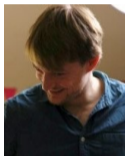


A Walk Along Models for Count Data in Microbial Ecology



M. Mariadassou, INRA-MaIAGE



joint work with Julien Chiquet and Stéphane Robin

RCAM'18, Athens, 2018, September 9th



Julien Chiquet, M.M., Stéphane Robin,
Variational inference for probabilistic Poisson PCA
<https://arxiv.org/abs/1703.06633> (in press at *Annals of Applied Statistics*)



PLNmodels package, development version on github
`devtools::install_github("jchiquet/PLNmodels", build_vignettes=TRUE)`



- 1 Motivation
- 2 Multinomial Models
- 3 Log-Normal Models
- 4 Applications





Data from [MBE⁺15].

- $n = 155$ samples (= 31 piglets at 5 times)
- $p = 1038$ bacterial species (OTUs) with prevalence ≥ 0.05
- Some covariates (sex, sire, etc)
- Offsets: $o_i =$ offset for sample i (sequencing depth)



Data from [MBE⁺15].

- $n = 155$ samples (= 31 piglets at 5 times)
- $p = 1038$ bacterial species (OTUs) with prevalence ≥ 0.05
- Some covariates (sex, sire, etc)
- Offsets: $o_i =$ offset for sample i (sequencing depth)

Aim: Study impact of weaning on gut microbiota

A look at the data

Metabarcoding data from [MBE⁺15]

- **count** matrix with $n = 155$ piglets, $p = 1038$ species

```
mach_counts[1:2, c(3, 9, 12, 15)]
```

```
##          5982 347 349 5854
## SF0901    0  23  3    0
## SF0902    8   0  4    0
```

- $d = 8$ **covariates** (sex, mother, weaning status, ...)

```
mach_covariates[1:2, ]
```

```
##          Run Project Time Bande sex      mere Weaned
## SF0901    3 Kinetic  D14  1105   1 17MAG101814  TRUE
## SF0902    3 Kinetic  D36  1105   1 17MAG101814  FALSE
```

- **Sampling effort** in each sample

```
mach_offsets[1:2, c(1:4, 48:51)]
```

```
##          16342  164 5982 5980 10413 6307 8949  346
## SF0901    3084 3084 3084 3084  3084 3084 3084 3084
## SF0902    2182 2182 2182 2182  2182 2182 2182 2182
```



Data from [JFS⁺16].

- $n = 116$ oak leaves = samples
- $p = 114$ microbial species
 - $p_1 = 66$ bacterial species (OTUs, based on the 16S)
 - $p_2 = 48$ fungal species (OTUs, based on the ITS)
- covariates: tree (resistant, intermediate, susceptible), height, distance to trunk, ...
- offsets: $o_{i1} \neq o_{i2}$ = offset for bacteria, fungi



Data from [JFS⁺16].

- $n = 116$ oak leaves = samples
- $p = 114$ microbial species
 - $p_1 = 66$ bacterial species (OTUs, based on the 16S)
 - $p_2 = 48$ fungal species (OTUs, based on the ITS)
- covariates: tree (resistant, intermediate, susceptible), height, distance to trunk, ...
- offsets: $o_{i1} \neq o_{i2}$ = offset for bacteria, fungi

```
offsets[1:2, c(1:4, 48:51)]
```

```
##      f_1  f_2  f_3  f_4 E_alphitoides b_1045 b_109 b_1093
## [1,] 2488 2488 2488 2488          2488   8315  8315   8315
## [2,] 2054 2054 2054 2054          2054    662   662    662
```




Data from [JFS⁺16].

- $n = 116$ oak leaves = samples
- $p = 114$ microbial species
 - $p_1 = 66$ bacterial species (OTUs, based on the 16S)
 - $p_2 = 48$ fungal species (OTUs, based on the ITS)
- covariates: tree (resistant, intermediate, susceptible), height, distance to trunk, ...
- offsets: $o_{i1} \neq o_{i2}$ = offset for bacteria, fungi

```
offsets[1:2, c(1:4, 48:51)]
```

```
##      f_1  f_2  f_3  f_4 E_alphitoides b_1045 b_109 b_1093
## [1,] 2488 2488 2488 2488          2488   8315  8315   8315
## [2,] 2054 2054 2054 2054          2054    662   662    662
```

Aim. Understand the interaction between the species, including the oak mildew pathogene *E. alphitoides*.

Data tables: $\mathbf{Y} = (Y_{ij}), n \times p$; $\mathbf{X} = (X_{ik}), n \times d$; $\mathbf{O} = (O_{ij}), n \times p$

- Y_{ij} = abundance (read counts) of species j in sample i
- X_{ik} = value of covariate k in sample i
- O_{ij} = offset (sampling effort) for species j in sample i

Data tables: $\mathbf{Y} = (Y_{ij}), n \times p$; $\mathbf{X} = (X_{ik}), n \times d$; $\mathbf{O} = (O_{ij}), n \times p$

- Y_{ij} = abundance (read counts) of species j in sample i
- X_{ik} = value of covariate k in sample i
- O_{ij} = offset (sampling effort) for species j in sample i

Need for multivariate analysis to help deciphering the ecosystem

- exhibit **patterns of diversity**
↪ summarize the information from \mathbf{Y} (PCA, clustering, ...)
- understand **between-species interactions**
↪ 'Network' inference (variable/covariance selection)
- correct for technical and **confounding effects**
↪ account for covariables and sampling effort

Data tables: $\mathbf{Y} = (Y_{ij}), n \times p$; $\mathbf{X} = (X_{ik}), n \times d$; $\mathbf{O} = (O_{ij}), n \times p$

- Y_{ij} = abundance (read counts) of species j in sample i
- X_{ik} = value of covariate k in sample i
- O_{ij} = offset (sampling effort) for species j in sample i

Need for multivariate analysis to help deciphering the ecosystem

- exhibit **patterns of diversity**
↪ summarize the information from \mathbf{Y} (PCA, clustering, ...)
- understand **between-species interactions**
↪ 'Network' inference (variable/covariance selection)
- correct for technical and **confounding effects**
↪ account for covariables and sampling effort

↪ need a generic framework to **model dependencies between count variables**

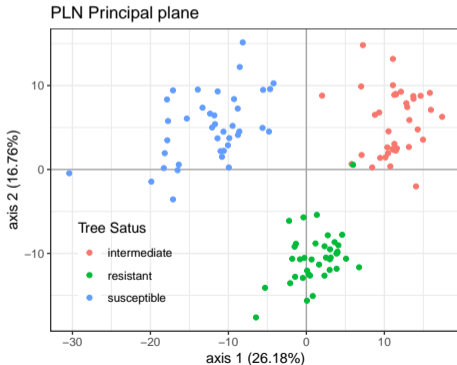
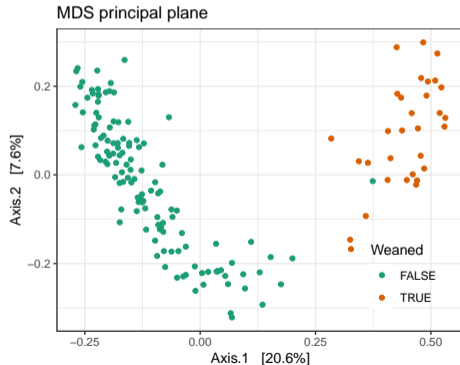
- 1 Apply your favorite **distance** (Jaccard, Bray-Curtis, UniFrac, weighted UniFrac, etc)

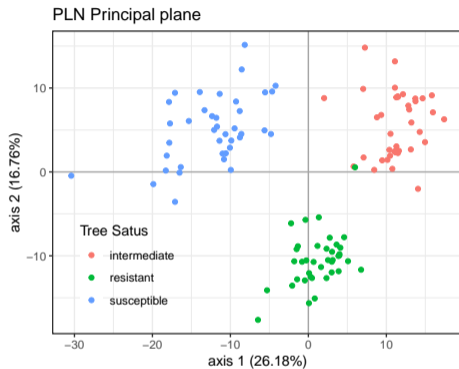
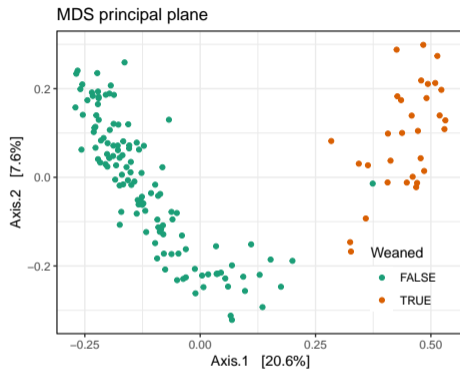
- 1 Apply your favorite **distance** (Jaccard, Bray-Curtis, UniFrac, weighted UniFrac, etc)
- 2 Apply your favorite **dimension reduction** technique (PCA, MDS/PCoA, NMDS, RDA, PLN, etc)

- 1 Apply your favorite **distance** (Jaccard, Bray-Curtis, UniFrac, weighted UniFrac, etc)
- 2 Apply your favorite **dimension reduction** technique (PCA, MDS/PCoA, NMDS, RDA, PLN, etc)
- 3 Plot resulting *graph*

Microbial Ecology 101

- 1 Apply your favorite **distance** (Jaccard, Bray-Curtis, UniFrac, weighted UniFrac, etc)
- 2 Apply your favorite **dimension reduction** technique (PCA, MDS/PCoA, NMDS, RDA, PLN, etc)
- 3 Plot resulting *graph*
- 4 *Et voilà!*



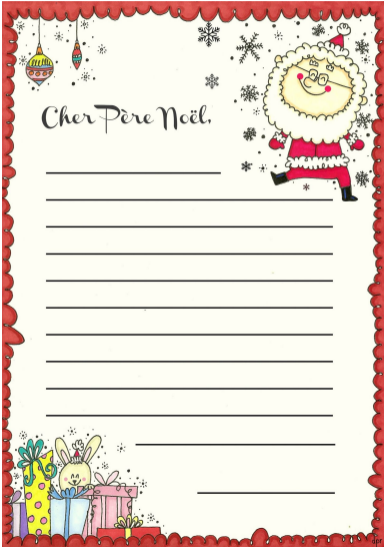


- 1 Perfect for **finding** structure...
- 2 But not for **modeling** it

What kind of generic models?

What kind of generic framework for multivariate count data?

What kind of generic models?



My Wish List to Santa

We want a family of **generative** models that are:

We want a family of **generative** models that are:

- **Flexible** enough to:
 - model average communities;
 - model dispersion (biological variability);
 - model interaction between OTUs (ecological networks);
 - accomodate heterogeneous communities;
 - integrate data from different sources (bacterial and fractions)

We want a family of **generative** models that are:

- **Flexible** enough to:
 - model average communities;
 - model dispersion (biological variability);
 - model interaction between OTUs (ecological networks);
 - accommodate heterogeneous communities;
 - integrate data from different sources (bacterial and fractions)
- yet as **parcimonious** as possible;

We want a family of **generative** models that are:

- **Flexible** enough to:
 - model average communities;
 - model dispersion (biological variability);
 - model interaction between OTUs (ecological networks);
 - accomodate heterogeneous communities;
 - integrate data from different sources (bacterial and fractions)
- yet as **parcimonious** as possible;
- **interpretable**;

We want a family of **generative** models that are:

- **Flexible** enough to:
 - model average communities;
 - model dispersion (biological variability);
 - model interaction between OTUs (ecological networks);
 - accomodate heterogeneous communities;
 - integrate data from different sources (bacterial and fractions)
- yet as **parcimonious** as possible;
- **interpretable**;
- **fast and easy** to fit to data;

We want a family of **generative** models that are:

- **Flexible** enough to:
 - model average communities;
 - model dispersion (biological variability);
 - model interaction between OTUs (ecological networks);
 - accomodate heterogeneous communities;
 - integrate data from different sources (bacterial and fractions)
- yet as **parcimonious** as possible;
- **interpretable**;
- **fast and easy** to fit to data;
- **good fits** to data (e.g. simulate **realistic** samples).

1 Motivation

2 Multinomial Models

- Multinomial
- Mixture of Multinomials
- (Mixture of) Dirichlet-Multinomial
- Latent Dirichlet Allocation

3 Log-Normal Models

4 Applications

1 Motivation

2 Multinomial Models

- **Multinomial**
- Mixture of Multinomials
- (Mixture of) Dirichlet-Multinomial
- Latent Dirichlet Allocation

3 Log-Normal Models

4 Applications



Intuition

- There are p species with proportions $\boldsymbol{\pi} = (\pi_1, \dots, \pi_p)$ in the species
- You pick N (sequencing depths) individuals with replacement

Multinomial Models

Intuition

- There are p species with proportions $\boldsymbol{\pi} = (\pi_1, \dots, \pi_p)$ in the species
- You pick N (sequencing depths) individuals with replacement

Mathematical Model

$$\mathbf{Y} \sim \mathcal{M}(N, \boldsymbol{\pi})$$

Multinomial Models

Intuition

- There are p species with proportions $\boldsymbol{\pi} = (\pi_1, \dots, \pi_p)$ in the species
- You pick N (sequencing depths) individuals with replacement

Mathematical Model

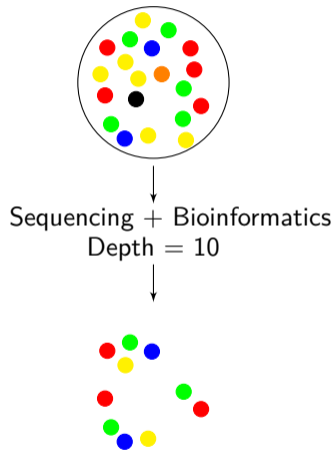
$$\mathbf{Y} \sim \mathcal{M}(N, \boldsymbol{\pi})$$

Inference is easy

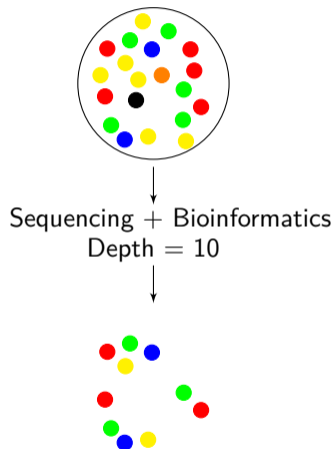
$$\hat{\pi}_j = \frac{\sum_{i=1}^n Y_{ij}}{\sum_{i=1}^n N_i}$$

with Y_{ij} the abundance of species j in sample i and N_i the depth of sample i .

Multinomial distribution: draw balls (with replacement) from a box

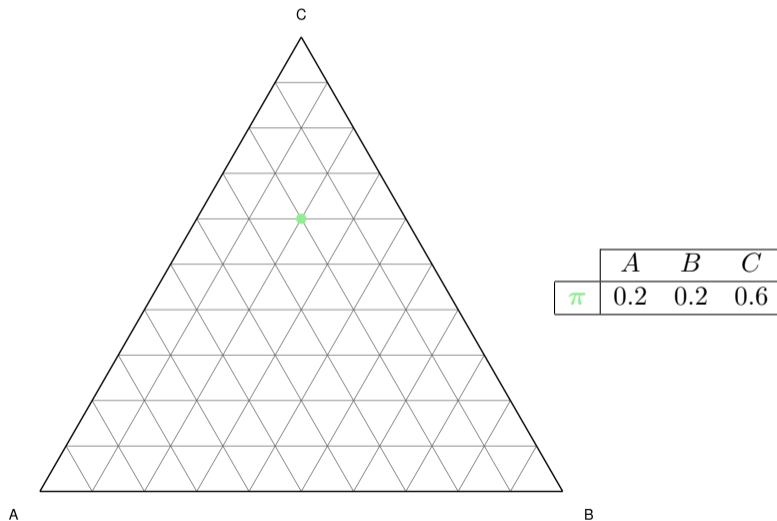


Multinomial distribution: draw balls (with replacement) from a box

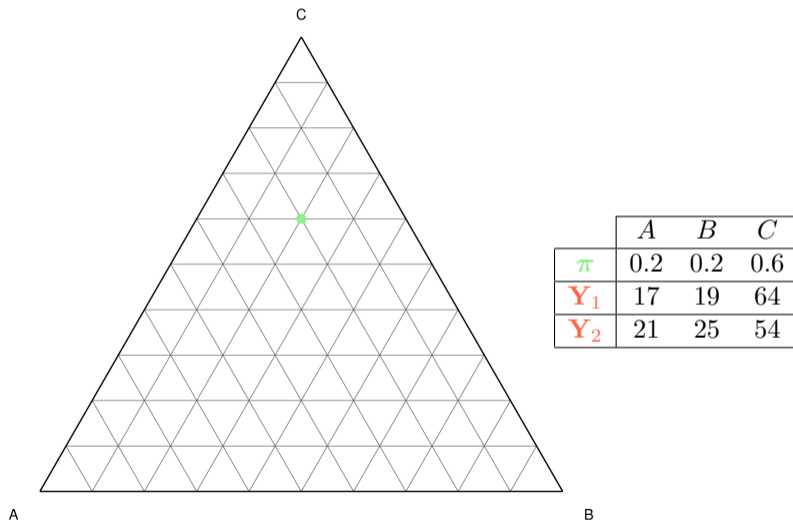


	●	●	●	●	●	●
Prop.	0.25	0.30	0.25	0.05	0.10	0.05
Counts	3	2	3	0	2	0
Obs. Prop.	0.3	0.2	0.3	0	0.2	0

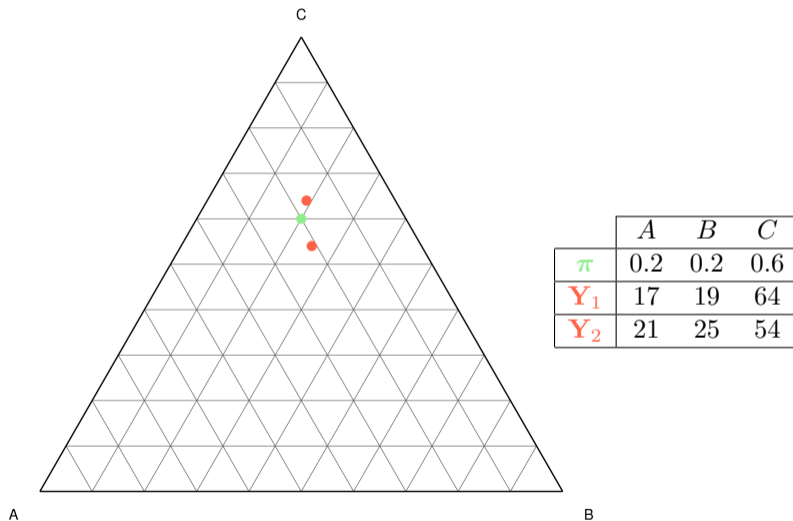
Multinomial Model



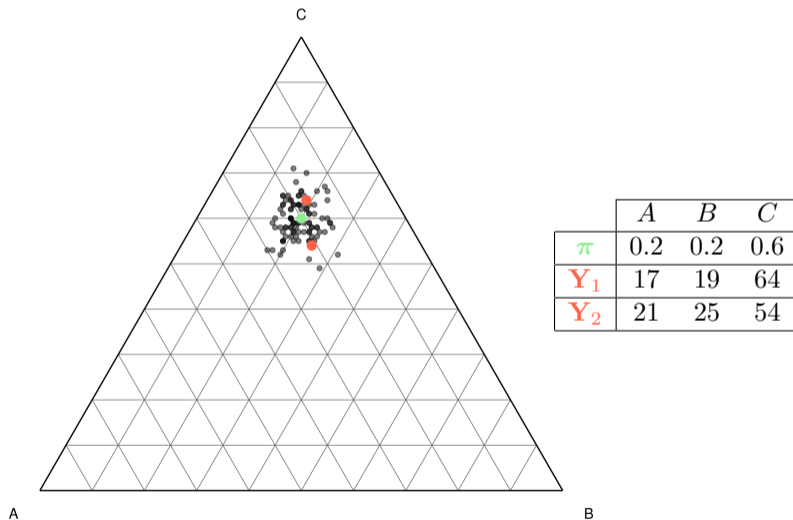
Multinomial Model



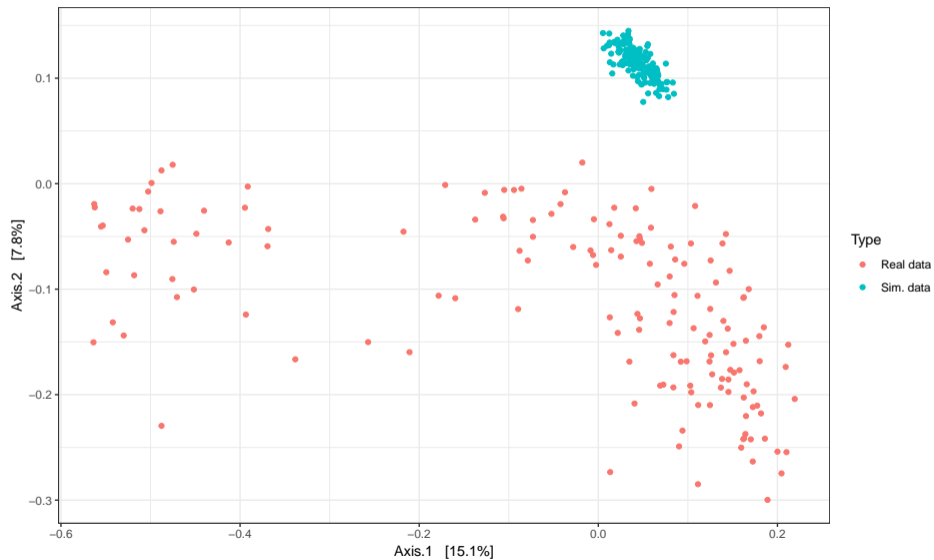
Multinomial Model



Multinomial Model

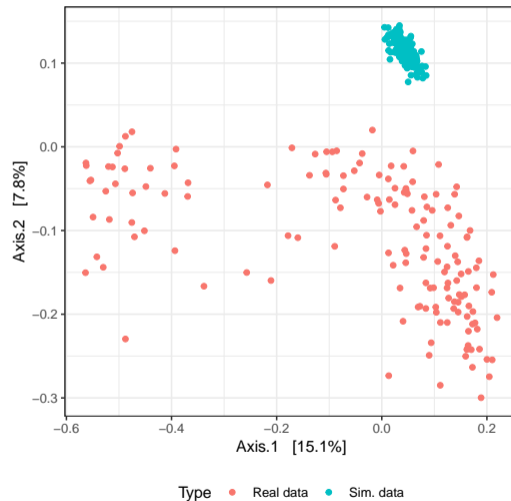


Example of Multinomial Model



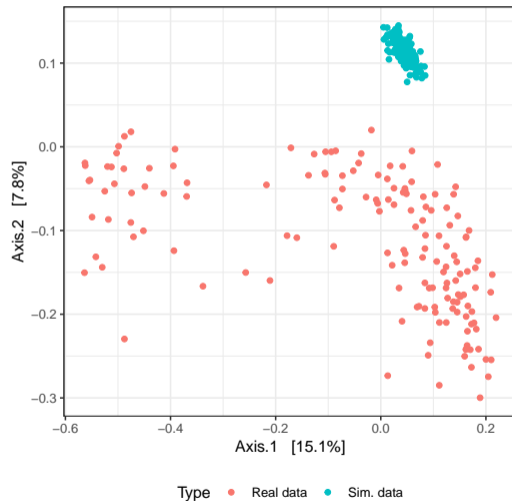
Heterogeneity

- Lack of heterogeneity



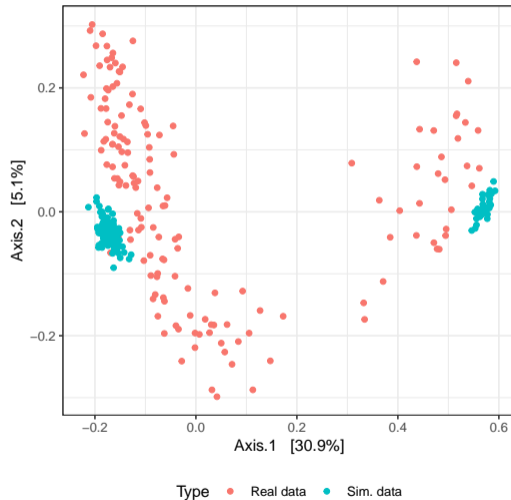
Heterogeneity

- Lack of heterogeneity
↳ Fit only part of the data



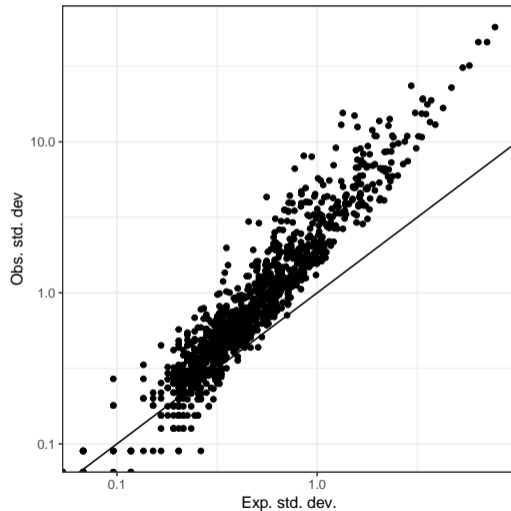
Heterogeneity

- Lack of heterogeneity
 ~> Fit only part of the data
- Lack of variance



Heterogeneity

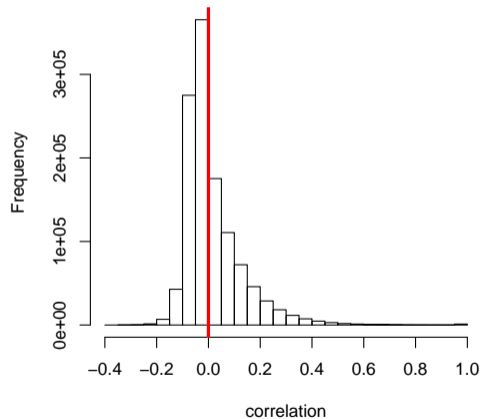
- Lack of heterogeneity
↪ Fit only part of the data
- Lack of variance
- Small dispersion



Heterogeneity

- Lack of heterogeneity
 ~> Fit only part of the data
- Lack of variance
- Small dispersion
- Wrong correlations

Correlation Histogram



Pros

- + Parsimonious model: $p - 1$ parameters to model p abundances
- + Easy to estimate
- + interpretable parameter

Pros

- + Parsimonious model: $p - 1$ parameters to model p abundances
- + Easy to estimate
- + interpretable parameter

Cons

- Bad for heterogeneity
- Bad for dispersion around average composition (\simeq biological variability)
- Bad for correlations between OTUs

1 Motivation

2 Multinomial Models

- Multinomial
- **Mixture of Multinomials**
- (Mixture of) Dirichlet-Multinomial
- Latent Dirichlet Allocation

3 Log-Normal Models

4 Applications



Intuition

- Each sample belongs to one of K groups
- Group k is characterized by its composition π_k
- A sample from group k has composition π_k
- Reads are sampled according to a multinomial process

Intuition

- Each sample belongs to one of K groups
- Group k is characterized by its composition π_k
- A sample from group k has composition π_k
- Reads are sampled according to a multinomial process

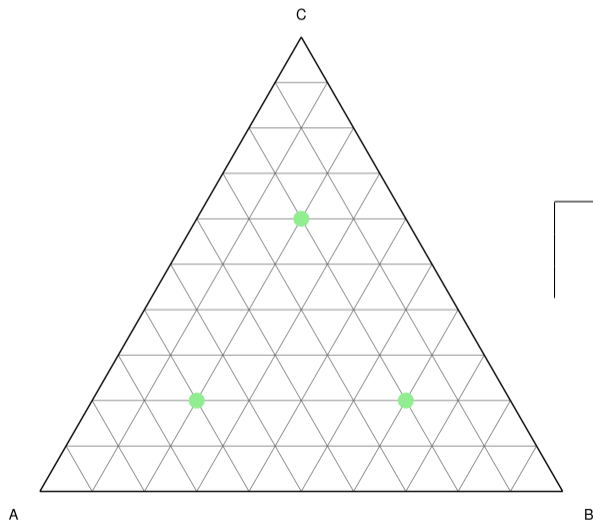
Hierarchical Model

$$Z \sim \mathcal{M}(1, \alpha)$$
$$Y|Z = k \sim \mathcal{M}(N, \pi_k)$$

where

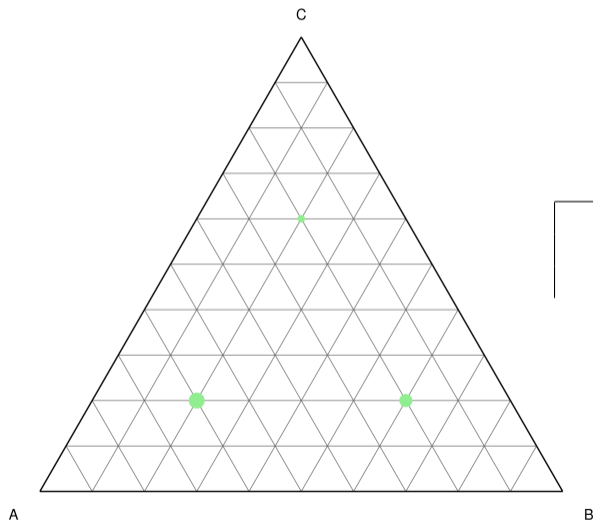
- $\alpha = (\alpha_1, \dots, \alpha_K)$ are the proportions of the K groups,

Mixture of Multinomial



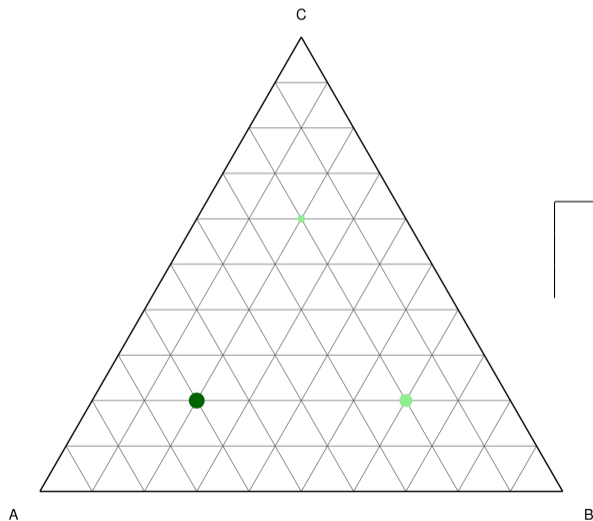
	<i>A</i>	<i>B</i>	<i>C</i>	α
π_1	0.6	0.2	0.2	
π_2	0.2	0.6	0.2	
π_3	0.2	0.2	0.6	

Mixture of Multinomial



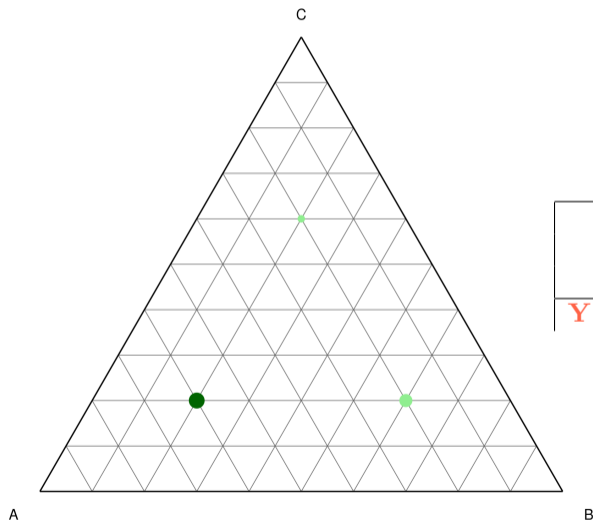
	<i>A</i>	<i>B</i>	<i>C</i>	α
π_1	0.6	0.2	0.2	0.5
π_2	0.2	0.6	0.2	0.4
π_3	0.2	0.2	0.6	0.1

Mixture of Multinomial



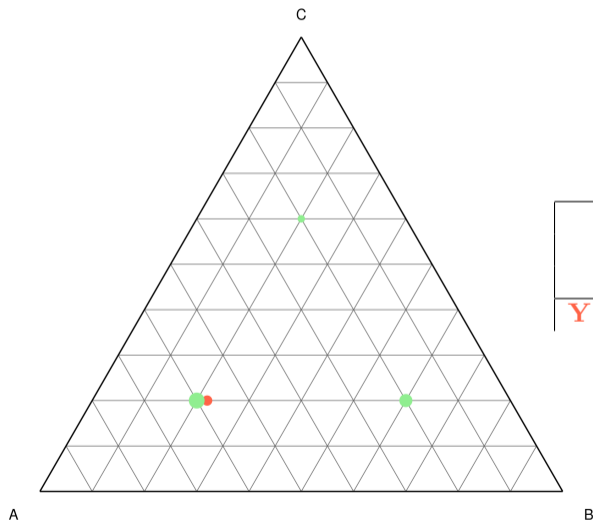
	<i>A</i>	<i>B</i>	<i>C</i>	α
π_1	0.6	0.2	0.2	0.5
π_2	0.2	0.6	0.2	0.4
π_3	0.2	0.2	0.6	0.1

Mixture of Multinomial



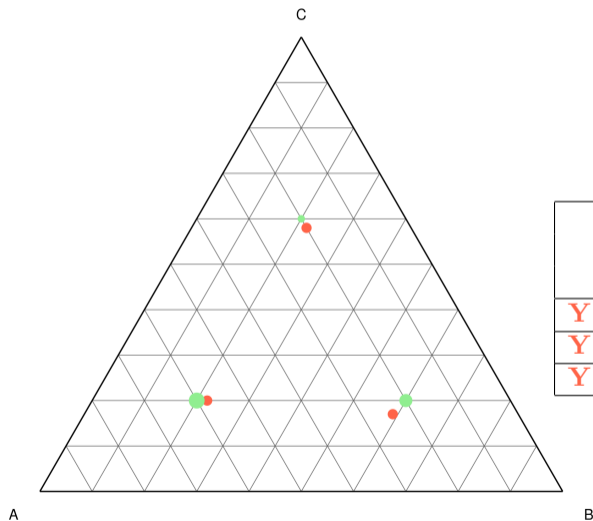
	A	B	C	α
π_1	0.6	0.2	0.2	0.5
π_2	0.2	0.6	0.2	0.4
π_3	0.2	0.2	0.6	0.1
$Y Z=1$	58	22	20	

Mixture of Multinomial



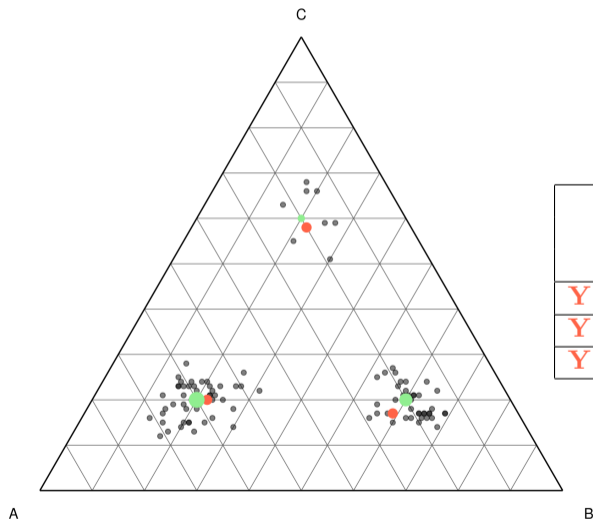
	A	B	C	α
π_1	0.6	0.2	0.2	0.5
π_2	0.2	0.6	0.2	0.4
π_3	0.2	0.2	0.6	0.1
$Y Z=1$	58	22	20	

Mixture of Multinomial



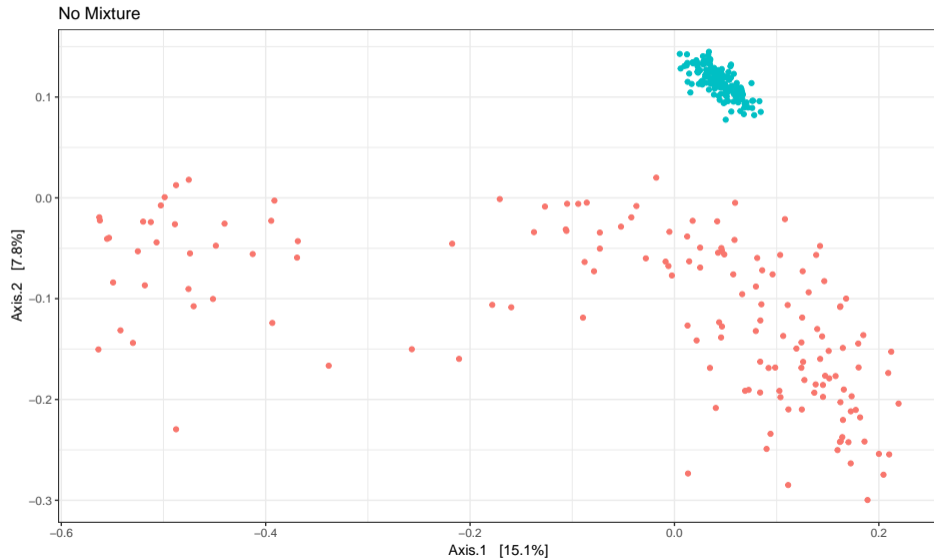
	A	B	C	α
π_1	0.6	0.2	0.2	0.5
π_2	0.2	0.6	0.2	0.4
π_3	0.2	0.2	0.6	0.1
$Y Z=1$	58	22	20	
$Y Z=2$	24	59	17	
$Y Z=3$	20	22	58	

Mixture of Multinomial

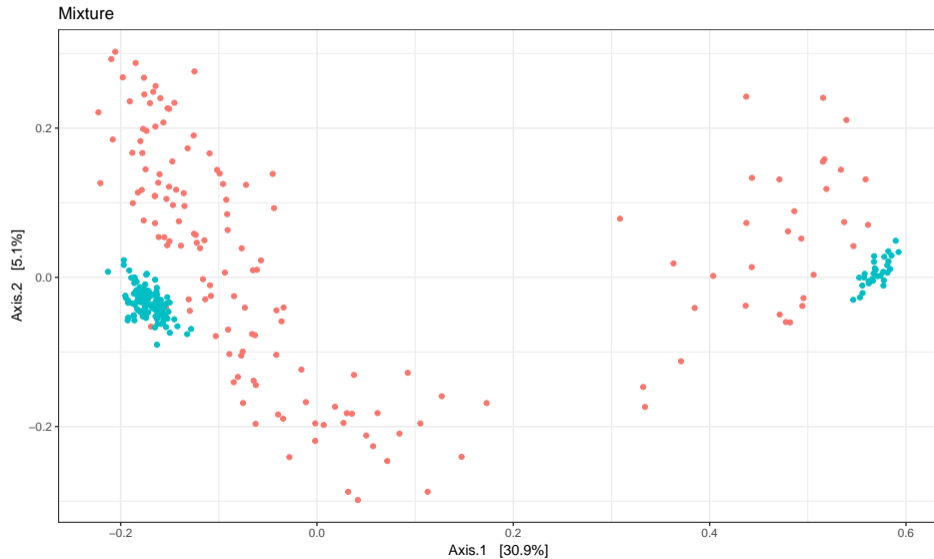


	A	B	C	α
π_1	0.6	0.2	0.2	0.5
π_2	0.2	0.6	0.2	0.4
π_3	0.2	0.2	0.6	0.1
$Y Z=1$	58	22	20	
$Y Z=2$	24	59	17	
$Y Z=3$	20	22	58	

Example of Mixture Models



Example of Mixture Models



Pros

- + Good for **heterogeneity**
- + **Parcimonious**: $Kp - 1$ parameters for K groups
- + Inference is easy when groups are known \rightsquigarrow simple averages

Pros

- + Good for **heterogeneity**
- + **Parcimonious**: $Kp - 1$ parameters for K groups
- + Inference is easy when groups are known \rightsquigarrow simple averages

Cons

- Inference is more involved when groups are unknown
 \rightsquigarrow iterative EM algorithm
- Bad for **dispersion**
- Bad for **correlations** between OTUs

1 Motivation

2 Multinomial Models

- Multinomial
- Mixture of Multinomials
- (Mixture of) Dirichlet-Multinomial
- Latent Dirichlet Allocation

3 Log-Normal Models

4 Applications



Intuition

- π is the **ecosystem-level** average composition

Intuition

- π is the **ecosystem-level** average composition
- Sample i has **own** composition π_i (**noisy version** of π) \rightsquigarrow **Biological** variability

Intuition

- π is the **ecosystem-level** average composition
- Sample i has **own** composition π_i (**noisy version** of π) \rightsquigarrow **Biological** variability
- N_i reads are sampled from π_i according to a multinomial \rightsquigarrow **Technical** / **Sampling** variability

Intuition

- π is the **ecosystem-level** average composition
- Sample i has **own** composition π_i (**noisy version** of π) \rightsquigarrow **Biological** variability
- N_i reads are sampled from π_i according to a multinomial \rightsquigarrow **Technical / Sampling** variability

Hierarchical Model

π	Ecosystem average composition
$\pi_i \sim \mathcal{D}(\kappa\pi)$	Sample average composition
$\mathbf{Y}_i \sim \mathcal{M}(N_i, \pi_i)$	Observed counts

where $1/\kappa$ models the **level of variability** (large $1/\kappa \rightsquigarrow$ large variability)

Dirichlet - Multinomial

Intuition

- π is the **ecosystem-level** average composition
- Sample i has **own** composition π_i (**noisy version** of π) \rightsquigarrow **Biological** variability
- N_i reads are sampled from π_i according to a multinomial \rightsquigarrow **Technical** / **Sampling** variability

Hierarchical Model

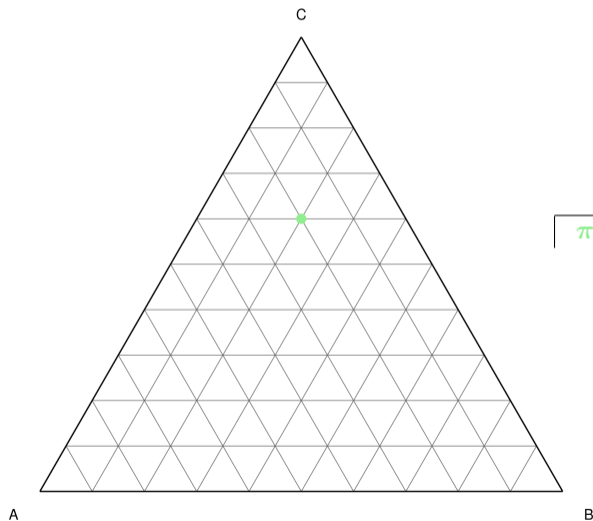
π	Ecosystem average composition
$\pi_i \sim \mathcal{D}(\kappa\pi)$	Sample average composition
$\mathbf{Y}_i \sim \mathcal{M}(N_i, \pi_i)$	Observed counts

where $1/\kappa$ models the **level of variability** (large $1/\kappa \rightsquigarrow$ large variability)

Mixture Layer

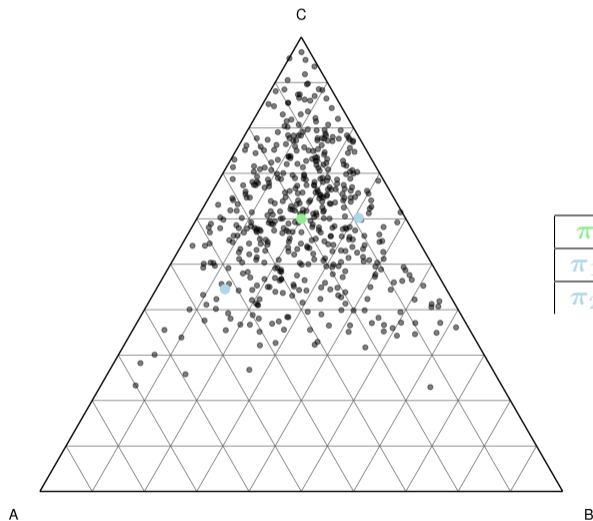
Can be **combined** with a mixture model

Dirichlet-Multinomial distribution



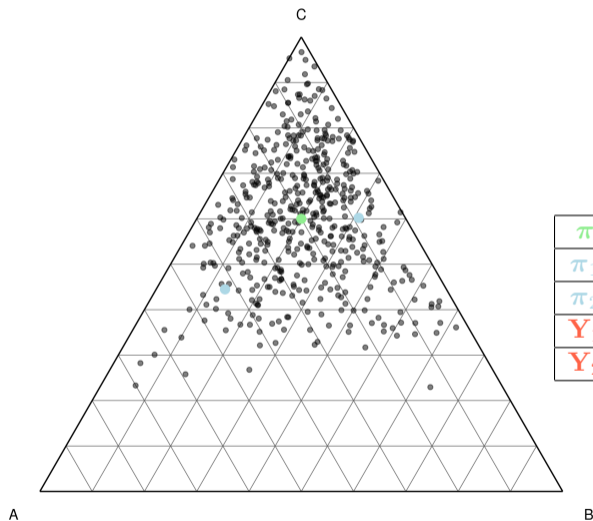
	<i>A</i>	<i>B</i>	<i>C</i>
π	0.2	0.2	0.6

Dirichlet-Multinomial distribution



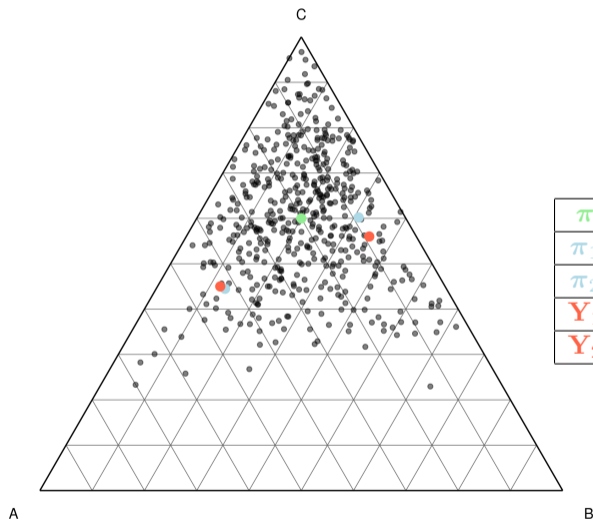
	<i>A</i>	<i>B</i>	<i>C</i>
π	0.2	0.2	0.6
π_1	0.089	0.309	0.602
π_2	0.423	0.132	0.445

Dirichlet-Multinomial distribution



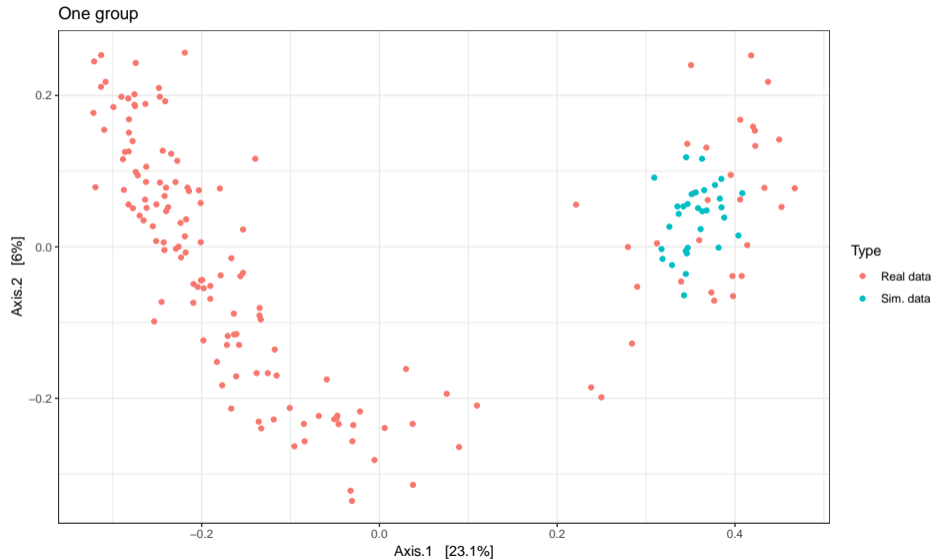
	A	B	C
π	0.2	0.2	0.6
π_1	0.089	0.309	0.602
π_2	0.423	0.132	0.445
Y_1	9	35	56
Y_2	43	12	45

Dirichlet-Multinomial distribution

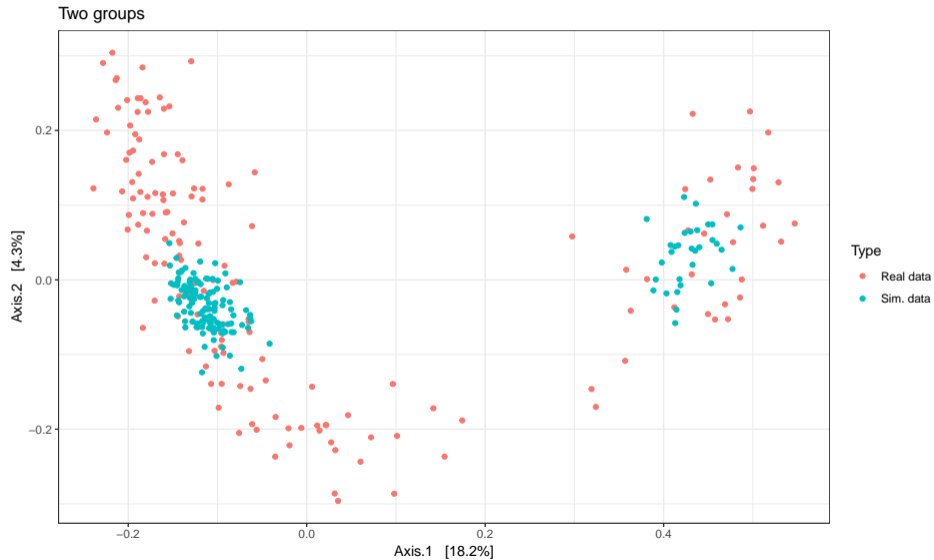


	A	B	C
π	0.2	0.2	0.6
π_1	0.089	0.309	0.602
π_2	0.423	0.132	0.445
Y_1	9	35	56
Y_2	43	12	45

Example of Dirichlet-Multinomial



Example of Dirichlet-Multinomial (Cont'd)



Pros

- + Good for **heterogeneity**
- + So-so of OK for **dispersion**
- + **Parcimonious**: $K(p + 1) - 1$ parameters for K groups

Pros

- + Good for **heterogeneity**
- + So-so of OK for **dispersion**
- + **Parcimonious**: $K(p + 1) - 1$ parameters for K groups

Cons

- **Inference** is more involved
Known groups \rightsquigarrow gradient descent
Unknown groups \rightsquigarrow Iterative EM algorithm + gradient descent
- Bad for **correlations** between OTUs

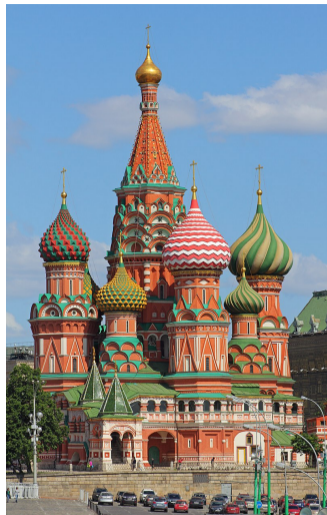
1 Motivation

2 Multinomial Models

- Multinomial
- Mixture of Multinomials
- (Mixture of) Dirichlet-Multinomial
- Latent Dirichlet Allocation

3 Log-Normal Models

4 Applications



Intuition

- There are K **archetype** ecosystems $1, \dots, K$

Intuition

- There are K **archetype** ecosystems $1, \dots, K$
- Each archetype has its own **composition** π_k

Intuition

- There are K **archetype** ecosystems $1, \dots, K$
- Each archetype has its own **composition** $\boldsymbol{\pi}_k$
- Each sample \mathbf{Y} is made-up of **several archetypes** in proportions $(\theta_1, \dots, \theta_K)$

Intuition

- There are K **archetype** ecosystems $1, \dots, K$
- Each archetype has its own **composition** π_k
- Each sample \mathbf{Y} is made-up of **several archetypes** in proportions $(\theta_1, \dots, \theta_K)$
- $\theta_k N$ reads are sampled from a **noisy version** of π_k

Intuition

- There are K **archetype** ecosystems $1, \dots, K$
- Each archetype has its own **composition** π_k
- Each sample \mathbf{Y} is made-up of **several archetypes** in proportions $(\theta_1, \dots, \theta_K)$
- $\theta_k N$ reads are sampled from a **noisy version** of π_k

Hierarchical Model

$$\pi_1, \dots, \pi_K$$

Archetypes average compositions

$$\boldsymbol{\theta} \sim \mathcal{D}(\kappa \boldsymbol{\alpha})$$

Proportion of archetypes in sample

$$\tilde{\pi}_k \sim \mathcal{D}(\kappa_k \pi_k)$$

Noisy version of π_k

$$z_i \sim \mathcal{M}(1, \boldsymbol{\theta})$$

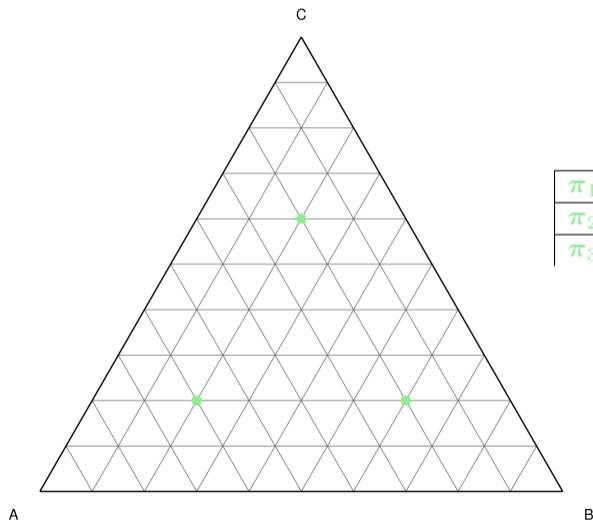
Archetype of origin of read i

$$w_i | z_i = k \sim \mathcal{M}(1, \tilde{\pi}_k)$$

OTU of read i

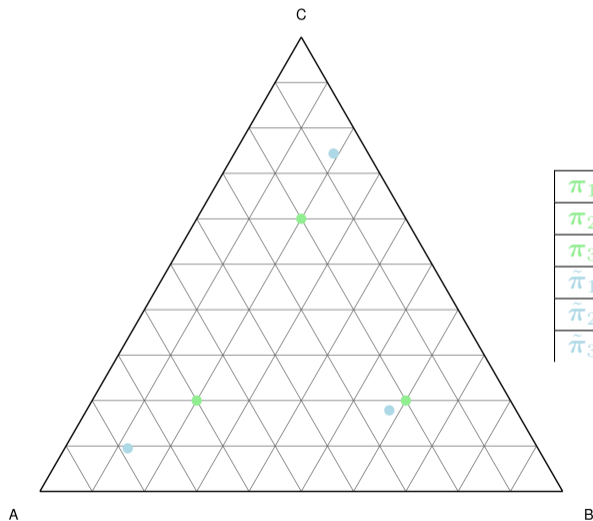
where κ and the κ_k control noise levels.

Latent Dirichlet Allocation



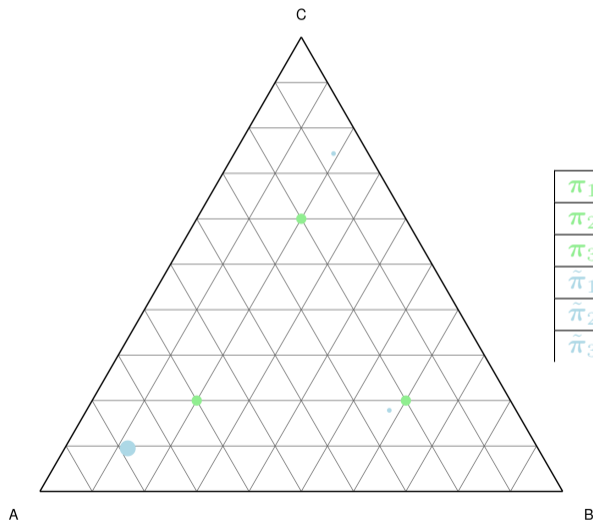
	<i>A</i>	<i>B</i>	<i>C</i>	θ
π_1	0.6	0.2	0.2	
π_2	0.2	0.6	0.2	
π_3	0.2	0.2	0.6	

Latent Dirichlet Allocation



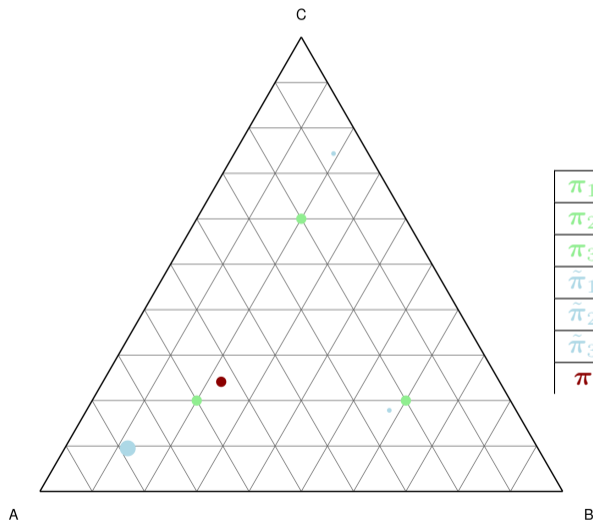
	<i>A</i>	<i>B</i>	<i>C</i>	θ
π_1	0.6	0.2	0.2	
π_2	0.2	0.6	0.2	
π_3	0.2	0.2	0.6	
$\tilde{\pi}_1$	0.784	0.121	0.095	
$\tilde{\pi}_2$	0.242	0.579	0.179	
$\tilde{\pi}_3$	0.423	0.132	0.445	

Latent Dirichlet Allocation



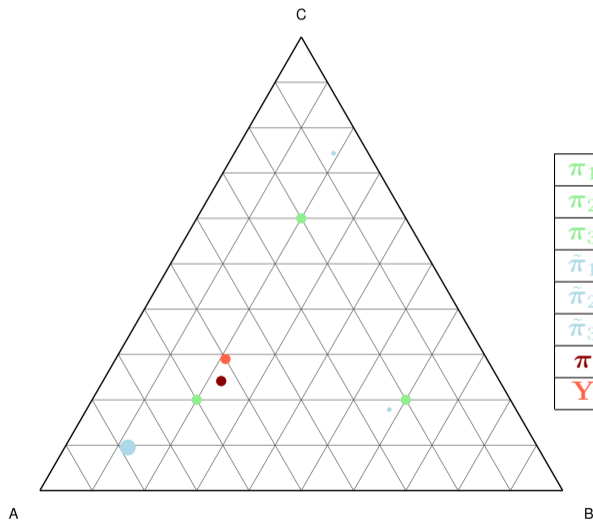
	<i>A</i>	<i>B</i>	<i>C</i>	θ
π_1	0.6	0.2	0.2	
π_2	0.2	0.6	0.2	
π_3	0.2	0.2	0.6	
$\tilde{\pi}_1$	0.784	0.121	0.095	0.6
$\tilde{\pi}_2$	0.242	0.579	0.179	0.2
$\tilde{\pi}_3$	0.423	0.132	0.445	0.2

Latent Dirichlet Allocation



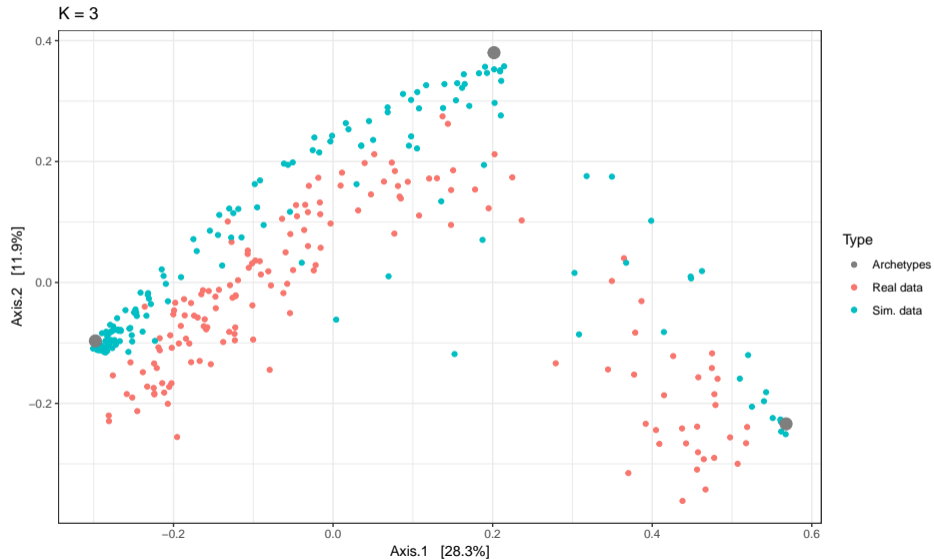
	A	B	C	θ
π_1	0.6	0.2	0.2	
π_2	0.2	0.6	0.2	
π_3	0.2	0.2	0.6	
$\tilde{\pi}_1$	0.784	0.121	0.095	0.6
$\tilde{\pi}_2$	0.242	0.579	0.179	0.2
$\tilde{\pi}_3$	0.423	0.132	0.445	0.2
π	0.532	0.226	0.241	

Latent Dirichlet Allocation

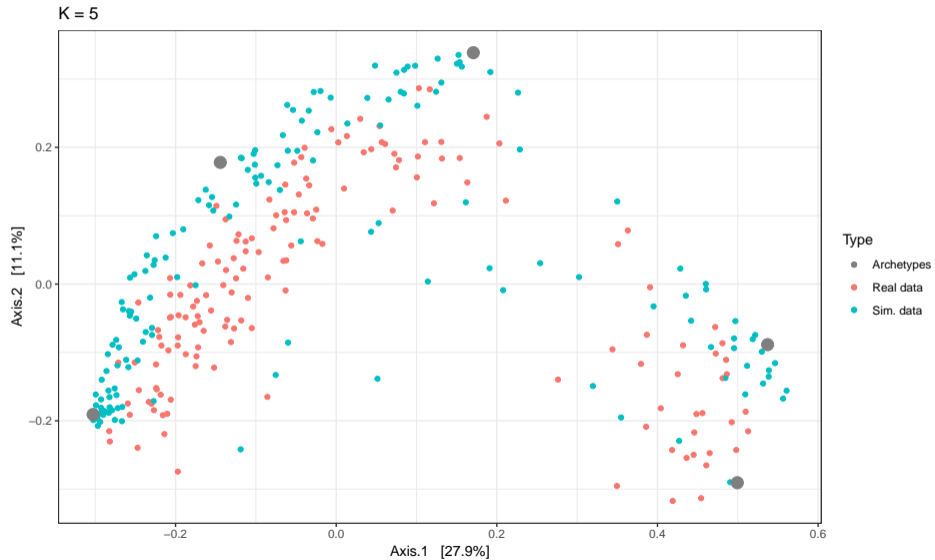


	A	B	C	θ
π_1	0.6	0.2	0.2	
π_2	0.2	0.6	0.2	
π_3	0.2	0.2	0.6	
$\tilde{\pi}_1$	0.784	0.121	0.095	0.6
$\tilde{\pi}_2$	0.242	0.579	0.179	0.2
$\tilde{\pi}_3$	0.423	0.132	0.445	0.2
π	0.532	0.226	0.241	
Y	54	18	28	

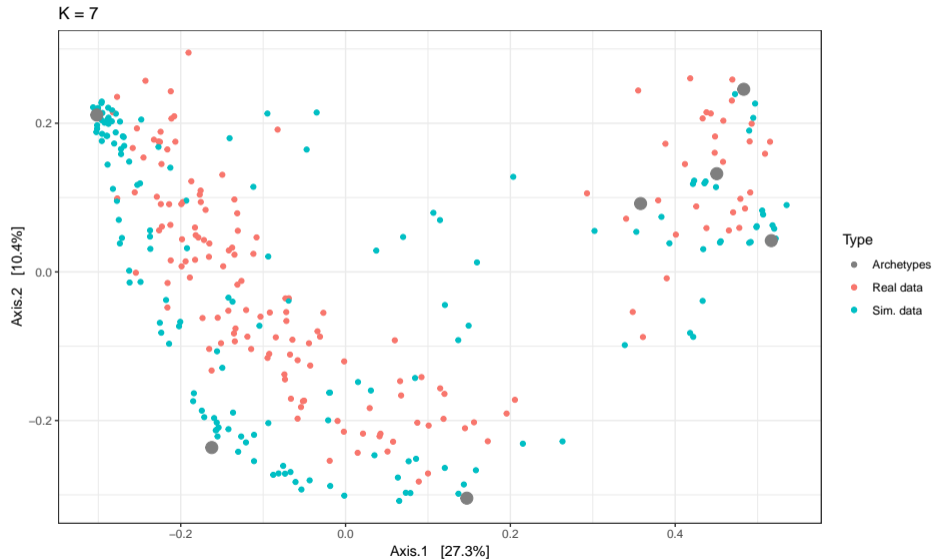
Example of Latent Dirichlet Allocation



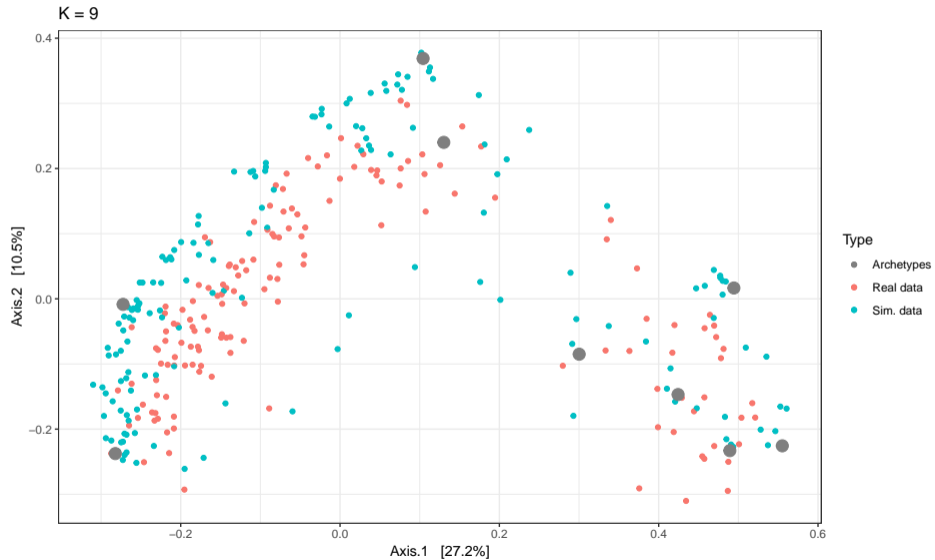
Example of Latent Dirichlet Allocation



Example of Latent Dirichlet Allocation



Example of Latent Dirichlet Allocation



Pros

- + Good for **heterogeneity**
- + Good for **dispersion**
- + **Parcimonious**: $K(p + 1)$ parameters for K archetypes

Pros

- + Good for **heterogeneity**
- + Good for **dispersion**
- + **Parcimonious**: $K(p + 1)$ parameters for K archetypes

Cons

- **Inference** is very involved
↔ gradient descent + EM algorithm / Gibbs sampling
- **Interpretation** is complex ↔ archetypes are **not groups**
- Bad for **correlations** between OTUs

Multinomial-based models are **good** at

- modeling **compositions**;
- modeling **dispersion** around average compositions;
- modeling **heterogeneity**;
- using (relatively) few parameters

Multinomial-based models are **good** at

- modeling **compositions**;
- modeling **dispersion** around average compositions;
- modeling **heterogeneity**;
- using (relatively) few parameters

Multinomial models are **bad** at

- modeling **interactions** between covariates;
- accounting for **covariates**;
- Integrating datasets from **different sources** (e.g. 16S, ITS)

- 1 Motivation
- 2 Multinomial Models
- 3 Log-Normal Models**
 - Multinomial Log-Normal
 - Poisson Log-Normal
- 4 Applications

Multivariate Gaussian models are the *de facto* distribution to model correlations.

Modeling Correlations

Multivariate Gaussian models are the *de facto* distribution to model correlations.

For continuous variables

- The p variables \mathbf{Y}_i (e.g. species abundances) are explained
- by the values of the d covariates \mathbf{X}_i and the p offsets \mathbf{O}_i

$$\mathbf{Y}_i = \underbrace{\mathbf{X}_i \mathbf{B}}_{\text{accounts for covariates}} + \underbrace{\mathbf{O}_i}_{\text{accounts for sampling effort}} + \varepsilon_i, \quad \varepsilon_i \sim \mathcal{N}(\mathbf{0}_p, \underbrace{\Sigma}_{\text{dependencies between species}})$$

+ null covariance \Leftrightarrow independence \rightsquigarrow uncorrelated species do not interact

Modeling Correlations

Multivariate Gaussian models are the *de facto* distribution to model correlations.

For continuous variables

- The p variables \mathbf{Y}_i (e.g. species abundances) are explained
- by the values of the d covariates \mathbf{X}_i and the p offsets \mathbf{O}_i

$$\mathbf{Y}_i = \underbrace{\mathbf{X}_i \mathbf{B}}_{\text{accounts for covariates}} + \underbrace{\mathbf{O}_i}_{\text{accounts for sampling effort}} + \varepsilon_i, \quad \varepsilon_i \sim \mathcal{N}(\mathbf{0}_p, \underbrace{\Sigma}_{\text{dependencies between species}})$$

+ ~~null covariance~~ \Leftrightarrow ~~independence~~ \rightsquigarrow ~~uncorrelated species do not interact~~

But abundances are not gaussian...

Modeling Correlations

Multivariate Gaussian models are the *de facto* distribution to model correlations.

For continuous variables

- The p variables \mathbf{Y}_i (e.g. species abundances) are explained
- by the values of the d covariates \mathbf{X}_i and the p offsets \mathbf{O}_i

$$\mathbf{Y}_i = \underbrace{\mathbf{X}_i \mathbf{B}}_{\text{accounts for covariates}} + \underbrace{\mathbf{O}_i}_{\text{accounts for sampling effort}} + \varepsilon_i, \quad \varepsilon_i \sim \mathcal{N}(\mathbf{0}_p, \underbrace{\Sigma}_{\text{dependencies between species}})$$

+ ~~null covariance~~ \Leftrightarrow ~~independence~~ \rightsquigarrow ~~uncorrelated species do not interact~~

But abundances are not gaussian...

Use a **latent variable models** with a gaussian latent layer and a count observed layer

1 Motivation

2 Multinomial Models

3 **Log-Normal Models**

- **Multinomial Log-Normal**
- Poisson Log-Normal

4 Applications



Intuition

- The latent layer models so-called **basis abundances** \mathbf{z}

Intuition

- The latent layer models so-called **basis abundances** \mathbf{z}
- Basis are **transformed** to an average **composition** π

Intuition

- The latent layer models so-called **basis abundances** \mathbf{z}
- Basis are **transformed** to an average **composition** π
- N reads are **sampled** from π according to a multinomial distribution

Multinomial Log-Normal

Intuition

- The latent layer models so-called **basis abundances** \mathbf{z}
- Basis are **transformed** to an average **composition** $\boldsymbol{\pi}$
- N reads are **sampled** from $\boldsymbol{\pi}$ according to a multinomial distribution

Hierarchical Model

$$\mathbf{z} \sim \mathcal{N}(\boldsymbol{\mu}, \boldsymbol{\Sigma})$$

Abundance basis

$$\boldsymbol{\pi} | \mathbf{z} = \left(\frac{e^{z_j}}{\sum_{j'} e^{z_{j'}}} \right)_j$$

Average composition

$$\mathbf{Y} \sim \mathcal{M}(N, \boldsymbol{\pi})$$

Observed composition

Multinomial Log-Normal

Intuition

- The latent layer models so-called **basis abundances** \mathbf{z}
- Basis are **transformed** to an average **composition** $\boldsymbol{\pi}$
- N reads are **sampled** from $\boldsymbol{\pi}$ according to a multinomial distribution

Hierarchical Model

$$\mathbf{z} \sim \mathcal{N}(\boldsymbol{\mu}, \boldsymbol{\Sigma})$$

Abundance basis

$$\boldsymbol{\pi} | \mathbf{z} = \left(\frac{e^{z_j}}{\sum_{j'} e^{z_{j'}}} \right)_j$$

Average composition

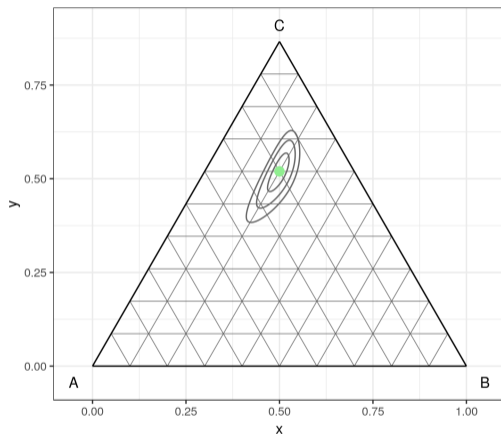
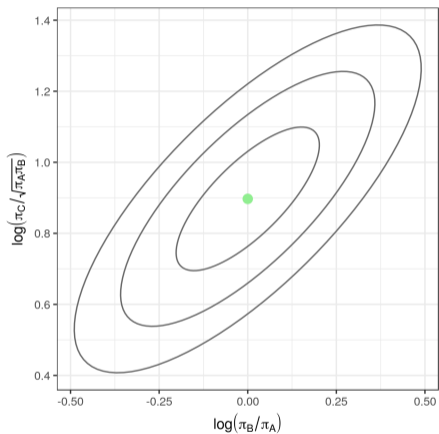
$$\mathbf{Y} \sim \mathcal{M}(N, \boldsymbol{\pi})$$

Observed composition

Mixture Layer

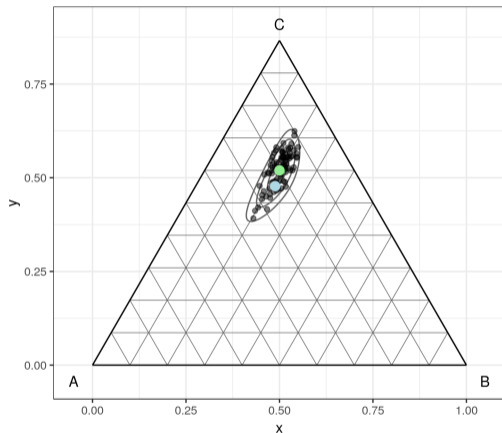
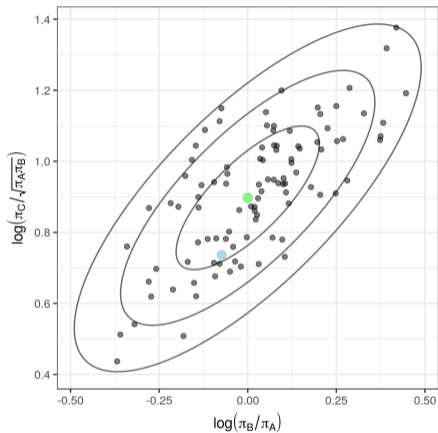
Can be combined with a mixture model

Multinomial Log-Normal



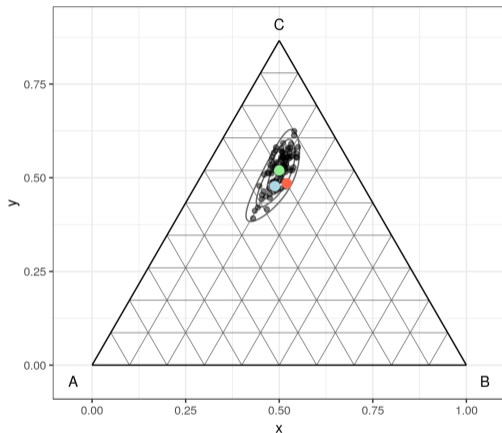
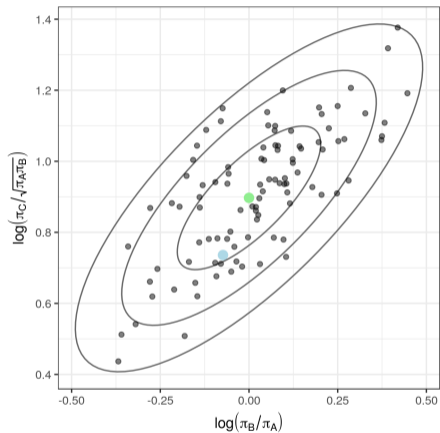
	A	B	C
π	0.2	0.2	0.6

Multinomial Log-Normal



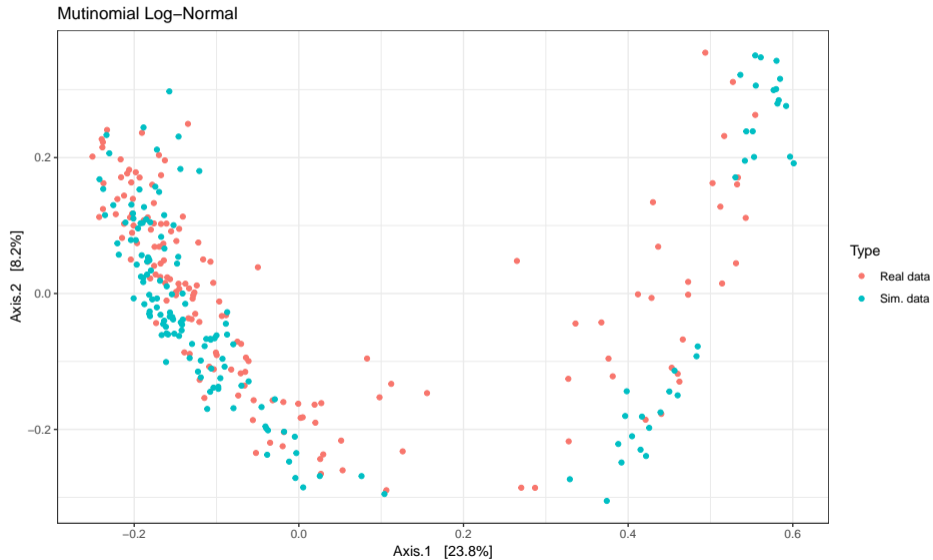
	A	B	C
π	0.2	0.2	0.6
π_1	0.235	0.213	0.552

Multinomial Log-Normal



	A	B	C
π	0.2	0.2	0.6
π_1	0.235	0.213	0.552
Y	20	24	56

Example of Multinomial Log-Normal



Pros

- + Good for **heterogeneity**
- + Good for **dispersion**
- + Good for **correlations** between OTUs

Pros and Cons

Pros

- + Good for **heterogeneity**
- + Good for **dispersion**
- + Good for **correlations** between OTUs

Cons

- The model is not **parsimonious**: $p(p + 3)/2$ parameters
- **Inference** is involved
 - ↪ iterative EM algorithm
- Modeling is done at the **proportion** level

- 1 Motivation
- 2 Multinomial Models
- 3 **Log-Normal Models**
 - Multinomial Log-Normal
 - **Poisson Log-Normal**
- 4 Applications



Intuition

- The latent layer models **basis** \mathbf{z}

Intuition

- The latent layer models **basis \mathbf{z}**
- Basis are *transformed* to average **counts**

Intuition

- The latent layer models **basis \mathbf{z}**
- Basis are *transformed* to average **counts**
- Reads are *sampled* according to Poisson distribution

Intuition

- The latent layer models **basis** \mathbf{z}
- Basis are *transformed* to average **counts**
- Reads are *sampled* according to Poisson distribution

Hierarchical Model

$$\mathbf{z} \sim \mathcal{N}(\boldsymbol{\mu}, \boldsymbol{\Sigma})$$

$$\lambda_j | \mathbf{z} = e^{z_j}$$

$$\mathbf{Y}_j | \mathbf{z} \sim \mathcal{P}(e^{z_j})$$

Basis

Average count of species j

Observed count of species j

Poisson-log normal (PLN) distribution [AH89]

Intuition

- The latent layer models **basis** \mathbf{z}
- Basis are *transformed* to average **counts**
- Reads are *sampled* according to Poisson distribution

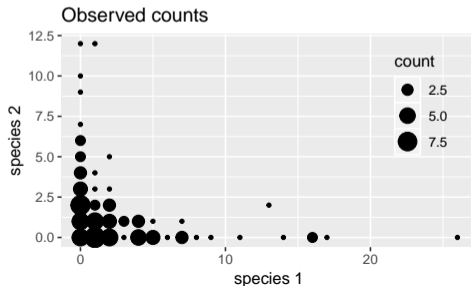
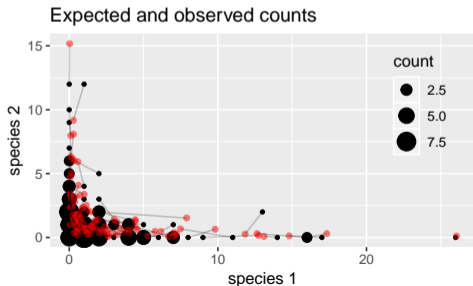
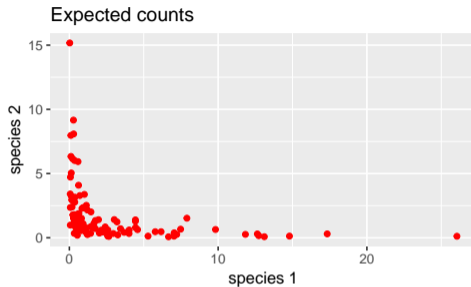
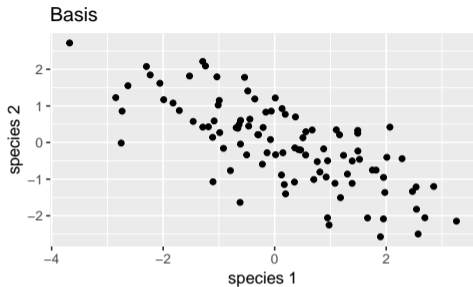
Hierarchical Model

$$\begin{array}{ll} \mathbf{z} \sim \mathcal{N}(\boldsymbol{\mu}, \boldsymbol{\Sigma}) & \text{Basis} \\ \lambda_j | \mathbf{z} = e^{z_j} & \text{Average count of species } j \\ \mathbf{Y}_j | \mathbf{z} \sim \mathcal{P}(e^{z_j}) & \text{Observed count of species } j \end{array}$$

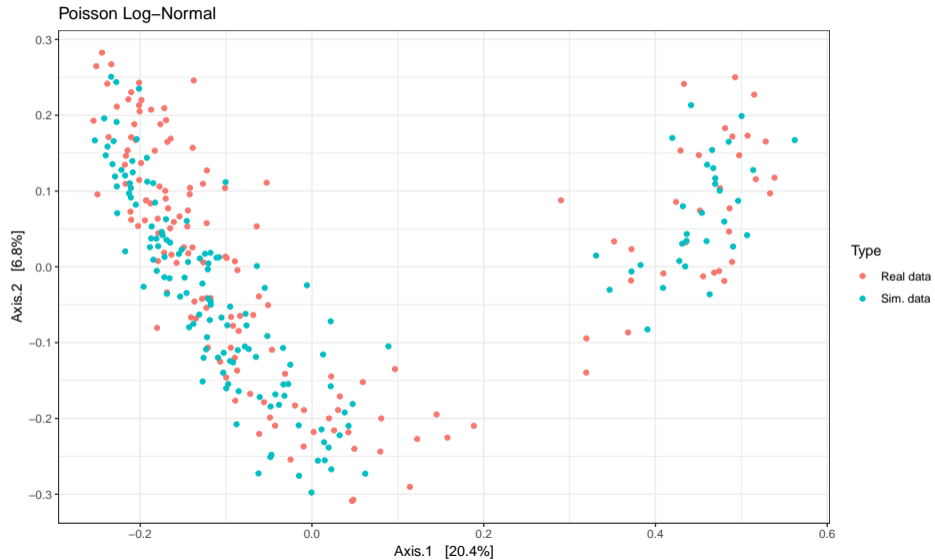
Mixture Layer

Can be combined with a mixture model

Geometrical view



Example of Poisson Log-Normal



Pros

- + Good for **heterogeneity**
- + Good for **dispersion**
- + Good for **correlations** between OTUs
- + Modeling done at the **count** level
 - ↪ counts can be on **different scales** and **come from different sources**

Pros and Cons

Pros

- + Good for **heterogeneity**
- + Good for **dispersion**
- + Good for **correlations** between OTUs
- + Modeling done at the **count** level
 - ↪ counts can be on **different scales** and **come from different sources**

Cons

- The model is not **parsimonious**: $p(p + 3)/2$ parameters
- **Inference** is quite involved
 - ↪ iterative EM algorithm + gradient descent
- Sequencing depths are only controlled on **average**

Log-Normal models are **good** at

- modeling **compositions**;
- modeling **dispersion** around average compositions;
- modeling **heterogeneity**;
- modeling **interactions** between OTUs;
- accounting for **covariates** through the linear model.

Partial Summary

Log-Normal models are **good** at

- modeling **compositions**;
- modeling **dispersion** around average compositions;
- modeling **heterogeneity**;
- modeling **interactions** between OTUs;
- accounting for **covariates** through the linear model.

Log-Normal models are **bad** at

- being **parsimonious**

Partial Summary

Log-Normal models are **good** at

- modeling **compositions**;
- modeling **dispersion** around average compositions;
- modeling **heterogeneity**;
- modeling **interactions** between OTUs;
- accounting for **covariates** through the linear model.

Log-Normal models are **bad** at

- being **parsimonious**

- MLN results are easier to interpret (proportions)
- PLN allows to mix data from **different sources** (16S, ITS, etc.)

- 1 Motivation
- 2 Multinomial Models
- 3 Log-Normal Models
- 4 Applications
 - PCA
 - Linear Discriminant Analysis

PLN: a flexible models accounting for:

- Heterogeneity and average compositions (\simeq first order moments)
- Dispersion and correlation between OTUs (\simeq second order moments)
- Structuring covariates
- Counts coming from different data sources

PLN: a flexible models accounting for:

- Heterogeneity and average compositions (\simeq first order moments)
- Dispersion and correlation between OTUs (\simeq second order moments)
- Structuring covariates
- Counts coming from dfferent data sources

Allows for *traditional* multivariate analysis:

Idea: put additional constraints in the model

- PCA \rightsquigarrow small rank Σ
- Linear Discriminant Analysis
- Network Inference \rightsquigarrow sparse/tree-like Σ^{-1}
- *Mixture Models*
- *etc.*

- 1 Motivation
- 2 Multinomial Models
- 3 Log-Normal Models
- 4 Applications
 - PCA
 - Linear Discriminant Analysis

PLN-PCA: summarize information

Dimension reduction and vizualization. Typical task in multivariate analysis

$$\begin{aligned} \mathbf{Z}_i \text{ iid } &\sim \mathcal{N}_p(\mathbf{0}_p, \Sigma), & \text{rank}(\Sigma) = q \ll p \\ \mathbf{Y}_i | \mathbf{Z}_i &\sim \mathcal{P}(\exp\{\mathbf{O}_i + \mathbf{X}_i\beta + \mathbf{Z}_i\}) \end{aligned}$$

↪ Find a low-dimensional base (PCA axes) to represent the latent covariance

Fit the PLNPCA models with offsets and various covariates.

```
Qmax = 30; Q <- 1:Qmax;

## Model with offset
models.offset <- PLNPCA(counts ~ 1 + offset(log(offsets)), ranks=Q)

## Models with offset and covariates (tree + orientation)
formula <- counts ~ 1 + covariates$tree + covariates$orientation + offset(log(offsets))
models.tree.orientation <- PLNPCA(formula, ranks=Q) # approx 10 mn
```

PCA: vizualization

PLN PCA separates well the kind of tree

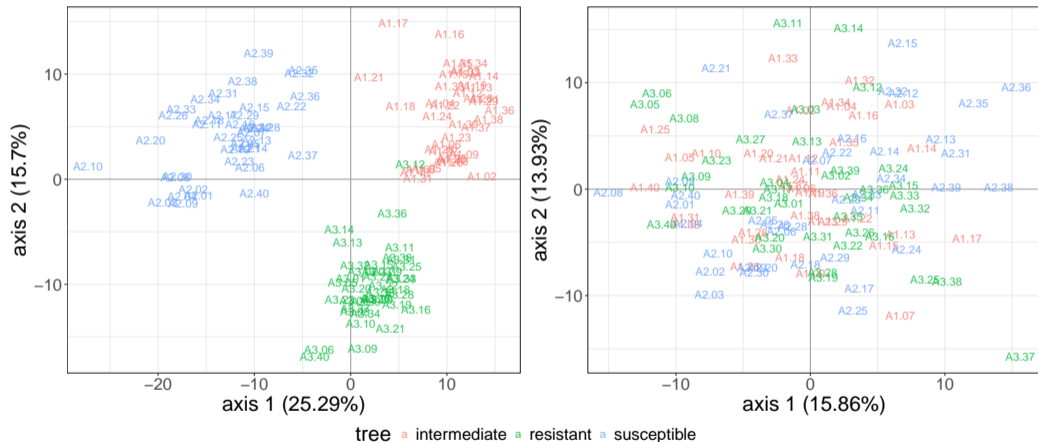


Figure: offset only

offset + covariates

PCA: vizualization II

Introduction of covariates unravels hidden patterns

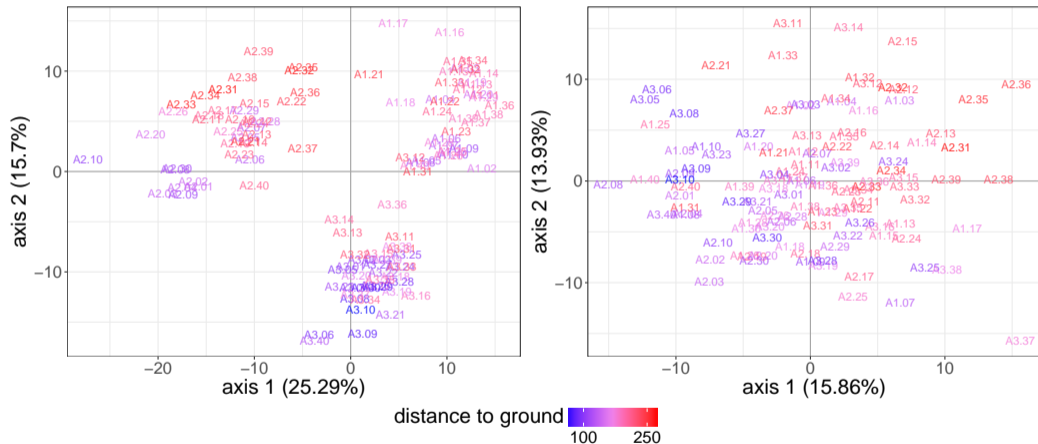


Figure: offset only

offset + covariates

- 1 Motivation
- 2 Multinomial Models
- 3 Log-Normal Models
- 4 Applications
 - PCA
 - Linear Discriminant Analysis

Fit the PLNLDA models

find the linear combination that separates the grouping

Fit the model with offsets, and various covariates

```
myLDA_tree <- PLNLDA(Y, grouping = treeStatus, 0 = log(0))
```

```
##  
## Initialization...  
## Adjusting the standard PLN model.  
## Performing Discriminant Analysis...  
## DONE!
```

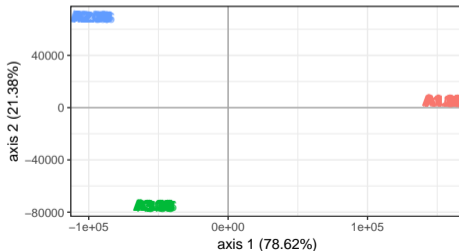
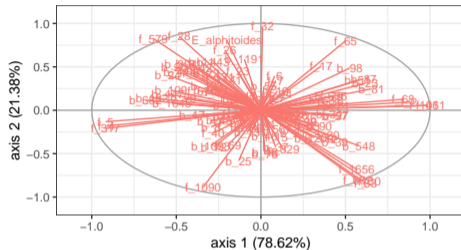
```
myLDA_tree$plot_LDA()
```


LDA on tree status

Axes contribution

axis 1 : 78.62%

axis 2 : 21.38%



Prediction error (10 fold cross-validation)

	intermediate	resistant	susceptible
intermediate	38	0	0
resistant	0	39	0
susceptible	0	0	39

Summary PLN = generic model for multivariate counts

- Corrects for covariates and offset (\simeq sequencing depths)
- Flexible statistical modeling
- `PLNmodels` R-package

Perspectives

- Add technical/biological "zeros" (zero-inflation)
- Extensions: sparse PCA, mixture models
- Confidence interval and tests
- Missing data...



John Aitchison and CH Ho.

The multivariate poisson-log normal distribution.

Biometrika, 76(4):643–653, 1989.



Boris Jakuschkin, Virgil Fievet, Loïc Schwaller, Thomas Fort, Cécile Robin, and Corinne Vacher.

Deciphering the pathobiome: Intra- and interkingdom interactions involving the pathogen *erysiphe alphitoides*.

Microbial Ecology, 72(4):870–880, Nov 2016.



Núria Mach, Mustapha Berri, Jordi Estellé, Florence Levenez, Gaëtan Lemonnier, Catherine Denis, Jean-Jacques Leplat, Claire Chevalyere, Yvon Billon, Joël Doré, and et al.

Early-life establishment of the swine gut microbiome and impact on host phenotypes.

Environmental Microbiology Reports, 7(3):554–569, May 2015.