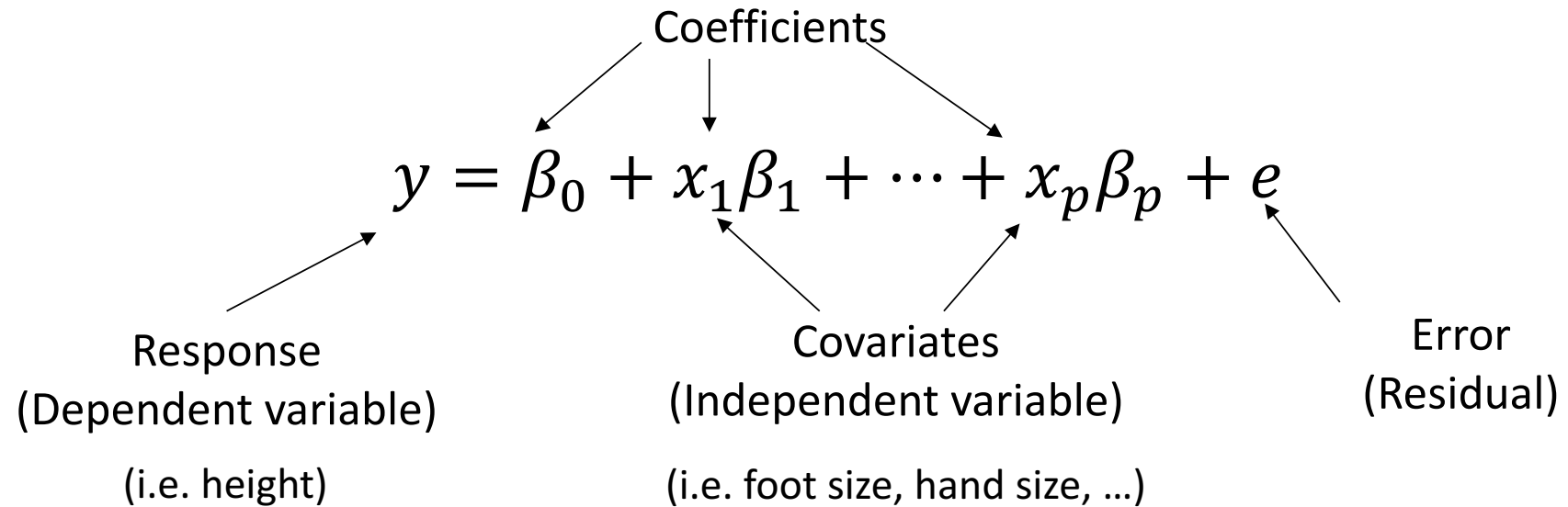


Insights into Cystic Fibrosis

Jamie Morton, University of California, San Diego

Review on Regression



Review on Regression

$$\begin{array}{ccccc} Y & = & X & B & + & e \\ \begin{array}{|c|} \hline y_0 \\ \hline y_1 \\ \hline \dots \\ \hline y_n \\ \hline \end{array} & = & \begin{array}{|c|c|c|c|} \hline 1 & x_{11} & \dots & x_{1p} \\ \hline \dots & \dots & \dots & \dots \\ \hline 1 & x_{n1} & \dots & x_{np} \\ \hline \end{array} & \begin{array}{|c|} \hline \beta_0 \\ \hline \beta_1 \\ \hline \dots \\ \hline \beta_p \\ \hline \end{array} & + & \begin{array}{|c|} \hline e_0 \\ \hline e_1 \\ \hline \dots \\ \hline e_n \\ \hline \end{array} \\ n \times 1 & & n \times p & p \times 1 & & n \times 1 \\ \text{Response} & & \text{Covariates} & \text{Coefficients} & & \text{Residuals} \end{array}$$

n = number of measurements

p = number of variables measured

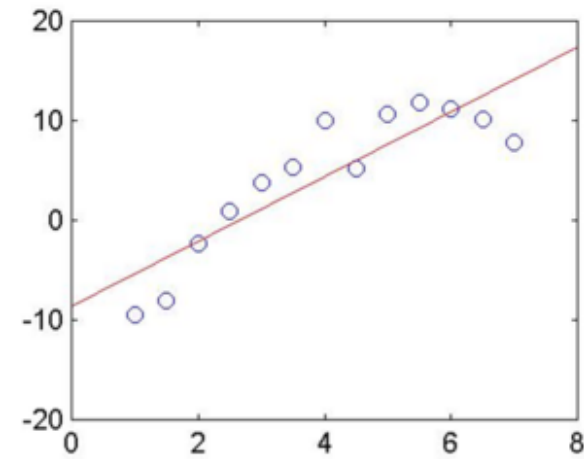
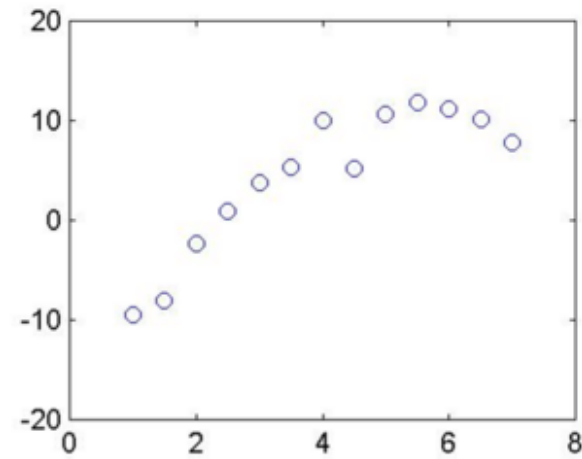
Assumptions

- Over-determined system ($n \gg p$)
- Independence between measurements
- Compositionality

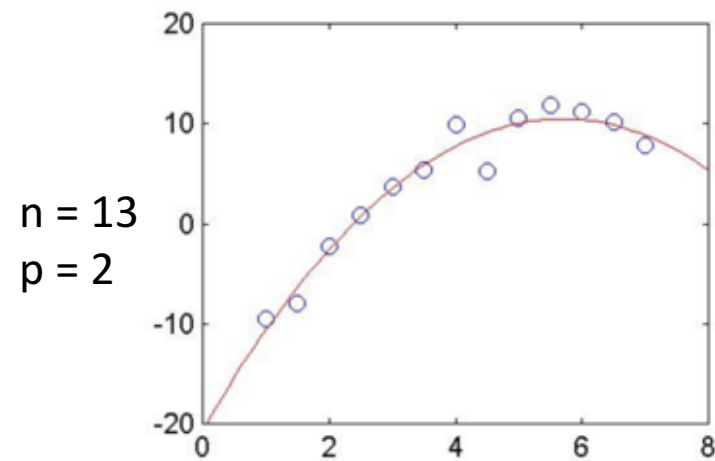
Assumptions

- Over-determined system ($n \gg p$)
- Independence between measurements
- Compositionality

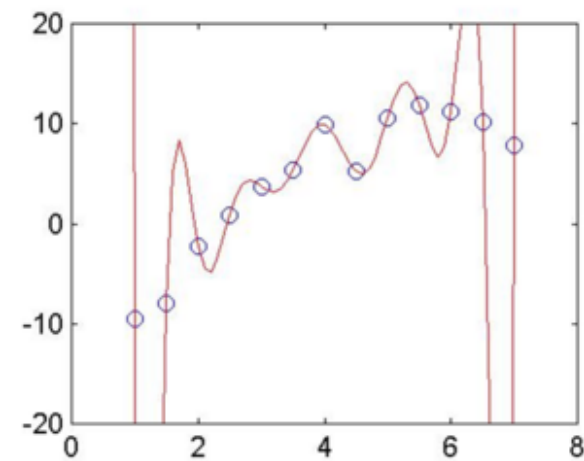
Occams Razor



$n = 13$
 $p = 1$



$n = 13$
 $p = 2$



$n = 13$
 $p = 13$

High dimensional variables

- Genetic data
 - Thousands of variables
 - Hundreds of samples
- Two solutions
 1. Regularization – wisely choose a subset of variables
 2. Multivariate response

High dimensional variables

- Genetic data
 - Thousands of variables
 - Hundreds of samples
- Two solutions
 1. Regularization – wisely choose a subset of variables
 2. Multivariate response

Univariate response

$$Y = X B$$

y_0
y_1
...
y_n

1	x_{11}	...	x_{1p}
...	
1	x_{n1}	...	x_{np}

β_0
β_1
...
β_p

$n \times 1$

Response

(BMI, age, sex, ...)

$n \times p$

Covariates

(gene abundances)

$p \times 1$

Coefficients

n = number of measurements

p = number of variables measured

Only one variable at a time!

Multivariate response

$$Y = XB$$

y_{00}	y_{01}	...	y_{0D}
y_{10}	y_{11}	...	y_{1D}
...
y_{n0}	y_{n1}	...	y_{nD}

$n \times 1$

Response
(gene abundances)

1	x_{11}	...	x_{1p}
...
1	x_{n1}	...	x_{np}

$n \times p$

Covariates
(BMI, age, sex, ...)

β_{00}	β_{01}	...	β_{0D}
β_{10}	β_{11}	...	β_{1D}
...
β_{p0}	β_{p1}	...	β_{pD}

$p \times D$

Coefficients

n = number of measurements

p = number of covariates measured

D = number of variables measured

Can encode categorical variables

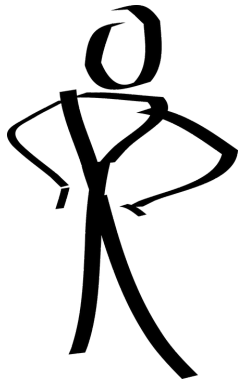
Advantages

- Effectively avoids over parameterization
 - $n \gg p$
 - But always try to run cross validation to confirm
 - Train model on subset of measurements
 - Try to predict the remaining measurements
- Build models with many covariates

Assumptions

- Over-determined system ($n \gg p$)
- Independence between measurements
- Compositionality

Independence



...



...



Time 1

Time 2

Time t

- Independence is violated
- Samples depend on person
 - Drawn from different distributions

Linear Mixed Effects Models

$$Y = X\beta + e \quad \text{Linear Regression}$$

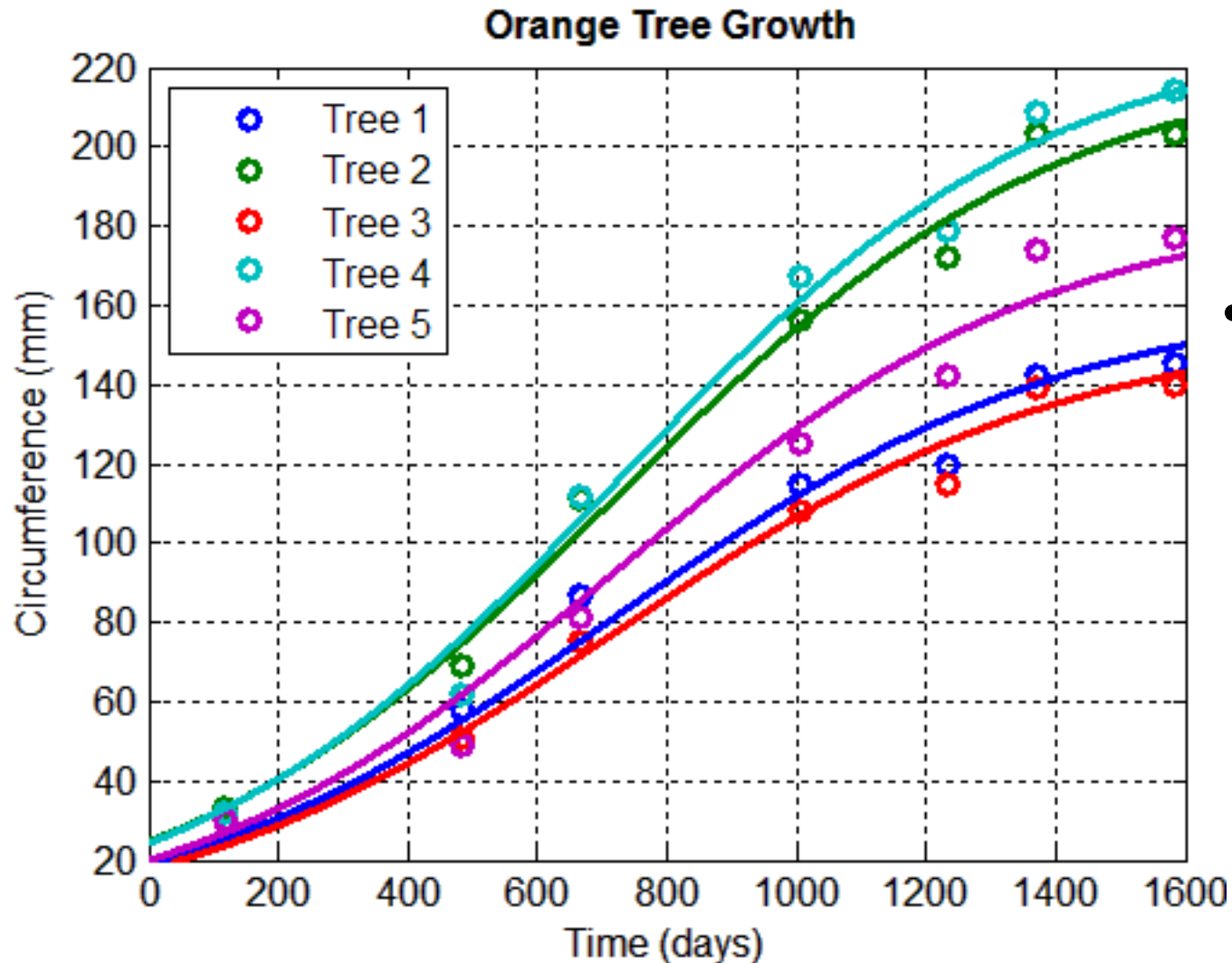
$$Y = X\beta + Z\mu + e \quad \text{Linear Mixed Effects Model}$$

Fixed effects
(i.e. sex)

Random effects
(i.e. test score)

Fixed effect: constant value
Random effect: value varies

Linear Mixed Effects Models



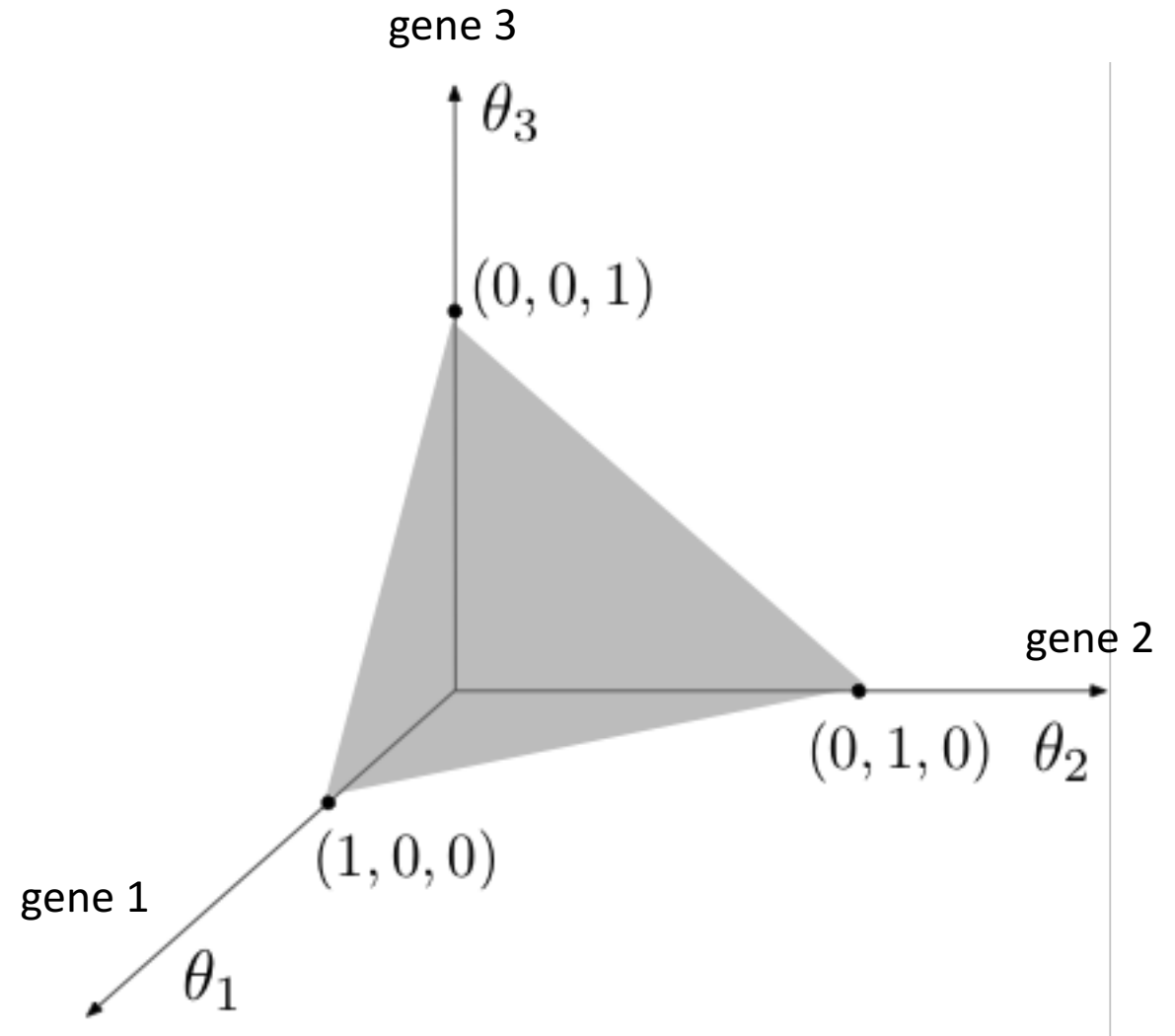
- Fit multiple regressions
 - One per individual

Assumptions

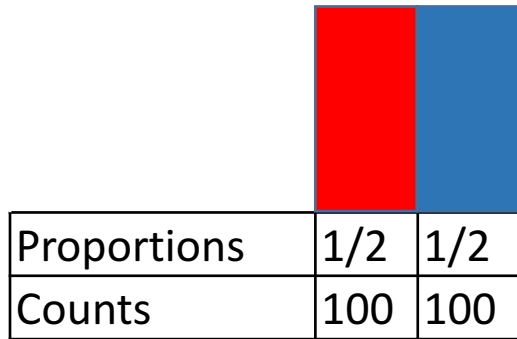
- Over-determined system ($n \gg p$)
- Independence between measurements
- Compositionality

Compositionality

- We're dealing with proportions
- Thus all of our samples live in the simplex

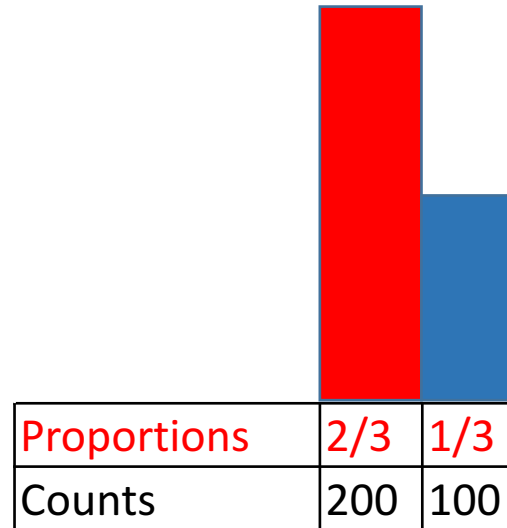


Compositionality

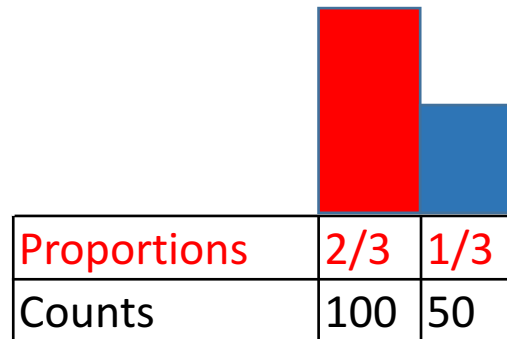
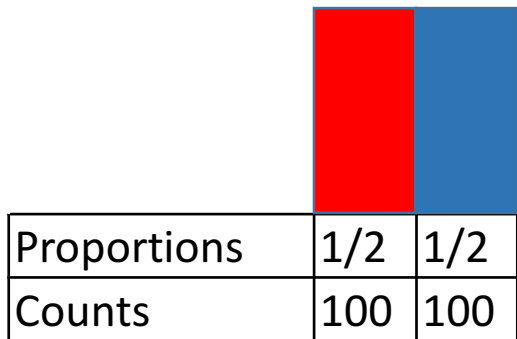


Time point 1

Red doubled



Blue halved

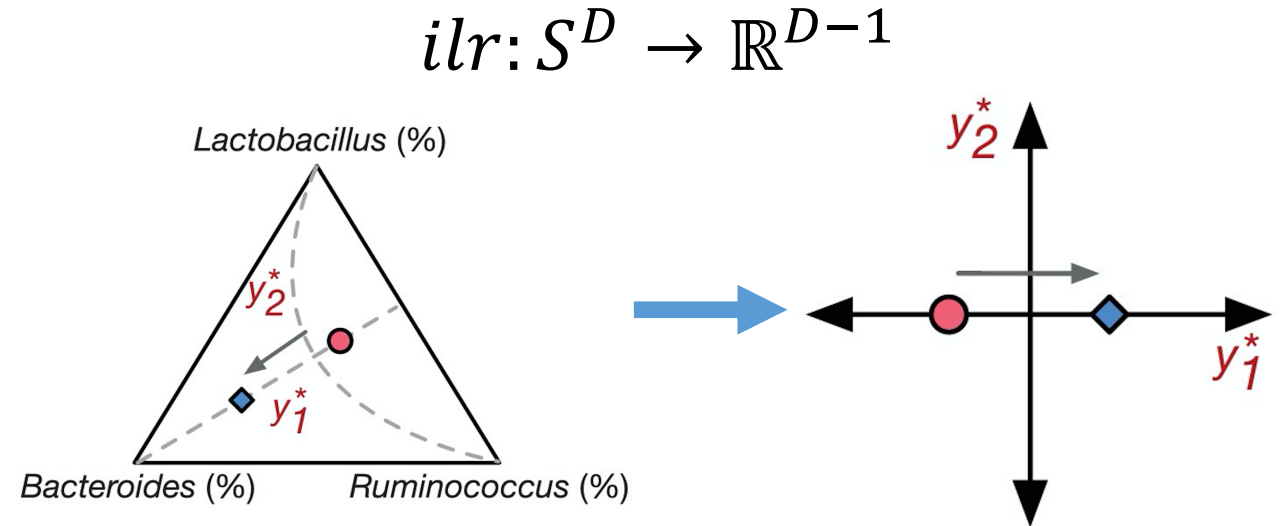
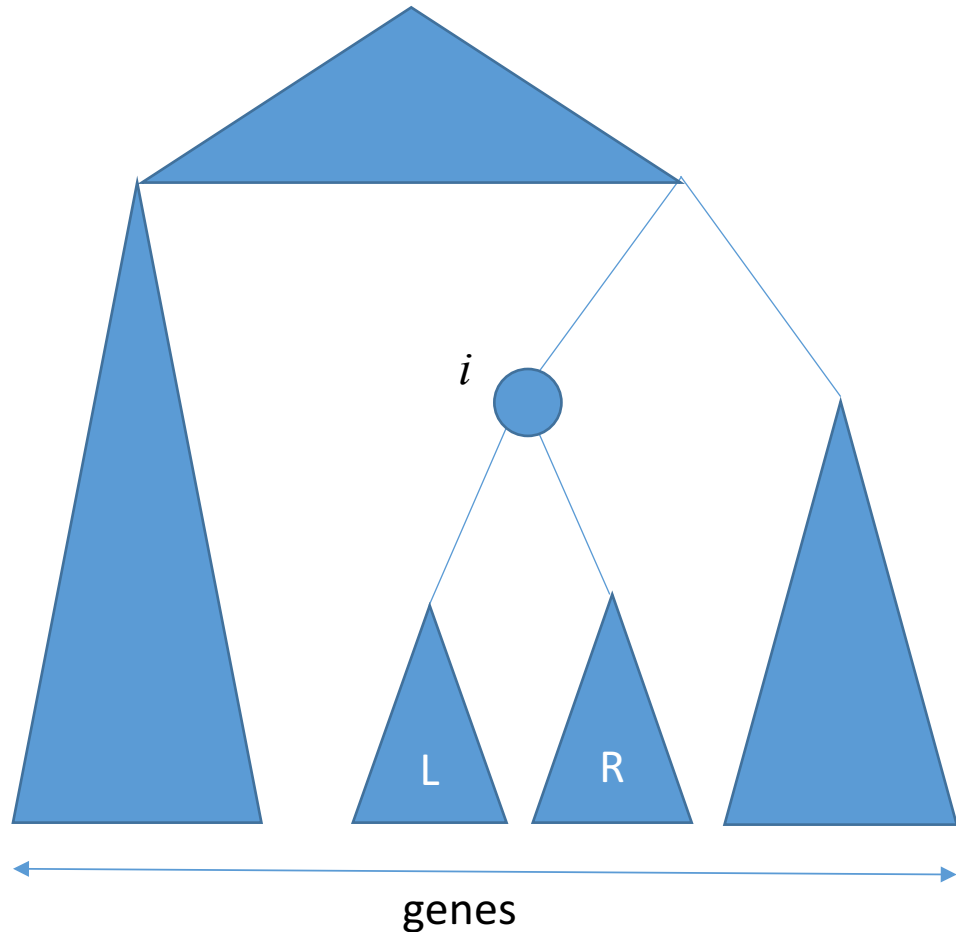


Time point 2

This breaks regression

- Cannot predict who's growing/dying
- Solutions must be restricted to the simplex

The isometric log ratio transform



$$b_i = \sqrt{\frac{|i_L||i_R|}{|i_L| + |i_R|}} \log \left(\frac{g(i_L)}{g(i_R)} \right)$$

$$g(x) = \sqrt[k]{x_1 \dots x_k} \quad \text{Geometric mean}$$

Recap

- Multivariate response resolves **Occam's razor**
- Linear mixed effects models resolves **sample dependence issue**
- ILR transform resolves **compositionality problem**

Lesson Plan

- Perform ILR transform
 - Build tree based on pH hierarchical clustering
 - i.e. group organisms that live in similar pH together
- Build Multivariate Response Linear Mixed Effects Model

$$Y = X\beta + Z\mu + e$$

Microbial proportions pH Patient id

Software



<https://biocore.github.io/gneiss/>



<https://qiime2.org/>



<http://scikit-bio.org/>