

# 2-Way $k$ -Means as a Model for Microbiome Samples

Weston J. Jackson<sup>1,\*</sup>, Ipsita Agarwal<sup>2</sup> and Itsik Pe'er<sup>1</sup>

<sup>1</sup>Department of Computer Science, Columbia University, New York, NY 10027 <sup>2</sup>Department of Biological Sciences, Columbia University, New York, NY 10027

## ABSTRACT

**Motivation:** Microbiome sequencing allows defining clusters of samples with shared composition. However, this paradigm poorly accounts for samples whose composition is a mixture of cluster-characterizing ones, and therefore lie in-between them in cluster space. This paper addresses unsupervised learning of 2-way clusters. It defines a mixture model that allows 2-way cluster assignment and describes a variant of generalized  $k$ -means for learning such a model. We demonstrate applicability to microbial 16S rDNA sequencing data from the Human Vaginal Microbiome Project. **Contact:** [wjj2106@columbia.edu](mailto:wjj2106@columbia.edu)

## 1 INTRODUCTION

Microbiome analysis [1] by sequencing of ubiquitous genes, most commonly 16S rRNA, is a standard, cost-effective way to characterize the composition of a microbial sample. Standard analysis tools facilitate quantifying the fraction of sequence reads from each bacterial species in a sample [2]. Interpretation of composition vectors across a collection of samples typically relies on dimensionality reduction followed by clustering in the lower-dimensionality space [3]. This allows identification of functionally-meaningful subsets of samples with characteristic microbiota. The Human Microbiome Project [4] and derivatives such as The Human Vaginal Microbiome Project [5] have collected and thus analyzed large numbers of samples towards elucidating the structure and composition of microbiota across physiological and pathological states.

Similar to variation in microbial genomes across different human individuals, variants along the nuclear genomes have been summarized by a small number of dimensions [6]. However, in contrast to analyses of microbiome samples, those of inherited genetic variation standardly assume and observe samples to be spread across a continuum in the reduced space, rather than be clustered [7]. Samples in between clusters are interpreted as originating from intermediate locales along a geographic cline [8], or as representing different levels of a mixture between cluster-specific populations.

In this paper, we formally tackle the problem of clustering while allowing elements to belong to two clusters. Specifically, we will describe in detail a model for clustering in  $\mathbb{R}^d$ . We construct a model that generalizes  $k$ -means clustering by allowing data points to be assigned to a point in the space along the line between two assigned clusters [9]. Each cluster is still modeled as a Gaussian with uniform, spherical covariances, the key difference is the presence of a parameter  $u \in [0, 1]$  for each 2-way-assigned data point  $x_i$ , which determines the proportional assignment of  $x_i$  between its two cluster representatives. We first describe the 2-way model's

inputs, parameters, and outputs. We then give the objective function, an algorithmic description, and a series of performance metrics. Next, we evaluate the performance on simulated data, describing benchmarks for optimal performance. Finally, we apply the model to real data of 16S rDNA sequencing from 1500 mid-vaginal bacterial samples by the Vaginal Human Microbiome Project.

## 2 METHODS

### 2-Way $k$ -means

The model characterizes a mixture where points are each sampled either from a  $k$ -mixture of uniform, spherical gaussian distributions, or from pairwise weighted averages of these Gaussians.

Formally, we describe a generative model for a set  $X$  of data points  $\{x_i\}_1^n \in \mathbb{R}^d$ . The model involves  $k \in \mathbb{Z}^+$  clusters. The  $j$ -th cluster is parametrized by its mean  $\mu_j \in \mathbb{R}^d$ . To simulate  $x_i$ , the model first chooses a pair of cluster indices  $j, j'$  along with a weighting  $u_i \in [0, 1]$ .  $x_i$  is drawn from a Gaussian distribution whose parameters are  $u_i$ -weighted averages of two representative clusters. Specifically,  $x_i \sim N(x_i; \tilde{\mu}_{ijj'}, \Sigma)$  such that  $\tilde{\mu}_{ijj'} = u_i\mu_j + (1-u_i)\mu_{j'}$ , and  $\Sigma \in \mathbb{R}^{d \times d}$  is the given uniform, spherical covariance matrix.

The inference problem involves the inputs of data  $X$  and number of clusters  $k$ , seeking output of the generative model parameters, i.e. the vectors of assignments  $C = (c_1, \dots, c_n)$  and weights  $U = (u_1, \dots, u_n)$ .

### Generalized $k$ -Means

Given input  $x_1 \dots x_n \in \mathbb{R}^d$  and cardinality  $k \in \mathbb{N}$ ,  $k$ -means traditionally provides us with the following objective:

$$\sum_{i=1}^n \min_{j \in [k]} \|x_j - c_j\|_2^2$$

where  $c_1 \dots c_j$  are the cluster representatives. The  $k$ -means objective can be generalized as the following:

$$\min_{C, \Phi} \sum_{i=1}^n \|x_i - C\phi_i\|_2^2 \quad (1)$$

where  $\Phi = [\phi_1 | \phi_2 | \dots | \phi_n] \in \{0, 1\}^{k \times n}$  are the cluster assignments and  $C = [c_1 | c_2 | \dots | c_k] \in \{0, 1\}^{d \times k}$  are the cluster representatives.

A common generalization of  $k$ -means is to permit each  $\phi_i$  to have  $s$  non-zero entries (in our case, we set  $s = 2$ ). An algorithm for this generalized objective is simply to hold  $C$  fixed while performing sparse regression on  $\Phi$ , then hold  $\Phi$  fixed and use Ordinary Least Squares (OLS) to find  $C$ .

In our case, because we only allow points  $x_i$  to lie uniformly between two cluster representatives, the two non-zero entries in a given  $\phi_i$  are restricted to some  $u_i \in [0, 1]$  and  $1 - u_i \in [0, 1]$ . Our problem is instead the following:

$$\min_{C, \Phi} \sum_{i=1}^n \|x_i - C\phi_i\|_2^2$$

subject to:

$$\|\phi_i\|_0 \leq 2, \|\phi_i\|_1 = 1, \phi_i \geq 0$$

## 2-Way $k$ -Means Algorithm

Our goal is to find a non-negative 2-sparse solution for each  $\phi_i$ . To do so, we can minimize over all  $\binom{k}{2}$  cluster representative possibilities. This 2-sparse solution gives us indices  $(j, j')$  which correspond with the two cluster representatives. This corresponds with the following objective:

$$\min_{u_i \in \mathbb{R}, c_j, c_{j'} \in C} \|x_i - (u_i c_j + u_i(1 - c_{j'}))\|_2^2$$

subject to:  $u_i \in [0, 1]$

For a given  $c_j$  and  $c_{j'}$ , minimizing with respect to  $u_{ijj'}$  reveals a global minimum at:

$$\frac{(c_{j'} - c_j)^T (c_{j'} - x_i)}{\|c_{j'} - c_j\|_2^2}$$

After minimizing with respect to  $u_{ijj'}$ , we project  $u_{ijj'}$  to the region  $[0, 1]$ . We set  $u_{ijj'} = 0$  if the minimizer is less than 0, and set  $u_{ijj'} = 1$  if the minimizer is greater than 1. This allows us to achieve the minimum value of  $u_i$  over the domain  $[0, 1]$  for  $x_i$ .

After minimizing the assignment  $\Phi$ , we then use *OLS* to pick optimal  $C$  as specified before. Formally, *OLS* produces a vector  $c_i^T$  that minimizes the squared residual error between an input matrix  $\Phi$  and vector  $x_i^T$ .

$$\min_{c_i^T} \|x_i^T - \Phi c_i^T\|_2^2$$

Taking the gradient and setting equal to zero yields the following formula:

$$c_i^T = (\Phi\Phi^T)^{-1}\Phi x_i$$

Thus, we perform *OLS* for all vectors  $c_i^T$  at once with matrix multiplication:

$$C^T = (\Phi\Phi^T)^{-1}\Phi X$$

Thus, this gives us representatives  $c_1 \dots c_k$  that minimize the residual error between the cluster representatives and data points subject to  $\Phi$ . We then alternate this process for  $r$  rounds until convergence.

## Performance Metrics

We use the 2-way  $k$ -means objective as a performance metric in measuring the accuracy of model in unsupervised examples.

$$obj(X, k, r) = \min_{C, \Phi} \sum_{i=1}^n \|x_i - C\phi_i\|_2^2 \quad (2)$$

Where  $\Phi$  has at most two non-zero entries with values  $u_i \in [0, 1]$  and  $1 - u_i \in [0, 1]$ .

Additionally, we also use four different error-rates to measure the accuracy of 2-Way  $k$ -means on test cases. Let  $c_i^*$ ,  $\mu_j^*$ , and  $u_{ijj'}^*$  be the ground truth instance parameters i.e., respectively, true 2-way cluster assignment of  $x_i$ , center of cluster  $j$ , and 2-way weighting for  $x_i$  between clusters  $(j, j')$ .

$err_{f(x)}$ : Defines the 0-1 Error rate for 2-way cluster assignment:

$$err_{f(x)} = \frac{\sum_{i=1}^n \mathbb{1}_{c_i \neq c_i^*}}{n}$$

$err_{\mu}$ : Defines the squared deviation from optimal  $\mu^*$ :

$$err_{\mu} = \sum_{\forall j} \|\mu_j - \mu_j^*\|_2$$

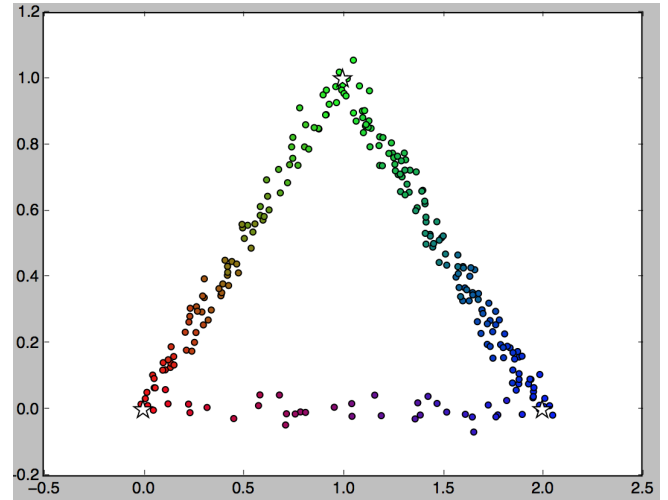
$err_u$ : Defines the squared deviation from optimal  $u_{ijj'}^*$ . WLOG, we assume  $u_{ijj'} = \max(u, 1 - u)$ , where  $u$  is the variable drawn from  $[0, 1]$ :

$$err_u = \frac{\sum_{i=1}^n \|u_{ijj'}^* - u_{ijj'}\|_2}{n}$$

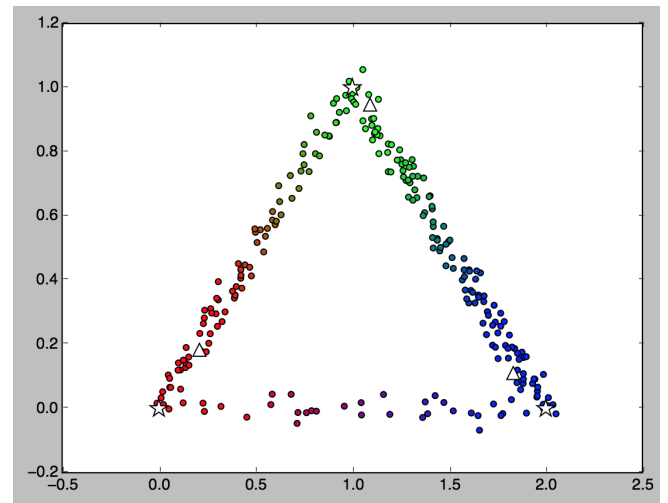
## 3 RESULTS

### Example Run for 2-Way $k$ -Means

We find it illuminating to demonstrate the performance of 2-way  $k$ -means vs vanilla  $k$ -means on a cartoon example.



**Fig. 1.**  $n=500$  simulated data points. The white stars are cluster centers for three simulated clusters (red cluster bottom left, green cluster top, blue cluster bottom right). Points are colored as a linear combination of the clusters they lie between (according to  $u$ ).



**Fig. 2.**  $n=500$  simulated data points after  $k$ -means. The white stars are cluster centers for three simulated clusters (red cluster bottom left, green cluster top, blue cluster bottom right). White triangles are cluster centers determined by  $k$ -means. Colors are  $u$  values determined by  $k$ -means.

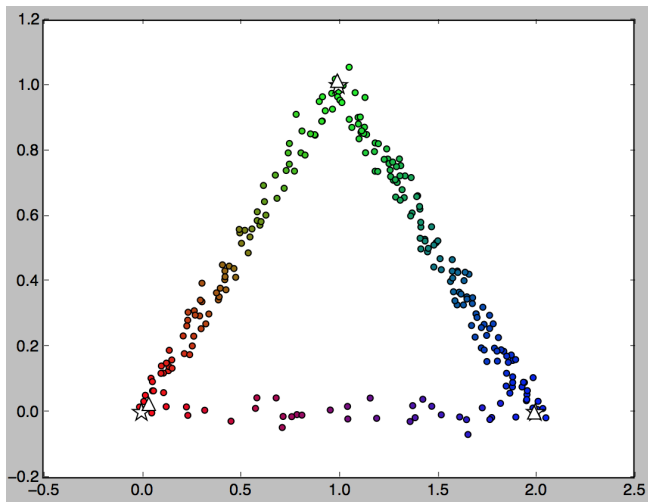
In Figure 1, we simulated  $n = 500$  data points in  $\mathbb{R}^2$  from three clusters, with respective means  $\mu_1 = [0, 0]$ ,  $\mu_2 = [1, 1]$ ,  $\mu_3 = [2, 0]$  and covariances matrices  $\Sigma = 0.001I$ . Data points are drawn into pairwise clusters by choosing two cluster representatives without replacement from the following prior probabilities:

$$\begin{aligned} P(c_1) &= 0.2 \\ P(c_2) &= 0.5 \\ P(c_3) &= 0.3 \end{aligned}$$

We initialize the cluster representatives with vanilla  $k$ -means. Vanilla  $k$ -means achieves the results in Figure 2. Statistics for vanilla  $k$ -means is given below:

$$\begin{aligned} obj &: 80.5795 \\ err_{f(x)} &: 0.406 \\ err_{\mu} &: 0.3677 \\ err_u &: 0.2422 \end{aligned}$$

$k$ -means predicts the cluster assignments of  $\approx 40\%$  points incorrectly (assuming many belong to just one cluster), and also skews cluster means toward the middle of the graph. 2-way  $k$ -means, however, avoids these problems. After 10 rounds of 2-way  $k$ -means, we achieve the results in Figure 3.

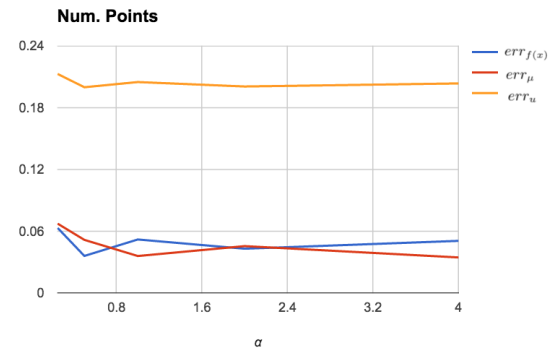


**Fig. 3.**  $n=500$  simulated data points after 10 rounds of 2-way  $k$ -means. White stars are cluster centers determined by 2-way  $k$ -means (10 rounds). Colors are  $u$  values determined by 2-way  $k$ -means (10 rounds).

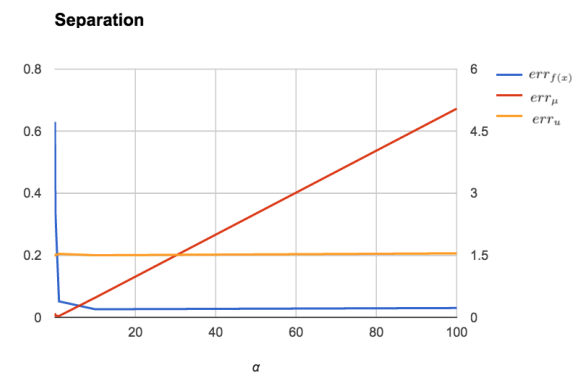
$$\begin{aligned} obj &: 13.237 \\ err_{f(x)} &: 0.038 \\ err_{\mu} &: 0.04335 \\ err_u &: 0.2031 \end{aligned}$$

For every statistic, the results are clearly an improvement on standard  $k$ -means. The 3% error-rate on cluster assignment still exists because 2-way  $k$ -means still allows some points behind cluster representatives to belong to one cluster.

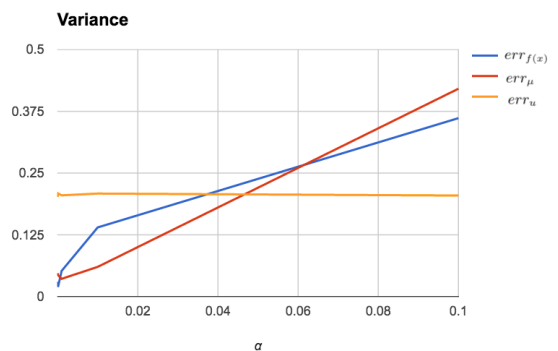
## Benchmarks



**Fig. 4.** Error rates when  $n = \alpha 500$ , and cluster priors and centers are fixed.



**Fig. 5.** Error rates as a function of the euclidean distances of  $\mu$ , where  $\mu_1 = \alpha[0, 0]$ ,  $\mu_2 = \alpha[1, 1]$ ,  $\mu_3 = \alpha[2, 0]$



**Fig. 6.** Error rates as a function of cluster variance  $\Sigma$ , where  $\Sigma = \alpha[[0, 0.0001], [0, 0.001]]$

### Sparsity (Avg. of 10 trials, 10 rounds each)

Our sparsity test was conducted by keeping cluster prior probabilities and cluster centers  $\mu$  constant while varying the number of data points (ratio of  $\alpha$  means  $n = 500\alpha$ ). From Figure 4, we see that the algorithm performs consistently well under a variety of conditions, but too few data points can hurt performance to an extent.

### Cluster Separation (Avg. of 10 trials, 10 rounds each)

We test the error rate as a function of the euclidean distances of  $\mu$  (ratio of  $\alpha$  means  $\mu_1 = \alpha[0, 0]$ ,  $\mu_2 = \alpha[1, 1]$ ,  $\mu_3 = \alpha[2, 0]$ ). From the results in Figure 5, we can see that a certain threshold is required for proper performance of the algorithm. This makes sense, as when  $\alpha = 0.01$ , the clusters are almost on top of each other, and difficult to distinguish. Additionally, as the cluster centers are moved farther apart, the  $\ell_2$  norm between the cluster representative determined by the algorithm and the actual cluster representative increases (but this is to be expected).

### Variance (Avg. of 10 trials, 10 rounds each)

We increase the variance of the clusters while fixing cluster prior probabilities, data points, and cluster centers (ratio of  $\alpha$  means  $\Sigma = \alpha[[0, 0.0001], [0, 0.001]]$ ). From the results in Figure 6, we can see that large variance hurts proper performance of the algorithm. Analogous to with cluster separation, as when  $\alpha = 100$ , the clusters are too close to distinguish.

### Real Data

Publicly available sequence data for the Human Microbiome Project (HMP) study SRP002462, described as Metagenomic sequencing of 16S rDNA from vaginal and related samples from clinical and twin subjects, was downloaded from the NCBI SRA database [10]. The downloaded sets of data correspond to two separate submissions: SRA169809 (1608/1608 samples were downloaded), and SRA273234 (34/133 samples were downloaded), for a total of 1642 samples.

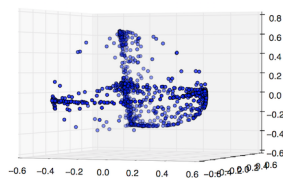


Fig. 7. 1500 data points graphed after PCoA.

The SFF files were processed and cleaned using the microbial community analysis software mothur [11], based on a standard protocol developed for 454 sequence data processing and quality control [12]. The dissimilarities between the samples were calculated using the Clayton-Yue dissimilarity measure. The data

was subsampled to 5000 sequences per sample (this step results in dropping out 136 samples that had less than 5000 reads in total) 500 times to produce the distance file, which was used to calculate principal coordinates. Figure 7 shows the graph of  $\sim 1500$  data points after PCoA. After implementing the 2-Way  $k$ -means algorithm [13], we initialized with  $k$ -means,  $k = 5$ , and ran 2-Way  $k$ -means for 5 rounds on the data.

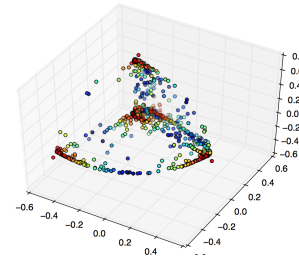


Fig. 8.  $\sim 1500$  data points after 5 rounds of 2-Way  $k$ -means. The red points are closer to cluster centers while the blue ones are viewed as between cluster centers. The cluster representatives are the red points which are slightly outside the clusters (compensating for the slightly non-linear arches between clusters).

Unfortunately, the non-linear arches between the clusters, pushed the cluster representatives slightly outside the clusters. Nonetheless, the algorithm was still an improvement over  $k$ -means. We note that after  $k$ -means, the 2-way objective had a value of 108.0 while our 2-way  $k$ -means algorithm converged on an objective of  $\approx 51.0$  after 5 rounds. Additionally, the algorithm gives us a characterization of the samples lying between two clusters. The results can be seen in the Figure 8.

### Discussion

We first get the most abundant operational taxonomic unit (OTU) in each sample (down to the genus level), and the closest cluster assignment for each sample. We use this to observe which OTUs are most common to each cluster. We can find the closest sample to each data point by simply taking the  $argmax(u)$  for each data point  $x_i$ .

From Table 1, we see that four of the five clusters have a unique most abundant OTU, while cluster  $c_3$  has a variety of abundant types. Aside from the top four OTUs, separating the data into discrete clusters obscures how the rest of the OTUs can be characterized.

By using each data point's cluster-pair assignment, we further separate the data into  $k^2 - k$  clusters. Let  $c_{jj'}$  designate the data points that are between clusters  $j$  and  $j'$ , but are nearer to cluster  $j$  than cluster  $j'$ . We take the most abundant OTUs in each sample, and the cluster pair for each sample. We can then find the most abundant OTUs for each cluster pair.

Table 2 shows the structure of the most abundant OTU types for each 2-way cluster  $c_{jj'}$  defined before. Once again, we find that clusters  $c_{1j}$ ,  $c_{2j}$ ,  $c_{4j}$ , and  $c_{5j}$  are all dominated by the same single OTU from before. Yet observing clusters  $c_{3j'}$  provides us with a more in-depth understanding of the diverse cluster  $c_3$ .

Interestingly, we see that the makeups of  $c_{31}$ ,  $c_{32}$ ,  $c_{34}$ , and  $c_{35}$  are remarkably different. We immediately see that

**Table 1.** Most abundant OTU per cluster. Because the 16S rDNA data maps multiple sequences to the same genus level, we use subscripts to denote different OTUs with the same genus.

Family;Genus	c <sub>1</sub>	c <sub>2</sub>	c <sub>3</sub>	c <sub>4</sub>	c <sub>5</sub>
Lactobacillaceae; <i>Lactobacillus</i> <sub>1</sub>	5	534	17	0	1
Lactobacillaceae; <i>Lactobacillus</i> <sub>2</sub>	0	0	3	269	0
Bifidobacteriaceae; <i>Gardenerella</i>	227	0	13	0	1
Lachnospiraceae; <i>Unclassified</i> <sub>1</sub>	0	0	4	0	150
Lactobacillaceae; <i>Lactobacillus</i> <sub>3</sub>	2	0	44	2	0
Lactobacillaceae; <i>Lactobacillus</i> <sub>4</sub>	2	4	32	1	1
Leptotrichiaceae; <i>Sneathia</i> <sub>1</sub>	8	0	33	0	0
Prevotellaceae; <i>Prevotella</i> <sub>1</sub>	2	0	30	0	0
Prevotellaceae; <i>Prevotella</i> <sub>2</sub>	3	0	16	0	1
Unclassified; <i>Unclassified</i> <sub>2</sub>	2	1	15	0	2
Prevotellaceae; <i>Prevotella</i> <sub>3</sub>	0	0	4	0	1
Leptotrichiaceae; <i>Sneathia</i> <sub>2</sub>	1	0	2	0	1
Lachnospiraceae; <i>Unclassified</i> <sub>3</sub>	0	0	1	0	3
Streptococcaceae; <i>Streptococcus</i> <sub>1</sub>	0	0	18	0	0
Veillonellaceae; <i>Unclassified</i> <sub>4</sub>	0	0	0	0	0
Streptococcaceae; <i>Streptococcus</i> <sub>2</sub>	0	1	15	0	0
Mycoplasmataceae; <i>Mycoplasma</i>	0	1	7	0	0
Bifidobacteriaceae; <i>Bifodobacterium</i>	0	0	9	0	0
Fusobacteriaceae; <i>Fusobacterium</i>	0	0	7	0	0
Enterobacteriaceae; <i>Unclassified</i> <sub>5</sub>	0	0	8	0	0

**Table 2.** Most abundant OTUs per cluster-pair.

Genus	c <sub>12</sub>	c <sub>13</sub>	c <sub>14</sub>	c <sub>15</sub>	c <sub>21</sub>	c <sub>23</sub>	c <sub>24</sub>	c <sub>25</sub>	c <sub>31</sub>	c <sub>32</sub>	c <sub>34</sub>	c <sub>35</sub>	c <sub>41</sub>	c <sub>42</sub>	c <sub>43</sub>	c <sub>45</sub>	c <sub>51</sub>	c <sub>52</sub>	c <sub>53</sub>	c <sub>54</sub>
<i>Lactobacillus</i> <sub>1</sub>	5	0	0	0	113	69	57	295	0	15	2	0	0	0	0	0	0	1	0	0
<i>Lactobacillus</i> <sub>2</sub>	0	0	0	0	0	0	0	0	0	0	3	0	14	60	27	168	0	0	0	0
<i>Gardenerella</i>	75	32	95	25	0	0	0	0	10	0	3	0	0	0	0	0	1	0	0	0
<i>Unclassified</i> <sub>1</sub>	0	0	0	0	0	0	0	0	0	0	1	3	0	0	0	0	54	23	23	50
<i>Lactobacillus</i> <sub>3</sub>	0	2	0	0	0	0	0	0	6	7	31	0	0	0	2	0	0	0	0	0
<i>Lactobacillus</i> <sub>4</sub>	1	1	0	0	0	4	0	0	4	4	23	1	0	1	0	0	1	0	0	0
<i>Sneathia</i> <sub>1</sub>	0	8	0	0	0	0	0	0	22	2	4	5	0	0	0	0	0	0	0	0
<i>Prevotella</i> <sub>1</sub>	0	2	0	0	0	0	0	0	9	4	17	0	0	0	0	0	0	0	0	0
<i>Prevotella</i> <sub>2</sub>	0	0	0	3	0	0	0	0	15	0	0	1	0	0	0	0	0	0	1	0
<i>Unclassified</i> <sub>2</sub>	2	0	0	0	0	1	0	0	8	6	0	1	0	0	0	0	1	0	1	0
<i>Prevotella</i> <sub>3</sub>	0	0	0	0	0	0	0	0	3	0	0	1	0	0	0	0	1	0	0	0
<i>Sneathia</i> <sub>2</sub>	0	1	0	0	0	0	0	0	1	1	0	0	0	0	0	0	0	0	1	0
<i>Unclassified</i> <sub>3</sub>	0	0	0	0	0	0	0	0	0	0	0	1	0	0	0	0	1	1	1	0
<i>Streptococcus</i> <sub>2</sub>	0	0	0	0	0	0	0	0	2	2	14	0	0	0	0	0	0	0	0	0
<i>Unclassified</i> <sub>4</sub>	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
<i>Streptococcus</i> <sub>2</sub>	0	0	0	0	0	1	0	0	2	2	11	0	0	0	0	0	0	0	0	0
<i>Mycoplasma</i>	0	0	0	0	1	0	0	0	5	2	0	0	0	0	0	0	0	0	0	0
<i>Bifodobacterium</i>	0	0	0	0	0	0	0	0	1	1	7	0	0	0	0	0	0	0	0	0
<i>Fusobacterium</i>	0	0	0	0	0	0	0	0	1	2	4	0	0	0	0	0	0	0	0	0
<i>Unclassified</i> <sub>5</sub>	0	0	0	0	0	0	0	0	0	2	5	1	0	0	0	0	0	0	0	0

the top four OTUs are all predominantly contained in cluster pairs that includes their single main cluster. In addition, we notice that the samples with abundant *Sneathia*<sub>1</sub>, *Prevotella*<sub>2</sub>, and *Unclassified* types are predominantly contained in c<sub>31</sub>. c<sub>32</sub> contains samples with a variety of abundant OTUs.

*Lactobacillus*<sub>3</sub>, *Lactobacillus*<sub>4</sub>, *Prevotella*<sub>1</sub>, *Streptococcus*<sub>1</sub>, *Streptococcus*<sub>2</sub>, and *Bifodobacterium* are abundant in samples that are predominantly contained in c<sub>34</sub>. Finally, almost no samples are in cluster pair c<sub>35</sub>, aside from a few *Sneathia*<sub>1</sub> types.

In this way, 2-way  $k$ -means also opens up a wealth of information on the relationships between samples. In particular, it now makes more sense to characterize the samples as being in 6 different clusters:  $c_1, c_2, c_{31}, c_{34}, c_5$ . We also see that certain clusters have mixed relationships, while others have almost no interaction. Without 2-way  $k$ -means this would not be immediately obvious.

## 4 CONCLUSION

The complexity of microbial populations is unfolding as microbiome data becomes increasingly available. Yet, standard methodologies oversimplify microbial compositions by pigeonholing them into discrete clusters. This paper further refines the models for microbial abundance across groups of samples. We allow samples to be presented as a weighted average of two clusters, rather than belonging to only one. This may be motivated biologically, as the sample often reflects a mixture of two sources of microbiota, each well represented by a cluster. An alternative explanation is that the averaged sample represents an intermediate, potentially temporary state of the microbial composition, between the more stable ones represented by the clusters themselves.

Technically, we formalize this model as a generalization of  $k$ -means. We derive a simple algorithm to infer such a structure, and validate its benchmarks on simulated data.

Applying our algorithm to real data from the Human Vaginal Microbiome Project provides empirical support to the 2-way model. We showed that while most of the samples lie in six clusters: four well-defined clusters and two subclusters. Furthermore, while previously, a sizable fraction of samples in-between clusters were ignored, the 2-way model characterized the entire distribution. Using 2-way  $k$ -means, we can tell that a large portion of the previously unclustered samples, which lie in-between two clusters, contain shared properties. In addition, we see that certain clusters have mixed relationships, while others have almost no interaction.

## 5 FURTHER RESEARCH

In addition, this paper leaves several open questions and opportunities for further research:

- How can we efficiently characterize a 2-way distribution with non-spherical covariance matrices?
- How can we efficiently characterize a  $k$ -way distribution?

- How can we efficiently characterize a 2-way distribution with non-linear paths between cluster representatives?

Addressing these questions will further help us understand the composition of microbial populations.

## 6 ACKNOWLEDGEMENTS

This work was supported by the National Science Foundation under CISE EAGER grant #1547120.

*The authors declare that there is no conflict of interest regarding the publication of this paper.*

## REFERENCES

- [1] Human Microbiome Project Consortium. Structure, function and diversity of the healthy human microbiome. *Nature*, 486(7402):207–214, 2012.
- [2] Raes J. Arumugam M. Burgdorf K.S. Manichanh C. Nielsen T. Qin J., Li R. et al. A human gut microbial gene catalog established by metagenomic sequencing. *Nature*, 464(7285):59–65, 2010.
- [3] Ramette A. and Buttigieg P. L. A guide to statistical analysis in microbial ecology: a community-focused, living review of multivariate data analyses. *FEMS Microbiol Ecol*, 90(3):543–550, 2014.
- [4] Human Microbiome Project Consortium. A framework for human microbiome research. *Nature*, 486(7402):215–221, 2012.
- [5] Serrano M.G. Fettweis J.M., Brooks J.P. et al. Differences in vaginal microbiome in african american women versus women of european ancestry. *Microbiology*, 160(10):2272–2282, 2014.
- [6] Plenge R. Weinblatt M. Shadick N. Price A., Patterson N. and Reich D. Principal components analysis corrects for stratification in genome-wide association studies. *Nat Genet*, 38(8):904–909, 2006.
- [7] Novembre J. Alexander D. H. and Lange K. Fast model-based estimation of ancestry in unrelated individuals. *Genome Research*, 19(9):1655–1664, 2009.
- [8] Bryc K. Kutalik Z. Boyko A. R. Auton A. Novembre J., Johnson T. et al. Genes mirror geography within europe. *Nature*, 456(7219):274, 2008.
- [9] Stuart P. Lloyd. Least squares quantization in pcm. *IEEE Transactions on Information Theory*, 28:129–137, 1982.
- [10] National Center for Biotechnology Information. <https://www.ncbi.nlm.nih.gov/sra/?term=SRP002462>, 2014.
- [11] Ryabin T. Hall J. R. Hartmann M. Schloss P. D., Westcott S. L. et al. Introducing mothur: Open-Source, Platform-Independent, Community-Supported Software for Describing and Comparing Microbial Communities. <http://aem.asm.org/content/75/23/7537.short?rss=1&ssource=mfc>, 2009.
- [12] Gevers D. Schloss, P. D. and Westcott S. L. Reducing the effects of pcr amplification and sequencing artifacts on 16s rna-based studies. *PLoS ONE*, 6(12), 2011.
- [13] Jackson W. J. 2-way cluster assignment. <https://github.com/westonjackson/2-Way-Cluster-Assignment>, 2016.