

# Tree Variable Selection for Paired Case-Control Studies with Application to Microbiome Data

Min Lu and Hemant Ishwaran

**Abstract** When case-control studies involve paired samples, tree analyses based on traditional splitting rules are suboptimal as they ignore the paired nature of the data. Paired samples occur in microbiome studies when they are collected from different locations of the same individual or when they are collected from paired individuals with familial ties. Borrowing concepts from tree splitting, we propose a novel approach that accommodates the paired structure in the data for fast and effective nonparametric variable ranking. Importantly this method allows detangling of different types of associations at play with structured correlated outcomes such as host genotype and environmental exposure effects. Another technique for variable selection are variable importance measures. We describe two types of measures useful for paired data analysis. The methodology is illustrated on the microbiota of paired samples from a case-control study of obesity.

## 1 Introduction

Paired samples occur in microbiome studies when they are collected from different locations of the same individual or from paired individuals with familial ties. Human microbiome can be shared among family members with variations in each individual's microbial community [4, 1]. Suppose an identifiable "core microbiome" exists at the microbial gene level and deviations from this core are associated with different physiologic states. It is of interest to study how family ties play a role in these deviations. For example, if deviations from a core gut microbiome are associated with body mass index (BMI), we can define "individual" and "family" outcomes

---

Min Lu

Division of Biostatistics, University of Miami, e-mail: luminwin@gmail.com

Hemant Ishwaran

Division of Biostatistics, University of Miami e-mail: hemant.ishwaran@gmail.com

with labels obese/lean, where for example obese family means the individual comes from a family containing at least one member who is obese, and lean family means the individual comes from a family whose members are all lean. By studying such outcomes, we can examine how each array of microbial genes is associated with obesity both at the family and individual level.

To illustrate our proposed methodology, we will use data from a cross-sectional study focusing on obesity in twins [4, 2, 6]. Data was collected from human stools of monozygotic or dizygotic twins or their mothers. We utilize 142 of these samples. The bacterial lineages present in the fecal microbiotas of these individuals were characterized by rRNA sequencing. Sequences were identified by assignment to taxonomic outcome groups using operational taxonomic units (OTUs). Specific details of how data was processed can be found in [4].

The original study found that obesity is associated with phylum-level changes in the microbiota and reduced bacterial diversity using linear approaches, such as PCA (Principal Components Analysis). Here we will focus on detecting which taxonomic outcome groups are the most informative for obesity risk at both the family and individual level using a novel approach that draws upon tree based concepts.

## 2 Gini Index

Consider a multiclass problem where  $Y$  is a categorical (factor) outcome such that  $Y \in \{1, \dots, J\}$  for  $J \geq 2$ . We call this the  $J$ -class problem and call  $\{1, \dots, J\}$  the  $J$  class labels for  $Y$ . Classification tree splitting is often based on the Gini index splitting rule. If  $\mathbf{p} = (p_1, \dots, p_J)$  are the data class proportions of  $Y$  for classes 1 through  $J$ , respectively, the Gini index of impurity is defined as

$$\phi(\mathbf{p}) = \sum_{j=1}^J p_j(1 - p_j) = 1 - \sum_{j=1}^J p_j^2.$$

Classification trees are grown using the Gini index by splitting features recursively into left and right daughter nodes, where tree splits are obtained by minimizing tree impurity. The Gini index split-statistic for a split  $s$  on a continuous feature  $x_m$  at a given tree node is

$$\theta(Y, x_m, s) = \frac{n_l}{n} \phi(\mathbf{p}_l) + \frac{n_r}{n} \phi(\mathbf{p}_r),$$

where the subscript  $l = \{x_m \leq s\}$  and  $r = \{x_m > s\}$  denote the left and right daughter nodes formed by the split on  $x_m$  at  $s$  ( $n_l$  and  $n_r$  are the sample sizes of the two daughter nodes where  $n = n_l + n_r$  is the parent sample size). To reduce tree impurity, the goal is to find  $x_m$  and  $s$  to *minimize*

$$\theta(Y, x_m, s) = \frac{n_l}{n} \left( 1 - \sum_{j=1}^J \frac{n_{j,l}^2}{n_l^2} \right) + \frac{n_r}{n} \left( 1 - \sum_{j=1}^J \frac{n_{j,r}^2}{n_r^2} \right),$$

where  $n_{j,l}$  and  $n_{j,r}$  are the number of cases of class  $j$  in the left and right daughters, respectively and  $n_j = n_{j,l} + n_{j,r}$  are the number of cases of class  $j$  and  $n = \sum_{j=1}^J n_j$ . With some algebra, it can be shown this is equivalent to *maximizing* the split-statistic

$$g(Y, x_m, s) = \frac{1}{n} \sum_{j=1}^J \frac{n_{j,l}^2}{n_l} + \frac{1}{n} \sum_{j=1}^J \frac{(n_j - n_{j,l})^2}{n - n_l}.$$

Although the Gini index is primarily used as a splitting rule for growing a classification tree, we note that it can be used as a fast preliminary variable ranking method. For each of the  $p$  predictors  $x_1, \dots, x_p$ , define

$$G(Y, x_m) = g(Y, x_m, s_{\max}),$$

where

$$s_{\max} = \arg \max_s g(Y, x_m, s).$$

For this analysis,  $g(Y, x_m, s)$  is taken to be the split statistic for the root node consisting of the entire data—thus  $n$  is the full sample size. Variables can be ranked in order of importance by the size of  $G(Y, x_m)$ . Notice this variable selection procedure is fully nonparametric and can be computed quickly even in big data settings. The following section provides a demonstration of how this approach works for our problem.

## 2.1 Simulation Study

Consider a binary class setting and denote the outcome as  $Y^I \in \{0, 1\}$ , where  $Y^I = 0$  represents a lean individual and  $Y^I = 1$  an obese individual. Family outcome is denoted as  $Y^F \in \{0, 1\}$ , where  $Y^F = 0$  signifies an individual from a family with all lean members and  $Y^F = 1$  indicates an individual from a family where at least one member is obese. Association with  $Y^I = 1$  reflects how host adiposity influences the gut microbiome, whereas association with  $Y^F = 1$  reflects environmental exposure influences. How the host genotype affects the gut microbiome under environmental exposure is reflected by an association with both  $Y^I = 1$  and  $Y^F = 1$ .

We use the following simulation where  $Y^F$  is specified according to

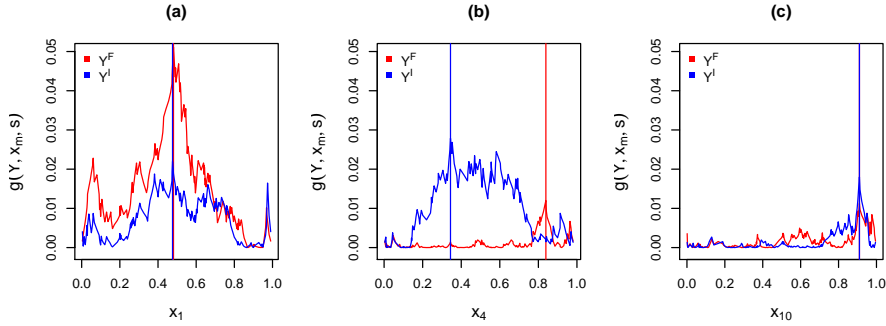
$$\mathbb{P}\{Y^F = 1 | \mathbf{X} = \mathbf{x}\} = \text{logistic}(-2 + x_1 + x_2 + x_3 + 2 \times \mathbf{1}_{\{x_1 < 0.5\}}) \quad (1)$$

and  $Y^I$  is specified by

$$\mathbb{P}\{Y^I = 1 | Y^F = 1, \mathbf{X} = \mathbf{x}\} = \text{logistic}(-2 + x_4 + x_5 + x_6 + 2 \times \mathbf{1}_{\{x_4 < 0.5\}}), \quad (2)$$

where  $\text{logistic}(\alpha) = 1/(1 + e^{-\alpha})$ . In this scenario,  $x_1, x_2$  and  $x_3$  are associated with environmental exposures that causes the presence of obesity. While,  $x_4, x_5$  and  $x_6$  are associated with host adiposity, given that the host is under these types of environmental exposures.

The feature space dimension was set to  $p = 10$ . Features were independently drawn from a uniform distribution  $U(0, 1)$ . Variables unrelated to outcome, representing noise variables, were also added to the design matrix. For  $Y^F$ , noise variables were  $x_4, \dots, x_{10}$ . For  $Y^I$ , noise variables were  $x_1, x_2, x_3$  and  $x_7, \dots, x_{10}$ . Split-statistics,  $g(Y, x_m, s)$ , are plotted in Figure 1 for features  $x_1, x_4$  and  $x_{10}$  and for both outcomes  $Y = Y^F$  and  $Y = Y^I$ . Red color represents the family level outcome  $Y^F$  and blue is used for the individual level outcome  $Y^I$ . Variable  $x_1$  in (a) predicts obesity at the family level, and is associated with  $Y^F$ , and the true optimal split point occurs at 0.5. We can see that the split-statistic of  $x_1$  is high for both  $Y^F$  and  $Y^I$  and both peak at around 0.5. Variable  $x_4$  in (b) is associated with  $\mathbb{P}\{Y^I = 1 | Y^F = 1\}$ , and therefore is associated with  $Y^I$ , and has a true optimal split point of 0.5. We can see that the split-statistic  $g(Y^I, x_4, s)$  is high for  $Y^I$  and reaches its peak near 0.5 (although not exactly at the true value—we will come back to this point later). In contrast, the split-statistic  $g(Y^F, x_4, s)$  for  $Y^F$  does not at all have an optimized value near 0.5 and its peak value occurs near its edge. This edge effect is typical of noisy variables and is a property of the Gini splitting rule called end-cut preference, ECP [13]. Variable  $x_{10}$  in (c) is a noise variable, and its split-statistic is low for both  $Y^F$  and  $Y^I$ . Observe that its optimal split points is close to the edge for both outcomes, which as stated is typical behavior of a noisy variable.



**Fig. 1** Univariate split-statistics for  $x_1, x_4$  and  $x_{10}$  from simulation (1)-(2). Values  $g(Y, x_m, s)$  are shown across different split values  $s$ . Red and blue display family level outcome  $Y^F$  and individual level outcomes  $Y^I$  respectively. Vertical lines mark the optimal split-statistic  $G(Y, x_m)$ . Variable  $x_1$  is associated with  $Y^F$  with true optimal split point of 0.5. Variable  $x_4$  is associated with  $\mathbb{P}\{Y^I = 1 | Y^F = 1\}$  with true optimal split point of 0.5. Variable  $x_{10}$  is a noise variable.

Comparing the results across Figure 1, it is clear that  $G(Y, x_m)$ , which is the highest point of  $g(Y, x_m, s)$ , is useful for variable ranking. However, focusing only on family level outcomes (red color) will ignore features like  $x_4$  that are related to the individual level outcome (blue color). Checking both split-statistics clearly helps better understand the underlying associations.

### 3 Multivariate Gini Index

Tang and Ishwaran [11] defined a multivariate Gini index split-statistic obtained by averaging univariate Gini split-statistics. For the bivariate outcome problem, this can be described as

$$\bar{g}_u(Y^F, Y^I, x_m, s) = \frac{1}{2} [g(Y^F, x_m, s) + g(Y^I, x_m, s)].$$

The subscript “ $u$ ” is used to emphasize that the split-statistic is unweighted. We can define

$$\bar{G}_u(Y^F, Y^I, x_m) = \bar{g}_u(Y^F, Y^I, x_m, s_{u_{\max}})$$

for ranking variables, where

$$s_{u_{\max}} = \arg \max_s \bar{g}_u(Y^F, Y^I, x_m, s).$$

Larger values of  $\bar{G}_u(Y^F, Y^I, x_m)$  identify informative variables and smaller values indicate noise variables.

#### 3.1 Conditional Gini Index

The problem with the split-statistic  $\bar{g}_u(Y^F, Y^I, x_m, s)$  is that by averaging across the outcomes it ignores the correlation between  $Y^F$  and  $Y^I$ . To resolve this issue, we introduce the following conditional Gini split-statistic.

Let  $\pi_c = \mathbb{P}\{Y^I = 1 | Y^F = 1\}$  be the population proportion of obese cases among individuals with at least one obese family member. The subscript “ $c$ ” is used to emphasize this is a conditional probability. Because there are only two classes, we have  $\mathbf{p}_c = (p_c, 1 - p_c)$  and  $\phi(\mathbf{p}_c) = 2p_c(1 - p_c)$  where  $p_c$  is the sample estimator of  $\pi_c$ . For a split  $s$  on variable  $x_m$ , the conditional Gini split-statistic is defined as

$$\theta_c(Y^F, Y^I, x_m, s) = \frac{\tilde{n}_l}{\tilde{n}} \phi(\mathbf{p}_c) + \frac{\tilde{n}_r}{\tilde{n}} \phi(\mathbf{p}_c),$$

where as before subscripts  $l$  and  $r$  denote left and right daughter nodes formed by the split. The numbers of cases  $Y^F = 1$  in the daughters are  $\tilde{n}_l$  and  $\tilde{n}_r$  where  $\tilde{n} = \tilde{n}_l + \tilde{n}_r$ . The numbers of these cases where  $Y^I = 1$  in the left and right daughters is denoted by  $\tilde{n}_{1,l}$  and  $\tilde{n}_{1,r}$  respectively. It can be shown that minimizing  $\theta_c(Y^F, Y^I, x_m, s)$  is equivalent to maximizing

$$g_c(Y^F, Y^I, x_m, s) = \frac{\tilde{n}_{1,l}^2}{\tilde{n}\tilde{n}_l} + \frac{\tilde{n}_{1,r}^2}{\tilde{n}\tilde{n}_r}.$$

We can define

$$G_c(Y^F, Y^I, x_m) = g_c(Y^F, Y^I, x_m, s_{c_{\max}})$$

for ranking variables, where  $s_{c_{\max}} = \arg \max_s g_c$ .

Now because  $g_c(Y^F, Y^I, x_m, s)$  conditions on  $Y^F = 1$ , it is not designed to identify signal affecting  $Y^F$ . To resolve this, define the conditional weighted split-statistic

$$\bar{g}_{cw}(Y^F, Y^I, x_m, s) = \frac{1}{w_F + w_I} [w_F \cdot g(Y^F, x_m, s) + w_I \cdot g_c(Y^F, Y^I, x_m, s)]$$

for detecting features that affect both  $Y^F$  and  $Y^I$ . Observe that when  $w_F = w_I = 1$ , this becomes an unweighted split-statistic and will be denoted by  $\bar{g}_{cu}(Y^F, Y^I, x_m, s)$ .

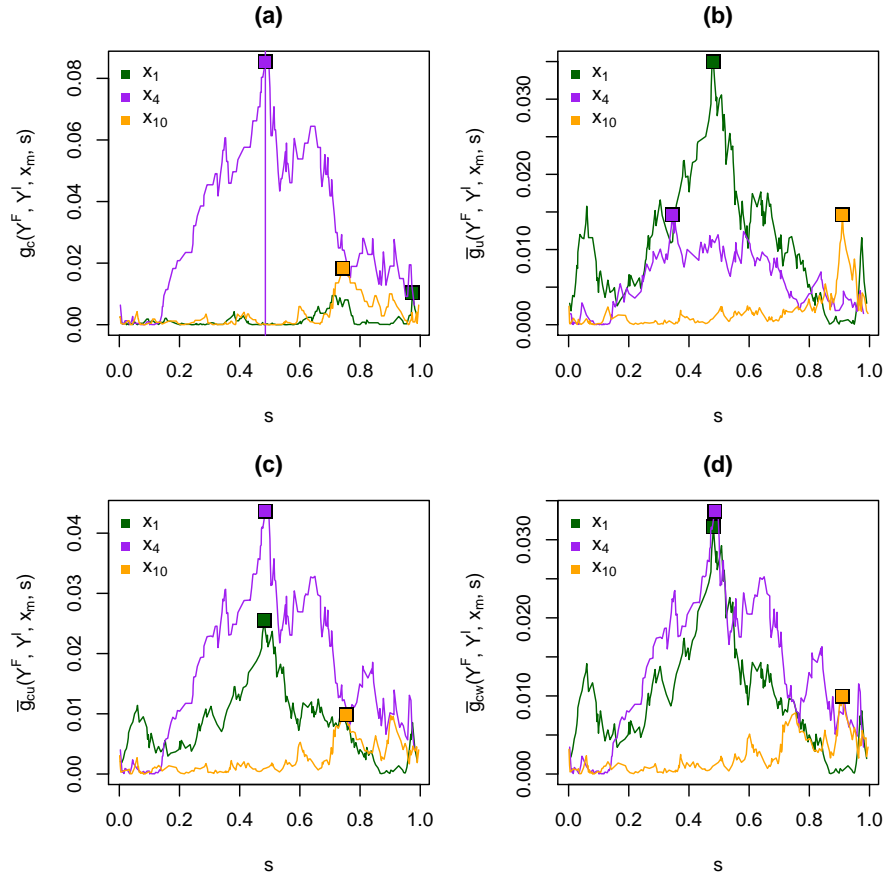
Weighted indices can be calculated as  $w_F = \sum_i^n \mathbf{1}_{\{Y_i^F=1\}}$  and  $w_I = \sum_i^n \mathbf{1}_{\{Y_i^I=1\}}$ , which adjusts for the fact that there are always more obese cases for  $Y^F$  than  $Y^I$ . The maximum value for the conditional weighted split-statistic is

$$\bar{G}_{cw}(Y^F, Y^I, x_m) = \bar{g}_{cw}(Y^F, Y^I, x_m, s_{cw_{\max}})$$

where  $s_{cw_{\max}} = \arg \max_s \bar{g}_{cw}$ . In a likewise fashion, define the maximum conditional unweighted split-statistic  $\bar{G}_{cu}(Y^F, Y^I, x_m)$ .

Figure 2 displays: (a)  $g_c(Y^F, Y^I, x_m, s)$ , (b)  $\bar{g}_u(Y^F, Y^I, x_m, s)$ , (c)  $\bar{g}_{cu}(Y^F, Y^I, x_m, s)$  and (d)  $\bar{g}_{cw}(Y^F, Y^I, x_m, s)$  for variables  $x_1, x_4$  and  $x_{10}$  from the simulation (1)-(2). Variable  $x_4$  affects the conditional probability  $\mathbb{P}(Y^I = 1 | Y^F = 1)$ , which is plotted in purple color. Returning to the point made earlier regarding Figure 1(b), when comparing Figure 2(a) to Figure 1(b), we find  $g_c(Y^F, Y^I, x_4, s)$  characterizes  $x_4$  better than  $g(Y^I, x_4, s)$  as the maximum value is closer to the true splitting point 0.5. Another point to observe is that the goal of  $\bar{g}_u(Y^F, Y^I, x_m, s)$  and  $\bar{g}_{cu}(Y^F, Y^I, x_m, s)$  is to detect features associated with  $Y^F$  and/or  $Y^I$ . However,  $\bar{g}_u(Y^F, Y^I, x_m, s)$  in (b) is less effective than  $\bar{g}_{cu}(Y^F, Y^I, x_m, s)$  in (c) because it ranks  $x_4$  similarly to noise variable  $x_{10}$  (shown in orange). In contrast,  $\bar{g}_{cu}(Y^F, Y^I, x_m, s)$  in (c) and the weighted  $\bar{g}_{cw}(Y^F, Y^I, x_m, s)$  in (d) properly rank  $x_4$  as more informative than  $x_{10}$ . In fact, the weighted split-statistic tends to do an even better job.

Figure 3 displays maximum Gini split-statistics for all  $p = 10$  variables averaged over 100 independent replications. For convenient calibration, the averaged split-statistic for the noise variable  $x_7$  is used as a selection cutoff. When comparing subfigure (c) with (b), we see that  $G_c(Y^F, Y^I, x_m)$  performs better in term of selecting the true signals,  $x_4, x_5$  and  $x_6$ , than  $G(Y^I, x_m)$ . When comparing subfigure (f) with (d), we observe that the weighted Gini split-statistic utilizing the conditional Gini index,  $\bar{G}_{cw}(Y^F, Y^I, x_m)$ , outperforms the simple averaged Gini split-statistic,  $\bar{G}_u(Y^F, Y^I, x_m)$ , in selecting the true signal variables  $x_1, \dots, x_6$  (in (d) the informative variable  $x_6$  is not selected whereas the noise variable  $x_{10}$  is selected). The performance of  $\bar{G}_{cu}(Y^F, Y^I, x_m)$  and  $\bar{G}_{cw}(Y^F, Y^I, x_m)$  are roughly similar except that noise variable  $x_{10}$  is less likely to be chosen using  $\bar{G}_{cw}(Y^F, Y^I, x_m)$ . Thus as before, the weighted split-statistic tends to do a better job. Finally, when comparing subpanel (f) to (a) notice that  $\bar{G}_{cw}(Y^F, Y^I, x_m)$  is as good as  $G(Y^F, x_m)$  in identifying variables  $x_1, x_2, x_3$  related to  $Y^F$ . However, this does not mean  $G(Y^F, x_m)$  is not useful, since when combined with  $\bar{G}_{cw}(Y^F, Y^I, x_m)$  it allows one to detangle variable relationships with the two outcomes.

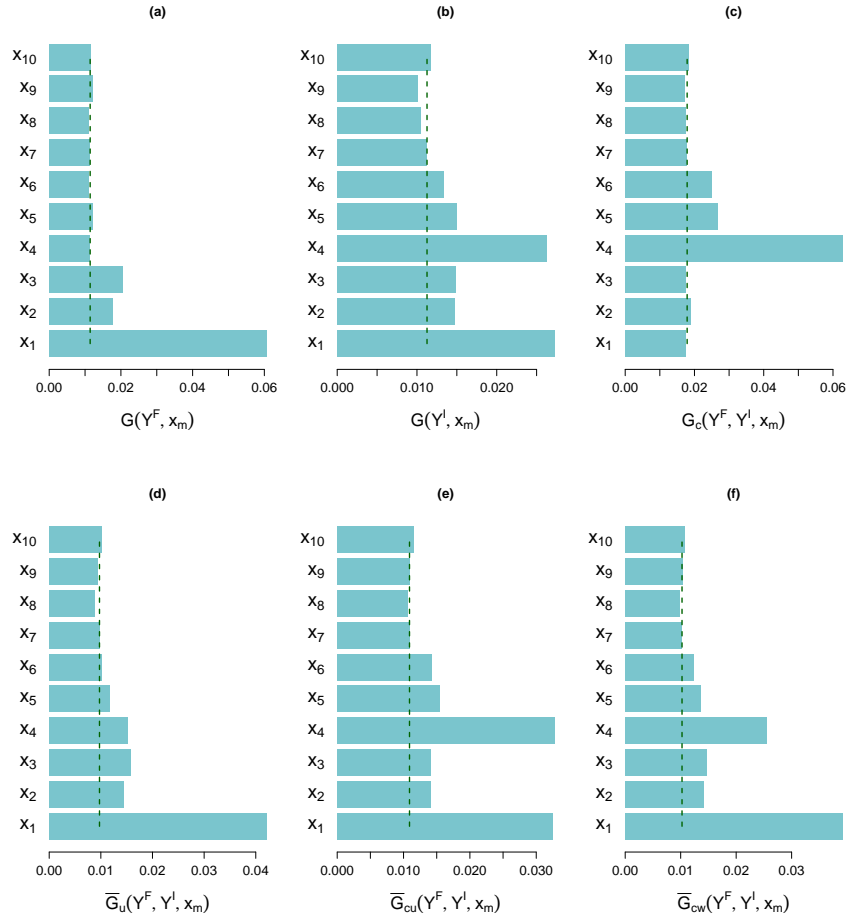


**Fig. 2** Multivariate split-statistics for  $x_1$ ,  $x_4$  and  $x_{10}$  from simulation (1)-(2). Curves displayed are: (a)  $g_c(Y^F, Y^I, x_m, s)$ , (b)  $\bar{g}_u(Y^F, Y^I, x_m, s)$ , (c)  $\bar{g}_{cu}(Y^F, Y^I, x_m, s)$  and (d)  $\bar{g}_{cw}(Y^F, Y^I, x_m, s)$  with maximum statistic marked by a square point.

## 4 Variable Importance

Another effective tool for variable selection is variable importance (VIMP). The permutation VIMP for a variable  $x_m$  is the prediction error for the model subtracted from the prediction error for the model using data that randomly permutes  $x_m$  [14]. This procedure can be implemented over independent bootstrap samples and the value averaged to obtain a more stable estimator [14]. More formally, let  $\hat{P}\hat{E}(Y)$  be the averaged out-of-sample (called out-of-bag and abbreviated as OOB) misclassification error for the original model. Let  $\hat{P}\hat{E}(Y, x_m^*)$  be the averaged OOB misclassification error when  $x_m$  is randomly permuted. The VIMP for  $x_m$  is

$$I(Y, x_m) = \hat{P}\hat{E}(Y, x_m^*) - \hat{P}\hat{E}(Y).$$



**Fig. 3** Variable ranking from maximum split-statistics for simulation (1)-(2) repeated 100 times independently. Dashed line is averaged value of maximum Gini split-statistic for noise variable  $x_7$  which represents a convenient cutoff value.

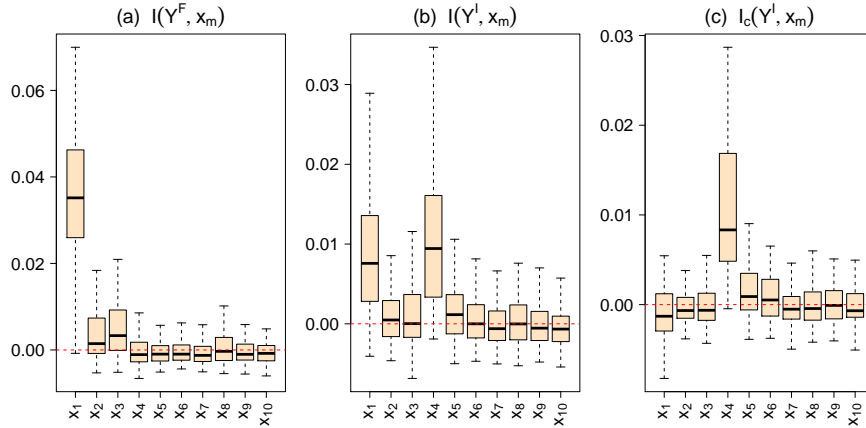
To determine if variables affect the conditional probability  $\mathbb{P}(Y^I = 1 | Y^F = 1)$ , we define a conditional VIMP analogous to the conditional Gini index. Conditional VIMP is calculated by restricting the data to those cases where  $Y^F = 1$ . The conditional VIMP index for  $x_m$  is

$$I_c(Y^I, x_m) = \hat{\mathbb{P}}\hat{\mathbb{E}}_c(Y^I, x_m^*) - \hat{\mathbb{P}}\hat{\mathbb{E}}_c(Y^I).$$

Figure 4 displays VIMP for all  $p$  features for our simulation. Values have been averaged over 100 independent replications. Unconditional VIMP,  $I(Y^F, x_m)$ , for  $Y^F$  displayed in subpanel (a) successfully ranks the true signal variables  $x_1, x_2$



and  $x_3$  as the most informative. When comparing subpanel (c) to (b), we see that conditional VIMP,  $I_c(Y^I, x_m)$ , is better at selecting true signal variables  $x_4, x_5$  and  $x_6$  than unconditional VIMP,  $I(Y^I, x_m)$ . In subfigure (b), VIMP for  $x_1$  is very large and would lead to incorrect selection compared with (c).



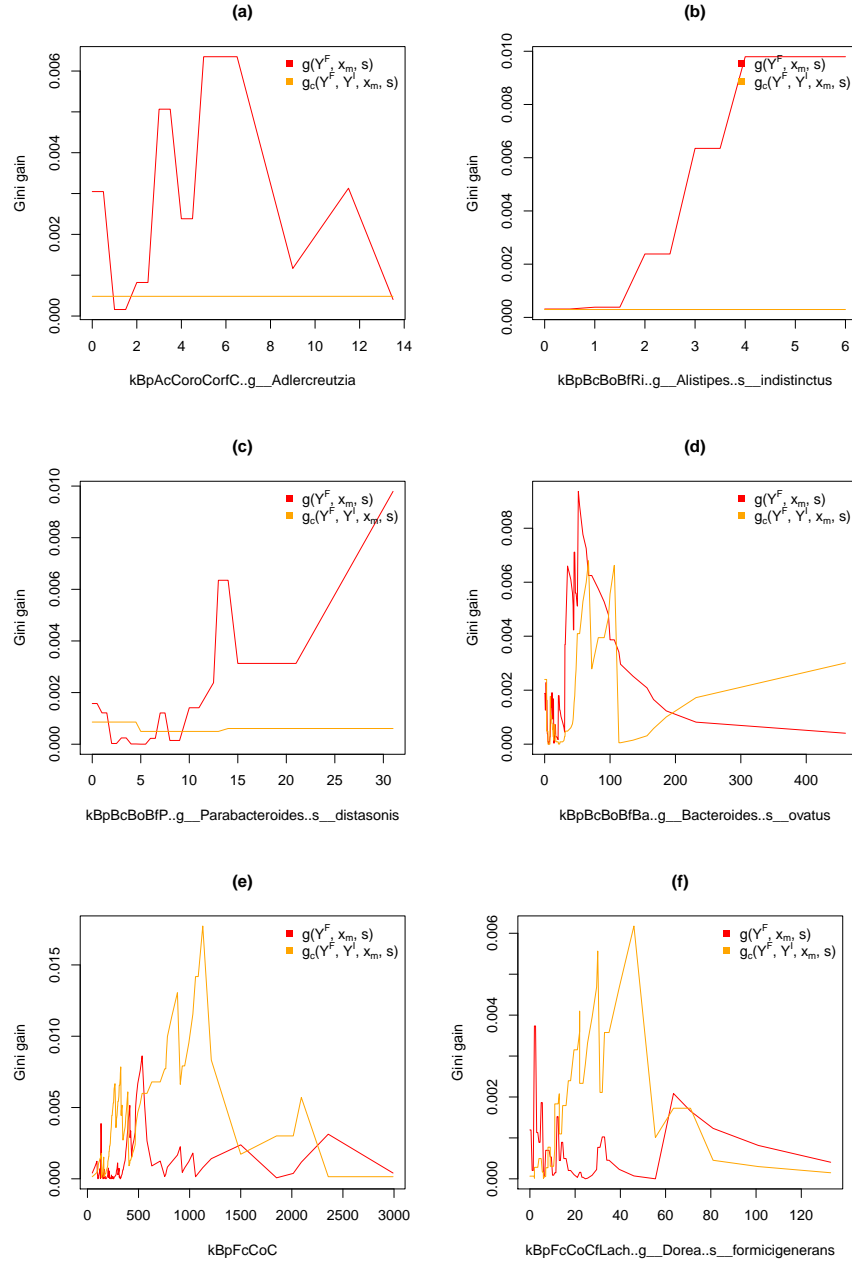
**Fig. 4** Variable importance from simulation (1)-(2) averaged over 100 independent replications.

## 5 Analysis of Obesity using Microbiome Data

Now we return to the microbiome obesity data described earlier ( $n = 142$  and  $p = 174$ ). Outcomes were coded as before:  $Y^I = 0$  represents a lean individual,  $Y^I = 1$  an obese individual,  $Y^F = 0$  signifies an individual from a family with all lean members, and  $Y^F = 1$  indicates an individual from a family where at least one member is obese. Table 1 of the Appendix provides convenient abbreviated names for features.

Figure 5 displays split-statistics for 6 representative features, chosen to illustrate how host and environmental factors affect the gut microbiome. Univariate split-statistics  $g(Y^F, x_m, s)$  for the family outcome  $Y^F$  are shown using red lines, and conditional split-statistics  $g_c(Y^F, Y^I, x_m, s)$  are displayed using orange lines. Bivariate split-statistics,  $\bar{g}_{cu}(Y^F, Y^I, x_m, s)$  and  $\bar{g}_{cw}(Y^F, Y^I, x_m, s)$ , lie between these two lines. Recall when optimal split points appear towards the edge of a features' range that this is a sign of a noisy feature (referred to as the ECP property [13]).

Subfigures (a), (b), (c) represent features informative for the environmental outcome  $Y^F$ . In all three figures  $g(Y^F, x_m, s)$  takes large values across the range of feature values. However, these three features are not informative for  $\mathbb{P}\{Y^I = 1|Y^F = 1\}$  as  $g_c(Y^F, Y^I, x_m, s)$  is near zero in all instances. Thus they do not reflect how host adiposity influences the gut microbiome under environmental exposure.



**Fig. 5** Split-statistics for microbiome obesity data. Shown are 6 representative variables illustrating how taxonomic outcome groups predict obesity risk at the family level (shown using the univariate Gini split-statistic on  $Y^F$ ,  $g(Y^F, x_m, c)$ , plotted in red) and at the individual level (shown using the univariate conditional Gini split-statistic  $g_c(Y^F, Y^I, x_m, c)$ , plotted in orange). Variable names are abbreviated according to Table 1.

Subfigures (d) and (e) represent features that are informative for both  $Y^F = 1$  and  $\mathbb{P}\{Y^I = 1|Y^F = 1\}$  as both  $g(Y^F, x_m, s)$  and  $g_c(Y^F, Y^I, x_m, s)$  assume relatively large values. These features identify influences from both environmental exposure and host adiposity. For (d), the two maximum split-statistics are nearly the same which suggests that effect of environmental exposure and host adiposity are roughly the same for this feature. For (e),  $g_c(Y^F, Y^I, x_m, s)$  attains a much larger maximum statistic than  $g(Y^F, x_m, s)$  at a higher feature value. This suggest the effect of environmental exposure and host adiposity depends on the feature value, for example whether the feature value is larger than 500 or 1000.

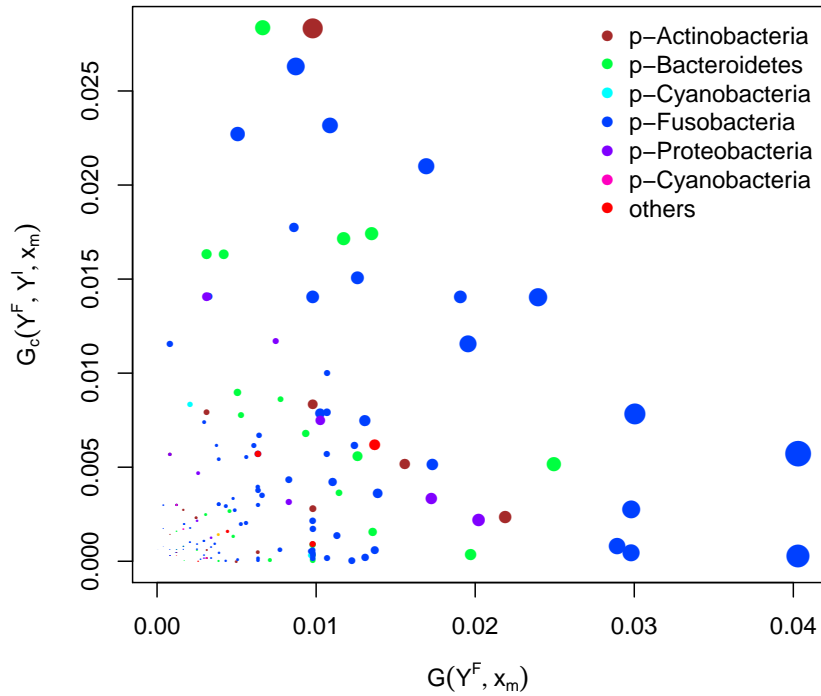
Subfigure (f) is a feature that mainly reflects the influence from host adiposity, rather than environmental exposure. This is because values of  $g(Y^F, x_m, s)$  are overall small and its optimal split point is close to the edge of its range, signaling that it is likely a noisy variable for  $Y^F$ .

Values of  $G(Y^F, x_m)$  and  $G_c(Y^F, Y^I, x_m)$  are given Figure 6. Size of circles are scaled proportional to  $\overline{G}_{cw}(Y^F, Y^I, x_m)$ . Phylum groups are used to color circles. It is interesting to note that features informative for  $Y^F = 1$  and  $\mathbb{P}\{Y^I = 1|Y^F = 1\}$  belong primarily to the Fusobacteria phylum. Generally values of  $G_c(Y^F, Y^I, x_m)$  are smaller than  $G(Y^F, x_m)$ . However when they are weighted to obtain  $\overline{G}_{cw}(Y^F, Y^I, x_m)$ , we can see that there is a nice balancing of values.

Finally, Figure 7 displays unconditional VIMP,  $I(Y^F, x_m)$ , and conditional VIMP,  $I_c(Y^I, x_m)$ , for all  $p$  variables. Many variables have small or negative values thus showing that VIMP can be used as an effective means to dimension.

## 6 Discussion

Fast nonparametric selection of features that accounts for correlation in paired data is a valuable tool for microbiome data analysis. Variable selection procedures can choose features that reflect influences from external effects (between pairs) and internal effects (within pairs), but without taking in account the paired structure of the data, they will be inefficient in separating the two types of effects. Our proposed conditional Gini split-statistic, when used alone or averaged with univariate Gini split-statistics, serves two purposes. First, the maximum value of the split-statistic can be used for variable ranking and variable selection. Conditional Gini is able to select variables reflecting how the microbiome is affected by host adiposity given the same environmental exposures. Second, how the value of the split-statistic varies within a feature provides useful insight into the magnitude of the external and/or internal effects. The optimal split point for conditional Gini represents the threshold that a feature can separate lean and obese individuals given the same environmental exposure. We demonstrated these two aspects in a systematic comparative simulation and through a real data application. We found that the paired structure of the data played a very strong role in performance of our methods. Without controlling for family level of obesity, features that only affect individual level of obesity are often noticeably masked.



**Fig. 6** Comparison between  $G(Y^F, x_m)$  and  $G_{Y^F=1}(Y^I, x_m)$ . The size of circles are proportional to  $G_{c_w}(Y^F, Y^I, x_m)$ . The color of circles identifies the phylum.

There are other variable selection procedures designed for multivariate outcomes. However, in big data settings, computational speed plays a key role. Practically speaking, the best method is not always optimal for the researcher because computational times can be too long. Our Gini split-statistics can be rapidly computed for a large number of features in big data settings and because the calculations are univariate the procedure could be parallelized to further reduce runtimes. Users can simulate a noise feature to determine the cutoff for screening noise variables. Potentially, our Gini indices can be used as tree splitting rules so that all the features can be taken into consideration together. Moreover, our approach could leverage powerful machine learning methods such as random forests and boosting to provide a direct approach to analyze paired data.

**Acknowledgements** This work was supported by the National Institutes of Health [grant numbers R01 CA200987 and R01 HL141892 to H.I.].

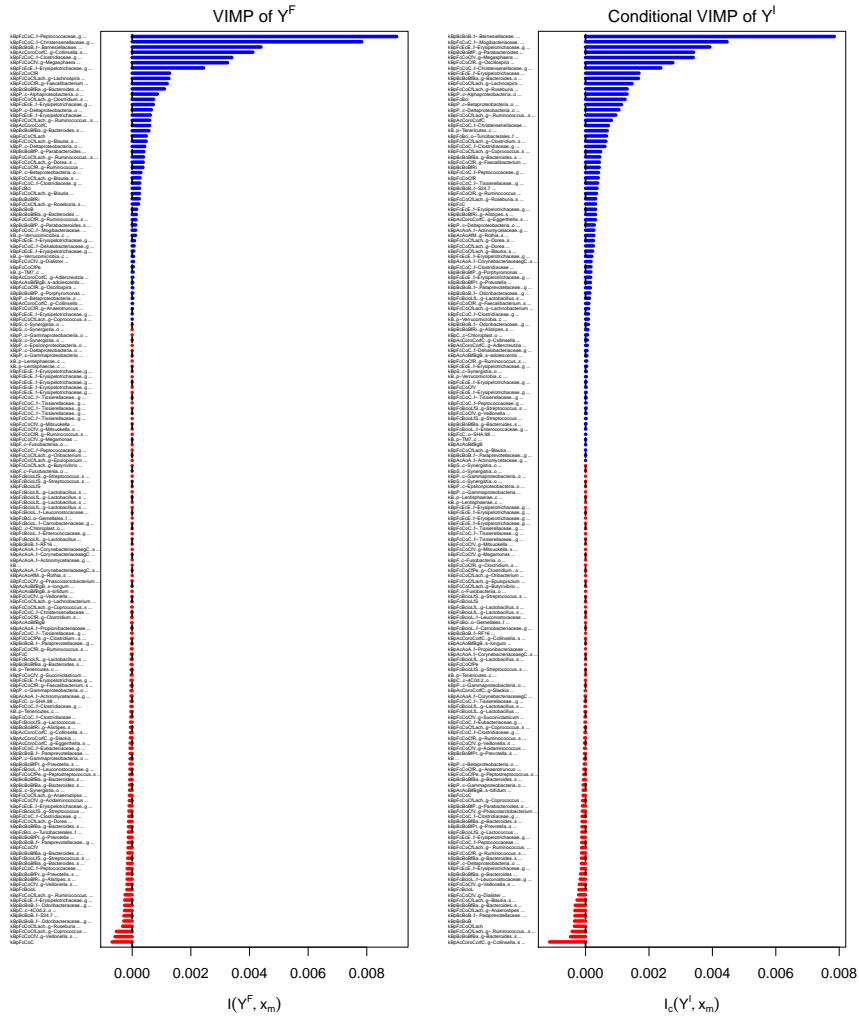
## References

1. Xia, Y., Sun, J., Chen, D. G.: Statistical analysis of microbiome data with R. Springer, Singapore (2018).
2. Knights-lab: Monozygotic or dizygotic twin pairs concordant for BMI class, and their mothers (Project Turnbaugh 2009).  
<https://github.com/knights-lab/MLRepo/tree/master/datasets/turnbaugh>. Cited 28 Mar 2020
3. Geddes, K.O., Czapor, S.R., Labahn, G.: Algorithms for Computer Algebra. Kluwer, Boston (1992)
4. Turnbaugh, P. J., Hamady, M., Yatsunenko, T., Cantarel, B. L., Duncan, A., Ley, R. E., ... Egholm, M.: A core gut microbiome in obese and lean twins. *Nature* **457**(7228), 480–484 (2009)
5. Tyler, A. D., Smith, M. I., Silverberg, M. S.: Analyzing the human microbiome: a “how to guide for physicians. *American Journal of Gastroenterology*, **109**(7), 983-993 (2014)
6. Vangay, P., Hillmann, B. M., Knights, D.: Microbiome Learning Repo (ML Repo): A public repository of microbiome regression and classification tasks. *GigaScience*, **8**(5), giz042. (2019)
7. Al-Ghalith, G., Knights, D.: BURST enables optimal exhaustive DNA alignment for big data. (2017) doi. org/10.5281/zenodo, 806850.
8. O’Leary, N. A., Wright, M. W., Brister, J. R., Ciufu, S., Haddad, D., McVeigh, R., ... Astashyn, A.: Reference sequence (RefSeq) database at NCBI: current status, taxonomic expansion, and functional annotation. *Nucleic acids research*, **44**(D1), D733–D745 (2016)
9. Breiman, L.: Random forests. *Machine Learning*. **45** 5–32 (2001)
10. Breiman, L.: Bagging predictors. *Machine Learning*. **24**(2), 123–140 (1996)
11. Tang, F., Ishwaran, H. Random forest missing data algorithms. *Statistical Analysis and Data Mining: The ASA Data Science Journal*, **10**(6), 363–377 (2017).
12. McDonald D., Price M.N., Goodrich J., et al: An improved Greengenes taxonomy with explicit ranks for ecological and evolutionary analyses of bacteria and archaea. *ISME J* **6**:610–8 (2012).
13. Ishwaran H. The effect of splitting on random forests. *Machine Learning*, **99**: 75-118 (2015).
14. Ishwaran, H. and Lu, M.: Standard errors and confidence intervals for variable importance in random forest regression, classification, and survival. *Statistics in Medicine*, **38**: 558-582, 2019.

## Appendix

**Table 1** Abbreviated feature names for microbiome obesity data.

Abbrev.	Full Form	Abbrev.	Full Form
kB	k-Bacteria	oBi	..o-Bifidobacteriales
pF	..p-Firmicutes	oE	..o-Erysipelotrichales
pA	..p-Actinobacteria	oL	..o-Lactobacillales
pB	..p-Bacteroidetes	fB	..f-Bifidobacteriaceae
pC	..p-Cyanobacteria	fBa	..f-Bacteroidaceae
pF	..p-Fusobacteria	fC	..f-Coriobacteriaceae
pP	..p-Proteobacteria	fL	..f-Lactobacillaceae
pS	..p-Synergistetes	fLach	..f-Lachnospiraceae
cA	..c-Actinobacteria	fM	..f-Micrococcaceae
cB	..c-Bacteroidia	fP	..f-Porphyrromonadaceae
cBci	..c-Bacilli	fPe	..f-Peptostreptococcaceae
cC	..c-Clostridia	fPr	..f-Prevotellaceae
cCor	..c-Coriobacteriia	fR	..f-Ruminococcaceae
cE	..c-Erysipelotrichi	fRi	..f-Rikenellaceae
oA	..o-Actinomycetales	fS	..f-Streptococcaceae
oC	..o-Clostridiales	fV	..f-Veillonellaceae
oCor	..o-Coriobacteriales	gB	..g-Bifidobacterium
oB	..o-Bacteroidales	gC	..g-Corynebacterium



**Fig. 7** Variable ranking using VIMP for microbiome data of obesity. Variables with higher value on the left reflect how the gut microbiome is influenced by environmental factors. Variables with higher values in the right reflect how gut microbiome is affected by host adiposity given the environmental exposures.