

Addressing Scale Uncertainty in Gene and Microbe Set Enrichment Analysis

Kyle McGovern

The Pennsylvania State University

kvm6065@psu.edu

GLBIO 2024
April 28, 2024

Review of Key Concepts

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

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

Review of Key Concepts

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

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

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 :

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

Review of Key Concepts

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

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

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 :

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

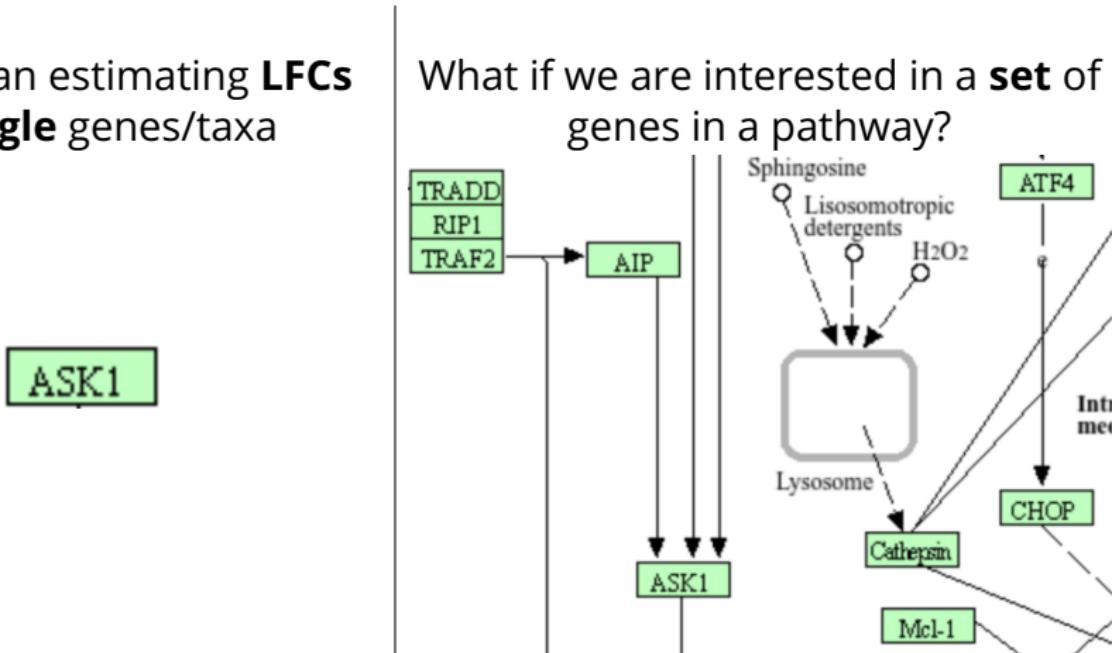
Estimates $\hat{\theta}^{\perp}$ come from normalization, for example:

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

Differential Set Analysis (DSA)

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

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):

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

GSEA with Gene Label Permutations

The GSEA Algorithm Step-by-Step

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

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

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

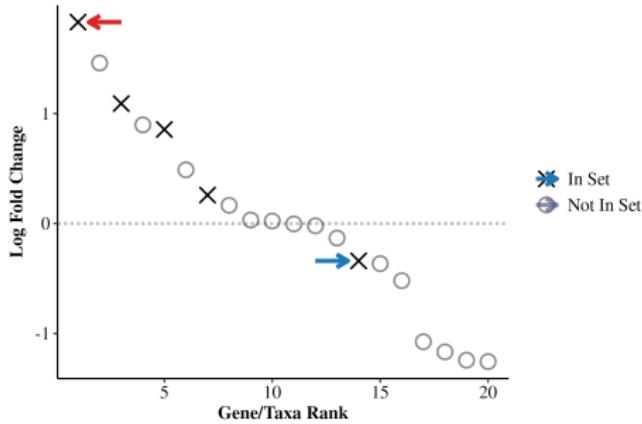
GSEA with Gene Label Permutations

The GSEA Algorithm Step-by-Step

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

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

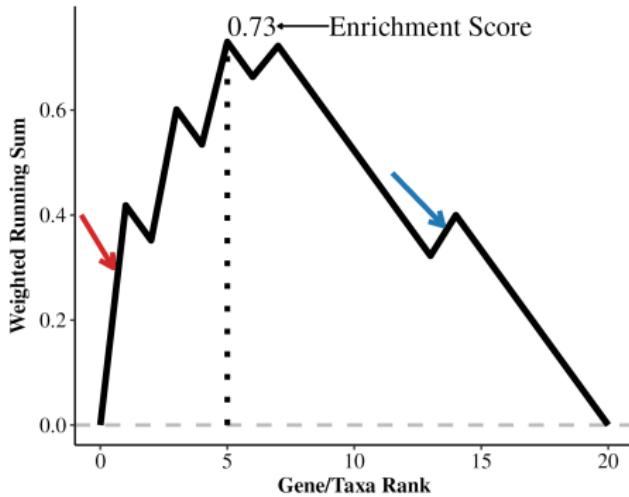
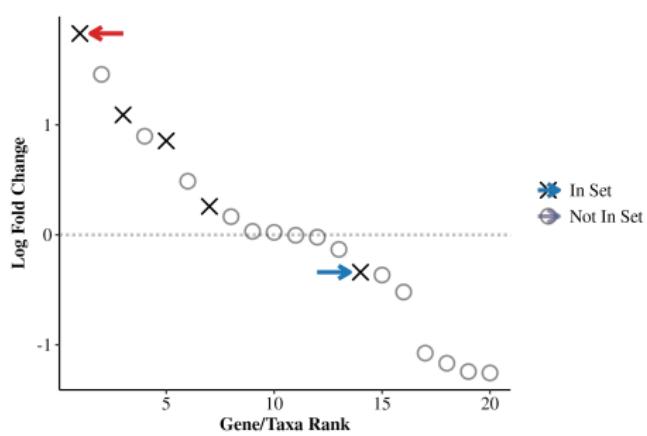
- ② Estimate LFCs $\hat{\theta} = f(Y)$ (i.e., with DESeq2, ALDEEx2, limma)
- ③ Order the LFCs from largest to smallest



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)

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

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

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

GSEA with Gene Label Permutations

The GSEA Algorithm Step-by-Step

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

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

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

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

- 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**:

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

(Normalization Assumption)

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**:

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

(Normalization Assumption)

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**:

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

(Normalization Assumption)

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

LFC Sensitivity Analysis

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

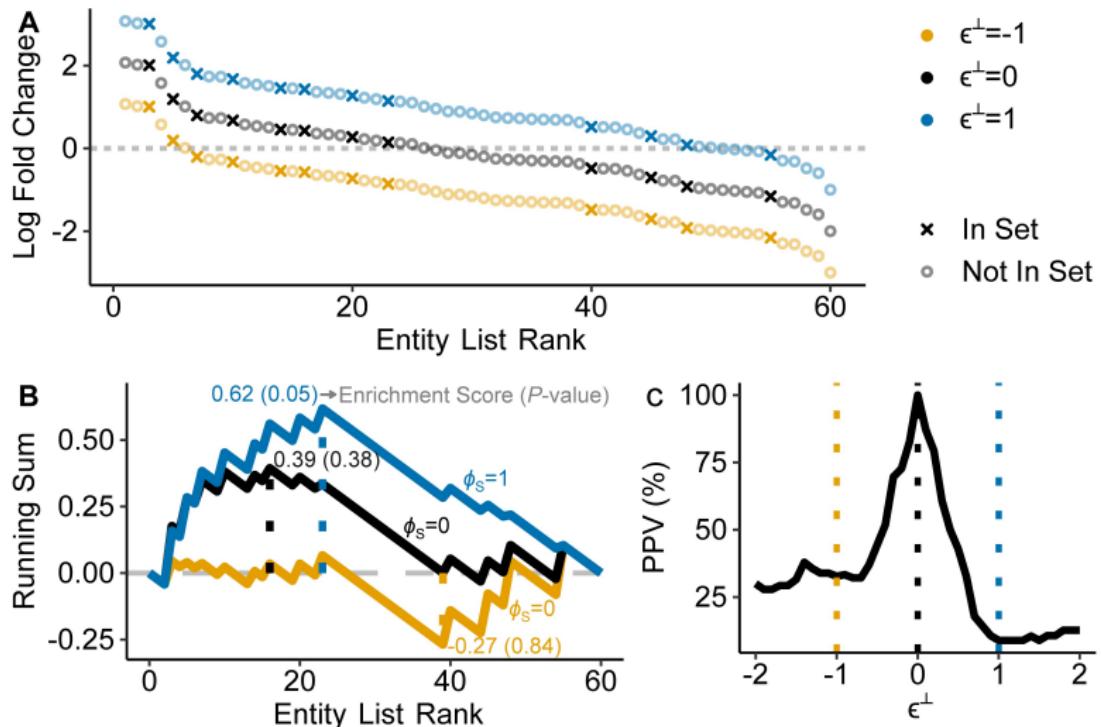
- ① This error corresponds to the true θ^\perp being $e^{0.5} = 1.65$ times lower/higher than $\hat{\theta}^\perp$
- ② Example results if a Gene set S is sensitive to error:

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

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

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

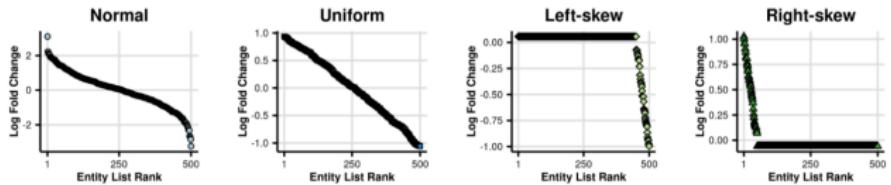
LFC Sensitivity Analysis Simulation



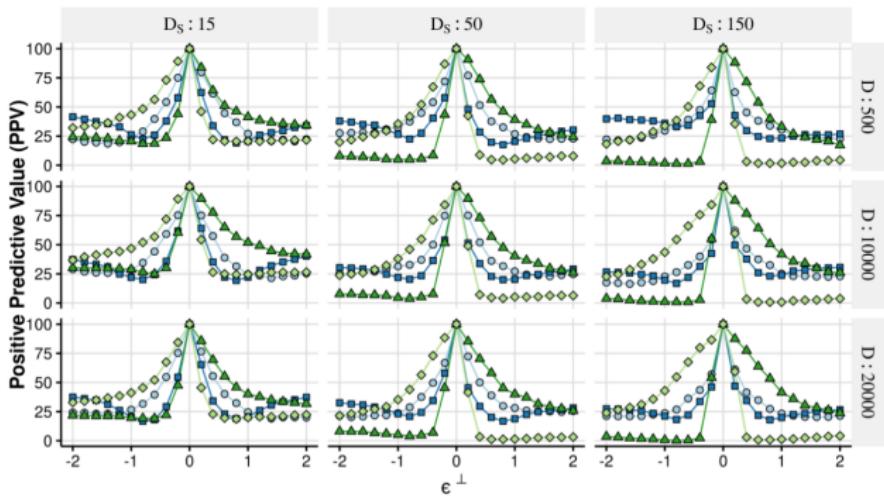
PPV (Positive Predictive Value) is the % of positives that are true positives

Multiple Simulations of Different LFC Distributions, Gene Set Sizes, and Total #'s of Genes

a



b



Real Data Analysis

RNA-seq: normal-adjacent-to-tumor vs. healthy thyroid tissue*

- ① LFCs estimated with Songbird (Morton et al., 2019)
- ② GSEA performed with fgsea (Korotkevich et al., 2021)

Real Data Analysis

RNA-seq: normal-adjacent-to-tumor vs. healthy thyroid tissue*

- ① LFCs estimated with Songbird (Morton et al., 2019)
- ② GSEA performed with fgsea (Korotkevich et al., 2021)
- ③ fgsea results for the Inflammatory Response Pathway:

```
## Run vanilla fgsea
simple_fgsea_res <- fgsea(stats=lfcs, pathways=gmt.file)
```

(Normalized) Enrichment Score	Adjusted p-value
1.53	0.003

Real Data Analysis

RNA-seq: normal-adjacent-to-tumor vs. healthy thyroid tissue*

- ① LFCs estimated with Songbird (Morton et al., 2019)
- ② GSEA performed with fgsea (Korotkevich et al., 2021)
- ③ fgsea results for the Inflammatory Response Pathway:

```
## Run vanilla fgsea
simple_fgsea_res <- fgsea(stats=lfcs, pathways=gmt.file)
```

(Normalized) Enrichment Score	Adjusted p-value
1.53	0.003

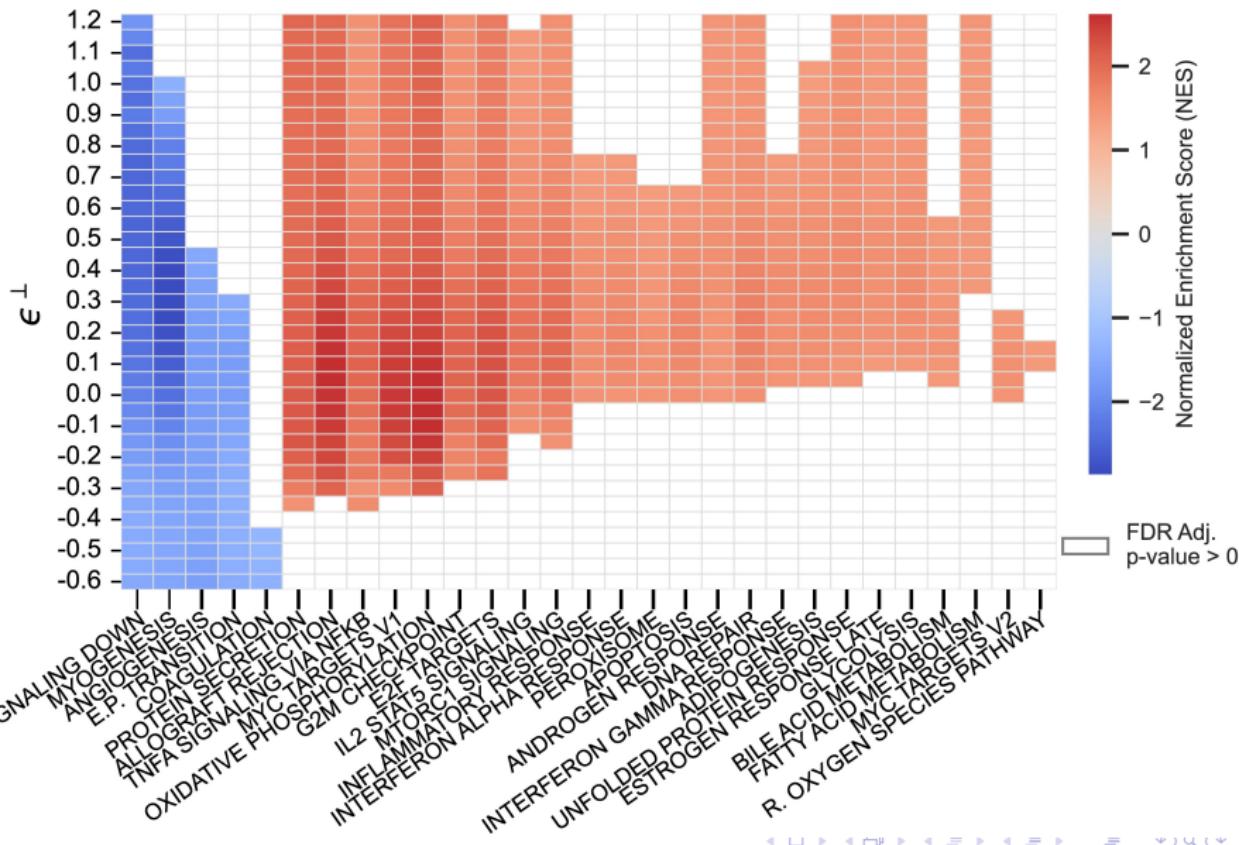
- ④ LFC Sensitivity Analysis results:

```
## Run LFC Sensitivity Analysis FGSEA
lfc_fgsea_res <- fgsea.error(lfcs, gmt.file,
                           epsilon=c(-0.5, 0, 0.5))
```

ϵ^\perp	(Normalized) Enrichment Score	Adjusted p-value
-0.5	-0.89	1
0	1.53	0.003
0.5	1.5	0

*(Aran et al., 2017)

Complete Results for Real Data Analysis



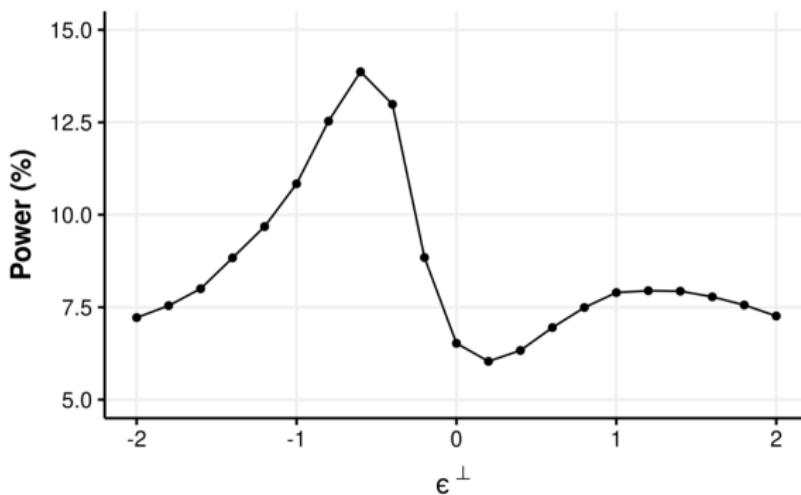
Two Additional Methods to Address False Positives

Two Additional Methods to Address False Positives

① The LFC Sensitivity Test

Let p_{ϵ^\perp} be the GSEA p-value at ϵ^\perp . Take the maximum GSEA p-value across all possible scale errors $\epsilon^\perp \in (-\infty, \infty)$:

$$p = \sup_{\epsilon^\perp \in (-\infty, \infty)} p_{\epsilon^\perp}$$



Two Additional Methods to Address False Positives

② GSEA with Compositional Weighting (GSEA-CW)

The GSEA target estimand (i.e., goal of inference) is

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

Two Additional Methods to Address False Positives

② GSEA with Compositional Weighting (GSEA-CW)

The GSEA target estimand (i.e., goal of inference) is

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

The GSEA-CW target estimand is different, and implies a different scientific question

$$\psi_S = g(\theta^{\parallel})$$

Two Additional Methods to Address False Positives

② GSEA with Compositional Weighting (GSEA-CW)

The GSEA target estimand (i.e., goal of inference) is

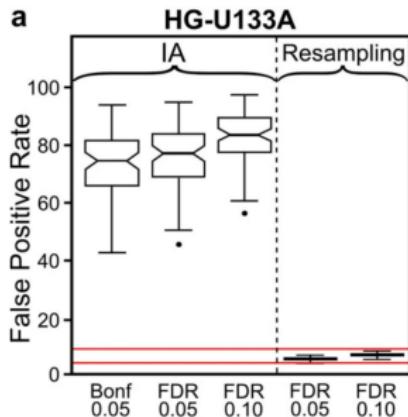
$$\phi_S = g(\theta)$$

The GSEA-CW target estimand is different, and implies a different scientific question

$$\psi_S = g(\theta^{\parallel})$$

DSA Methods that Account for Inter-gene/inter-taxa Correlations

Inter-gene/taxa correlations can massively inflate the false positive rate of GSEA with gene label permutations (Gatti et al., 2010)



Two methods that handle inter-gene/taxa correlations:

- ① GSEA with **sample label** permutations
- ② limma's CAMERA method

The CAMERA Method

Two methods that handle inter-gene/taxa correlations:

- ① GSEA with **sample label** permutations

Permute the sample labels (e.g., whether sample n came from healthy or tumor tissue), then re-estimate LFCs $\hat{\theta}$

The CAMERA Method

Two methods that handle inter-gene/taxa correlations:

- ① GSEA with **sample label** permutations

Permute the sample labels (e.g., whether sample n came from healthy or tumor tissue), then re-estimate LFCs $\hat{\theta}$

- ② limma's CAMERA method

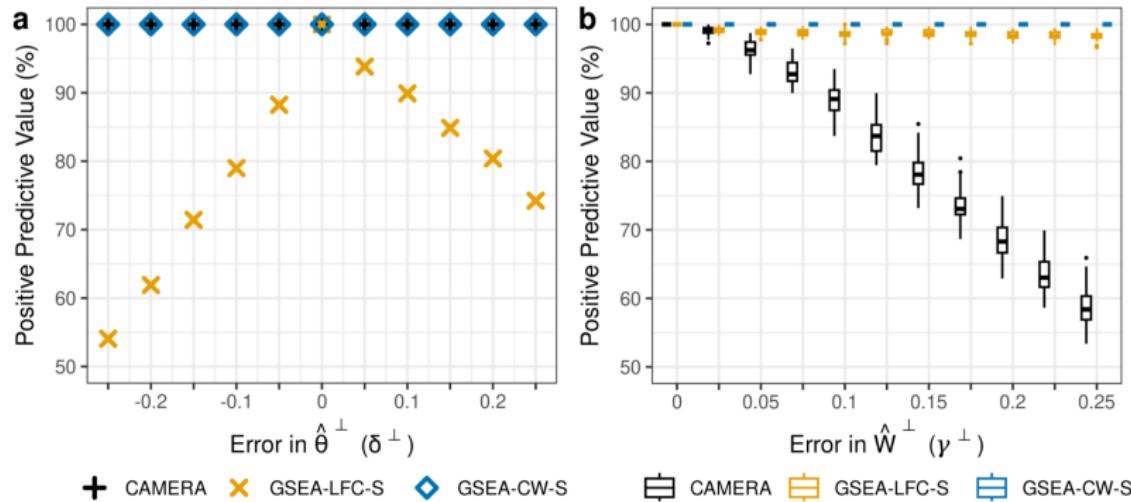
A two-sample t-test for the set S that uses a Variance Inflation Factor (VIF) to account for correlation:

$$\frac{\hat{\theta}_{\in S} - \hat{\theta}_{\notin S}}{s_p \sqrt{VIF/m_1 + 1/m_0}}$$

Scale Sensitivity Analyses

Here I consider two real data studies:

- ① Constant Error δ^\perp added to just samples from normal-adjacent-to-tumor
- ② Variable Error simulated from Uniform $[-\gamma^\perp, \gamma^\perp]$ added to all samples



* GSEA-LFC-S is GSEA with sample label permutations, GSEA-CW-S is GSEA with compositional weighting and sample label permutations

Future Directions

The LFC Sensitivity Test considers all possible errors $\epsilon^\perp \in (-\infty, \infty)$:

$$p = \sup_{\epsilon^\perp \in (-\infty, \infty)} p_{\epsilon^\perp}$$

Future Directions

The LFC Sensitivity Test considers all possible errors $\epsilon^\perp \in (-\infty, \infty)$:

$$p = \sup_{\epsilon^\perp \in (-\infty, \infty)} p_{\epsilon^\perp}$$

But why consider *all* possible errors? For instance $\theta^\perp = 10$ might imply total transcription is 22,000 times higher in tumor than healthy tissue!

Future Directions

The LFC Sensitivity Test considers all possible errors $\epsilon^\perp \in (-\infty, \infty)$:

$$p = \sup_{\epsilon^\perp \in (-\infty, \infty)} p_{\epsilon^\perp}$$

But why consider *all* possible errors? For instance $\theta^\perp = 10$ might imply total transcription is 22,000 times higher in tumor than healthy tissue!

Rather assume $\epsilon^\perp \in [-1, 1]$:

$$p = \sup_{\epsilon^\perp \in [-1, 1]} p_{\epsilon^\perp}$$

Future Directions

The LFC Sensitivity Test considers all possible errors $\epsilon^\perp \in (-\infty, \infty)$:

$$p = \sup_{\epsilon^\perp \in (-\infty, \infty)} p_{\epsilon^\perp}$$

But why consider *all* possible errors? For instance $\theta^\perp = 10$ might imply total transcription is 22,000 times higher in tumor than healthy tissue!

Rather assume $\epsilon^\perp \in [-1, 1]$:

$$p = \sup_{\epsilon^\perp \in [-1, 1]} p_{\epsilon^\perp}$$

Future Directions

In the previous presentation, SSRVs were a Bayesian approach to LFC estimation

$$\begin{aligned}\theta^\perp &= \hat{\theta}^\perp + \epsilon^\perp \\ \epsilon^\perp &\sim P\end{aligned}$$

Where P is a user-defined distribution of scale uncertainty.

Future Directions

In the previous presentation, SSRVs were a Bayesian approach to LFC estimation

$$\begin{aligned}\theta^\perp &= \hat{\theta}^\perp + \epsilon^\perp \\ \epsilon^\perp &\sim P\end{aligned}$$

Where P is a user-defined distribution of scale uncertainty.

What if we take a *purely Frequentist* approach? Rather than define P we just assume:

$$\theta^\perp \in [\hat{\theta}_L^\perp, \hat{\theta}_U^\perp]$$

Let $p_{\hat{\theta}^\perp}$ be the p-value under a single normalization assumption $\hat{\theta}^\perp$, we can define a new p-value:

$$p = \sup_{\theta^\perp \in [\hat{\theta}_L^\perp, \hat{\theta}_U^\perp]} p_{\hat{\theta}^\perp}$$

References

- ① Gatti, et al. Heading down the wrong pathway: on the influence of correlation within gene sets. *BMC Genomics.* 2010 Oct 18;11:574. doi: 10.1186/1471-2164-11-574.
- ② McGovern, et al. Addressing erroneous scale assumptions in microbe and gene set enrichment analysis. *PLoS Comput Biol.* 2023 Nov 20;19(11):e1011659. doi: 10.1371/journal.pcbi.1011659.
- ③ Subramanian, et al. Gene set enrichment analysis: a knowledge-based approach for interpreting genome-wide expression profiles. *Proc Natl Acad Sci U S A.* 2005 Oct 25;102(43):15545-50. doi: 10.1073/pnas.0506580102. Epub 2005 Sep 30.
- ④ Wu, et al. Camera: a competitive gene set test accounting for inter-gene correlation. *Nucleic Acids Res.* 2012 Sep 1;40(17):e133. doi: 10.1093/nar/gks461. Epub 2012 May 25.
- ⑤ Aran, et al. Comprehensive analysis of normal adjacent to tumor transcriptomes. *Nat Commun.* 2017 Oct 20;8(1):1077. doi: 10.1038/s41467-017-01027-z. PMID: 29057876; PMCID: PMC5651823.