



# PROLONG

Penalized Regression on Longitudinal multi-Omics Data with Network and Group Lasso Constraints

Steve Broll <sup>1</sup>

Advised by Sumanta Basu, <sup>1</sup> Martin Wells <sup>1</sup>, and Myung Hee Lee <sup>2</sup>

<sup>1</sup>Cornell University

<sup>2</sup>Weill Cornell Medicine



Clinician gives you a longitudinal clinical outcome, along with hundreds of longitudinal -omics variables, and asks

Which variables co-vary with the outcome?



We have:

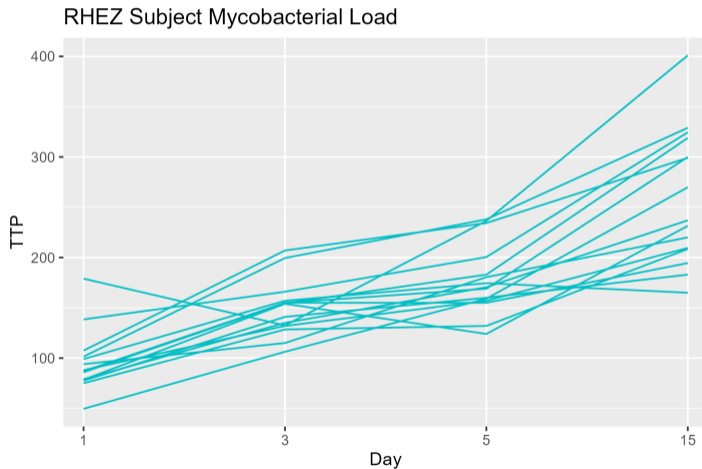
- Longitudinal measurements for some continuous phenotype and for -omics variables with only a few time points
- Large amount of variables with relatively small number of subjects

We want to:

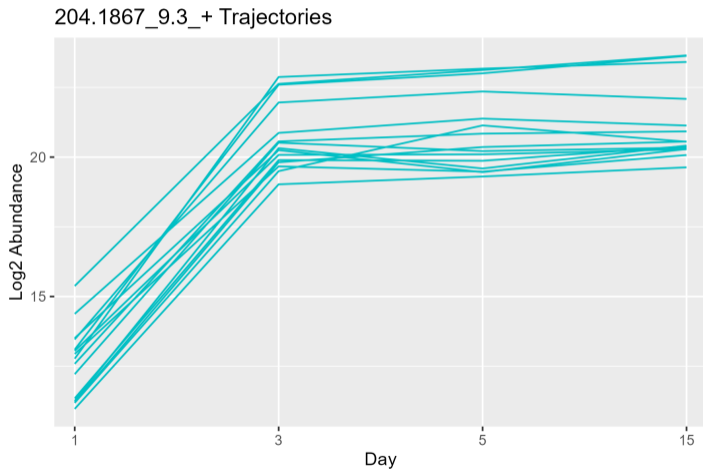
- Identify -omics variables that co-vary with the phenotype
- Overcome time dependence, low signal, and high subject variability
- Incorporate correlation of the variables



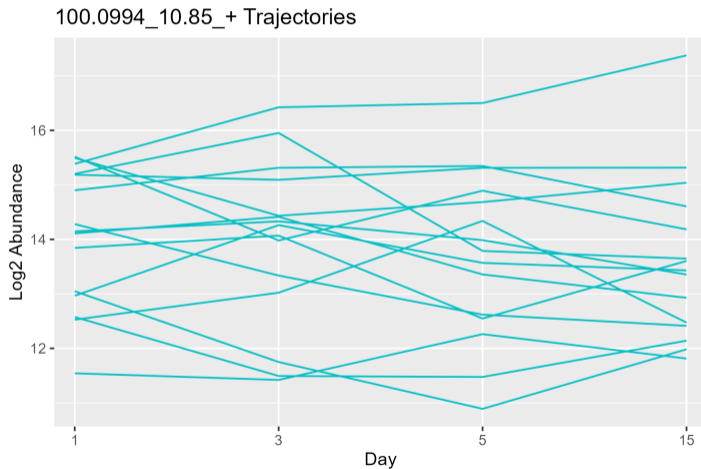
- 15 subjects, TB patients treated with RHEZ [rifampin (R), isoniazid (H), ethambutol (E), and pyrazinamide (Z)]
- Mycobacterial load measured by Time to Positivity (TTP)
- 352 metabolites with complete measurements for  $>80\%$  of subjects, softImpute used for missing values
- 4 time points, days 1, 3, 5, 15
- Additionally, we have microbiome and RNAseq data [1] for days 1 and 15 - more on this later



# TB Example Metabolite 1



# TB Example Metabolite 2





- Take first difference of the data to deal with observed temporal dependence
- Stack our  $t - 1$  first differenced value of  $X$  and  $Y$  so we have

$$Y = |Y_4 - Y_3 \quad Y_3 - Y_2 \quad Y_2 - Y_1|^T$$

And for each variable  $j$  we have

$$X_j = |X_{j4} - X_{j3} \quad X_{j3} - X_{j2} \quad X_{j2} - X_{j1}|^T$$

- Set up design matrix so that each first differenced  $Y$  value is regressed on all prior first differenced values of  $X$  to account for potential lags
- Apply network and group lasso penalties to induce sparsity while utilizing correlation and inherent group structure



$$\tilde{Y} = \begin{bmatrix} \tilde{Y}_{11} & \cdots & \tilde{Y}_{1T} \\ \vdots & & \\ \tilde{Y}_{n1} & \cdots & \tilde{Y}_{nT} \end{bmatrix}_{n \times T} \rightarrow \begin{bmatrix} \Delta \tilde{Y}_{11} & \cdots & \Delta \tilde{Y}_{1(T-1)} \\ \vdots & & \\ \Delta \tilde{Y}_{n1} & \cdots & \Delta \tilde{Y}_{n(T-1)} \end{bmatrix}_{n \times (T-1)}$$

$$\rightarrow Y = \begin{bmatrix} \Delta \tilde{Y}_{11} \\ \vdots \\ \Delta \tilde{Y}_{n1} \\ \Delta \tilde{Y}_{1(T-1)} \\ \vdots \\ \Delta \tilde{Y}_{n(T-1)} \end{bmatrix}_{n(T-1) \times 1}$$

# Moving $X$ from Tensor to Matrix



$$\tilde{X} = \begin{bmatrix} \tilde{X}_{11}^{[j]} & \cdots & \tilde{X}_{1T}^{[j]} \\ \vdots & & \vdots \\ \tilde{X}_{n1}^{[j]} & \cdots & \tilde{X}_{nT}^{[j]} \end{bmatrix}_{n \times T} \rightarrow \begin{bmatrix} \Delta \tilde{X}_{11}^{[j]} & \cdots & \Delta \tilde{X}_{1(T-1)}^{[j]} \\ \vdots & & \vdots \\ \Delta \tilde{X}_{n1}^{[j]} & \cdots & \Delta \tilde{X}_{n(T-1)}^{[j]} \end{bmatrix}_{n \times (T-1)}$$

$$\rightarrow X^{[j]} = \left[ \begin{array}{c|cc|c|ccc} \Delta \tilde{X}_{11}^{[j]} & & & 0 & & & \\ \vdots & & & 0 & & & \\ \Delta \tilde{X}_{n1}^{[j]} & & & 0 & & & \\ \hline 0 & \Delta \tilde{X}_{11}^{[j]} & \Delta \tilde{X}_{12}^{[j]} & 0 & & & 0 \\ & \vdots & & & & & \\ 0 & \Delta \tilde{X}_{n1}^{[j]} & \Delta \tilde{X}_{n2}^{[j]} & 0 & & & \\ \hline 0 & 0 & & \ddots & & & 0 \\ \hline 0 & 0 & 0 & 0 & \Delta \tilde{X}_{11}^{[j]} & \cdots & \Delta \tilde{X}_{1(T-1)}^{[j]} \\ & & & & \vdots & & \\ & & & & \Delta \tilde{X}_{n1}^{[j]} & \cdots & \Delta \tilde{X}_{n(T-1)}^{[j]} \end{array} \right]_{n(T-1) \times T(T-1)/2}$$



Given our stacked response vector  $Y$  and design matrix  $X$  we seek to minimize

$$(Y - X\beta)^T(Y - X\beta) + \lambda_1 \sum_{j=1}^p \|\beta_{(j)}\|_2 + \lambda_2 \beta^T L \beta,$$

- $\lambda_1$  is the tuning parameter for our group lasso penalty, where each group  $j$  corresponds to all of the representations in the design matrix of the  $j$ th variable
- $\lambda_2$  is the tuning parameter for the network penalty
- $L$  is the Laplacian matrix for the weighted graph where the edge weights between each pair of variables are their absolute correlation



- Each variable is represented multiple times in the model, but the group lasso penalty results in either all zero or all non-zero coefficients for the representations of each variable, helping interpretability
- If two variables are highly correlated, and one is a strong enough predictor to be selected, the other variable is likely to be selected as well
- If two variables are identical, either both will be selected and have the same coefficient or neither will be selected



- Linear Mixed Effects Model
- Wald test on the  $\Delta$  scale with each  $X^{[j]}$
- PROLONG

In the following simulations, the univariate models are evaluated at different FDR thresholds and compared to PROLONG



$$x_1 \sim N(\mu, \Sigma_X); \quad \mu \sim U(10, 20), \Sigma_X = \text{diag}(\sigma_1, \dots, \sigma_p), \sigma_j \sim U(1, 2)$$

$$x_2 \sim x_1 + N(d\mu, \Sigma_X); \quad d\mu = (5, \dots, 10, 0, \dots, 0)$$

$$x_t \sim x_{t-1} + N(0, \Sigma_X) \quad t \in 3, 4$$

$$y_1 \sim N(15, 5); \quad y_2 = N(y_1 + \beta(x_2 - x_1), 5)$$

$$y_3 \sim N(y_2 + \beta(x_3 - x_2) + \beta(x_2 - x_1), 5)$$

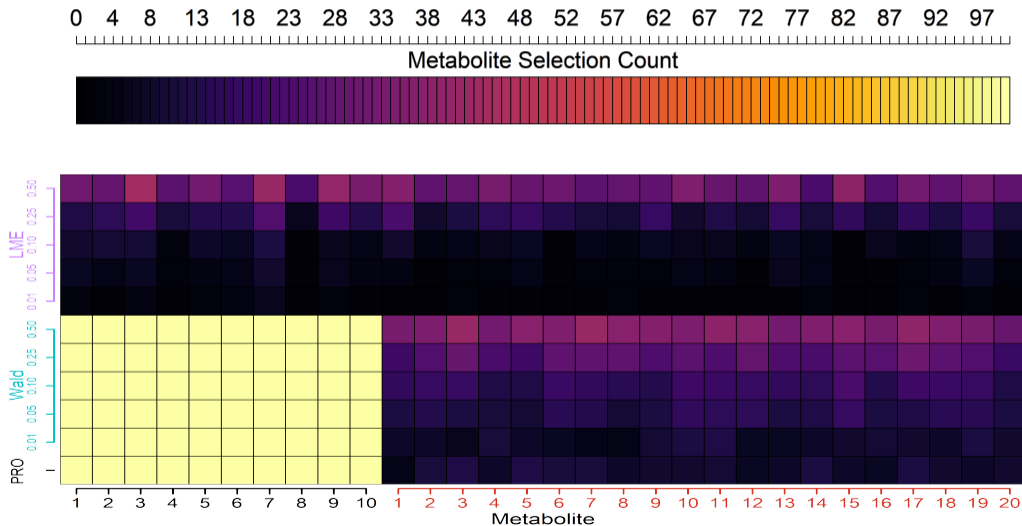
$$y_4 \sim N(y_3 + \beta(x_4 - x_3) + \beta(x_3 - x_2) + \beta(x_2 - x_1), 5)$$

$$\beta = (1/3, 1/3, \dots, 0, \dots, 0)$$

SNR ranging from 1 to 2 in targets

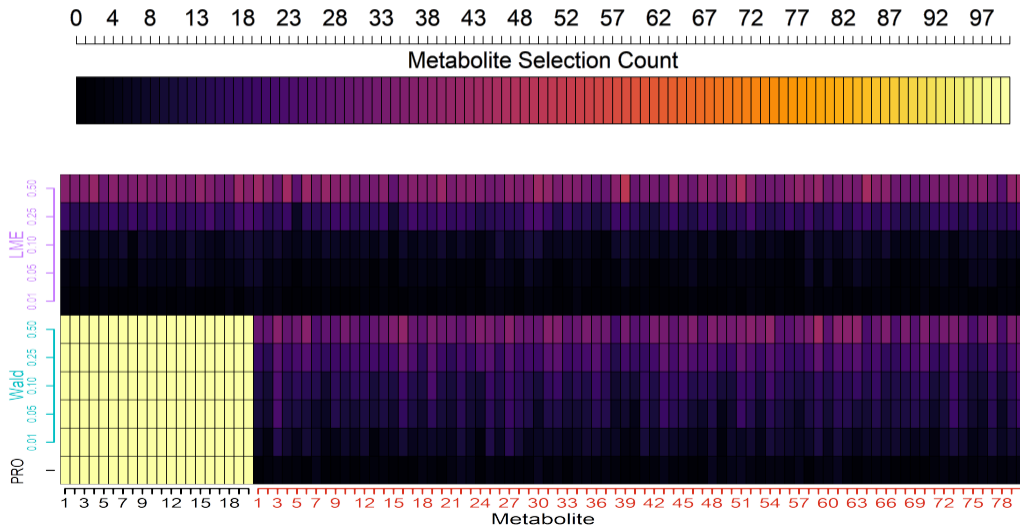
# Performance in Simulations

## Uncorrelated Simulated Variables



# Performance in Simulations

## Uncorrelated Simulated Variables







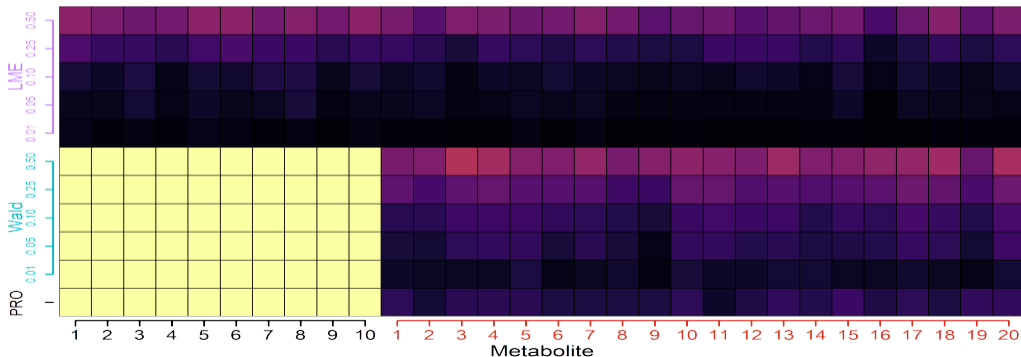
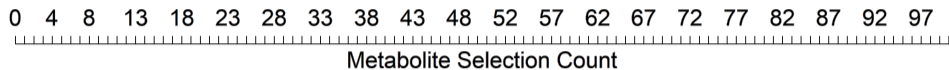
Same as previous scenario, but with

$$\Sigma_X = \begin{bmatrix} \Sigma_C & 0 \\ 0 & \Sigma_\epsilon \end{bmatrix}$$

where  $\Sigma_C$  generated so that the variances are in the same range as in  $\Sigma_\epsilon$  and the covariances correspond to correlations uniformly drawn from  $(-1, 1)$

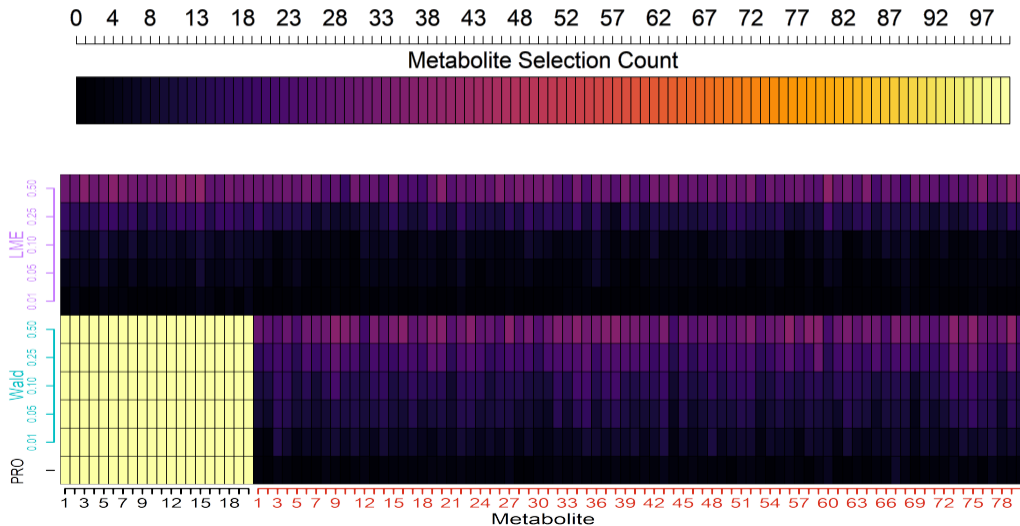
# Performance in Simulations

## Correlated Simulated Variables



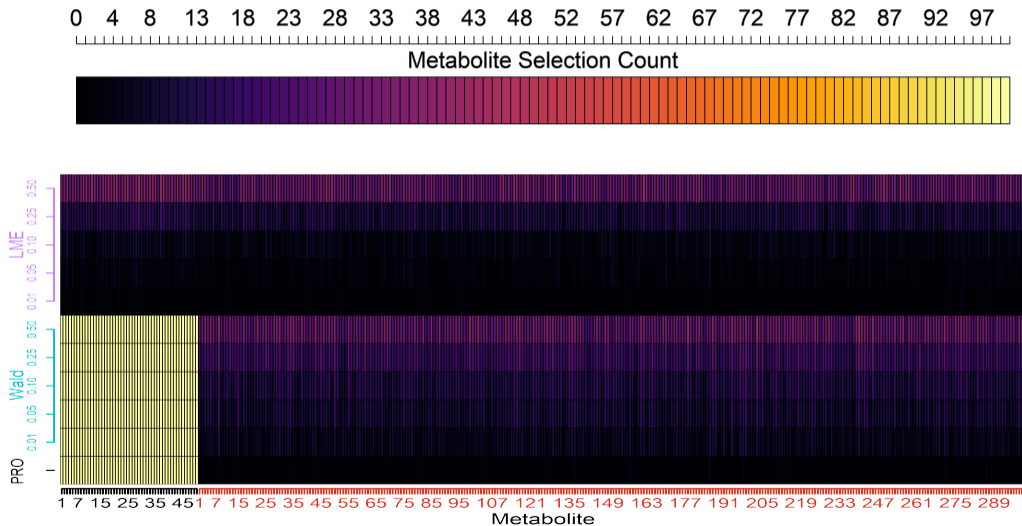
# Performance in Simulations

## Correlated Simulated Variables



# Performance in Simulations

## Correlated Simulated Variables





- Univariate methods don't pick up a single metabolite from our 352 even with an FDR of 0.5
- PROLONG selects 45 metabolites, including targets identified by our clinician collaborators and during our EDA



- High sensitivity and specificity in simulations
- Group lasso + network penalty model is slightly less sensitive at some  $\lambda_2$  values but much more specific than regular lasso + network penalty
- Limited preprocessing necessary
- Stable across choice of  $\lambda_2$ ,  $\lambda_1$  can be chosen with usual MSE cross-validation for lasso and group lasso or with a grid search using AIC/BIC, Mallows's  $C_p$ , etc



- Extension to other continuous omics variables is immediate
- Our current work is incorporating the relative abundances of 282 microbiome species measured at the first and last time points



- Zero Inflation
- Compositional data - relative abundances are used instead of raw counts
- Estimating correlation within microbiome and between microbiome and metabolites
- Subset of time points for clinical outcome and metabolomic variables
- High between-subject variation



# Microbiome Composition at Class Level





We propose incorporating the compositional data directly into the same model framework along with the metabolomic variables by using the radial transformation [2]

$$\frac{x}{\|x\|_2}$$

Additional investigation is needed to determine if Pearson's correlation using the radial transformed data is adequate for the purposes of our network constraint.



- [1] Wipperman, M.F., Bhattarai, S.K., Vorkas, C.K. et al. Gastrointestinal microbiota composition predicts peripheral inflammatory state during treatment of human tuberculosis. *Nat Commun* 12, 1141 (2021).
- [2] Park, Junyoung, et al. "Kernel Methods for Radial Transformed Compositional Data with Many Zeros." *International Conference on Machine Learning*. PMLR, (2022).