



**Compositional analysis: a valid approach to analyze
microbiome high throughput sequencing data**

| | |
|-------------------------------|--|
| Journal: | <i>Canadian Journal of Microbiology</i> |
| Manuscript ID | cjm-2015-0821.R2 |
| Manuscript Type: | Article |
| Date Submitted by the Author: | 30-Mar-2016 |
| Complete List of Authors: | Gloor, Gregory; The University of Western Ontario, Biochemistry Reid, Gregor; The Lawson Research Institute |
| Keyword: | microbiome, compositional data, correlation, multivariate statistics, multiple test correction |
| | |



1 Compositional analysis: a valid approach to analyze microbiome
2 high throughput sequencing data

3

4 Gregory B. Gloor (1,2)*, Gregor Reid (2, 3)

5

6

7 1. Department of Biochemistry, Western University, London, Ontario, Canada

8 2. Canadian Center for Human Microbiome and Probiotic Research, Lawson Health
9 Research Institute, London, Ontario, Canada

10 3. Departments of Microbiology and Immunology, and Surgery, Western University,
11 London, Ontario, Canada

12

13

14

15 * Address for Correspondence: Gregory B. Gloor, E-mail: ggloor@uwo.ca

16

17

18

19

20

21

22

23

24 **Abstract**

25 A workshop held at the 2015 annual meeting of the Canadian Society of
26 Microbiologists highlighted compositional data analysis methods, and the importance of
27 exploratory data analysis, for the analysis of microbiome datasets generated by high
28 throughput DNA sequencing. A summary of the content of that workshop, a review of
29 new methods of analysis, and information on the importance of careful analyses are
30 presented herein. The workshop focussed on explaining the rationale behind the use of
31 compositional data analysis, and a demonstration of these methods for the examination
32 of two microbiome datasets. A clear understanding of bioinformatics methodologies and
33 the type of data being analyzed is essential given the growing number of studies
34 uncovering the critical role of the microbiome in health and disease, and the need to
35 understand alterations to its composition and function following intervention with fecal
36 transplant, probiotics, diet and pharmaceutical agents.

37

38

39

40 **Key Words:** microbiome, compositional data, correlation, multiple test correction

41

42

43

44

45

46

47

48

49

50 **Introduction**

51 Human microbiome studies have shown a major link between microbial
52 composition and health and disease and dysbiosis (Fremont et al. 2013; Lourenço et al.
53 2014; Urbaniak et al. 2014). High throughput DNA sequencing methodologies have
54 made this possible, along with breakthroughs in culturing techniques. The former has
55 used approaches such as 16S rRNA gene sequencing, metagenomics, transcriptomics
56 and meta-transcriptomics, leading to vast datasets that must be simplified and analyzed
57 (Di Bella et al. 2013). Indeed, each sample may have tens of thousands to millions of
58 sequence reads associated with it, and the entire dataset across all samples can easily
59 exceed many hundreds of millions of reads. Such has been the rapidity of these
60 developments that some studies appear to have been published using methods that are
61 potentially. The result can be papers with serious deficiencies that are publicized as
62 major advances or breakthroughs (Reardon 2013), when in some cases the data are far
63 from sufficient for such claims. We will examine the evidence for one of these papers
64 below (Hsiao et al. 2013).

65 Data for microbiome analysis are collected by the following general workflow.
66 The sample (swab, stool, saliva, urine or other type) is collected, the DNA is isolated
67 and used in a polymerase chain reaction with primers specific to one or more variable
68 regions of the 16S rRNA gene. It is also possible to target other conserved genes such
69 as the *cpn60* gene (Schellenburg et al. 2009). However, analysis problems are the

70 same regardless of the amplification target chosen, and Walker et al. (2015) present a
71 good summary of how choices taken upstream of data analysis affect the results.
72 Following amplification, a random sample of the product is used to make a sequencing
73 library, and it is common to multiplex many samples in the library. A small aliquot of the
74 library is processed on the high throughput DNA sequencing instrument. As outlined
75 below, this workflow imposes constraints on the resulting data.

76 It should be recognized that the investigator is sequencing a random sample of
77 the DNA in the library, which is itself a random sample of the DNA in the environment.
78 Thus, it is important to ensure that any analysis takes this random component into
79 account (Fernandes et al. 2013).

80 Perhaps less obvious is that the number of sequencing reads obtained for a
81 sample bears no relationship to the number of molecules of DNA in the environment,
82 because the number of reads obtained for a sample is determined by the capacity of the
83 instrument. For example, the same library sequenced on an Illumina MiSeq or HiSeq
84 would return approximately 20 million or 200 million reads. That there is no information
85 in the actual read numbers per sample is implicitly acknowledged by the common use of
86 'relative abundance' values for analysis of microbiome datasets. Such datasets are
87 referred to as compositional and there is a long history of the development of proper
88 analysis techniques for such data in other fields (Pawlowsky-Glahn et al. 2015).

89 Compositional data is a term used to describe a dataset in which the parts in
90 each sample have an arbitrary or non-informative sum (Aitchison 1986), such as data
91 obtained from high throughput DNA sequencing (Friedman and Alm 2012, Fernandes et
92 al. 2013, 2014). These data have long been known to be problematic (Pearson 1896),

93 and we now understand that multivariate data analysis approaches such as ordination
94 and clustering and univariate methods that measure differential abundance are invalid
95 (Aitchison 1986, Warton et al. 2012, Friedman and Alm 2012, Fernandes et al. 2013
96 Pawlowsky-Glahn et al. 2015).

97 The essential problem is illustrated in Figure 1 where we set up an artificial
98 example and count the number of molecules in the environment. We allow one part
99 (shown as solid black) to increase 10-fold between samples 1 and 2, while the
100 abundance of the other 49 parts (in open circles) remain unchanged. The proportion
101 panel shows how the data are distorted when we convert it to relative abundances or
102 proportions, or as happens when the sequencing instrument imposes a constant sum.
103 The black part still appears to become more abundant, although it is less than a 10-fold
104 change. However, the 49 other parts appear to become less abundant. This property
105 leads to the *negative correlation bias* observed in compositional data, and renders
106 invalid any type of correlation or covariance based analysis such as correlation
107 networks, principle component analysis, and others (Pearson 1896, Aitchison 1986).
108 Note that this distortion will also lead to false univariate inferences as well (Fernandes
109 et al. 2013,2014).

110 The original issue with compositional data identified by Pearson (1896) was that
111 of spurious correlation. That is, two or more variables can appear to be correlated
112 simply because the data are transformed to have a constant sum. Spurious correlation
113 also causes the correlations observed in these data to depend on the membership of
114 the sample. For example, consider the simple case of three samples (a, b and c) with

115 four taxonomic variables measured to have the following absolute counts in three
 116 environmental samples (i.e., samples are in rows, taxa are in columns):

$$117 \quad abc = \begin{bmatrix} 470 & 66 & 839 & 751 \\ 541 & 569 & 787 & 512 \\ 167 & 906 & 959 & 504 \end{bmatrix}, \text{cor}(abc) = \begin{bmatrix} & -0.68 & \mathbf{-0.99} & 0.36 \\ -0.77 & & \mathbf{0.59} & -0.93 \\ \mathbf{-0.30} & \mathbf{-0.37} & & -0.25 \\ 0.55 & -0.95 & 0.62 & \end{bmatrix}.$$

118 The Pearson correlation for the numerical values is in the upper triangle of the
 119 right hand matrix, and we see that taxon 1 and taxon 3 have a near perfect negative
 120 correlation of -0.99 (shown in bold), and taxon 2 and taxon 3 have a positive correlation
 121 of 0.59. The lower triangle on the right hand matrix shows the Pearson correlation
 122 values that are found when these are converted to relative abundances by dividing by
 123 the total sum of counts in each sample. Now, the correlations between the same taxa
 124 have changed. The correlation between 1 and 3 is now moderately negative at -0.30,
 125 and between 2 and 3 is now -0.37. Thus, the correlation observed in compositional data
 126 is not the same as the correlation for the counts, and the correlations measured can
 127 even change sign.

128 There is a further complication: the correlations observed in compositional data
 129 depend on the membership in the sample. So, for example, when the last value is
 130 dropped from each sample, the correlations between taxa 1 and 2 is positive (0.43), and
 131 the correlation between 2 and 3 is even more strongly negative at -0.79. Thus,
 132 correlation determined from compositional data has the potential to be wildly wrong, and
 133 normal approaches to determine correlation cannot be used (Friedman and Alm 2012,
 134 Lovell et al. 2015, Kurtz et al. 2015). It is worth noting that any method of determining
 135 correlation (including Spearman, Kendall, etc) will suffer from the same problems. Thus
 136 the current tools used to examine the analysis goals give results that may be

137 inconsistent, difficult to interpret and in many cases completely wrong (Filmoser et al.
138 2009, Friedman and Alm 2012, Fernandes et al 2013, Fernandes et al. 2014, Lovell et
139 al. 2015, Kurtz et al. 2015).

140 The essential first step of proper compositional data analysis is to convert the
141 relative abundances of each part, or the values in the table of counts for each part, to
142 ratios between all parts. This can be accomplished in several ways (Aitchison 1986), but
143 the most widely used and most convenient for our purposes is to convert the data using
144 the centred log-ratio (clr) transformation. So if X is a vector of numbers that contains D
145 parts:

$$146 X = [x_1, x_2, \dots, x_D],$$

147 the centered log-ratio of X can be computed as:

$$148 X_{\text{clr}} = [\log[x_1/g_X], \log[x_2/g_X], \dots, \log[x_D/g_X],$$

149 where g_X is the geometric mean of all values in vector X (Aichison 1986). This
150 simple transformation renders valid all standard multivariate analysis techniques
151 (Aitchison 1986, van den Boogaart 2013, Pawlowsky-Glahn et al. 2015), and as shown
152 in the Ratios panel of Figure 1, can reconstitute the shape of the data so that univariate
153 analyses are also more likely to be valid. This transformation is also the starting point
154 for essentially all compositional data analysis (CoDa) based assessments of the
155 datasets.

156 A CoDa approach would be robust if microbiome datasets were not sparse, that
157 is, they did not contain any 0 values. However a frequent criticism of the CoDa
158 approach is that the geometric mean cannot be computed if any of the values in the
159 vector are 0. It is here we reiterate that our data represent the counts per taxon through

160 the process of random sampling (Fernandes et al. 2013, 2014). Thus, some 0 values
161 could arise simply by random chance, while others arise because of true absence of the
162 taxon in the environment. Fortunately, we can couple Bayesian approaches to estimate
163 the likelihood of 0 values with the compositional analysis approach (Fernandes et al.
164 2013, 2014, Gloor et al. 2016). With this paradigm we dispose of taxa with 0 counts in
165 all or most samples (Palarea-Albaladejo and Martin-Fernandez 2015), and assign an
166 estimate of the likelihood of the 0 being a sampling artifact to the remainder. When
167 performing univariate tests or correlation analyses, it is often convenient to keep many
168 such estimates of 0 and to determine the expected value of test statistics to reduce
169 false positive inferences (Friedman and Alm 2012, Fernandes et al. 2013, Fernandes et
170 al. 2014).

171 **Microbiome analysis tools that account for compositional data**

172 Fortunately, the compositional data analysis problem of microbiome datasets is starting
173 to be examined by several groups and there are now an increasing number of tools
174 available as outlined below.

175 These tools can be applied to address three major objectives of many microbiome
176 analyses:

- 177 1. Do the data show any structure? That is, do the data partition into groups?
- 178 2. What is the difference between groups? This can be between groups identified
179 beforehand, or following the exploratory data analysis.
- 180 3. What is the correlation structure of the taxonomic groups? Do any of these taxa
181 correlate with the metadata?

182 These analyses are usually done using either the mothur (Schloss et al. 2009) or the
183 QIIME (Kuczynski et al. 2012) aggregated toolsets, containing approaches adapted
184 from the field of ecology. However, the use of an analysis paradigm based on
185 compositional data analysis (Aitchison 1986), or CoDa, offers a number of advantages
186 over these tools, as explained below.

187 The first objective is to determine if there is structure in the dataset. In the
188 microbiome field this is generally described as beta-diversity analysis. Beta-diversity as
189 currently used requires a distance or dissimilarity measure, and popular ones include
190 the unweighted or weighted Unifrac distance metrics (Lozopone and Knight 2005) or the
191 Bray-Curtis dissimilarity metric. These methods are included in both the mothur and
192 QIIME toolkits. The distance metrics from these tools can be used to generate Principle
193 Co-ordinate (PCoA) plots that can be used to assess similarities and differences
194 between samples and groups. Unfortunately, distance-based tools can confuse location
195 (difference) and dispersion (variance) effects (Warton et al. 2012), and so additional
196 approaches based on a compositional paradigm should be used for exploratory data
197 analysis.

198 The CoDa analysis analog to PCoA is a principle component analysis (PCA) of
199 center-log ratio transformed data that has been modified to either remove taxa with 0
200 observed counts, or to adjust 0 values to an estimated value (Palarea-Albaladejo and
201 Martin-Fernandez 2015). PCA has the advantage of being a more interpretable metric
202 than PCoA, since it directly assesses the variance in the data and because both the
203 locations of the samples and the contribution of each taxon to the total variance can be
204 shown on the so-called compositional biplot (Aitchison and Greenacre 2002). The ability

205 to examine variation of both the samples and the taxa on the same plot provides
206 powerful insights into which taxa are compositionally associated and which taxa are
207 driving (or not) the location of particular samples. Thus, the biplot can serve as a
208 summary of the entire dataset, and it is up to the investigator to attach numerical
209 significance to the qualitative results observed. The example usage of compositional
210 biplots is explained in detail below.

211 The second major objective is often to determine which taxa are driving the
212 difference observed between groups. Several methods are in widespread use to assess
213 the difference in abundance of taxa between groups. These include microbiome specific
214 methods such as Metastats (White et al. 2009) or LEfSe (Segata et al. 2011), and more
215 general t-tests or nonparametric tests. However, all use as input a table of proportional
216 abundances. As shown in Figure 1, examination of proportions can result in a gross
217 distortion of the data, such that some taxa can appear to change in abundance when
218 measured by proportion, when in fact, their true abundance in the environment may be
219 unchanged. This effect can be ameliorated by the center-log ratio transformation.

220 There are two approaches that assess differential abundance in a compositional
221 data analysis framework. The simplest approach is the ANCOM tool (Mandal et al.
222 2015), which assesses statistical significance on log-ratio transformed data. This is
223 more robust than both traditional t-tests and more sophisticated approaches such as
224 zero-inflated Gaussian methods. It should be noted that the software is not currently
225 deposited into a public repository, and that the 0-replacement value used is fixed in the
226 software.

227 A slightly more complex approach is used by the ALDEx2 package, available
228 from Bioconductor (Fernandes et al 2013, Fernandes et al 2014). Like ANCOM,
229 ALDEx2 centre log-ratio transforms the data prior to the assessment of statistical
230 significance, however ALDEx2 differs greatly in how values of 0 are handled. ALDEx2
231 estimates a large number of possible values for 0 (and any other count for a taxon in a
232 sample), conducts significance tests on all estimated values, and takes the average
233 significance test value as the most representative for that taxon. In essence, ALDEx2
234 determines which taxa are significantly different between groups after accounting for the
235 random sampling that occurs when the DNA is extracted and loaded onto the
236 sequencing instrument. In either case, both ANCOM and ALDEx2 explicitly
237 acknowledge the multivariate compositional nature of the data, and control for false
238 positive identifications much better than do the usual approaches.

239 The third objective is to determine if there are taxa in the dataset with correlated
240 abundances. As noted above, spurious correlation is a very large problem in
241 microbiome datasets. Therefore, analyses that report correlations using traditional
242 methods, such as Pearson's or Spearman's correlations, Kendall's Tau or Partial
243 correlations are likely to be wrong (Friedman and Alm 2012, Lovell et al. 2015, Kurtz et
244 al 2015). However, there are a number of approaches that use a compositional data
245 analytic approach to correlation. In a compositional approach, the variance between
246 ratios of two taxa should be 0 or nearly so for two taxa to be counted as correlated
247 (Aitchison 1986, Lovell et al. 2015). The difficulty comes when placing this approach
248 into a familiar null hypothesis test framework, or when applying a consistent scale to the
249 measure. The simplest approach is to calculate the phi statistic for two taxa X and Y,

250 which is the $\text{var}(\log(X/Y))/\text{var}(\log(X))$ (Lovell et al. 2015), where $\log()$ is meant to imply
251 the clr values of X or Y. This measure has the advantage of being easily calculated and
252 of strictly enforcing the compositional data analysis approach. The SparCC method
253 (Friedman and Alm, 2012) uses Bayesian estimates of the value of X and Y but
254 calculates a mean value of a measure similar to the concordance correlation coefficient.
255 The SPIEC-EASI approach (Kurtz et al. 2015) uses clr-transformed values and infers a
256 graphical model under the assumption of a sparse correlation network. Both of the latter
257 approaches make strong assumptions about the sparsity of the data, and so are less
258 rigorous for estimating correlations in compositional data than is the calculation of phi.
259 However, they both offer the advantage of using a full or partial Bayesian approach,
260 which is generally more powerful than point-estimate based approaches.

261 **Application of CoDa to Two Case Studies**

262 Having introduced the issue of compositional data analysis, we now present the
263 results of two worked examples presented at the Bioinformatics Workshop was held on
264 June 16, 2015 in Regina at the Annual Scientific Meeting of the Canadian Society of
265 Microbiologists. This illustrates how these approaches can be applied to two different
266 16S rRNA gene sequencing datasets from the recent literature. A full description of the
267 methodology, the datasets and the code used to generate the figures is given in the
268 Supplementary file workshop.Rnw (Gloor 2016). Downloading and running this file in R
269 (R Core Team 2015) or RStudio will generate the associated workshop.pdf. The .Rnw
270 document contains both the code and annotation for the code, and the .pdf document
271 contains the code and the resulting figures.

272 The first worked example is a vaginal microbiome dataset. This dataset is from
273 an experiment that examined the effect of treating women suffering from bacterial
274 vaginosis (BV) with antibiotics and placebo or antibiotics plus a probiotic supplement
275 (Macklaim et.al, 2015). For this example, we extracted only the 'before' (samples
276 labeled as BXXX) and 'after' (AXXX) treatment samples, which were further identified by
277 their Nugent status, a Gram stain scoring system that acts as a rough indicator of
278 whether the subject had BV or was healthy (normal, n), or whose status was
279 indeterminate (labeled as 'i' for intermediate). In addition, individual taxa were
280 aggregated to genus level using QIIME (Kuczynski et al. 2012), except for *Lactobacillus*
281 *iners* and *Lactobacillus crispatus*, which remained as separate species in the tables.
282 This relatively simple dataset will be used to introduce and explain the CoDa analysis
283 methods.

284 The compositional biplot is the essential initial tool for exploratory compositional
285 data analysis and replaces ordinations based on Unifrac or Bray-Curtis metrics.
286 Compositional biplots are principle component plots of the singular value decomposition
287 of the data. This approach displays the major axes of variance (or change) in a dataset
288 (Aitchison and Greenacre 2002). Properly made and interpreted, these plots summarize
289 all the essential results of an experiment. However, it should be remembered that they
290 are descriptive and exploratory, not quantitative. Quantitative tools can be applied later
291 to support the conclusions derived from the biplot.

292 For simplicity, we filtered the dataset to include only those taxa that were at least
293 0.1% abundant in any sample. One of the desirable properties of compositional data
294 analysis is that subsets of the dataset are expected to give essentially the same answer

295 as the entire dataset *for the taxa in common* between the whole and the subset dataset
296 (Aitchison 1986).

297 Figure 2 shows the compositional biplot for this dataset along with the associated
298 scree plot that displays the percentage of variance explained by each sample or
299 component. The sample names (labeled in red for BV, blue for Normal or purple for
300 Intermediate) illustrate the variance of the samples, and the taxa values (represented by
301 the black rays) illustrate the variance between the taxa. In fact, the length of the arrow
302 for each taxon is proportional to the standard deviation of the ratio of each taxon to all
303 other taxa. There are many interpretation rules for biplots of compositional data
304 (Aitchison and Greenacre 2002), but these rules are dependent on remembering that
305 only the *ratios* between taxa can be examined. Thus, the links between the tips of the
306 rays, or between samples contain the most information. Keeping this in mind, we can
307 see the following:

308 First, the proportion of variance explained in the first component is very good,
309 being 47%, then falling to 13% on component 2, and decreasing rapidly thereafter. This
310 indicates that the major difference between samples can be captured in essentially one
311 direction along component 1. While the amount of variance explained on the first
312 component is relatively large in this dataset, a rule of thumb is that PCA plots that
313 display less than 80% of the variance on the first two components are not necessarily
314 accurate projections of the data. Thus, some of the quantitative results are expected to
315 be somewhat different than is displayed in the qualitative PCA projection.

316 Second, the longest link from the center to a taxon is the one to *Lactobacillus*
317 *iners*. This indicates that the ratio of this taxon to all others is the most variable across

318 all samples. Likewise, the shortest link is to *Gardnerella*, implying that the ratio of this
319 taxon to all others is the least variable.

320 Third, the longest link is between *L. iners* and *Leptotrichia* (*Sneathia*). This
321 means we can infer that these two taxa likely have the strongest reciprocal ratio
322 relationship. That is, when one becomes more abundant relative to everything else, the
323 other becomes less abundant relative to everything else.

324 Fourth, the shortest link observed in the plot is between *Megasphaera* and
325 BVAB2. From this we conclude that the ratio of these two taxa is relatively constant
326 across all samples. That is, their ratio abundance is highly correlated. These two taxa
327 should be seen to have a low value of phi, but we must keep in mind the limit of the
328 projection of the data.

329 Fifth, the link between *Prevotella* and *Lactobacillus crispatus* passes directly
330 through *Atopobium*. This indicates that these three taxa are linearly related. In this case,
331 it is clear when *L. crispatus* increases, the other two will decrease. Likewise, this
332 property can be extended to any linear relationships containing three or more links.

333 Sixth, the link between *L. iners* and *Megasphaera*, and the link between
334 *Leptotrichia* (*Sneathia*) and *Lactobacillus* cross at approximately 90°. The cosine of the
335 angle approximates the correlation between the connected log ratios. Thus, we can
336 conclude that the abundance relationship between the former pair of taxa is poorly
337 correlated with that of the latter two taxa. In other words, these two pairs vary
338 independently in the dataset.

339 Some samples (A312_bv, B312_i, A282_n at the bottom), are tightly grouped,
340 indicating that they contain similar sets of taxa at similar ratio abundances. We can see

341 from the biplot that these samples contain an abundance of *Lactobacillus* and are
342 depleted in *Leptotrichia* (*Sneathia*). Furthermore, we can see that the samples divide
343 into two fairly clear groups, with most of the before or “B” samples on the left, and most
344 of the after or “A” samples on the right. We further observe that the majority of the B
345 samples are colored red indicating a diagnosis of BV, and the majority of the A samples
346 are colored blue indicating a diagnosis of non-BV.

347 The result of the biplot suggested that there were two main groups that could be
348 defined with this set of data. With a few exceptions, there appears to be a fairly strong
349 separation between the samples containing a majority of *Lactobacillus* sp., and those
350 lacking them. We can explore this by performing an unsupervised cluster analysis on
351 the log-ratio transformed data. In traditional microbiome evaluation methodologies,
352 clustering is based on the weighted or unweighted unifrac distances or on the Bray-
353 Curtis dissimilarity metric, for example see the standard workflow in QIIME (Kuczynski
354 et al. 2012). These metrics are much more sensitive to the abundance of community
355 members than is the Aitchison distance used in compositional data analysis (Martin
356 Fernandez 1998). Thus, here we used the Aitchison distance metric that fulfills the
357 criteria required for compositional data. In particular, by using a compositional approach,
358 it is appropriate to examine a defined sub-composition of the data (i.e., a subset of the
359 taxa).

360 The results of unsupervised clustering of the dataset are shown in Figure 3.
361 Again, it is important to remember that all distances are calculated from the ratios
362 between taxa, and not on the taxa abundances themselves. For this figure, we used the
363 ward.D2 method which clusters groups together by their squared distance from the

364 geometric mean distance of the group. There are many other options, and the user
365 should choose one that best represents the data, although Ward.D and Ward.D2 are
366 usually the most appropriate (Martin-Fernandez 1998).

367 The cluster analysis supports the results of the biplot and shows the split
368 between two types of samples rather clearly. Samples containing an abundance of
369 *Lactobacillus* sp. are grouped together on the right, and samples with an abundance of
370 other taxa are grouped together on the left. The cluster analysis helps explain and
371 clarify the compositional biplot. For example, the four samples in the middle lower part
372 of the biplot in Figure 2 labelled A/B312 and A/B282, group together in both the biplot
373 and the cluster plot. These samples are atypical for both the N and BV groups,
374 containing substantially more of the *Lactobacillus* taxon, and somewhat more of the
375 taxa normally found in BV than in the other N samples. Based on these two results it
376 would be appropriate to exclude these four samples from further analysis because of
377 their atypical makeup.

378 Next, a univariate comparison between the B and A groups was performed. For
379 simplicity of coding, we kept the outlier samples, but the reader is encouraged to
380 remove them and see how the results change. For this, we used the ALDEx2 tool
381 (Fernandes et al. 2013, 2014) that incorporates a Bayesian estimate of taxon
382 abundance into a compositional framework, with the results shown in Table 1 and the
383 effect plot (Gloor et al. 2016) shown in Figure 4. Of note, ALDEx2 examines differential
384 abundance by estimating the measurement error inherent in high throughput DNA
385 sequencing experiments, including the measurement error associated with 0 count taxa,

386 and uses the assumptions of compositional data analysis to normalize the data for the
387 differing number of reads in each sample (Fernandes et al. 2013, Lovell et al. 2015).

388 When interpreting these results, it is important to remember that we are actually
389 examining ratios between values, rather than abundances. Thus, we are examining the
390 change in abundance of a taxon *relative to all others* in the dataset. The user should
391 also remember that all values reported are the means or medians over the number of
392 Dirichlet instances as given by the mc.samples variable in the aldex.clr function and
393 explained more fully in the supplementary material and the original papers (Fernandes
394 et al. 2013, 2014).

395 In the examples given in Table 1, we filtered to show only those taxa where the
396 expected Bejamini-Hochberg (1995) adjusted P value was less than 0.05, meaning that
397 the expected likelihood of a false positive identification per taxon is less than 5%, with
398 the actual value per taxon given in the wi.eBH column. Using *L. iners*, we note that the
399 absolute difference between groups can be up to -2.25. Thus, the absolute fold change
400 in the ratio between *L. iners* and all other taxa between groups for this organism is on
401 average 4.76 fold ($1/2^{-2.25}$): being more abundant in the A samples than in the B
402 samples. However, the difference within the groups (roughly equivalent to the standard
403 deviation) is even larger, giving an effect size of -0.79. Thus, the difference between
404 groups is less than the variability within a group, a result that is typical for microbiome
405 studies.

406 These quantitative results are largely congruent with the biplot, which showed
407 that the taxa represented here were the ones that best explained the variation between
408 groups, and that the *Leptotrichia (Sneathia)* and *Lactobacillus* taxa were not contributing

409 to the separation of the two large groups and so would not be expected to be
410 significantly different, despite being highly variable.

411 The left panel of Figure 4 shows a plot of the within (diff.win) to between (diff.btw)
412 condition differences, with the large black dots representing those that have a BH
413 adjusted P value of 0.05 or less. Taxa that are more abundant than the mean in the B
414 samples have positive y values, and those that are more abundant than the mean in the
415 A samples have negative y values. These are referred to as 'effect size' plots, and they
416 summarize the data in an intuitive way (Gloor et al. 2015). The grey lines represent the
417 line of equivalence for the within and between group values. Small black dots represent
418 taxa that are less abundant than the mean taxon abundance: here it is clear that the
419 abundance of rare taxa, are generally difficult to estimate with any precision.

420 The middle plot in Figure 4 shows a plot of the effect size vs. the BH adjusted P
421 value, with a strong correspondence between these two measures. In general, an effect
422 size cutoff is preferred because it is more robust than P values. The right plot in this
423 figure shows a volcano plot for reference.

424 Finally, we can determine which taxa are most correlated or compositionally
425 associated. As noted above, correlation is especially problematic, and the only way to
426 avoid false positive associations is to identify those taxa that have constant or nearly
427 constant ratios in all samples: this is the underlying basis of the phi measure (Lovell et
428 al. 2015). In the example shown in the supplementary material, we calculate the mean
429 phi using the same philosophy as outlined above for univariate statistical tests.

430 In the context of microbiome datasets, the phi metric (Lovell et al. 2015) seeks to
431 identify those pairs of taxa that have a near constant ratio abundance across all

432 samples. Applying this approach to the dataset shows that the two most compositionally
433 associated taxa are *Prevotella* sp. and *Megasphaera* sp. Note, that these taxa do not
434 have the shortest links in the compositional biplot, indicating that the amount of variance
435 explained is not high enough to provide an accurate projection of the dataset.

436 For the second worked example we include in the workshop.Rnw document a
437 second example based on the data of Hsiao et al. (2013) that examined the effect of
438 *Bacteriodes fragilus* supplementation on the microbiome composition of a mouse model
439 of autism. This paper determined that there was a strong functional association between
440 *B. fragilus* supplementation and mouse behavior. One of the major conclusions was that
441 this functional change in behavior was associated with changes in abundance of a
442 number of bacteria that composed the mouse gut microbiome. We will focus our
443 analysis only on the conclusions derived from the analysis of the microbiome data that
444 were presented in Figure 4 of the paper.

445 Figure 5 shows a compositional biplot of this dataset, and it is obvious that there
446 is little evidence of difference between the poly-IC treated control (IC) and poly-IC
447 treated mice supplemented with *B. fragilus* (Bf) groups when analyzed using this
448 approach. This is in accordance with their conclusions when analyzing the data using
449 an unweighted Unifrac distance based approach. Interestingly, the compositional biplot
450 shows that the Bf samples are generally closer to the origin of the plot than are the IC
451 samples, suggesting that the Bf samples have lower dispersion than the IC samples.

452 Since the authors concluded that there was no evidence for multivariate
453 differences between groups, and the CoDa approach agrees, it is generally not advised

454 to conduct a univariate analysis since it is likely that only false positive results would be
455 obtained (Hubert and Wainer 2012).

456 However, these authors went on to identify a number of univariate differences in
457 taxon abundance between groups using the LEfSe and Metastats tools that are
458 standard in the field (White et al. 2009, Segata et al. 2012), but that do not assume the
459 data are multivariate compositions. When examining univariate differences with the
460 ALDEx2 tool, we found that none of the univariate differences reported in the original
461 paper were supported by subsequent analysis. In particular, the authors indicated that
462 the largest differences between groups were found for six taxa labeled as 53, 145, 638,
463 836, 837, and 956 in Figure 4 of the paper. The reason for this discrepancy is that
464 inspection of the original paper reveals that raw, and not Benjamini-Hochberg adjusted
465 P values were reported. Thus it is likely that the majority, if not all, of the taxa different
466 between the control and treatment groups are false positive identifications. This result is
467 congruent with the multivariate results found in both the original paper, and by the
468 compositional biplot. Finally, in support of this assertion, we observe that all of these
469 predicted differences become insignificant following a multiple test correction using
470 either the P values reported in the paper, or P values calculated using the ALDEx2
471 software.

472 While we have been critical of the microbiome analysis methods used in this
473 paper, we must acknowledge that other published papers exhibit many of the same
474 flaws: namely an over-reliance on tools that do not treat the data as compositions, the
475 identification of extremely rare taxa as the most 'significantly different' taxa between
476 groups, and a general lack of corrections for multiple hypothesis testing.

477 **Summary**

478 Because the total number of reads is uninformative in high throughput DNA
479 sequencing datasets, the only information available is the ratio of abundances between
480 components: thus these data are compositional. Using two 16S rRNA gene sequencing
481 datasets, we have illustrated that microbiome data can be examined using a
482 multivariate CoDa approach where the data are ratios between the OTU count in a
483 sample and geometric mean for that sample. Dirichlet Monte-Carlo replicates coupled
484 with the centered log-ratio transformation can ameliorate the sparse data problem
485 inherent in microbiome datasets.

486 In essence, we argue here that 16S rRNA gene sequencing datasets are not
487 special and do not need their own unique statistical analysis approaches. The data
488 generated can be examined by a general multivariate approach after accounting for the
489 compositional nature of the data, and such an analysis is comparable or superior to
490 domain-specific approaches, such as those used in the second example paper (Hsiao
491 et al. 2013).

492 With the human body associated with a large number and diversity of bacteria,
493 we need to understand the evolution of this association and how and when this intimate
494 association develops. Such understanding will in turn lead us to robust approaches
495 focussed on when and how to influence the microbiome by probiotic supplementation or
496 by nutrient or antimicrobial means. More and more studies are exploring how the
497 microbiome can predict outcomes, including following fecal transplant, probiotic, dietary
498 and drug treatment (David et al. 2014; Kwak et al. 2014; Seekatz et al. 2014; Rajca et
499 al. 2014). Such work will require carefully designed studies with high quality clinical

500 documentation, and samples that are processed using some of the methods described
501 herein. As the compositional toolkit for microbiome analysis evolves, these studies will
502 reveal aspects of human life not previously envisaged. In order to have confidence in
503 such findings, datasets must be interrogated with rigour. The public is thirsty for
504 knowledge and the media anxious to attract attention. Reliance on pharmaceutical
505 agents is longer acceptable, and the ability to manipulate the microbiome is not only
506 appealing but actually feasible. Thus, studies that help to understand how such
507 manipulations occur, what communication is taking place between microbes and the
508 host, will allow for more precisely targeted interventions, even to some extent
509 personalized. In particular for the latter, as precise knowledge of microbiome
510 components and activity will be critical.

511 Interested readers wishing to progress beyond this demonstration should consult
512 the compositional data literature, but in particular the original book by Aitchison (1986)
513 and a comprehensive book by Pawlowsky-Glahn et al. (2015) that outlines the essential
514 geometric problem of compositional data as it is understood at present. For a guide that
515 goes beyond the introduction given here and in the supplementary material, a book
516 outlining how to use the compositions R package by Van den Boogaart and Tolosana-
517 Delgado (2013) is particularly helpful, although none of the examples are drawn from
518 the biological literature. For others wishing to understand bioinformatics and data
519 analysis of sequencing data in general terms, hopefully this paper will prove helpful, and
520 encourage people to enroll in specialized courses. The temptation may be to rely on
521 proprietary third party systems, even at a cost, but the 'devil is in the details' and for

522 thoroughness we recommend developing the highest level of skill possible, especially to
523 continue to create new analytical tools.

524 We hope that this report will help researchers to better understand their data and
525 thereby conduct analyses that are more likely to be robust, and more importantly to
526 bring badly needed breakthroughs in prevention, treatment and cure of disease.
527

Draft

528 **Funding:** Financial support for this study was provided by a joint Canadian Institutes of
529 Health Research (CIHR) Emerging Team Grant and a Genome British Columbia (GBC)
530 grant awarded on which GR was a co-PI and GG and ML were co-investigators (grant
531 reference #108030). Work in the lab of GG is also supported by an NSERC Discovery
532 Grant. The funders had no role in study design, data collection and analysis, decision to
533 publish, or preparation of the manuscript.

534

Draft

535 **References**

- 536 Aitchison, J. 1986. The statistical analysis of compositional data, Chapman and Hall,
537 London England. ISBN 1-930665-78-4
- 538 Aitchison, J and Greenacre, M. 2002. Biplots of compositional data. J. Royal Stat. Soc:
539 Series C. 51:375-92
- 540 Benjamini, Y., Hochberg, Y. 1995. Controlling the false discovery rate: a practical and
541 powerful approach to multiple testing. . Royal Stat. Soc: Series B (Methodological), 289-
542 300.
- 543 David, L.A., Maurice, C.F., Carmody, R.N., Gootenberg, D.B., Button, J.E., Wolfe, B.E.,
544 Ling, A.V., Devlin, A.S., Varma, Y., Fischbach, M.A., Biddinger, S.B., Dutton, R.J.,
545 Turnbaugh, P.J. 2014. Diet rapidly and reproducibly alters the human gut microbiome.
546 Nature. 505(7484):559-63.
- 547 Di Bella, J.M., Bao, Y., Gloor, G.B., Burton, J.P., Reid, G. 2013. High throughput
548 sequencing methods and analysis for microbiome research. J Microbiol Methods. 2013
549 Dec;95(3):401-14.
- 550 Fernandes, A. D., Macklaim, J. M., Linn, T. G., Reid, G., Gloor, G. B. 2013. ANOVA-like
551 differential expression (ALDEx) analysis for mixed population RNA-Seq. PloS One, 8(7),
552 e67019.
- 553 Fernandes, A.D., Reid, J.N., Macklaim, J.M., McMurrough, T.A., Edgell, D.R., Gloor,
554 G.B. 2014. Unifying the analysis of high-throughput sequencing datasets: characterizing
555 RNA-seq, 16S rRNA gene sequencing and selective growth experiments by
556 compositional data analysis. Microbiome. 2:15.

- 557 Filzmoser, P., Hron, K., Reimann, C. 2009. Univariate statistical analysis of
558 environmental (compositional) data: problems and possibilities. *Sci Total Environ.*
559 407:6100-8.
- 560 Frémont, M., Coomans, D., Massart, S., De Meirleir, K.. 2013. High-throughput 16S
561 rRNA gene sequencing reveals alterations of intestinal microbiota in myalgic
562 encephalomyelitis/chronic fatigue syndrome patients. *Anaerobe.* 22:50-6.
- 563 Friedman, J., Alm, E. J. 2012. Inferring correlation networks from genomic survey data.
564 *PLoS Comput. Biol.* 8(9): e1002687
- 565 Gentleman, R. C., Carey, V. J., Bates, D. M., Bolstad, B., Dettling, M., Dudoit, S., Ellis,
566 B. et al. 2004. Bioconductor: open software development for computational biology
567 and bioinformatics. *Gen. Biol.* 5 (10): R80.
- 568 Gloor, G.B., Macklaim, J.M., Fernandes, A.F. 2016. Displaying variation in large
569 datasets: a visual summary of effect sizes. *J. Comput. Graph. Stat.* (in press),
570 **DOI:10.1080/10618600.2015.1131161.**
- 571 Gloor, G.B., Macklaim, J.M., Vu, M, Fernandes, A.F. 2016. Compositional uncertainty
572 should not be ignored in high-throughput sequencing data analysis. *Austrian Journal of*
573 *Statistics* (in press).
- 574 Gloor, G.B. 2016. Compositional data analysis for high throughput sequencing: an
575 example from 16S rRNA gene sequencing. Supplementary Information at:
576 http://github.com/ggloor/CJM_Supplement. DOI:10.5281/zenodo.49579.
- 577 Hsiao, E. Y., McBride, S.W., Hsien, S., Sharon, G., Hyde, E.R., McCue, T., Codelli, J.A.,
578 Chow, J., Reisman, S.E., Petrosino, J.F., Patterson, P.H., Mazmanian, S.K. 2013.

579 Microbiota modulate behavioral and physiological abnormalities associated with
580 neurodevelopmental disorders. *Cell*. 155(7):1451-63

581 Hubert, L., Wainer, H. 2012. A statistical guide for the ethically perplexed. CRC Press,
582 London, UK.

583 Kuczynski, J., Stombaugh, J., Walters, W. A., González, A., Caporaso, J. G., Knight, R.
584 2012. Using QIIME to analyze 16S rRNA gene sequences from microbial communities.
585 *Curr. Prot. Microbiol.* 1E-5.

586 Kurtz, Zachary D and Müller, Christian L and Miraldi, Emily R and Littman, Dan R and
587 Blaser, Martin J and Bonneau, Richard A 2015. Sparse and compositionally robust
588 inference of microbial ecological networks. *PLoS Comp. Bio.* 11:e1004226

589 Kwak, D.S., Jun, D.W., Seo, J.G., Chung, W.S., Park, S.E., Lee, K.N., Khalid-Saeed,
590 W., Lee, H.L., Lee, O.Y., Yoon, B.C., Choi, H.S. 2014. Short-term probiotic therapy
591 alleviates small intestinal bacterial overgrowth, but does not improve intestinal
592 permeability in chronic liver disease. *Eur J Gastroenterol Hepatol.* 26(12):1353-9.

593 Lourenço, T.G., Heller, D., Silva-Boghossian, C.M., Cotton, S.L., Paster, B.J., Colombo,
594 A.P. 2014. Microbial signature profiles of periodontally healthy and diseased patients. *J*
595 *Clin Periodontol.* 41(11):1027-36.

596 Lovell, D., Pawlowsky-Glahn, V., Egozcue, J. J. Marguerat, S., Bähler, J. 2015.
597 Proportionality: a valid alternative to correlation for relative data. *PLoS Comput Biol*
598 11:e1004075.

599 Lozopone, C., Knight, R. 2005. Unifrac: a new phylogenetic method for comparing
600 microbial communities. *Applied Env. Micro.* 71:8228-8235.

- 601 Macklaim, J.M., Clemente, J.C., Knight, R., Gloor, G.B., Reid, G. 2015. Changes in
602 vaginal microbiota following antimicrobial and probiotic therapy. *Microb Ecol Health Dis.*
603 26:27799.
- 604 Mandal, S., Van Treuren, W., White, RA., and Eggesbø, M., Knight, R., Peddada, S. D.
605 2015. Analysis of composition of microbiomes: a novel method for studying microbial
606 composition. *Microl. Ecol. Health Dis.* 26:27663.
- 607 Martín-Fernández, J. A., Barceló-Vidal, C., & Pawlowsky-Glahn, V. 1998. Measures of
608 difference for compositional data and hierarchical clustering methods. In A. Buccianti,
609 G. Nardi, & R. Potenza (Eds.), *Proc. IAMG* (Vol. 98, pp. 526-531).
- 610 Palarea-Albaladejo J., Antoni Martín-Fernández, J. 2015. zCompositions --- R package
611 for multivariate imputation of left-censored data under a compositional approach.
612 *Chemometrics and Intelligent Laboratory Systems.* 143:85-96
- 613 Pawlowsky-Glahn, V., Egozcue, J. J., Tolosana-Delgado, R. 2015. Modeling and
614 Analysis of Compositional Data. John Wiley & Sons. Springer. 258 pg, London, UK.
- 615 Pearson, K. 1896. Mathematical contributions to the theory of evolution. -- on a form of
616 spurious correlation which may arise when indices are used in the measurement of
617 organs. *Proc. Royal Soc. Lond.* 60:489-498
- 618 R Core Team 2015. R: A language and environment for statistical computing. R
619 Foundation for Statistical Computing, Vienna, Austria. URL <https://www.R-project.org/>.
- 620 Rajca, S., Grondin, V., Louis, E., Vernier-Massouille, G., Grimaud, J.C., Bouhnik, Y.,
621 Laharie, D., Dupas, J.L., Pillant, H., Picon, L., Veyrac, M., Flamant, M., Savoye, G.,
622 Jian, R., Devos, M., Paintaud, G., Piver, E., Allez, M., Mary, J.Y., Sokol, H., Colombel,
623 J.F., Seksik, P. 2014. Alterations in the intestinal microbiome (dysbiosis) as a predictor

- 624 of relapse after infliximab withdrawal in Crohn's disease. *Inflamm Bowel Dis.* 20(6):978-
625 86.
- 626 Reardon, S. 2013, Bacterium can reverse autism-like behaviour in mice. *Nature.*
627 doi:10.1038/nature.2013.14308.
- 628 Schellenberg, J., Links, M. G., Hill, J. E., Dumonceaux, T. J., Peters, G. A., Tyler, S.,
629 Ball, T. B., Severini, A., Plummer, F. A. 2009. Pyrosequencing of the chaperonin-60
630 universal target as a tool for determining microbial community composition. *Appl*
631 *Environ Microbiol.* 75: 2889-98.
- 632 Schloss, P.D, Westcott, S.L, Ryabin, T., Hall, J.R., Hartmann, M., Hollister, E.B.,
633 Lesniewski, R.A., Oakley, B.B., Parks, D.H., Robinson, C.J., Sahl, J.W., Stres, B.,
634 Thallinger, G.G., and Van Horn, D.J., Weber, C.F. 2009. Introducing mothur: open-
635 source, platform-independent, community-supported software for describing and
636 comparing microbial communities
- 637 Seekatz, A.M., Aas, J., Gessert, C.E., Rubin, T.A., Saman, D.M., Bakken, J.S., Young,
638 V.B. 2014. Recovery of the gut microbiome following fecal microbiota transplantation.
639 *MBio.* 5(3):e00893-14.
- 640 Segata, N., Izard, J., Waldron, L., Gevers, D., Miropolsky, L., Garrett, W.S.,
641 Huttenhower, C. 2011. Metagenomic biomarker discovery and explanation. *Genome*
642 *Biol.* 12:R60
- 643 Urbaniak, C., Cummins, J., Brackstone, M., Macklaim, J.M., Gloor, G.B., Baban, C.K.,
644 Scott, L., O'Hanlon, D.M., Burton, J.P., Francis, K.P., Tangney, M., Reid, G. 2014.
645 Microbiota of human breast tissue. *Appl Environ Microbiol.* 80(10):3007-14. Van den

- 646 Boogaart, K. G., Tolosana-Delgado, R. 2013. Analyzing compositional data with R.
647 Heidelberg: Springer. Heidelberg 258 pages.
- 648 Walker, A. W., and Martin, J.C., Scott, P., Parkhill, J., Flint, H. J. Scott, K. P. 2015. 16S
649 rRNA gene-based profiling of the human infant gut microbiota is strongly influenced by
650 sample processing and PCR primer choice. *Microbiome*. 3:26
- 651 Warton, D.I., Wright, S.T., Wang, Y. 2012. Distance-based multivariate analyses
652 confound location and dispersion effects. *Methods Ecol. Evol.* 3:89-101.\
- 653 White, J.R., Nagarajan, N., Pop, M. 2009. Statistical methods for detecting differentially
654 abundant features in clinical metagenomic samples. *PLoS Comput. Biol.* 5:e1000352
- 655
- 656

Draft

657 **Figure Legends**

658 **Figure 1:** The difference between counting, proportions and ratios. The 'Counts' panel
659 shows a scatter plot of a simulated dataset with two samples composed of 49 invariant
660 taxa in open circles, and 1 taxon that changes in count 10-fold (black-filled circle). This
661 is the type of data that most current analysis tools in the microbiome field expect is
662 being analyzed. The 'Proportions' panel shows the same samples after they have been
663 sequenced and so constrained to have a constant sum. With such a constraint, their
664 representation is the same whether the sum is 1 (as shown here) or an arbitrarily larger
665 number (such as would be obtained from a sequencing instrument). The distortion in the
666 data is obvious: the black-filled circle still appears to be more abundant, but the open
667 circles appear to have become less abundant! It is obvious that we would draw incorrect
668 inferences regarding abundance changes in these data, yet these are the data as used
669 by existing tools. The third panel shows that much of this distortion can be removed
670 using a ratio transformation where each count (or proportion) is divided by the
671 geometric mean of the 50 taxa in the sample. Examination of the data after this
672 transformation can thus provide more robust inferences.

673 **Figure 2:** The left figure shows a covariance biplot of the abundance-filtered dataset,
674 the right figure shows a scree plot of the same data. This exploratory analysis is
675 encouraging, but not definitive, because the amount of variance explained is substantial
676 with 0.469 of the variance being explained by component 1, and 0.139 being explained
677 by component 2. The numbers on the left and right indicated unit-scaled variance of the
678 taxa, the numbers on the top and right indicate unit scaled variances of the samples.
679 Samples are colored in red if diagnosed as BV, blue if healthy, and purple if

680 intermediate. The scree plot also shows that the majority of the variability is on
681 component 1. We can interpret this biplot with some confidence, although it is likely that
682 any associations will be found to have large variation.

683 **Figure 3:** Unsupervised clustering of the reduced dataset. The top figure shows a
684 dendrogram of relatedness generated by unsupervised clustering of the Aitchison
685 distances, which is a distance that is robust to perturbations and sub-compositions of
686 the data (Aitchison 1986). The bottom figure shows a stacked bar plot of the samples in
687 the same order. The legend indicating the colour scheme for the taxa is on the right side.

688 **Figure 4:** An effect plot showing the univariate differences between groups (Gloor et al.
689 2015). The left plot shows a plot of the maximum variance within the B or A group vs.
690 the difference between groups. Large black points indicate those that have a mean
691 Benjamini-Hochberg adjusted P-value of 0.05 or less using P values calculated with the
692 Wilcoxon rank test. The middle plot shows a plot of the effect size vs. the adjusted P
693 value. In general, effect size measures are more robust than are P values and are
694 preferred. For a large sample size such as this one, an effect size of 0.5 or greater will
695 likely correspond to biological relevance. The right plot shows a volcano plot where the
696 difference between groups is plotted vs the adjusted P value.

697 **Figure 5:** A form biplot of the Hsiao et al. (2013) dataset that best represents the
698 distances between samples. Here we can see that the control and experimental
699 samples are intermingled, suggesting no separation between the groups. Furthermore,
700 the proportion of variance explained in the first component is not large when compared
701 to the other components. The evidence of structure within this dataset is thus weak.

702

Draft

703 Table 1: List of significantly different taxa.

| Taxon | diff.btw | diff.win | effect | overlap | wi.ep | wi.eBH |
|----------------------|----------|----------|--------|---------|-------|--------|
| <i>Atopobium</i> | 0.86 | 1.51 | 0.53 | 0.30 | 0.007 | 0.037 |
| <i>Prevotella</i> | 1.41 | 1.77 | 0.75 | 0.22 | 0.000 | 0.002 |
| <i>L. crispatus</i> | -1.07 | 1.78 | -0.49 | 0.23 | 0.000 | 0.004 |
| <i>L. iners</i> | -2.25 | 2.68 | -0.79 | 0.20 | 0.000 | 0.001 |
| <i>Streptococcus</i> | -1.14 | 2.38 | -0.37 | 0.30 | 0.008 | 0.041 |
| <i>Dialister</i> | 0.89 | 1.38 | 0.59 | 0.25 | 0.001 | 0.009 |
| <i>Megasphaera</i> | 1.56 | 2.31 | 0.63 | 0.28 | 0.002 | 0.015 |

704 diff.btw: median difference between groups on a log base 2 scale

705 diff.win: largest median variation within group H or BV

706 effect: effect size of the difference, median of diff.btw/diff.win

707 overlap: confusion in assigning an observation to H or BV group. Smaller is better

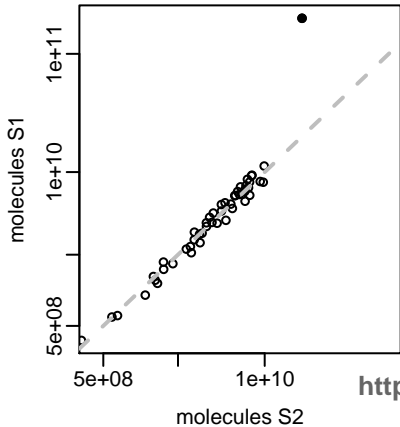
708 wi.ep: expected value of the Wilcoxon Rank Test P-value

709 wi.eBH: expected value of the Benjamini-Hochberg corrected P-value

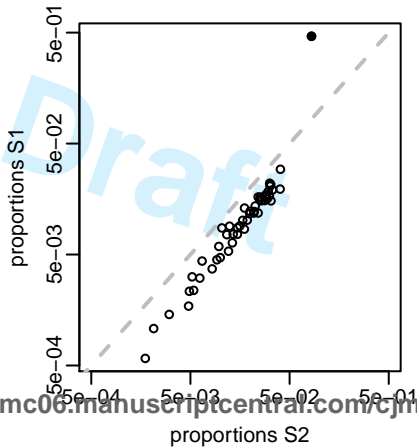
710

711

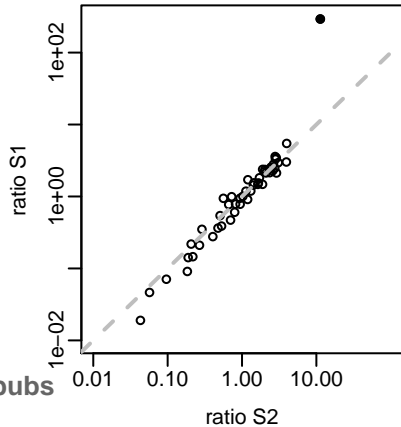
Counts



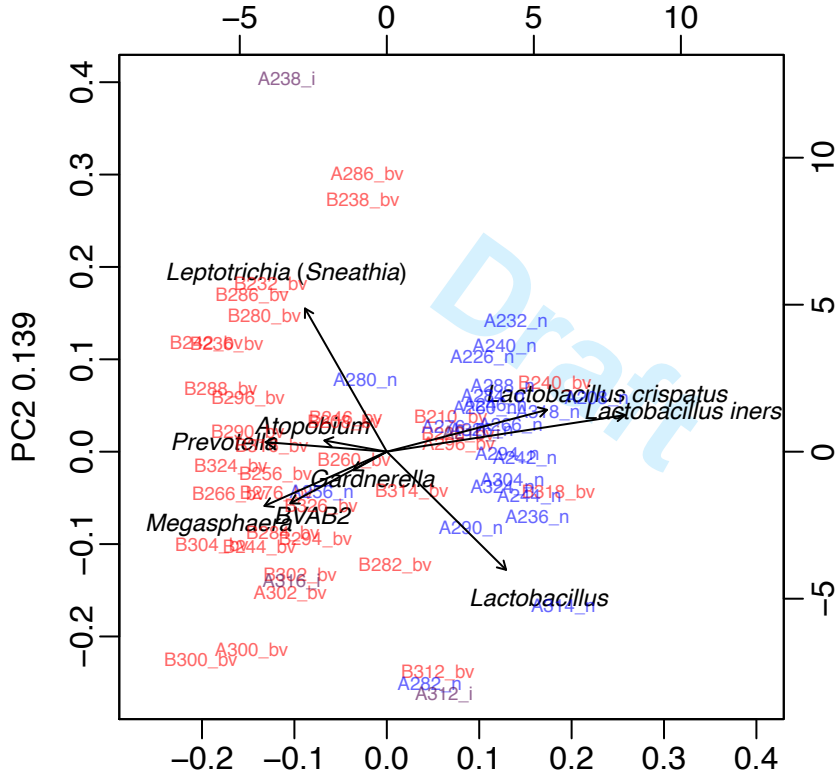
Proportions



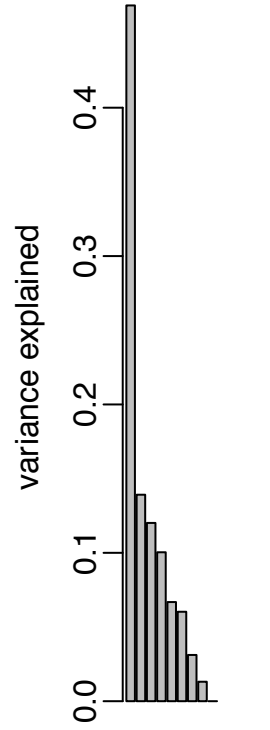
Ratios

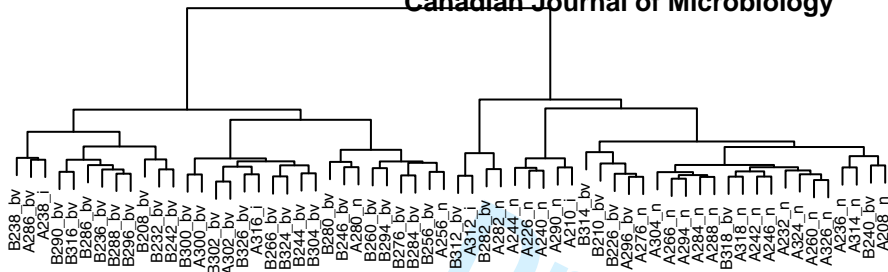


Biplot



Scree plot





- Firmicutes:*Lactobacillus iners*
- Actinobacteria:*Gardnerella*
- Bacteroidetes:*Prevotella*
- Firmicutes:*Megasphaera*
- Firmicutes:*Lactobacillus crispatus*
- Fusobacteria:*Leptotrichia*
- Actinobacteria:*Atopobium*
- Firmicutes:*Lactobacillus*
- Firmicutes:BVAB2

