

Addressing Scale Uncertainty in Gene and Microbe Set Enrichment Analysis

Kyle McGovern

The Pennsylvania State University

kvm6065@psu.edu

GLBIO 2024

May 2, 2024

Review of Key Concepts

Consider as an example an 16S rRNA-seq experiment measuring D taxa in the colons of N patients:

$$\underbrace{W_{dn}}_{\substack{\text{Absolute Abundance} \\ \text{Taxa } d, \text{ Patient } n \\ \text{(Unmeasured)}}} = \underbrace{W_{dn}^{\parallel}}_{\substack{\text{Composition} \\ \text{Taxa } d, \text{ Patient } n \\ \text{(Measured)}}} \times \underbrace{W_n^{\perp}}_{\substack{\text{Scale} \\ \text{(e.g., total \# of microbes in} \\ \text{patient } n\text{'s colon)} \\ \text{(Unmeasured)}}}$$

Review of Key Concepts

Consider as an example an 16S rRNA-seq experiment measuring D taxa in the colons of N patients:

$$\underbrace{W_{dn}}_{\substack{\text{Absolute Abundance} \\ \text{Taxa } d, \text{ Patient } n \\ \text{(Unmeasured)}}} = \underbrace{W_{dn}^{\parallel}}_{\substack{\text{Composition} \\ \text{Taxa } d, \text{ Patient } n \\ \text{(Measured)}}} \times \underbrace{W_n^{\perp}}_{\substack{\text{Scale} \\ \text{(e.g., total \# of microbes in} \\ \text{patient } n\text{'s colon)} \\ \text{(Unmeasured)}}}$$

Further consider as an example estimation of the LFC (Log Fold Change) of taxa d in patients with and without Ulcerative Colitis:

$$\underbrace{\theta_d}_{\text{LFC in Absolute Abundance}} = \underbrace{\theta_d^{\parallel}}_{\text{LFC in Composition}} + \underbrace{\theta_d^{\perp}}_{\text{LFC in Scale}}.$$

Review of Key Concepts

Methods like ALDEx2, DESeq2, Limma, etc. estimate LFCs using sequence count data Y :

$$f(Y) = \hat{\theta}_d$$
$$= \underbrace{\hat{\theta}_d^{\parallel}}_{\substack{\text{Estimated LFC in the} \\ \text{measured composition}}} + \underbrace{\hat{\theta}_d^{\perp}}_{\substack{\text{Estimated LFC in the} \\ \text{unmeasured scale}}}.$$

Review of Key Concepts

Methods like ALDEx2, DESeq2, Limma, etc. estimate LFCs using sequence count data Y :

$$f(Y) = \hat{\theta}_d = \underbrace{\hat{\theta}_d^{\parallel}}_{\text{Estimated LFC in the measured composition}} + \underbrace{\hat{\theta}_d^{\perp}}_{\text{Estimated LFC in the unmeasured scale}}.$$

Estimates $\hat{\theta}^{\perp}$ come from normalization, for example:

- Total Sum Scaling (TSS): $\hat{\theta}^{\perp} = 0$

Review of Key Concepts

Methods like ALDEx2, DESeq2, Limma, etc. estimate LFCs using sequence count data Y :

$$f(Y) = \hat{\theta}_d = \underbrace{\hat{\theta}_d^{\parallel}}_{\substack{\text{Estimated LFC in the} \\ \text{measured composition}}} + \underbrace{\hat{\theta}_d^{\perp}}_{\substack{\text{Estimated LFC in the} \\ \text{unmeasured scale}}}.$$

Estimates $\hat{\theta}^{\perp}$ come from normalization, for example:

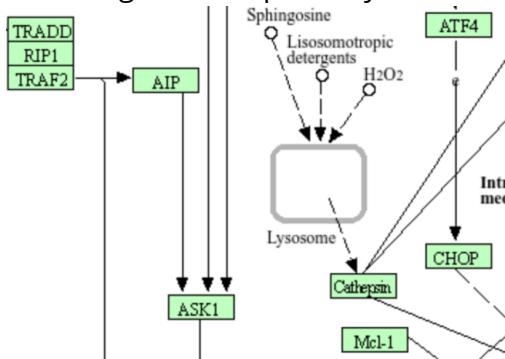
- Total Sum Scaling (TSS): $\hat{\theta}^{\perp} = 0$
- Centered Log Ratio (CLR): $\hat{\theta}^{\perp} = -\text{mean}(\hat{\theta}^{\parallel})$

Differential Set Analysis (DSA)

Rather than estimating **LFCs** of **single** genes/taxa

ASK1

What if we are interested in a **set** of genes in a pathway?



Differential Set Analysis (DSA) is used to estimate enrichment or depletion of a gene/taxa set

Key Points of this Talk

- 1 Errors in scale assumptions (i.e., estimates $\hat{\theta}^\perp$, \hat{W}^\perp) inflate false positive rates in DSA

Key Points of this Talk

- 1 Errors in scale assumptions (i.e., estimates $\hat{\theta}^\perp$, \hat{W}^\perp) inflate false positive rates in DSA
- 2 Errors in DSA estimates are a **non-linear** function of scale errors

Key Points of this Talk

- 1 Errors in scale assumptions (i.e., estimates $\hat{\theta}^\perp$, \hat{W}^\perp) inflate false positive rates in DSA
- 2 Errors in DSA estimates are a **non-linear** function of scale errors
- 3 We have developed three solutions to these errors:
 - 1 LFC Sensitivity Analysis
 - 2 LFC Sensitivity Testing
 - 3 Compositional Weighting Methods

Three Methods for DSA

In this presentation 3 common DSA methods will be considered

- 1 **Gene Set Enrichment Analysis (GSEA) with Gene Label permutations**
- 2 Gene Set Enrichment Analysis (GSEA) with Sample Label permutations
- 3 CAMERA

GSEA with Gene Label Permutations

The GSEA Algorithm Step-by-Step

- 1 Pick a set of genes S (e.g., the apoptosis signaling pathway):

$$S = \{\text{ASK1}, \text{CHOP}, \text{TRAF2}\}$$

GSEA with Gene Label Permutations

The GSEA Algorithm Step-by-Step

- 1 Pick a set of genes S (e.g., the apoptosis signaling pathway):

$$S = \{\text{ASK1}, \text{CHOP}, \text{TRAF2}\}$$

- 2 Estimate LFCs $\hat{\theta} = f(Y)$ (i.e., with DESeq2, ALDEx2, limma)

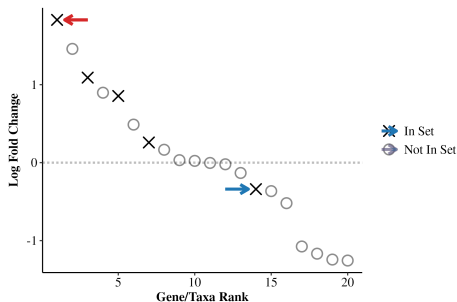
GSEA with Gene Label Permutations

The GSEA Algorithm Step-by-Step

- 1 Pick a set of genes S (e.g., the apoptosis signaling pathway):

$$S = \{\text{ASK1, CHOP, TRAF2}\}$$

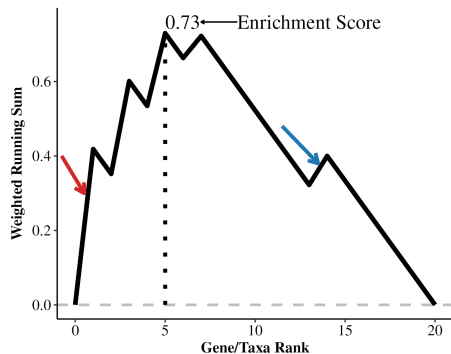
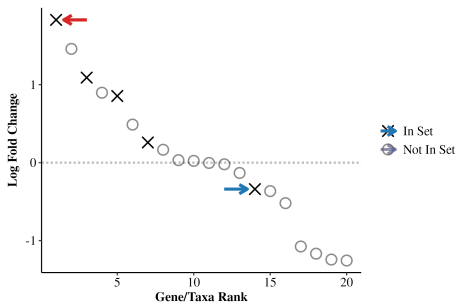
- 2 Estimate LFCs $\hat{\theta} = f(Y)$ (i.e., with DESeq2, ALDEx2, limma)
- 3 Order the LFCs from largest to smallest



GSEA with Gene Label Permutations

The GSEA Algorithm Step-by-Step

- 3 Calculate a running sum **weighted** by the LFC
- 4 Calculate an enrichment score (max distance from 0 of **weighted** running sum)



GSEA with Gene Label Permutations

The GSEA Algorithm Step-by-Step

- 5 Calculate a null distribution of Enrichment Scores (ESs)

$$S = \{\mathbf{ASK1}, \mathbf{CHOP}, \mathbf{TRAF2}\} \implies \text{ES}$$

$$S_1^* = \{\mathbf{ASK1}, \mathbf{CHOP}, \text{B2M}\} \implies \text{ES}_1^*$$

$$S_2^* = \{\text{BRCA1}, \text{EGFR}, \text{XRCC4}\} \implies \text{ES}_2^*$$

GSEA with Gene Label Permutations

The GSEA Algorithm Step-by-Step

- 5 Calculate a null distribution of Enrichment Scores (ESs)

$$S = \{\mathbf{ASK1}, \mathbf{CHOP}, \mathbf{TRAF2}\} \implies \text{ES}$$

$$S_1^* = \{\mathbf{ASK1}, \mathbf{CHOP}, \text{B2M}\} \implies \text{ES}_1^*$$

$$S_2^* = \{\text{BRCA1}, \text{EGFR}, \text{XRCC4}\} \implies \text{ES}_2^*$$

- 6 Use null distribution to calculate a p-value

DSA Target Estimand

The goal of DSA is to estimate a **target estimand** ϕ_S :

$$\phi_S = \begin{cases} 1 & \text{Gene Set } S \text{ is significantly enriched} \\ -1 & \text{Gene Set } S \text{ is significantly depleted} \\ 0 & \text{Gene Set } S \text{ is not significantly changing.} \end{cases}$$

In GSEA the target estimand is a function of the **true** LFCs:

$$\phi_S = g(\theta)$$

LFC Sensitivity Analysis

In GSEA the target estimand is a function of the **true** LFCs:

$$\phi_S = g(\theta)$$

But we don't know the true LFCs, we only have **estimates**:

$$\begin{aligned}\hat{\phi}_S &= g(\hat{\theta}) \\ &= g(\hat{\theta}^{\parallel} + \underbrace{\hat{\theta}^{\perp}}_{\substack{\text{Estimated LFC in Scale} \\ \text{(Normalization Assumption)}}}).\end{aligned}$$

LFC Sensitivity Analysis

In GSEA the target estimand is a function of the **true** LFCs:

$$\phi_S = g(\theta)$$

But we don't know the true LFCs, we only have **estimates**:

$$\begin{aligned}\hat{\phi}_S &= g(\hat{\theta}) \\ &= g(\hat{\theta}^{\parallel} + \underbrace{\hat{\theta}^{\perp}}_{\substack{\text{Estimated LFC in Scale} \\ \text{(Normalization Assumption)}}}).\end{aligned}$$

Our DSA estimate $\hat{\phi}_S$ **depends on our scale estimate $\hat{\theta}^{\perp}$!**

LFC Sensitivity Analysis

In GSEA the target estimand is a function of the **true** LFCs:

$$\phi_S = g(\theta)$$

But we don't know the true LFCs, we only have **estimates**:

$$\begin{aligned}\hat{\phi}_S &= g(\hat{\theta}) \\ &= g(\hat{\theta}^{\parallel} + \underbrace{\hat{\theta}^{\perp}}_{\substack{\text{Estimated LFC in Scale} \\ \text{(Normalization Assumption)}}}).\end{aligned}$$

Our DSA estimate $\hat{\phi}_S$ **depends on our scale estimate $\hat{\theta}^{\perp}$!**

A sensitivity analysis of how error in $\hat{\theta}^{\perp}$ affects ϕ_S

LFC Sensitivity Analysis

Error ϵ^\perp in our estimate of the unmeasured scale θ^\perp :

$$\underbrace{\theta^\perp}_{\text{True LFC in Scale}} = \underbrace{\hat{\theta}^\perp}_{\text{Estimate}} + \underbrace{\epsilon^\perp}_{\text{Estimation Error}}$$

LFC Sensitivity Analysis

Error ϵ^\perp in our estimate of the unmeasured scale θ^\perp :

$$\underbrace{\theta^\perp}_{\text{True LFC in Scale}} = \underbrace{\hat{\theta}^\perp}_{\text{Estimate}} + \underbrace{\epsilon^\perp}_{\text{Estimation Error}}$$

How does the **true** ϕ_S change with error ϵ^\perp ?

$$\begin{aligned}\phi_S &= g(\hat{\theta}^\parallel + \hat{\theta}^\perp + \epsilon^\perp) \\ &= g(\hat{\theta} + \epsilon^\perp)\end{aligned}$$

LFC Sensitivity Analysis

Error ϵ^\perp in our estimate of the unmeasured scale θ^\perp :

$$\underbrace{\theta^\perp}_{\text{True LFC in Scale}} = \underbrace{\hat{\theta}^\perp}_{\text{Estimate}} + \underbrace{\epsilon^\perp}_{\text{Estimation Error}}$$

How does the **true** ϕ_S change with error ϵ^\perp ?

$$\begin{aligned}\phi_S &= g(\hat{\theta}^\parallel + \hat{\theta}^\perp + \epsilon^\perp) \\ &= g(\hat{\theta} + \epsilon^\perp)\end{aligned}$$

LFC Sensitivity Analysis Algorithm:

- 1 Get estimated LFCs $\hat{\theta}$ (e.g., from ALDEx2, limma, DESeq2, etc.)

LFC Sensitivity Analysis

Error ϵ^\perp in our estimate of the unmeasured scale θ^\perp :

$$\underbrace{\theta^\perp}_{\text{True LFC in Scale}} = \underbrace{\hat{\theta}^\perp}_{\text{Estimate}} + \underbrace{\epsilon^\perp}_{\text{Estimation Error}}$$

How does the **true** ϕ_S change with error ϵ^\perp ?

$$\begin{aligned}\phi_S &= g(\hat{\theta}^\parallel + \hat{\theta}^\perp + \epsilon^\perp) \\ &= g(\hat{\theta} + \epsilon^\perp)\end{aligned}$$

LFC Sensitivity Analysis Algorithm:

- 1 Get estimated LFCs $\hat{\theta}$ (e.g., from ALDEx2, limma, DESeq2, etc.)
- 2 Run GSEA with $\epsilon^\perp = 0$ (i.e., $\hat{\phi}_S = g(\hat{\theta})$)

LFC Sensitivity Analysis

Error ϵ^\perp in our estimate of the unmeasured scale θ^\perp :

$$\underbrace{\theta^\perp}_{\text{True LFC in Scale}} = \underbrace{\hat{\theta}^\perp}_{\text{Estimate}} + \underbrace{\epsilon^\perp}_{\text{Estimation Error}}$$

How does the **true** ϕ_S change with error ϵ^\perp ?

$$\begin{aligned}\phi_S &= g(\hat{\theta}^\parallel + \hat{\theta}^\perp + \epsilon^\perp) \\ &= g(\hat{\theta} + \epsilon^\perp)\end{aligned}$$

LFC Sensitivity Analysis Algorithm:

- 1 Get estimated LFCs $\hat{\theta}$ (e.g., from ALDEx2, limma, DESeq2, etc.)
- 2 Run GSEA with $\epsilon^\perp = 0$ (i.e., $\hat{\phi}_S = g(\hat{\theta})$)
- 3 Rerun GSEA with $\epsilon^\perp \neq 0$ and compare to $\epsilon^\perp = 0$ (i.e., $\phi_S = g(\hat{\theta} + \epsilon^\perp)$)

Interpreting Error ϵ^\perp and LFC Sensitivity Analysis Results

Consider error $\epsilon^\perp = \pm 0.5$:

- 1 This error corresponds to the true θ^\perp being $e^{0.5} = 1.65$ times lower/higher than $\hat{\theta}^\perp$

Interpreting Error ϵ^\perp and LFC Sensitivity Analysis Results

Consider error $\epsilon^\perp = \pm 0.5$:

- 1 This error corresponds to the true θ^\perp being $e^{0.5} = 1.65$ times lower/higher than $\hat{\theta}^\perp$
- 2 Example results if a Gene set S is sensitive to error:

$\epsilon^\perp = -0.5$	$\epsilon^\perp = 0$	$\epsilon^\perp = 0.5$
$\phi_S = 0$	$\phi_S = 1$	$\phi_S = 0$

Interpreting Error ϵ^\perp and LFC Sensitivity Analysis Results

Consider error $\epsilon^\perp = \pm 0.5$:

- 1 This error corresponds to the true θ^\perp being $e^{0.5} = 1.65$ times lower/higher than $\hat{\theta}^\perp$
- 2 Example results if a Gene set S is sensitive to error:

$\epsilon^\perp = -0.5$	$\epsilon^\perp = 0$	$\epsilon^\perp = 0.5$
$\phi_S = 0$	$\phi_S = 1$	$\phi_S = 0$

- 3 Example results if a Gene set S is not sensitive to error:

$\epsilon^\perp = -0.5$	$\epsilon^\perp = 0$	$\epsilon^\perp = 0.5$
$\phi_S = 1$	$\phi_S = 1$	$\phi_S = 1$