

Evaluation of machine learning methods for 16S rRNA gene data

Running title: Machine learning methods in microbiome studies

Begüm D. Topçuoğlu¹, Jenna Wiens², Patrick D. Schloss^{1†}

† To whom correspondence should be addressed: pschloss@umich.edu

1. Department of Microbiology and Immunology, University of Michigan, Ann Arbor, MI 48109

2. Department of Computer Science and Engineering, University of Michigan, Ann Arbor, MI 49109

1 Abstract

2 Introduction

Advances in sequencing technology and decreasing costs of generating 16S rRNA gene sequences have allowed rapid exploration of human associated microbiome and its health implications. Currently, the human microbiome field is growing at an unprecedented rate and as a result, there is an increasing demand for methods that identify associations between members of the microbiome and human health. However, this is an undertaking as human associated microbial communities are remarkably complex and uneven. It is unlikely that a single species can explain a disease. Instead, subsets of those communities, in relation to one another and to their host, account for the differences in the health outcomes. Therefore, researchers have started to explore the utility of machine learning (ML) models that use microbiota associated biomarkers to predict human health and to understand the microbial ecology of diseases such as liver cirrhosis, colorectal cancer, inflammatory bowel diseases (IBD), obesity, and type 2 diabetes (1–11). ML methods are effective at recognizing and highlighting patterns in complex microbial datasets. However, currently the field's use of ML lacks clarity and consistency on which methods are used and how these methods are implemented. Moreover, there is a lack of deliberation on why a particular ML model is utilized. Recently, there is a trend towards using more complex ML models such as random forest, extreme gradient boosting and neural networks without a discussion on if and how much model interpretability is necessary for the study (11–14). The lack of transparency on methodology and modeling decisions can have negative consequences on model reproducibility and reliability which we need avoid by (1) implementing consistent, accessible and transparent machine learning practices; (2) selecting ML models that reflect the goal of the study as it will inform the expectations of model accuracy, complexity, interpretability and computational efficiency.

To showcase transparent ML methodologies and to shed light on how much ML model selection can affect modeling results, we performed an empirical analysis comparing several different ML models using the same dataset. We used a previously published colorectal cancer (CRC) study (3) which had fecal 16S rRNA gene sequences and human hemoglobin concentrations from 490 patients. We built ML models using fecal 16S rRNA gene sequences and human hemoglobin concentrations to predict patients with normal or screen relevant neoplasias (SRN) disease status. The study

had 261 normal and 229 SRN samples. We established modeling pipelines for L2-regularized logistic regression, L1 and L2 support vector machines (SVM) with linear and radial basis function kernels, a decision tree, random forest and XGBoost which increase in complexity and decrease in interpretability respectively. Our ML pipeline performed 100 data-splits and utilized held-out test data to evaluate generalization and prediction performance of each ML model. The mean AUROC varied from 0.5 (std \pm 0.01) to 0.82 (std \pm 0.04). Random Forest and XGBoost had the highest mean AUROC for detecting SRN. Despite the simplicity, the L1-regularized linear kernel SVM followed Random Forest and XGBoost in performance. In terms of computational efficiency, L2 SVM with linear kernel trained the fastest (0.123 hours, std \pm 0.014), while XGBoost took the longest (106.568 hours, std \pm 6.839). We found that mean cross-validation and testing AUROC could vary by as much as 0.238, which highlights the importance of a separate held-out test set for evaluation. Aside from evaluating generalization and classification performances for each of these models, this study established standards for modeling pipelines of microbiome-associated machine learning models.

Results

Model selection and construction

We established modeling pipelines for binary classification with L2-regularized logistic regression, L1 and L2 support vector machines (SVM) with linear and radial basis function kernels, a decision tree, random forest and XGBoost to emphasize the differences in model accuracy, complexity, interpretability, computational efficiency and scalability due to model selection. We randomly split the data into stratified training and test sets 100 times with a 80/20 split [Figure 1]. The training set consisted of 393 patients (209 SRN), while the test set was composed of 97 patients (52 SRN). The training data was used for training purposes and validation of parameter selection, and the test set was used for evaluation purposes. Validation of parameter selection was repeated 100 times to get a robust parameter selection for each model [Figure 1].

The prediction and generalization performance of classifiers during cross-validation and

when applied to the held-out test data.

We evaluated the prediction performance of seven binary classification models when applied to held-out test data over 100 data-splits using AUROC as classification performance metric. Random Forest and XGBoost had the highest mean AUROC for detecting SRN, 0.816 (std \pm 0.039) and 0.814 (std \pm 0.04) respectively [Figure 2]. L1 linear SVM and decision tree had significantly lower AUROC values, 0.503 (std \pm 0.027) and 0.741 (std \pm 0.038) [Figure 2]. However, they had significantly higher performances than L2 linear SVM, RBF SVM and L2 logistic regression which had mean AUROC values of 0.501 (std \pm 0.013), 0.68 (std \pm 0.05) and 0.672 (std \pm 0.05) respectively [Figure 2]. We also evaluated the generalization performance of each classifier by comparing their mean cross-validation AUROC and mean testing AUROC. We found a statistically significant difference for classifiers L2 support vector machines (SVM) with linear and radial basis function kernels and L2 logistic regression (p-values = 2.4e-37, 1.3e-14 and 2.6e-21 respectively). These differences were 0.222, 0.05 and 0.064, respectively [Figure 2] .

The complexity and interpretability of each classifier.

We interpreted the feature weights of L1 and L2 SVM with linear kernel and regression coefficients of L2 logistic regression using the training data. To get a sense of the importance of each feature we randomly subsampled 80% of the data 100 times and trained 100 models; this resulted in 100 different sets of feature weights. We reported the mean weights of the 10 most important features [Figure 4].

The computational efficiency of each classifier.

Linear models trained faster than non-linear models. L2 logistic Regression and L1 and L2 SVM with linear kernel had training times of 7.2 min, (std \pm 0.6), 10.2 min, (std \pm 0.6) and 13.2 min, (std \pm 1.8) respectively. Whereas, SVM with radial basis function kernel, decision tree, random forest and xgboost had training times of 2.01 hrs, (std \pm 0.28), 4.2 hrs, (std \pm 0.49), 106.57 hrs, (std \pm 6.84) and 116.83 hrs, (std \pm 17.24), respectively.

81 **Conclusions**

82 **Materials and Methods**

83 **Data collection**

84 The data used for this analysis are stool bacterial abundances, stool hemoglobin levels and clinical
85 information of the patients recruited by Great Lakes-New England Early Detection Research
86 Network study. These data were obtained from Sze et al (15). The stool samples were provided by
87 recruited adult participants who were undergoing scheduled screening or surveillance colonoscopy.
88 Colonoscopies were performed and fecal samples were collected from participants in four locations:
89 Toronto (ON, Canada), Boston (MA, USA), Houston (TX, USA), and Ann Arbor (MI, USA).
90 Patients' colonic disease status was defined by colonoscopy with adequate preparation and tissue
91 histopathology of all resected lesions. Patients with an adenoma greater than 1 cm, more than
92 three adenomas of any size, or an adenoma with villous histology were classified as advanced
93 adenoma. Study had 172 patients with normal colonoscopies, 198 with adenomas and 120 with
94 carcinomas. Of the 198 adenomas, 109 were identified as advanced adenomas. Stool provided by
95 the patients was used for Fecal Immunological Tests (FIT) which measure human hemoglobin
96 concentrations and for 16S rRNA gene sequencing to measure bacterial population abundances.
97 The bacterial abundance data was generated by Sze et al, by processing 16S rRNA sequences
98 in Mothur (v1.39.3) using the default quality filtering methods, identifying and removing chimeric
99 sequences using VSEARCH and assigning to OTUs at 97% similarity using the OptiClust algorithm
100 (16–18).

101 **Data definitions and pre-processing**

102 The colonic disease status is re-defined as two encompassing classes; Normal or Screen Relevant
103 Neoplasias (SRNs). Normal class includes patients with non-advanced adenomas or normal colons
104 whereas SRN class includes patients with advanced adenomas or carcinomas. Colonic disease
105 status is the label predicted with each classifier. The bacterial abundances and FIT results are the

features used to predict colonic disease status. Bacterial abundances are discrete data in the form of Operational Taxonomic Unit (OTU) counts. There are 6920 OTUs for each sample. FIT levels are continuous data present for each sample. Because the data are in different scales, features are transformed by scaling each feature to a [0-1] range (Table 1).

Learning the Classifier

To train and validate our model, labeled data is randomly split 80/20 into a training set and testing set. Then, seven binary class classifiers, L2 logistic regression, L1 and L2 linear support vector machines (SVM), radial basis function SVM, decision tree, random forest and XGBoost, are learned. The training set is used for training purposes and validation of hyperparameter selection, and the test set is used for evaluation purposes. Hyperparameters are selected using 5-fold cross-validation with 100-repeats on the training set. Since the colonic disease status are not uniformly represented in the data, 5-fold splits are stratified to maintain the overall label distribution on the training set.

Classifier Performance

The classification performance of learned classifier is evaluated on the labeled held-out testing set. The optimal classifier with optimal hyperparameters selected in the cross-validation step is used to produce a prediction for the testing set. The performance of this prediction is measured in terms of the sensitivity and specificity, in addition to Area Under the Curve (AUC) metrics. This process of splitting the data, learning a classifier with cross-validation, and testing the classifier is repeated on 100 different splits. In the end cross-validation AUC and testing AUC averaged over the 100 different training/test splits are reported. Hyperparameter budget and performance for each split is also reported.

Figure 1. Generalization and classification performance of modeling methods AUC values of all cross validation and testing performances. The boxplot shows quartiles at the box ends and the statistical median as the horizontal line in the box. The whiskers show the farthest points that are not outliers. Outliers are data points that are not within $3/2$ times the interquartile ranges.

References

1. **Zeller G, Tap J, Voigt AY, Sunagawa S, Kultima JR, Costea PI, Amiot A, Böhm J, Brunetti F, Habermann N, Hercog R, Koch M, Luciani A, Mende DR, Schneider MA, Schrotz-King P, Tournigand C, Tran Van Nhieu J, Yamada T, Zimmermann J, Benes V, Kloor M, Ulrich CM, Knebel Doeberitz M von, Sobhani I, Bork P.** 2014. Potential of fecal microbiota for early-stage detection of colorectal cancer. *Mol Syst Biol* **10**. doi:10.15252/msb.20145645.
2. **Zackular JP, Rogers MAM, Ruffin MT, Schloss PD.** 2014. The human gut microbiome as a screening tool for colorectal cancer. *Cancer Prev Res* **7**:1112–1121. doi:10.1158/1940-6207.CAPR-14-0129.
3. **Baxter NT, Koumpouras CC, Rogers MAM, Ruffin MT, Schloss PD.** 2016. DNA from fecal immunochemical test can replace stool for detection of colonic lesions using a microbiota-based model. *Microbiome* **4**. doi:10.1186/s40168-016-0205-y.
4. **Baxter NT, Ruffin MT, Rogers MAM, Schloss PD.** 2016. Microbiota-based model improves the sensitivity of fecal immunochemical test for detecting colonic lesions. *Genome Medicine* **8**:37. doi:10.1186/s13073-016-0290-3.
5. **Hale VL, Chen J, Johnson S, Harrington SC, Yab TC, Smyrk TC, Nelson H, Boardman LA, Druliner BR, Levin TR, Rex DK, Ahnen DJ, Lance P, Ahlquist DA, Chia N.** 2017. Shifts in the fecal microbiota associated with adenomatous polyps. *Cancer Epidemiol Biomarkers Prev* **26**:85–94. doi:10.1158/1055-9965.EPI-16-0337.
6. **Pasolli E, Truong DT, Malik F, Waldron L, Segata N.** 2016. Machine learning meta-analysis of large metagenomic datasets: Tools and biological insights. *PLoS Comput Biol* **12**. doi:10.1371/journal.pcbi.1004977.
7. **Sze MA, Schloss PD.** 2016. Looking for a signal in the noise: Revisiting obesity and the microbiome. *mBio* **7**. doi:10.1128/mBio.01018-16.
8. **Walters WA, Xu Z, Knight R.** 2014. Meta-analyses of human gut microbes associated with

obesity and IBD. *FEBS Lett* **588**:4223–4233. doi:10.1016/j.febslet.2014.09.039.

9. **Vázquez-Baeza Y, Gonzalez A, Xu ZZ, Washburne A, Herfarth HH, Sartor RB, Knight R.** 2018. Guiding longitudinal sampling in IBD cohorts. *Gut* **67**:1743–1745. doi:10.1136/gutjnl-2017-315352.

10. **Qin N, Yang F, Li A, Prifti E, Chen Y, Shao L, Guo J, Le Chatelier E, Yao J, Wu L, Zhou J, Ni S, Liu L, Pons N, Batto JM, Kennedy SP, Leonard P, Yuan C, Ding W, Chen Y, Hu X, Zheng B, Qian G, Xu W, Ehrlich SD, Zheng S, Li L.** 2014. Alterations of the human gut microbiome in liver cirrhosis. *Nature* **513**:59–64. doi:10.1038/nature13568.

11. **Geman O, Chiuchisan I, Covasa M, Doloc C, Milici M-R, Milici L-D.** 2018. Deep learning tools for human microbiome big data, pp. 265–275. *In* Balas, VE, Jain, LC, Balas, MM (eds.), *Soft computing applications*. Springer International Publishing.

12. **Galkin F, Aliper A, Putin E, Kuznetsov I, Gladyshev VN, Zhavoronkov A.** 2018. Human microbiome aging clocks based on deep learning and tandem of permutation feature importance and accumulated local effects. *bioRxiv*. doi:10.1101/507780.

13. **Reiman D, Metwally A, Dai Y.** 2017. Using convolutional neural networks to explore the microbiome, pp. 4269–4272. *In* 2017 39th annual international conference of the IEEE engineering in medicine and biology society (EMBC).

14. **Fioravanti D, Giarratano Y, Maggio V, Agostinelli C, Chierici M, Jurman G, Furlanello C.** 2017. Phylogenetic convolutional neural networks in metagenomics. *arXiv:170902268 [cs, q-bio]*.

15. **Sze MA, Schloss PD.** 2018. Leveraging existing 16S rRNA gene surveys to identify reproducible biomarkers in individuals with colorectal tumors. *mBio* **9**:e00630–18. doi:10.1128/mBio.00630-18.

16. **Schloss PD, Westcott SL, Ryabin T, Hall JR, Hartmann M, Hollister EB, Lesniewski RA, Oakley BB, Parks DH, Robinson CJ, Sahl JW, Stres B, Thallinger GG, Van Horn DJ, Weber CF.** 2009. Introducing mothur: Open-Source, Platform-Independent, Community-Supported Software for Describing and Comparing Microbial Communities. *Appl Environ Microbiol*

181 **75:7537–7541.**

182 17. **Westcott SL, Schloss PD.** 2017. OptiClust, an Improved Method for Assigning

183 Amplicon-Based Sequence Data to Operational Taxonomic Units. *mSphere* **2**. doi:10.1128/mSphereDirect.00073-17

184 18. **Rognes T, Flouri T, Nichols B, Quince C, Mahé F.** 2016. VSEARCH: A versatile open source

185 tool for metagenomics. *PeerJ* **4**:e2584. doi:10.7717/peerj.2584.