

Une promenade parmi les modèles pour comptages multivariés



M. Mariadassou, INRAE-MaiAGE

Travail avec Julien Chiquet and Stéphane Robin

CNRS Villa Clythia, Fréjus, 6 octobre 2022



Julien Chiquet, M.M., Stéphane Robin,
Variational inference for probabilistic Poisson PCA
<http://doi.org/10.1214/18-AOAS1177> (*Annals of Applied Statistics*, 2019)



Julien Chiquet, M.M., Stéphane Robin,
Variational inference for network inference with count data
<http://proceedings.mlr.press/v97/chiquet19a/chiquet19a.pdf> (*ICML19*)



Julien Chiquet, M.M., Stéphane Robin,
The Poisson-lognormal model as a versatile framework for the joint analysis of species abundances
<http://doi.org/10.3389/fevo.2021.588292> (*Frontiers in Ecol. and Evol.*, 2021)



PLNmodels package, development version on github
`devtools::install_github("pln-team/PLNmodels", build_vignettes=TRUE)`
<https://pln-team.github.io/PLNmodels/>



- 1 Motivation
- 2 Modèles multinomiaux
- 3 Modèles Log-Normaux
- 4 Applications





Données issues de [MBE⁺15].

- $n = 155$ échantillons (= 31 porcelets à 5 points de temps)
- $p = 1038$ espèces bactériennes (OTUs) avec prévalence ≥ 0.05
- Certaines covariables (sexe, portée, etc)
- Offsets: $o_i =$ offset de l'échantillon i (profondeur de séquençage)



Données issues de [MBE⁺15].

- $n = 155$ échantillons (= 31 porcelets à 5 points de temps)
- $p = 1038$ espèces bactériennes (OTUs) avec prévalence ≥ 0.05
- Certaines covariables (sexe, portée, etc)
- Offsets: $o_i =$ offset de l'échantillon i (profondeur de séquençage)

But: Étudier l'impact du **sevrage** sur le microbiote intestinal

Données de métagénétique de [MBE⁺15]

- **count** matrice de taille $n = 155$ échantillons, $p = 1038$ espèces

```
mach_counts[1:2, c(3, 9, 12, 15)]  
##           5982 347 349 5854  
## SF0901    0 23  3  0  
## SF0902    8  0  4  0
```

- $d = 8$ **covariables** (sexe, portée, statut de sevrage, ...)

```
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
```

- Effort d'observation pour chaque échantillon

```
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
```



Données issues de [JFS⁺16].

- $n = 116$ feuilles de bouleau = échantillons
- $p = 114$ espèces microbiennes
 - $p_1 = 66$ espèces microbiennes (OTUs, issues du marqueur 16S)
 - $p_2 = 48$ espèces fongiques (OTUs, issues du marqueur ITS)
- covariables: arbre (résistant, intermédiaire, susceptible), hauteur, distance au tronc, ...
- offsets: $o_{i1} \neq o_{i2}$ = offset pour les bactéries, espèces fongiques



Données issues de [JFS⁺16].

- $n = 116$ feuilles de bouleau = échantillons
- $p = 114$ espèces microbiennes
 - $p_1 = 66$ espèces microbiennes (OTUs, issues du marqueur 16S)
 - $p_2 = 48$ espèces fongiques (OTUs, issues du marqueur ITS)
- covariables: arbre (résistant, intermédiaire, susceptible), hauteur, distance au tronc, ...
- offsets: $o_{i1} \neq o_{i2}$ = offset pour les bactéries, espèces fongiques

```
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
```



Données issues de [JFS⁺16].

- $n = 116$ feuilles de bouleau = échantillons
- $p = 114$ espèces microbiennes
 - $p_1 = 66$ espèces microbiennes (OTUs, issues du marqueur 16S)
 - $p_2 = 48$ espèces fongiques (OTUs, issues du marqueur ITS)
- covariables: arbre (résistant, intermédiaire, susceptible), hauteur, distance au tronc, ...
- offsets: $o_{i1} \neq o_{i2}$ = offset pour les bactéries, espèces fongiques

```
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
```

But. Comprendre les interactions entre les espèces, y compris le pathogène *E. alphitoides*.

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

- Y_{ij} = abondance (nombre de lectures) de l'espèce j dans l'échantillon i
- X_{ik} = valeur de la covariable k dans l'échantillon i
- O_{ij} = offset (effort d'observation) pour l'espèce j dans l'échantillon i

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

- Y_{ij} = abondance (nombre de lectures) de l'espèce j dans l'échantillon i
- X_{ik} = valeur de la covariable k dans l'échantillon i
- O_{ij} = offset (effort d'observation) pour l'espèce j dans l'échantillon i

Besoin de méthodes d'analyses multivariées pour comprendre les écosystèmes microbiens

- Mettre en évidence des **patrons de diversité**
↪ Résumer l'information de \mathbf{Y} (ACP, classification, ...)
- Comprendre **les interactions entre espèces**
↪ Inférence de réseau (sélection de covariance)
- Corriger les effets techniques et **confondants**
↪ Prise en compte de covariables et d'offsets.

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

- Y_{ij} = abondance (nombre de lectures) de l'espèce j dans l'échantillon i
- X_{ik} = valeur de la covariable k dans l'échantillon i
- O_{ij} = offset (effort d'observation) pour l'espèce j dans l'échantillon i

Besoin de méthodes d'analyses multivariées pour comprendre les écosystèmes microbiens

- Mettre en évidence des **patrons de diversité**
↪ Résumer l'information de \mathbf{Y} (ACP, classification, ...)
- Comprendre **les interactions entre espèces**
↪ Inférence de réseau (sélection de covariance)
- Corriger les effets techniques et **confondants**
↪ Prise en compte de covariables et d'offsets.

↪ Besoin d'un cadre générique pour **modéliser les données de comptages multivariées**

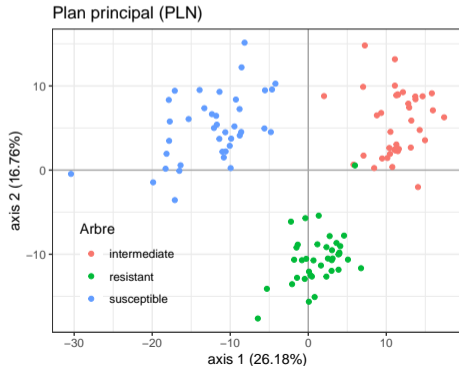
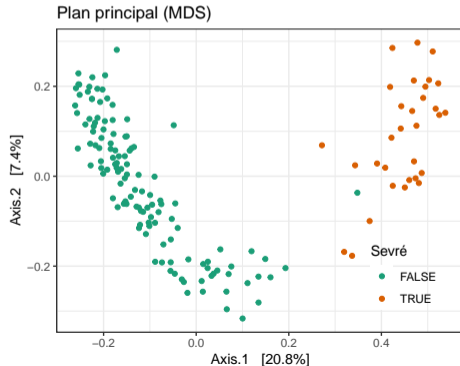
- 1 Appliquer votre **distance** préférée (Jaccard, Bray-Curtis, UniFrac, weighted UniFrac, Aitchison)

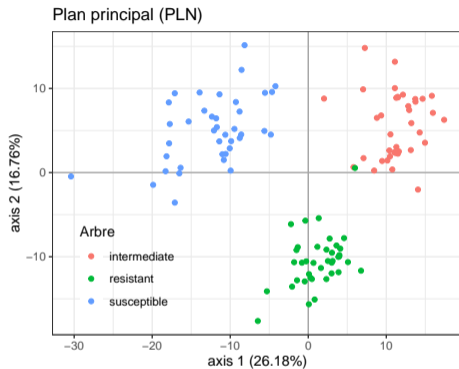
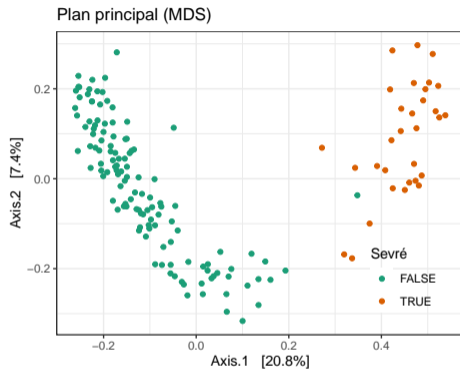
- 1 Appliquer votre **distance** préférée (Jaccard, Bray-Curtis, UniFrac, weighted UniFrac, Aitchison)
- 2 Appliquer votre méthode de **réduction de dimension** préférée (ACP, MDS/PCoA, NMDS, RDA, PLN, etc)

- 1 Appliquer votre **distance** préférée (Jaccard, Bray-Curtis, UniFrac, weighted UniFrac, Aitchison)
- 2 Appliquer votre méthode de **réduction de dimension** préférée (ACP, MDS/PCoA, NMDS, RDA, PLN, etc)
- 3 **Tracer** le résultat

Microbial Ecology 101

- 1 Appliquer votre **distance** préférée (Jaccard, Bray-Curtis, UniFrac, weighted UniFrac, Aitchison)
- 2 Appliquer votre méthode de **réduction de dimension** préférée (ACP, MDS/PCoA, NMDS, RDA, PLN, etc)
- 3 **Tracer** le résultat
- 4 *Et voilà!*





- 1 Parfait pour **trouver** des structures...
- 2 Mais peu idéal pour la **modéliser**

Quel type de modèle générique?

Quel cadre générique pour les données de comptages multivariées?

Quel type de modèle générique?



Pour Noël, je voudrais une famille de modèles **génératifs** qui sont:

Pour Noël, je voudrais une famille de modèles **génératifs** qui sont:

- Suffisamment **flexibles** pour:
 - modéliser les compositions moyennes;
 - modéliser la dispersion (variabilité biologique);
 - modéliser les interactions entre espèces (réseaux écologiques);
 - prendre en compte l'hétérogénéité des communautés;
 - intégrer les données issues de sources différentes (bactéries et champignons)

Pour Noël, je voudrais une famille de modèles **génératifs** qui sont:

- Suffisamment **flexibles** pour:
 - modéliser les compositions moyennes;
 - modéliser la dispersion (variabilité biologique);
 - modéliser les interactions entre espèces (réseaux écologiques);
 - prendre en compte l'hétérogénéité des communautés;
 - intégrer les données issues de sources différentes (bactéries et champignons)
- aussi **parcimonieux** que possible;

Pour Noël, je voudrais une famille de modèles **génératifs** qui sont:

- Suffisamment **flexibles** pour:
 - modéliser les compositions moyennes;
 - modéliser la dispersion (variabilité biologique);
 - modéliser les interactions entre espèces (réseaux écologiques);
 - prendre en compte l'hétérogénéité des communautés;
 - intégrer les données issues de sources différentes (bactéries et champignons)
- aussi **parcimonieux** que possible;
- **interprétables**;

Pour Noël, je voudrais une famille de modèles **génératifs** qui sont:

- Suffisamment **flexibles** pour:
 - modéliser les compositions moyennes;
 - modéliser la dispersion (variabilité biologique);
 - modéliser les interactions entre espèces (réseaux écologiques);
 - prendre en compte l'hétérogénéité des communautés;
 - intégrer les données issues de sources différentes (bactéries et champignons)
- aussi **parcimonieux** que possible;
- **interprétables**;
- **rapides et faciles** à estimer à partir de données;

Pour Noël, je voudrais une famille de modèles **génératifs** qui sont:

- Suffisamment **flexibles** pour:
 - modéliser les compositions moyennes;
 - modéliser la dispersion (variabilité biologique);
 - modéliser les interactions entre espèces (réseaux écologiques);
 - prendre en compte l'hétérogénéité des communautés;
 - intégrer les données issues de sources différentes (bactéries et champignons)
- aussi **parcimonieux** que possible;
- **interprétables**;
- **rapides et faciles** à estimer à partir de données;
- **réalistes** (e.g. simulent des échantillons **réalistes**).

1 Motivation

2 Modèles multinomiaux

- Multinomiale
- Mélange de multinomiales
- (Mélange de) Dirichlet-Multinomiale
- Latent Dirichlet Allocation

3 Modèles Log-Normaux

4 Applications

1 Motivation

2 Modèles multinomiaux

- **Multinomiale**
- Mélange de multinomiales
- (Mélange de) Dirichlet-Multinomiale
- Latent Dirichlet Allocation

3 Modèles Log-Normaux

4 Applications



Intuition

- Il y a p espèces, présentes en proportions $\boldsymbol{\pi} = (\pi_1, \dots, \pi_p)$ et en quantité infinie dans la population
- On échantillonne N (profondeur de séquençage) bactéries dans la population

Intuition

- Il y a p espèces, présentes en proportions $\boldsymbol{\pi} = (\pi_1, \dots, \pi_p)$ et en quantité infinie dans la population
- On échantillonne N (profondeur de séquençage) bactéries dans la population

Modèle mathématiques

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

Intuition

- Il y a p espèces, présentes en proportions $\boldsymbol{\pi} = (\pi_1, \dots, \pi_p)$ et en quantité infinie dans la population
- On échantillonne N (profondeur de séquençage) bactéries dans la population

Modèle mathématiques

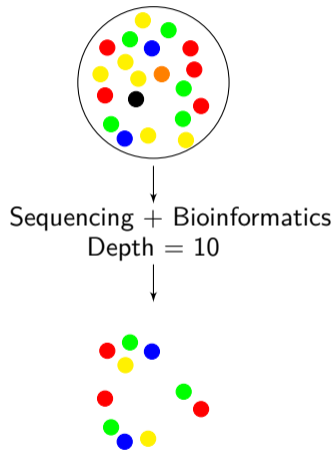
$$\mathbf{Y} \sim \mathcal{M}(N, \boldsymbol{\pi})$$

L'inférence est facile

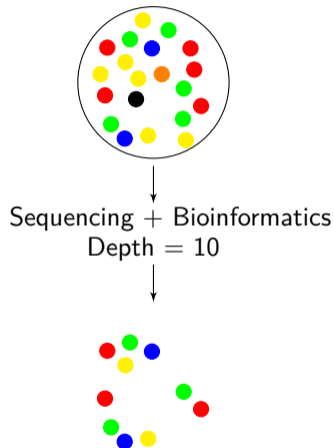
$$\hat{\pi}_j = \frac{\sum_{i=1}^n Y_{ij}}{\sum_{i=1}^n N_i}$$

avec Y_{ij} l'abundance de l'espèce j dans l'échantillon i et N_i la profondeur de l'échantillon i .

Distribution multinomiale: tirage dans une urne (avec remise)

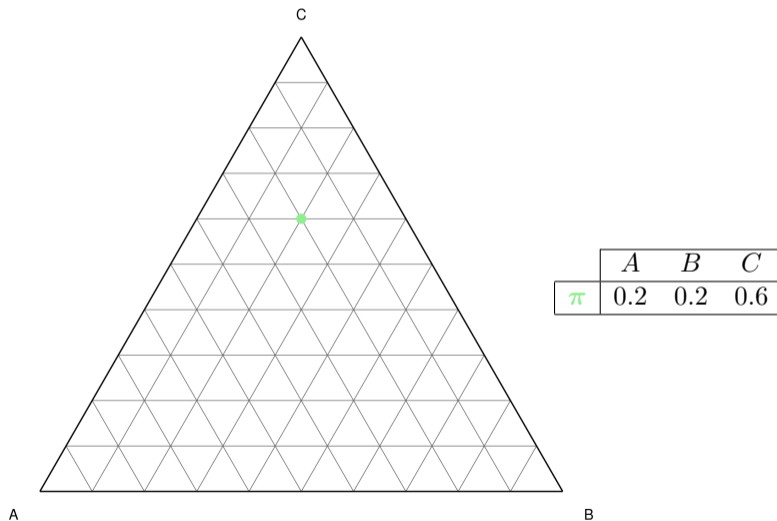


Distribution multinomiale: tirage dans une urne (avec remise)

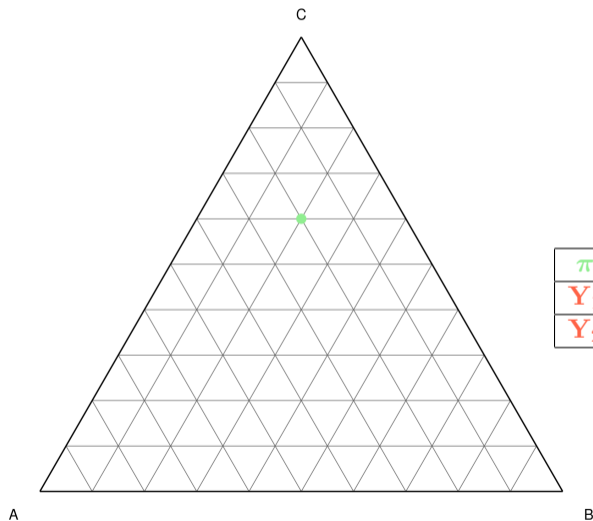


	●	●	●	●	●	●
Prop.	0.25	0.30	0.25	0.05	0.10	0.05
Counts	3	2	3	0	2	0
Obs. Prop.	0.3	0.2	0.3	0	0.2	0

Modèle multinomial

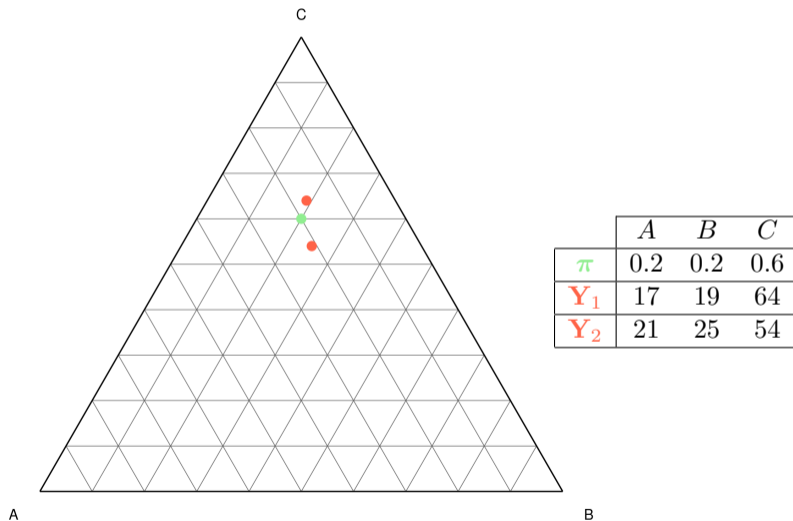


Modèle multinomial

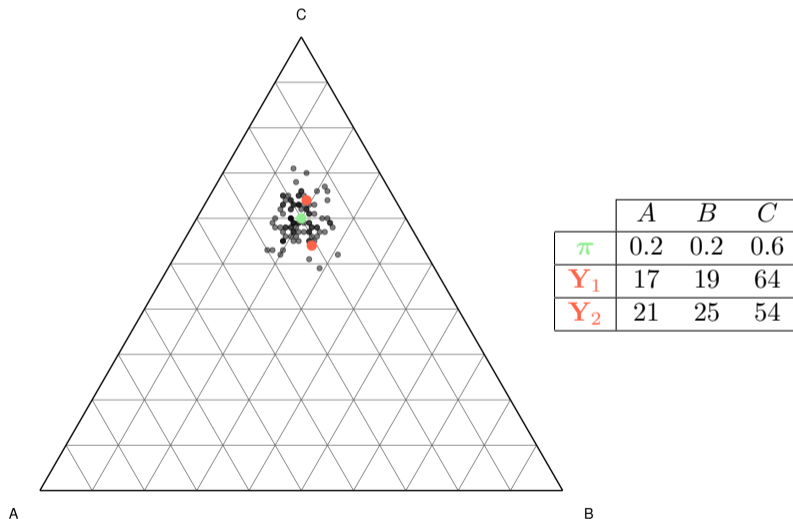


	A	B	C
π	0.2	0.2	0.6
Y_1	17	19	64
Y_2	21	25	54

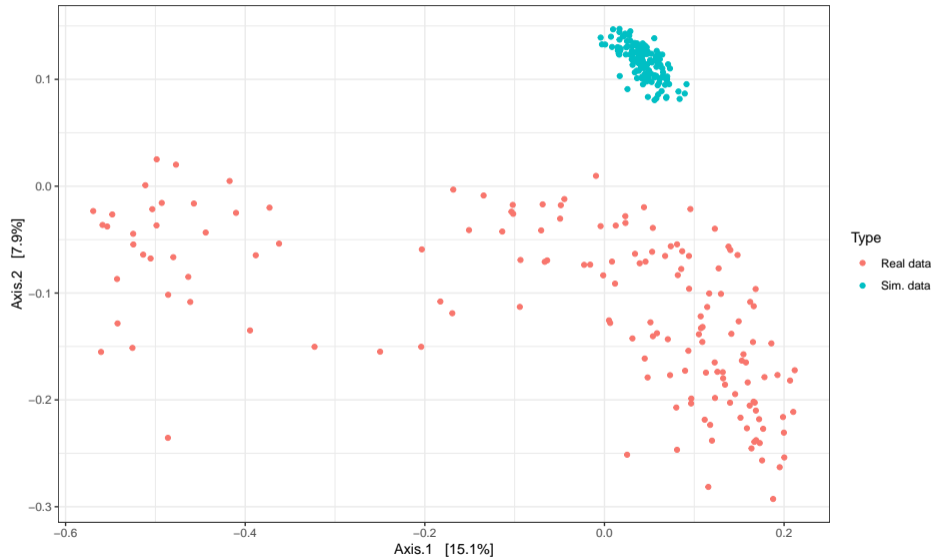
Modèle multinomial



Modèle multinomial

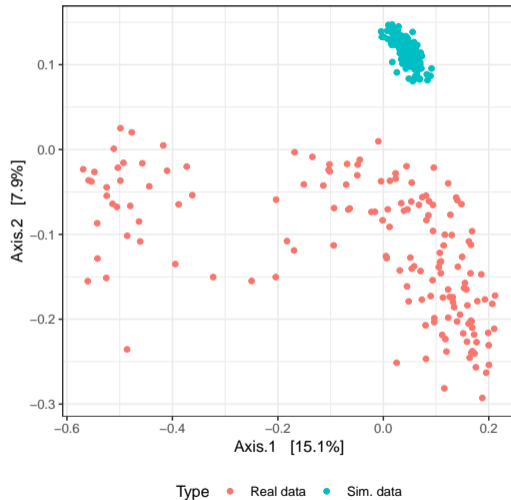


Exemple de modèle multinomial



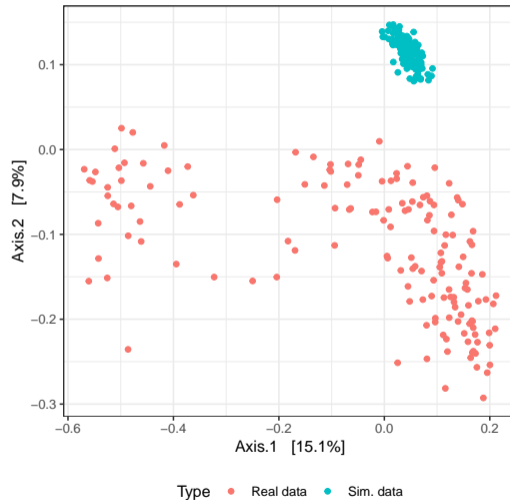
Hétérogénéité

- Pas assez d'hétérogénéité



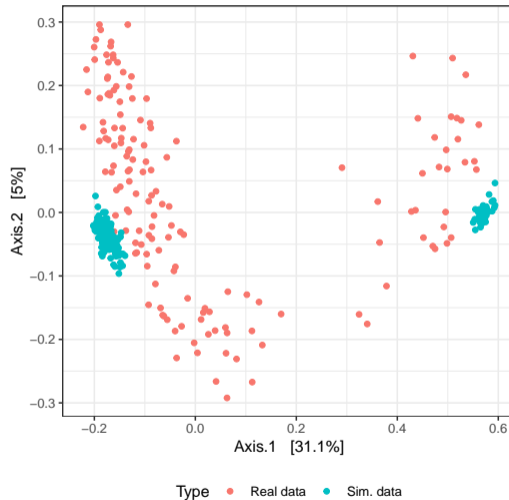
Hétérogénéité

- Pas assez d'hétérogénéité
↪ Ne s'ajuste qu'à une partie des données



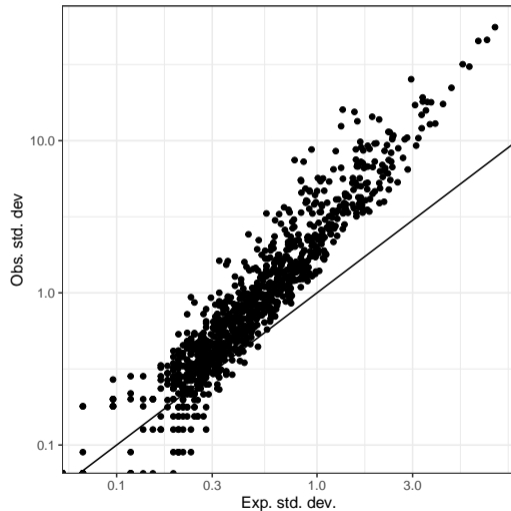
Hétérogénéité

- Pas assez d'hétérogénéité
↪ Ne s'ajuste qu'à une partie des données
- Variance faible



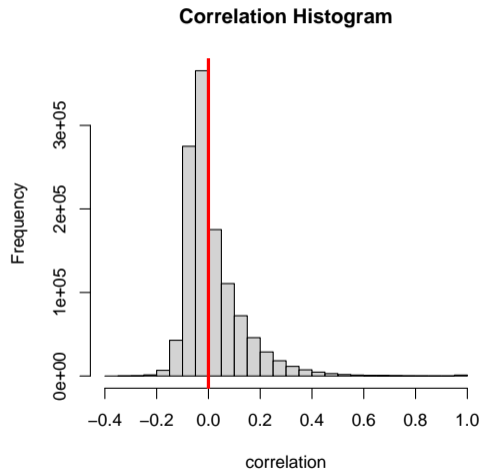
Hétérogénéité

- Pas assez d'hétérogénéité
↪ Ne s'ajuste qu'à une partie des données
- Variance faible
- Dispersion faible



Hétérogénéité

- Pas assez d'hétérogénéité
↪ Ne s'ajuste qu'à une partie des données
- Variance faible
- Dispersion faible
- Corrélations du mauvais signe



Avantages

- + **Parcimonie**: $p - 1$ paramètres pour p abondances
- + **Facile à estimer**
- + **Interprétable**

Avantages

- + **Parcimonie**: $p - 1$ paramètres pour p abondances
- + **Facile à estimer**
- + **Interprétable**

Inconvénients

- Mauvais pour **hétérogénéité**
- Mauvais pour la **dispersion** autour de la composition moyenne (\simeq variabilité biologique)
- Mauvais pour les **corrélations** entre espèces

1 Motivation

2 Modèles multinomiaux

- Multinomiale
- **Mélange de multinomiales**
- (Mélange de) Dirichlet-Multinomiale
- Latent Dirichlet Allocation

3 Modèles Log-Normaux

4 Applications



©Manfred Heyde

Intuition

- Chaque échantillon appartient à un parmi K groupes
- Le groupe k est caractérisé par sa composition π_k
- Un échantillon du groupe k a pour composition π_k
- Les comptages sont échantillonnés suivant une loi multinomiale

Intuition

- Chaque échantillon appartient à un parmi K groupes
- Le groupe k est caractérisé par sa composition π_k
- Un échantillon du groupe k a pour composition π_k
- Les comptages sont échantillonnés suivant une loi multinomiale

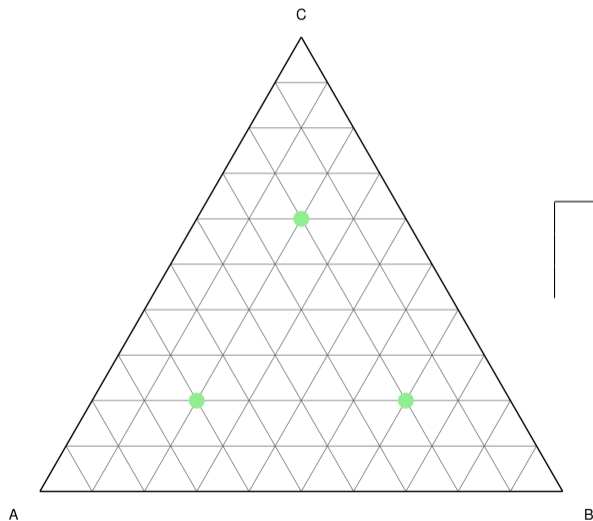
Modèle hiérarchique

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

où

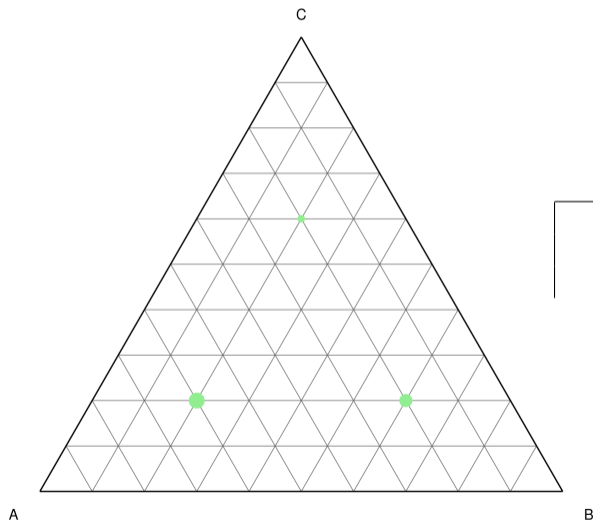
- $\alpha = (\alpha_1, \dots, \alpha_K)$ sont les proportions des K groupes,

Mélange de multinomiales



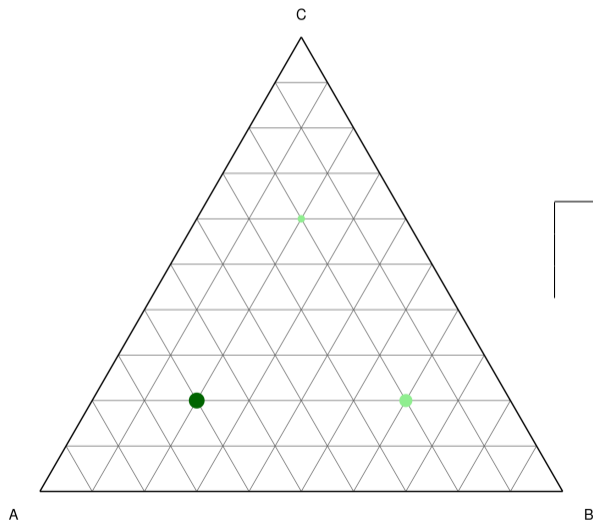
	A	B	C	α
π_1	0.6	0.2	0.2	
π_2	0.2	0.6	0.2	
π_3	0.2	0.2	0.6	

Mélange de multinomiales



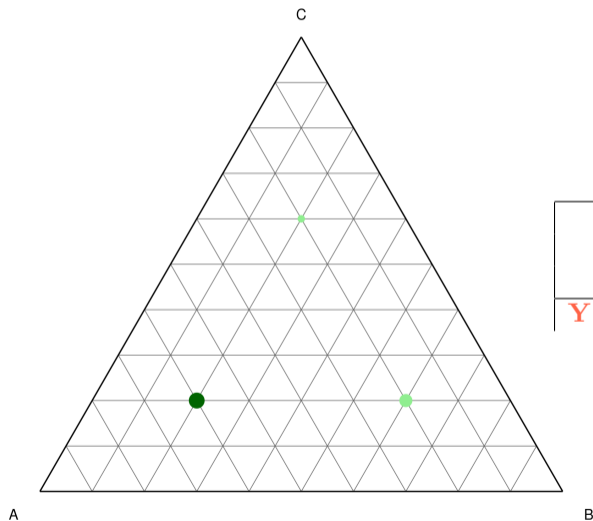
	<i>A</i>	<i>B</i>	<i>C</i>	α
π_1	0.6	0.2	0.2	0.5
π_2	0.2	0.6	0.2	0.4
π_3	0.2	0.2	0.6	0.1

Mélange de multinomiales



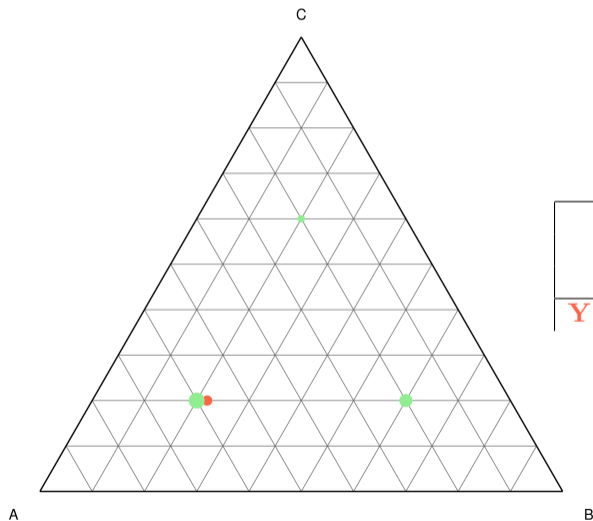
	<i>A</i>	<i>B</i>	<i>C</i>	α
π_1	0.6	0.2	0.2	0.5
π_2	0.2	0.6	0.2	0.4
π_3	0.2	0.2	0.6	0.1

Mélange de multinomiales



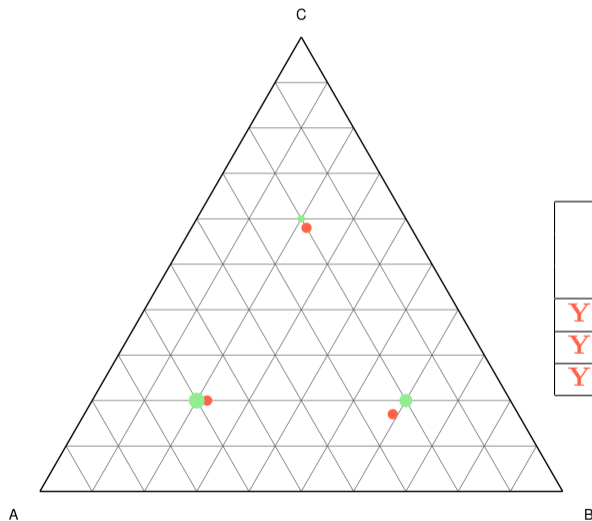
	A	B	C	α
π_1	0.6	0.2	0.2	0.5
π_2	0.2	0.6	0.2	0.4
π_3	0.2	0.2	0.6	0.1
$Y Z=1$	58	22	20	

Mélange de multinomiales



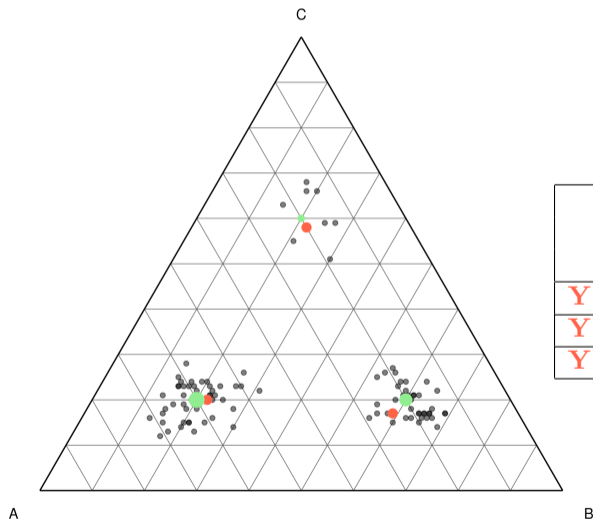
	A	B	C	α
π_1	0.6	0.2	0.2	0.5
π_2	0.2	0.6	0.2	0.4
π_3	0.2	0.2	0.6	0.1
$Y Z=1$	58	22	20	

Mélange de multinomiales



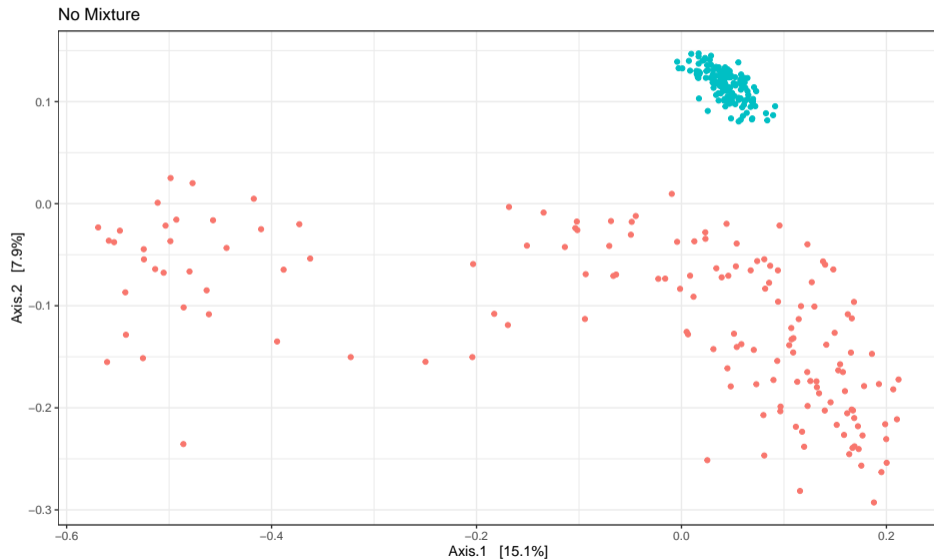
	A	B	C	α
π_1	0.6	0.2	0.2	0.5
π_2	0.2	0.6	0.2	0.4
π_3	0.2	0.2	0.6	0.1
$Y Z = 1$	58	22	20	
$Y Z = 2$	24	59	17	
$Y Z = 3$	20	22	58	

Mélange de multinomiales

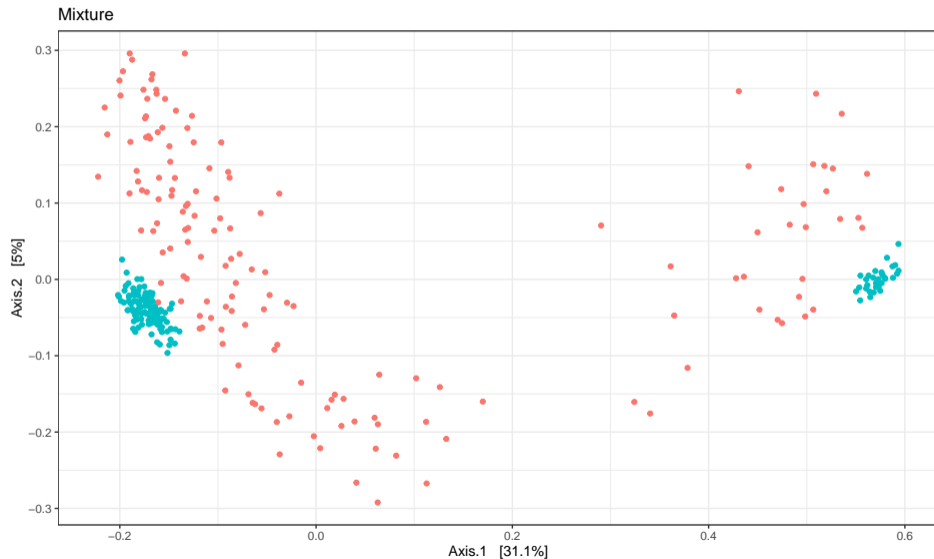


	A	B	C	α
π_1	0.6	0.2	0.2	0.5
π_2	0.2	0.6	0.2	0.4
π_3	0.2	0.2	0.6	0.1
$Y Z=1$	58	22	20	
$Y Z=2$	24	59	17	
$Y Z=3$	20	22	58	

Exemple de mélange de multinomial



Exemple de mélange de multinomial



Avantages

- + Bon pour **hétérogénéité**
- + **Parcimonieux**: $Kp - 1$ paramètres pour K groups
- + L'inférence est facile quand les groupes sont connus \rightsquigarrow simples moyennes

Avantages

- + Bon pour **hétérogénéité**
- + **Parcimonieux**: $Kp - 1$ paramètres pour K groups
- + L'inférence est facile quand les groupes sont connus \rightsquigarrow simples moyennes

Inconvénients

- Inférence moins facile quand les groupes sont inconnus
 \rightsquigarrow algorithme itératif EM
- Mauvais pour la **dispersion**
- Mauvais pour les **corrélations** entre OTUs

1 Motivation

2 Modèles multinomiaux

- Multinomiale
- Mélange de multinomiales
- (Mélange de) Dirichlet-Multinomiale
- Latent Dirichlet Allocation

3 Modèles Log-Normaux

4 Applications



Intuition

- π est la composition moyenne au **niveau de l'écosystème**

Intuition

- π est la composition moyenne au **niveau de l'écosystème**
- Échantillon i a sa **propre** composition π_i (**version bruitée** de π) \rightsquigarrow variabilité **biologique**

Intuition

- π est la composition moyenne au **niveau de l'écosystème**
- Échantillon i a sa **propre** composition π_i (**version bruitée** de π) \rightsquigarrow variabilité **biologique**
- N_i microbes tirés suivant $\mathcal{M}(1, \pi_i)$ \rightsquigarrow variabilité **technique** / **échantillonnage**

Intuition

- π est la composition moyenne au **niveau de l'écosystème**
- Échantillon i a sa **propre** composition π_i (**version bruitée** de π) \rightsquigarrow variabilité **biologique**
- N_i microbes tirés suivant $\mathcal{M}(1, \pi_i)$ \rightsquigarrow variabilité **technique** / **échantillonnage**

Modèle Hiérarchique

π	Composition de l'écosystème
$\pi_i \sim \mathcal{D}(\kappa\pi)$	Composition de l'échantillon
$Y_i \sim \mathcal{M}(N_i, \pi_i)$	Comptages observés

où $1/\kappa$ modélise le **niveau de variabilité** (petit $\kappa \rightsquigarrow$ grande dispersion)

Dirichlet - Multinomiale

Intuition

- π est la composition moyenne au **niveau de l'écosystème**
- Échantillon i a sa **propre** composition π_i (**version bruitée** de π) \rightsquigarrow variabilité **biologique**
- N_i microbes tirés suivant $\mathcal{M}(1, \pi_i)$ \rightsquigarrow variabilité **technique** / **échantillonnage**

Modèle Hiérarchique

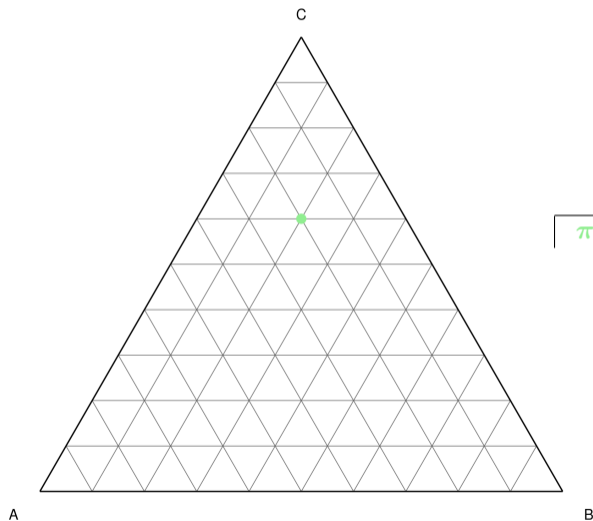
π	Composition de l'écosystème
$\pi_i \sim \mathcal{D}(\kappa\pi)$	Composition de l'échantillon
$\mathbf{Y}_i \sim \mathcal{M}(N_i, \pi_i)$	Comptages observés

où $1/\kappa$ modélise le **niveau de variabilité** (petit $\kappa \rightsquigarrow$ grande dispersion)

Couche de mélange

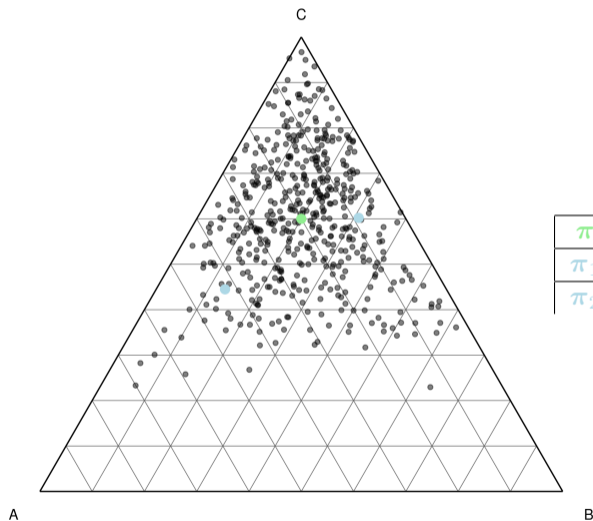
Ce modèle peut se **combiner** avec une couche de mélange

Distribution Dirichlet-Multinomiale



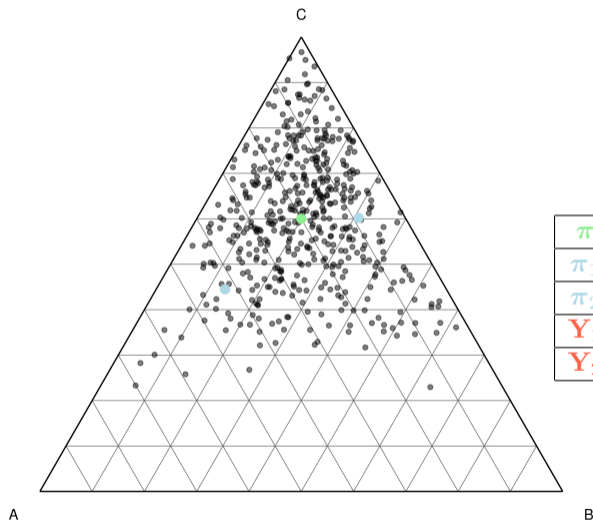
	<i>A</i>	<i>B</i>	<i>C</i>
π	0.2	0.2	0.6

Distribution Dirichlet-Multinomiale



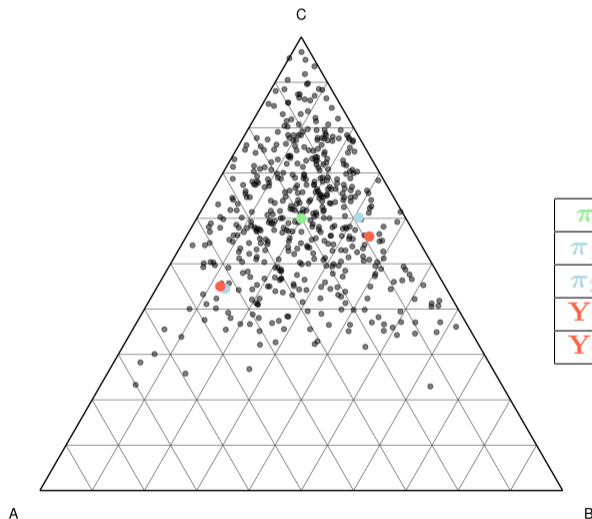
	A	B	C
π	0.2	0.2	0.6
π_1	0.089	0.309	0.602
π_2	0.423	0.132	0.445

Distribution Dirichlet-Multinomiale



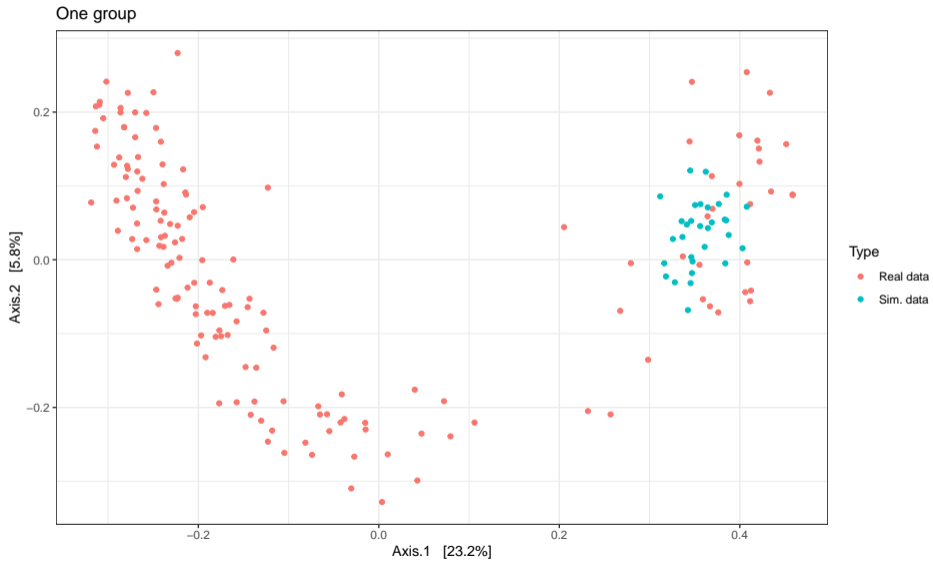
	A	B	C
π	0.2	0.2	0.6
π_1	0.089	0.309	0.602
π_2	0.423	0.132	0.445
Y_1	9	35	56
Y_2	43	12	45

Distribution Dirichlet-Multinomiale

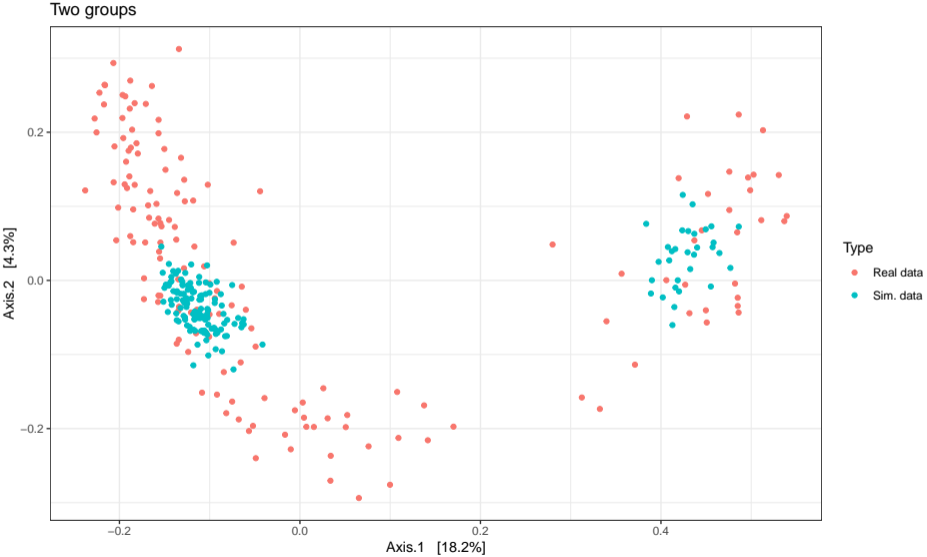


	A	B	C
π	0.2	0.2	0.6
π_1	0.089	0.309	0.602
π_2	0.423	0.132	0.445
Y_1	9	35	56
Y_2	43	12	45

Exemple de Dirichlet-Multinomiale



Exemple de Dirichlet-Multinomiale (II)



Avantages

- + Bon pour l'**hétérogénéité**
- + Moyen pour la **dispersion**
- + **Parcimonieux**: $K(p + 1) - 1$ paramètres pour K groupes

Avantages

- + Bon pour l'**hétérogénéité**
- + Moyen pour la **dispersion**
- + **Parcimonieux**: $K(p + 1) - 1$ paramètres pour K groupes

Inconvénients

- **Inférence** non triviale
 - Groupes connus \rightsquigarrow descente de gradient
 - Groupes inconnus \rightsquigarrow Algorithme itérative EM + descente de gradient
- Mauvais pour les **corrélations** entre OTUs

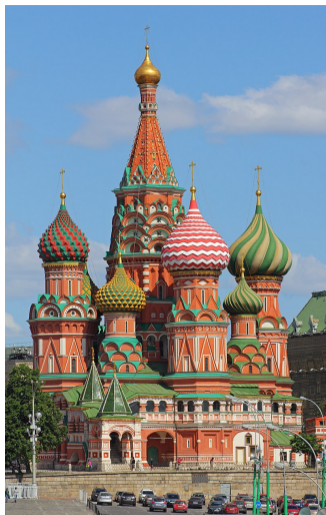
1 Motivation

2 Modèles multinomiaux

- Multinomiale
- Mélange de multinomiales
- (Mélange de) Dirichlet-Multinomiale
- Latent Dirichlet Allocation

3 Modèles Log-Normaux

4 Applications



Intuition

- Il y a K archétypes d'écosystèmes $1, \dots, K$

Intuition

- Il y a K **archétypes** d'écosystèmes $1, \dots, K$
- Chaque archétype a sa **propre composition** π_k

Intuition

- Il y a K **archétypes** d'écosystèmes $1, \dots, K$
- Chaque archétype a sa **propre composition** π_k
- Chaque échantillon \mathbf{Y} est le mélange de **plusieurs archétypes** en proportions $(\theta_1, \dots, \theta_K)$

Intuition

- Il y a K **archétypes** d'écosystèmes $1, \dots, K$
- Chaque archétype a sa **propre composition** π_k
- Chaque échantillon \mathbf{Y} est le mélange de **plusieurs archétypes** en proportions $(\theta_1, \dots, \theta_K)$
- $\theta_k N$ microbes sont échantillonnés à partir d'une **version bruitée** de π_k

Intuition

- Il y a K **archétypes** d'écosystèmes $1, \dots, K$
- Chaque archétype a sa **propre composition** π_k
- Chaque échantillon \mathbf{Y} est le mélange de **plusieurs archétypes** en proportions $(\theta_1, \dots, \theta_K)$
- $\theta_k N$ microbes sont échantillonnés à partir d'une **version bruitée** de π_k

Modèle Hiérarchique

π_1, \dots, π_K

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

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

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

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

Composition des archétypes

Proportion des archétypes dans un échantillon

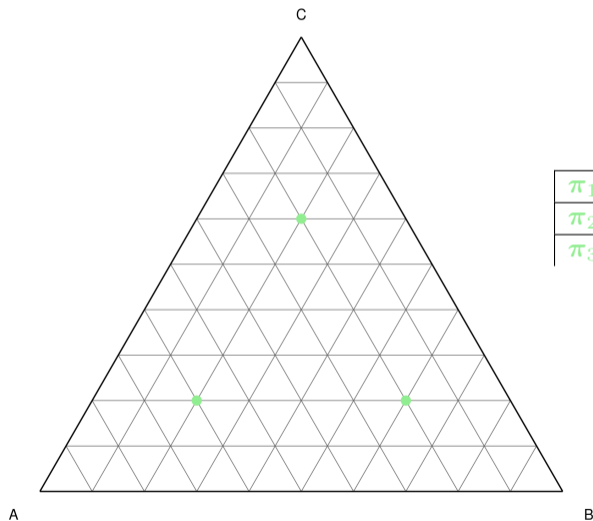
Version bruitée de π_k

Archétype d'origine du comptage i

espèce du comptage i

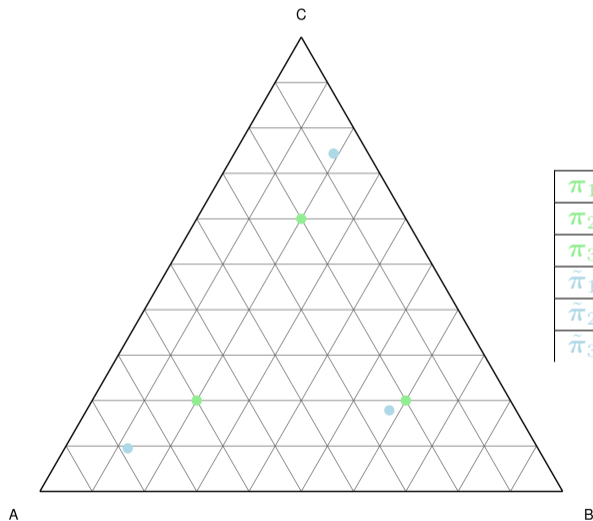
où κ et les κ_k contrôlent le niveau de variabilité.

Latent Dirichlet Allocation



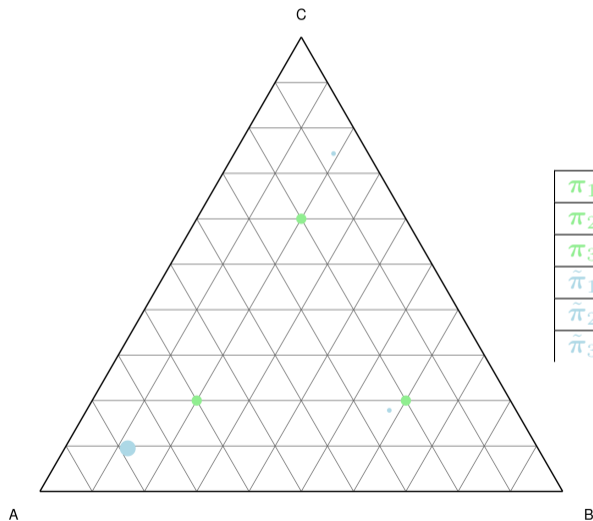
	<i>A</i>	<i>B</i>	<i>C</i>	θ
π_1	0.6	0.2	0.2	
π_2	0.2	0.6	0.2	
π_3	0.2	0.2	0.6	

Latent Dirichlet Allocation



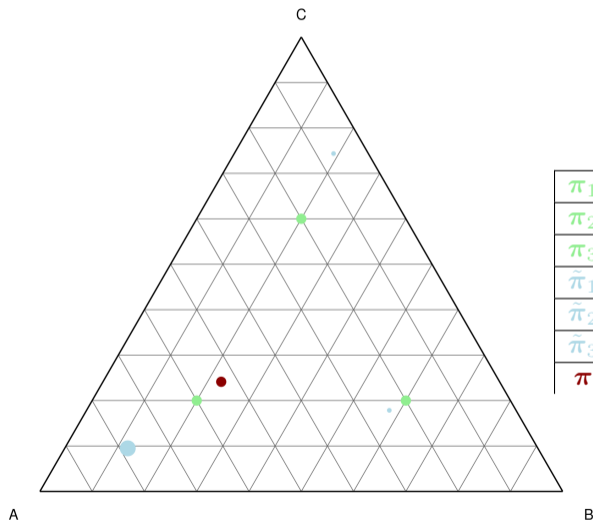
	A	B	C	θ
π_1	0.6	0.2	0.2	
π_2	0.2	0.6	0.2	
π_3	0.2	0.2	0.6	
$\tilde{\pi}_1$	0.784	0.121	0.095	
$\tilde{\pi}_2$	0.242	0.579	0.179	
$\tilde{\pi}_3$	0.423	0.132	0.445	

Latent Dirichlet Allocation



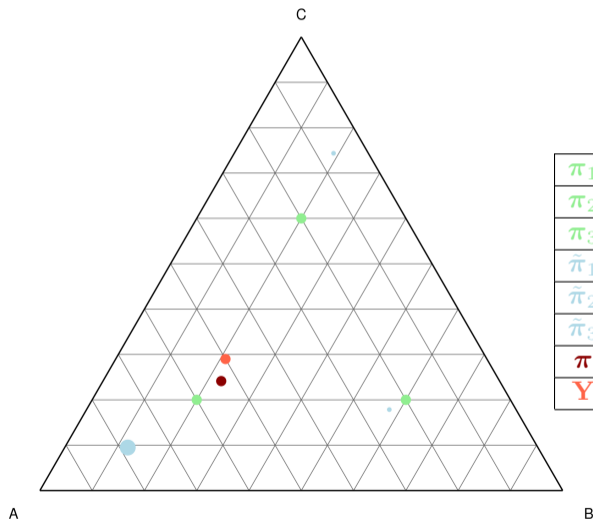
	A	B	C	θ
π_1	0.6	0.2	0.2	
π_2	0.2	0.6	0.2	
π_3	0.2	0.2	0.6	
$\tilde{\pi}_1$	0.784	0.121	0.095	0.6
$\tilde{\pi}_2$	0.242	0.579	0.179	0.2
$\tilde{\pi}_3$	0.423	0.132	0.445	0.2

Latent Dirichlet Allocation



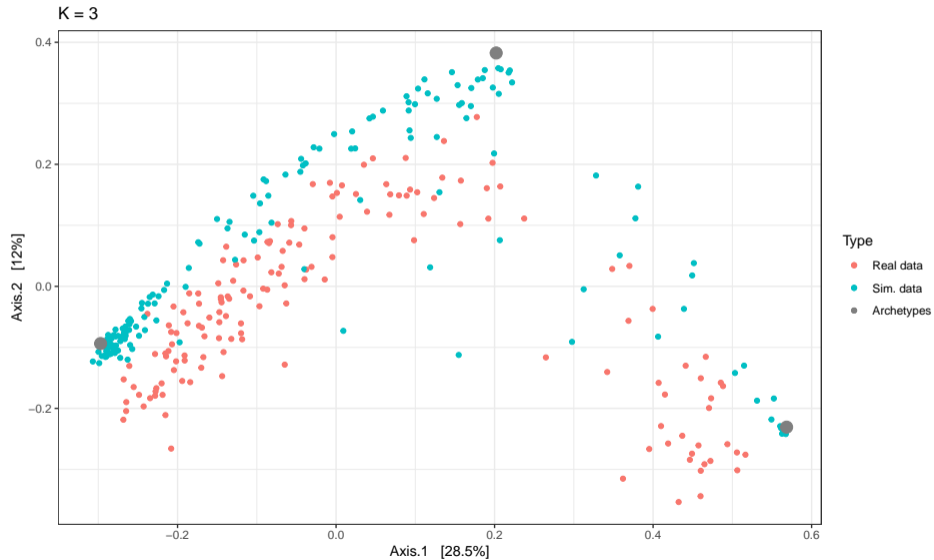
	A	B	C	θ
π_1	0.6	0.2	0.2	
π_2	0.2	0.6	0.2	
π_3	0.2	0.2	0.6	
$\tilde{\pi}_1$	0.784	0.121	0.095	0.6
$\tilde{\pi}_2$	0.242	0.579	0.179	0.2
$\tilde{\pi}_3$	0.423	0.132	0.445	0.2
π	0.532	0.226	0.241	

Latent Dirichlet Allocation

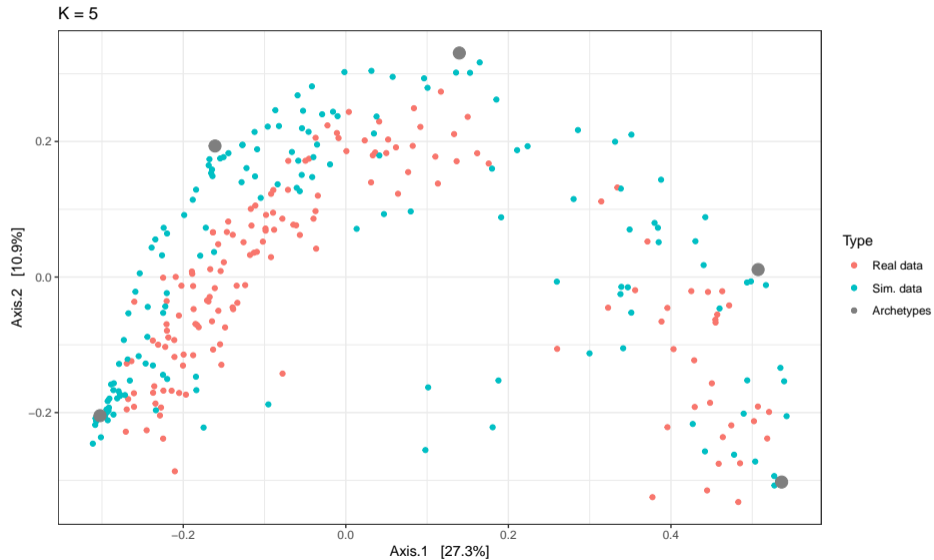


	A	B	C	θ
π_1	0.6	0.2	0.2	
π_2	0.2	0.6	0.2	
π_3	0.2	0.2	0.6	
$\tilde{\pi}_1$	0.784	0.121	0.095	0.6
$\tilde{\pi}_2$	0.242	0.579	0.179	0.2
$\tilde{\pi}_3$	0.423	0.132	0.445	0.2
π	0.532	0.226	0.241	
Y	54	18	28	

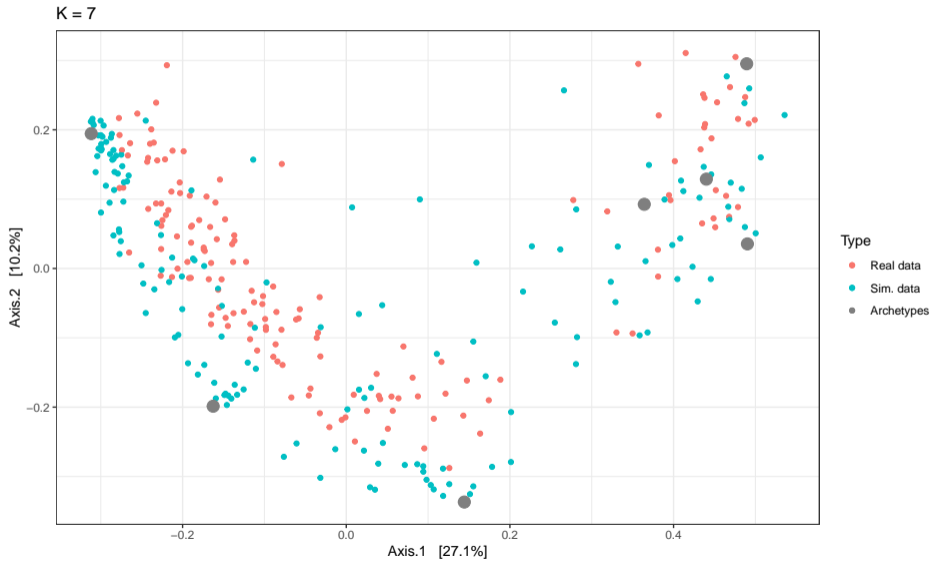
Exemple de Latent Dirichlet Allocation



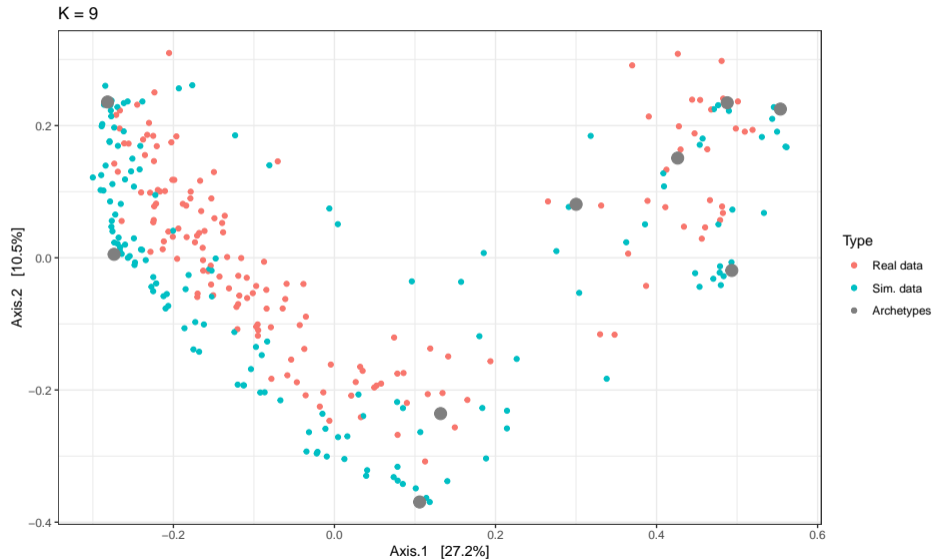
Exemple de Latent Dirichlet Allocation



Exemple de Latent Dirichlet Allocation



Exemple de Latent Dirichlet Allocation



Avantages

- + Bon pour l'**hétérogénéité**
- + Bon pour la **dispersion**
- + **Parcimonieux**: $K(p + 1)$ paramètres pour K archétypes

Avantages

- + Bon pour l'**hétérogénéité**
- + Bon pour la **dispersion**
- + **Parcimonieux**: $K(p + 1)$ paramètres pour K archétypes

Inconvénients

- **Inférence** moyennement complexe
↪ algorithme EM + descente de gradient / Échantillonnage de Gibbs
- **Interprétation** complexe ↪ les archétypes ne **sont pas des groupes**
- Mauvaises **corrélations** entre espèces

Les modèles multinomiaux sont **bons** pour

- modéliser les **compositions moyennes**;
- modéliser la **dispersion** autour de ces moyennes;
- modéliser l'**hétérogénéité**;
- en utilisant (relativement) peu de paramètres

Les modèles multinomiaux sont **bons** pour

- modéliser les **compositions moyennes**;
- modéliser la **dispersion** autour de ces moyennes;
- modéliser l'**hétérogénéité**;
- en utilisant (relativement) peu de paramètres

Les modèles multinomiaux sont **mauvais** pour

- modéliser les **interactions** entre espèces;
- prendre en compte les **covariables**;
- intégrer les données issues de **différentes sources** (e.g. 16S, ITS)

1 Motivation

2 Modèles multinomiaux

3 Modèles Log-Normaux

- Multinomiale Log-Normale
- Poisson Log-Normale

4 Applications

Modéliser les corrélations

Les modèles gaussiens multivariés sont le standard *de facto* pour modéliser les corrélations.

Modéliser les corrélations

Les modèles gaussiens multivariés sont le standard *de facto* pour modéliser les corrélations.

Pour des variables continues

- Les p variables \mathbf{Y}_i (e.g. abondances d'espèces) sont expliquées
- par les valeurs de d covariables \mathbf{X}_i et de p offsets \mathbf{O}_i

$$\mathbf{Y}_i = \underbrace{\mathbf{X}_i \mathbf{B}}_{\text{capture les covariables}} + \underbrace{\mathbf{O}_i}_{\text{captures les efforts d'observation}} + \varepsilon_i, \quad \varepsilon_i \sim \mathcal{N}(\mathbf{0}_p, \underbrace{\Sigma}_{\text{dépendances entre espèces}})$$

+ covariance nulle \Leftrightarrow indépendance \rightsquigarrow les espèces non-corrélées n'interagissent pas.

Modéliser les corrélations

Les modèles gaussiens multivariés sont le standard *de facto* pour modéliser les corrélations.

Pour des variables continues

- Les p variables \mathbf{Y}_i (e.g. abondances d'espèces) sont expliquées
- par les valeurs de d covariables \mathbf{X}_i et de p offsets \mathbf{O}_i

$$\mathbf{Y}_i = \underbrace{\mathbf{X}_i \mathbf{B}}_{\text{capture les covariables}} + \underbrace{\mathbf{O}_i}_{\text{captures les efforts d'observation}} + \varepsilon_i, \quad \varepsilon_i \sim \mathcal{N}(\mathbf{0}_p, \underbrace{\Sigma}_{\text{dépendances entre espèces}})$$

+ ~~covariance nulle~~ \Leftrightarrow ~~indépendance~~ \rightsquigarrow ~~les espèces non-corrélées n'interagissent pas.~~

Mais les abondances ne sont pas gaussiennes...

Modéliser les corrélations

Les modèles gaussiens multivariés sont le standard *de facto* pour modéliser les corrélations.

Pour des variables continues

- Les p variables \mathbf{Y}_i (e.g. abondances d'espèces) sont expliquées
- par les valeurs de d covariables \mathbf{X}_i et de p offsets \mathbf{O}_i

$$\mathbf{Y}_i = \underbrace{\mathbf{X}_i \mathbf{B}}_{\text{capture les covariables}} + \underbrace{\mathbf{O}_i}_{\text{captures les efforts d'observation}} + \varepsilon_i, \varepsilon_i \sim \mathcal{N}(\mathbf{0}_p, \underbrace{\Sigma}_{\text{dépendances entre espèces}})$$

+ ~~covariance nulle~~ \Leftrightarrow ~~indépendance~~ \rightsquigarrow ~~les espèces non-corrélées n'interagissent pas.~~

Mais les abondances ne sont pas gaussiennes...

Modèle à variable latente avec une couche *latente* gaussienne et une couche *observée* comptage

1 Motivation

2 Modèles multinomiaux

3 Modèles Log-Normaux

- Multinomiale Log-Normale
- Poisson Log-Normale

4 Applications



Intuition

- La couche latente modélise des **abondances de bases** \mathbf{z}

Intuition

- La couche latente modélise des **abondances de bases** \mathbf{z}
- Ces abondances sont **transformées** en une **composition** moyenne $\boldsymbol{\pi}$

Intuition

- La couche latente modélise des **abondances de bases** \mathbf{z}
- Ces abondances sont **transformées** en une **composition** moyenne π
- N microbes sont **tirés** suivant une multinomiale de paramètre π

Intuition

- La couche latente modélise des **abondances de bases** \mathbf{z}
- Ces abondances sont **transformées** en une **composition** moyenne $\boldsymbol{\pi}$
- N microbes sont **tirés** suivant une multinomiale de paramètre $\boldsymbol{\pi}$

Modèle hiérarchique

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

abondances de base

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

composition moyenne

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

comptages observés

Multinomiale Log-Normale

Intuition

- La couche latente modélise des **abondances de bases** \mathbf{z}
- Ces abondances sont **transformées** en une **composition** moyenne $\boldsymbol{\pi}$
- N microbes sont **tirés** suivant une multinomiale de paramètre $\boldsymbol{\pi}$

Modèle hiérarchique

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

abondances de base

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

composition moyenne

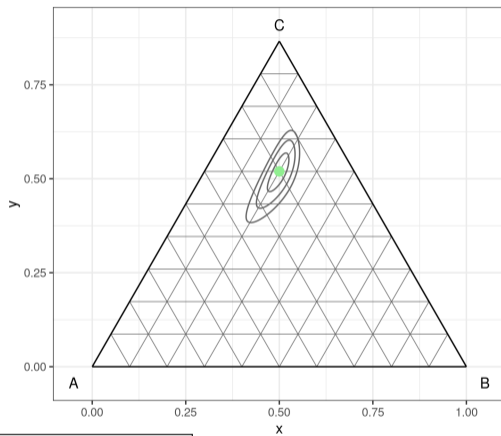
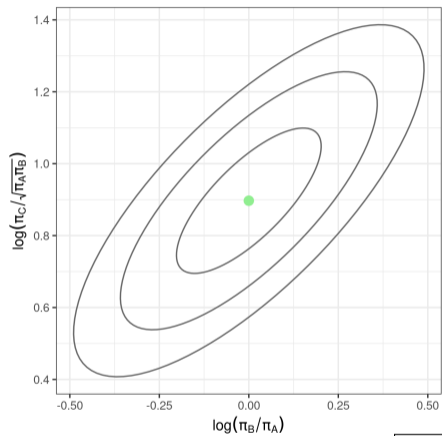
$$\mathbf{Y} \sim \mathcal{M}(N, \boldsymbol{\pi})$$

comptages observés

Couche de mélange

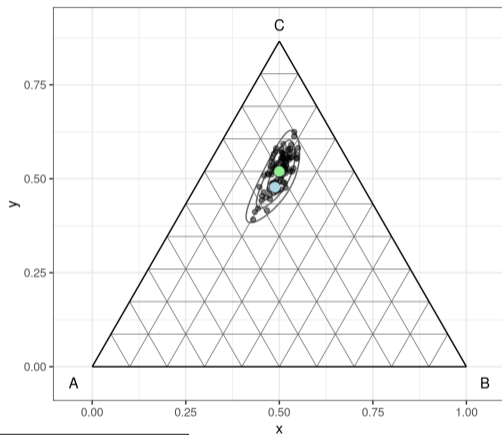
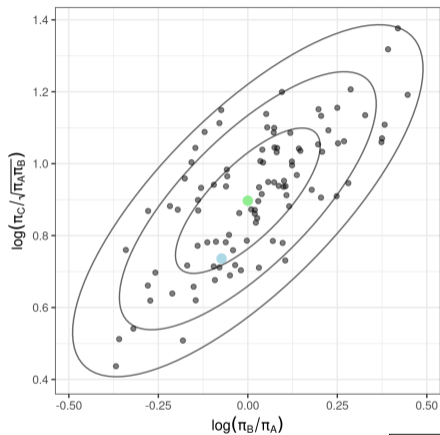
On peut rajouter une couche de mélange au modèle pour l'hétérogénéité.

Multinomiale Log-Normale



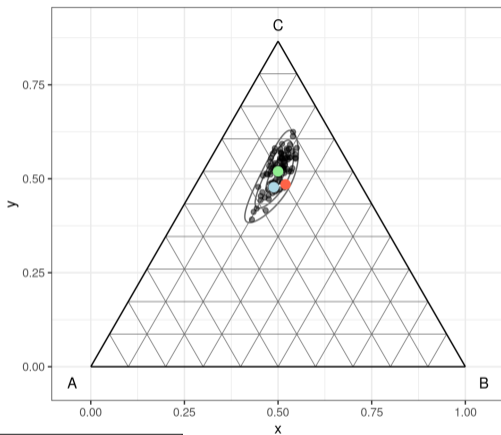
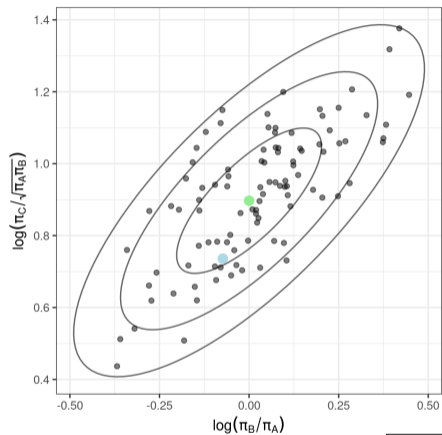
	A	B	C
π	0.2	0.2	0.6

Multinomiale Log-Normale



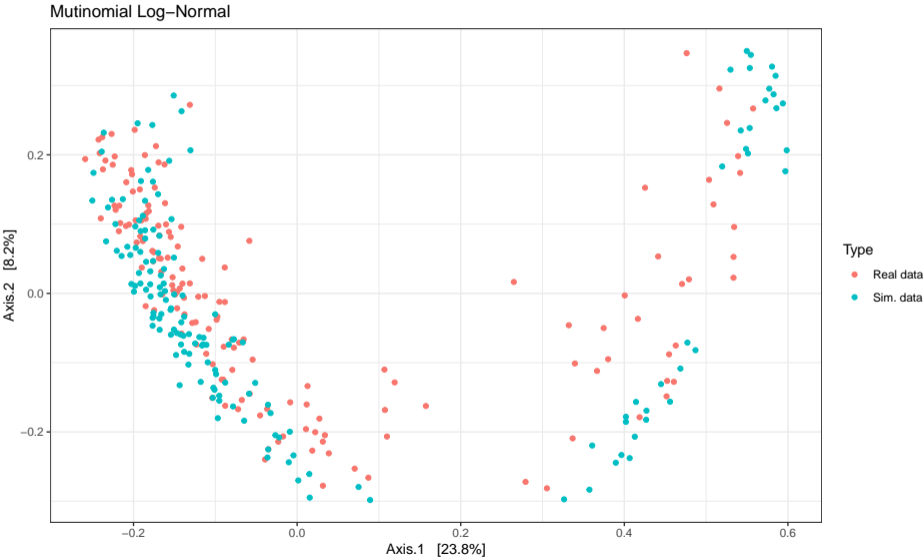
	A	B	C
π	0.2	0.2	0.6
π_1	0.235	0.213	0.552

Multinomiale Log-Normale



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

Exemple de Multinomiale Log-Normale



Avantages

- + Bon pour l'**hétérogénéité**
- + Bon pour la **dispersion**
- + Bon pour les **corrélations** entre espèces

Avantages

- + Bon pour l'**hétérogénéité**
- + Bon pour la **dispersion**
- + Bon pour les **corrélations** entre espèces

Inconvénients

- Modèle non-**parcimonieux**: $p(p + 3)/2$ paramètres
- **Inférence** complexe
 - ↔ Algorithme EM itératif / Échantillonnage de Gibbs
- Modélisation faite au niveau des **compositions** et pas des comptages.

1 Motivation

2 Modèles multinomiaux

3 Modèles Log-Normaux

- Multinomiale Log-Normale
- Poisson Log-Normale

4 Applications



Intuition

- La couche latente modélise des bases z

Intuition

- La couche latente modélise des **bases** z
- Les bases sont *transformées* en **comptages** moyens

Intuition

- La couche latente modélise des **bases** z
- Les bases sont *transformées* en **comptages** moyens
- Les comptages des microbes sont *tirés* suivant une distribution de Poisson

Intuition

- La couche latente modélise des **bases** \mathbf{z}
- Les bases sont *transformées* en **comptages** moyens
- Les comptages des microbes sont *tirés* suivant une distribution de Poisson

Modèle hiérarchique

$\mathbf{z} \sim \mathcal{N}(\boldsymbol{\mu}, \boldsymbol{\Sigma})$	Bases
$\lambda_j \mathbf{z} = e^{z_j}$	Comptage moyen de l'espèce j
$Y_j z_j \sim \mathcal{P}(e^{z_j})$	Comptage observé de l'espèce j

Intuition

- La couche latente modélise des **bases** \mathbf{z}
- Les bases sont *transformées* en **comptages** moyens
- Les comptages des microbes sont *tirés* suivant une distribution de Poisson

Modèle hiérarchique

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

Bases

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

Comptage moyen de l'espèce j

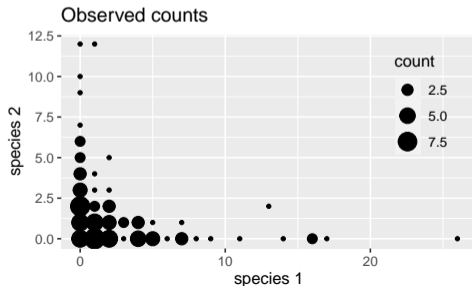
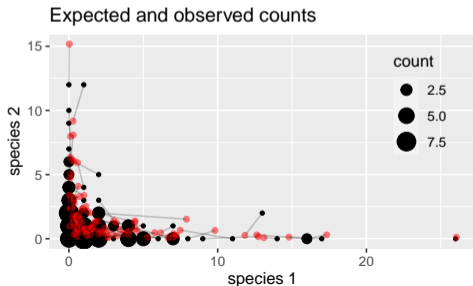
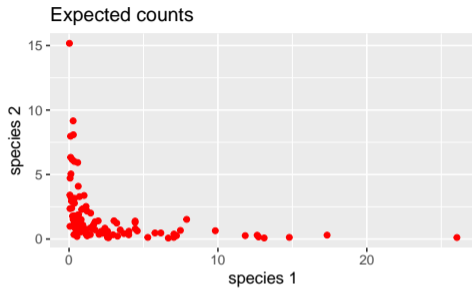
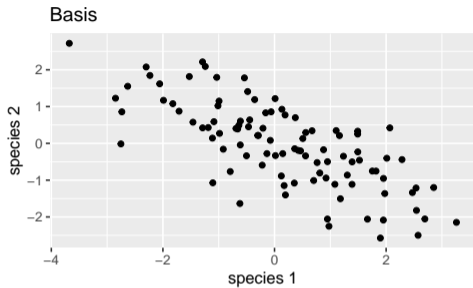
$$Y_j | z_j \sim \mathcal{P}(e^{z_j})$$

Comptage observé de l'espèce j

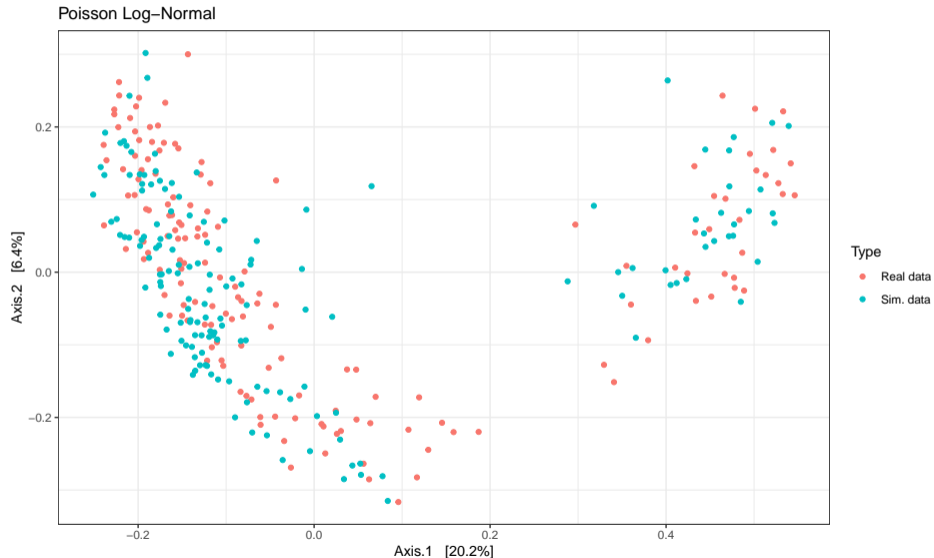
Couche mélange

On peut rajouter une couche de mélange au modèle pour l'hétérogénéité.

Interprétation géométrique



Exemple de Poisson Log-Normale



Avantages

- + Bon pour l'**hétérogénéité**
- + Bon pour la **dispersion**
- + Bon pour les **corrélations** entre espèces
- + Modélisation faite au niveau des **comptages**
 - ↪ les comptages peuvent être à des **échelles différentes** et provenir de **différentes sources**

Avantages et inconvénients

Avantages

- + Bon pour l'**hétérogénéité**
- + Bon pour la **dispersion**
- + Bon pour les **corrélations** entre espèces
- + Modélisation faite au niveau des **comptages**
 - ↪ les comptages peuvent être à des **échelles différentes** et provenir de **différentes sources**

Inconvénients

- Modèle non-**parcimonieux**: $p(p + 3)/2$ paramètres
- **Inférence** complexe
 - ↪ Algorithme EM itératif + descente de gradient
- Comptage total seulement contrôlé en **moyenne**

Les modèles Log-Normaux sont **bons** pour

- modéliser les **compositions moyennes**;
- modéliser la **dispersion** autour de ces moyennes;
- modéliser l'**hétérogénéité**;
- modéliser les **interactions** entre espèces;
- corriger les effets de **covariables** à l'aide d'un modèle linéaire

Les modèles Log-Normaux sont **bons** pour

- modéliser les **compositions moyennes**;
- modéliser la **dispersion** autour de ces moyennes;
- modéliser l'**hétérogénéité**;
- modéliser les **interactions** entre espèces;
- corriger les effets de **covariables** à l'aide d'un modèle linéaire

Les modèles log-normaux sont **mauvais** pour

- leur **grand nombre** de paramètres
- la **complexité** de leurs méthodes d'inférence

Résumé intermédiaire

Les modèles Log-Normaux sont **bons** pour

- modéliser les **compositions moyennes**;
- modéliser la **dispersion** autour de ces moyennes;
- modéliser l'**hétérogénéité**;
- modéliser les **interactions** entre espèces;
- corriger les effets de **covariables** à l'aide d'un modèle linéaire

Les modèles log-normaux sont **mauvais** pour

- leur **grand nombre** de paramètres
- la **complexité** de leurs méthodes d'inférence

- Les modèles MLN sont plus faciles à **interpréter** (compositions)
- Les modèles PLN permettent de mixer les données de **différentes sources** (16S, ITS, etc.)

1 Motivation

2 Modèles multinomiaux

3 Modèles Log-Normaux

4 Applications

- ACP
- Analyse discriminante
- Inférence de réseaux

PLN: un modèle flexible pour prendre en compte:

- l'hétérogénéité et les comptages moyens (\simeq moment de premier ordre)
- la dispersion et les corrélations entre espèces (\simeq moment de second ordre)
- des covariables et des facteurs de confusion
- des comptages issus de différentes sources

PLN: un modèle flexible pour prendre en compte:

- l'hétérogénéité et les comptages moyens (\simeq moment de premier ordre)
- la dispersion et les corrélations entre espèces (\simeq moment de second ordre)
- des covariables et des facteurs de confusion
- des comptages issus de différentes sources

Permet de se ramener à des analyses multivariées *classiques*:

Idée: Mettre des contraintes sur le modèle

- **ACP** \rightsquigarrow faible rang sur Σ
- **Analyse discriminante** \rightsquigarrow structure (connue) de groupe sur μ
- **Inférence de réseaux** $\rightsquigarrow \Sigma^{-1}$ parcimonieuse
- **Modèles de mélange** \rightsquigarrow structure (inconnue) de groupe sur μ
- *etc.*

1 Motivation

2 Modèles multinomiaux

3 Modèles Log-Normaux

4 Applications

- ACP
- Analyse discriminante
- Inférence de réseaux

Réduction de dimension et visualisation. Tâches typiques en analyse multivariée

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

↪ Trouver une base de faible dimension (axes de l'ACP) pour représenter la covariance latente

Ajustement d'un modèle PLNPCA avec offset et covariables.

```
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: visualisation

PLN-PCA séparent les échantillons suivant l'arbre d'origine

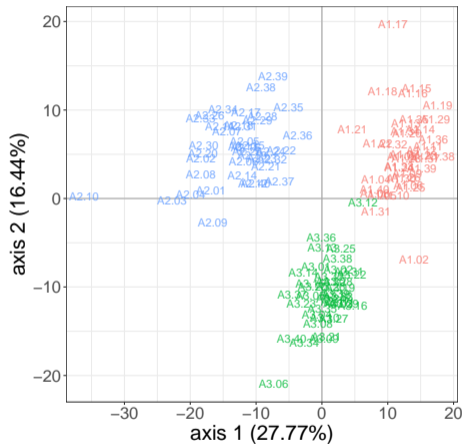


Figure: Offset seul



Offset + covariables

- tree
- a intermediate
 - a resistant
 - a susceptible

PCA: visualisation II

L'introduction de covariables fait émerger des structures

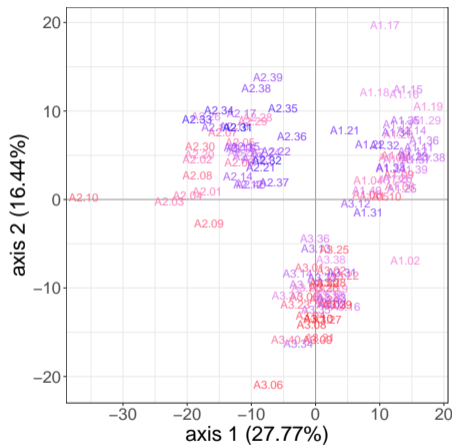
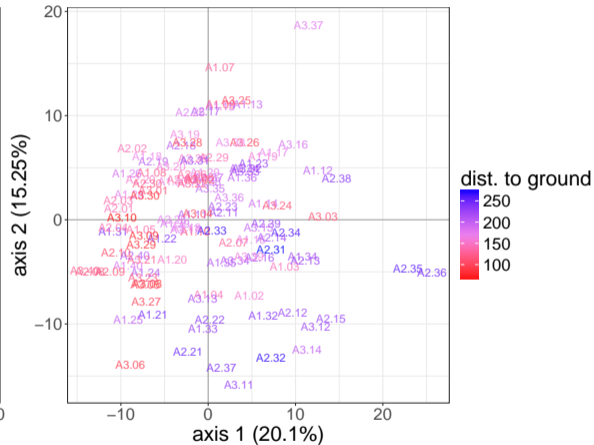


Figure: Offset seul



Offset + covariables

- 1 Motivation
- 2 Modèles multinomiaux
- 3 Modèles Log-Normaux
- 4 Applications**
 - ACP
 - Analyse discriminante
 - Inférence de réseaux

Estimation des modèles PLN-LDA

Trouver la combinaison linéaire qui sépare les groupes

Estimation du modèle, avec offsets et covariables.

```
myLDA_tree <- PLNLDA(Abundance ~ offset(log(Offset)), grouping = tree, data = oaks)
```

```
##  
## Performing discriminant Analysis...  
## DONE!
```

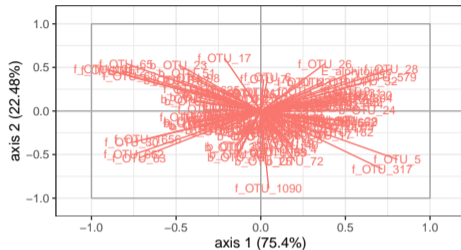
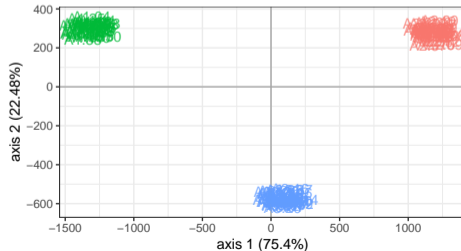
```
myLDA_tree$plot_LDA()
```

Analyse discriminante sur l'arbre d'origine

Axes contribution

axis 1 : 75.4%

axis 2 : 22.48%



classification

- a susceptible
- a intermediate
- a resistant

Erreur de prédiction (10-fold validation croisée)

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

- 1 Motivation
- 2 Modèles multinomiaux
- 3 Modèles Log-Normaux
- 4 Applications
 - ACP
 - Analyse discriminante
 - Inférence de réseaux

Même cadre:

$$\begin{array}{lll} \text{Espace des paramètres } \mathbb{R}^p & \mathbb{E}[\mathbf{Z}_i] = \mu_i & \mathbf{Z}_i \sim \mathcal{N}(\mu_i, \Sigma = \mathbf{B}^\top \mathbf{B}) \\ \text{Espace des observations } \mathbb{N}^p & Y_{ij} | Z_{ij} \text{ indep.} & Y_{ij} | Z_{ij} \sim \mathcal{P}(e^{Z_{ij}}) \end{array}$$

où \mathbf{B} est de (faible-)rang $q \ll p$.

Même cadre, contraintes différentes:

$$\begin{array}{lll} \text{Espace des paramètres } \mathbb{R}^p & \mathbb{E}[\mathbf{Z}_i] = \mu_i & \mathbf{Z}_i \sim \mathcal{N}(\mu_i, \Sigma = \Omega^{-1}) \\ \text{Espace des observations } \mathbb{N}^p & Y_{ij}|Z_{ij} \text{ indep.} & Y_{ij}|Z_{ij} \sim \mathcal{P}(e^{Z_{ij}}) \end{array}$$

où Ω est creuse et reflète la topologie du réseau.

EM variationnel Maximiser une borne inf J de la vraisemblance sous contrainte de parcimonie sur Ω :

$$\arg \max_{\mathbf{M}, \mathbf{S}, \Theta, \Omega} J(\mathbf{M}, \mathbf{S}, \Theta, \Omega) - \lambda \|\Omega\|_1$$

EM variationnel Maximiser une borne inf J de la vraisemblance sous contrainte de parcimonie sur Ω :

$$\arg \max_{\mathbf{M}, \mathbf{S}, \Theta, \Omega} J(\mathbf{M}, \mathbf{S}, \Theta, \Omega) - \lambda \|\Omega\|_1$$

Itérer jusqu'à convergence entre

- **step 1**: Mise à jour des paramètres $(\mathbf{M}, \mathbf{S}, \Theta)$ par descente de gradient

EM variationnel Maximiser une borne inf J de la vraisemblance sous contrainte de parcimonie sur Ω :

$$\arg \max_{\mathbf{M}, \mathbf{S}, \Theta, \Omega} J(\mathbf{M}, \mathbf{S}, \Theta, \Omega) - \lambda \|\Omega\|_1$$

Itérer jusqu'à convergence entre

- **step 1**: Mise à jour des paramètres $(\mathbf{M}, \mathbf{S}, \Theta)$ par descente de gradient
- **step 2**: Mise à jour du réseau Ω par lasso graphique (*glasso*)

EM variationnel Maximiser une borne inf J de la vraisemblance sous contrainte de parcimonie sur Ω :

$$\arg \max_{\mathbf{M}, \mathbf{S}, \Theta, \Omega} J(\mathbf{M}, \mathbf{S}, \Theta, \Omega) - \lambda \|\Omega\|_1$$

Itérer jusqu'à convergence entre

- **step 1**: Mise à jour des paramètres $(\mathbf{M}, \mathbf{S}, \Theta)$ par descente de gradient
- **step 2**: Mise à jour du réseau Ω par lasso graphique (*glasso*)

λ sélectionné en utilisant la méthode de rééchantillonnage StARS [LRW10].

Un exemple miroir de celui de Thibaud

Données

Suffrages de 63 242 bureaux de vote lors des élections présidentielles de 2017 (premier tour)

Objectif

Reconstruire les corrélations partielles entre candidats.

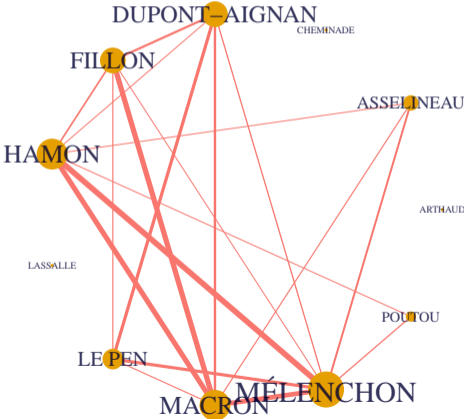
Méthode

On ajuste deux modèles:

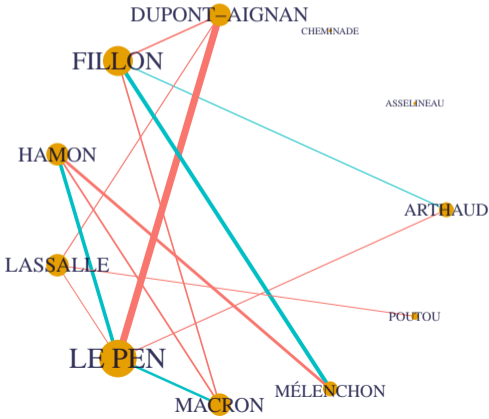
- Un modèle naïf sans offset (pas de prise en compte de la compositionnalité)
- Un modèle moins naïf avec la **taille du bureau de vote** comme offset.

La compositionnalité compte

No offset



Offset

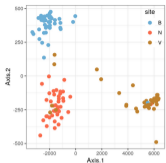


Résumé PLN = modèle générique pour les données de comptages

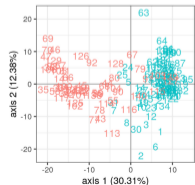
- Prise en compte des effets compositionnels (via les *offset*)
- Prise en compte des covariables
- Modélisation statistique flexible
- PLNmodels R-package

Extension

- Ajout des zéros essentiels ("zéro-inflation")
- Extensions: ACP parcimonieuse
- Intervalle de confiance et tests
- Données manquantes. . .



Classification accuracy: 94.3%
(work with S. Even)

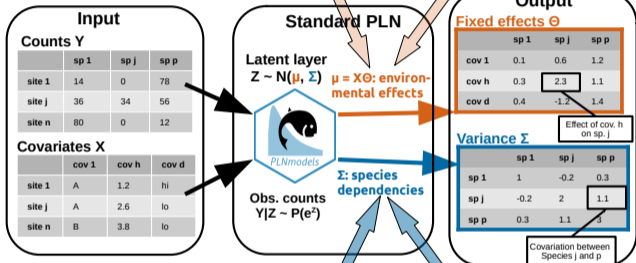


Work with N. Peyrard and M.-J. Cros

PLN-LDA: compare sites
Goal: find **systematic differences** between sites in different classes.
Constraint: $\mu = \mu_k$ if site in known class k

PLN-mixture: find groups
Goal: cluster sites into **homogeneous groups**
Constraint: $\mu = \mu_k$ if site in unknown group k

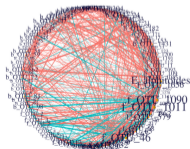
Constrain **species abundances** μ



Constrain **species dependencies** Σ

PLN-PCA: find structure
Goal: find **few structuring factors** governing species dependencies
Model: force Σ to have **low rank**

PLN-network: find interactions
Goal: find pairs of species in **direct interaction**
Model: force $\Omega^{-1} = \Sigma$ to be **sparse**



Work with C. Vacher



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.



Han Liu, Kathryn Roeder, and Larry Wasserman.

Stability approach to regularization selection (stars) for high dimensional graphical models.

In *Proceedings of the 23rd International Conference on Neural Information Processing Systems - Volume 2*, NIPS'10, pages 1432–1440, USA, 2010. Curran Associates Inc.



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.