# Scale Uncertainty in ALDEx2

Michelle Nixon

May 13, 2024
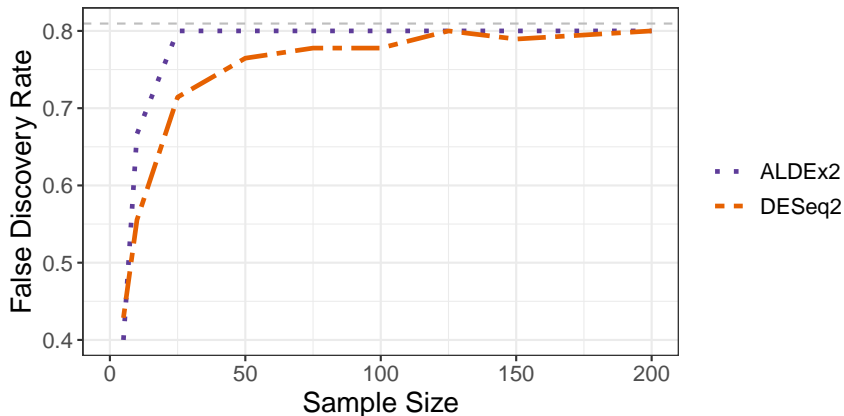
## Recap: Sequencing depth can confound conclusions.

| Observed data (Y) | Sample 1 | Sample 2 | Sample 3 |            |
|-------------------|----------|----------|----------|------------|
| Condition         | Health   | Health   | Disease  | Conclusion |
| Entity 1          | 5        | 10       | 100      | Increase   |
| Entity 2          | 10       | 25       | 3        | Decrease   |
| Entity 3          | 0        | 1        | 8        | Increase   |
| Entity 4          | 0        | 0        | 19       | Increase   |
| Sequencing Depth  | 15       | 36       | 130      |            |

## This can mislead analyses.

| System data (W) | Sample 1 | Sample 2 | Sample 3 | |
|---|---|---|---|---|
| Condition | Health | Health | Disease | Conclusion |
| Entity 1 | 227 | 351 | 154 | Decrease |
| Entity 2 | 684 | 891 | 3 | Decrease |
| Entity 3 | 48 | 32 | 15 | Decrease |
| Entity 4 | 43 | 39 | 27 | Decrease |
| Scale ($W^\perp$) | 1,002 | 1,313 | 200 | |

# . . . and lead to unacknowledged bias.

Section 1

## Problem Set-Up

# Observed Data as a Sample from the System



**Information Loss from System to Data:**

sampling

data processing

W

Y

$W^\perp$ = 14 (scale)

$Y^\perp$ = 4 (sequencing depth)

# Observed Data as a Sample from the System



Information Loss from System to Data:

W

sampling

data processing

Y

$W^{\perp} = 14$ (scale)

$Y^{\perp} = 4$ (sequencing depth)

non-scaled research questions + analyses

A
B
C

Relative abundance

scaled research questions + analyses

A
B
C

Absolute abundance

## Notation

- **Y**: a measurement of the underlying system **W**.

$$\mathbf{W}_{dn} = \underbrace{\mathbf{W}_{dn}^{\parallel}}_{\text{composition}} \times \underbrace{W_n^{\perp}}_{\text{scale}}$$

## Notation

- **Y**: a measurement of the underlying system $W$.

$$\mathbf{W}_{dn} = \underbrace{\mathbf{W}_{dn}^{\parallel}}_{\text{composition}} \times \underbrace{W_n^{\perp}}_{\text{scale}}$$

- **Compostion:** $\mathbf{W}_{dn}^{\parallel} = \frac{\mathbf{W}_{dn}}{\sum_{d=1}^{D} \mathbf{W}_{dn}}$

- **Scale:** $W_n^{\perp} = \sum_{d=1}^{D} \mathbf{W}_{dn}$

## Notation

- **Y**: a measurement of the underlying system $W$.

$$\mathbf{W}_{dn} = \underbrace{\mathbf{W}_{dn}^{\parallel}}_{\text{composition}} \times \underbrace{W_n^{\perp}}_{\text{scale}}$$

- **Compostion:** $\mathbf{W}_{dn}^{\parallel} = \frac{\mathbf{W}_{dn}}{\sum_{d=1}^{D} \mathbf{W}_{dn}}$

- **Scale:** $W_n^{\perp} = \sum_{d=1}^{D} \mathbf{W}_{dn}$

- $\boldsymbol{\theta}$: what we want to estimate.

## Example: Notation

| System data ($W^{\parallel}$) | Sample 1 | Sample 2 | Sample 3 |
|---|---|---|---|
| Condition | Health | Health | Disease |
| Entity 1 | 0.27 | 0.27 | 0.77 |
| Entity 2 | 0.68 | 0.68 | 0.02 |
| Entity 3 | 0.05 | 0.02 | 0.08 |
| Entity 4 | 0.04 | 0.03 | 0.13 |
| Scale ($W^{\perp}$) | 1,002 | 1,313 | 200 |

Differential Abundance/Expression Analysis

- **Question:** How do entities (e.g., taxa or genes) change between conditions?

- In this case, $\theta$ is the log-fold change (LFC):

$$\theta_d = \text{mean}_{\text{case}}(\log \mathbf{W}_{dn}) - \text{mean}_{\text{control}}(\log \mathbf{W}_{dn})$$

## The Original ALDEx2 Model

**Step 1: Model Sampling Uncertainty**

$$\mathbf{Y}_{\cdot n} \sim \text{Multinomial}(\mathbf{W}_{\cdot n}^{\parallel})$$

$$\mathbf{W}_{\cdot n}^{\parallel} \sim \text{Dirichlet}(\alpha)$$

**Step 2: Centered Log-Ratio Transformation**

$$\log \mathbf{W}_{\cdot n} = \left[ \log \mathbf{W}_{1n}^{\parallel} - \text{mean}(\log \mathbf{W}_{\cdot n}^{\parallel}), ..., \log \mathbf{W}_{Dn}^{\parallel} - \text{mean}(\log \mathbf{W}_{\cdot n}^{\parallel}) \right]$$

**Step 3: Calculate LFCs and Test if Different from Zero.**

$$\theta_d = \text{mean}_{\text{case}}(\log \mathbf{W}_{dn}) - \text{mean}_{\text{control}}(\log \mathbf{W}_{dn})$$
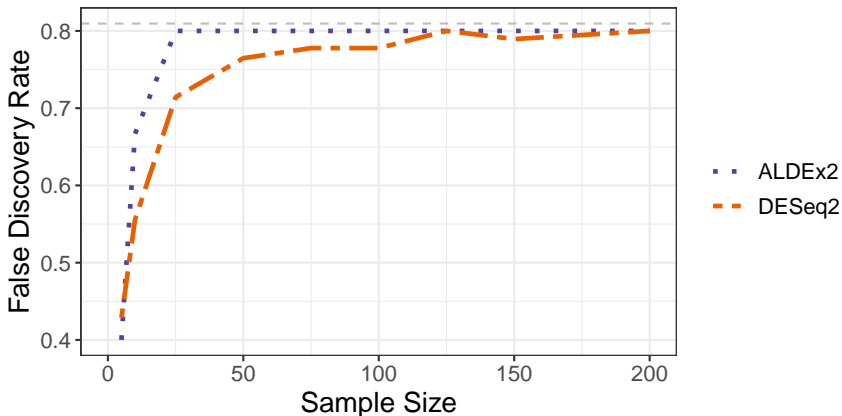
## Implied Assumptions about Scale

**Step 1: Model Sampling Uncertainty**

$$\mathbf{Y}_{\cdot n} \sim \text{Multinomial}(\mathbf{W}^{\|}_{\cdot n})$$
$$\mathbf{W}^{\|}_{\cdot n} \sim \text{Dirichlet}(\alpha)$$

**Step 2: Centered Log-Ratio Transformation**

$$\log \mathbf{W}_{\cdot n} = \left[\log \mathbf{W}^{\|}_{1n} - \text{mean}(\log \mathbf{W}^{\|}_{\cdot n}), ..., \log \mathbf{W}^{\|}_{Dn} - \text{mean}(\log \mathbf{W}^{\|}_{\cdot n})\right]$$

**Step 3: Calculate LFCs and Test if Different from Zero.**

$$\theta_d = \text{mean}_{\text{case}}(\log \mathbf{W}_{dn}) - \text{mean}_{\text{control}}(\log \mathbf{W}_{dn})$$

## Implied Assumptions about Scale, cont.

Using the relationship $\mathbf{W}_{dn} = \mathbf{W}^{\parallel}_{dn} W^{\perp}_n$ and some math, the CLR normalization implies:
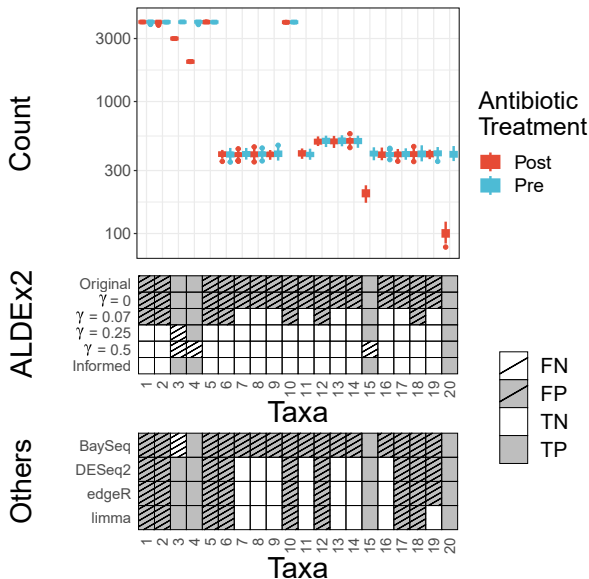
$$\log W^{\perp}_n = -\text{mean}(\log \mathbf{W}^{\parallel}_{\cdot n}).$$

What does this mean?

# Unacknowledged bias!

# Adding Uncertainty in Scale can Help.

Section 2

Scale Reliant Inference

## Scale Reliant Inference: The Basics

- **The CoDA perspective:** Research questions that depend on $W^{\perp}$ (scale) are not possible.

- **The Normalization perspective:** Research questions that depend on $W^{\perp}$ (scale) can be answered after normalization.

- Who is right?

## Scale Reliant Inference: The Basics

- **The CoDA perspective:** Research questions that depend on $W^{\perp}$ (scale) are not possible.

- **The Normalization perspective:** Research questions that depend on $W^{\perp}$ (scale) can be answered after normalization.

- Who is right?

- **The CoDA perspective:** Yes, but this is limiting in practice.

- **The Normalization perspective:** Not correct, but attempting to answer relevant questions.

## Scale Reliant Inference: The Basics

- What happens if $\theta$ depends on $W^\perp$?

- Consider LFCs: how are taxa changing between two conditions?

$$
\begin{aligned}
\theta_d &= \text{mean}_{\text{case}}(\log \mathbf{W}_{dn}) - \text{mean}_{\text{control}}(\log \mathbf{W}_{dn}) \\
&= ... \\
&= \underbrace{\text{mean}_{\text{case}}(\log \mathbf{W}_{dn}^{\parallel}) - \text{mean}_{\text{control}}(\log \mathbf{W}_{dn}^{\parallel})}_{\theta^{\parallel}} \\
&\quad - \underbrace{\text{mean}_{\text{case}}(\log W_n^{\perp}) - \text{mean}_{\text{control}}(\log W_n^{\perp})}_{\theta^{\perp}}
\end{aligned}
$$

Don't we need $\theta^{\perp}$?

## Scale Reliant Inference: Theory Intro

Recall for LFCs:

$$\theta_d = \text{mean}_{\text{case}}(\log \mathbf{W}_{dn}) - \text{mean}_{\text{control}}(\log \mathbf{W}_{dn})$$
$$= \theta^{\|} + \theta^{\perp}$$

- What can we say about $\theta$ from $\theta^{\|}$ alone?

## Scale Reliant Inference: Theory Intro

Recall for LFCs:

$$\theta_d = \text{mean}_{\text{case}}(\log \mathbf{W}_{dn}) - \text{mean}_{\text{control}}(\log \mathbf{W}_{dn})$$
$$= \theta^{\parallel} + \theta^{\perp}$$

- What can we say about $\theta$ from $\theta^{\parallel}$ alone?

- Statistical perspective: $\theta$ is not identifiable without $\theta^{\perp}$.

- Practical issues: unbiased estimators, calibrated confidence sets, and type-I error control **NOT** possible!

- See Nixon et al. (2023) for details.

## Scale Simulation Random Variables

**Goal:** Estimate $\theta = f(\mathbf{W}^{\parallel}, W^{\perp})$.

1. Draw samples of $\mathbf{W}^{\parallel}$ from a measurement model (can depend on $\mathbf{Y}$).

2. Draw samples of $W^{\perp}$ from a scale model (can depend on $\mathbf{W}^{\parallel}$).

3. Estimate samples of $\theta = f(\mathbf{W}^{\parallel}, W^{\perp})$.

Section 3

# The Updated ALDEx2 Software

## ALDEx2 as an SSRV

**Step 1: Model Sampling Uncertainty**

$$\mathbf{Y}_{.n} \sim \text{Multinomial}(\mathbf{W}^{\parallel}_{.n})$$
$$\mathbf{W}^{\parallel}_{.n} \sim \text{Dirichlet}(\alpha)$$

**Step 2: Draw Samples from a Scale Model**

$$\log W^{\perp}_{n} = -\text{mean}(\log \mathbf{W}^{\parallel}_{.n}) + \epsilon,\ \epsilon \sim N(0, \gamma^2)$$
$$\log \mathbf{W}_{.n} = \log \mathbf{W}^{\parallel}_{.n} + \log W^{\perp}_{n}$$

**Step 3: Calculate LFCs and Test if Different from Zero.**

$$\theta_d = \text{mean}_{\text{case}}(\log \mathbf{W}_{dn}) - \text{mean}_{\text{control}}(\log \mathbf{W}_{dn})$$

# Benefits of Moving Past Normalizations to Scale

## Intro to Scale Models

Normalizations are replaced by a scale model:

$$\log W_n^{\perp} = -\text{mean}(\log \mathbf{W}_{\cdot n}^{\parallel}) + \epsilon$$
$$\epsilon \sim N(0, \gamma^2)$$

What about other options?

## Intro to Scale Models, cont.

There are no restrictions on what scale models can be, although there are some helpful options:

1. Based on normalizations. (Stochastic normalizations)
2. Based on biological knowledge.
3. Based on outside measurements.

## Scale Models based on Biological Knowledge

What do past studies or biological mechanisms tell about the scale of the system?

## Scale Models based on Biological Knowledge

What do past studies or biological mechanisms tell about the scale of the system?

1. You are confident that taking an antibiotic will kill at least some microbes in the gut.

2. A past study showed that a certain disease (e.g., Crohn's disease) leads to lower microbial load in the gut.

3. You believe the total microbial load in the mouth changes after brushing your teeth.

This type of information can be used in scale model building.

## Scale Models based on Outside Measurements

How can outside measurements be used to quantify scale?

## Scale Models based on Outside Measurements

How can outside measurements be used to quantify scale?

1. These measurements can be used *if* they relate to your scale of interest.

2. Examples include flow cytometry, qPCR, etc.

3. Scale models can incorporate measurement uncertainty.

Section 4

## Coding Changes to ALDEx2

## Including scale

**The new ALDEx2 model removes normalizations in lieu of scale models.**

## Including scale

**The new ALDEx2 model removes normalizations in lieu of scale models.**

Major updates:

1. A new argument `gamma` which makes it easy to incorporate scale uncertainty.

2. A new function `aldex.senAnalysis` to see how analysis results change as a function of scale uncertainty.

## The gamma argument

- Added as argument to the aldex and aldex.clr function.

- gamma can either be a single numeric or a matrix.

    1. Single numeric: controls the noise on the default scale model.
    2. Matrix: A $N \times S$ matrix of samples of $W^{\perp}$.

- gamma = NULL returns the original behavior of ALDEx2.

## Option 1: Default Scale Model

The default scale model is based on errors in the CLR normalization.

$$\log \hat{W}_n^{\perp(s)} = -\text{mean}\left(\log \hat{W}_{\cdot n}^{\parallel(s)}\right) + \Lambda^{\perp} x_n$$

$$\Lambda^{\perp} \sim \ N(0, \gamma^2).$$

## Advantages of the Default Scale Model

1. It is built off the status quo for ALDEx2.

2. Any value of $\gamma > 0$ will reduce false positives compared to the CLR normalization.

3. It has a concrete interpretation to contextualize scale assumptions.

## Interpreting the Default Scale Model

$$\log \hat{\mathcal{W}}_n^{\perp(s)} = -\text{mean}\left(\log \hat{\mathcal{W}}_{\cdot n}^{\parallel(s)}\right) + \Lambda^{\perp} x_n$$

$$\Lambda^{\perp} \sim N(0, \gamma^2).$$

**Empirical Rule:** 95% of the samples of $\Lambda^{\perp}$ fall within a factor of $\pm 2\gamma$ from zero.

## Interpreting the Default Scale Model, cont.

$$\log \hat{W}_n^{\perp(s)} = -\text{mean}\left(\log \hat{W}_{.n}^{\parallel(s)}\right) + \Lambda^\perp x_n$$

$$\Lambda^\perp \sim N(0, \gamma^2).$$

**For case/control experiments:**

1. If $x_n = 1$: 95% of samples of $\log \hat{W}_n^{\perp(s)}$ fall within a factor of $\pm 2\gamma$ of the negative geometric mean.

2. If $x_n = 0$: $\log \hat{W}_n^{\perp(s)}$ is equal to the negative geometric mean.

## Interpreting the Default Scale Model, cont.

Recall that with the CLR normalization:

$$\log W_n^{\perp} = -\text{mean}(\log \mathbf{W}_{\cdot n}^{\|}) = -\text{GM}(\mathbf{W}_{\cdot n}^{\|}).$$

Thus, when using the CLR normalization:

$$\theta^{\perp} = \text{mean}_{\text{case}}(-\text{GM}(\mathbf{W}_{\cdot n}^{\|})) - \text{mean}_{\text{control}}(-\text{GM}(\mathbf{W}_{\cdot n}^{\|}))$$

This is same mean that the default scale model is centered on.

## Interpreting the Default Scale Model, cont.

Taken together, the default scale model implies that:

1. The value of $\theta^{\perp}$ is within $\pm 2\gamma$ of the value of $\theta^{\perp}_{\mathsf{CLR}}$ implied by the CLR normalization.

2. With 95% certainty, the true difference in scales falls within the the range $2^{\theta^{\perp}_{\mathsf{CLR}} \pm 2\gamma}$.

## Option 2: More Complex Scale Models

Alternatively, can pass a matrix of scale samples to gamma so long as:

1. The dimension is $N \times S$.
2. They are samples of $W^{\perp}$ not log $W^{\perp}$.

Reasons to do this:

1. **Biological beliefs:** Scale is guided by the biological system or the researcher's prior beliefs.

2. **Outside Measurements:** These can be used in building a scale model *if* they are informative on the scale of interest (e.g., qPCR, flow cytometry).

## Sensitivity Analyses

- Recall that the default scale model has a parameter $\gamma$ controlling the amount of noise added.

- Instead of picking $\gamma$, why not test over a range instead?

- Enter sensitivity analyses.

## Sensitivity Analyses

**Step 1: Model Sampling Uncertainty**

$$\mathbf{Y}_{\cdot n} \sim \text{Multinomial}(\mathbf{W}_{\cdot n}^{\parallel})$$
$$\mathbf{W}_{\cdot n}^{\parallel} \sim \text{Dirichlet}(\alpha)$$

**Step 2: Draw Samples from a Scale Model** For a given $\gamma$:

$$\log W_n^{\perp,\gamma} = -\text{mean}\left(\log \hat{W}_{\cdot n}^{\parallel(s)}\right) + \Lambda^{\perp} x_n$$
$$\Lambda^{\perp} \sim \ N(0, \gamma^2)$$
$$\log \mathbf{W}^{\gamma}_{\cdot n} = \log \mathbf{W}_{\cdot n}^{\parallel} + \log W_n^{\perp,\gamma}$$

**Step 3: Calculate LFCs and Test if Different from Zero.**

**Step 4: Repeat for all desired values of $\gamma$.**

Section 5

# Data Examples

## Simulation Study

Consider a simple study of the microbiome pre/post antibiotic administration.

- **Research question:** Which taxa change in absolute abundance after taking an antibiotic?

- 100 study participants, 50 in each condition (pre/post antibiotics).

- 20 taxa total with 4 taxa truly changing (decreasing)

## Data

## Adding Scale is Easy

```
## Adding noise via the default scale model
mod.ss.high <- aldex(Y, conds, gamma = 0.5)
```

## Investigating Assumptions about Scale

```
## Looking at the implied scale
clr <- aldex.clr(Y, conds, gamma = 0.001)
clr@scaleSamps[1:6, 1:4]

##            [,1]     [,2]     [,3]     [,4]
## [1,] 5.174279 5.124890 5.199780 5.175163
## [2,] 5.175705 5.144470 5.184953 5.167715
## [3,] 5.178751 5.171188 5.130795 5.100749
## [4,] 5.158594 5.195139 5.164371 5.145696
## [5,] 5.120674 5.175533 5.189581 5.171154
## [6,] 5.208741 5.273464 5.207085 5.162631
```

## Investigating Assumptions about Scale, cont.

# Scale Model based on Biology

```
## Creating an informed model using biological
↪   reasoning
scales <- c(rep(1, 50), rep(0.9, 50))
scale_samps <- aldex.makeScaleMatrix(gamma = 0.15, mu
↪   = scales,
    conditions = conds, log = FALSE)

mod.know <- aldex(Y, conds, gamma = scale_samps)
```

## Scale Model based on Outside Measurements

```
flow_data_collapse <- flow_data %>%
    group_by(sample) %>%
    mutate(mean = mean(flow)) %>%
    mutate(stdev = sd(flow)) %>%
    dplyr::select(-flow) %>%
    ungroup() %>%
    unique()
scale_samps <- matrix(NA, nrow =
↪  nrow(flow_data_collapse), ncol = 128)
for (i in 1:nrow(scale_samps)) {
    scale_samps[i, ] <- rnorm(n = 128, mean =
↪  flow_data_collapse$mean[i],
        sd = flow_data_collapse$stdev[i])
}
mod.flow <- aldex(Y, conds, gamma = scale_samps)
```
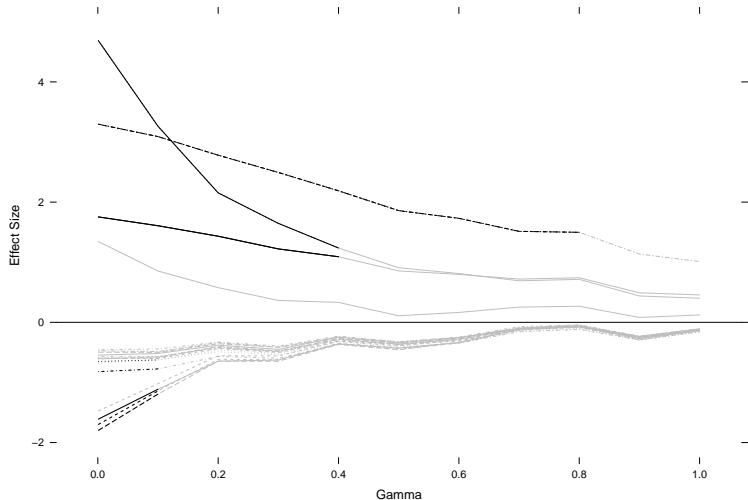
# Plotting Results

## Sensitivity Analyses

```
## First, specifying different values for the noise
↪  in the
## scale
gamma_to_test <- c(0.001, seq(0.1, 1, by = 0.1))

## Run the CLR function
clr <- aldex.clr(Y, conds)

## Run sensitivity analysis function
sen_res <- aldex.senAnalysis(clr, gamma =
↪  gamma_to_test)
plotGamma(sen_res, thresh = 0.1, blackWhite = TRUE,
↪  taxa_to_label = 3)
```
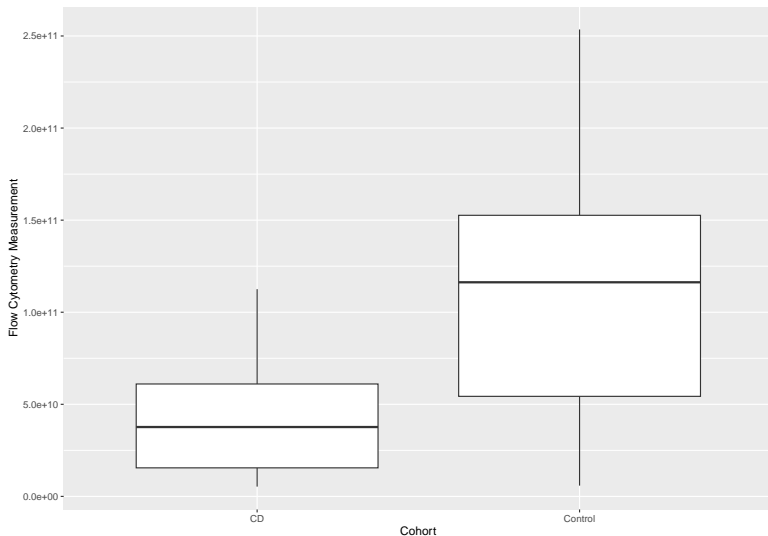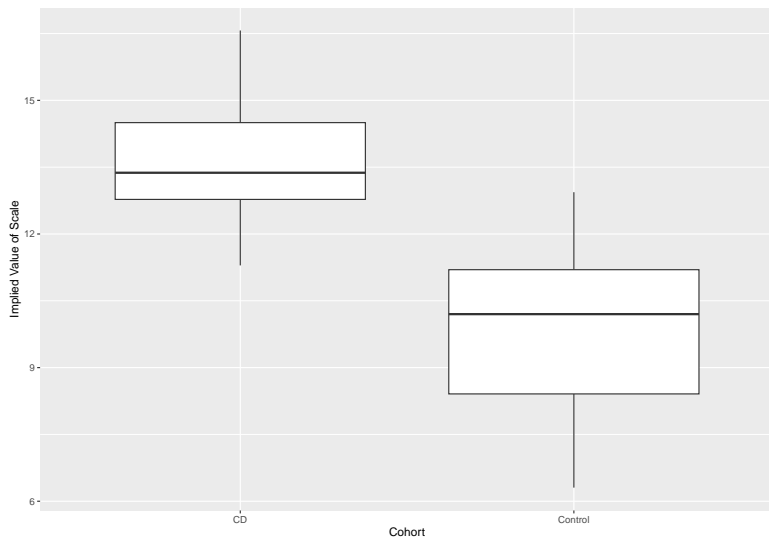
# Sensitivity Analyses, cont.

## Real Example: Vandputte

1. Comparison study of 29 Crohn's disease patients and 66 healthy controls.

2. For each patient, they sequenced the fecal sample and obtained flow cytometry measurements.

3. Proposed an approach that supplemented sequence count data with flow cytometry measurements.

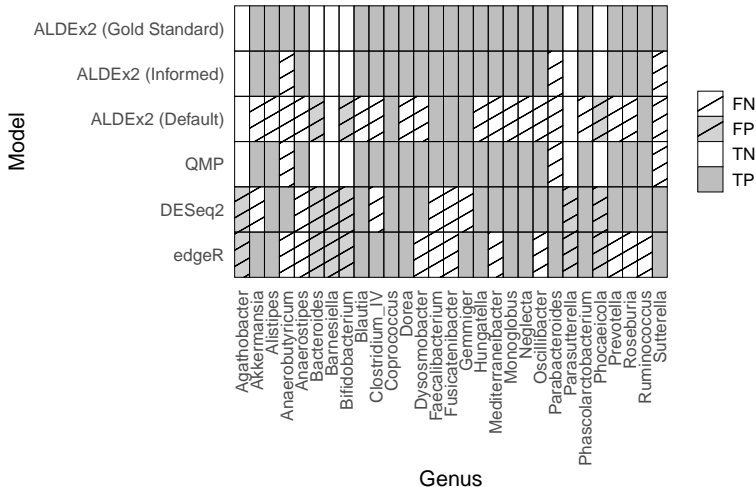## Difference in Scale Implied by Flow Cytometry

## Difference in Scale Implied by CLR

## Creating a Gold Standard Model

```r
scale_mean <- log2(sample_data(phylo)$CellCount)
scale_var <- rep(0.7, 95)

scale_samples <- matrix(NA, nrow = 95, ncol = 1000)
for (i in 1:95) {
    scale_samples[i, ] <- 2^rnorm(1000,
↪   scale_mean[i], scale_var[i])
}
```
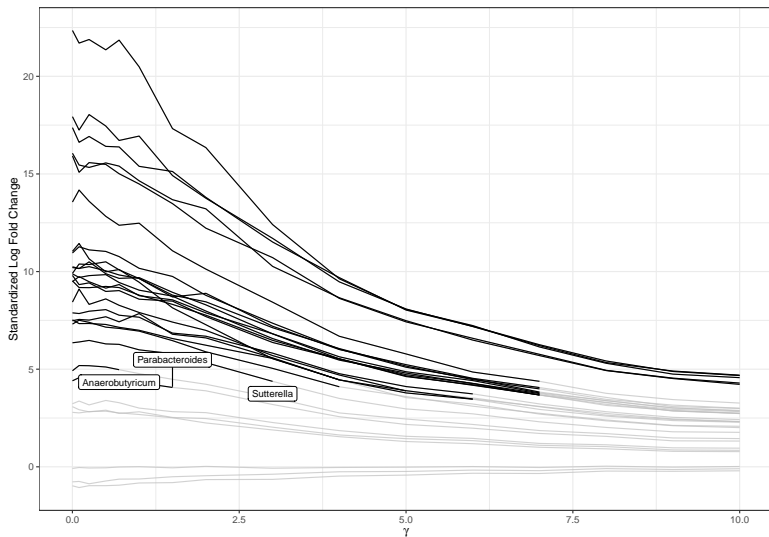
## Creating an Informed Model

```r
scale.cd <- 2^matrix(rnorm(1000 * 29, mean =
↪ log2(0.7), sd = 0.125),
    nrow = 29)
scale.control <- 2^matrix(rnorm(1000 * 66, mean =
↪ log2(1), sd = 0.125),
    nrow = 66)

scale.inf <- rbind(scale.cd, scale.control)
aldex_informed <- aldex(Y, X, mc.samples = 1000,
↪ gamma = scale.inf)
```

## Comparing to Other Methods

## Sensitivity Analyses

## References

**Scale Reliant Inference/Updates to ALDEx2:**

- Nixon, et. al. (2023) "Scale Reliant Inference." *ArXiv Preprint 2201.03616*.

- Gloor, Nixon, and Silverman. (2023) "Scale is Not What You Think; Explicit Scale Simulation in ALDEx2." *BioRXiv Preprint 2023.10.21.563431*.

- Nixon, Gloor, and Silverman. (2024) "Beyond Normalizations: Incorporating Scale Uncertainty in ALDEx2." *BioRXiv Preprint 2024.04.01.587602*.

- Fernandes et. al. (2014). "Unifying the analysis of high-throughput sequencing datasets: characterizing RNA-seq, 16S rRNA gene sequencing and selective growth experiments by compositional data analysis." *Microbiome*.

## References

**Data Sources:**

- McMurrough et. al. (2014)."Control of catalytic efficiency by a coevolving network of catalytic and noncatalytic residues." *PNAS*.

- Vandputte et. al. (2017). "Quantitative microbiome profiling links gut community variation to microbial load." *Nature*.