E, et al. Estimation of polio infection prevalence from environmental surveillance data. *Sci Transl Med*. 2017; 9(383). <https://doi.org/10.1126/scitranslmed.aaf6786> PMID: 28356510
3. Newton RJ, McLellan SL, Dila DK, Vineis JH, Morrison HG, Eren AM, et al. Sewage Reflects the Microbiomes of Human Populations. *mBio*. 2015; 6(2). <https://doi.org/10.1128/mBio.02574-14> PMID: 25714718
4. Matus M, Duvallet C, Soule MK, Kearney SM, Endo N, Ghaeli N, et al. 24-hour multi-omics analysis of residential sewage reflects human activity and informs public health. *bioRxiv*. 2019; p. 728022.
5. Medema G, Heijnen L, Elsinga G, Italiaander R, Brouwer A. Presence of SARS-Coronavirus-2 RNA in Sewage and Correlation with Reported COVID-19 Prevalence in the Early Stage of the Epidemic in The Netherlands. *Environ Sci Technol Lett*.
6. Manor Y, Shulman LM, Kaliner E, Hindiyeh M, Ram D, Sofer D, et al. Intensified environmental surveillance supporting the response to wild poliovirus type 1 silent circulation in Israel, 2013. *Eurosurveillance*. 2014; 19(7):20708. <https://doi.org/10.2807/1560-7917.ES2014.19.7.20708> PMID: 24576473
7. Murakami M, Hata A, Honda R, Watanabe T. Letter to the Editor: Wastewater-Based Epidemiology Can Overcome Representativeness and Stigma Issues Related to COVID-19. *Environ Sci Technol*. 2020; 54(9):5311. <https://doi.org/10.1021/acs.est.0c02172> PMID: 32323978
8. CDC. National Wastewater Surveillance System; 2020. Available from: <https://www.cdc.gov/coronavirus/2019-ncov/cases-updates/wastewater-surveillance.html>.
9. Daughton CG. Real-time estimation of small-area populations with human biomarkers in sewage. *Sci Total Environ*. 2012; 414:6–21. <https://doi.org/10.1016/j.scitotenv.2011.11.015> PMID: 22137478
10. Fazel-Zarandi MM, Feinstein JS, Kaplan EH. The number of undocumented immigrants in the United States: Estimates based on demographic modeling with data from 1990 to 2016. *PLoS One*. 2018; 13(9):e0201193. <https://doi.org/10.1371/journal.pone.0201193> PMID: 30240392
11. Zhang T, Breitbart M, Lee WH, Run JQ, Wei CL, Soh SWL, et al. RNA Viral Community in Human Feces: Prevalence of Plant Pathogenic Viruses. *PLOS Biol*. 2005; 4(1):e3. <https://doi.org/10.1371/journal.pbio.0040003>
12. Yang Z, Xu G, Reboud J, Kasprzyk-Hordern B, Cooper JM. Monitoring Genetic Population Biomarkers for Wastewater-Based Epidemiology. *Anal Chem*. 2017; 89(18):9941–9945. <https://doi.org/10.1021/acs.analchem.7b02257> PMID: 28814081
13. Šizling AL, Storch D, Šizlingová E, Reif J, Gaston KJ. Species abundance distribution results from a spatial analogy of central limit theorem. *Proc Natl Acad Sci USA*. 2009; 106(16):6691–6695. <https://doi.org/10.1073/pnas.0810096106> PMID: 19346488
14. Shoemaker WR, Locey KJ, Lennon JT. A macroecological theory of microbial biodiversity. *Nat Ecol Evol*. 2017; 1(5):1–6. <https://doi.org/10.1038/s41559-017-0107> PMID: 28812691
15. Zhernakova A, Kurilshikov A, Bonder MJ, Tigchelaar EF, Schirmer M, Vatanen T, et al. Population-based metagenomics analysis reveals markers for gut microbiome composition and diversity. *Science*. 2016; 352(6285):565. <https://doi.org/10.1126/science.aad3369> PMID: 27126040
16. Hotelling H, Frankel LR. The Transformation of Statistics to Simplify their Distribution. *Ann Math Stat*. 1938; 9(2):87–96. <https://doi.org/10.1214/aoms/1177732330>
17. Bai Z, Saranadasa H. Effect of high dimension: by an example of a two sample problem. *Stat Sin*. 1996; p. 311–329.
18. Chen SX, Qin YL, et al. A two-sample test for high-dimensional data with applications to gene-set testing. *Ann Stat*. 2010; 38(2):808–835. <https://doi.org/10.1214/09-AOS716>
19. Xu M, Zhang D, Wu WB. L2 Asymptotics for High-Dimensional Data. *arXiv preprint arXiv:14057244*. 2014;.
20. Heaton KW, Radvan J, Cripps H, Mountford RA, Braddon FE, Hughes AO. Defecation frequency and timing, and stool form in the general population: a prospective study. *Gut*. 1992; 33(6):818–824. <https://doi.org/10.1136/gut.33.6.818> PMID: 1624166
21. David LA, Materna AC, Friedman J, Campos-Baptista MI, Blackburn MC, Perrotta A, et al. Host lifestyle affects human microbiota on daily timescales. *Genome Biol*. 2014; 15(7):R89. <https://doi.org/10.1186/gb-2014-15-7-r89> PMID: 25146375
22. Lozupone CA, Stombaugh JI, Gordon JI, Jansson JK, Knight R. Diversity, stability and resilience of the human gut microbiota. *Nature*. 2012; 489(7415):220–230. <https://doi.org/10.1038/nature11550> PMID: 22972295

23. Mehta RS, Abu-Ali GS, Drew DA, Lloyd-Price J, Subramanian A, Lochhead P, et al. Stability of the human faecal microbiome in a cohort of adult men. *Nat Microbiol.* 2018; 3(3):347–355. <https://doi.org/10.1038/s41564-017-0096-0> PMID: 29335554
24. Johnson AJ, Vangay P, Al-Ghalith GA, Hillmann BM, Ward TL, Shields-Cutler RR, et al. Daily Sampling Reveals Personalized Diet-Microbiome Associations in Humans. *Cell Host Microbe.* 2019; 25(6):789–802.e5. <https://doi.org/10.1016/j.chom.2019.05.005> PMID: 31194939
25. He Y, Wu W, Zheng HM, Li P, McDonald D, Sheng HF, et al. Regional Variation Limits Applications of Healthy Gut Microbiome Reference Ranges and Disease Models. *Nat Med.* 2018; 24(10):1532–1535. <https://doi.org/10.1038/s41591-018-0164-x> PMID: 30150716
26. Miller ET, Svanbäck R, Bohannan BJM. Microbiomes as Metacommunities: Understanding Host-Associated Microbes through Metacommunity Ecology. *Trends Ecol Evol.* 2018; 33(12):926–935. <https://doi.org/10.1016/j.tree.2018.09.002> PMID: 30266244
27. Leibold MA, Chase JM. *Metacommunity Ecology*, Volume 59. Princeton University Press; 2017.
28. Bivins A, North D, Ahmad A, Ahmed W, Alm E, Been F, et al. Wastewater-Based Epidemiology: Global Collaborative to Maximize Contributions in the Fight Against COVID-19. *Environ Sci Technol.* 2020; 54(13):7754–7757. <https://doi.org/10.1021/acs.est.0c02388> PMID: 32530639
29. Burkholder DL. Sharp inequalities for martingales and stochastic integrals. *Astérisque.* 1988; 157(158):75–94.
30. Catoni O. Challenging the empirical mean and empirical variance: a deviation study. In: *Annales de l'IHP Probabilités et statistiques.* vol. 48; 2012. p. 1148–1185.
31. Callahan BJ, McMurdie PJ, Rosen MJ, Han AW, Johnson AJA, Holmes SP. DADA2: high-resolution sample inference from Illumina amplicon data. *Nat Methods.* 2016; 13(7):581–583. <https://doi.org/10.1038/nmeth.3869> PMID: 27214047
32. Rognes T, Flouri T, Nichols B, Quince C, Mahé F. VSEARCH: A Versatile Open Source Tool for Metagenomics. *PeerJ.* 2016; 4:e2584. <https://doi.org/10.7717/peerj.2584> PMID: 27781170
33. Quast C, Pruesse E, Yilmaz P, Gerken J, Schweer T, Yarza P, et al. The SILVA Ribosomal RNA Gene Database Project: Improved Data Processing and Web-Based Tools. *Nucleic Acids Res.* 2013; 41(Database issue):D590–596. <https://doi.org/10.1093/nar/gks1219> PMID: 23193283
34. Bokulich NA, Kaehler BD, Rideout JR, Dillon M, Bolyen E, Knight R, et al. Optimizing Taxonomic Classification of Marker-Gene Amplicon Sequences with QIIME 2's Q2-Feature-Classifer Plugin. *Microbiome.* 2018; 6(1):90. <https://doi.org/10.1186/s40168-018-0470-z> PMID: 29773078
35. Bolyen E, Rideout JR, Dillon MR, Bokulich NA, Abnet CC, Al-Ghalith GA, et al. Reproducible, interactive, scalable and extensible microbiome data science using QIIME 2. *Nat biotechnol.* 2019; 37(8):852–857. <https://doi.org/10.1038/s41587-019-0209-9> PMID: 31341288
36. Prado PI, Miranda MD, Chalom A. *sads: Maximum Likelihood Models for Species Abundance Distributions*; 2018. Available from: <https://CRAN.R-project.org/package=sads>.
37. Preheim SP, Perrotta AR, Friedman J, Smilie C, Brito I, Smith MB, et al. Computational methods for high-throughput comparative analyses of natural microbial communities. *Methods Enzymol.* 2013; 531:353–370. <https://doi.org/10.1016/B978-0-12-407863-5.00018-6> PMID: 24060130
38. Martin M. Cutadapt Removes Adapter Sequences from High-Throughput Sequencing Reads. *EMBnet J.* 2011; 17(1):10–12. <https://doi.org/10.14806/ej.17.1.200>
39. Kuhn M. Building Predictive Models in R Using the caret Package. *J Stat Softw.* 2008; 028(i05).