

GigaScience

Microbiome Learning Repo (ML Repo): A public repository of microbiome regression and classification tasks

--Manuscript Draft--

Manuscript Number:	GIGA-D-18-00316	
Full Title:	Microbiome Learning Repo (ML Repo): A public repository of microbiome regression and classification tasks	
Article Type:	Technical Note	
Funding Information:	National Institute of Allergy and Infectious Diseases (R01AI121383)	Not applicable
Abstract:	<p>The use of machine learning in high-dimensional biological applications, such as the human microbiome, has grown exponentially in recent years. Unfortunately, challenges still exists for machine learning algorithm developers who often lack domain expertise required for interpretation and curation of the heterogeneous microbiome datasets. We present Microbiome Learning Repo (ML Repo), a public, web-based repository of 33 curated classification and regression tasks from 15 published human microbiome datasets. We highlight the use of ML Repo in several use cases to demonstrate its wide application, and expect it to be an important resource for algorithm developers.</p>	
Corresponding Author:	Dan Knights UNITED STATES	
Corresponding Author Secondary Information:		
Corresponding Author's Institution:		
Corresponding Author's Secondary Institution:		
First Author:	Pajau Vangay, Ph.D.	
First Author Secondary Information:		
Order of Authors:	Pajau Vangay, Ph.D. Dan Knights, Ph.D. Benjamin M. Hillmann	
Order of Authors Secondary Information:		
Additional Information:		
Question	Response	
Are you submitting this manuscript to a special series or article collection?	No	
Experimental design and statistics	Yes	
Full details of the experimental design and statistical methods used should be given in the Methods section, as detailed in our Minimum Standards Reporting Checklist . Information essential to interpreting the data presented should be made available in the figure legends.		

<p>Have you included all the information requested in your manuscript?</p>	
<p>Resources</p> <p>A description of all resources used, including antibodies, cell lines, animals and software tools, with enough information to allow them to be uniquely identified, should be included in the Methods section. Authors are strongly encouraged to cite Research Resource Identifiers (RRIDs) for antibodies, model organisms and tools, where possible.</p> <p>Have you included the information requested as detailed in our Minimum Standards Reporting Checklist?</p>	<p>Yes</p>
<p>Availability of data and materials</p> <p>All datasets and code on which the conclusions of the paper rely must be either included in your submission or deposited in publicly available repositories (where available and ethically appropriate), referencing such data using a unique identifier in the references and in the “Availability of Data and Materials” section of your manuscript.</p> <p>Have you have met the above requirement as detailed in our Minimum Standards Reporting Checklist?</p>	<p>Yes</p>

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65

1 **Microbiome Learning Repo (ML Repo): A public repository of microbiome regression**
2 **and classification tasks**

3
4 Pajau Vangay¹, Benjamin M. Hillmann², Dan Knights^{12*}

5
6 ¹Bioinformatics and Computational Biology, University of Minnesota, Minneapolis, MN, 55455,
7 vanga015@umn.edu

8 ²Department of Computer Science and Engineering, University of Minnesota, Minneapolis, MN,
9 55455, hillm096@umn.edu

10
11 *Corresponding author: dknights@umn.edu

12
13
14 **Abstract**

15 The use of machine learning in high-dimensional biological applications, such as the human
16 microbiome, has grown exponentially in recent years. Unfortunately, challenges still exists for
17 machine learning algorithm developers who often lack domain expertise required for
18 interpretation and curation of the heterogeneous microbiome datasets. We present Microbiome
19 Learning Repo (ML Repo), a public, web-based repository of 33 curated classification and
20 regression tasks from 15 published human microbiome datasets. We highlight the use of ML
21 Repo in several use cases to demonstrate its wide application, and expect it to be an important
22 resource for algorithm developers.

23
24 **Keywords**

25 Microbiome, machine learning, repository, database

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65

27 **Findings**

28

29 **Background**

30 Machine learning is widely used as a method for classification and prediction, with a growing
31 number of applications in human health [1]. The use of machine learning in biological fields
32 [2,3], and more specifically the microbiome research field [4–7], has grown exponentially due to
33 the robustness of these algorithms to high dimensional data. However, challenges exist for
34 large-scale meta-analysis as they often require manual curation of metadata and standardized
35 processing of raw sequence data, resulting in variation in the results derived from chosen
36 datasets across studies [8,9]. In addition, microbiome research data can be challenging to
37 access and analyze for expert machine learning algorithm developers, who often do not have
38 the domain expertise required to parse the data and metadata in complex microbiome studies.
39 There exist general resources with curated classification tasks from variety of domains. The
40 University of California Irvine (UCI) Machine Learning Repository [10] revolutionized machine
41 learning methods development by giving developers access to many curated datasets; its
42 widespread usage and impact can be seen from its thousands of resulting citations. Currently,
43 we are unaware of any machine learning repository specifically for microbiome classification
44 tasks. We constructed a complementary database to address this deficiency, in order to
45 promote the development of and usage of improved machine learning methods for the
46 microbiome community.

47

48 **Workflow**

49 We present the Microbiome Learning Repo (ML Repo), a repository of 33 curated classification
50 and regression tasks using human microbiome data. Our 33 tasks are curated from 15 publicly
51 available human microbiome datasets, which include 12 amplicon-based and 3 shotgun
52 sequencing datasets [Table 1]. These datasets vary across sequencing technology platforms,

1
2
3
4 53 16s hypervariable regions, and study design, in order to help developer ensure robustness of
5
6 54 algorithms across data types. We streamlined the microbiome data using a single post-
7
8 55 processing workflow [Fig 1A]. We downloaded trimmed and quality filtered sequencing reads for
9
10 56 n=8 datasets from QIITA [11], and raw sequences for n=7 datasets from public repositories. We
11
12 57 preprocessed raw sequences using SHI7 [12] or QIIME [13] according to individual technologies
13
14 58 and characteristics of each study. Full details regarding the data preprocessing are provided for
15
16 59 each data set in the repository. We picked Operational Taxonomic Units (OTUs) from all quality
17
18 60 filtered sequences using a closed-reference method with the BURST [14] aligner against both
19
20 61 the NCBI RefSeq 16S ribosomal RNA project [15] and the Greengenes 97 database [16].
21
22 62 Samples with depths lower than 1000 sequences per sample were dropped for n=10 datasets,
23
24 63 while we applied a lower threshold of 100 sequences per sample for n=5 datasets which had
25
26 64 lower expected bacterial load. As a result, for each dataset we generated RefSeq-based OTU
27
28 65 and taxa abundance counts, and Greengenes-based OTU and taxa abundance counts. We
29
30 66 excluded additional post-processing filtering and normalization steps so that these parameters
31
32 67 can be included in future benchmarking use cases as needed. We also limit our data to OTU
33
34 68 and taxa tables as other metrics such as alpha and beta diversity can be subsequently
35
36 69 generated as needed.
37
38
39
40
41
42
43

44 71 Sample metadata from individual studies were manually curated to generate viable prediction
45
46 72 tasks. When available, published study exclusion criteria was applied accordingly and
47
48 73 confounders were removed by dropping samples or stratification. Studies that were cross-
49
50 74 sectional by design but contained several samples per subject were filtered to contain one
51
52 75 sample per subject. Well-known confounders, such as geography, were accounted for when
53
54 76 constructing prediction tasks for other human-associated conditions. Longitudinal studies were
55
56 77 reduced to single time points of interest to minimize the effect of high intra-individual similarities.
57
58 78 Hence, each prediction task is made available as an individual, compartmentalized metadata file
59
60
61
62
63
64
65

1
2
3
4 79 that contains sample identifiers, responses to predict, and optionally, confounder variables to
5
6 80 control for. As a result, we generated 33 distinct tasks for predicting human-associated
7
8 81 responses.
9

10 82

13 83 **Publicly available web-based interface**

15 84 We expect two types of users: (1) machine-learning algorithm developers with limited
16
17 85 knowledge of microbiome study designs and (2) microbiome researchers interested in obtaining
18
19 86 additional datasets for meta-analysis. Generally, we expect that methods developers will be
20
21 87 most interested in sweeping through the full set of prediction tasks for benchmarking, and hence
22
23 88 would prefer to download a single compressed file containing all tasks and data. On the other
24
25 89 hand, we expect that microbiome researchers will be more selective in downloading specific
26
27 90 datasets and tasks depending on their research domain. Hence, researchers may prefer to
28
29 91 browse specific details about tasks and datasets prior to downloading.
30
31

32 92

35 93 Based on these expected use cases, we created a publicly available web-interface for MLRepo
36
37 94 hosted by GitHub Pages and available at: <https://knights-lab.github.io/MLRepo>. Tasks are
38
39 95 organized by relevant response categories [Fig 2A]. Task pages contain descriptive details such
40
41 96 as Sample Size and Response Type that are specific to the selected prediction task, as well as
42
43 97 links for downloading OTU tables, taxa tables, and sample metadata [Fig 2B]. Dataset pages
44
45 98 contain important details about the entire dataset, including links to the original research study,
46
47 99 as well as original metadata files and quality filtered sequences [Fig 2C]. We also provide a
48
49 100 single compressed file containing the entire set of available tasks (OTU tables, taxa tables, and
50
51 101 relevant metadata) for download from the main home page.
52
53

54 102

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65

103 **Benefits of curated microbiome-based prediction tasks**

104 We expect MLRepo to be beneficial for both the machine-learning community as well as the
105 microbiome research community. MLRepo will be a powerful complement to UCI's machine
106 learning repository, as it will allow for benchmarking curated classification tasks with high-
107 dimensional data, and hence enable the subsequent development of novel algorithms for these
108 complex datasets. Our streamlined approach in generating OTU and taxa tables offers a rich set
109 of 15 datasets that microbiome researchers can use directly for further comparison with their
110 own studies, for teaching and learning purposes, or for large meta-analyses. We expect that our
111 provided OTU and taxa tables will also be beneficial for researchers with limited access to high-
112 performance computing resources or bioinformatics skills necessary for processing raw
113 sequencing data. In addition, we expect microbiome-specific methods development will also
114 benefit from our repository prediction tasks. The subsetted samples found in each prediction
115 task metadata file replaces the work of rigorously deciphering metadata and nuances from
116 individual research studies. Hence, new methods, such as OTU-picking algorithms, can be
117 evaluated not only on metrics such as speed and accuracy, but also based on overall impact to
118 study findings.

120 **Comparison to similar databases**

121 Although a number of microbiome repositories exist, many are intended as data archival
122 repositories [17,18] or function as resources for aggregating across studies [19]. Resources
123 such as QIITA [11] offer an extensive collection of datasets, and mock-community-based
124 Mockrobiota [20] is well-suited for benchmarking upstream methods, but neither offer support
125 for the metadata interpretation necessary for predicting high-level phenotypes. MLRepo differs
126 from all of these resources in that we provide well-defined tasks for predicting responses from
127 manually curated metadata and standardized data from published microbiome research studies.

1
2
3
4 **129 Case studies**

5
6 130 We compare the performance of three machine learning models: a random forest [21], and a
7
8 131 support vector machine [22] (SVM) with either a radial or linear kernel. Sweeping through
9
10 132 available tasks with binary responses, we compare our models by examining receiver operating
11
12 133 curves (ROCs) and areas under the curve (AUC) [Fig 3]. Through comparison of ROCs, we can
13
14 134 see that random forest outperforms or ties the other two models in 21 out of the 28 tasks. The
15
16 135 choice of kernels for SVM appears to have limited impact on overall mean accuracy, yet a linear
17
18 136 kernel can perfectly classify penicillin-treated and vancomycin-treated mouse cecal contents
19
20 137 when the other models could not; further examination of the microbial features in these samples
21
22 138 may be warranted to better understand the strengths of this kernel. We also performed pairwise
23
24 139 comparisons of random forest against the other models across all tasks. When evaluated by
25
26 140 AUC, considered the standard method for machine learning model evaluation [23,24], random
27
28 141 forest performs significantly better than both SVM with a linear kernel ($P=0.0014$) and with a
29
30 142 radial kernel ($P=0.00032$) [Fig 4A]. We found that random forest accuracy improvements were
31
32 143 moderate when compared with SVM-Linear ($P=0.083$) and SVM-Radial ($P=0.03$) [Fig 4B]. Our
33
34 144 results support the broad usage [4,5,8,25] and acceptance of random forest as a robust
35
36 145 classifier [6] with high-dimensional microbiome data.
37
38
39
40
41
42
43

44 146
45
46 147 To assess the impact of reference database choice on classification accuracies, we also used
47
48 148 the classification tasks to compare random forest using OTUs picked with the Greengenes 97
49
50 149 database or the NCBI RefSeq Targeted Loci Project 16s project. We find that there is limited
51
52 150 impact of database choice to overall classification accuracies [Fig 4C, Fig 5]. This may be due
53
54 151 to (1) large effect sizes that are driven mainly by several well-characterized bacterial taxa
55
56 152 present in both databases (e.g. stool versus tongue samples), or (2) small effect sizes such that
57
58 153 classification is difficult regardless of the database (e.g. male versus female stool).
59
60
61
62
63
64
65

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65

155 Future work

156 We expect and hope that the broader microbiome research community will add new datasets
157 and prediction tasks to MLRepo. We provide instructions on our GitHub repository to guide
158 users to create a fork from our repository, add the appropriate data and files, and update the
159 master task and dataset lists. Researchers can then submit a pull request for our review, and if
160 properly formatted, will be accepted and merged into the repository. We expect that data
161 submissions will come from either the original researchers or those well-acquainted with the
162 datasets, and hence will expect that sample selection and subsetting will have undergone
163 rigorous review for prediction tasks.

164

165 Methods

166 Pre-processing of sequencing reads

167 When available, preprocessed FASTA files were downloaded from QIITA (or previously, the
168 QIIME database). For all other datasets, raw FASTQ files were downloaded from sources listed
169 in Supplemental Table 1. Sequences were trimmed and quality filtered using SHI7 [12] or QIIME
170 [13]. OTUs were picked from processed FASTA files using BURST [26] with Greengenes [16]
171 97 or the NCBI RefSeq Targeted Loci Project 16s project [15] (accessed on 17-07-04). Samples
172 with sequencing depth lower than 1000 sequences per sample were dropped for all studies,
173 except for five datasets [27–31], where the minimum threshold was 100 sequences per sample.

174

175 Selection of classification tasks

176 Classification tasks were selected based on reported study results, biologically relevant high-
177 level phenotypes, and sufficient sample sizes. Original metadata files and research methods
178 were rigorously and manually curated in order to subset samples with minimal confounders. For
179 confounders that were inherent to the study, we include an additional variable to control for in

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65

180 the task metadata files. Presence of control variables can be found by examining “control_vars”
181 in the Tasks table.

182

183 **Website generation**

184 Website templating was developed using Jinja2 [32] and custom Python scripts. Individual
185 webpages were generated by iterating through items in the Tasks and Datasets tables, and
186 dynamically populating templates in order to generate individual Markdown [33] pages. The
187 resulting Markdown pages are hosted as GitHub Pages.

188

189 **Case Study Benchmarking**

190 Case study results were generated with custom R [34] scripts, which can be found in the
191 */example* folder in the MLRepo Github repository. To compare machine learning models, we
192 iterated through tasks with binary responses. OTU counts were converted to relative
193 abundances, filtered at a minimum of 10% prevalence across samples, and collapsed at a
194 complete-linkage correlation of 95%. We then constructed a 5-fold cross-validation for tasks
195 containing more than 100 samples, or a leave-one-out cross-validation for tasks with smaller
196 sample sizes. For n-fold cross validation, samples were assigned to folds such that classes
197 were equally balanced within each fold (e.g. if our task contained 40% healthy and 60%
198 diseased samples, our folds would also be selected to represent this distribution). For tasks that
199 contained control variables, we selected folds such that samples with the same control variable
200 value were contained within the same fold. For example, for a task dataset containing matching
201 stool and oral samples from subjects, the Subject Identifier would be listed as the control
202 variable and we should assign samples to folds such that all samples from a specific subject
203 were contained within a fold. This step is crucial to avoid biasing or overfitting the training
204 model; test folds should contain not only new samples, but also samples that are independent
205 from those in the training set. Models were constructed using the ‘caret’ package [35]. This

1
2
3
4 206 process was bootstrapped 100 times, and the mean class probabilities were used to calculate
5
6 207 the resulting AUCs and ROCs. To compare classification accuracies using different reference
7
8 208 databases, we used a similar procedure but held the model constant and predicted using
9
10 209 different base OTU tables. This framework enables comparison of a myriad of machine learning
11
12 210 models available in the 'caret' package, and can be easily expanded to compare different OTU-
13
14 211 picking algorithms, or normalization and filtering techniques.
15
16
17
18
19

20 213 **Availability of supporting source code and requirements**

21
22 214
23
24 215 Project name: Microbiome Learning Repo
25
26 216 Project home page: <https://knights-lab.github.io/MLRepo/>
27
28 217 Operating system: Platform independent
29
30
31 218 Programming language: Python, R
32
33 219 License: MIT License
34
35
36
37

38 221 **Declarations**

39
40 222
41
42 223 **Authors' contributions**
43
44 224 Conceptualization: P.V. and D.K; Data curation: P.V.; Formal analyses: P.V.; Methodology:
45
46 225 P.V., B.H., D.K.; Software: P.V.; Writing - original draft: P.V.; Writing - review and editing: B.H.
47
48 226 and D.K.
49
50

51 227
52
53 228 **Competing Interests**

54
55 229 D.K. serves as CEO and holds equity in CoreBiome, a company involved in the
56
57 230 commercialization of microbiome analysis. The University of Minnesota also has financial
58
59 231 interests in CoreBiome under the terms of a license agreement with CoreBiome. These interests
60
61
62
63
64
65

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65

232 have been reviewed and managed by the University of Minnesota in accordance with its
233 Conflict-of-Interest policies.

234

235 **References**

236 1. Jordan MI, Mitchell TM. Machine learning: Trends, perspectives, and prospects. *Science*.
237 2015;349:255–60.

238 2. Shipp MA, Ross KN, Tamayo P, Weng AP, Kutok JL, Aguiar RCT, et al. Diffuse large B-cell
239 lymphoma outcome prediction by gene-expression profiling and supervised machine learning.
240 *Nat Med*. 2002;8:68–74.

241 3. Furey TS, Cristianini N, Duffy N, Bednarski DW, Schummer M, Haussler D. Support vector
242 machine classification and validation of cancer tissue samples using microarray expression
243 data. *Bioinformatics*. 2000;16:906–14.

244 4. Aagaard K, Riehle K, Ma J, Segata N, Mistretta T-A, Coarfa C, et al. A metagenomic
245 approach to characterization of the vaginal microbiome signature in pregnancy. *PLoS One*.
246 2012;7:e36466.

247 5. Yatsunenko T, Rey FE, Manary MJ, Trehan I, Dominguez-Bello MG, Contreras M, et al.
248 Human gut microbiome viewed across age and geography. *Nature*. 2012;486:222–7.

249 6. Knights D, Costello EK, Knight R. Supervised classification of human microbiota. *FEMS*
250 *Microbiol Rev*. 2011;35:343–59.

251 7. Smith MI, Yatsunenko T, Manary MJ, Trehan I, Mkakosya R, Cheng J, et al. Gut microbiomes
252 of Malawian twin pairs discordant for kwashiorkor. *Science*. 2013;339:548–54.

253 8. Pasolli E, Truong DT, Malik F, Waldron L, Segata N. Machine Learning Meta-analysis of

1
2
3
4 254 Large Metagenomic Datasets: Tools and Biological Insights. PLoS Comput Biol. Public Library
5
6 255 of Science; 2016;12:e1004977.
7
8
9
10 256 9. Sze MA, Schloss PD. Looking for a Signal in the Noise: Revisiting Obesity and the
11
12 257 Microbiome. MBio [Internet]. 2016;7. Available from: <http://dx.doi.org/10.1128/mBio.01018-16>
13
14
15 258 10. Asuncion A, Newman D. UCI machine learning repository [Internet]. 2007. Available from:
16
17 259 <https://ergodicity.net/2013/07/>
18
19
20 260 11. Qiita Development Team. QIITA [Internet]. Available from: <http://qiita.microbio.me/>
21
22
23 261 12. Al-Ghalith GA, Hillmann B, Ang K, Shields-Cutler R, Knights D. SHI7 Is a Self-Learning
24
25 262 Pipeline for Multipurpose Short-Read DNA Quality Control. mSystems [Internet]. 2018;3.
26
27 263 Available from: <http://dx.doi.org/10.1128/mSystems.00202-17>
28
29
30
31 264 13. Caporaso JG, Kuczynski J, Stombaugh J, Bittinger K, Bushman FD, Costello EK, et al.
32
33 265 QIIME allows analysis of high-throughput community sequencing data. Nat Methods.
34
35 266 2010;7:335–6.
36
37
38
39 267 14. Al-Ghalith G, Knights D. BURST enables optimal exhaustive DNA alignment for big data
40
41 268 [Internet]. Zenodo; 2017. Available from: <http://dx.doi.org/10.5281/ZENODO.806850>
42
43
44 269 15. O’Leary NA, Wright MW, Brister JR, Ciufu S, Haddad D, McVeigh R, et al. Reference
45
46 270 sequence (RefSeq) database at NCBI: current status, taxonomic expansion, and functional
47
48 271 annotation. Nucleic Acids Res. 2016;44:D733–45.
49
50
51
52 272 16. McDonald D, Price MN, Goodrich J, Nawrocki EP, DeSantis TZ, Probst A, et al. An
53
54 273 improved Greengenes taxonomy with explicit ranks for ecological and evolutionary analyses of
55
56 274 bacteria and archaea. ISME J. 2012;6:610–8.
57
58
59
60 275 17. Hunter S, Corbett M, Denise H, Fraser M, Gonzalez-Beltran A, Hunter C, et al. EBI
61
62
63
64
65

1
2
3
4 276 metagenomics—a new resource for the analysis and archiving of metagenomic data. *Nucleic*
5
6 277 *Acids Res.* Oxford University Press; 2014;42:D600–6.
7
8
9
10 278 18. Leinonen R, Sugawara H, Shumway M, International Nucleotide Sequence Database
11
12 279 Collaboration. The sequence read archive. *Nucleic Acids Res.* 2011;39:D19–21.
13
14
15 280 19. Forster SC, Browne HP, Kumar N, Hunt M, Denise H, Mitchell A, et al. HPMCD: the
16
17 281 database of human microbial communities from metagenomic datasets and microbial reference
18
19 282 genomes. *Nucleic Acids Res.* 2016;44:D604–9.
20
21
22
23 283 20. Bokulich NA, Rideout JR, Mercurio WG, Shiffer A, Wolfe B, Maurice CF, et al. mockrobiota:
24
25 284 a Public Resource for Microbiome Bioinformatics Benchmarking. *mSystems* [Internet]. 2016;1.
26
27 285 Available from: <http://dx.doi.org/10.1128/mSystems.00062-16>
28
29
30 286 21. Breiman L. Random Forests. *Mach Learn.* 2001;45:5–32.
31
32
33
34 287 22. Cortes C, Vapnik V. Support-vector networks. *Mach Learn.* 1995;20:273–97.
35
36
37 288 23. Huang J, Ling CX. Using AUC and accuracy in evaluating learning algorithms. *IEEE Trans*
38
39 289 *Knowl Data Eng.* 2005;17:299–310.
40
41
42 290 24. Ling CX, Huang J, Zhang H, Others. AUC: a statistically consistent and more discriminating
43
44 291 measure than accuracy. *IJCAI.* 2003. p. 519–24.
45
46
47
48 292 25. Karlsson FH, Tremaroli V, Nookaew I, Bergström G, Behre CJ, Fagerberg B, et al. Gut
49
50 293 metagenome in European women with normal, impaired and diabetic glucose control. *Nature.*
51
52 294 2013;498:99–103.
53
54
55 295 26. Al-Ghalith G, Knights D. BURST enables optimal exhaustive DNA alignment for big data
56
57 296 [Internet]. Zenodo; 2017. Available from: <http://dx.doi.org/10.5281/ZENODO.806850>
58
59
60
61
62
63
64
65

1
2
3
4 297 27. Gevers D, Kugathasan S, Denson LA, Vázquez-Baeza Y, Van Treuren W, Ren B, et al. The
5
6 298 treatment-naive microbiome in new-onset Crohn's disease. *Cell Host Microbe*. 2014;15:382–92.
7
8
9
10 299 28. Human Microbiome Project Consortium. A framework for human microbiome research.
11
12 300 *Nature*. 2012;486:215–21.
13
14
15 301 29. Kostic AD, Gevers D, Pedamallu CS, Michaud M, Duke F, Earl AM, et al. Genomic analysis
16
17 302 identifies association of *Fusobacterium* with colorectal carcinoma. *Genome Res*. 2012;22:292–
18
19 303 8.
20
21
22
23 304 30. David LA, Maurice CF, Carmody RN, Gootenberg DB, Button JE, Wolfe BE, et al. Diet
24
25 305 rapidly and reproducibly alters the human gut microbiome. *Nature*. 2014;505:559–63.
26
27
28 306 31. Turnbaugh PJ, Hamady M, Yatsunencko T, Cantarel BL, Duncan A, Ley RE, et al. A core gut
29
30 307 microbiome in obese and lean twins. *Nature*. 2009;457:480–4.
31
32
33
34 308 32. Ronacher A. Jinja2 [Internet]. 2017. Available from: <http://jinja.pocoo.org/>
35
36
37 309 33. Gruber J, Swartz A, Others. Markdown. 2004.
38
39
40 310 34. Team RC, Others. R: A language and environment for statistical computing. Citeseer; 2013;
41
42 311 Available from:
43
44 312 <http://citeseerx.ist.psu.edu/viewdoc/download?doi=10.1.1.470.5851&rep=rep1&type=pdf>
45
46
47
48 313 35. Kuhn M, Others. Caret package. *J Stat Softw*. 2008;28:1–26.
49
50
51 314
52
53 315
54
55 316
56
57
58 317
59
60 318
61
62
63
64
65

319 **Tables**320 **Table 1. Microbiome datasets with available classification tasks in ML Repo.**

Project Name	V Region	Target size	Num samples	Num subjects	Area	Description	Sequencing Technology	Study Design
Cho 2012	V3	177	95	47	Antibiotics	Mouse fecal and cecal samples, Control vs. 4 kinds of antibiotics	454	Cross-Sectional
Claesson 2012	V4	221	168	168	Age	Elderly and young adults	454	Cross-Sectional
David 2014	V4	282	235	11	Diet	Plant-based vs. Animal-based diet, Cross-over study	Illumina MiSeq	Longitudinal
Gevers 2014	V4	173	1321	668	IBD	Biopsies from IBD patients prior to treatment	Illumina MiSeq	Cross-Sectional
HMP 2012	V35	527	6407	242	Body Habitat, Gender	Up to 18 body sites across 242 healthy subjects at 1-2 time points	454	Cross-Sectional
Kostic 2012	V35	569	190	95	Colorectal Cancer	Adjacent Healthy vs. Tumor Colon Biopsy Tissues	454	Paired
Montassier 2016	V56	280	28	28	Bacteremia	Patients prior to chemotherapy who did or did not develop bacteremia	454	Cross-Sectional
Morgan 2012	V35	569	231	231	IBD	Healthy, Crohn's Disease, or Ulcerative Colitis patients	454	Cross-Sectional
Turnbaugh 2009	V2	230	281	154	Obesity	Monozygotic or dizygotic twin pairs concordant for BMI class, and their mothers	454	Cross-Sectional
Wu 2011	V12	244	95	10	Diet	Controlled HighFat or LowFat feeding on 10 subjects over 10 days	454	Longitudinal
Yatsunencko 2012	V4	282	531	531	Geography, Age, Gender	Humans of varying ages from the USA, Malawi, and Venezuela	Illumina MiSeq	Cross-Sectional
Ravel 2011	V12	240	396	396	Bacterial Vaginosis	Vaginal samples from four ethnic groups nugen scores for bacterial vaginosis	454	Cross-Sectional
Karlsson 2013	NA	NA	144	144	Diabetes	Patients with normal, impaired, or type 2 diabetes glucose tolerance categories	Illumina HiSeq	Cross-Sectional
Qin 2012	NA	NA	134	134	Diabetes	Healthy vs type 2 diabetes Chinese patients	Illumina HiSeq	Cross-Sectional
Qin 2014	NA	NA	130	130	Cirrhosis	Cirrhosis versus healthy	Illumina HiSeq	Cross-Sectional

321 ML Repo contains 33 classification and regression tasks from 15 publicly available human
322 microbiome datasets shown here.

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65

323 Figure Legends

324 Figure 1. Data processing workflow and website generation.

325 (A) Quality-filtered sequences were obtained from either the QIITA or from another public
326 repository and trimmed and filtered using SHI7. Reference-based OTUs were picked
327 using BURST with the NCBI RefSeq and Greengenes 97 databases.

328 (B) Individual GitHub Markdown pages were generated from dataset and task lists with a
329 custom Python script and Jinja2 template, then uploaded to GitHub to be hosted.

330 Figure 2. Screenshots of ML Repo web interface.

331 (A) Available classification and regression tasks are listed by high level phenotype
332 categories for browsing.

333 (B) Individual task webpages contain links to files for classifying a specific task, as well as
334 relevant task-specific metadata.

335 (C) Individual dataset webpages contain relevant metadata pertaining to the entire dataset,
336 as well as links to raw metadata files and sequencing data.

337 Figure 3. ROCs comparing random forest and SVM with different kernels.

338 Sweeping across all binary classification tasks available in MLRepo (n=28), we compare ROCs
339 of random forest, SVM with a radial kernel, and SVM with a linear kernel. AUCs are listed within
340 plots and are colored respective to each model.

341 Figure 4. Summary statistics of framework and database comparisons.

342 (A) AUCs random forest (rf) to SVM-Linear (left) and random forest to SVM-Radial (right).

343 Paired t-tests reveal that random forest results in significantly higher AUC than both
344 SVM-Linear (P=0.0014) and SVM-Radial (P=0.00032).

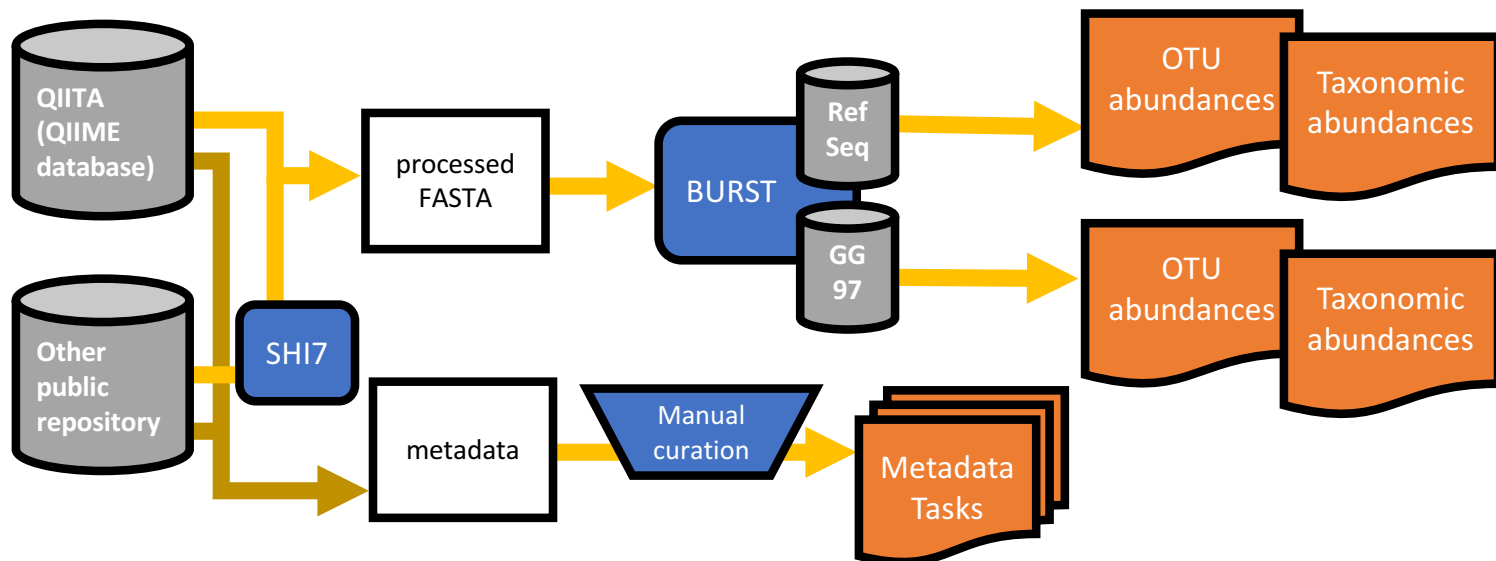
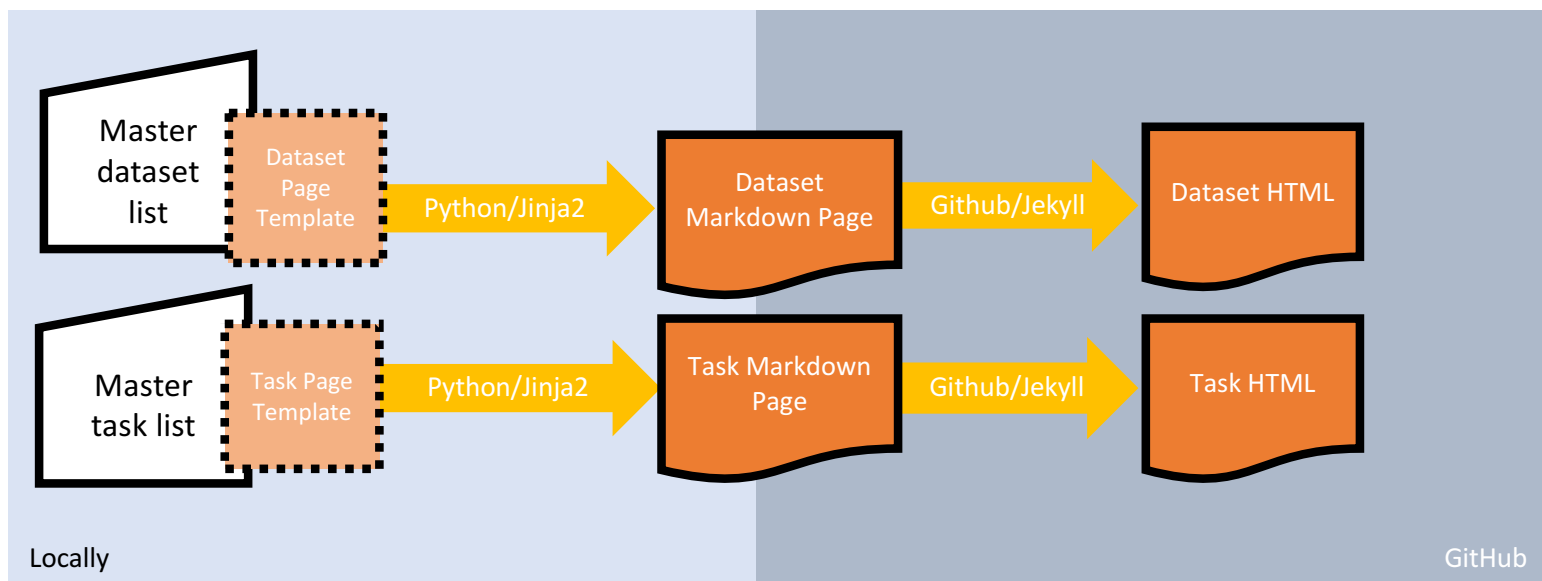
345 (B) Accuracies of random forest to SVM-Linear (left) and random forest to SVM-Radial
346 (right). Paired t-tests reveal that random forest results in significantly better accuracy
347 than SVM-Radial (P=0.03), but not SVM-Linear (P=0.083).

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65

348 (C) AUCs (left) and accuracies (right) of random forest classifications of n=24 tasks using
349 OTUs picked with NCBI RefSeq database or Greengenes database as predictors.
350 Student's t-test reveals that reference database choice has limited impact on
351 classification AUC or accuracy.
352 Lines are colored by the top model for each classification task.

Figure 5. ROCs comparing NCBI RefSeq and Greengenes 97 databases.

354 Sweeping across 16s-based binary classification tasks available in MLRepo (n=24), we
355 compare ROCs of random forest with genus-level taxonomic summaries as predictors from
356 OTU-picking strategies with the NCBI RefSeq prokaryote reference database and the
357 Greengenes 97 reference database. AUCs are listed within plots and are colored respective to
358 each database.

A**B**

A

MLRepo
Machine learning repository for microbiome datasets

[View On GitHub](#)

Available Tasks

- Bacteremia
- Diet
- Antibiotics
- Age
- IBD
- Gender
- Vaginal
- Geography
- Body Habitat
- Cancer
- Obesity
- Diabetes
- Cirrhosis

Available Tasks
Download a single file containing all available tasks

Bacteremia

- bacteremia vs no bacteremia

Diet

- high fat vs low fat diet
- animal vs plant diet, last diet day

Antibiotics

- chlortetracycline vs control, cecal
- chlortetracycline vs control, fecal

B

Task: bacteremia vs no bacteremia
Patients prior to chemotherapy who did or did not develop bacteremia

Project	Montassier 2016
Topic area	Bacteremia
Sample type	human stool
Number of samples	28
Response type	binary
Additional task details	
Multiple samples per subject?	No
Task mapping file	task.txt
OTU file <i>gg97</i>	otutable.txt
Taxa file <i>gg97</i>	taxatable.txt
OTU file <i>RefSeq</i>	otutable.txt
Taxa file <i>RefSeq</i>	taxatable.txt

[back to task index](#)

C

Montassier 2016
Patients prior to chemotherapy who did or did not develop bacteremia

Overview

Description	Patients prior to chemotherapy who did or did not develop bacteremia
Study design	Cross-Sectional
Topic area	Bacteremia
Attributes	Treatment: NObact, bact
Dataset notes	
Number of samples	28
Number of subjects	28

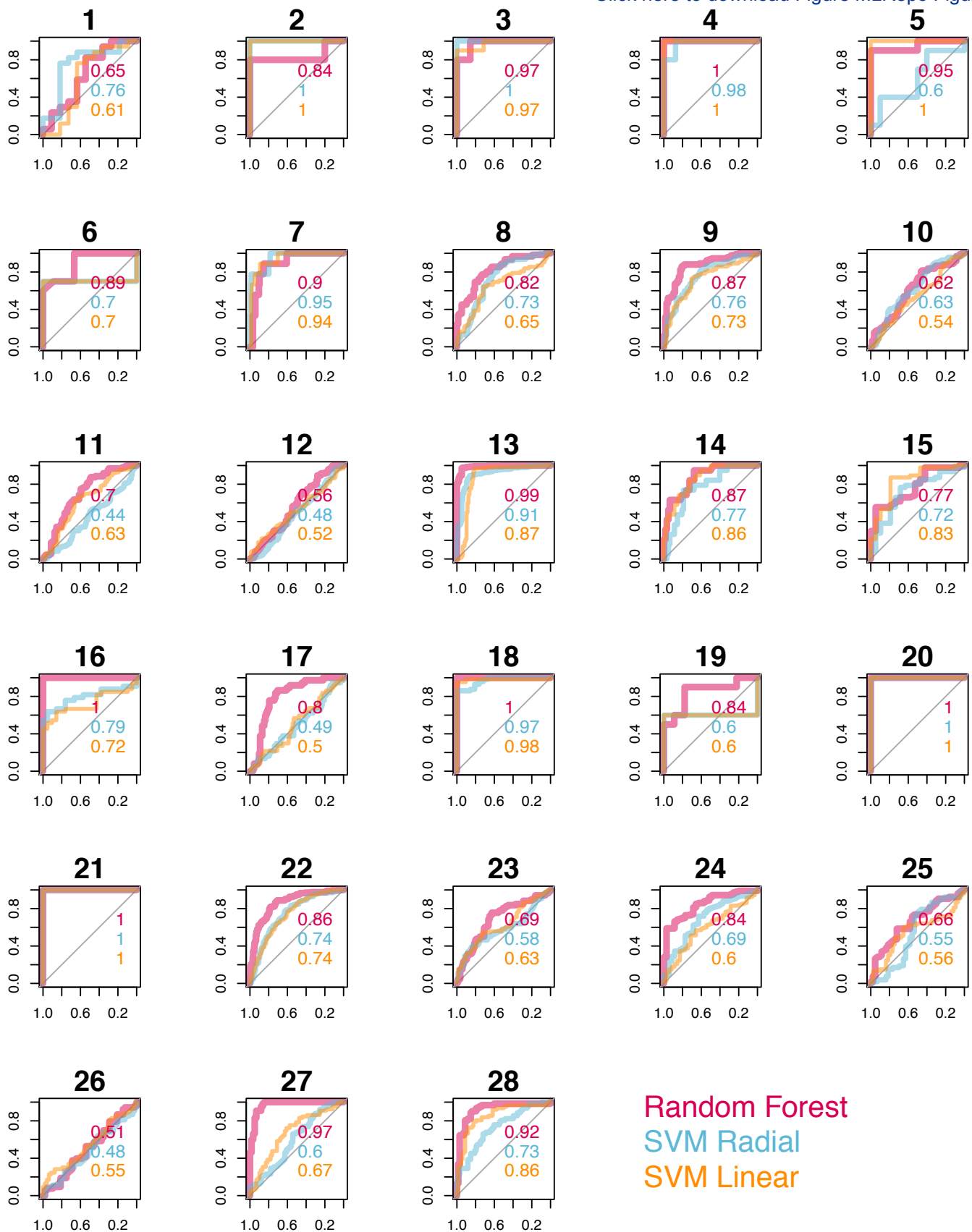
Other Details

16s hypervariable region	V56
Targeted amplicon size	280
Sequencing technology	454
Fraction of sequences mapped to database	
Processed sequences	montassier2016.fasta.gz
Raw metadata file	mapping-orig.txt
Raw sequence source	https://www.ncbi.nlm.nih.gov/sra/SRX733464
Literature source	https://www.ncbi.nlm.nih.gov/pubmed/27121964

[back to task index](#)

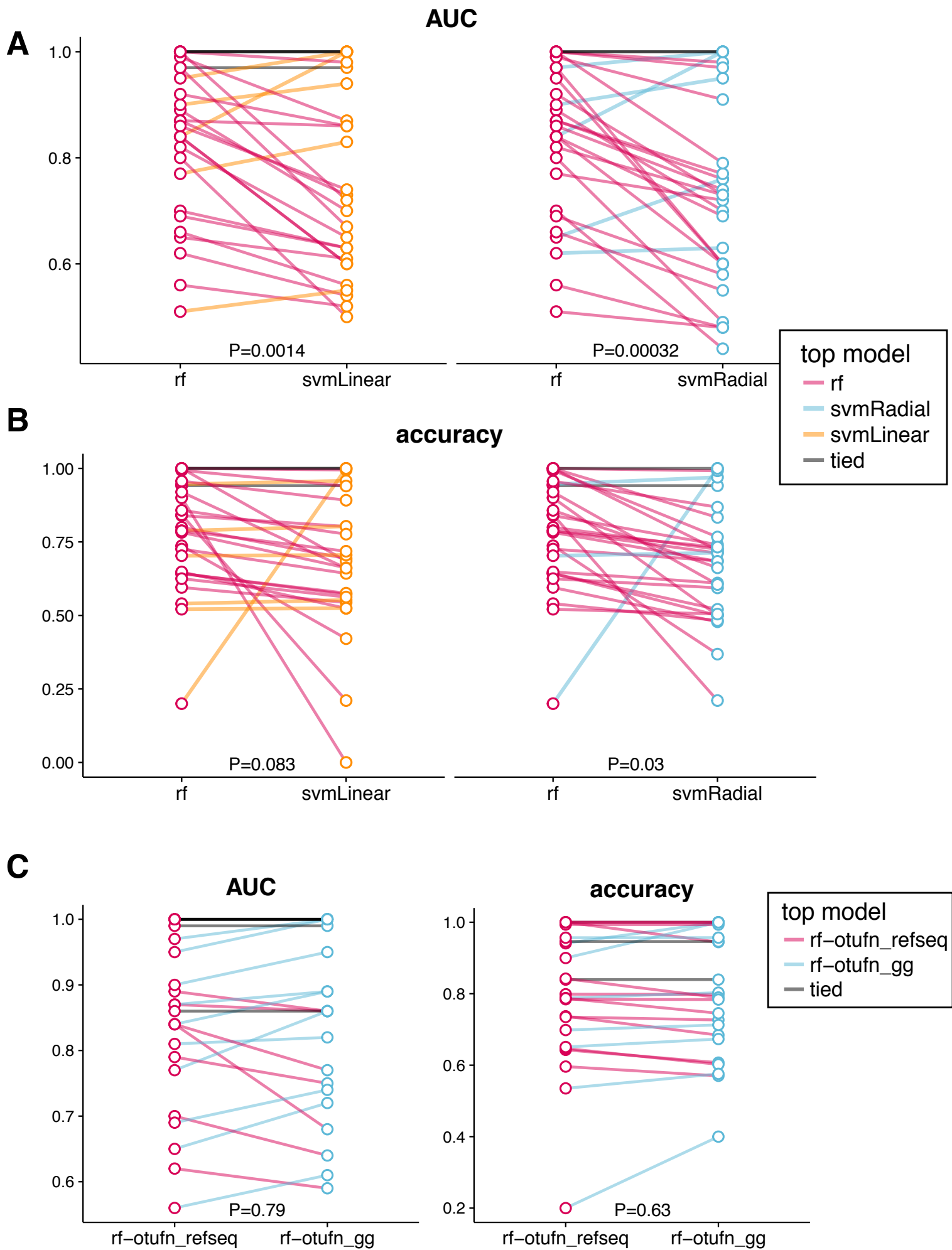
Figure 3

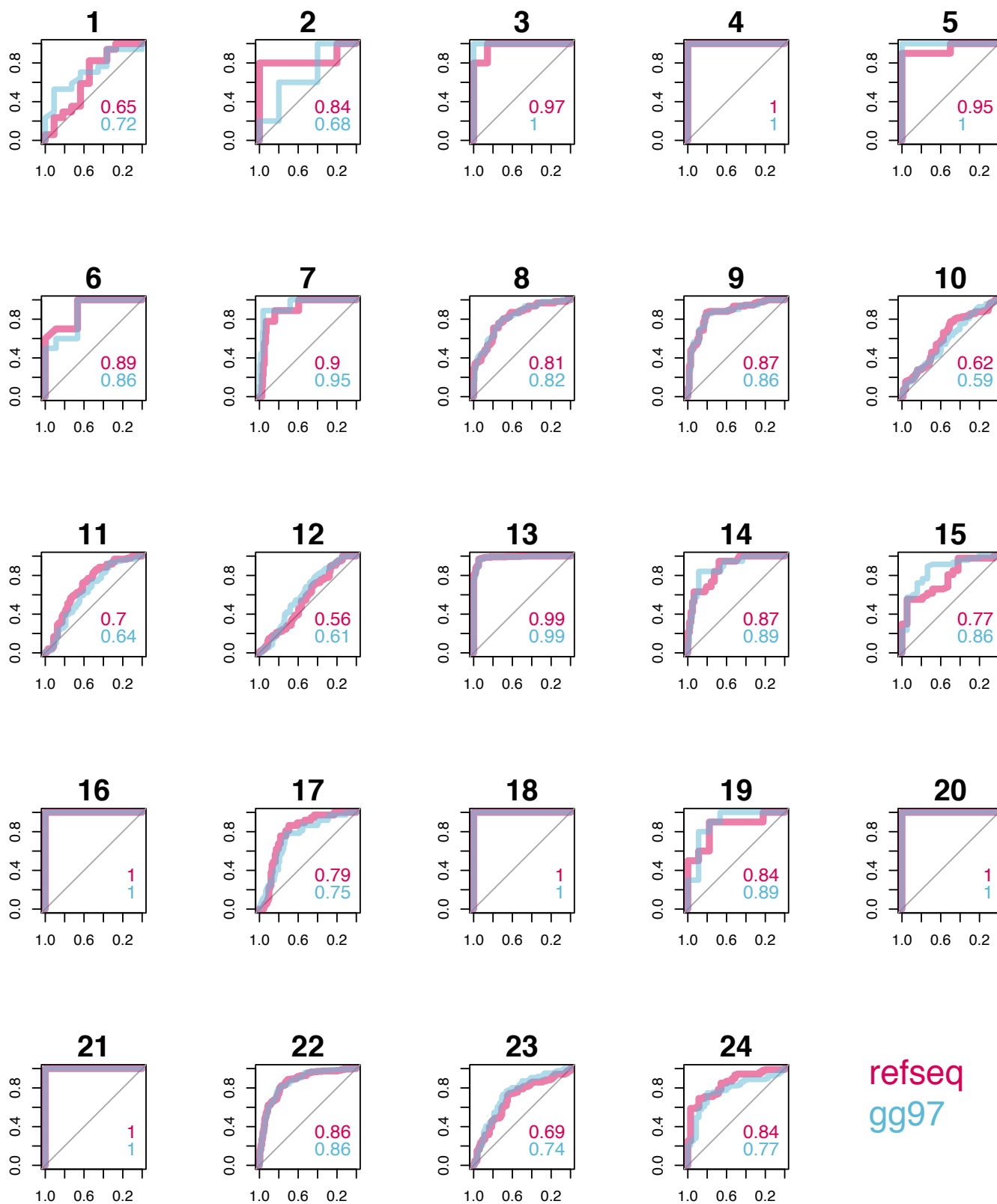
[Click here to download Figure MLRepo Figure 3.pdf](#)



Random Forest
SVM Radial
SVM Linear

- | | | |
|----------------------------------------------|-----------------------------------------------|--------------------------------------------------|
| 1 bacteremia vs no bacteremia | 11 white vs black, vaginal | 21 stool vs tongue |
| 2 high fat vs low fat diet | 12 black vs hispanic, vaginal | 22 subgingival vs supragingival plaque |
| 3 chlortetracycline vs control, cecal | 13 low vs high nugent category | 23 healthy vs tumor biopsy, paired |
| 4 chlortetracycline vs control, fecal | 14 healthy vs cd, stool | 24 lean vs obese, mz/dz/mom |
| 5 penicillin vs vancomycin, cecal | 15 healthy vs uc, stool | 25 normal vs diabetes glucose tolerance |
| 6 penicillin vs vancomycin, fecal | 16 malawi vs venezuela, adults only | 26 impaired vs diabetes glucose tolerance |
| 7 elderly vs young | 17 male vs female, usa | 27 healthy vs type 2 diabetes |
| 8 control vs cd, ileum | 18 us vs malawi, adults only | 28 healthy vs cirrhosis |
| 9 control vs cd, rectum | 19 animal vs plant diet, last diet day | |
| 10 male vs female, stool | 20 gastrointestinal vs oral | |





refseq
gg97

- | | |
|---------------------------------------|----------------------------------------|
| 1 bacteremia vs no bacteremia | 13 low vs high nugent category |
| 2 high fat vs low fat diet | 14 healthy vs cd, stool |
| 3 chlortetracycline vs control, cecal | 15 healthy vs uc, stool |
| 4 chlortetracycline vs control, fecal | 16 malawi vs venezuela, adults only |
| 5 penicillin vs vancomycin, cecal | 17 male vs female, usa |
| 6 penicillin vs vancomycin, fecal | 18 us vs malawi, adults only |
| 7 elderly vs young | 19 animal vs plant diet, last diet day |
| 8 control vs cd, ileum | 20 gastrointestinal vs oral |
| 9 control vs cd, rectum | 21 stool vs tongue |
| 10 male vs female, stool | 22 subgingival vs supragingival plaque |
| 11 white vs black, vaginal | 23 healthy vs tumor biopsy, paired |
| 12 black vs hispanic, vaginal | 24 lean vs obese, mz/dz/mom |