

Supplementary Material:
Feature selection and causal analysis for microbiome studies in the presence
of confounding using standardization

Emily Goren

May 2, 2021

Contents

S1 Additional simulation results	2
S2 Real data analysis: correlation structure and assumption checks	9

List of Tables

S1 Simulation results: TPR and FPR, $n = 100$, Poisson features	2
----------------------------------------------------------------------------	---

List of Figures

S1 Simulation results: AUC, $n = 50$, Poisson features	3
S2 Simulation results: AUC, $n = 50$, negative binomial features	4
S3 Simulation results: AUC, $n = 100$, negative binomial features	5
S4 Simulation results: FDP, $n = 50$, Poisson features	6
S5 Simulation results: FDP, $n = 50$, negative binomial features	7
S6 Simulation results: FDP, $n = 100$, negative binomial features	8
S7 Sorghum data: pairwise OTU correlations	9
S8 Sorghum data: conditional residuals vs. predicted	9
S9 Sorghum data: conditional residual Q-Q plot	10

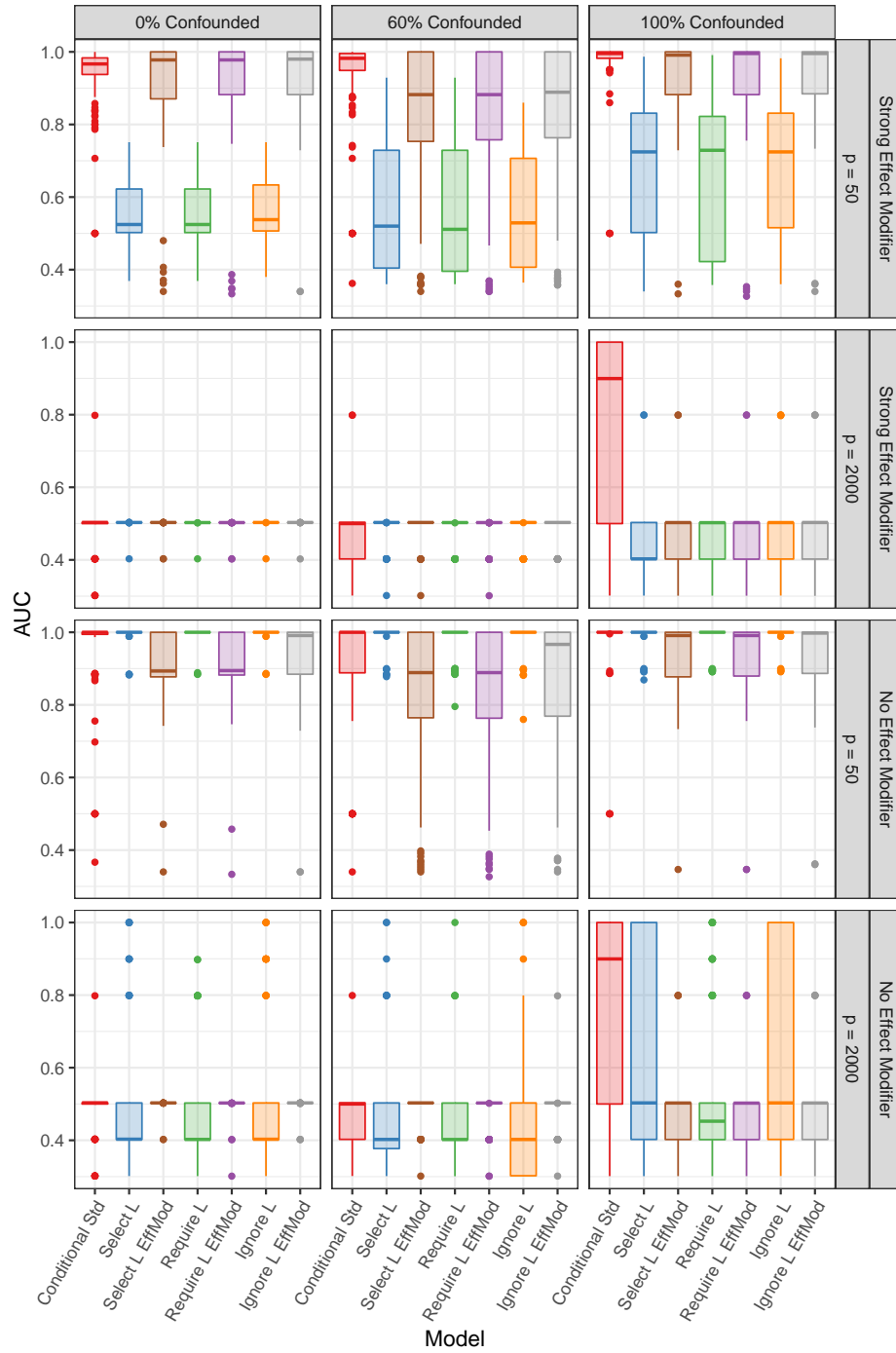


Figure S1: Simulation results: box plots of the area under the curve (AUC) from 100 simulation replications for $n = 50$ and Poisson features using p -values based on the debiased LASSO estimate following iterative sure independence screening (iterative SIS).

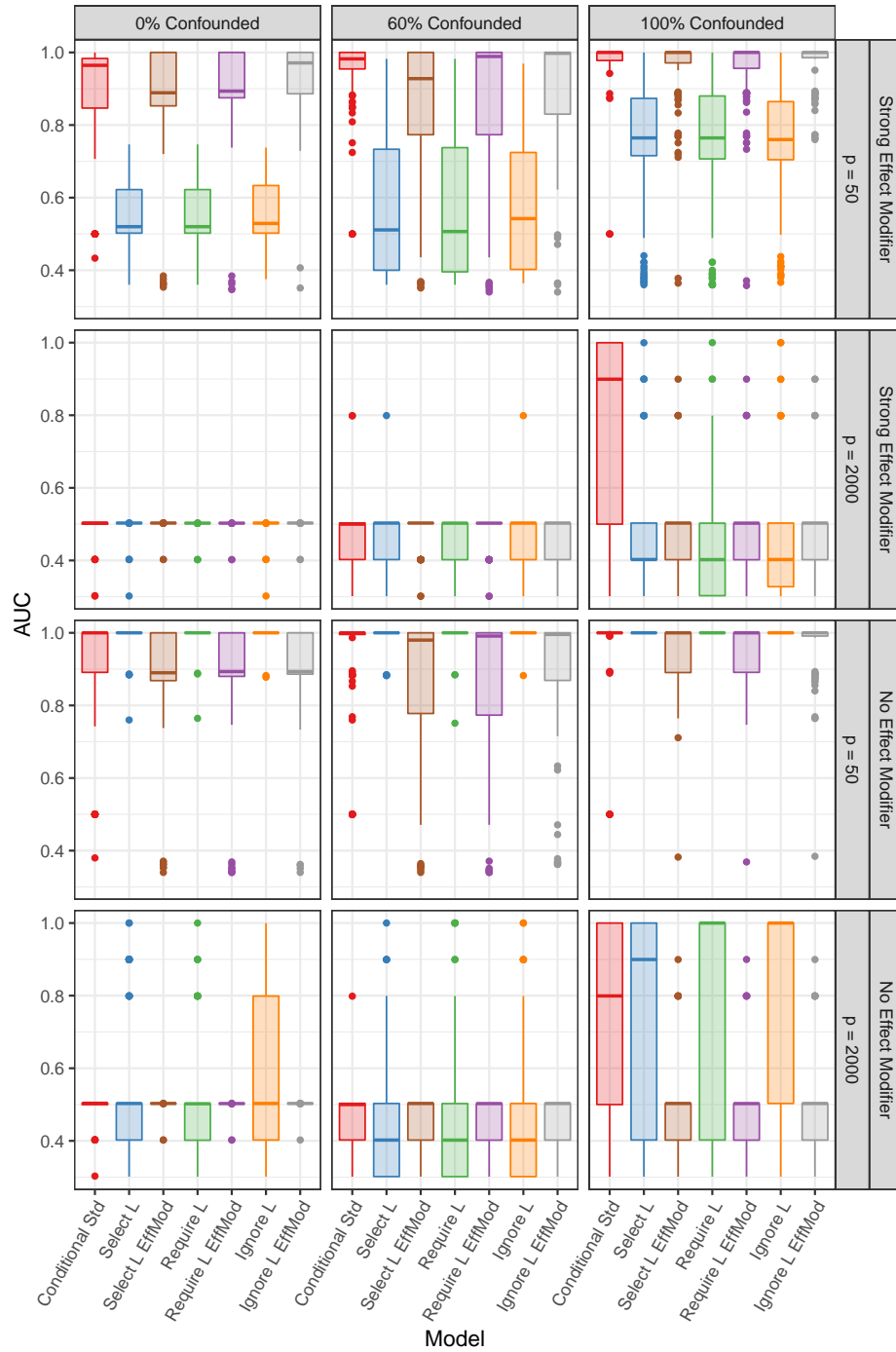


Figure S2: Simulation results: box plots of the area under the curve (AUC) from 100 simulation replications for $n = 50$ and negative binomial features using p -values based on the debiased LASSO estimate following iterative sure independence screening (iterative SIS).

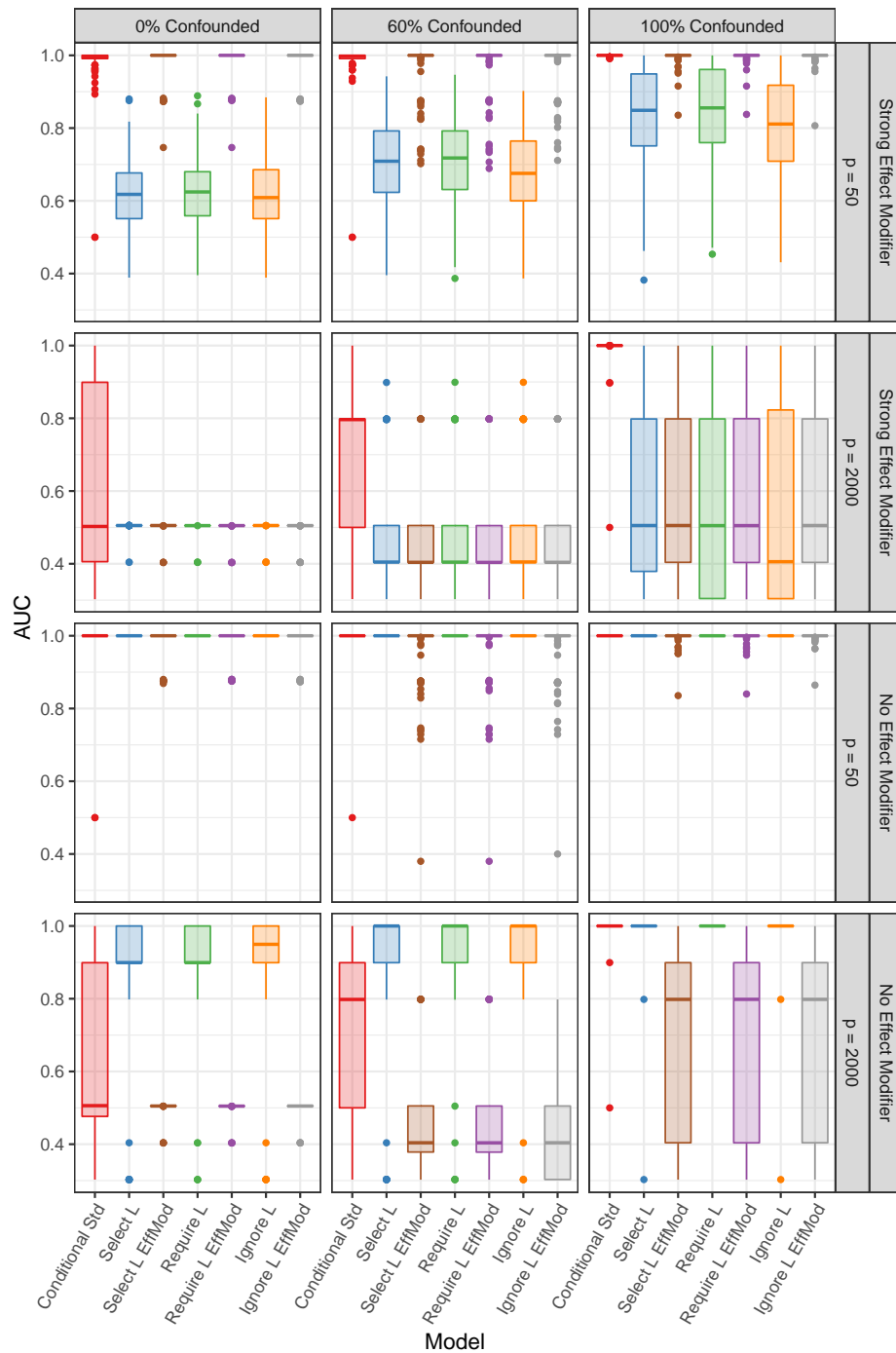


Figure S3: Simulation results: box plots of the area under the curve (AUC) from 100 simulation replications for $n = 100$ and negative binomial features using p -values based on the debiased LASSO estimate following iterative sure independence screening (iterative SIS).

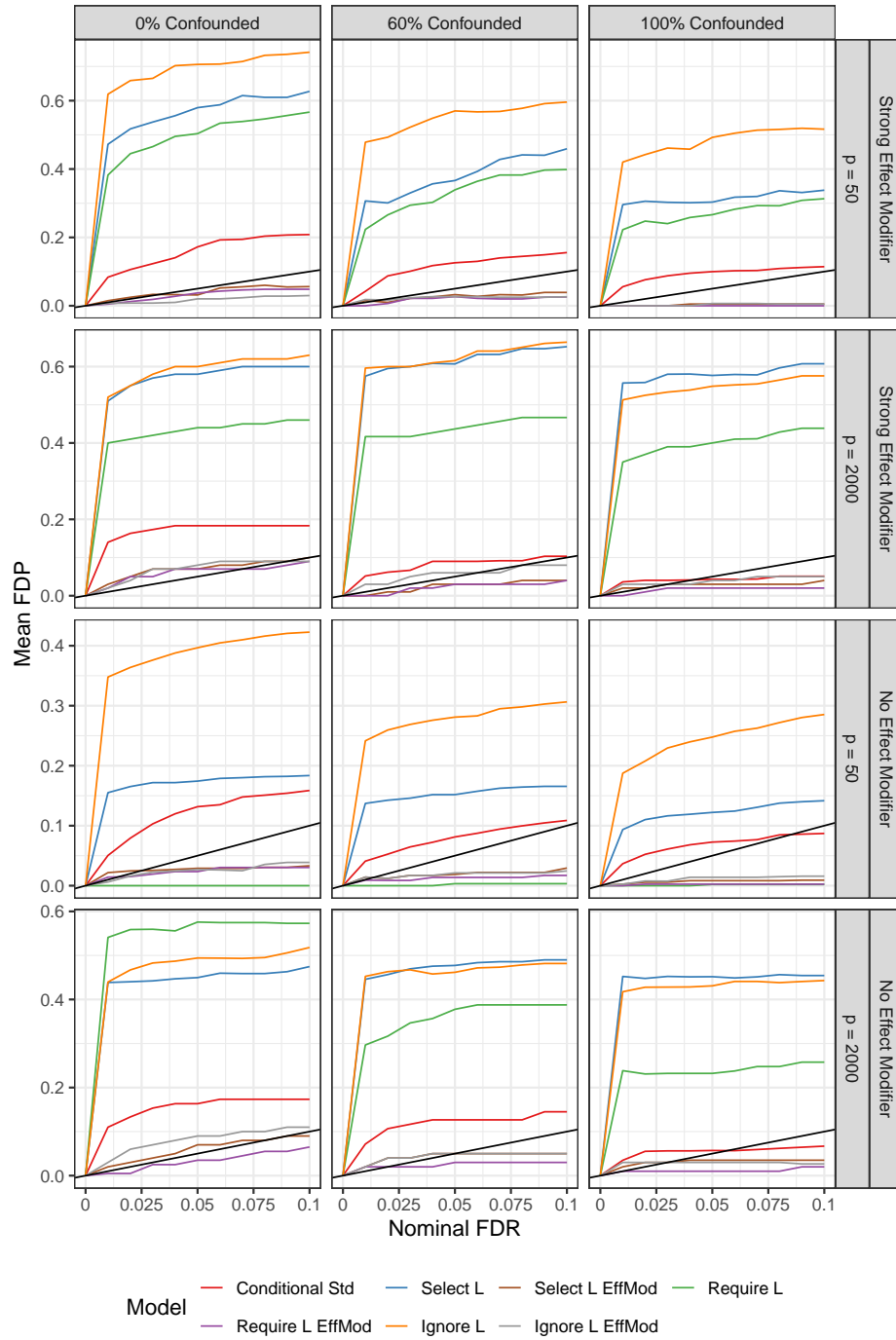


Figure S4: Simulation results: mean estimated false discovery proportion (FDP) for $n = 50$ and Poisson features at varying nominal false discovery rate (FDR) values using Benjamini-Hochberg adjusted p -values based on the debiased LASSO estimate following iterative sure independence screening (iterative SIS). The $y = x$ line is shown in black; any values above this line indicate lack of FDR control.

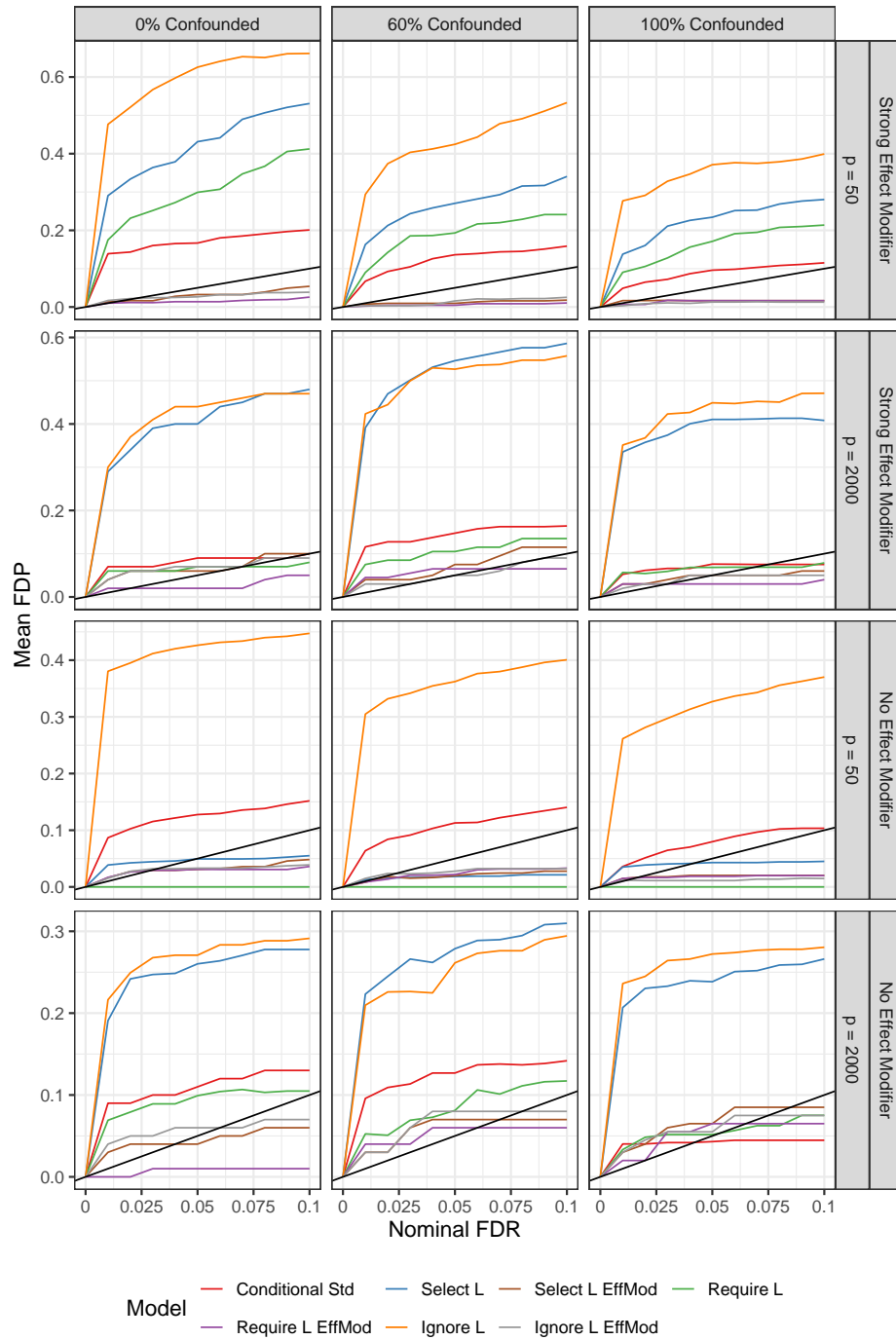


Figure S5: Simulation results: mean estimated false discovery proportion (FDP) for $n = 50$ and negative binomial features at varying nominal false discovery rate (FDR) values using Benjamini-Hochberg adjusted p -values based on the debiased LASSO estimate following iterative sure independence screening (iterative SIS). The $y = x$ line is shown in black; any values above this line indicate lack of FDR control.

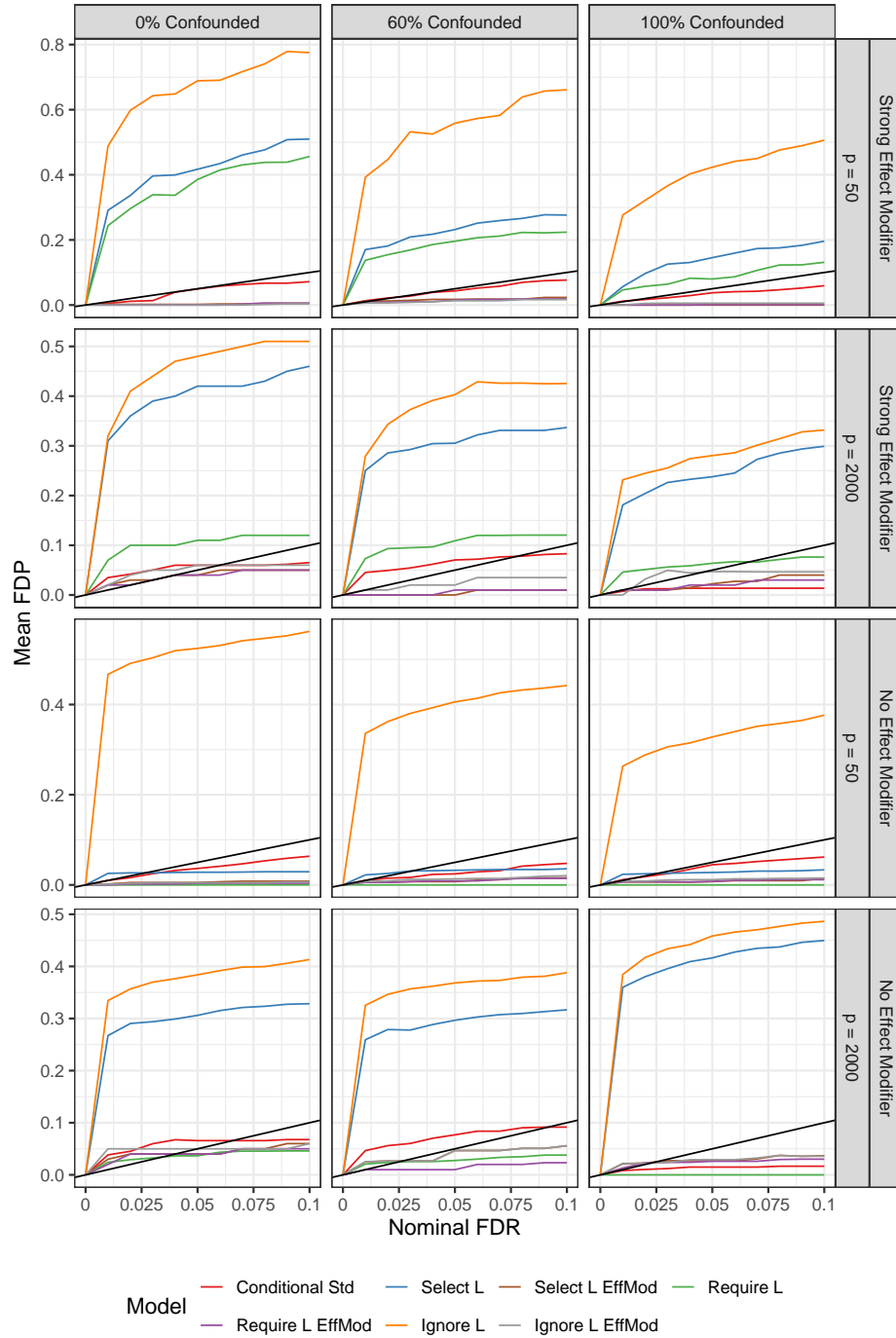


Figure S6: Simulation results: mean estimated false discovery proportion (FDP) for $n = 100$ and negative binomial features at varying nominal false discovery rate (FDR) values using Benjamini-Hochberg adjusted p -values based on the debiased LASSO estimate following iterative sure independence screening (iterative SIS). The $y = x$ line is shown in black; any values above this line indicate lack of FDR control.

S2 Real data analysis: correlation structure and assumption checks

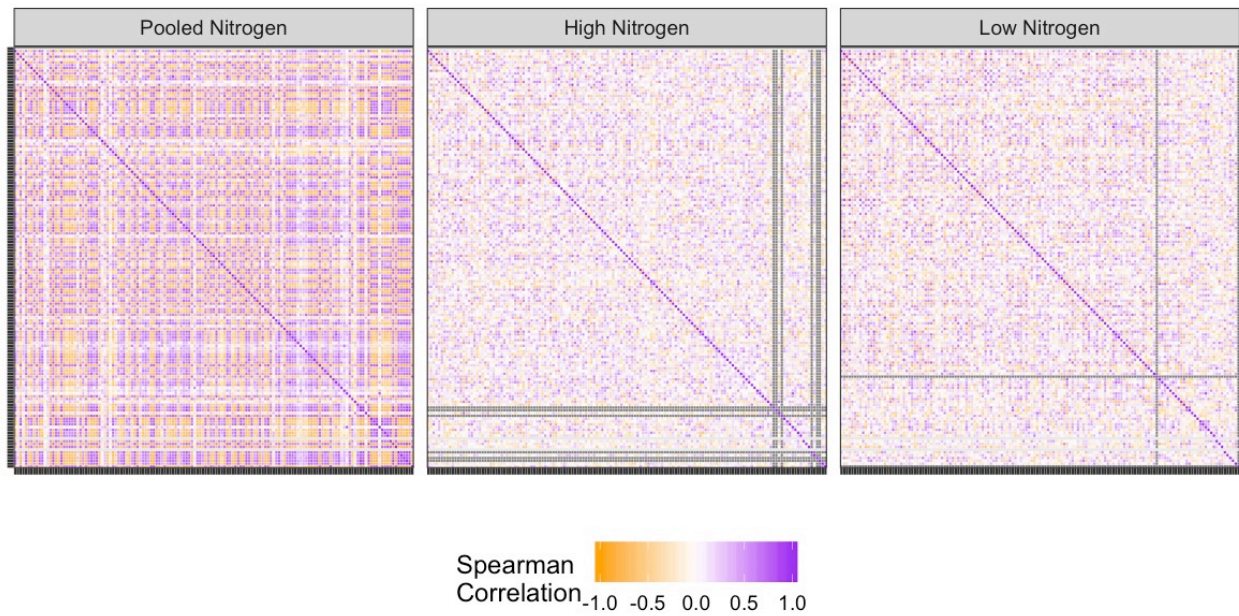


Figure S7: Sorghum microbiome data: Spearman's correlation between the top 150 marginally correlated OTUs for all samples, pooled across both nitrogen conditions (left) and stratified by nitrogen application (center and right).

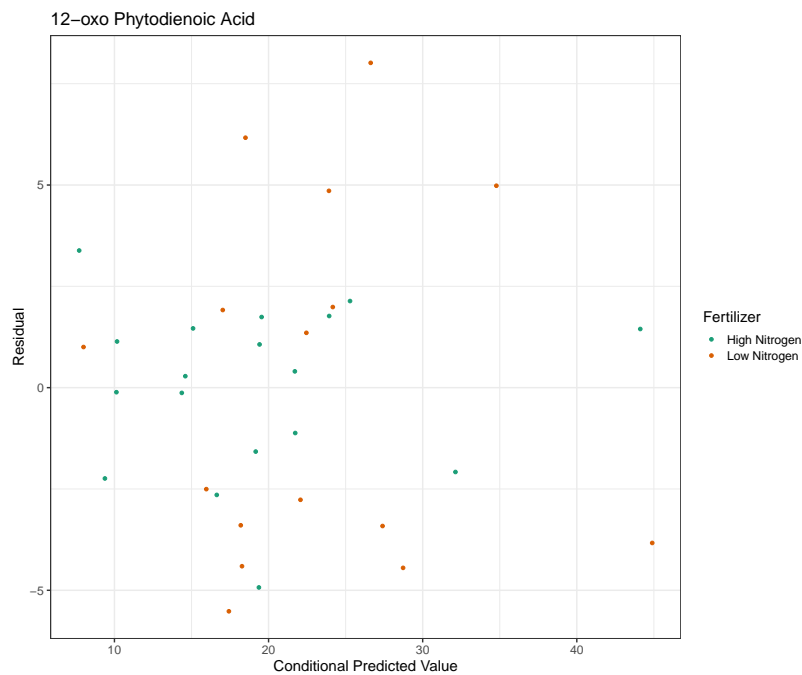


Figure S8: Real data analysis: residuals versus conditional (nitrogen stratum-specific) predicted values from the conditional debiased iterative SIS-LASSO estimates.

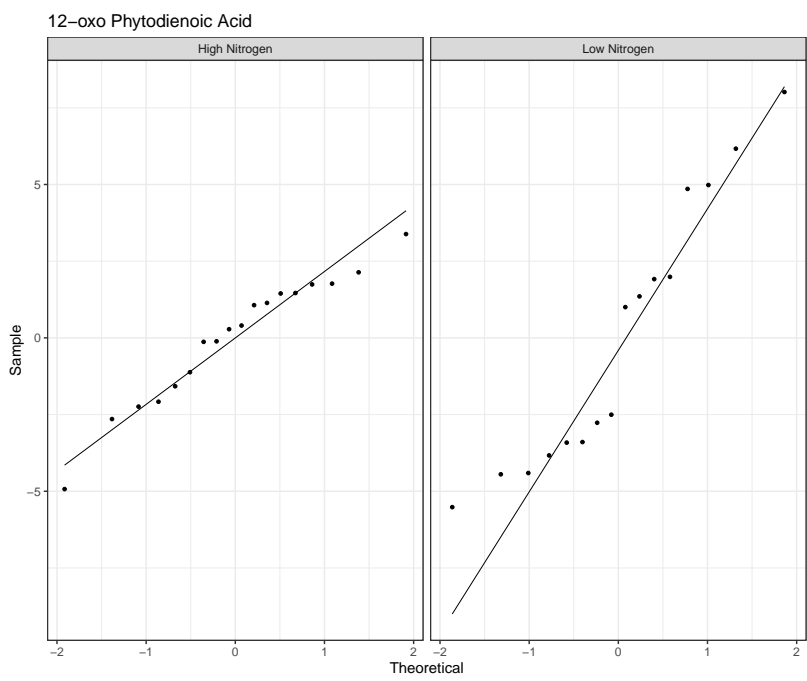


Figure S9: Real data analysis: Q-Q plot of residuals based on conditional (nitrogen stratum-specific) predicted values from the conditional debiased iterative SIS-LASSO estimates.