

# Estimation of Bacterial Functions in Samples Based on 16S rRNA

Bc. Michaela Bielíková\*



## Abstract

Humans are hosts to an enormous variety of microbes, bacterial, archaeal, fungal, and viral. Unfortunately, science knows only little about them. Since most of the bacteria has not been studied yet, the main question for a given sample is not only which species of bacteria a specific sample contains, but also what can the bacteria in this sample do (i.e. lipid digestion or resistance to antibiotics). This task is called functional profile prediction and it will be the main focus of this paper. In this paper, I introduce methods for functional analysis, describe existing tools and then design a new tool inspired by them, which implements different methods for the prediction. The results of the experiments imply, that the implemented tool is accurate and useful when using the same method for experimental evaluation as existing tools. However, I propose a new approach to evaluation, that concerns only the most specific bacterial functions, where the results differ from the classic one. In the end, I discuss possible implications of this difference.

**Keywords:** bioinformatics — metagenomics — bacterial functional profile — KO profile — 16S rRNA — PiCRUST

**Supplementary Material:** [Github repository of the created tool](#)

\*[xbieli06@stud.fit.vutbr.cz](mailto:xbieli06@stud.fit.vutbr.cz), Faculty of Information Technology, Brno University of Technology

## 1. Introduction

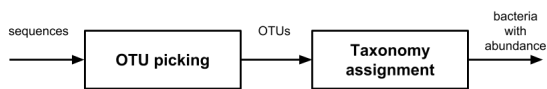
Humans are hosts to an enormous variety of microbes. Some of these are invaders that can cause serious diseases, but there is a lot of microbes that are essential to human life. Particularly gut microbiome is crucial for the regular function of the digestion tract. In the last years, it was proven that irregularities in gut microbiome are linked to many conditions ranging from digestion tract diseases like inflammatory bowel disease to antibiotic resistant infections [1]. Unfortunately, because of a big variety of present bacterial species and the impossibility to cultivate most of them in laboratories, the gut microbiome is not well described. Modern approaches in microbiology, specifically high-

throughput sequencing and metagenomics, seem to be able to solve these problems and allow us to study microbiome thoroughly and understand how it is connected to human health [1, 2].

Since most of the bacteria present in gut microbiome has not been studied yet, the main question is not which species of bacteria a specific sample contains, as we lack named species for most of the bacteria present in the sample, but instead what can the bacteria do (i.e. lipid digestion or resistance to antibiotics). This task is called functional profile prediction and it will be the main focus of this paper. Functional profile prediction is based on the observation, that bacteria species with similar RNA sequence tend to have sim-

15  
16  
17  
18  
19  
20  
21  
22  
23  
24  
25  
26  
27  
28

**Figure 1.** Diagram showing steps of bacterial composition analysis



- It is present in every organism we want to study 71
- It is unique for every species 72
- It is similar for closely related species and different for non-related species 73 74

ilar functions, whereas between species with small RNA similarity the functional profile differs [1].

In this paper, we will introduce existing bioinformatics tools for functional profile prediction, namely PiCRUST, and Tax4Fun. We will discuss the different methods they use for prediction. Then we will design a new tool inspired by them, that will implement the classic methods for functional profile prediction, but should also include a new approach to this problem based on linear regression.

In Chapter 2, I will define theoretical background needed to understand this paper. In Chapter 3, I will discuss the created tool. Chapter 4 contains experimental evaluation of the tool. In the last Chapter, I will talk about the future work and possible extension of the created tool.

## 2. Theoretical background

In this section, we will define the field this paper relates to - metagenomics. Then, we will discuss the details of the functional composition analysis of a given sample. This section is based mostly on papers from Xochitl C. Morgan [1], Jay-Hyun Jo [3] and Andreas Hiergeist [2], where more detailed information can be found.

### 2.1 Metagenomics

Metagenomics is a study of genetic material recovered directly from samples. It does not require isolating the DNA of individual species, neither cultivating in laboratories.

There are two main types of analysis often performed in metagenomics. The first one is taxonomic, where the main question is: Which bacteria are present in the given sample? The second one, the main focus of this thesis, is functional: What can the bacteria in this sample do?

The steps of taxonomic analysis can be seen in Figure 1.

### 2.2 16S rRNA

To minimize the length of the DNA sequence that must be processed to determine the species and functional profile, only a part of genetic information, called marker gene, is used. Marker gene needs to have the following attributes:

For bacteria, a commonly used marker gene is 16S RNA. It contains conserved regions, that are consistent among all species, and variable regions, that are different. In the taxonomic analysis, we study and compare the variable regions to determine which species we are dealing with.

### 2.3 Functional profile analysis

In functional analysis, we want to find the different metabolic functions of organisms in the sample, as well as to estimate their abundance — how many organisms in the sample have this function. The process of functional analysis is shown in Figure 2. Functional profiles have the form of KO identifiers with abundance in the corresponding sequences. KO identifiers refer to molecular functions and can be found in the Kegg Orthology database [4]. The functional profiles represent which bacterial functions in what quantity are present in the given sample.

The input of functional analysis can be either raw sequences or already preprocessed and clustered data (OTU table — a matrix that gives the number of reads per sample per OTU). For the purpose of explaining all steps of the analysis, we will suppose we start with raw DNA (or RNA in microbiome research) sequences.

In the data preprocessing, DNA sequences that are very similar (95-98% similarity) are clustered into groups called OTUs. OTU stands for the operational taxonomic unit and is used as a synonym to species since we lack a real, named species corresponding to most of the clusters. Each OTU has an identifier, but these are generic and are not consistent among different samples, so to represent OTUs we use the nucleotide sequences.

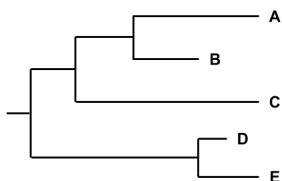
The basic principle of functional analysis is to compare representative sequences of OTUs to a reference database that contains the functional profile of previously studied organisms and find the best match. For the genetic material that cannot be paired with a known organism, we can search for the most similar organisms and deduce the functional profile from there.

A big question is, how to find the most similar organisms to a given sentence. Various methods for solving this problem exist. The naive algorithm is based on analyzing the sequences and finding the most similar one - these methods will be further called distance based. More advanced algorithms are based on

**Figure 2.** Diagram showing steps of functional analysis of a sample



**Figure 3.** Example of a phylogenetic tree



122 constructing a phylogenetic tree, that represents evolu-  
123 tionary relationships between the species, of all OTUs  
124 and deduces the estimated functional profile from the  
125 tree structure.

### 126 2.3.1 Functional analysis methods

127 In this section, I will describe two distinctive groups  
128 of methods for functional analysis.

#### 129 **Phylogenetic tree based algorithms**

130 This group of methods is commonly used in bioin-  
131 formatic tools for functional analysis. It is based on  
132 constructing a phylogenetic tree which is a graph that  
133 represents evolutionary relations between organisms.  
134 Each node of such a tree represents a species. Some of  
135 them, specifically the leaves, are living species, while  
136 the others may be extinct or only theoretical. The  
137 common parent of two nodes is their most probable  
138 evolutionary ancestor.

139 An example of a phylogenetic tree can be seen in  
140 Figure 3. This is a tree where the lengths of individual  
141 lines between nodes represent the estimated time of  
142 evolution. If the line is short, the nodes it connects  
143 are very similar, since the time for evolution is short  
144 which implies fewer changes in the genome compared  
145 to the long lines.

146 From the phylogenetic tree, we can estimate the  
147 evolutionary relationship between different species.  
148 Then it is possible to infer a correct combination of  
149 known functional profiles for all species for which  
150 the functional profile was not found in the reference  
151 database.

152 The inference of unknown functional profiles can  
153 be done by finding the nearest nodes with known pro-  
154 files. We can search for a certain number of known  
155 profiles, or limit the search by sequence similarity to  
156 the investigated. After we have a set of nodes with  
157 known profiles, we compute a consensus profile based  
158 on the distance to the investigated node — it can be  
159 a simple average, or closer nodes may have a bigger

weight than the more distant ones. 160

#### 161 **Distance based algorithms**

162 The basic idea used in these algorithms is my origi-  
163 nal work that I introduced in my masters thesis. It  
164 is based on analyzing the representative sequences  
165 of given OTUs and comparing them to reference se-  
166 quences with known functional profiles. The resulting  
167 functional profile is then inferred from the most similar  
168 reference OTUs.

169 To speed up the search, the similarity between  
170 sequences is usually precomputed and stored in a dis-  
171 tance matrix. The rows and columns of the distance  
172 matrix represent the OTUs and the numbers in the  
173 matrix represent the distance of OTUs in the corre-  
174 sponding row and column.

175 To compute the similarity between sequences, dif-  
176 ferent methods can be used. One of them is to simply  
177 count the number of equal characters in their sequence  
178 alignments. Others punish the differences according to  
179 their evolutionary probability — because of the differ-  
180 ent chemical nature of the nucleotides in RNA, certain  
181 changes in the sequences are more probable than the  
182 others. There are various matrices that express the  
183 probability of interchange between the nucleotides [5].

### 184 2.4 Existing tools

185 The most used tools for functional analysis are Picrust  
186 [6] and Tax4Fun [7]. They both implement the phylo-  
187 genetic tree approach. The main difference between  
188 them is their reference database. Picrust uses Green-  
189 genes [8], which is outdated, while Tax4Fun uses Silva  
190 [9], which is a newer database that is still frequently  
191 updated. To eliminate this disadvantage, the creators  
192 of Picrust developed Picrust 2, which is not dependent  
193 on reference database [10]. Unfortunately, Picrust 2 is  
194 still in beta version.

195 Picrust is a bioinformatic software package imple-  
196 mented in Python and R, while Tax4Fun is an open-  
197 source package for R. Tax4Fun is newer and tries to  
198 alter the prediction method of Picrust to make it more  
199 accurate. It is also easier to use and faster.

200 The work-flow of Picrust can be divided into two  
201 parts, Gene content inference, and Metagenome infer-  
202 ence. In Gene content inference, Picrust takes the  
203 reference OTU table from Greengenes database and  
204 gene content table from IMG, which is a table contain-  
205 ing functional profiles for known genomes. It creates a  
206 tree featuring all OTUs from the reference database us-  
207 ing ancestral state reconstruction algorithm. For OTUs  
208 with an unknown functional profile, an estimated pro-  
209 file is computed, using the position of the given OTU  
210 in the phylogenetic tree and the closest OTUs with a

211 known functional profile. This step is independent on  
212 the sample, so it is computed only once. [11]

213 In Metagenome inference, Picrust takes an user-  
214 provided table of OTUs, and using the gene content  
215 table from the previous step, predicts metagenomic  
216 content of the given sample. The prediction is done  
217 by summing up the functional profiles (obtained in  
218 the previous step) corresponding to OTUs in the input  
219 table while taking into account their abundance. [11]

### 220 3. Created tool

221 In my masters thesis, I have created a new tool for  
222 functional profile prediction. It is not dependent on a  
223 reference database, as Picrust is. It implements various  
224 methods for dealing with OTUs with the unknown  
225 functional profile.

226 The dataflow of the designed tool can be seen in  
227 Figure 4. The yellow modules (Input parser, Known  
228 profile resolver, Output generator) will be the same for  
229 every sample and method for dealing with unknown  
230 OTUs, the pink one (Unknown profile resolver) differs  
231 and its accuracy is the target of the experiments.

232 The Input parser loads the data from the input  
233 sample. Then the Known profile resolver determines,  
234 which OTUs have known functional profiles, and which  
235 do not. The unknown ones are then processed by  
236 the Unknown profile resolver, which tries to estimate  
237 the most probable KO profile composition. Both the  
238 known and the estimated profiles are then merged to-  
239 gether in the Output generator.

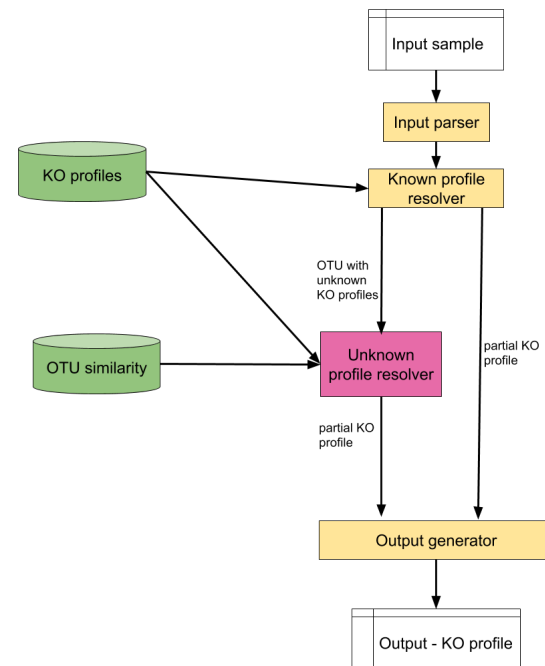
240 The Unknown profile resolver currently imple-  
241 ments three methods for functional profile prediction,  
242 two are distance based and one uses a phylogenetic  
243 tree. They will be described in detail in the following  
244 sections. The results of the experiments will also be  
245 given.

246 The green modules (KO profiles, OTU similarity)  
247 represent data sources. They are precomputed and  
248 saved in files, but the tool contains code for the pre-  
249 computation so it is possible to repeat it for other re-  
250 ference databases. KO profiles data source is a table of  
251 known species with corresponding KO profiles. OTU  
252 similarity is a data source of similarity between OTUs  
253 with known and with unknown functional profiles. It  
254 can either be a similarity matrix or a phylogenetic tree,  
255 or anything else, that somehow represents similarity  
256 between OTUs.

### 257 4. Experiments

258 To simulate the real-life situation, where we do not  
259 have information about many functional profiles, I  
260 have created a version of the reference table with a

**Figure 4.** Design of the created tool for functional profile prediction based on 16S rRNA data



261 fraction of rows missing. For 0% missing, the predic-  
262 tion should be 100% accurate, since we have all the  
263 data and no estimation is needed.

264 For each method, the accuracy was tested on 10  
265 artificial samples. To simulate missing functional pro-  
266 files, a part of reference KO profile table was randomly  
267 deleted. The ratio of the deleted table was incremen-  
268 tally increased, from 0% to 90%, to see how much the  
269 accuracy drops with more profiles missing. Since the  
270 deletion from the reference table was randomized, this  
271 step was performed 10 times. To summarize, for each  
272 ratio of missing functional profiles, I performed 100  
273 tests.

274 In the visualization of the results of the experi-  
275 ments I always show the correlation between the ex-  
276 pected and the computed functional profile in a box-  
277 plot. The y-axis will be the correlation and the x-axis  
278 the ratio of the known functional profiles. This way  
279 we can see how much the accuracy drops when we do  
280 not have enough reference data.

281 The correlation coefficients are computed as Pear-  
282 son Product-Moment Correlation, which shows the  
283 linear association between two vectors. The values  
284 of the coefficients can range from -1 to 1, where 0  
285 means no association between the vectors, values big-  
286 ger than 0 show positive association and values smaller  
287 than 0 show negative association. The proximity to  
288 -1 and 1 show the strength of the association. Gen-

289 erally, values bigger than 0.5 are considered a strong  
290 association [12].

#### 291 4.1 Distance based methods

292 I have experimented with two methods for functional  
293 prediction based on distance. I have used the Green-  
294 genes database, for which a global alignment of all  
295 sequences is available. Greengenes database was cho-  
296 sen so I can compare my results with Picrust in the  
297 future. However, the tool is not dependent on the  
298 database, and if another multiple sequence alignment  
299 data are used, the prediction will work.

300 Both methods look for the most similar OTUs with  
301 a known functional profile. The similarity is measured  
302 by the distance of the multiple sequence alignment.  
303 The number of similar sequences that are taken into  
304 account is an attribute I tried to experiment with. The  
305 best results were achieved when the estimated profile  
306 was computed as the average of profiles of 4 most  
307 similar sequences.

308 The difference in the implemented methods is the  
309 computation of the distance metric. The first one  
310 simply counts the number of similar symbols in the  
311 multiple sequence alignment. The second one uses a  
312 transition/transversion scoring matrix for nucleotides.  
313 It takes into account the chemical attributes of RNA  
314 bases and the probability of exchange of the nucleotide  
315 pairs.

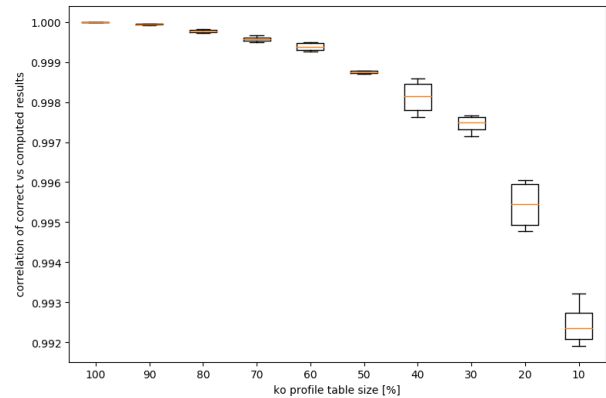
316 Surprisingly, better results were achieved with the  
317 simpler method, that does not respect the chemical  
318 properties of nucleotides. This might be caused by a  
319 wrong approach to gaps in the scoring matrix, or by  
320 the fact that the distance matrix was computed from  
321 the global alignment of all Greengenes database se-  
322 quences, so the impact of alignment errors might be  
323 significant. Therefore, with respect to the limited scope  
324 of this paper, this is the only method which will be  
325 discussed further. We can see the results of the average-  
326 finding method in Figure 5.

#### 327 4.2 Phylogenetic tree based methods

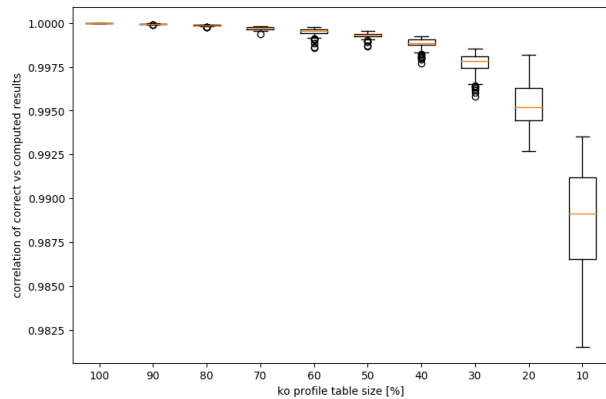
328 I have also implemented a method based on a phyloge-  
329 netic tree. It uses a UPGMA (short for "unweighted  
330 pair group method with arithmetic mean", an algorithm  
331 for creating a phylogenetic tree based on sequence sim-  
332 ilarity) method to compute a phylogenetic tree of all  
333 the OTUs from the Greengenes database, both the ones  
334 with known and unknown KO profiles. Then the un-  
335 known OTUs are determined as the average profile of  
336 OTUs which are connected to it in the tree.

337 The tree is computed before the functional predic-  
338 tion. The script for the tree creation is a part of the

**Figure 5.** Evaluation of a distance based method for functional prediction on complete functional profiles. The x-axis represents the ratio of known profiles in the sample, the y-axis represents the correlation between expected and computed result.



**Figure 6.** Evaluation of a phylogenetic tree based method for functional prediction on complete functional profiles. The x-axis represents the ratio of known profiles in the sample, the y-axis represents the correlation between expected and computed result.



339 tool, so the process can be repeated for any reference  
340 data.

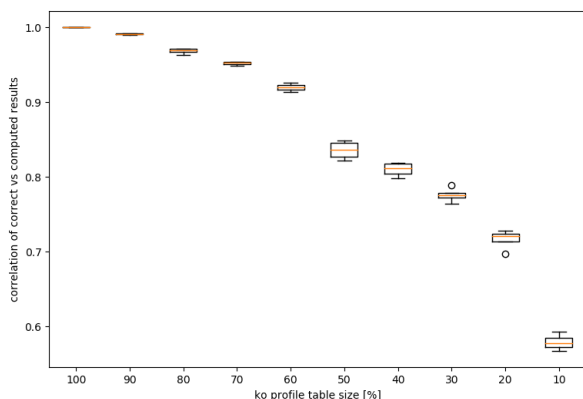
341 In UPGMA, each item is paired with the closest  
342 item by a given distance matrix. The pair is connected  
343 and is assigned a parent node in the resulting tree. In  
344 the name of the node, I store names of all its children.  
345 That means that the root has names of all the OTUs in  
346 the reference database. The closest OTUs to any OTU  
347 are the ones with which it was connected.

348 When looking for the most similar OTU, I search  
349 the tree and find, in which point the searched OTU  
350 was connected with some other node. The name of the  
351 node with which it was connected are the ones that are  
352 the most similar, and the estimated profile is inferred  
353 from them.

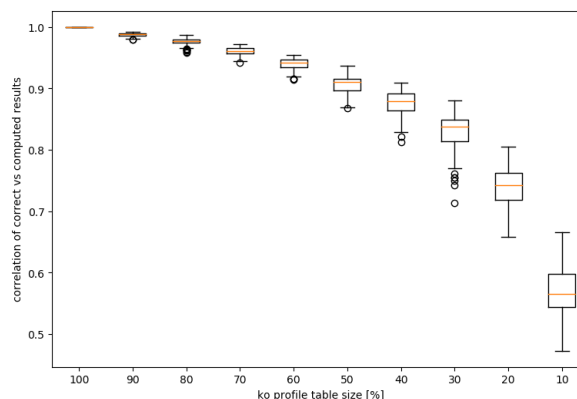
#### 354 4.3 Comparison

355 As we can see, the method based on average is slightly  
356 better than the one based phylogenetic tree. However,

**Figure 7.** Evaluation of a distance based method for functional prediction on the rarest functions. The x-axis represents the ratio of known profiles in the sample, the y-axis represents the correlation between expected and computed result.



**Figure 8.** Evaluation of a phylogenetic tree based method for functional prediction on the rarest functions. The x-axis represents the ratio of known profiles in the sample, the y-axis represents the correlation between expected and computed result.



357 the correlation between expected and computed pro-  
 358 files is surprisingly high, over 99% even when only  
 359 10% of the profiles are known. Picrust and Tax4Fun  
 360 also report a high prediction accuracy [11, 13].

361 The high correlation can be caused by a number  
 362 of reasons, the most probable is **common metabolic**  
 363 **functions** — each bacteria must have basic functions  
 364 for translation, transcription, and processing of com-  
 365 mon metabolites. This is the part of the functional  
 366 profile, that is the same for every species of bacteria,  
 367 independent on the sample. A number of KO specific  
 368 for a certain species of bacteria is much smaller, so it  
 369 might not have that much of an effect for the correla-  
 370 tion.

371 Since every bacteria has common metabolic func-  
 372 tions, it is expected that they will be present in every  
 373 sample. More interesting are the functions that are  
 374 exclusive for a certain species. From this reason, I  
 375 have altered the correlation so that it takes into account  
 376 only those KOs, that are present in less than 1% of the  
 377 species in the reference table.

378 With this approach, the results between the meth-  
 379 ods are more different, so it is more useful for compar-  
 380 ing different algorithms. We can see the new results in  
 381 Figures 7 and 8. Now, we can see that the correlations  
 382 drop more significantly, and with only 10% of the pro-  
 383 files known, we are at the 60% match. The difference  
 384 between the methods is still minimal, which might  
 385 indicate, that while 16S rRNA is immensely useful as  
 386 a species identifier, the connection between the 16S  
 387 sequence and functional profile is not very significant  
 388 and to predict functional profile more accurately, we  
 389 need to look at the whole genome. To confirm this  
 390 hypothesis, more analysis and evaluation have to be  
 391 performed.

## 5. Conclusion 392

In this paper, I have introduced functional profile anal- 393  
 394 ysis, which is an important part of metagenomic re- 395  
 396 search. I have discussed the most used methods and 397  
 398 created a tool that implements them. 399

The focus of this paper is the comparison of differ- 397  
 398 ent methods for functional analysis. I have shown, 399  
 400 that the classic approach to experimental evaluation — 401  
 402 when we look at the whole functional profile — gives 403  
 404 a different result than the evaluation that looks only at 405  
 406 the more specific functions. 407

The significance of the evaluation of the specific 403  
 404 functions is that it gives us more detailed information 405  
 406 about a sample. The common metabolic functions are 407  
 408 a part of every bacteria species, so it is not a surprise 409  
 410 to find them in every sample. More rare functions are 411  
 412 more informative and harder to predict. 413

In the future, I would like to compare my tool with 409  
 410 Picrust and Tax4Fun. It would be interesting to see 411  
 412 how they stand in the tests of only specific KOs since 413  
 414 they are widely used in bioinformatics research. I will 415  
 416 also add some more methods for functional analysis, 417  
 418 including one based on linear regression, to try a more 419  
 420 computer-sciences inspired approach. 421

## Acknowledgements 416

I would like to thank my supervisor Ing. Stanislav 417  
 418 Smatana for his help. 419

## References 419

[1] Xochitl C. Morgan and Curtis Huttenhower. 420  
 421 Chapter 12: Human microbiome analysis. In 422  
 423 *PLoS Computational Biology*, 2012. 424

- 423 [2] Andreas Hiergeist, Udo Reischl, and André Gess-  
424 ner. Multicenter quality assessment of 16s ribo-  
425 somal dna-sequencing for microbiome analyses  
426 reveals high inter-center variability. *International*  
427 *journal of medical microbiology : IJMM*, 306  
428 5:334–342, 2016.
- 429 [3] Jay-Hyun Jo, Elizabeth A. Kennedy, and Heidi  
430 Kong. Research techniques made simple: Bacte-  
431 rial 16s ribosomal rna gene sequencing in cuta-  
432 neous research. *Journal of Investigative Derma-*  
433 *tology*, 136:e23–e27, 03 2016.
- 434 [4] KEGG Orthology database. [Online:  
435 <https://www.genome.jp/kegg/ko.html>]. [Visited  
436 4.12.2018].
- 437 [5] William R. Pearson. An introduction to sequence  
438 similarity (“homology”) searching. *Curr Protoc*  
439 *Bioinformatics*, page Chapter 3:Unit3.1, 2013.
- 440 [6] Picrust webpage. [Online:  
441 <http://picrust.github.io/picrust/>]. [Visited  
442 4.12.2018].
- 443 [7] Tax4Fun webpage. [Online:  
444 <http://tax4fun.gobics.de/>]. [Visited 4.12.2018].
- 445 [8] Greengenes database webpage. [Online:  
446 <http://greengenes.secondgenome.com/>]. [Visited  
447 28.12.2018].
- 448 [9] Silva database webpage. [Online:  
449 <https://www.arb-silva.de/>]. [Visited 28.12.2018].
- 450 [10] Picrust 2 github repository. [Online:  
451 <https://github.com/picrust/picrust2>]. [Vis-  
452 ited 4.12.2018].
- 453 [11] Morgan G I Langille, Jesse Zaneveld, J Gre-  
454 gory Caporaso, Daniel Mcdonald, Dan Knights,  
455 Joshua A Reyes, Jose C Clemente, Deron  
456 Burkepille, Rebecca Vega Thurber, Rob Knight,  
457 Robert G Beiko, and Curtis Huttenhower. Predic-  
458 tive functional profiling of microbial communi-  
459 ties using 16s rna marker gene sequences. *Na-*  
460 *nure biotechnology*, 31, 08 2013.
- 461 [12] Pearson Product-Moment Correlation. [On-  
462 line: [https://statistics.laerd.com/statistical-](https://statistics.laerd.com/statistical-guides/pearson-correlation-coefficient-statistical-guide.php)  
463 [guides/pearson-correlation-coefficient-](https://statistics.laerd.com/statistical-guides/pearson-correlation-coefficient-statistical-guide.php)  
464 [statistical-guide.php](https://statistics.laerd.com/statistical-guides/pearson-correlation-coefficient-statistical-guide.php)]. [Visited 1.4.2019].
- 465 [13] Kathrin P ABhauer, Bernd Wemheuer, Rolf  
466 Daniel, and Peter Meinicke. Tax4fun: Predicting  
467 functional profiles from metagenomic 16s rna  
468 data. *Bioinformatics (Oxford, England)*, 31, 05  
469 2015.