

Taxonomic Resolution Matters for Microbiome-Based Classification of Colorectal Cancer

Other title options:

Optimal Taxonomic Resolution for Microbiome-Based Classification of Colorectal Cancer ASV level resolution is unnecessary for prediction of SRNs from microbiome data . . .

Courtney R. Armour, Begüm D.Topçuoğlu, Andrea Garretto, Patrick D. Schloss [†]

[†] To whom correspondence should be addressed:

pschloss@umich.edu

Department of Microbiology

University of Michigan

Ann Arbor, MI 48109

observation format - max 1200 words, 2 figures, 25 ref

¹³ **Abstract (max 250 words)**

¹⁴ **Importance (max 150 words)**

Introduction (~250 words)

Colorectal cancer is one of the most common cancers in men and women and a leading cause of cancer related deaths in the United States {}. Early detection and treatment are essential to increase survival rates, but for a variety of reasons including the invasiveness and high cost of screening (i.e. colonoscopy), many people do not comply with recommended screening guidelines {} prompting a need for low cost, non-invasive detection methods. A growing body of research points to the gut microbiome as a promising target for non-invasive screening to detect screen relevant neoplasias (neoplasms?) (SRNs) consisting of pre-cancerous polyps (i.e. advanced adenomas) and cancer. Efforts to realize the diagnostic potential of the gut microbiome in detecting SRNs have focused on machine learning (ML) methods using abundances from operational taxonomic unit (OTU) classifications based on the amplicon sequencing of the 16S rRNA gene {}. However, whether this is the optimal taxonomic resolution for classifying SRNs from microbiome data is unknown. Additionally, recent work has pushed for the use of amplicon sequence variants (ASVs) to replace OTUs for marker-gene analysis because of the improved resolution with ASVs {}. However, whether the additional resolution provided by ASVs is useful for ML classification is unclear. Since ML classification relies on consistent differences between groups, its possible that the resolution at the ASV level is too individualized to accurately differentiate groups. Topçuoğlu *et al* {mBio 2020} recently demonstrated effective application of machine learning (ML) to microbiome based classification problems and developed a framework for applying ML practices in a more reproducible way (mikropml). This analysis utilizes the reproducible framework developed by Topçuoğlu *et al* to quantify which ML method and taxonomic level produce the best performing classifier for detecting SRNs from microbiome data.

Results (~700 words)

Utilizing publicly available 16S rRNA sequence data from patients with cancer and healthy controls, we generated abundance tables with Mothur {} annotated to Phylum, Class, Order, Family, Genus, OTU and ASV levels. Using the taxonomic abundance data, we quantified how accurately samples could be classified as “normal” or “SRN” (i.e. advanced adenoma or cancer) using machine learning (methods). Across the five machine learning methods tested, model performance tended to increase with taxonomic level usually peaking around genus/OTU level before dropping off slightly with ASVs (Supplemental Figure 1). Random forest (RF) was consistently the top performer at most taxonomic levels. [TODO: dig into literature on this: RF might be more appropriate for microbiome analysis since its more suitable for zero inflated data]. Within the RF model, the highest AUROCs were observed for family (median AUROC: 0.687), genus (median AUROC: 0.686), and OTU (median AUROC: 0.698) level data with no significant difference between the

three (Figure 1). Performance with ASVs (median AUROC: 0.676) was significantly lower than OTU and genus ($p < 0.05$) but not family level ($P = 0.06$). These results were consistent with ASVs generated with DADA2 {} (median AUROC: 0.66). These results suggest that finer taxonomic resolution is not necessarily better for differentiating a persons cancer status from the microbiome.

One hypothesis for the observation that model performance increases from phylum to OTU level then drops slightly at ASV level is that at higher taxonomic levels (e.g. phylum) there are too few taxa and too much overlap to reliably differentiate between cases and controls. At the level of genus/OTU data there is enough data and variation but at the ASV level, the data is too specific to individuals and doesn't overlap enough. Examination of the prevalence of taxa in samples at each level supports this idea. A majority of taxa are present in greater than 75% of samples at the phylum (67% of taxa) and class (63% of taxa) levels). The opposite is observed at the OTU and ASV level where 60% and 53% of taxa respectively are only present in less than 25% of samples (Supplemental Figure 2). **SUMMARY SENTENCE**

TODO: what are important taxa, patterns along levels (e.g. fusobacterium always in top)

Of note, the ML pipeline includes a pre-processing step that occurs prior to training and classifying the ML models (methods). As part of this step, features are removed that will not provide useful information to build the model. For example, strongly correlated features provide the same information to build the model and thus can be collapsed. Additionally, features with zero or near-zero variance will not help the model differentiate groups and thus can be removed. Interestingly, despite starting with 104106 features at the ASV level, only 478 (0.5%) remained after pre-processing. At the OTU level, 20079 of the 705 features (3.5%) remained after preprocessing (Table 1).

TODO: check why ASVs were removed in preprocessing (nzv?)

(depending on answer can discuss why features were removed - e.g. only in one or a few samples)

Discussion/Conclusions (~250 words)

These results demonstrate that consideration of the appropriate taxonomic resolution for utilizing the microbiome as a predictive tool. In general, we found that finer resolution (e.g. ASVs) does not add additional sensitivity to model prediction and at the ASV level actually impedes model performance due to the sparsity of shared taxa. Additionally, utilizing genus or family level data could allow for merging data generated from different 16S regions or sequencing platforms.

- Overall, the fine resolution of ASVs is unnecessary for ML classification with microbiome data since the vast majority of the ASVs are removed during preprocessing and the model performance is diminished

76 compared to other taxonomic levels.

- 77 • since family/genus level is just as good as otu could we use data from different 16S regions?

78 **Materials and Methods**

79 **Dataset.** Raw 16S rRNA gene amplicon sequence data isolated from human gut samples {Baxter} was
80 downloaded from NCBI SRA (accession #). This dataset contains stool samples from 490 subjects. Based
81 on the available metadata, samples categorized as normal, high risk normal, or adenoma were labeled
82 “normal” for this analysis and samples categorized as advanced adenoma or carcinoma were labeled as
83 “screen relevant neoplasia” (SRN). This resulted in a total of 261 “normal” samples and 229 “SRN” samples.

84 **Data processing.** Sequence data was processed with Mothur (1.44.3) using the SILVA reference database
85 (v132) to produce count tables for phylum, class, order, family, genus, OTU, and ASV following the Schloss
86 Lab MiSeq SOP described on the Mothur website {}. ASV level data was also produced using DADA2 to
87 ensure consistent results with a different pipeline. Data was processed following the DADA2 pipeline, but
88 setting pool=TRUE to infer ASVs from the whole dataset rather than per sample. The resulting ASV table
89 was subsampled for consistency with the Mothur data.

90 **Machine Learning.** Machine learning models were run with the R package mikropml (v0.0.2) {} to predict
91 the diagnosis category (normal vs SRN) of each sample. Data was preprocessed to normalize values
92 (scale/center? range?), remove values with zero or near-zero variance, and collapse collinear features using
93 default parameters. Initially the models were run with default hyperparameters, but were expanded if the
94 peak performance was at the edge of the hyperparameter range. Each taxonomic model taxonomic level
95 combination (e.g. Random Forest on genus) was run with 100 different seeds (1-100). Each seed split the
96 data into a training (80%) and testing (20%) set, and output performance of the training and testing as area
97 under the receiver operating curve (AUROC).

98 To compare performance between taxonomic levels and models, P values were calculated as previously
99 described {begum}. Another output from the mikropml package is the permuted feature importance which is
100 quantified by iteratively permuting each feature in the model and assessing the change in model performance.
101 Features are presumed to be important if the performance of the model, measured by the AUROC, decreases
102 when that feature is permuted. Ranking of feature importance was determined by ordering the features
103 based on the change in AUROC where features with a larger decrease in AUROC are ranked higher in
104 importance.

105 To quantify prevalence of the features, the number of samples with non-zero abundance was divided by the

106 total number of samples resulting in values ranging from 0 to 1 where 0 indicates the feature is not found in
107 any samples, 0.5 indicates the feature is found in half of the samples, an 1 indicating the feature is found in
108 all of the samples.

109 All code is available at: **TODO: link to code**

110 **Acknowledgements**

Figures

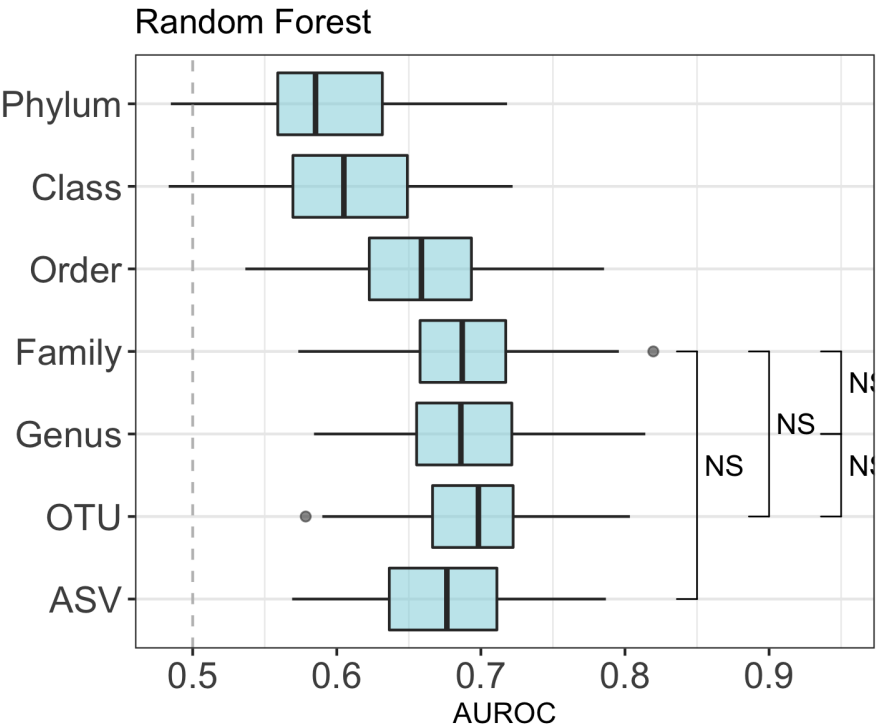


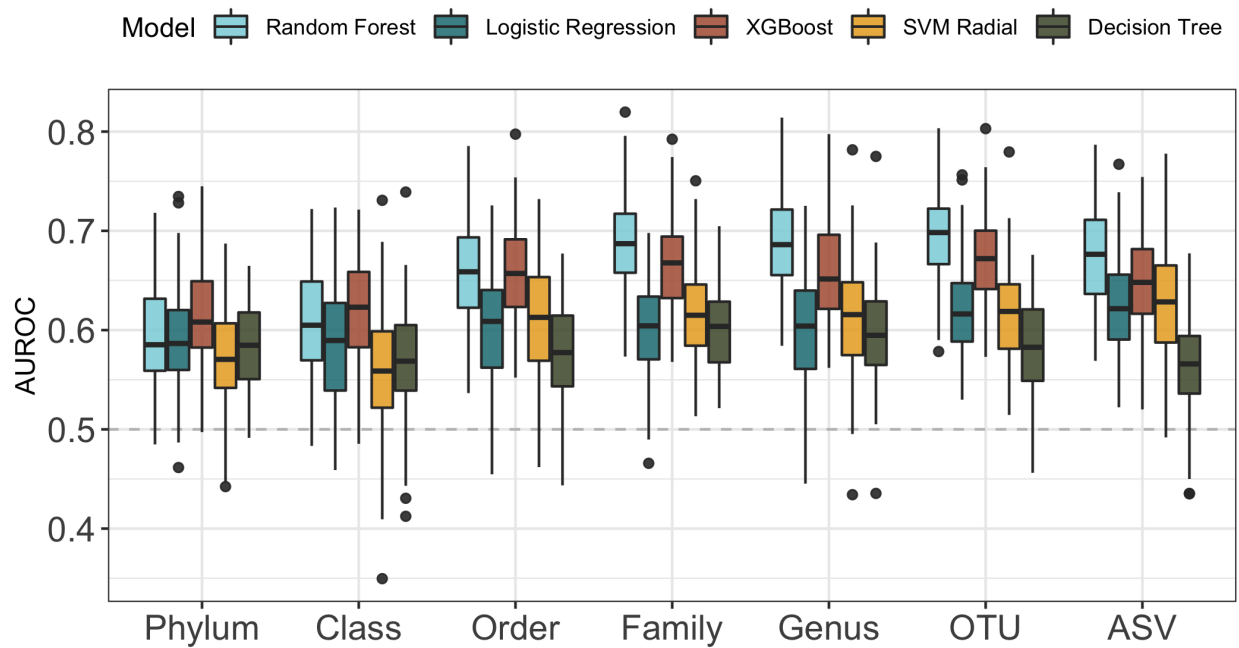
Figure 1: Random Forest Model Performance with Significance. Boxplots of AUROC values on the test dataset for 100 seeds predicting SRNs using a Random Forest model. Dashed line denotes AUROC of 0.5 which is equivalent to random classification. Differences in AUROC values between taxonomic levels are significant unless otherwise specified with “NS” (not significant)

117 Tables

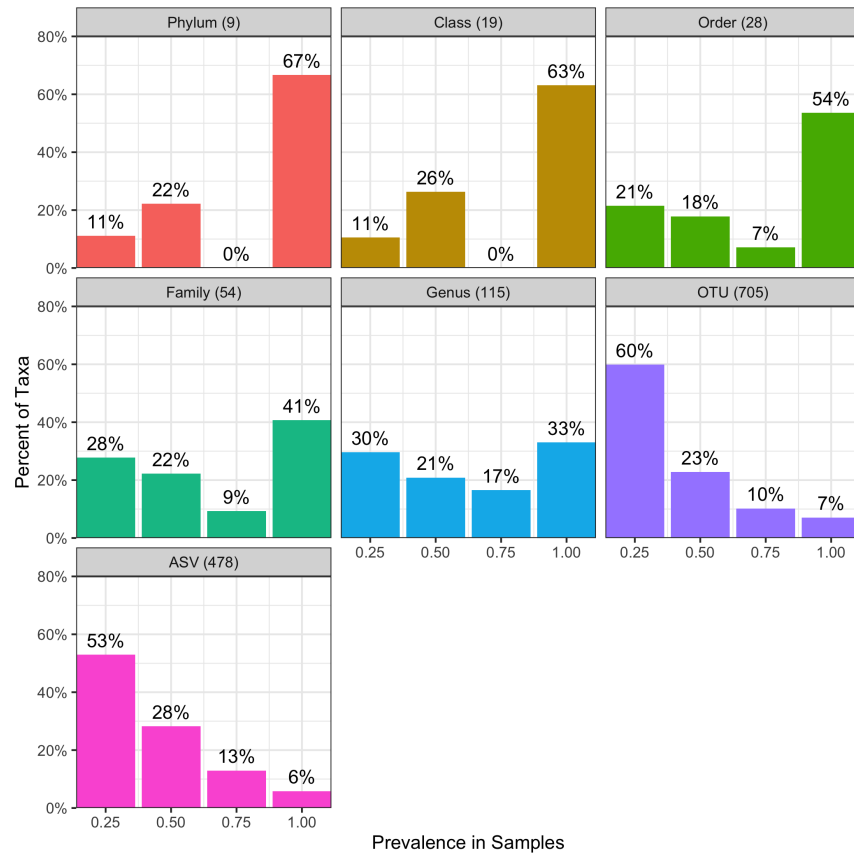
Taxonomic Level	Number of Features	Number of Features After Preprocessing	Percent of Features Kept After Preprocessing
Phylum	19	9	47.4 %
Class	36	19	52.8 %
Order	65	28	43.1 %
Family	124	54	43.5 %
Genus	316	115	36.4 %
OTU	20079	705	3.5 %
ASV	104106	478	0.5 %

118 **Table 1: Summary of Features.** Overview of the number of features at each taxonomic level before and
119 after preprocessing as described in the methods.

Supplemental Figures



Supplemental Figure 1: Model Performance across Taxonomy. Boxplot of AUROC values on the test dataset for 100 seeds for each model type across all taxonomic levels.



Supplemental Figure 2: Prevalence of Taxa in Samples. Distribution of taxa prevalence in samples at each taxonomic level.