

1 MB-SupCon: Microbiome-based predictive models via Supervised 2 Contrastive Learning

3
4 Sen Yang (Sen.Yang@UTSouthwestern.edu)¹,
5 Shidan Wang (Shidan.Wang@UTSouthwestern.edu)²,
6 Yiqing Wang (lucy@mail.smu.edu)¹,
7 Ruichen Rong (Ruichen.Rong@UTSouthwestern.edu)²,
8 Jiwoong Kim (Jiwoong.Kim@UTSouthwestern.edu)²,
9 Bo Li (Bo.Li@UTSouthwestern.edu)^{3,4},
10 Andrew Y. Koh (Andrew.Koh@UTSouthwestern.edu)^{3,5,6},
11 Guanghua Xiao (Guanghua.Xiao@UTSouthwestern.edu)^{2,4},
12 Qiwei Li (Qiwei.Li@UTDallas.edu)⁷,
13 Dajiang Liu (dajiang.liu@psu.edu)^{8,*},
14 Xiaowei Zhan (Xiaowei.Zhan@UTSouthwestern.edu)^{1,2,9*}

15
16 ¹Department of Statistical Science, Southern Methodist University, Dallas, TX 75275
17 ²Quantitative Biomedical Research Center, Department of Population and Data Sciences, University of Texas
18 Southwestern Medical Center, Dallas, TX, 75390
19 ³Harold C. Simmons Cancer Center, University of Texas Southwestern Medical Center, Dallas, Texas, 75390
20 ⁴Department of Bioinformatics, University of Texas Southwestern Medical Center, Dallas, Texas, 75390
21 ⁵Department of Microbiology, University of Texas Southwestern Medical Center, Dallas, Texas, 75390
22 ⁶Department of Paediatrics, University of Texas Southwestern Medical Center, Dallas, Texas, 75390
23 ⁷Department of Mathematical Sciences, University of Texas at Dallas, Richardson, TX 75080
24 ⁸Department of Public Health Sciences, Pennsylvania State University, Hershey, PA, 17033
25 ⁹Center for Genetics of Host Defense, University of Texas Southwestern Medical Center, Dallas, Texas, 75390

26
27 *To whom correspondence should be addressed

28
29

30

Abstract

31 Human microbiome consists of trillions of microorganisms. Microbiota can modulate the host
32 physiology through molecule and metabolite interactions. Integrating microbiome and metabolomics
33 data have the potential to predict different diseases more accurately. Yet, most datasets only measure
34 microbiome data but without paired metabolome data. Here, we propose a novel integrative modeling
35 framework, Microbiome-based Supervised Contrastive Learning Framework (MB-SupCon). MB-
36 SupCon integrates microbiome and metabolome data to generate microbiome embeddings, which can
37 be used to improve the prediction accuracy in datasets that only measure microbiome data. As a proof
38 of concept, we applied MB-SupCon on 720 samples with paired 16S microbiome data and
39 metabolomics data from patients with type 2 diabetes. MB-SupCon outperformed existing prediction
40 methods and achieves high average prediction accuracies for insulin resistance status (84.62%), sex
41 (78.98%), and race (80.04%). Moreover, the microbiome embeddings form separable clusters for
42 different covariate groups in the lower-dimensional space, which enhances data visualization. We
43 also applied MB-SupCon on a large inflammatory bowel disease study and observed similar
44 advantages. Thus, MB-SupCon could be broadly applicable to improve microbiome prediction
45 models in multi-omics disease studies.

46

47 **1 Introduction**

48 The human microbiome is a collection of living microorganisms cohabitating in distinct body niches
49 [1, 2]. The microbiome significantly impacts human health, including diseases and treatments [3].
50 Accordingly, it is possible to use microbiome measurements to predict host physiologic conditions
51 non-invasively. Creating microbiome-based prediction models has great benefits for medical research
52 [4].

53 Earlier work on microbiome-based prediction models using microbiome abundances includes random
54 forest, support vector machines models[5]. While identification and quantification of microbiome
55 taxa using microbiome data alone lead to associative and correlative insights, multi-omics can offer
56 mechanistic insights and potentially improve prediction accuracy over models based on microbiomes
57 alone. For example, in colorectal cancer, specific bacterial species has been associated with increased
58 disease risk [6]. Follow-up mechanistic studies further elucidated the functions of the pathogenic
59 species through multi-omics data analysis [7, 8]. Similar multi-omics approaches, especially in
60 microbiome and metabolomics, have been applied to other diseases [9, 10]. To leverage multi-omics
61 data features and unleash the potential of non-invasive microbiome biomarkers, we aim to develop a
62 general framework for phenotype prediction using microbiome data.

63 Statistical learning and artificial intelligence research have advanced microbiome-based prediction
64 models. Earlier work utilized taxonomic abundance data and linear or logistic regression models with
65 penalties (e.g., LASSO model, and elastic net model) [11]. More recent approaches integrate multi-
66 omics data using partial least squares (PLS), partial least squares-discriminant analysis (PLS-DA), or
67 canonical correlation analysis (CCA) [12]. These models rely on linear transformations of original
68 features in supervised or unsupervised learning. Recently, contrastive learning has been introduced
69 in the analysis of the multi-omics data [13] that can capture non-linear relationships between features.
70 For example, a simple framework for unsupervised contrastive learning (simCLR) achieves state-of-
71 the-art prediction performance [14]. Supervised contrastive learning (SupCon) in computer vision
72 tasks also demonstrated superior robustness and prediction accuracies [15], and these advantages have
73 solid theoretical support [16]. Inspired by the success of these approaches, we propose a novel
74 supervised-learning framework (MB-SupCon) based on non-linear transformations of multi-omics
75 datasets, which achieve robust and accurate prediction performance. Our method architecture is
76 intuitive and requires only modest-sized multi-omics data. We demonstrate MB-SupCon's utility
77 using data from a published type 2 diabetes study where MB-SupCon-based model improves
78 prediction accuracies by a large margin; Another independent application of MB-SupCon to an
79 Inflammatory Bowel Disease (IBD) study also produced consistent improvements. Moreover, we
80 demonstrated that the microbiome embeddings from MB-SupCon can better separate different
81 phenotype groups and lead to more informative visualizations of the data. We posit that our
82 microbiome-based prediction model can easily be applied to other disease types and used to integrate
83 data from a variety of omics technologies.

84 **2 Results**

85 **2.1 MB-SupCon: Microbiome-based prediction model via supervised contrastive** 86 **learning**

87 The main goal of MB-SupCon is to improve the prediction of phenotype or clinical covariates via
88 supervised contrastive learning. An overall workflow is shown in **Figure 1**. The model input includes
89 gut microbiome and metabolome data, phenotype information and/or clinical covariates. We then

90 train a supervised contrastive learning model to obtain the weights of the encoder networks. Finally,
91 we apply the predictive model to independent test datasets to assess its accuracy. The microbiome
92 embedding is critically useful for downstream analysis tasks, including 1) predicting phenotypic
93 outcomes and covariates and 2) visualizing the lower-dimensional representation. We show that
94 approaches using microbiome embedding from MB-SupCon often have better performance than
95 approaches using raw microbiome abundances.

96 **2.2 MB-SupCon improved categorical outcome prediction in type 2 diabetes study**

97 We trained MB-SupCon using real human gut microbiome and metabolome data obtained in a host-
98 microbe dynamics study by Zhou, Sailani [17]. The omics data were collected longitudinally from
99 subjects with prediabetes over approximately four years. Gut microbiome data were obtained from
100 stool samples, and host metabolome data was obtained from blood samples at each visit of subjects.
101 We subset both datasets and retained 720 samples with both 16s gut microbiome and metabolome
102 data. Microbiome data is encoded as a matrix of 720 x 96 dimension with entries having values
103 between $[0,1)$, (i.e., $[0, 1)^{720 \times 96}$), and each of the 96 features represents the relative abundance of one
104 microbial taxon from 5 taxonomic levels - phylum, class, order, family, and genus. Metabolome data
105 is encoded as a matrix of dimension 720 x 724, with each entry taking values from non-negative real
106 numbers, (i.e., $\mathbb{R}_+^{720 \times 724}$), and each of the 724 features represents the abundance of one metabolite.
107 Standardization was applied to both datasets before model fitting so that each feature has a mean
108 value of zero and unit variance. In addition, at each visit, demographic or clinical covariates (e.g., sex,
109 age, insulin resistant/insulin sensitive, BMI, etc.) were also recorded for all subjects. We also
110 attempted to predict the covariates using microbiome and metabolome data to evaluate different
111 predictive models. To evaluate the predictive performance for each machine learning model, we
112 applied 12 random splitting of training (70%), validation (15%), and testing (15%) to the data. For
113 each split, the training and validation sets were used for model fitting and hyperparameter tuning
114 (**Supplementary texts: Training and tuning procedure**), and the testing set was used for
115 benchmarking.

116 To illustrate the advantage of MB-SupCon, we used

- 117 • a logistic regression with elastic net regularization (EN),
- 118 • a multi-layer perceptron (MLP),
- 119 • a support vector machine classifier (SVM),
- 120 • a random forest classifier (RF)

121 to analyze and compare their performance on

- 122 • the original microbiome abundances,
- 123 • the embedding of supervised contrastive learning (MB-SupCon).

124 We also compared MB-SupCon with a method that uses a logistic regression model to analyze
125 unsupervised embeddings (MB-simCLR).

126 To distinguish analyses using original abundance and embeddings, we denote methods that analyze
127 embeddings with prefix “MB-SupCon” e.g., MB-SupCon+MLP represents using MLP to analyze
128 MB-SupCon embeddings.

129 We listed the details in **Supplementary texts: Calculation of the microbiome embedding** on
130 obtaining microbiome embeddings in unsupervised or supervised learning. To evaluate prediction
131 accuracy, we compute the fraction of correctly predicted labels for each model. Since we create

132 multiple splits of the data for training, validation, and testing, the average prediction accuracy using
133 different test folds are reported.

134 MB-SupCon embeddings, compared with the original data, lead to improved prediction accuracies in
135 logistic regression with an elastic net penalty, SVM, MB-simCLR. The methods using MB-SupCon
136 embedding almost always outperform RF and MLP models using raw microbiome abundance, which
137 are two of the most accurate methods (**Table 1, Figure 3**). For the prediction of insulin resistance,
138 methods using MB-SupCon embeddings achieved 84.62% average accuracy (MB-SupCon+Logistic,
139 MB-SupCon+SVM, MB-SupCon+RF, and MB-SupCon+MLP), which is better than methods that
140 uses raw abundances, i.e., the elastic net logistic regression (76.69%), SVM (79.46%), MB-simCLR
141 (65.67%), and similar to RF (83.93%) and MLP (83.73%). Similarly, for predicting sex, MB-SupCon
142 also has good average prediction accuracy (78.98%). For predicting race, a four-category outcome,
143 approaches using MB-SupCon embeddings reaches the lead average accuracy (80.04%), and their
144 advantage is consistent over the other methods, including RF (77.90%) and MLP (75.60%). More
145 importantly, MB-SupCon embedding leads to a near-best prediction accuracy regardless of the choice
146 of machine learning algorithms, which demonstrated its utility and robustness.

147 **2.3 MB-SupCon better visualized embeddings in independent datasets**

148 In addition to improving prediction accuracy, MB-SupCon embeddings in the lower dimensional
149 space can be useful for visualizations. In **Figure 2A**, we applied PCA on 1) raw abundance data, 2)
150 embeddings from MB-simCLR, and 3) embeddings from MB-SupCon in an independent test data.
151 We placed the samples of test datasets onto the principal component 2 (PC2) vs 1 (PC1) scatterplot
152 using a random seed of 1. In addition, we also compared MB-SupCon to three other methods, i.e.,
153 Sparse Partial Least Squares Discriminant Analysis (sPLS-DA) [18], Sparse Partial Least Squares
154 (sPLS)[19], and Data Integration Analysis for Biomarker discovery using Latent cOmponents
155 (DIABLO) [20], for their capability to distinguish different groups of covariates. sPLS-DA [18]
156 predicts covariates using microbiome data only; the other two methods are based on integrative
157 modeling of both microbiome and metabolome data. sPLS [19] uses microbiome data as predictors
158 and metabolome data as responses. DIABLO [20] uses multiple omics data from the same samples to
159 be blocks and covariate values to be the outcome. All three methods can be implemented under the
160 “mixOmics” [12] framework. In **Figure 2B**, we compared the lower-dimensional scatterplots
161 (Component 2 vs. Component 1) on the same testing data for each method to those of MB-SupCon
162 in **Figure 2A**. Only the embedding from MB-SupCon leads to separable clusters from distinct
163 covariate groups, whereas the other established methods failed to separate different categories of
164 covariates. This result confirms that the improvements in prediction accuracy of MB-SupCon can be
165 attributable to better feature embeddings.

166 **2.4 MB-SupCon analysis of an inflammatory bowel disease study**

167 To further evaluate the performance of MB-SupCon, we applied it to another independent multi-
168 omics Inflammatory Bowel Disease (IBD) study with both metagenomics and metabolomics data [9]
169 (detailed in **Supplementary texts: Network architecture and training of MB-SupCon model for
170 IBD study**). With “diagnosis” of IBD status as the covariate, we trained, validated and tested our
171 model using 12 different random splits similar to the diabetes study. For each model, we evaluated
172 the model performance on testing data. As shown in **Table 2** and **Figure 4**, the results remained
173 consistent with the T2D study. Approaches using MB-SupCon embeddings achieved significantly
174 better average prediction accuracies (74.04%) compared to approaches using original data directly,
175 including logistic regression (67.79%) and SVM (52.70%). When RF or MLP is used, predictions

176 based on MB-SupCon embedding was comparable to the predictions using original abundance
177 information, although MB-SupCon+RF had a slightly smaller variance compared to RF and has a
178 marginal advantage compared to MLP. This validated the reliability and extensive applicability of
179 MB-SupCon.

180

181 **3 Discussion**

182 A reliable microbiome-based prediction model could have immediate values in disease diagnosis and
183 treatment responses prediction [6, 21, 22]. Here, we propose a novel method, MB-SupCon, to improve
184 those models by utilizing increasingly accessible multi-omics datasets. The method leverages the
185 strengths of contrastive learning, which were first established in computer vision tasks [14-16, 23].
186 MB-SupCon performs the nonlinear transformation of microbiome abundance and produces useful
187 embeddings, which lead to improved prediction accuracies and more informative visual
188 representations. We demonstrate these advantages of MB-SupCon utilizing existing published data
189 from a diabetes study and an inflammatory bowel disease study. We showed that the improved
190 microbiome prediction model using MB-SupCon embeddings is more accurate than elastic net
191 logistic regression, support vector machine, and unsupervised contrastive learning model, and can
192 achieve comparable or better performance of random forest and multi-layer perceptron.

193 Like all other deep learning models, MB-SupCon has limitations. One drawback is that it does not
194 explicitly offer biological interpretations between the microbiome and metabolomics. This “black-
195 box” nature of the deep learning model often leads to criticisms. Developing more interpretable
196 machine learning models can potentially address the emerging biological questions. Another
197 limitation is that MB-SupCon does not explicitly model sample relatedness. Specifically, as paired
198 longitudinal data is relatively infrequent, MB-SupCon does not incorporate features that could
199 account for correlations among longitudinal samples. A better solution to explore in the future is to
200 change the current MLP encoders to mixed effect neural networks [24, 25] so that variation within
201 subjects for longitudinal data could be better modeled and explained.

202 There are numerous future applications and extensions of MB-SupCon. MB-SupCon is not restricted
203 to the microbiome and metabolomic data analysis. It can be applied to any omics technology (e.g.,
204 proteomics, host transcriptomics, host methylome, etc.). Moreover, MB-SupCon can be extended to
205 integrate more than two types of omics data. This can be achieved by adding pair-wise supervised
206 contrastive losses.

207 In summary, we believe MB-SupCon and encoder-based on the neural network in general have
208 advantage in approximating non-linear functions and modeling high-dimensional data. MB-SupCon
209 framework can be applicable in broad multi-omics settings and improves microbiome-based
210 prediction models.

211 **4 Methods**

212 Contrastive learning aims to maximize the similarities between microbiome embedding and
213 metabolome embedding from a pair of samples. Let X^g and X^m be the standardized microbiome and
214 metabolome data. Suppose there are n samples in a minibatch. For a single sample i ($i = 1, 2, \dots, n$),
215 we denote the associated microbial and metabolic data as x_i^g and x_i^m , respectively. Let the
216 microbiome (or metabolome) encoder network be a multi-layer perceptron $f^g(\cdot)$ (or $f^m(\cdot)$). The

217 encoded features (embeddings) of microbiome and metabolome for sample i are $z_i^g = f^g(x_i^g)$ and
 218 $z_i^m = f^m(x_i^m)$, respectively. We define the similarity between the encoded vectors z_i^g, z_j^m for $i, j \in$
 219 $S = \{1, 2, \dots, n\}$ in the latent space by the cosine similarity,

$$220 \quad \text{sim}(z_i^g, z_j^m) = \frac{z_i^g \cdot z_j^m}{\|z_i^g\| \|z_j^m\|} \text{ for } i, j \in S = \{1, 2, \dots, n\}$$

221 where \cdot denotes the dot product of two vectors and $\| \cdot \|$ denotes the Euclidean norm of a vector.

222 We first introduce MB-SimCLR, an unsupervised contrastive learning approach: if a pair of
 223 microbiome and metabolome samples are from the same sample, we define the corresponding data
 224 $\{x_i^g, x_i^m\}$ as a “positive pair”. Otherwise, we define the pair of data $\{x_i^g, x_j^m\}$ ($i \neq j$) as a “negative
 225 pair”. Given n pairs of microbiome and metabolome samples, if we set the embedding vector of
 226 microbiome z_i^g as an anchor, the loss of unsupervised contrastive learning is

$$227 \quad \text{Loss}_{\text{unsup}}^{g,m} = -E_{i \in S} \left[\log \frac{\exp\{\text{sim}(z_i^g, z_i^m)/\tau\}}{\sum_{j=1}^n \exp\{\text{sim}(z_i^g, z_j^m)/\tau\}} \right]$$

228 where $i \in S = \{1, 2, \dots, n\}$, $\tau \in \mathbb{R}_+$ is the temperature parameter.

229 Symmetrically, by anchoring the embedding of the metabolome we can get loss $\text{Loss}_{\text{unsup}}^{m,g}$. The total
 230 loss will be the sum of these two parts: $\text{Loss}_{\text{unsup}} = \text{Loss}_{\text{unsup}}^{g,m} + \text{Loss}_{\text{unsup}}^{m,g}$

231 Improved upon MB-SimCLR, we describe a supervised contrastive learning method, MB-SupCon,
 232 where we incorporate labels in calculating the loss function. Given a specific categorical label y_i from
 233 sample i , $P(y_i)$ denotes the index set of samples with label y_i . Any pairs of microbiome and
 234 metabolome vectors $\{x_k^g, x_l^m\}$ with $k, l \in P(y_i)$ are treated as “positive pairs”. Otherwise, they are
 235 “negative pairs”. Suppose we set microbiome embedding z_i^g for $i \in S$ with label y_i as an anchor.
 236 Then supervised contrastive loss [15] is defined as

$$237 \quad \text{Loss}^{g,m} = -E_{i \in S} \left[\frac{1}{|P(y_i)|} \sum_{l \in P(y_i)} \log \frac{\exp\{\text{sim}(z_i^g, z_l^m)/\tau\}}{\sum_{j=1}^n \exp\{\text{sim}(z_i^g, z_j^m)/\tau\}} \right]$$

238 where $|P(y_i)|$ is the cardinality of index set $P(y_i)$, $\tau \in \mathbb{R}_+$ is the temperature parameter.

239 By anchoring metabolome embedding, we can get $\text{Loss}^{m,g}$. The total loss is still the sum of $\text{Loss}^{g,m}$
 240 and $\text{Loss}^{m,g}$.

241 In all, the difference between supervised contrastive learning and unsupervised contrastive learning
 242 is the definition of positive and negative sample pairs. Once the loss is determined, we can update the
 243 weights of encoder networks using the stochastic gradient descent (SGD) method. Embedding can be
 244 calculated as the network outputs. Details are provided in the **Supplemental Texts: Network**
 245 **Architecture and Training**.

246

247 **Figure and table legend**

248

249 **Figure 1. Overview of the MB-SupCon framework.** Step 1 - Data Collection: Microbiome,
250 metabolome, phenotype/covariates are collected; Step 2 – Contrastive Learning – MB-SupCon is
251 applied, and two encoder networks are trained; Step 3 – Predictive Model – microbiome encoder
252 network can be applied to new microbiome data to obtain microbiome embeddings. The embeddings
253 lead to an improved microbiome-based prediction model and lower-dimensional representation.

254

255 **Figure 2. Scatter plots of test data on lower-dimensional space (T2D study).**

256 **Panel A: Scatter plots of test data (random seed 1) on a 2-dimensional space by PCA.** 1st row:
257 the first two principal components for the embeddings learned from MB-SupCon; 2nd row: the first
258 two principal components for the original data; 3rd row: the first two principal components for the
259 embeddings learned from MB-simCLR.

260 Acronyms: PCA - Principal component analysis.

261 **Panel B: Scatter plots of test data (random seed 1) on 2-dimensional space by other methods.**

262 1st row: the first two components learned from sPLSDA on original data; 2nd row: the first two
263 components learned from sPLS on original data; 3rd row: the first two principal components
264 learned from DIABLO on original data.

265 Acronyms: PCA - Principal component analysis. sPLS-DA - Sparse Partial Least Squares
266 Discriminant Analysis; sPLS - Sparse Partial Least Squares; DIABLO - Data Integration Analysis
267 for Biomarker discovery using Latent cOmponents.

268

269 **Figure 3. Scatter plot of average prediction accuracies on test data from 12 random training-
270 validation-testing splits, by using different methods for categorical covariates (T2D study).**

271 Green triangles and red points represent predictions based on MB-SupCon embeddings. Orange
272 squares and blue points represent predictions based on original microbiome data. **Panel A:** Insulin
273 resistant/sensitive; **Panel B:** Sex; **Panel C:** Race.

274 Acronyms: LOGISTIC - logistic regression with elastic net penalty; SVM - support vector machine
275 classifier; RF - random forest classifier; MLP - multi-layer perceptron.

276

277 **Figure 4. Scatter plots of average prediction accuracies for diagnosis on testing data from 12
278 random training-validation-testing splits, by using different methods for categorical
279 covariates (IBD study).** Green triangles and red points represent predictions based on MB-SupCon
280 embeddings. Orange squares and blue points represent predictions based on original microbiome
281 data.

282 Acronyms: LOGISTIC - logistic regression with elastic net penalty; SVM - support vector machine
283 classifier; RF - random forest classifier; MLP - multi-layer perceptron.

284 **Table 1. Average prediction accuracies on testing data from 12 random training-validation-**
285 **testing splits, by using different methods for categorical covariates (T2D study).**

286 Acronyms: Logistic - logistic regression with elastic net penalty using original data; SVM - support
287 vector machine classifier using original data; RF - random forest classifier using original data; MLP
288 multi-layer perceptron using original data; MB-simCLR - logistic regression model with elastic net
289 penalty using microbiome embeddings learned from unsupervised contrastive learning; MB-
290 SupCon + Logistic - logistic regression model with elastic net penalty using microbiome
291 embeddings learned from supervised contrastive learning. MB-SupCon + SVM: support vector
292 machine classifier using microbiome embeddings learned from supervised contrastive learning;
293 MB-SupCon + RF: random forest classifier using microbiome embeddings learned from supervised
294 contrastive learning; MB-SupCon + MLP: multi-layer perceptron using microbiome embeddings
295 learned from supervised contrastive learning; Avg. Acc. based on MB-SupCon: average accuracies
296 among MB-SupCon + Logistic, MB-SupCon + SVM, MB-SupCon + RF and MB-SupCon + MLP.

297 **Table 2. Average prediction accuracies on testing data from 12 random training-validation-**
298 **testing splits, by using different methods for categorical covariates (IBD study).** Acronyms are
299 defined the same as those from Table 1.

300 **Supplementary Figure 1. Structure of the microbiome and metabolome encoder network.** Only
301 dense layers are visualized where numbers represent neuron counts. Batch normalized layer,
302 Activation layer, Dropout layer are appended after each dense layer but not shown.

303 **Supplementary Figure 2. Hyperparameter tuning for MB-SupCon on T2D study.** Panel A1 –
304 A3: hyperparameter tuning result for covariate Insulin resistant/sensitive by logistic regression with
305 an elastic net penalty, SVM, and RF; Panel B1 - B3: hyperparameter tuning results for covariate
306 Sex by logistic regression with an elastic net penalty, SVM, and RF; Panel C1 – C3:
307 hyperparameter tuning result for covariate Race by logistic regression with an elastic net penalty,
308 SVM and RF.

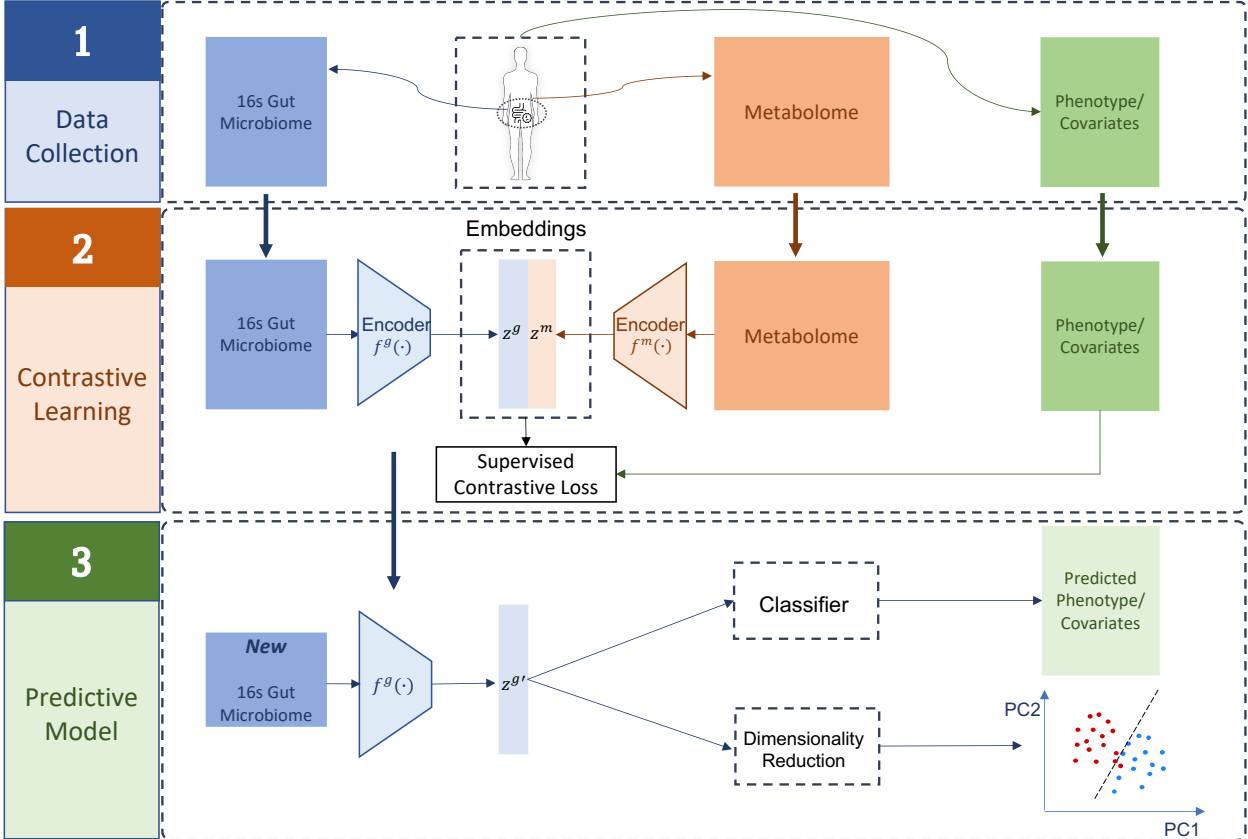
309 Acronyms: Dropout Rate: dropout rate of the encoders from MB-SupCon; weight_decay: weight
310 decay value (l2 regularization) of the stochastic gradient descent (SGD) optimizer; temperature:
311 temperature hyperparameter when calculating contrastive losses.

312

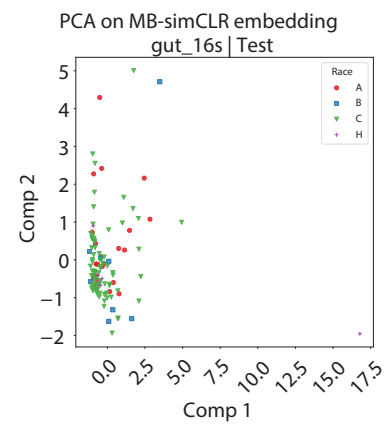
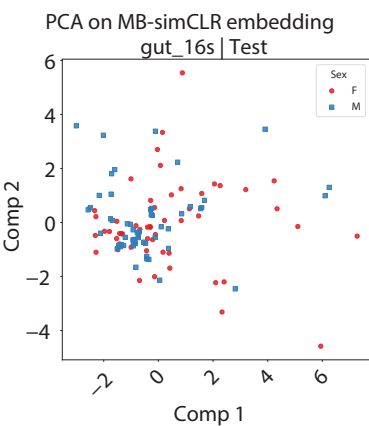
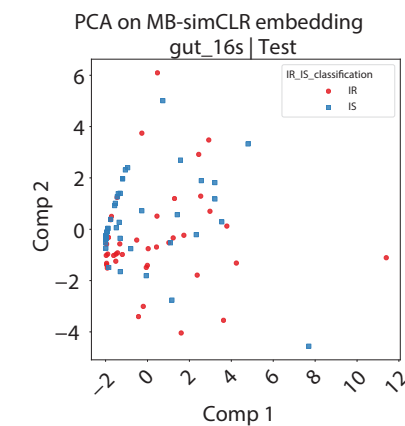
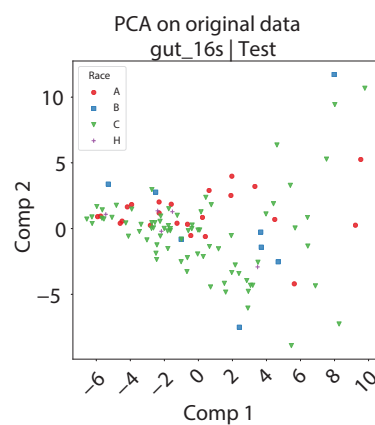
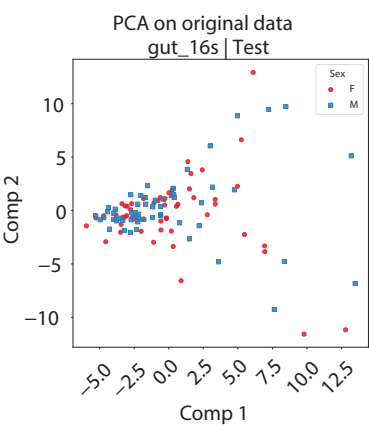
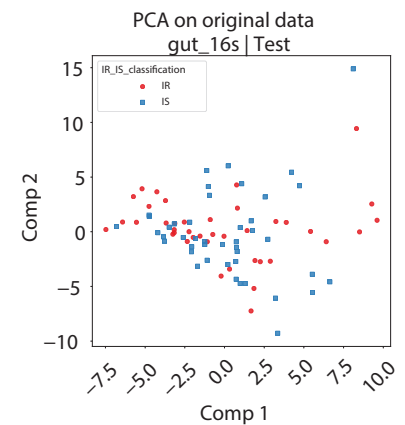
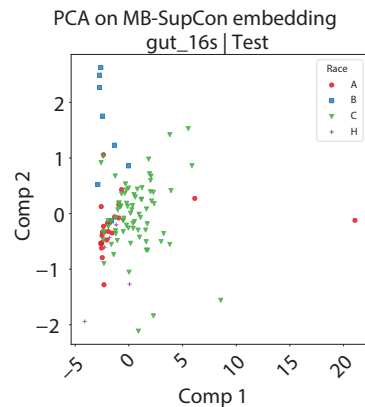
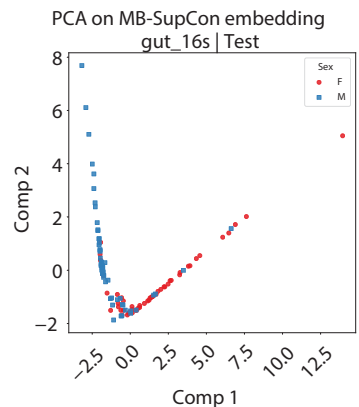
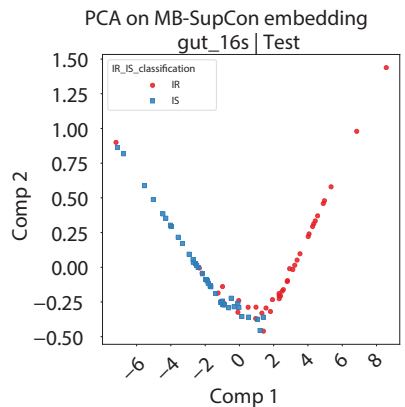
313 Reference

- 314 [1] Human Microbiome Project C. A framework for human microbiome research. *Nature*.
315 2012;486:215-21.
316 [2] Sender R, Fuchs S, Milo R. Revised Estimates for the Number of Human and Bacteria Cells in
317 the Body. *PLoS Biol*. 2016;14:e1002533.
318 [3] Cho I, Blaser MJ. The human microbiome: at the interface of health and disease. *Nature reviews*
319 *Genetics*. 2012;13:260-70.
320 [4] Claesson MJ, Clooney AG, O'Toole PW. A clinician's guide to microbiome analysis. *Nat Rev*
321 *Gastroenterol Hepatol*. 2017;14:585-95.
322 [5] Pasolli E, Truong DT, Malik F, Waldron L, Segata N. Machine Learning Meta-analysis of Large
323 Metagenomic Datasets: Tools and Biological Insights. *PLoS Comput Biol*. 2016;12:e1004977.
324 [6] Yu J, Feng Q, Wong SH, Zhang D, Liang QY, Qin Y, et al. Metagenomic analysis of faecal
325 microbiome as a tool towards targeted non-invasive biomarkers for colorectal cancer. *Gut*.
326 2017;66:70-8.

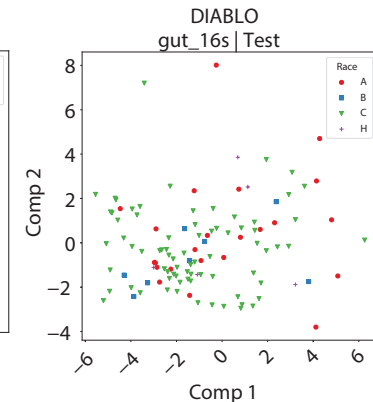
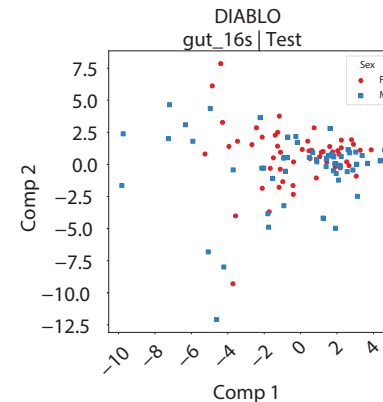
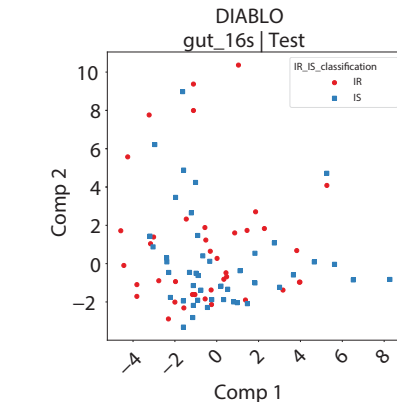
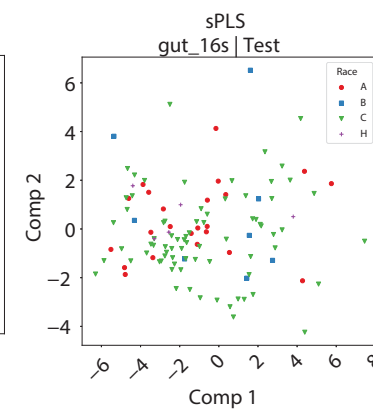
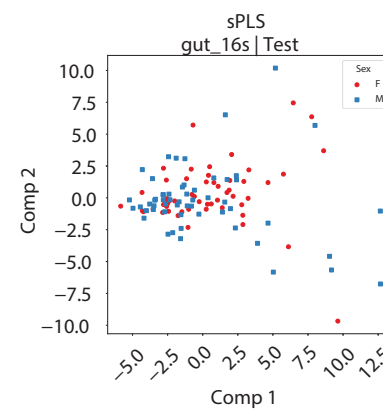
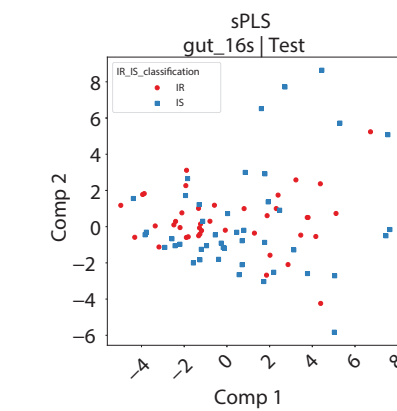
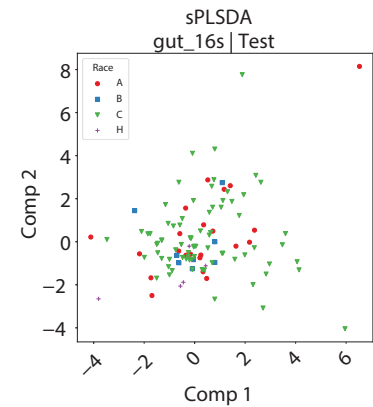
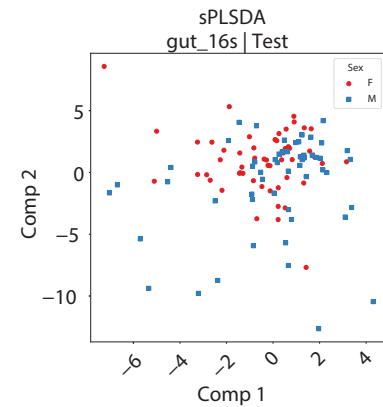
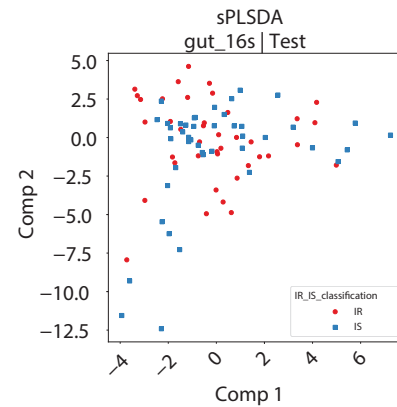
- 327 [7] Chen F, Dai X, Zhou CC, Li KX, Zhang YJ, Lou XY, et al. Integrated analysis of the faecal
328 metagenome and serum metabolome reveals the role of gut microbiome-associated metabolites in
329 the detection of colorectal cancer and adenoma. *Gut*. 2021.
- 330 [8] Wang Q, Ye J, Fang D, Lv L, Wu W, Shi D, et al. Multi-omic profiling reveals associations
331 between the gut mucosal microbiome, the metabolome, and host DNA methylation associated gene
332 expression in patients with colorectal cancer. *BMC Microbiol*. 2020;20:83.
- 333 [9] Lloyd-Price J, Arze C, Ananthakrishnan AN, Schirmer M, Avila-Pacheco J, Poon TW, et al.
334 Multi-omics of the gut microbial ecosystem in inflammatory bowel diseases. *Nature*. 2019;569:655-
335 62.
- 336 [10] Heintz-Buschart A, May P, Laczny CC, Lebrun LA, Bellora C, Krishna A, et al. Integrated
337 multi-omics of the human gut microbiome in a case study of familial type 1 diabetes. *Nat*
338 *Microbiol*. 2016;2:16180.
- 339 [11] Friedman J, Hastie T, Tibshirani R. *The elements of statistical learning: Springer series in*
340 *statistics Springer, Berlin; 2001.*
- 341 [12] Rohart F, Gautier B, Singh A, Lê Cao K-A. mixOmics: An R package for ‘omics feature
342 selection and multiple data integration. *PLOS Computational Biology*. 2017;13:e1005752.
- 343 [13] Tian Y, Krishnan D, Isola P. Contrastive multiview coding. *European conference on computer*
344 *vision: Springer; 2020. p. 776-94.*
- 345 [14] Chen T, Kornblith S, Norouzi M, Hinton G. A simple framework for contrastive learning of
346 visual representations. *International conference on machine learning: PMLR; 2020. p. 1597-607.*
- 347 [15] Khosla P, Teterwak P, Wang C, Sarna A, Tian Y, Isola P, et al. Supervised contrastive
348 learning. *Advances in Neural Information Processing Systems*. 2020;33:18661--73.
- 349 [16] Tian Y, Krishnan D, Isola P. Contrastive Multiview Coding. *European conference on computer*
350 *vision. 2020:776--94.*
- 351 [17] Zhou W, Sailani MR, Contrepois K, Zhou Y, Ahadi S, Leopold SR, et al. Longitudinal multi-
352 omics of host-microbe dynamics in prediabetes. *Nature*. 2019;569:663-71.
- 353 [18] Le Cao KA, Boitard S, Besse P. Sparse PLS discriminant analysis: biologically relevant feature
354 selection and graphical displays for multiclass problems. *BMC Bioinformatics*. 2011;12:253.
- 355 [19] Lê Cao K-A, Rossouw D, Robert-Granié C, Besse P. A sparse PLS for variable selection when
356 integrating omics data. *Statistical applications in genetics and molecular biology*. 2008;7.
- 357 [20] Singh A, Shannon CP, Gautier B, Rohart F, Vacher M, Tebbutt SJ, et al. DIABLO: an
358 integrative approach for identifying key molecular drivers from multi-omics assays. *Bioinformatics*.
359 2019;35:3055-62.
- 360 [21] Andrews MC, Duong CPM, Gopalakrishnan V, Iebba V, Chen WS, Derosa L, et al. Gut
361 microbiota signatures are associated with toxicity to combined CTLA-4 and PD-1 blockade. *Nat*
362 *Med*. 2021;27:1432-41.
- 363 [22] Zeller G, Tap J, Voigt AY, Sunagawa S, Kultima JR, Costea PI, et al. Potential of fecal
364 microbiota for early-stage detection of colorectal cancer. *Mol Syst Biol*. 2014;10:766.
- 365 [23] Wu Z, Xiong Y, Yu S, Lin D. Unsupervised Feature Learning via Non-Parametric Instance-
366 level Discrimination. *Proceedings of the IEEE conference on computer vision and pattern*
367 *recognition. 2018:3733--42.*
- 368 [24] Xiong Y, Kim HJ, Singh V. Mixed effects neural networks (menets) with applications to gaze
369 estimation. *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern*
370 *Recognition2019. p. 7743-52.*
- 371 [25] Tandon R, Adak S, Kaye JA. Neural networks for longitudinal studies in Alzheimer’s disease.
372 *Artificial intelligence in medicine. 2006;36:245-55.*

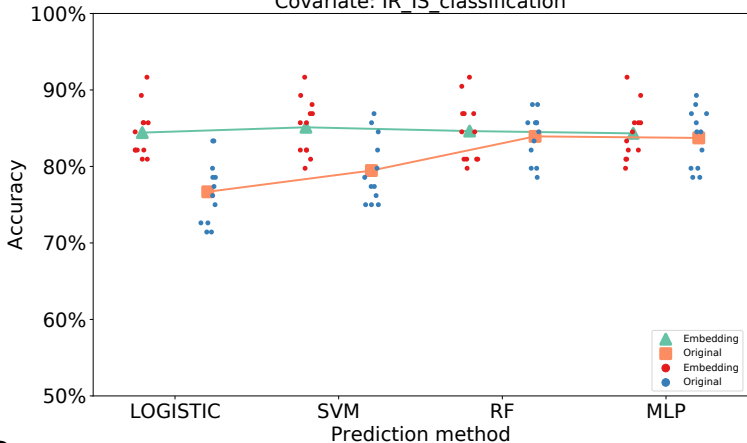
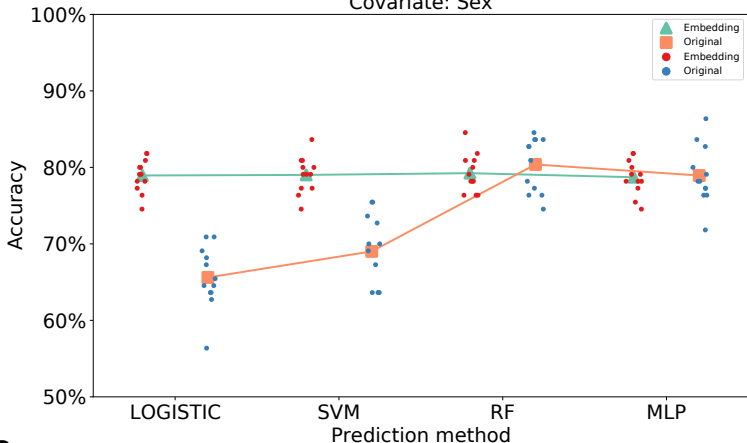
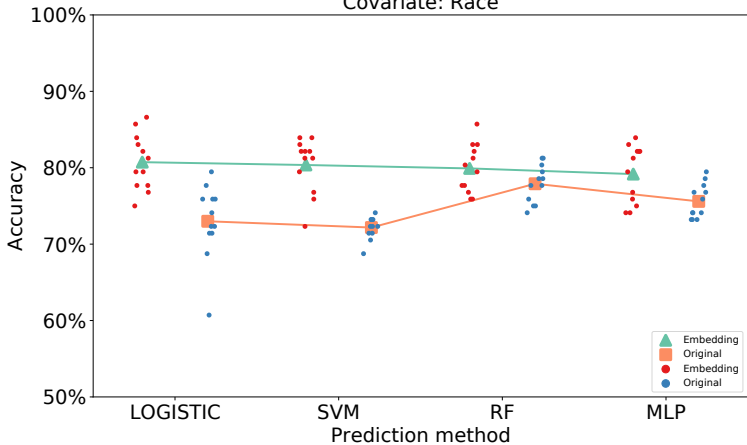


A)



B)



AT2D Study - Prediction accuracy on testing data
Covariate: IR_IS_classification**B**T2D Study - Prediction accuracy on testing data
Covariate: Sex**C**T2D Study - Prediction accuracy on testing data
Covariate: Race

IBD Study - Prediction accuracy on testing data
Covariate: diagnosis

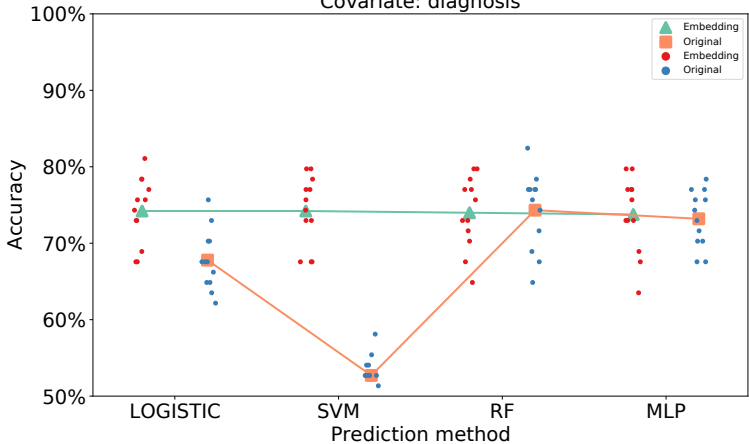


Table 1. Average prediction accuracies on testing data from 12 random training-validation-testing splits, by using different methods for categorical covariates (T2D study).

Acronyms: Logistic - logistic regression with elastic net penalty using original data; SVM - support vector machine classifier using original data; RF - random forest classifier using original data; MLP - multi-layer perceptron using original data; MB-simCLR - logistic regression model with elastic net penalty using microbiome embeddings learned from unsupervised contrastive learning; MB-SupCon + Logistic - logistic regression model with elastic net penalty using microbiome embeddings learned from supervised contrastive learning. MB-SupCon + SVM: support vector machine classifier using microbiome embeddings learned from supervised contrastive learning; MB-SupCon + RF: random forest classifier using microbiome embeddings learned from supervised contrastive learning; MB-SupCon + MLP: multi-layer perceptron using microbiome embeddings learned from supervised contrastive learning; Avg. Acc. based on MB-SupCon: average accuracies among MB-SupCon + Logistic, MB-SupCon + SVM, MB-SupCon + RF and MB-SupCon + MLP.

Prediction Task	Logistic	SVM	RF	MLP	MB-simCLR
Insulin resistance	76.69%	79.46%	83.93%	83.73%	65.67%
Sex	65.61%	69.02%	80.38%	78.94%	59.85%
Race	72.99%	72.17%	77.90%	75.60%	68.38%

Prediction Task	MB-SupCon + Logistic	MB-SupCon + SVM	MB-SupCon + RF	MB-SupCon + MLP	Avg. Acc. based on MB-SupCon
Insulin resistance	84.42%	85.12%	84.62%	84.33%	84.62%
Sex	78.94%	79.02%	79.24%	78.71%	78.98%
Race	80.73%	80.36%	79.91%	79.17%	80.04%

Table 2. Average prediction accuracies on testing data from 12 random training-validation-testing splits, by using different methods for categorical covariates (IBD study). Acronyms are defined the same as those from Table 1.

Prediction Task	Logistic	SVM	RF	MLP
diagnosis	67.79%	52.70%	74.32%	73.20%

Prediction Task	MB-SupCon + Logistic	MB-SupCon + SVM	MB-SupCon + RF	MB-SupCon + MLP	Avg. Acc. based on MB-SupCon
diagnosis	74.21%	74.21%	73.99%	73.76%	74.04%