# FRAMING

# WHAT IS SEQUENCE COUNT DATA?

Multivariate count data **Y**<sub>ij</sub> representing the number of transcripts of type **j** sequenced in sample **i** 

## **Examples:**

- 16s rRNA sequencing
- RNA-seq (± Single Cell)
- T-cell receptor sequencing
- [hand waving] tumor subtypes

## **Extended Applications**

[Beyond Sequencing]:

- Multiparametric Flow
   Cytometry
- Political Polling







Adapted from Hamady. et al., *Nature Methods*, 2008





## **KEY POINT**

Sequencing depth does not seem to correlate with microbial load.





% Blue % Orange % Green







# **COUNTING INTRODUCES UNCERTAINTY**



# **COUNTING INTRODUCES UNCERTAINTY**



## THE PROBLEM WITH FEW COUNTS OR SMALL COUNTS



## THE PROBLEM WITH FEW COUNTS OR SMALL COUNTS





# MICROBIOME DATA IS SPARSE



Silverman, et al., <u>eLife</u> 2017

## TAKE HOME MESSAGES

- Counts (with many zero or small counts)
- Compositional Information Only
- The data is multivariate not univariate!
- Lots of variation (technical and biological)
  - Variation is not always random (then we call it "bias")

# MODELING

## **GENERATIVE MODELING**

Sequencing



## $Y_i \sim \text{Multinomial}(n_i, \pi_i)$

?

DNA Extraction PCR Amplification



Sample Collection and Storage



Adapted from Hamady. et al., *Nature Methods*, 2008

The Aitchison Geometry in the Simplex is the relevant space to model our systems in

- The Aitchison Geometry in the Simplex is the relevant space to model our systems in
- "Our Systems Multiply" + The information in our data is relative

#### **GROUP THEORY / VECTOR SPACE**

- The Aitchison Geometry in the Simplex is the relevant space to model our systems in
- "Our Systems Multiply" + The information in our data is relative
- Conclusions should be drawn from [Log]-Ratios



INTUITIVE

- The Aitchison Geometry in the Simplex is the relevant space to model our systems in
- "Our Systems Multiply" + The information in our data is relative
- Conclusions should be drawn from [Log]-Ratios
- The Logistic-Normal Distribution is the CLT for our unobserved system(s)



- The Aitchison Geometry in the Simplex is the relevant space to model our systems in
- "Our Systems Multiply" + The information in our data is relative
- Conclusions should be drawn from [Log]-Ratios
- The Logistic-Normal Distribution is the CLT for our unobserved system(s)
- All methods for analyzing relative data should adhere to three principles (1) Scale Invariance,
   (2) Permutation Invariance, (3)
   Subcompositional Coherence





## **GENERATIVE MODELING (LIKELIHOOD ONLY)**

Sequencing



 $Y_i \sim \text{Multinomial}(n_i, \pi_i)$ 

DNA Extraction PCR Amplification



Sample Collection and Storage



 $\mu_i \sim \text{Logistic Normal}(\alpha, V)$ 

 $\pi_i \sim \text{Logistic Normal}(\mu_i, V)$ 

Adapted from Hamady. et al., Nature Methods, 2008

## WHAT IS MISSING?

Sequencing



DNA Extraction PCR Amplification



Sample Collection and Storage



Assign Sequences to Samples

| >AGTGAGAGAAGCAGGGTCGTAATGTT |   |   |
|-----------------------------|---|---|
| >AGTGCGATGCGTAGGGTCGTAATGCG |   |   |
| >AGTGCGATGCGTAGGGTCGTAATG7A | - |   |
| >AGTGGATGCTCTAGGGTCGTAATGCA |   |   |
| >AGTGTCACGGTGAGGGTCGTAATGGG | - | - |
| >AGTGGATGCTCTAGGGTCGTAATGTT |   |   |
| >AGTGTCACGGTGAGGGTCGTAATGCC |   |   |
| >AGTGAGAGAAGCAGGGTCGTAATCAC |   |   |
|                             |   |   |

Denoise Reads or Cluster

|          | Species 1 | Species 2 | Species 3 | Sp |
|----------|-----------|-----------|-----------|----|
| Sample 1 | 23        | 53        | 2         |    |
| Sample 2 | 69        | 64        | 70        |    |
| Sample 3 | 33        | 100       | 68        |    |
| Sample 4 | 5         | 63        | 57        |    |
| Sample 5 | 76        | 80        | 46        |    |
| Sample 6 | 58        | 7         | 37        |    |
| Sample 7 | 10        | 87        | 32        |    |
| comple 0 | 01        | 00        | 70        |    |

Adapted from Hamady. et al., *Nature Methods*, 2008

## MULTINOMIAL-LOGISTIC NORMAL (OR NORMAL ON THE SIMPLEX)

```
Y \sim \text{Multinomial}(\pi)
\pi \sim \text{Logistic Normal}(\rho, \Xi)
Y \sim \text{Multinomial}(\pi)
\pi = ILR^{-1}(\eta)
\eta \sim \text{Multivariate Normal}(\mu, \Sigma)
```

## **ISOMETRIC LOGRATIO TRANSFORM – AN ORTHONORMAL BASIS**

### **UNOBSERVED ABSOLUTE ABUNDANCES**



## **ISOMETRIC LOGRATIO TRANSFORM – AN ORTHONORMAL BASIS**

### **UNOBSERVED ABSOLUTE ABUNDANCES**





# ISOMETRIC LOGRATIO TRANSFORM

## **ORTHONORMAL BASIS IN SIMPLEX**



## DATA PROJECTED ONTO BASIS



Silverman, et al., <u>eLife</u> 2017

## PHYLOGENIC ISOMETRIC LOGRATIO (PHILR) TRANSFORM

## **PHYLOGENETIC BALANCES**



## **ORTHONORMAL BASIS IN SIMPLEX**



## DATA PROJECTED ONTO BASIS



Silverman, et al., <u>eLife</u> 2017

## WHY AN ORTHONORMAL BASIS? CURRENT STATISTICAL STANDARD "IDENTIFIED SOFTMAX" (AKA INVERSE ALR)



## WHY AN ORTHONORMAL BASIS? CURRENT STATISTICAL STANDARD "IDENTIFIED SOFTMAX" (AKA INVERSE ALR)



## WHY AN ORTHONORMAL BASIS? CURRENT STATISTICAL STANDARD "IDENTIFIED SOFTMAX" (AKA INVERSE ALR)



# WHY AN ORTHONORMAL BASIS?

## **ORTHONORMAL BASIS**



# WHY AN ORTHONORMAL BASIS?

## **ORTHONORMAL BASIS**



## WHY AN ORTHONORMAL BASIS?

- > Ability to rotate your view requires orthonormal basis.
- Interpretability of low-dimensional projections requires orthonormal basis
- "Objects change when you look at them differently with non-orthonormal bases"
- "Units" Require an orthonormal Basis (Evidence Information)

## **PHILR BASIS**



## **VARIATION ARRAY**



## **PRINCIPLE BALANCE ANALYSIS**



## MANUAL CURATION



# **EXAMPLE APPLICATIONS**

True State with Biological Noise



True State with Biological Noise



Addition of Technical Noise

True State with Biological Noise



Addition of Technical Noise

True State with Biological Noise











Observed Counts
$$\boldsymbol{Y}_t \sim \operatorname{Multinomial}(\boldsymbol{\pi}_t)$$
 $\uparrow$  $\boldsymbol{\pi}_t = \operatorname{ILR}^{-1}(\boldsymbol{\eta}_t)$ Addition of Technical Noise $\boldsymbol{\eta}_t = \boldsymbol{F}'_t \boldsymbol{\theta}_t + \boldsymbol{\nu}_t$  $\uparrow$  $\boldsymbol{\eta}_t = \boldsymbol{F}'_t \boldsymbol{\theta}_t + \boldsymbol{\nu}_t$ True State with Biological Noise $\boldsymbol{\theta}_t = \boldsymbol{G}_t \boldsymbol{\theta}_{t-1} + \boldsymbol{\omega}_t$  $\boldsymbol{\omega}_t \sim N(0, \boldsymbol{W}_t)$ 

## **A SIMPLE SIMULATION**

$$Y_t \sim \text{Multinomial}(\pi_t)$$
  

$$\pi_t = \text{ILR}^{-1}(\eta_t)$$
  

$$\eta_t = \mu_t + v_t \qquad v_t \sim N(0, V)$$
  

$$\mu_t = \mu_{t-1} + \omega_t \qquad \omega_t \sim N(0, W)$$

## **A SIMPLE SIMULATION**



$$Y_t \sim \text{Multinomial}(\pi_t)$$
  

$$\pi_t = \text{ILR}^{-1}(\eta_t)$$
  

$$\eta_t = \mu_t + v_t \qquad v_t \sim N(0, V)$$
  

$$\mu_t = \mu_{t-1} + \omega_t \qquad \omega_t \sim N(0, W)$$

Posterior Estimate for Eta with 95% Credible Interval







#### **STANDARD LONGITUDINAL MODEL**

 $\begin{aligned} Y_t &\sim \mathsf{Multinomial}(\pi_t) \\ \pi_t &= \mathsf{ILR}^{-1}(\eta_t) \\ \eta_t &= \mu_t + v_t \qquad v_t \sim \mathit{N}(\mathsf{0}, \mathit{V}) \\ \mu_t &= \mu_{t-1} + \omega_t \qquad \omega_t \sim \mathit{N}(\mathsf{0}, \mathit{W}) \end{aligned}$ 



#### **STANDARD LONGITUDINAL MODEL**

 $\begin{aligned} Y_t &\sim \mathsf{Multinomial}(\pi_t) \\ \pi_t &= \mathsf{ILR}^{-1}(\eta_t) \\ \eta_t &= \mu_t + v_t \qquad v_t \sim \mathit{N}(\mathsf{0}, \mathit{V}) \\ \mu_t &= \mu_{t-1} + \omega_t \qquad \omega_t \sim \mathit{N}(\mathsf{0}, \mathit{W}) \end{aligned}$ 



## ESTIMATING "SIGNAL-TO-NOISE" RATIO

**Biological Noise to Technical Noise Ratio** 

$$\frac{\operatorname{Tr}(W)}{\operatorname{Tr}(V)}$$

Percent of Noise (Excluding Counting) Attributable to Biology Tr(W) $\overline{Tr(V) + Tr(W)}$ 

## ESTIMATING "SIGNAL-TO-NOISE" RATIO

**Biological Noise to Technical Noise Ratio** 

 $\frac{\operatorname{Tr}(W)}{\operatorname{Tr}(V)}$ 

Percent of Noise (Excluding Counting) Attributable to Biology Tr(W) $\overline{Tr(V) + Tr(W)}$ 

#### **STANDARD LONGITUDINAL MODEL**

 $\begin{aligned} Y_t &\sim \mathsf{Multinomial}(\pi_t) \\ \pi_t &= \mathsf{ILR}^{-1}(\eta_t) \\ \eta_t &= \mu_t + v_t \\ \mu_t &= \mu_{t-1} + \omega_t \end{aligned} \qquad \begin{aligned} v_t &\sim \mathsf{N}(\mathsf{0}, \mathsf{V}) \\ \omega_t &\sim \mathsf{N}(\mathsf{0}, \mathsf{W}_t) \end{aligned}$ 

#### **CONDITION TO HANDLE REPLICATES**

$$W_t = \begin{cases} 0 & \text{if } t \text{ is a replicate of } t-1; \\ W & \text{otherwise;} \end{cases}$$

## DOMINATED BY TECHNICAL NOISE AT HOURLY INTERVALS

**Biological Noise to Technical Noise Ratio** 

 $\frac{\operatorname{Tr}(W)}{\operatorname{Tr}(V)}$ 

Percent of Noise (Excluding Counting) Attributable to Biology Tr(W) $\overline{Tr(V) + Tr(W)}$ 





## **BIOLOGICAL AND TECHNICAL VARIATION CAN HAVE DIFFERENT SHAPES**

