1  **IDENTIFICATION OF DYNAMIC MICROBIAL SIGNATURES IN**
2  **LONGITUDINAL STUDIES**

3  Short running title: Microbiome signatures in longitudinal studies

4  M.Luz Calle[1]*, Antoni Susin[2]

5  [1] Biosciences Department, Faculty of Sciences and Technology, University of Vic -

6  Central University of Catalonia, Vic, Spain

7  [2] Mathematical Department, UPC-Barcelona Tech, Barcelona, Spain

8  *Corresponding author: M.Luz Calle, Biosciences Department, Faculty of Sciences and

9  Technology, University of Vic - Central University of Catalonia, Carrer de la Laura, 13,

10  08500-Vic, Spain

11  Tel: +34938861222  Email: malu.calle@uvic.cat

12

## Abstract

14 The study of microbiome dynamics is key for unveiling the role of the microbiome in

15 human health. Addressing the compositional structure of microbiome data is

16 particularly critical in longitudinal studies where compositions measured at different

17 times can yield to different subcompositions.

18 We propose a new compositional data analysis (CoDA) algorithm for inferring dynamic

19 microbial signatures. The algorithm performs penalized regression over the summary of

20 the log-ratio trajectories (the area under these trajectories) and the inferred microbial

21 signature is expressed as a log-contrast model. Graphical representations of the results

22 are provided to facilitate the interpretation of the analysis: plot of the log-ratio

23 trajectories, plot of the signature and plot of the prediction accuracy of the model. The

24 new proposal is illustrated with data on the developing microbiome of infants.

25 The algorithm is implemented in the R package "code4microbiome" (https://cran.r-

26 project.org/web/packages/coda4microbiome/) that is accompanied with a vignette with

27 a detailed description of the functions. The website of the project contains several

28 tutorials: https://malucalle.github.io/coda4microbiome/

29

30 **Key-words:** log-ratio analysis, longitudinal studies, microbiome analysis, microbiome

31 signatures, penalized regression

32

33

34

35

## 1. Introduction

Microbiome composition is dynamic and the study of microbiome changes over time is of primary importance for understanding the relationship between microbiome and human phenotypes. Longitudinal studies are costly, both economically and logistically, but there is growing evidence of the limitations of cross-sectional studies for providing a full picture of the role of the microbime in human health. Microbiome longitudinal studies can be very valuable in this context, provided appropriate methods of analysis are used (Schmidt et al. 2018)

Microbiome data analysis is challenging because, among other things, the compositional nature of the data (Susin et al. 2020, Calle 2019, Gloor et al. 2016, 2017, Gloor and Reid, 2016). This is particularly critical in the context of microbiome longitudinal studies where compositions measured at different times can be affected by distinct batch effects and similar quality control or filtering protocols may yield to different subcompositions at each time point.

The log-ratio approach (Aitchison 1986), that consists in analyzing the abundances of some taxa relative to the abundances of other taxa, is subcompositionally coherent and provides an especially interesting standpoint for exploring microbiome dynamics. In longitudinal studies, the log-ratio between two groups of taxa measured at different time points gives a curve profile or trajectory for each sample. We propose to explore the association between the phenotype of interest and the shape of the log-ratio microbiome trajectories.

Among the questions outstanding about microbiome dynamics, we focus on inferring dynamic microbial signatures and propose a novel algorithm to identify a set of microbial taxa whose joint dynamics is associated with the phenotype of interest. For

60    binary outcomes, such as disease status, we aim to identify two groups of taxa with

61    clearly different log-ratio trajectories for cases and controls.

62    The algorithm performs variable selection through penalized regression over the

63    summary of the log-ratio trajectories (the area under these trajectories). The inferred

64    microbial signature is expressed as a log-contrast model (Aitchison, J. and Bacon-

65    Shone,J. 1984), i.e. a log-linear model with the constraint that the sum of the

66    coefficients is equal to zero. The zero-sum constraint ensures the invariance principle

67    required for compositional data analysis.

68    The interpretability of results is of major importance in the context of microbiome

69    studies. We provide several graphical representations of the results that facilitate the

70    interpretation of the analysis: plot of the log-ratio trajectories, plot of the signature

71    (selected taxa and coefficients) and plot of the prediction accuracy of the model.

72    The methodology is illustrated with data from the "Early childhood and the microbiome

73    (ECAM) study" (Bokulich et al. 2016).

74    The algorithm is implemented in the R package "code4microbiome" (https://cran.r-

75    project.org/web/packages/coda4microbiome/) that is accompanied with a vignette with

76    a detailed description of the functions. The website of the project contains several

77    tutorials: https://malucalle.github.io/coda4microbiome/

78

79

## 2. Materials and methods

We first describe the analysis of log-ratios between two taxa A and B in longitudinal studies, which involves the summary of the log-ratio trajectories. Then we explain how to generalize the analysis of pairwise log-ratios to identify microbial signatures involving more than two taxa.

**Log-ratio analysis and taxa prioritization**

Assume $n$ subjects with fixed phenotype $Y = (Y_1, \ldots, Y_n)$. Subject $i$ has been observed in $L_i$ time points, $(t_{i1}, t_{i2}, \ldots, t_{iL_i})$. We denote by $X_i(t_{ij}) = (X_{i1}(t_{ij}), X_{i2}(t_{ij}), \ldots, X_{iK}(t_{ij}))$ the microbiome composition of subject $i$ at time $t_{ij}$, where $K$ is the number of taxa which is assumed to be the same for all the individuals and all the time points. $X_i(t_{ij})$ can represent either relative abundances (proportions) or raw counts. We denote by $logX_i(t_{ij})$ the logarithm transformation of microbiome abundances after zero imputation. The log-abundance trajectory of component A for individual $i$ is denoted by $logX_{iA} = (logX_{iA}(t_{i1}), logX_{i2A}(t_{i2}), \ldots, logX_{iA}(t_{iL_i}))$ and the log-ratio trajectory between components A and B for individual $i$ is given by:

$$logX_{iA} - logX_{iB} = (logX_{iA}(t_{i1}) - logX_{iB}(t_{i1}),$$

$$logX_{i2A}(t_{i2}) - logX_{iB}(t_{i2}), \ldots, logX_{iA}(t_{iL_i}) - logX_{iB}(t_{iL_i}))$$

We summarize the log-ratio trajectories within two time points $l_1$ and $l_2$ as the integral of the log-ratio trajectory:

$$s_i(A, B) = \int_{l_1}^{l_2} (logX_{iA}(t) - logX_{iB}(t)) \, dt \tag{1}$$

where the values of the log-ratio for $t \notin (t_{i1}, t_{i2}, \ldots, t_{iL_i})$ are linearly interpolated.

101 We do not take the absolute value in equation (1) because the sign of the integral is

102 informative: Positive values of $s_i(A, B)$ correspond to trajectories of component A

103 above trajectories of component B, that is, larger relative abundances of A with respect

104 to B, while negative values represent the opposite. Values of $s_i(A, B)$ around zero can

105 represent similar abundances between A and B over time or a non-homogeneous trend

106 between A and B within the observed region.

107 Another advantage of the summary $s_i(A, B)$ is computational. Since the integral is

108 linear, $s_i(A, B)$ is equal to the difference between the integrals of log-transformed

109 microbiome abundances of taxa A and taxa B:

110
$$s_i(A, B) = \int_{l_1}^{l_2} logX_{iA}(t)\, dt - \int_{l_1}^{l_2} logX_{iB}(t)\, dt$$

111 Thus, the number of integrals to be calculated is of the order of $K$, the number of taxa,

112 instead of $K(K - 1)/2$, the number of pairwise log-ratios.

113 The log-ratio summary for the $n$ subjects, $s(A, B) = (s_1(A, B), \dots, s_n(A, B))$, can be

114 tested for association with the phenotype $Y = (Y_1, \dots, Y_n)$ with a generalized linear

115 model (glm) adjusted for some covariates Z:

116
$$g\big(E(Y_i)\big) = \beta_0 + \beta_1 \cdot s_i(A, B) + \gamma' \cdot Z_i \qquad (2)$$

117 where $\beta_0$ is the intercept, $\beta_1$ is the regression coefficient for the log-ratio summary

118 between components $A$ and $B$, $Z = (Z_1, Z_2, \dots, Z_r)$ are non-compositional covariates and

119 $\gamma$ is the vector of regression coefficients for Z.

120 **Microbiome signature based on log-ratio analysis**

121 To identify those log-ratios that are most associated with the outcome $Y$, we implement

122 glm penalized regression on the log-ratio summaries for all pairs of taxa:

123
$$g\big(E(Y)\big) = \beta_0 + \sum_{j \in J} \beta_j \cdot S(j) \qquad (3)$$

124  where $J = \{1, \dots, K(K-1)/2\}$ and $S(j) = s(j_1, j_2)$ is the log-ratio summary of

125  components $j_1$ and $j_2$ with $(j_1, j_2) \in J_{12}$, the set of all possible combinations of pairs of

126  components.

127  The regression coefficients in equation (3) are estimated to minimize a loss function

128  $L(\beta)$ subject to a penalization on the regression coefficients, $P(\beta)$:

129
$$\hat{\beta} = \underset{\beta}{\mathrm{argmin}}\{L(\beta) + P(\beta)\} \qquad (4)$$

130  For the penalty term we consider the elastic-net, which combines the L1 and L2 norms:

131  $P(\beta) = \lambda_1 \|\beta\|_2^2 + \lambda_2 \|\beta\|_1$. A common reparameterization of $P(\beta)$ is $\lambda_1 = \lambda(1-\alpha)/2$

132  and $\lambda_2 = \lambda\alpha$ where $\lambda$ controls the amount of penalization and $\alpha$ the mixing between the

133  two norms.

134  For the linear regression model the loss function is given by the residual sum of squares

135
$$\hat{\beta} = \underset{\beta}{\mathrm{argmin}}\{\|Y - S\beta\|_2^2 + \lambda_1 \|\beta\|_2^2 + \lambda_2 \|\beta\|_1\},$$

136  where $S$ is the matrix of all log-ratio summaries and has dimension $n$ by $K(K-1)/2$.

137  The expression of the optimization problem (4) for other models, like the logistic

138  regression and the multinomial regression models, can be found in Friedman et al.

139  (2010). Non-compositional covariates $Z$ are previously modeled with Y and the fitted

140  values are considered as "offset" in the penalized regression.

141  The result of the penalized optimization provides a set of selected pairs of taxa, those

142  with a non-null estimated coefficient. For each individual, $i \in \{1, \dots, n\}$, the linear

143  predictor of the generalized linear model (3), $M_i = \sum_{j \in J, (j_1,j_2)=J_{12}(j)} \widehat{\beta_J} \cdot s_i(j_1, j_2)$, is the

144  microbiome signature which is associated with the phenotype $Y_i$. Because of the

145  linearity of the integrals used as summaries of the log-ratio trajectories, the microbiome

146  signature $M$ can be rewritten in terms of the selected single taxa which is more

147  interpretable than the selected pairs of components:

148
$$M = \sum_{j \in J, (j_1, j_2) = J_{12}(j)} \widehat{\beta}_J \cdot s(j_1, j_2) =$$

149
$$= \sum_{j \in J, (j_1, j_2) = J_{12}(j)} \widehat{\beta}_J \cdot \int_{l_1}^{l_2} logX_{j_1}(t) \, dt - \sum_{j \in J, (j_1, j_2) = J_{12}(j)} \widehat{\beta}_J \cdot \int_{l_1}^{l_2} logX_{j_2}(t) \, dt =$$

150
$$= \sum_{k=0}^{K} \widehat{\theta}_k \cdot \int_{l_1}^{l_2} logX_k(t) \, dt =$$

151
$$= \int_{l_1}^{l_2} \left( \sum_{k=0}^{K} \widehat{\theta}_k \cdot logX_k(t) \right) dt$$

152
$$(5)$$

153    where $\widehat{\theta}_k = \sum_{j:k \in J_{12}(j)} \widehat{\beta}_J$ , that is, the sum of the coefficients $\widehat{\beta}_J$ corresponding to a log-

154    ratio that involves component $k$.

155    It can be proved that $\sum_{k=0}^{K} \widehat{\theta}_k = 0$ and thus, the microbiome signature $M$ is the integral

156    of the trajectory of a log-contrast function involving the selected taxa (those with $\widehat{\theta}_k \neq$

157    0). This ensures the invariance principle required for proper compositional data analysis

158    and it facilitates the interpretation of the microbiome signature: Expression $\sum_{k=0}^{K} \widehat{\theta}_k \cdot$

159    $logX_k(t)$ in (5) can be interpreted as a weighted balance between two groups of taxa,

160    $G_1$ and $G_2$, the taxa with a positive coefficient vs those with a negative coefficient

161    (Susin et al. 2020).

162    The package "coda4microbiome" (Calle and Susin, 2022) contains several functions

163    that implement the proposed algorithms. The two main functions are

164    explore_lr_longitudinal(), that implements the simple generalized linear model

165    (equation 2), and coda_glmnet_longitudinal(), that performs penalized regression for

166    the multivariable generalized linear model (equation 4). Additional functions are

167    available like function plot_signature_curves() that provides a plot of the signature

168    trajectories or `filter_longitudinal()` that filters those individuals and taxa with

169    enough longitudinal information.

170    To illustrate the proposed approach and the R implementation we use data from the

171    early childhood and the microbiome (ECAM) study (Bokulich et al. 2016). Metadata

172    and microbiome data were downloaded from https://github.com/caporaso-

173    lab/longitudinal-notebooks. Microbiome data, corresponding to 16S rRNA gene

174    microbiota compositions sampled at regular intervals, were available in QIIME 2 qza

175    file format (file ecam-table-genus.qza) and were transformed to R objects with function

176    `read_qza()` of the R library qiime2R: https://github.com/jbisanz/qiime2R. Metadata

177    (file ecam-sample-metadata.tsv) were in long format: multiple rows for individual, one

178    for each time-point observation. Initially the data contained information on 43 child and

179    445 taxa at the genus level. We filtered those individuals and taxa with enough

180    information for time-course profiling: we removed individuals with only one time-point

181    observation and those taxa with less than 30 children (70% of individuals) with at least

182    3 non-zero observations over the follow-up period. After filtering, the data reduced to

183    42 children and 37 taxa.

184    **3.  Results**

185    We demonstrate the proposed methodology with data from the "Early childhood and the

186    microbiome (ECAM) study" that followed a cohort of 43 U.S. infants during the first 2

187    years of life for the study of their microbial development and its association with early-

188    life antibiotic exposures, cesarean section, and formula feeding (Bokulich et al. 2016).

189    Microbiome data were available for 43 child and 445 taxa at the genus level (Bokulich

190    et al. 2018). After filtering those individuals and taxa with enough information for time-

191    course profiling, the data were reduced to 42 child and 37 taxa. We focus on the effects

192     of the diet on the early modulation of the microbiome by comparing microbiome

193     profiles between children with breastmilk diet (bd) vs. formula milk diet (fd) in their

194     first 3 months of life.

195     **Most important taxa**

196     By implementing the pairwise log-ratio approach for longitudinal data (function

197     `explore_lr_ongitudinal()`), we identified which taxa have more different dynamics

198     between bd and fd children in the first three months of life. Table 1 provides the top 15

199     taxa with more discriminative dynamics between both diets.

200     Table 1. Taxa with most different abundances between the two diets groups during the first
201                                     three months of life.
202

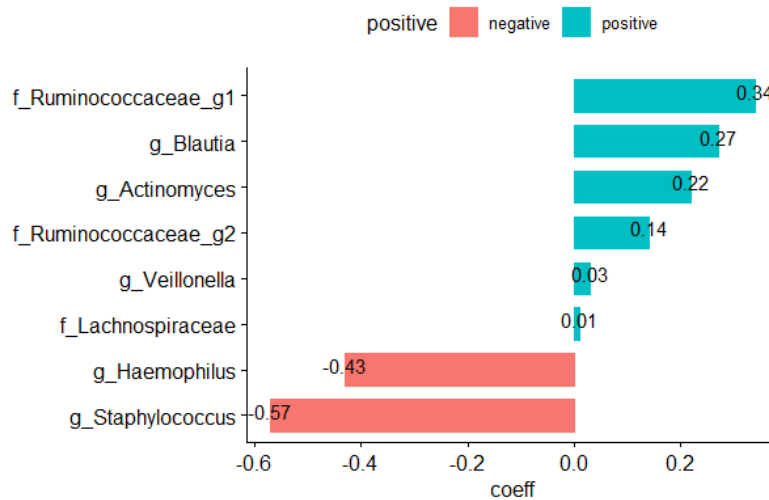| Taxanomic assignment | More abundant group |
|---|---|
| "p_Proteobacteria;c_Gammaproteobacteria;o_Pasteurellales;f_Pasteurellaceae;g_Haemophilus" | bd |
| "p_Firmicutes;c_Bacilli;o_Bacillales;f_Staphylococcaceae;g_Staphylococcus" | bd |
| "p_Firmicutes;c_Clostridia;o_Clostridiales;f_Veillonellaceae;g_Veillonella" | fd |
| "p_Firmicutes;c_Clostridia;o_Clostridiales;f_Lachnospiraceae;g_Blautia" | fd |
| "p_Firmicutes;c_Clostridia;o_Clostridiales;f_Ruminococcaceae;g_1" | fd |
| "p_Firmicutes;c_Clostridia;o_Clostridiales" | fd |
| "p_Firmicutes;c_Clostridia;o_Clostridiales;f_Clostridiaceae" | fd |
| "p_Actinobacteria;c_Actinobacteria;o_Bifidobacteriales;f_Bifidobacteriaceae;g_Bifidobacterium" | bd |
| "p_Firmicutes;c_Clostridia;o_Clostridiales;f_Lachnospiraceae" | fd |
| "p_Bacteroidetes;c_Bacteroidia;o_Bacteroidales;f_Bacteroidaceae;g_Bacteroides" | bd |
| "p_Firmicutes;c_Bacilli;o_Lactobacillales;f_Lactobacillaceae;g_Lactobacillus" | bd |
| "p_Firmicutes;c_Erysipelotrichi;o_Erysipelotrichales;f_Erysipelotrichaceae;g_[Eubacterium]" | fd |
| "p_Firmicutes;c_Clostridia;o_Clostridiales;f_Lachnospiraceae;g_Coprococcus" | fd |
| "p_Firmicutes;c_Clostridia;o_Clostridiales;f_Lachnospiraceae;g_Dorea" | fd |
| "p_Firmicutes;c_Bacilli;o_Lactobacillales;f_Enterococcaceae;g_Enterococcus" | fd |

203


204

**Microbiome signature**

206  The application of the proposed algorithm (with function `coda_glmnet_longitudinal()`)

207  identified a microbiome signature with maximum discrimination accuracy between the

208  two diet groups. The signature is defined by the relative abundances of two groups of

209  taxa, $G_1$ and $G_2$, where $G_1$ is composed of 6 taxa (those with a positive coefficient in the

210  regression model) and $G_2$ is composed of 2 taxa (those with a negative coefficient)

211  (Table 1 and Figure 1). Group $G_1$ is mainly dominated by three taxa within the order

212  *Clostridiales* (family *Ruminococcaceae* (2) and gender *Blautia*) and one taxon within

213  the gender *Actinomyces*. Two taxa (*g_Veillonella* and *f_Lachnospiraceae*) have a

214  coefficient close to zero and will have a very small contribution to the signature. Group

215  $G_2$ is composed by two taxa within the genders *Haemophilus* and *Staphylococcus*. Note

216  that the selected taxa within the microbial signature are among most important taxa

217  according to the results of the pairwise log-ratio analysis (Table 1).

218

219  Table 2. Taxa included in the microbiome signature that best discriminates between the two diet
220                                                             groups

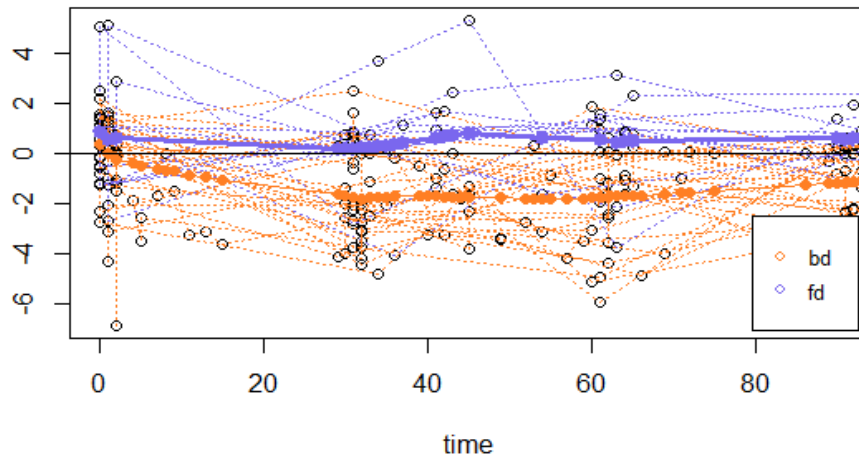| Balance group | Coefficient | Taxanomic assignment |
|---|---|---|
| $G_1$ | 0.3359 | *p_Firmicutes;c_Clostridia;o_Clostridiales;f_Ruminococcaceae;g_1* |
| | 0.2730 | *p_Firmicutes;c_Clostridia;o_Clostridiales;f_Lachnospiraceae;g_Blautia* |
| | 0.2159 | *p_Actinobacteria;c_Actinobacteria;o_Actinomycetales;f_Actinomycetaceae;g_Actinomyces* |
| | 0.1358 | *p_Firmicutes;c_Clostridia;o_Clostridiales;f_Ruminococcaceae;g_2* |
| | 0.0337 | *p_Firmicutes;c_Clostridia;o_Clostridiales;f_Veillonellaceae;g_Veillonella* |
| | 0.0055 | *p_Firmicutes;c_Clostridia;o_Clostridiales;f_Lachnospiraceae;g_* |
| $G_2$ | −0.4327 | *p_Proteobacteria;c_Gammaproteobacteria;o_Pasteurellales;f_Pasteurellaceae;g_Haemophilus* |
| | −0.5672 | *p_Firmicutes;c_Bacilli;o_Bacillales;f_Staphylococcaceae;g_Staphylococcus* |

221

222

Fig 1. Taxa composing the microbiome signature that best discriminates between the two diet groups (green: positive coefficient and red: negative coefficient)
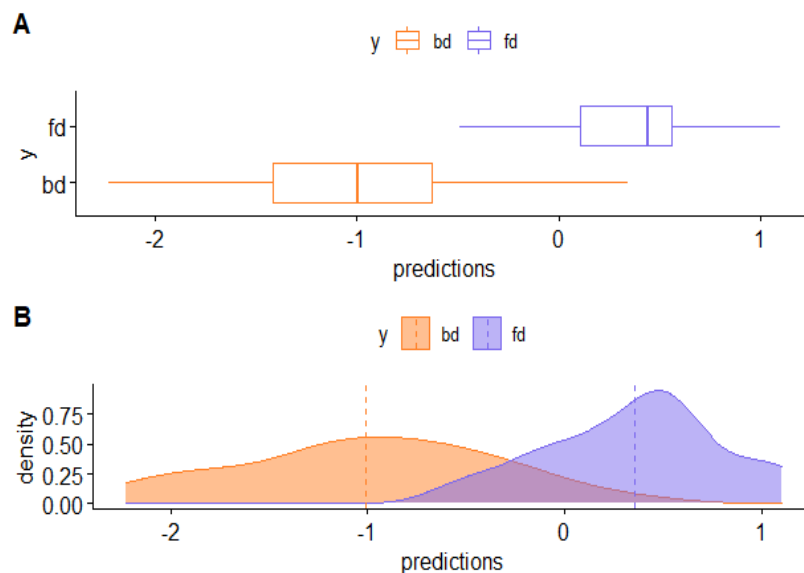
223
224
225

226    The trajectories of the microbial signature over the observed period are represented in

227    Figure 2, where the color of the curves corresponds to the diet group. Each trajectory

228    represents the relative mean abundances between the two taxa groups for each child. We

229    can see that the two groups are clearly separated. Those children under breastmilk diet

230    (in orange) usually have trajectories below zero, which means they have more relative

231    mean abundance of *g_Haemophilus* and *g_Staphylococcus* with respect to the relative

232    abundance of taxa in group $G_1$, while children with formula milk diet (in blue) have

233    more relative abundance of taxa in group $G_1$ relative to $G_2$.

234

Fig 2. Relative abundance between group $G_1$ and $G_2$ during the first three months of life. Highlighted curves represent the mean value of the signature for each diet group (orange: breast milk diet, blue: formula milk diet)

Figure 3 displays the distribution of the microbial signature scores for the two diet groups and offers a visual assessment of the (apparent) discrimination accuracy of the signature. Quantitatively, the apparent discrimination accuracy of the signature (i. e. the AUC of the signature applied to the same data that was used to generate the model) is 0.96 and the mean cross-validation AUC is 0.74 (sd=0.10).

244



245

Fig 3. Distribution of the microbial signature scores for the two diet groups (orange: breast milk diet, blue: formula milk diet)

248
249

13

250 Both results, the pairwise analysis and the taxa selected in the microbial signature, are

251 consistent with previous studies on the association of the infant gut microbiome

252 composition and breastmilk feeding practices. In Fehr et al. (2020), *Haemophilus*

253 *parainfluenzae* and *Staphylococcus* were found to be enriched with exclusive breastmilk

254 feeding together with lower prevalence of *Actinomyces* at 3 months. *Lachnospiraceae*

255 (*Blautia*) was enriched among infants who were no longer fed breastmilk. Similar

256 results are reported in Laursen et al. (2016) where the duration of exclusive

257 breastfeeding was negatively correlated with genera within *Lachnospiraceae* (e.g.,

258 *Blautia*) and genera within *Ruminococcaceae*. Positive correlations with exclusive

259 breastfeeding were observed for *g_Bifidobacterium* and *Pasteurellaceae*

260 (*Haemophilus*).

261 **4. Discussion**

262 Longitudinal microbiome studies, especially those focused on the human microbiome,

263 have usually low resolution: the number of individuals is small, each individual has few

264 observation times, the observations of the different individuals are not made at exactly

265 the same time, the data are very variable, the expected behavior of the abundance

266 trajectories is not linear or quadratic, etc. This makes it difficult to justify and

267 implement a parametric modeling of trajectories and limits the use of models for

268 longitudinal data (time series, mixed models). In this context, a description of the

269 trajectories such as the one we propose, although less precise, allows to extract valuable

270 information from the data as we have shown in the example. Other longitudinal data

271 modeling strategies (Gerberg et al. 2012, Park et al. 2020, Silverman et al. 2018, Äijö et

272 al. 2017) could be used in longitudinal microbiome studies with higher resolution such

273 as laboratory or animal experimental studies.

274 The applicability of CoDA methods in microbiome studies has been limited by the

275 difficulty in interpreting the obtained results. We hope that this work and the R package

276 "coda4microbiome" will help to increase the use of these methods in this field.

**Acknowledgements**

**Data Accessibility**

281 The filtered data from the ECAM study is available as a data object in the

282 "coda4microbiome" package.

**References**

284 1. Äijö T, Müller CL and Bonneau R. Temporal probabilistic modeling of bacterial

285 compositions derived from 16S rRNA sequencing. Bioinformatics, 34(3), 2018,

286 372–380 doi: 10.1093/bioinformatics/btx549

287 2. Aitchison J. The Statistical Analysis of Compositional Data. London: Chapman &

288 Hall, 1986.

289 3. Aitchison J. and Bacon-Shone J. Log contrast models for experiments with

290 mixtures. Biometrika. 1984, 71: 323–330.

291 4. Bokulich NA, Chung J, Battaglia T, Henderson N, Jay M, Li H, Lieber AD, Wu F,

292 Perez-Perez GI, Chen Y, Schweizer W, Zheng X, Contreras M, Dominguez-Bello

293 MG, Blaser MJ. Antibiotics, birth mode, and diet shape microbiome maturation

294 during early life. Sci Transl Med. 2016, 8:343ra82.

295 https://doi.org/10.1126/scitranslmed.aad7121

296 5. Bokulich NA, Dillon MR, Zhang Y, Rideout JR, Bolyen E, Li H, Albert PS,

297 Caporaso JG. q2-longitudinal: longitudinal and paired-sample analyses of

298      microbiome data. mSystems. 2018, 3:e00219-18.

299      https://doi.org/10.1128/mSystems.00219-18.Microbiome Data

300  6.  Calle ML. Statistical analysis of metagenomics data. Genomics Inform. 2019, 17(1):

301      e6

302  7.  Calle ML, Susin A, coda4microbiome R-package (CRAN). 2022

303  8.  Fehr K, Moossavi S, Sbihi H, Finlay B, Turvey SE, Azad MB. Breastmilk Feeding

304      Practices Are Associated with the Co-Occurrence of Bacteria in Mothers' Milk and

305      the Infant Gut: the CHILD Cohort Study. Cell Host & Microbiome. 2020,

306      28(2):285-297.e4 https://doi.org/10.1016/j.chom.2020.06.009

307  9.  Gerber GK, Onderdonk AB, Bry L. Inferring Dynamic Signatures of Microbes in

308      Complex Host Ecosystems. PLoS Comput Biol. 2012, 8(8): e1002624.

309      https://doi.org/10.1371/journal.pcbi.1002624

310  10. Gloor, Gregory B and Wu, Jia Rong and Pawlowsky-Glahn, Vera and Egozcue,

311      Juan José It's all relative: analyzing microbiome data as compositions. Annals.

312      Epidemiology. 2016, 26(5):322-9.

313      http://dx.doi.org/10.1016/j.annepidem.2016.03.003

314  11. Gloor GB. and Reid G. Compositional analysis: a valid approach to analyze

315      microbiome high throughput sequencing data. Can J Microbiol. 2016, 62(8):692–

316      703. http://dx.doi.org/10.1139/cjm-2015-0821

317  12. Laursen MF, Andersen LBB, Michaelsen KF, Mølgaard C, Trolle E, Bahl MI, Licht

318      TR. Infant gut microbiota development is driven by transition to family foods

319      independent of maternal obesity. MSphere. 2016, 1(1): e00069-15.

320      doi:10.1128/mSphere.00069-1

321   13. Park Y, Ufondu A, Lee K, Jayaraman A, Emerging computational tools and models

322       for studying gut microbiota composition and function, Current Opinion in

323       Biotechnology. 2020, 66: 301-311. https://doi.org/10.1016/j.copbio.2020.10.005.

324   14. Schmidt T, Raes J., Bork P. The Human Gut Microbiome: From Association to

325       Modulation, Cell. 2018, 172: 1198-1215. https://doi.org/10.1016/j.cell.2018.02.044

326   15. Silverman JD, Durand HK, Bloom RJ, Mukherjee S, David LA. Dynamic linear

327       models guide design and analysis of microbiota studies within artificial human guts.

328       Microbiome. 2018, 6:202.  https://doi.org/10.1186/s40168-018-0584-3

329   16. Susin A, Wang Y, Lê Cao KA, Calle ML. Variable selection in microbiome

330       compositional data analysis. NAR Genomics and Bioinformatics. 2020, 2 (2):

331       lqaa029, https://doi.org/10.1093/nargab/lqaa029