

# Regression Methods for GSEA

AJ Fagan

June 14, 2024

# Gene Set Enrichment Analysis

Goal: from some gene-level data, determine which of some gene sets are of interest in the present data

# GSEA in DE Analysis

- DE output: list of scores for each gene (insert example here)
- Fed as input into GSEA methods

# Permutation Methods

- Construct binary scores for each gene
- Determine, for each gene set, if it is enriched for 1's compared to a random set of the same size, or to the remaining genes
- Simplest example is Fisher's Exact Test to see if a score of 1 is associated with inclusion in each gene set

# Problem - Gene Set Overlap

Many gene sets, such as those from the Gene Ontology (GO), have considerable overlap. As a result, methods like those above will often return hundreds of highly correlated gene sets.

## Solution - Multiset Methods

Instead of considering only one gene set at a time, methods like SetRank and the Rolemodel consider every gene set when determining significance of any other. This enables overlap to be addressed, reducing the sizes of generated lists substantially.

# Problem - Gene-Gene Correlation

Another problem with permutation-based methods is that they fail to account for inter-gene correlation. As an example, if a gene set has 5 perfectly correlated genes found as hits, and 5 uncorrelated genes found as misses, this set should be treated closer to  $1/6$ , then  $5/10$ .

# Solution - Multivariate Regression Methods

Methods such as ROAST utilize multivariate regression techniques to account for inter-gene correlation. These methods cannot operate solely on gene-level output of some pre-existing DE analysis, as such data would lose all inter-gene information.



# Proposed Work

As far as I can find, no method currently exists that enables multiset, multivariate regression techniques that can adequately account for both the inter-gene correlation, and the gene set-overlap problem. Such a method would permit both multiset methods' cleanliness of returned gene sets, as well as the robustness to inter-gene correlation offered by the multivariate regression methodology.

# Data Model

Let  $g = 1, \dots, G$  denote the genes present in our analysis, and let  $s = 1, \dots, S$  denote the gene sets. We model, for  $i = 1, \dots, n$ :

$$y_i = X_i\alpha + Z_i\beta + \varepsilon_i,$$

where

- $y_i$  is the length  $G$  vector of gene expression data for sample  $i$
- $X_i$  is the size  $G \times k_0$  design matrix for sample  $i$  under the null of non-differential expression
- $Z_i$  is the size  $G \times k_A$  design matrix for sample  $i$  under the alternative of differential expression
- and  $\varepsilon_i$  is an error term.

# The Rolemodel

For each gene  $g$ , let  $A_g \in \{0, 1\}$  denote that gene's "activity". Similarly, for each gene set  $s$ , let  $T_s \in \{0, 1\}$  denote that gene set's "activity".

The Rolemodel asserts that, for each gene  $g$ ,

$$A_g = 1 \iff \exists s, g \in s, T_s = 1.$$

Traditionally, it operates on a set of observed gene-level binary responses  $\hat{A}_g$ , informing the likelihood of  $P(\hat{A}_g | A_g)$ , forcing

$$P(\hat{A}_g | A_g = 1) > P(\hat{A}_g | A_g = 0).$$

From this, it constructs the posterior probability that  $T_s = 1$  for each gene set  $s$ .

# The Rolemodel - Expanded

To extend the Rolemodel to function in the multivariate regression context, we simply alter the role of  $A_g$  in the model. Here, we use the latent  $A_g$  to inform the prior on  $\beta_g$ :

$$\beta_g \sim (1 - A_g)F_0(\beta_g|\lambda_0) + A_gF_1(\beta_g|\lambda_1),$$

where  $F_0(\beta_g|\lambda_0)$  is a null distribution, and  $F_1(\beta_g|\lambda_0)$  is a distribution indicative of DE.

# Example - Atovaquone vs DMSO

Design - three batches ( $j = 1, 2, 3$ ), each containing:

- Three times ( $t = 1, 8, 24$  hours)
- 3 treatments ( $i = \text{DMSO}, \text{IC50}, \text{IC75}$ )
- 3 replicates ( $k = 1, 2, 3$ )

Goal:

Explain DE in terms of GO functions.

## Example (cont'd)

Let

$$y_{itjk} = \mu + \alpha_j + \gamma_t + \beta_i + (\gamma\beta)_{it} + \varepsilon_{itjk},$$

where

- $\alpha_j$  indicates the batch mean for each gene,
- $\gamma_t$  indicates the (discrete) time effect on each gene,
- $\beta_i$  indicates the treatment effect on each gene, and
- $(\gamma\beta)_{it}$  indicates the time-treatment combination effect on each gene.

Then, under the null of non-DE, for each  $i$  and  $t$ ,

$$\beta_i = (\gamma\beta)_{it} = 0.$$

## Example - Data Likelihood

At this point, DE models such as DESeq2 or multiset GSEA models such as [Cao and Zhang, 2013] employ some discrete data distribution (negative binomial and non-central hypergeometric, respectively) to model  $y_{itjk}$ .

However, such discrete distributions are unable to account for the covariance structure of the gene expression. Therefore, we employ log normalized-counts as our dependent variable, and model,

$$y_{itjk} - \hat{y}_{itjk} = \varepsilon_{itjk} \sim N(0, \Sigma).$$

## Problem - Variance

The human genome contains approximately 20,000 *genes*. Therefore,  $\Sigma$  would be approximately  $20,000 \times 20,000$ . Assuming  $\Sigma$  was populated with 64-bit floating point numbers, this matrix would occupy 3.2GB worth of memory, and would be computationally infeasible to perform many basic operations on.



# A proposed solution - Marginalize

Suppose we want to fit

$$\mathbf{Y}_{n \times G} | \mathbf{X}, \mathbf{Z}, \mathbf{B}, \mathbf{A}, \boldsymbol{\Sigma}, \delta = \mathbf{X}_{n \times k_0} (\mathbf{B}_\delta)_{k_0 \times |\delta|} (\mathbf{R}_\delta)_{|\delta| \times G} + \mathbf{Z}_{n \times k_1} \mathbf{A}_{k_1 \times G} + \mathbf{E},$$

where,

- $\mathbf{E} \sim IW(\boldsymbol{\Sigma}, \nu),$
- $\mathbf{A} \sim MN_{n,G}(\mathbf{A}_0, \boldsymbol{\Lambda}_{0A}, \frac{1}{\kappa_{0,A}} \boldsymbol{\Sigma}),$
- $\mathbf{B}_\delta \sim MN_{n,|\delta|}(\mathbf{B}_{0\delta}, \boldsymbol{\Lambda}_{0B}, \frac{1}{\kappa_{0,B}} \boldsymbol{\Sigma}_\delta),$
- $\delta$  is a length  $G$  binary vector
- $R_\delta \in \mathbb{R}^{|\delta| \times G}$  is the matrix that inserts zero columns upon right-multiplication

## Marginalization (cont'd).

Then (TODO: NEED TO VERIFY)

$$\mathbf{Y}|\delta \sim T_{n,G}(\nu, \mathbf{Z}\mathbf{A}_0 + \mathbf{X}\mathbf{B}_{0,\delta}\mathbf{R}_\delta, \mathbf{I} + \mathbf{X}\mathbf{\Lambda}_{0B}\mathbf{X}' + \mathbf{Z}\mathbf{\Lambda}_{0A}\mathbf{Z}', \frac{1 + \kappa_{0A}}{\kappa_{0A}}\mathbf{\Sigma} + \frac{1}{\kappa_{0B}}\mathbf{R}_\delta'\mathbf{\Sigma}_\delta\mathbf{R}_\delta)$$

This is useful, as we can avoid sampling coefficients and variances.

$$P(\delta_J = 1|Y, \delta_{\setminus J}) = \left( 1 + \frac{P(\mathbf{Y}|\delta_J = 0, \delta_{\setminus J})}{P(\mathbf{Y}|\delta_J = 1, \delta_{\setminus J})} \times \frac{1 - P(\delta_J = 1|\delta_{\setminus J})}{P(\delta_J = 1|\delta_{\setminus J})} \right)^{-1}$$

The above equation enables us to perform Gibbs sampling on activity indicators.

To avoid the large variance matrix, I propose the integration of gene regulatory networks into the analysis. Such methods have been used before [Wang et al., 2024, Alexeyenko et al., 2012]. However, to my knowledge, such methods, again, do not account for the problem of gene set overlap.

GRNs can enable

- smaller models
  - [Fang et al., 2021] lists 254\_792 TF-Target pairs, as opposed to the 400\_000\_000 needed in the previous formulation
- directionality, as in [Wang et al., 2024]
- re-use of existing structures and packages to simplify analysis

# Possible Applications - DAGs

## ① With DAG's

- add, as regressors to each gene, the expression level of its parents from the same sample

# Possible Applications - non-DAGs



Alexeyenko, A., Lee, W., Pernemalm, M., Guegan, J., Dessen, P., Lazar, V., Lehtiö, J., and Pawitan, Y. (2012).

Network enrichment analysis: extension of gene-set enrichment analysis to gene networks.

*BMC Bioinformatics*, 13(1):226.



Cao, J. and Zhang, S. (2013).

A bayesian extension of the hypergeometric test for functional enrichment analysis.

*Biometrics*, 70(1):84–94.



Fang, L., Li, Y., Ma, L., Xu, Q., Tan, F., and Chen, G. (2021).

GRNdb: decoding the gene regulatory networks in diverse human and mouse conditions.

*Nucleic Acids Res*, 49(D1):D97–D103.



Wang, B., van der Kloet, F., Kes, M. B. M. J., Luirink, J., and Hamoen, L. W. (2024).

Improving gene set enrichment analysis (GSEA) by using regulation directionality.

*Microbiol Spectr*, 12(3):e0345623.