

Inference de réseaux écologiques à partir de données de comptage

J. Chiquet², J. Lao¹, M. Mariadassou¹, S. Robin², S. Schbath¹

¹MAIAGE, INRA, Université Paris-Saclay, 78350 Jouy-en-Josas, France

²MIA 518, INRA, AgroParisTech, 16 Rue Claude Bernard, 75005 Paris, France

Journées NETBIO
10 Novembre 2017

- 1 Introduction
 - Métagénomique
 - Interactions écologiques
 - Réseaux de co-occurrences

- 2 Méthodes de Reconstruction

- 3 Modèles pour Données Compositionnelles

- 4 Modèles pour Données de Comptages

- 5 Conclusion et Perspectives

Principe

Caractérisation de la **globalité** d'un écosystème microbien,
sans isolation au préalable de microorganismes qui le compose

Principe

Caractérisation de la **globalité** d'un écosystème microbien, **sans isolation au préalable** de microorganismes qui le compose

Pourquoi ?

- < 1% des microorganismes de la biosphère sont cultivables *in vitro*

Principe

Caractérisation de la **globalité** d'un écosystème microbien, **sans isolation au préalable** de microorganismes qui le compose

Pourquoi ?

- < 1% des microorganismes de la biosphère sont cultivables *in vitro*
- Perte des interactions écologiques avec les analyses par isolation

Principe

Caractérisation de la **globalité** d'un écosystème microbien, **sans isolation au préalable** de microorganismes qui le compose

Pourquoi ?

- < 1% des microorganismes de la biosphère sont cultivables *in vitro*
- Perte des interactions écologiques avec les analyses par isolation

Fort intérêt

Microbiote intestinal humain (MetaHit, HMP), écosystèmes océaniques (Tara Expeditions), et de la terre (Earth Microbiome Project), Produits Fermentés (FoodMicrobiome, MétaPDOCheese)

- Identification d'un mélange de bactéries via leurs séquences génétiques;

- Identification d'un mélange de bactéries via leurs séquences génétiques;
- Utilisation d'un marqueur génétique **universel** : ARNr 16S pour les prokaryotes

- Identification d'un mélange de bactéries via leurs séquences génétiques;
- Utilisation d'un marqueur génétique **universel** : ARNr 16S pour les prokaryotes
- Caractéristiques importantes :
 - Succession de régions variables et conservées



CONSERVED REGIONS: unspecific applications

VARIABLE REGIONS: group or species-specific applications

- Identification d'un mélange de bactéries via leurs séquences génétiques;
- Utilisation d'un marqueur génétique **universel** : ARNr 16S pour les prokaryotes
- Caractéristiques importantes :
 - Succession de régions variables et conservées

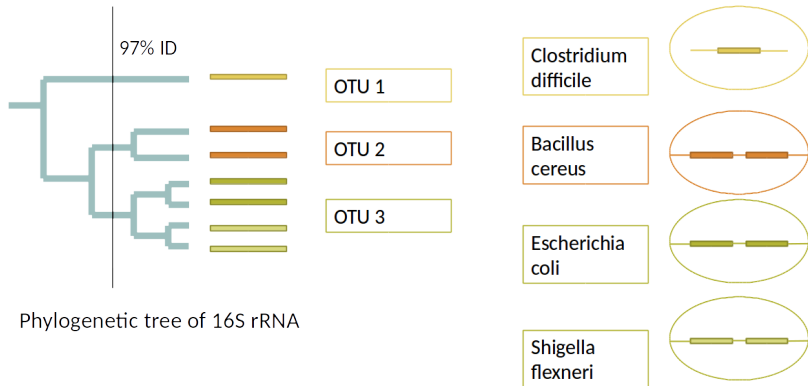


CONSERVED REGIONS: unspecific applications

VARIABLE REGIONS: group or species-specific applications

- Amorces PCR possibles sur les régions flanquantes car conservées
- Régions variables sont de longueur "séquençable"
- Bases de données exhaustives et (à peu près) fiables

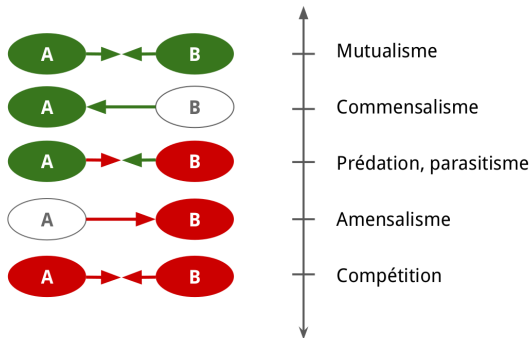
Principe des OTUs



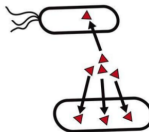
Données Types

##		Sample_a	Sample_b	Sample_c	Sample_d	Sample_e
##	OTU_1	3	3	3	3	6
##	OTU_2	5	2	3	3	1
##	OTU_3	2	0	2	4	3
##	OTU_4	4	3	0	4	3
##	OTU_5	0	2	2	1	4
##	OTU_6	0	2	2	1	3

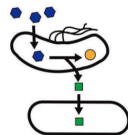
Interactions écologiques



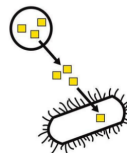
Compétition



Syntrophie

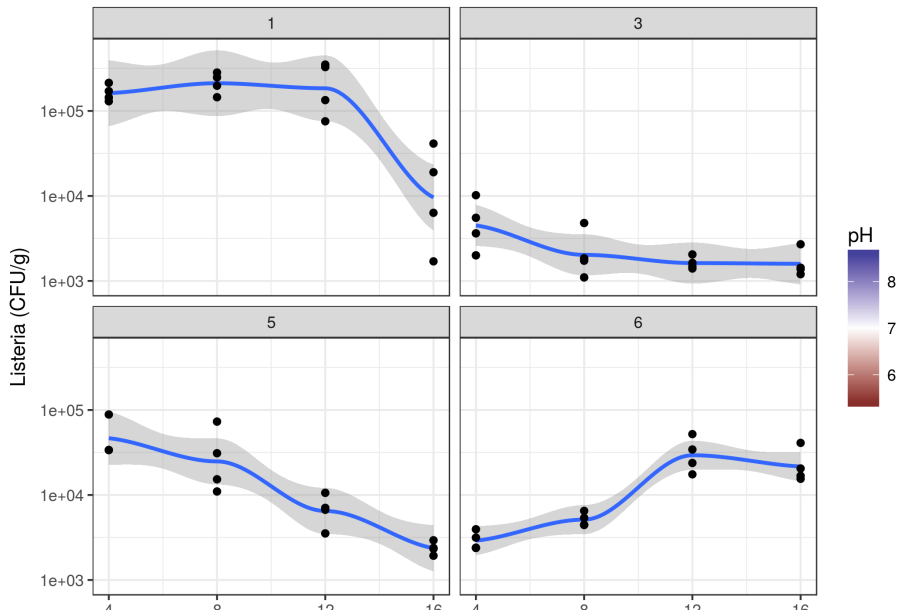


Alimentation croisée

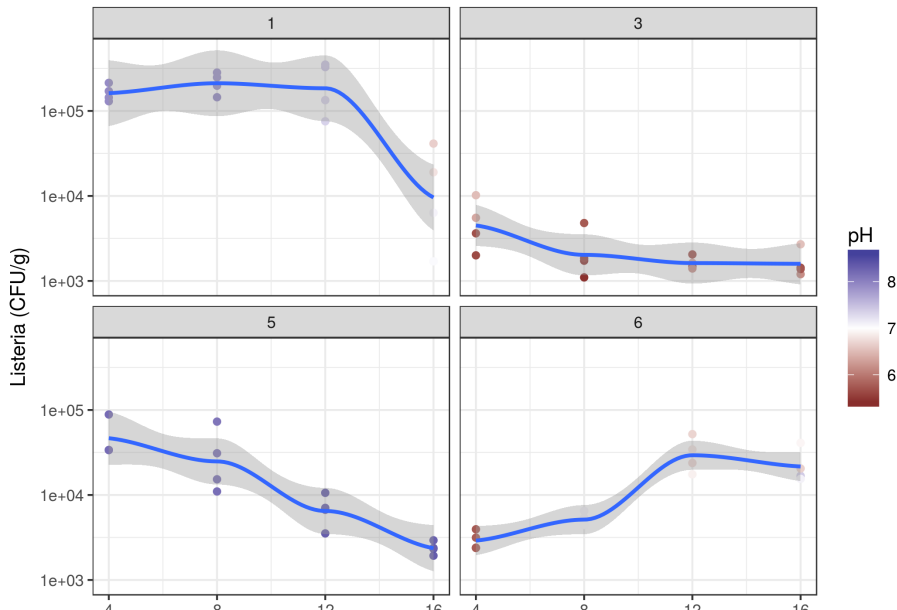


Seth EC, Taga ME (2014). Nutrient cross-feeding in the microbial world. *Front. Microbiol.* 5(350). doi:10.3389/fmicb.2014.00350

Exemple: Fortress (E. Dugat-Bony)



Exemple: Fortress (E. Dugat-Bony)



Exemple: Millup (V. Gagnaire/A. Thierry)

Maximal population reached (log cfu or nb copies/ml)



- On se limite à des interactions entre **deux** taxa;
- On s'intéresse surtout à des interactions de type
 - *commensales*,
 - *mutualistes*,
 - *compétitives*

Le reste étant trop difficile à détecter

- On se limite à des interactions entre **deux** taxa;
- On s'intéresse surtout à des interactions de type
 - *commensales*,
 - *mutualistes*,
 - *compétitives*

Le reste étant trop difficile à détecter

- Ces 3 types induisent par de la **co-présence** (commensalisme, mutualisme) ou de la **co-exclusion** (compétition)
- On passe par la co-occurrence pour identifier des interactions

Interaction \simeq co-occurrence

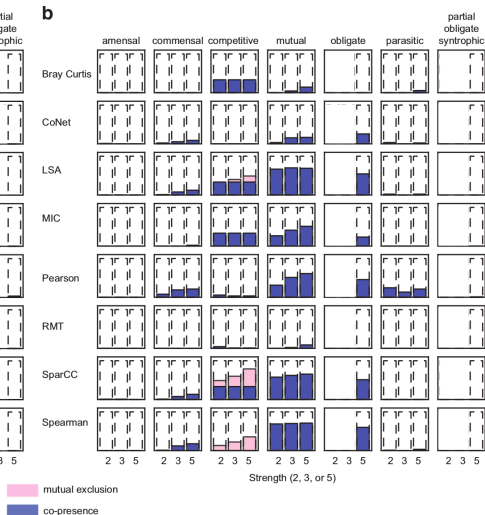
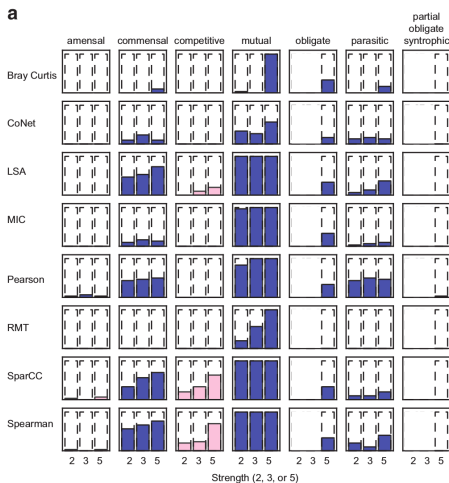
- On se limite à des interactions entre **deux** taxa;
- On s'intéresse surtout à des interactions de type
 - *commensales*,
 - *mutualistes*,
 - *compétitives*

Le reste étant trop difficile à détecter

- Ces 3 trois types induisent par de la **co-présence** (commensalisme, mutualisme) ou de la **co-exclusion** (compétition)
- On passe par la co-occurrence pour identifier des interactions

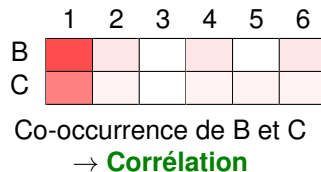
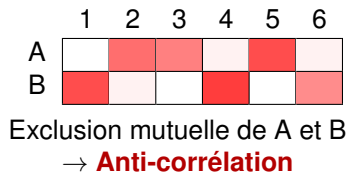
Interaction $\begin{matrix} \implies \\ \nleftarrow \end{matrix}$ Co-occurrence

Autres interactions (d'après Weiss et al., 2016)

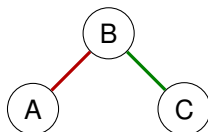


Réseaux de co-occurrences

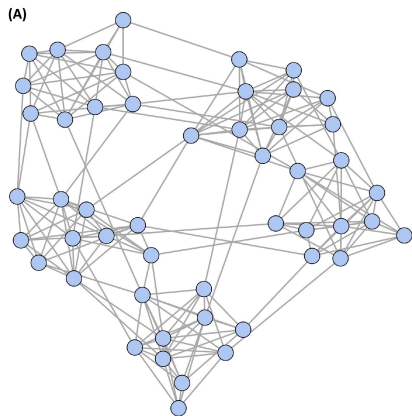
Observation des abondances



Inférence du réseau

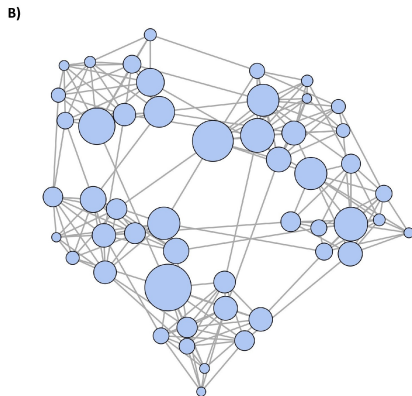
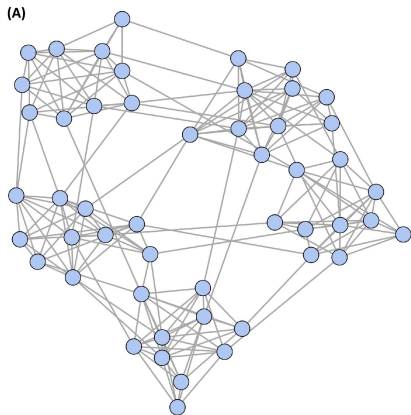


1. Avoir une jolie figure



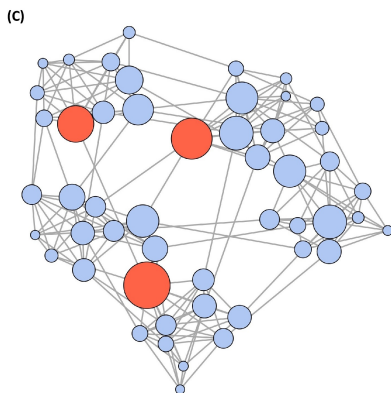
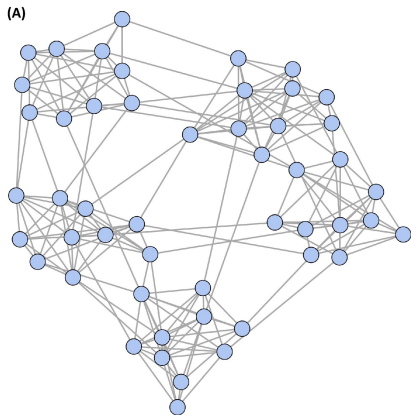
Des réseaux pour quoi faire? (Layeghifard et al. 2017)

1. Avoir une jolie figure
2. Trouver des taxa **systemiques**



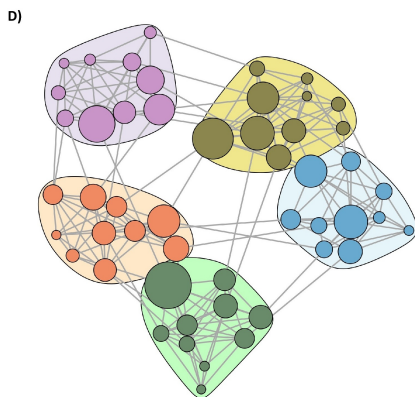
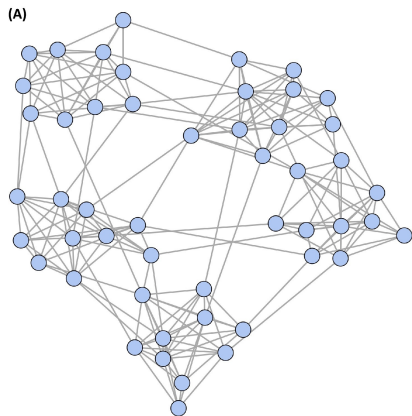
Des réseaux pour quoi faire? (Layeghifard et al. 2017)

1. Avoir une jolie figure
2. Trouver des taxa **systemiques**
3. Trouver des taxa **clé de route** (*keystone species*)



Des réseaux pour quoi faire? (Layeghifard et al. 2017)

1. Avoir une jolie figure
2. Trouver des taxa **systemiques**
3. Trouver des taxa **clé de route** (*keystone species*)
4. Trouver des **communautés** de cluster de taxa en interaction



Problématique

Évaluer la présence/absence d'une arête pour chaque **paire** de taxa possible

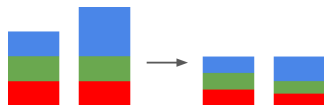
Problématique

Évaluer la présence/absence d'une arête pour chaque **paire** de taxa possible

Sparsité Beaucoup de 0 dans les tables d'abondances
→ Détection de *co-absences* au détriment de taxa rares

Dimensionnalité Nombre d'OTUs \gg Nombre d'échantillons
→ Manque de puissance

Compositionnalité Les abondances peuvent être biaisées : profondeurs de séquençage différentes, effet de la normalisation, ...



Lorsque les abondances de ces 2 communautés sont représenté de façon relative, les OTUs vert et rouge ont des abondances différentes

- 1 Introduction
- 2 **Méthodes de Reconstruction**
 - Panorama des Méthodes
 - Corrélation pour données de comptages
- 3 Modèles pour Données Compositionnelles
- 4 Modèles pour Données de Comptages
- 5 Conclusion et Perspectives

Par ordre (approximatif) de complexité

1. Corrélation (Pearson, Spearman)

Une grande variété de méthodes...

Par ordre (approximatif) de complexité

1. Corrélation (Pearson, Spearman)
2. (Dis)similarité des profils d'occurrence (Bray-Curtis, Kullback-Leibler, etc)

Par ordre (approximatif) de complexité

1. Corrélation (Pearson, Spearman)
2. (Dis)similarité des profils d'occurrence (Bray-Curtis, Kullback-Leibler, etc)
3. Modèles (Graphiques) Gaussiens
 - sur données *transformées* (log, arcsin, etc)
 - avec correction *compositionnelle* (sparCC, Rebacka, SPIEC-EASI)
 - pour données de *comptages* (Mint, PLN-network)

Par ordre (approximatif) de complexité

3. Modèles (Graphiques) Gaussiens

- sur données *transformées* (log, arcsin, etc)
- avec correction *compositionnelle* (sparCC, Rebacca, SPIEC-EASI)
- pour données de *comptages* (Mint, PLN-network)

Interaction \simeq Co-occurrence \simeq Corrélation (directe/**partielle**)

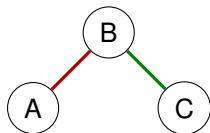
Interaction \simeq Co-occurrence \simeq Corrélation (directe/**partielle**)

Réseau d'interaction $G \simeq$ Matrice de covariance directe (Σ) ou partielle (Σ^{-1})

Interaction \simeq Co-occurrence \simeq Corrélation (directe/**partielle**)

Réseau d'interaction $G \simeq$ Matrice de covariance directe (Σ) ou partielle (Σ^{-1})

Réseau G



Matrice Σ ou Σ^{-1}

	A	B	C
A		red	
B	red		green
C		green	

Approche compositionnelle

- Modéliser les proportions (à la Aitchinson)
- Ajouter une couche d'échantillonnage (Multinomiale)

Approche directe

- Modéliser directement les comptages (avec un paramètre pour la profondeur)
- Modèle graphique Poisson (Inouye et al., 2017)
- Espace latent et émission poissonnienne

- 1 Introduction
- 2 Méthodes de Reconstruction
- 3 Modèles pour Données Compositionnelles**
 - Définition
 - Inférence
 - Métrique d'évaluation
 - Schéma de simulation
 - Résultats
- 4 Modèles pour Données de Comptages
- 5 Conclusion et Perspectives

Definition

1. Données positives (ou nulles) sommant à 1

$$\sum_{j=1}^p x_j = 1; \quad \mathbf{x} = (x_1, \dots, x_p) \in \mathbb{S}^{p-1}$$

2. Description des parties d'un tout

Definition

1. Données positives (ou nulles) sommant à 1

$$\sum_{j=1}^p x_j = 1; \quad \mathbf{x} = (x_1, \dots, x_p) \in \mathbb{S}^{p-1}$$

2. Description des parties d'un tout

- information **relative** entre les parties
- Si les x_j sont des **proportions** correspondants à des **abondances absolues** b_j , on a

$$\log(x_j/x_{j'}) = \log(b_j/b_{j'})$$

La transformation centered log-ratio (*clr*) permet de se ramener à $\mathbb{H}^{p-1} = \{\mathbf{x}^\top \in \mathbb{R}^p : \mathbf{x}\mathbf{1} = 0\}$

Definition

$$clr : \begin{cases} \mathbb{S}^{p-1} & \rightarrow \mathbb{H}^{p-1} \subset \mathbb{R}^p \\ \mathbf{x} & \mapsto \mathbf{y} = clr(\mathbf{x}) = (x_1/m_g(\mathbf{x}), \dots, x_p/m_g(\mathbf{x})) \end{cases}$$

où $m_g(\mathbf{x})$ est la moyenne géométrique de \mathbf{x} .

Si les x_j sont des **proportions** correspondants à des **abondances absolues** b_j , on a

$$clr(\mathbf{x}) = clr(\mathbf{b})$$

Modèle Multinomial lognormal

Abondances	\mathbb{R}^p	\mathbf{b}_i i.i.d.	$\log(\mathbf{b}_i) \sim \mathcal{N}_p(\boldsymbol{\mu}, \boldsymbol{\Sigma})$
Proportions	\mathbb{S}^{p-1}	$\mathbf{x}_i = \text{clr}^{-1}(\text{clr}(\mathbf{b}_i))$	$\text{clr}(\mathbf{p}_i) \sim \mathcal{N}_p(\text{clr}(\boldsymbol{\mu}), \boldsymbol{\Gamma} = \mathbf{G}\boldsymbol{\Sigma}\mathbf{G})$
Comptages	\mathbb{N}^p	\mathbf{n}_i	$\mathbf{n}_i \sim \mathcal{M}(N_i, \mathbf{x}_i)$

Où $\mathbf{G} = \mathbf{I}_p - \frac{1}{p}\mathbf{1}\mathbf{1}^\top$ est la matrice de centrage.

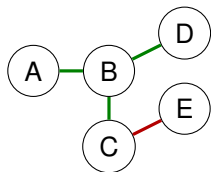
Modèle Multinomial lognormal

Abondances	\mathbb{R}^p	\mathbf{b}_i i.i.d.	$\log(\mathbf{b}_i) \sim \mathcal{N}_p(\boldsymbol{\mu}, \boldsymbol{\Sigma})$
Proportions	\mathbb{S}^{p-1}	$\mathbf{x}_i = \text{clr}^{-1}(\text{clr}(\mathbf{b}_i))$	$\text{clr}(\mathbf{p}_i) \sim \mathcal{N}_p(\text{clr}(\boldsymbol{\mu}), \boldsymbol{\Gamma} = \mathbf{G}\boldsymbol{\Sigma}\mathbf{G})$
Comptages	\mathbb{N}^p	\mathbf{n}_i	$\mathbf{n}_i \sim \mathcal{M}(N_i, \mathbf{x}_i)$

Où $\mathbf{G} = \mathbf{I}_p - \frac{1}{p}\mathbf{1}\mathbf{1}^\top$ est la matrice de centrage.

Le réseau est donnée par $\boldsymbol{\Sigma} / \boldsymbol{\Sigma}^{-1}$

Réseau réel



Abondances
de bases b

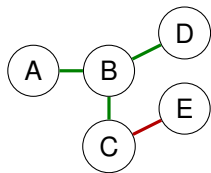
$$\begin{pmatrix} 8 \\ 8 \\ 8 \\ 9 \\ 2 \end{pmatrix}$$

$$\begin{pmatrix} 6 \\ 6 \\ 6 \\ 8 \\ 1 \end{pmatrix}$$

- Les abondances dans la communauté i sont gouvernées par des abondances de bases b_i log-normales $\log(\mathbf{b}_i) \sim \mathcal{N}(\mu, \Sigma)$

Modèle Multinomial lognormal (II)

Réseau réel



Abondances
de bases b

$$\begin{pmatrix} 8 \\ 8 \\ 8 \\ 9 \\ 2 \end{pmatrix}$$

→

Abondances
relatives x

$$\begin{pmatrix} 0.17 \\ 0.17 \\ 0.17 \\ 0.47 \\ 0.02 \end{pmatrix}$$

$$\begin{pmatrix} 6 \\ 6 \\ 6 \\ 8 \\ 1 \end{pmatrix}$$

→

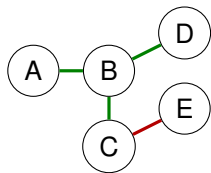
$$\begin{pmatrix} 0.09 \\ 0.09 \\ 0.09 \\ 0.71 \\ 0.02 \end{pmatrix}$$

- L'abondance relative (ou proportion) x_{ij} du taxon j dans l'échantillon i se déduit de b_i par une transformation logistique:

$$x_{ij} = \frac{e^{b_{ij}}}{\sum_j e^{b_{ij}}}$$

Modèle Multinomial lognormal (II)

Réseau réel



Abondances
de bases b

$$\begin{pmatrix} 8 \\ 8 \\ 8 \\ 9 \\ 2 \end{pmatrix}$$

→

Abondances
relatives x

$$\begin{pmatrix} 0.17 \\ 0.17 \\ 0.17 \\ 0.47 \\ 0.02 \end{pmatrix}$$

→

Tableau
comptage

	1	2
A	3	8
B	3	6
C	5	2
D	8	30
E	1	4
N	20	50

de

- Les comptages $\mathbf{n}_i = (n_{i1}, \dots, n_{ip})$ des différents taxa dans l'échantillon i suivent alors une loi multinomiale de paramètres N_i (profondeur de séquençage) et $\mathbf{x}_i = (x_{i1}, \dots, x_{ip})$.

$$\mathbf{n}_i \sim \mathcal{M}(N_i, \mathbf{x}_i)$$

Definition

Matrice de variation On considère la *matrice de variation* T définie par

$$T_{jj'} = \frac{1}{2} \text{Var}(\log(x_j) - \log(x_{j'})) = \frac{1}{2} \text{Var}(\log(b_j) - \log(b_{j'}))$$

Remarques

- T peut-être estimée par sa version empirique \hat{T} .
- $T_{ij} = \frac{1}{2}(\sigma_{ii} + \sigma_{jj}) - \sigma_{ij}$
- Si on connaissait les σ_{ii} , on aurait un estimateur de σ_{ij} .

Description

- **Objectif:** Reconstruire Σ
- **Méthode:** Méthode des moments glorifiée

Description

- **Objectif:** Reconstruire Σ
 - **Méthode:** Méthode des moments glorifiée
1. Transformation des n_i en x_i via clr (avec pseudocomptage éventuel)

Description

- **Objectif:** Reconstruire Σ
 - **Méthode:** Méthode des moments glorifiée
1. Transformation des n_i en x_i via clr (avec pseudocomptage éventuel)
 2. Calcul de \hat{T} à partir des x_i

Description

- **Objectif:** Reconstruire Σ
 - **Méthode:** Méthode des moments glorifiée
1. Transformation des n_i en x_i via clr (avec pseudocomptage éventuel)
 2. Calcul de \hat{T} à partir des x_i
 3. Écriture d'un système de $p(p-1)/2$ équations linéaires à $p(p-1)/2$ inconnues mais de rang $p(p-1)/2 - p$ reliant les σ_{ij} et les \hat{T}_{ij} :

$$\text{Avech}(\Sigma) = \text{vech}(T)$$

Description

- **Objectif:** Reconstruire Σ
 - **Méthode:** Méthode des moments glorifiée
1. Transformation des n_i en x_i via clr (avec pseudocomptage éventuel)
 2. Calcul de \hat{T} à partir des x_i
 3. Écriture d'un système de $p(p-1)/2$ équations linéaires à $p(p-1)/2$ inconnues mais de rang $p(p-1)/2 - p$ reliant les σ_{ij} et les \hat{T}_{ij} :
$$\text{Avech}(\Sigma) = \text{vech}(T)$$
 4. Résolution du système sous l'hypothèse que

Description

- **Objectif:** Reconstruire Σ
- **Méthode:** Méthode des moments glorifiée

1. Transformation des n_i en x_i via clr (avec pseudocomptage éventuel)
2. Calcul de \hat{T} à partir des x_i
3. Écriture d'un système de $p(p-1)/2$ équations linéaires à $p(p-1)/2$ inconnues mais de rang $p(p-1)/2 - p$ reliant les σ_{ij} et les \hat{T}_{ij} :

$$\text{Avech}(\Sigma) = \text{vech}(T)$$

4. Résolution du système sous l'hypothèse que
SparCC $\sum_{j \neq i} \sigma_{ij} \ll (p-1)\sigma_{ii} + \sum_{j \neq i} \sigma_{jj}$ (faibles corrélations)
 \implies estimateur facile des σ_{ii} et plug-in pour les σ_{ij}

Description

- **Objectif:** Reconstruire Σ
- **Méthode:** Méthode des moments glorifiée

1. Transformation des n_i en x_i via clr (avec pseudocomptage éventuel)
2. Calcul de \hat{T} à partir des x_i
3. Écriture d'un système de $p(p-1)/2$ équations linéaires à $p(p-1)/2$ inconnues mais de rang $p(p-1)/2 - p$ reliant les σ_{ij} et les \hat{T}_{ij} :

$$\text{Avech}(\Sigma) = \text{vech}(T)$$

4. Résolution du système sous l'hypothèse que

SparCC $\sum_{j \neq i} \sigma_{ij} \ll (p-1)\sigma_{ii} + \sum_{j \neq i} \sigma_{jj}$ (faibles corrélations)
 \implies estimateur facile des σ_{ii} et plug-in pour les σ_{ij}

Rebacca Σ est sparse (moins de $p/4$ coefficients non-nuls par ligne), via une régularisation ℓ_1 du système:

$$\| \arg \min_{\Sigma} \text{Avech}(\Sigma) - \text{vech}(T) \| + \lambda \| \text{vech}(\Sigma) \|_1$$

Avantages

- + Méthode très simple et très rapide
- + Simple résolution d'un système linéaire

Inconvénients

- Focus sur Σ et non Σ^{-1}
- Oubli de la couche multinomiale
- Complexité en espace $O(p^4)$ (pour REBACCA)
- Choix de λ (REBACCA) ou du nombre d'itérations (SparCC)

Description

- **Objectif:** Reconstruire Σ^{-1}
- **Méthode:** Se ramener à un lasso graphique

Description

- **Objectif:** Reconstruire Σ^{-1}
- **Méthode:** Se ramener à un lasso graphique

1. Transformation des n_i en x_i via clr (avec pseudocomptage éventuel)

Description

- **Objectif:** Reconstruire Σ^{-1}
- **Méthode:** Se ramener à un lasso graphique

1. Transformation des n_i en x_i via clr (avec pseudocomptage éventuel)
2. Supposer que p est grand donc que $G \simeq I_p$ et $\Gamma = G\Sigma G \simeq \Sigma$

Description

- **Objectif:** Reconstruire Σ^{-1}
 - **Méthode:** Se ramener à un lasso graphique
1. Transformation des n_i en x_i via clr (avec pseudocomptage éventuel)
 2. Supposer que p est grand donc que $G \simeq I_p$ et $\Gamma = G\Sigma G \simeq \Sigma$
 3. Utiliser $\hat{\Gamma}$ en lieu et place de $\hat{\Sigma}$

Description

- **Objectif:** Reconstruire Σ^{-1}
- **Méthode:** Se ramener à un lasso graphique

1. Transformation des n_i en x_i via clr (avec pseudocomptage éventuel)
2. Supposer que p est grand donc que $G \simeq I_p$ et $\Gamma = G\Sigma G \simeq \Sigma$
3. Utiliser $\hat{\Gamma}$ en lieu et place de $\hat{\Sigma}$
4. Insérer dans un lasso graphique¹ pour estimer $\Omega = \Sigma^{-1}$

$$\hat{\Omega} = \arg \min_{\Omega \in \mathbb{S}_{++}} -\log \det \Omega + \text{Tr}(\Omega \hat{\Gamma}) + \lambda \|\Omega\|_1$$

¹Variante possible avec l'estimation de voisinage à la Meinshausen et Bühlmann

Avantages

- + Focus sur $\Omega = \Sigma^{-1}$
- + Basée sur une méthode de vraisemblance pénalisée

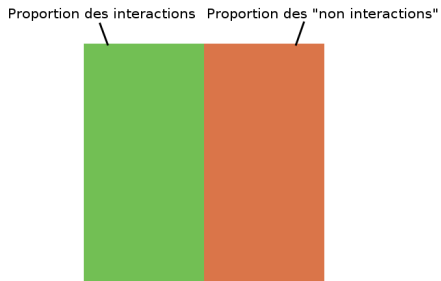
Inconvénients

- Oubli de la couche multinomiale
- Choix de λ (REBACCA)
- Relativement long

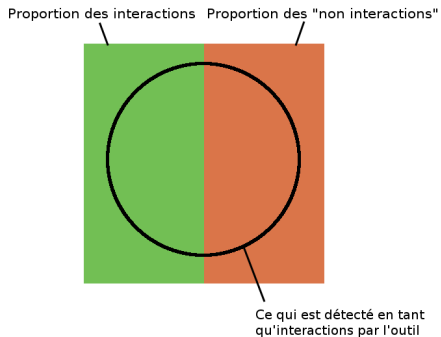
Proportion des interactions



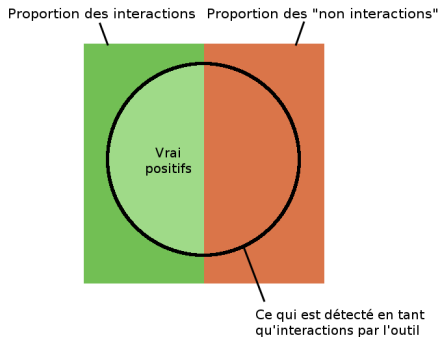
Mesures de la performance



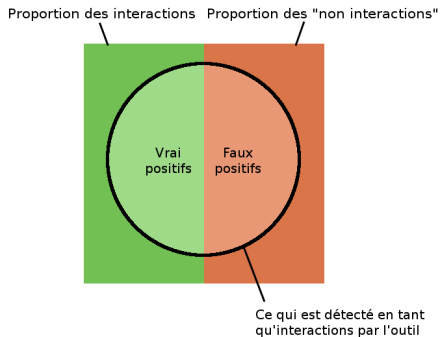
Mesures de la performance



Mesures de la performance



Mesures de la performance

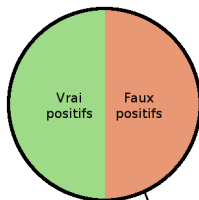


Proportion des interactions



$$\frac{\textit{Vrai positifs}}{\textit{Positifs}}$$

Sensibilité/Rappel Proportion des interactions du réseau ayant été détectées par l'outil

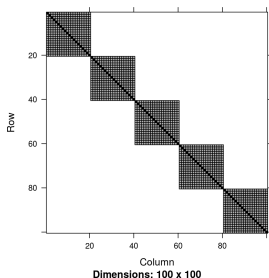


Ce qui est détecté en tant qu'interactions par l'outil

$$\frac{\text{Vrai positifs}}{\text{Vrai positifs} + \text{Faux positifs}}$$

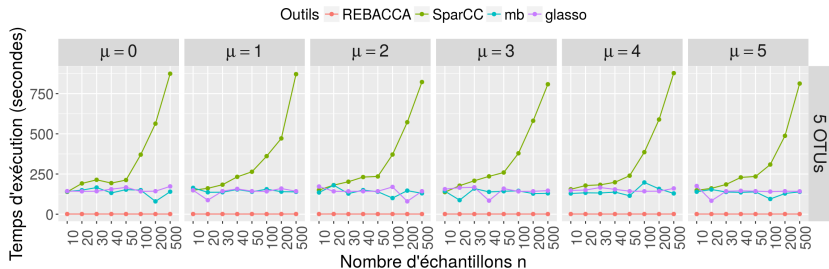
Précision Proportion des interactions ayant été correctement détectées parmi ce que l'outil a détecté

Schéma de simulation

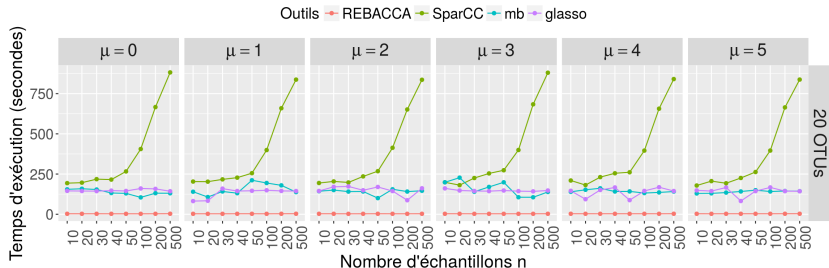


- p Nombre d'OTUs
→ Réseaux de taille variable, passage à l'échelle
- μ Gamme des abondances des OTUs
→ Détection des taxa rares
- ρ Force moyenne des corrélations
→ Détection d'interactions faibles (puissance de l'outil)
- k Matrice diagonale par blocs de k blocs : Sparsité du réseau
→ Robustesse aux données non sparses

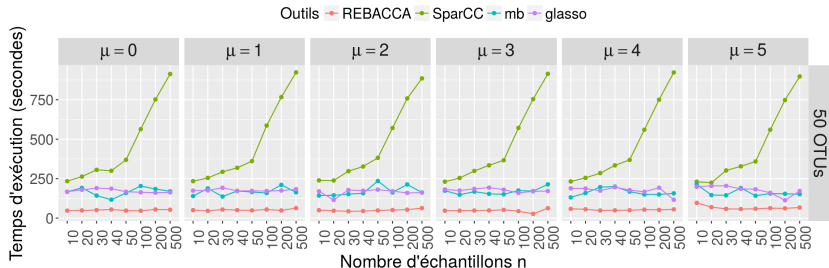
SparCC est le plus lent pour des petits réseaux (12 minutes contre 4 minutes)



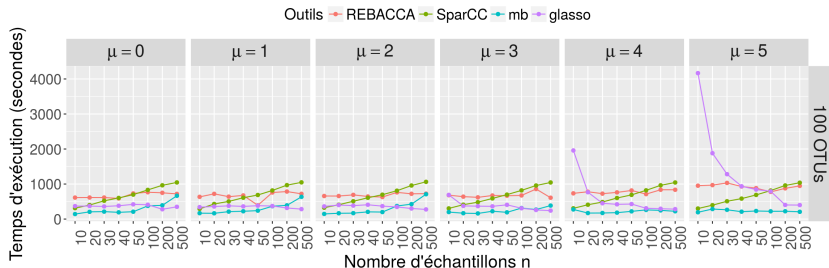
SparCC est le plus lent pour des petits réseaux (12 minutes contre 4 minutes)



SparCC est le plus lent pour des petits réseaux (12 minutes contre 4 minutes)

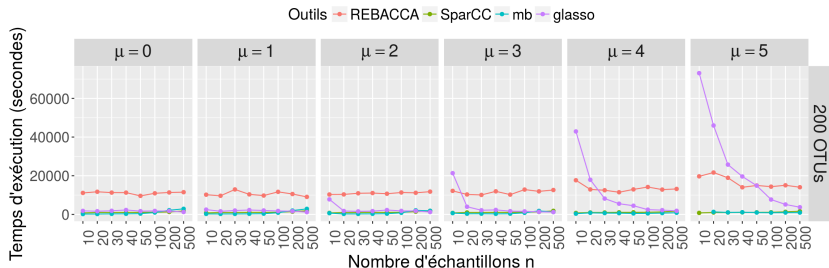


SparCC possède les meilleurs temps sur les larges réseaux
glasso manque de puissance lorsque n est petit (10...50) :
2,5x plus lent que les autres méthodes



Temps d'exécution

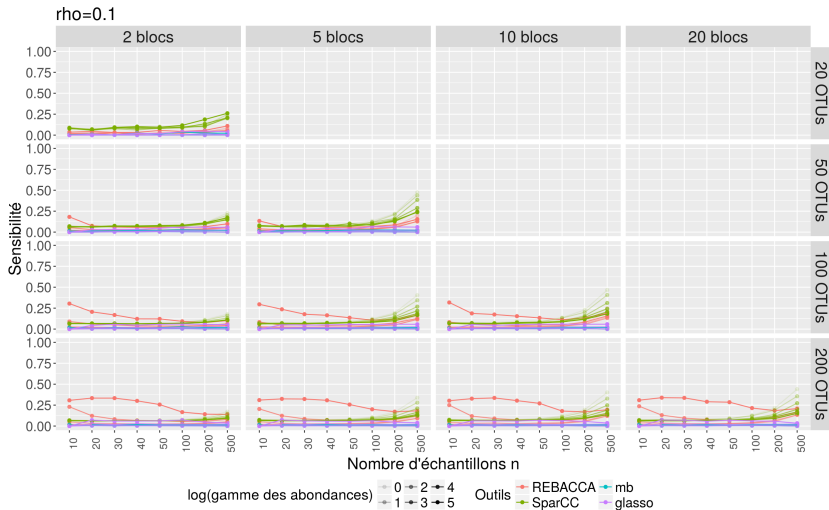
SparCC possède les meilleurs temps sur les larges réseaux
glasso manque de puissance lorsque n est petit (10...50) :
Presque **5x plus lent** que les autres méthodes (19 heures contre 4 heures)



- REBACCA est limité à 250 OTUs
→ Complexité en mémoire en $O(D^4)$
- Les simulations très difficiles n'aboutissent pas avec la méthode SPIEC-EASI (MB)
→ Compensation du manque de puissance pour détecter les taxa rares

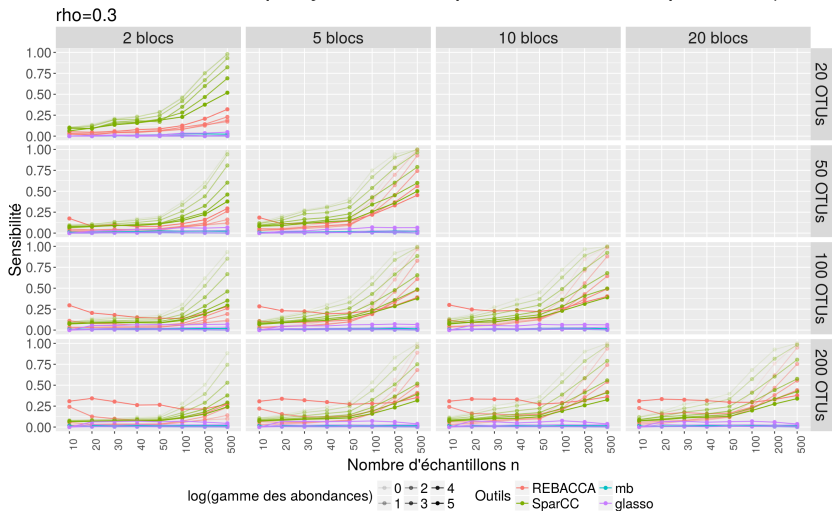
Réseaux de densité variable : Sensibilité

Mauvaise sensibilité pour de **faibles** interactions ($\rho = 0.1$)



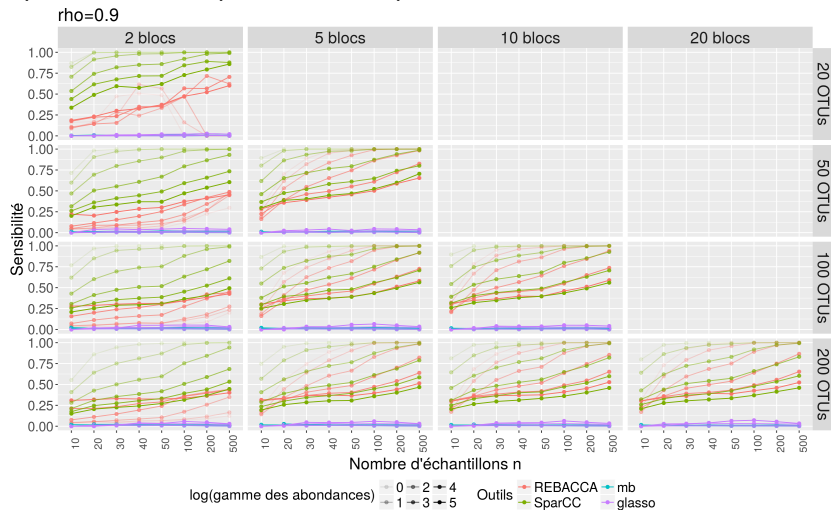
Réseaux de densité variable : Sensibilité

Bonne sensibilité lorsqu'il y a beaucoup d'échantillons à partir de $\rho = 0.3$



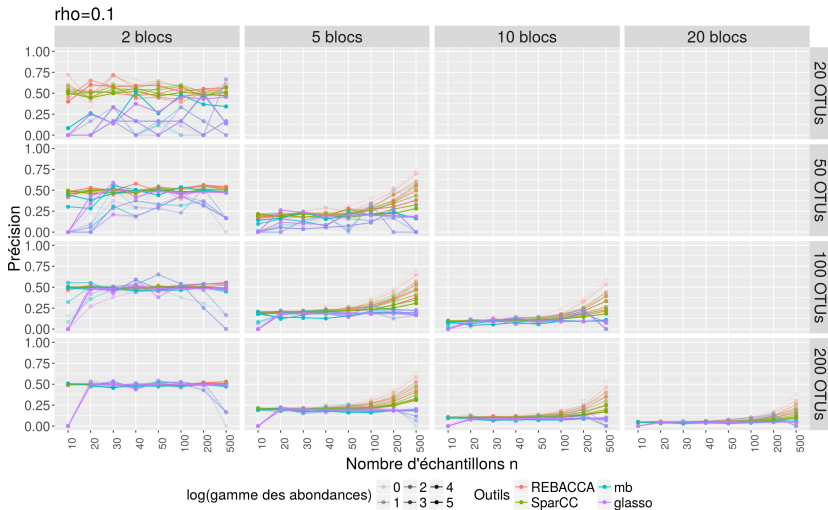
Réseaux de densité variable : Sensibilité

SparCC meilleur que REBACCA pour des réseaux denses



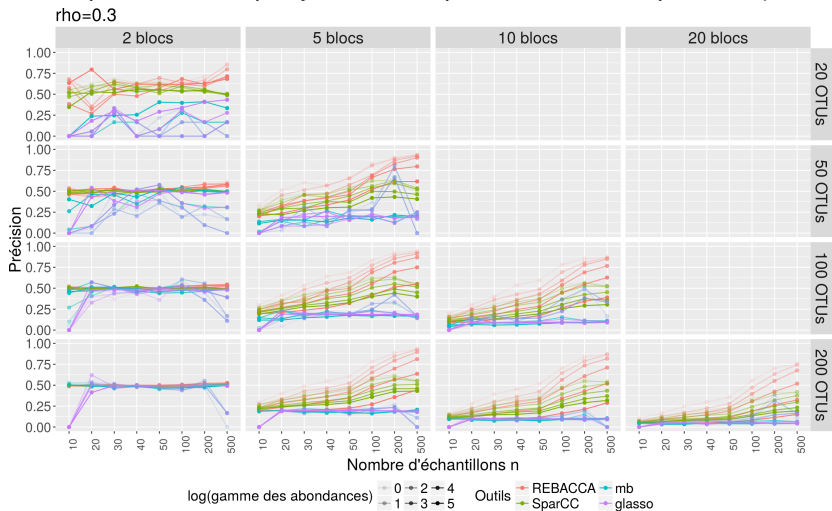
Réseaux de densité variable : Précision

Mauvaise précision pour de **faibles** interactions ($\rho = 0.1$)



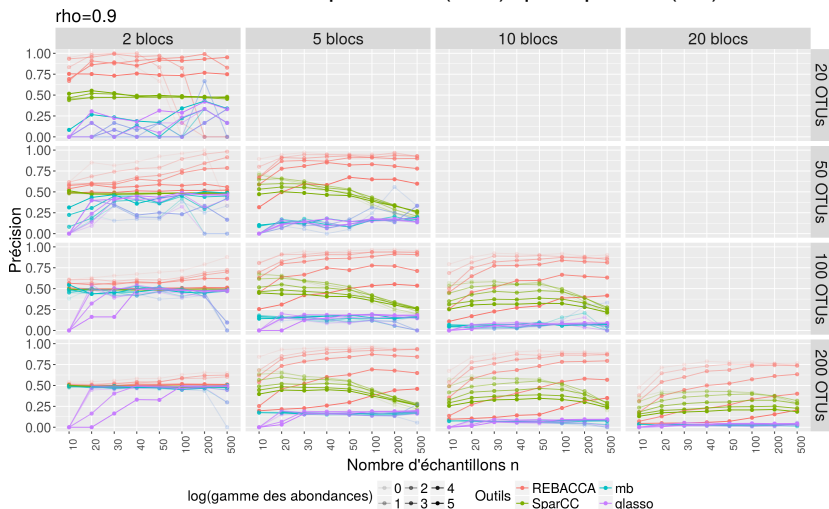
Réseaux de densité variable : Précision

Bonne précision lorsqu'il y a beaucoup d'échantillons à partir de $\rho = 0.3$



Réseaux de densité variable : Précision

REBACCA a une meilleure précision (0.75) que SparCC (0.5)



- REBACCA: (Assez) **peu** d'interactions mais plutôt **fiables**
sensibilité moyenne, bonne précision
 - Bonnes performances, mais limité à 250 OTUs.
 - Solutions possibles : Agglomération d'OTUs en méta-OTUs selon leur assignation taxonomique.

- REBACCA: (Assez) **peu** d'interactions mais plutôt **fiables**
sensibilité moyenne, bonne précision
 - Bonnes performances, mais limité à 250 OTUs.
 - Solutions possibles : Agglomération d'OTUs en méta-OTUs selon leur assignation taxonomique.
- SparCC: **Beaucoup** d'interactions mais **moins** fiables
bonne sensibilité, précision moyenne

- REBACCA: (Assez) **peu** d'interactions mais plutôt **fiables**
sensibilité moyenne, bonne précision
 - Bonnes performances, mais limité à 250 OTUs.
 - Solutions possibles : Agglomération d'OTUs en méta-OTUs selon leur assignation taxonomique.
- SparCC: **Beaucoup** d'interactions mais **moins** fiables
bonne sensibilité, précision moyenne
- SPIEC-EASI: Mauvais résultats
mauvaise sensibilité, mauvaise précision
 - Résultats très inférieurs à ceux de la publi.
 - Modèle de simulation (dans la publi) non compositionnel.

Conclusions Partielles

- REBACCA: (Assez) **peu** d'interactions mais plutôt **fiables**
sensibilité moyenne, bonne précision
 - Bonnes performances, mais limité à 250 OTUs.
 - Solutions possibles : Agglomération d'OTUs en méta-OTUs selon leur assignation taxonomique.
- SparCC: **Beaucoup** d'interactions mais **moins** fiables
bonne sensibilité, précision moyenne
- SPIEC-EASI: Mauvais résultats
mauvaise sensibilité, mauvaise précision
 - Résultats très inférieurs à ceux de la publi.
 - Modèle de simulation (dans la publi) non compositionnel.
- Plus **difficile** pour des gammes d'abondances **variables**
- Aucun des 3 outils ne gère les co-variables (confusion entre préférence d'habitat et interaction).

- 1 Introduction
- 2 Méthodes de Reconstruction
- 3 Modèles pour Données Compositionnelles
- 4 **Modèles pour Données de Comptages**
 - Modèles Graphiques Poissonniens
 - Modèles Graphiques Poisson Log-Normal
 - Simulation
- 5 Conclusion et Perspectives

- Peut-on modéliser directement la loi jointe de n ?

- Peut-on modéliser directement la loi jointe de \mathbf{n} ?
- Par analogie au cas Gaussien, on aimerait avoir une loi de la forme

$$P(\mathbf{n}) \propto \exp \left(\mathbf{n}^\top \boldsymbol{\Omega} \mathbf{n} + \boldsymbol{\eta}^\top \mathbf{n} - \sum_j \log(n_j!) - B(\boldsymbol{\eta}, \boldsymbol{\Omega}) \right)$$

- Peut-on modéliser directement la loi jointe de \mathbf{n} ?
- Par analogie au cas Gaussien, on aimerait avoir une loi de la forme

$$P(\mathbf{n}) \propto \exp \left(\mathbf{n}^\top \boldsymbol{\Omega} \mathbf{n} + \boldsymbol{\eta}^\top \mathbf{n} - \sum_j \log(n_j!) - B(\boldsymbol{\eta}, \boldsymbol{\Omega}) \right)$$

qui garantit que

- les lois **conditionnelles** de chaque comptage sont Poisson
- $\Omega_{ij} = 0$ sont conditionnellement indépendantes

- Peut-on modéliser directement la loi jointe de \mathbf{n} ?
- Par analogie au cas Gaussien, on aimerait avoir une loi de la forme

$$P(\mathbf{n}) \propto \exp \left(\mathbf{n}^\top \boldsymbol{\Omega} \mathbf{n} + \boldsymbol{\eta}^\top \mathbf{n} - \sum_j \log(n_j!) - B(\boldsymbol{\eta}, \boldsymbol{\Omega}) \right)$$

qui garantit que

- les lois **conditionnelles** de chaque comptage sont Poisson
- $\Omega_{ij} = 0$ sont conditionnellement indépendantes
- Possible uniquement si
 - $\Omega_{jj} = 0$ pour tout j
 - $\Omega_{ij} \leq 0$ pour tout $i \neq j$

- Peut-on modéliser directement la loi jointe de \mathbf{n} ?
- Par analogie au cas Gaussien, on aimerait avoir une loi de la forme

$$P(\mathbf{n}) \propto \exp \left(\mathbf{n}^\top \boldsymbol{\Omega} \mathbf{n} + \boldsymbol{\eta}^\top \mathbf{n} - \sum_j \log(n_j!) - B(\boldsymbol{\eta}, \boldsymbol{\Omega}) \right)$$

qui garantit que

- les lois **conditionnelles** de chaque comptage sont Poisson
- $\Omega_{ij} = 0$ sont conditionnellement indépendantes
- Possible uniquement si
 - $\Omega_{jj} = 0$ pour tout j
 - $\Omega_{ij} \leq 0$ pour tout $i \neq j$
- \implies Uniquement des interactions **négatives** ☹

Modèle Poisson Log-Normal (PLN)

Modèle

(log)-abondances	\mathbb{R}^p	\mathbf{b}_i	i.i.d.	$\mathbf{b}_i \sim \mathcal{N}_p(\boldsymbol{\mu}, \boldsymbol{\Sigma})$
Comptages	\mathbb{N}^p	$n_{ij} b_{ij}$	indep.	$n_{ij} b_{ij} \sim \mathcal{P}(e^{b_{ij}})$

Modèle Poisson Log-Normal (PLN)

Modèle

(log)-abondances	\mathbb{R}^p	\mathbf{b}_i	i.i.d.	$\mathbf{b}_i \sim \mathcal{N}_p(\boldsymbol{\mu}, \boldsymbol{\Sigma})$
Comptages	\mathbb{N}^p	$n_{ij} b_{ij}$	indep.	$n_{ij} b_{ij} \sim \mathcal{P}(e^{b_{ij}})$

Properties:

$$\mathbb{E}(n_{ij}) = e^{\mu_j + \Sigma_{jj}/2} =: \lambda_j$$

$$\mathbb{V}(n_{ij}) = \lambda_j + \lambda_j^2 (e^{\Sigma_{jj}} - 1) \quad (\text{surdispersion})$$

$$\text{Cov}(n_{ij}, n_{ik}) = \lambda_j \lambda_k (e^{\Sigma_{jk}} - 1) \quad (\text{signe arbitraire})$$

Modèle

(log)-abondances	\mathbb{R}^p	\mathbf{b}_i	i.i.d.	$\mathbf{b}_i \sim \mathcal{N}_p(\boldsymbol{\mu}, \boldsymbol{\Sigma})$
Comptages	\mathbb{N}^p	$n_{ij} b_{ij}$	indep.	$n_{ij} b_{ij} \sim \mathcal{P}(e^{b_{ij}})$

Modèle

(log)-abondances \mathbb{R}^p \mathbf{b}_i i.i.d. $\mathbf{b}_i \sim \mathcal{N}_p(\boldsymbol{\mu}, \boldsymbol{\Sigma})$

Comptages \mathbb{N}^p $n_{ij}|b_{ij}$ indep. $n_{ij}|b_{ij} \sim \mathcal{P}(e^{b_{ij}})$

Propriétés:

$$\mathbb{E}(n_{ij}) = e^{\mu_j + \Sigma_{jj}/2} =: \lambda_j$$

$$\mathbb{V}(n_{ij}) = \lambda_j + \lambda_j^2 (e^{\Sigma_{jj}} - 1) \quad (\text{surdispersion})$$

$$\text{Cov}(n_{ij}, n_{ik}) = \lambda_j \lambda_k (e^{\Sigma_{jk}} - 1) \quad (\text{signe arbitraire})$$

PLN Covariables et Offset

Modèle

(log)-abondances	\mathbb{R}^p	\mathbf{b}_i	i.i.d.	$\mathbf{b}_i \sim \mathcal{N}_p(\boldsymbol{\mu}, \boldsymbol{\Sigma})$
Comptages	\mathbb{N}^p	$n_{ij} b_{ij}$	indep.	$n_{ij} b_{ij} \sim \mathcal{P}(e^{b_{ij}})$

Propriétés:

$$\mathbb{E}(n_{ij}) = e^{\mu_j + \Sigma_{jj}/2} =: \lambda_j$$

$$\mathbb{V}(n_{ij}) = \lambda_j + \lambda_j^2 (e^{\Sigma_{jj}} - 1) \quad (\text{surdispersion})$$

$$\text{Cov}(n_{ij}, n_{ik}) = \lambda_j \lambda_k (e^{\Sigma_{jk}} - 1) \quad (\text{signe arbitraire})$$

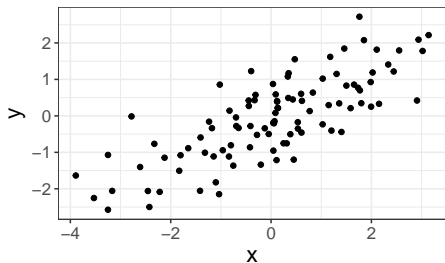
Covariables et Offset

En présence de covariables \mathbf{X}_i et d'offset \mathbf{O}_i , on peut modéliser $\boldsymbol{\mu}_i$ par

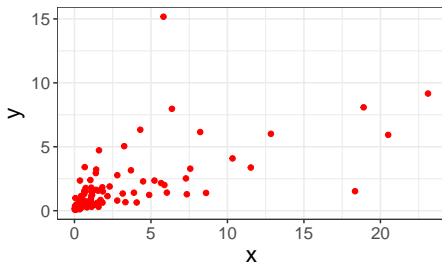
$$\boldsymbol{\mu}_i = \mathbf{O}_i + \mathbf{X}_i \boldsymbol{\Theta}$$

Vue géométrique

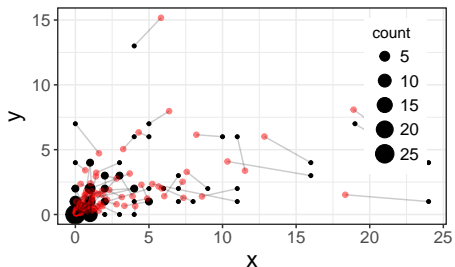
Latent Space (B)



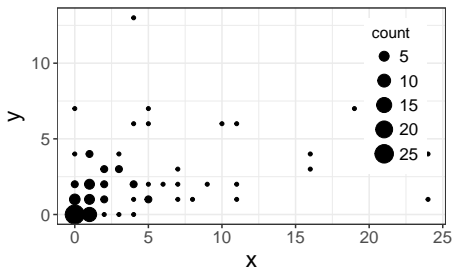
Observation Space ($\exp(B)$)



Observation Space ($N = P(\exp(B))$)



Observation Space (N)



$$\begin{aligned}\log p(\mathbf{N}, \mathbf{B}; \Theta, \Sigma) &= \sum_{i=1}^n \log p(\mathbf{n}_i | \mathbf{b}_i; \Theta, \Sigma) + \log p(\mathbf{b}_i) \\ &= \mathbf{1}_n^\top [\mathbf{N} \odot (\mathbf{X}\Theta + \mathbf{B}) - \exp(\mathbf{X}\Theta + \mathbf{B})] \mathbf{1}_p \\ &\quad - \frac{\|\Sigma^{-1/2} \mathbf{B}\|_F^2}{2} - \frac{nq}{2} \log(2\pi) - K(\mathbf{N})\end{aligned}$$

où \exp est appliquée par composante et \odot le produit de Hadamard.

$$\begin{aligned}\log p(\mathbf{N}, \mathbf{B}; \Theta, \Sigma) &= \sum_{i=1}^n \log p(\mathbf{n}_i | \mathbf{b}_i; \Theta, \Sigma) + \log p(\mathbf{b}_i) \\ &= \mathbf{1}_n^\top [\mathbf{N} \odot (\mathbf{X}\Theta + \mathbf{B}) - \exp(\mathbf{X}\Theta + \mathbf{B})] \mathbf{1}_p \\ &\quad - \frac{\|\Sigma^{-1/2} \mathbf{B}\|_F^2}{2} - \frac{nq}{2} \log(2\pi) - K(\mathbf{N})\end{aligned}$$

où \exp est appliquée par composante et \odot le produit de Hadamard.

- Pas de forme simple pour $\log p(\mathbf{N}; \Theta, \Sigma) \Rightarrow$ Pas de solution exacte pour Θ and Σ
- Pas de forme simple pour $p(\mathbf{n}_i | \mathbf{b}_i; \Theta, \Sigma) \Rightarrow$ Pas d'algorithme EM standard

On veut calculer la vraisemblance **intractable** $\log p(\mathbf{N}; \Theta, \Sigma)$.

On veut calculer la vraisemblance **intractable** $\log p(\mathbf{N}; \Theta, \Sigma)$.

On veut calculer la vraisemblance **intractable** $\log p(\mathbf{N}; \Theta, \Sigma)$.

On peut le développer en

$$\log p(\mathbf{N}; \Theta, \Sigma) = \mathbb{E}_{\mathbf{B}|\mathbf{N}}[\log p(\mathbf{N}, \mathbf{N}; \Theta, \Sigma) - \log p(\mathbf{B}|\mathbf{N})]$$

On veut calculer la vraisemblance **intractable** $\log p(\mathbf{N}; \Theta, \Sigma)$.

On peut le développer en

$$\log p(\mathbf{N}; \Theta, \Sigma) = \mathbb{E}_{\mathbf{B}|\mathbf{N}}[\log p(\mathbf{N}, \mathbf{B}; \Theta, \Sigma) - \log p(\mathbf{B}|\mathbf{N})]$$

- $\log p(\mathbf{N}, \mathbf{B}; \Theta)$ est simple ☺;
- $\log p(\mathbf{B}|\mathbf{N})$ est complexe ☹.

On veut calculer la vraisemblance **intractable** $\log p(\mathbf{N}; \Theta, \Sigma)$.

On peut le développer en

$$\log p(\mathbf{N}; \Theta, \Sigma) = \mathbb{E}_{\mathbf{B}|\mathbf{N}}[\log p(\mathbf{N}, \mathbf{B}; \Theta, \Sigma) - \log p(\mathbf{B}|\mathbf{N})]$$

- $\log p(\mathbf{N}, \mathbf{B}; \Theta)$ est simple ☺;
- $\log p(\mathbf{B}|\mathbf{N})$ est complexe ☹.

Idée: Remplacer $\log p(\mathbf{B}|\mathbf{N})$ par une approximation **agréable** $\tilde{p}(\mathbf{B})$.

Approximation variationnelle

Soit $\tilde{p} = \tilde{p}_{M,S} = \otimes_i \mathcal{N}_p(\mathbf{m}_i, \text{Diag}(\mathbf{s}_i)^2)$ une **approximation** de $p(\mathbf{B}|\mathbf{N})$.

On a la **borne inférieure variationnelle** de la vraisemblance

$$\log p(\mathbf{N}; \Theta, \Sigma) \geq J(\mathbf{M}, \mathbf{S}, \Theta, \Sigma) := J(\tilde{p}, \Theta, \Sigma)$$

avec

$$J(\tilde{p}, \Theta, \Sigma) := \log p(\mathbf{N}; \Theta, \Sigma) - KL(\tilde{p}(\mathbf{B}) || p(\mathbf{B}|\mathbf{N}; \Theta, \Sigma))$$

Approximation variationnelle

Soit $\tilde{p} = \tilde{p}_{M,S} = \otimes_i \mathcal{N}_p(\mathbf{m}_i, \text{Diag}(\mathbf{s}_i)^2)$ une **approximation** de $p(\mathbf{B}|\mathbf{N})$.

On a la **borne inférieure variationnelle** de la vraisemblance

$$\log p(\mathbf{N}; \Theta, \Sigma) \geq J(\mathbf{M}, \mathbf{S}, \Theta, \Sigma) := J(\tilde{p}, \Theta, \Sigma)$$

avec

$$\begin{aligned} J(\tilde{p}, \Theta, \Sigma) &:= \log p(\mathbf{N}; \Theta, \Sigma) - KL(\tilde{p}(\mathbf{B}) || p(\mathbf{B}|\mathbf{N}; \Theta, \Sigma)) \\ &= \mathbb{E}_{\tilde{p}}[\log p(\mathbf{N}, \mathbf{B}; \Theta, \Sigma) - \log \tilde{p}(\mathbf{B})] \end{aligned}$$

Approximation variationnelle

Soit $\tilde{p} = \tilde{p}_{M,S} = \otimes_i \mathcal{N}_p(\mathbf{m}_i, \text{Diag}(\mathbf{s}_i)^2)$ une **approximation** de $p(\mathbf{B}|\mathbf{N})$.

On a la **borne inférieure variationnelle** de la vraisemblance

$$\log p(\mathbf{N}; \Theta, \Sigma) \geq J(\mathbf{M}, \mathbf{S}, \Theta, \Sigma) := J(\tilde{p}, \Theta, \Sigma)$$

avec

$$\begin{aligned} J(\tilde{p}, \Theta, \Sigma) &:= \log p(\mathbf{N}; \Theta, \Sigma) - KL(\tilde{p}(\mathbf{B}) || p(\mathbf{B}|\mathbf{N}; \Theta, \Sigma)) \\ &= \mathbb{E}_{\tilde{p}}[\log p(\mathbf{N}, \mathbf{B}; \Theta, \Sigma) - \log \tilde{p}(\mathbf{B})] \\ &= \mathbf{1}_n^\top [\mathbf{N} \odot (\mathbf{X}\Theta + \mathbf{M}) - \mathbf{A}] \mathbf{1}_p \\ &\quad - \frac{1}{2} \mathbf{1}_n^\top [\mathbf{M} \odot \mathbf{M} + \mathbf{S} \odot \mathbf{S} - 2 \log(\mathbf{S}) - \mathbf{1}_{n,p}] \mathbf{1}_p - K(\mathbf{N}). \end{aligned}$$

où

$$\mathbf{A} = \exp \left(\mathbf{X}\Theta + \mathbf{M} + \frac{1}{2}(\mathbf{S} \odot \mathbf{S}) \right)$$

Pour reconstruire le réseau, on ajoute de la structure sur $\Omega = \Sigma^{-1}$ via une contrainte de **sparsité** $\|\Omega\|_1$

Variational EM On maximise la borne J avec la structure sur $\|\Omega\|_1$

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

Pour reconstruire le réseau, on ajoute de la structure sur $\Omega = \Sigma^{-1}$ via une contrainte de **sparsité** $\|\Omega\|_1$

Variational EM On maximise la borne J avec la structure sur $\|\Omega\|_1$

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

On alterne entre les étapes suivantes

- Mise à jour de $(M, S, \Theta) = \arg \max_{M, S, \Theta} J(M, S, \Theta, \Omega)$
- Mise à jour de $\Omega = \arg \max_{\Omega} J(M, S, \Theta, \Omega) - \lambda \|\Omega\|_1$

Propriétés de J_q

- J n'est pas globalement concave;
- $J(\cdot, \cdot, \cdot, \cdot, \Omega) - \lambda \|\Omega\|_1$ concave en Ω : résolution via un lasso
- J concave in (M, S, Θ) : résolution via une descente de gradient

Propriétés de J_q

- J n'est pas globalement concave;
- $J(\cdot, \cdot, \cdot, \cdot, \Omega) - \lambda \|\Omega\|_1$ concave en Ω : résolution via un glasso
- J concave in (M, S, Θ) : résolution via une descente de gradient

En pratique

Tests de plusieurs procédures:

- Procédure itérative / Optimisation de J
- Remplacer glasso par Meinshausen et Bühlman

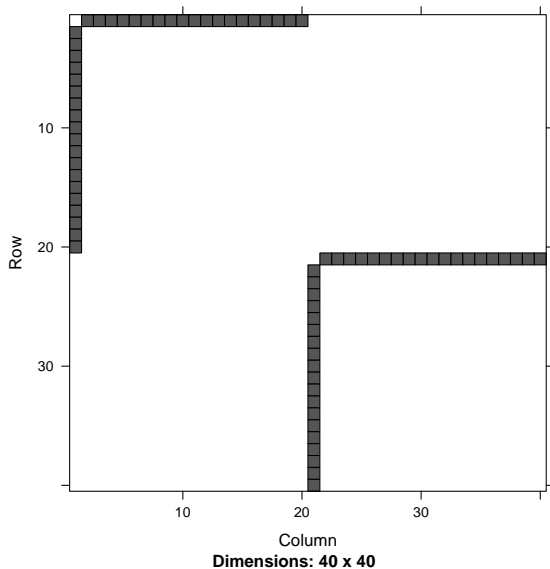
Avantages

- + Focus sur $\Omega = \Sigma^{-1}$
- + Basée sur une méthode de vraisemblance pénalisée
- + Modélisation des comptages et pas des fréquences

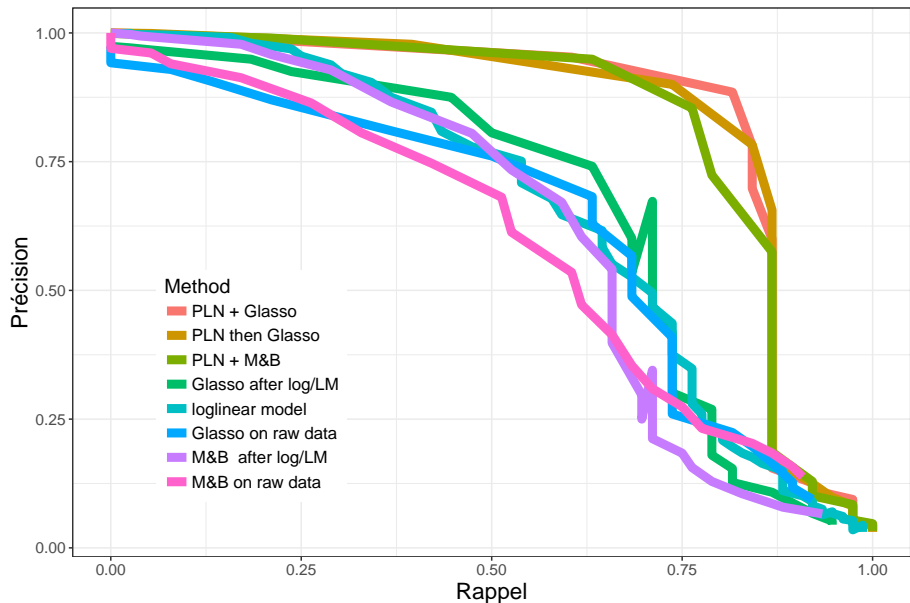
Inconvénients

- Inférence variationnelle (approchée)
- Choix de λ

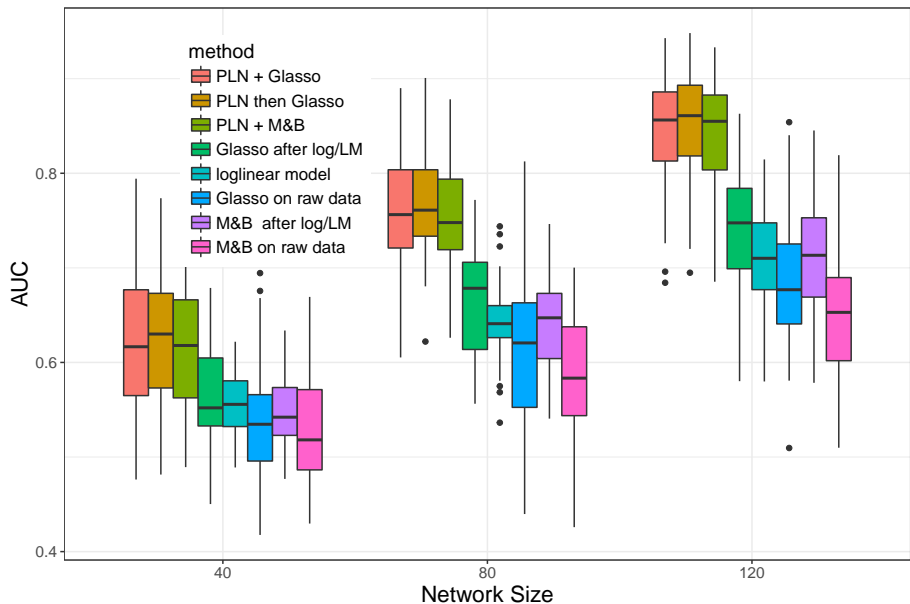
Un graphe simple (Hub)



Comparaison de différentes méthodes



(Petite) Étude de simulations



- 1 Introduction
- 2 Méthodes de Reconstruction
- 3 Modèles pour Données Compositionnelles
- 4 Modèles pour Données de Comptages
- 5 Conclusion et Perspectives

Interaction en écologie

- Définition assez floue
- Beaucoup de méthodes (SparCC, REBACCA, SPIEC-EASI, etc)
- Mais chacune avec ses limitations

Interaction en écologie

- Définition assez floue
- Beaucoup de méthodes (SparCC, REBACCA, SPIEC-EASI, etc)
- Mais chacune avec ses limitations

Modèle PLN

- Modèle générique pour données de comptages
- Adaptation *facile* aux familles exponentielles
- Algorithme variationnel pour l'inférence

Interaction en écologie

- Définition assez floue
- Beaucoup de méthodes (SparCC, REBACCA, SPIEC-EASI, etc)
- Mais chacune avec ses limitations

Modèle PLN

- Modèle générique pour données de comptages
- Adaptation *facile* aux familles exponentielles
- Algorithme variationnel pour l'inférence

Extensions

- Ajouter du *Zero Inflated* pour distinguer zéros structuraux / de sous-échantillonnage
- Rajouter une structuration spatiale / temporelle dans l'espace latent.

- Code **R/C++** disponible sur <https://github.com/jchiquet/PLNmodels>;
- Programmation OO avec des classes R6;

```
Y <- read.csv("myCounts.csv")
meta <- read.csv("myCovariates.csv")
## offset
O <- log(matrix(rep(rowSums(Y), ncol(Y)), ncol = ncol(Y)))
## Fit many models
PLNnets.models <- PLNnetwork(dat$Y)
## Choose best model
best.model <- PLNnets.models$getBestModel()
## Reconstruct and plot network
best.model$latentNetwork()
best.model$plot_network()
```

1. J. T. Staley and a. A. Konopka, "Measurement of in Situ Activities of Nonphotosynthetic Microorganisms in Aquatic and Terrestrial Habitats," *Annual Review of Microbiology*, vol. 39, no. 1, pp. 321–346, 1985.
2. W. Wade, "Unculturable bacteria—the uncharacterized organisms that cause oral infections," *Journal of the Royal Society of Medicine*, vol. 95, pp. 81–83, Feb. 2002.
3. G. Lima-Mendez, K. Faust, et al., "Determinants of community structure in the global plankton interactome," *Science*, vol. 348, p. 1262073, May 2015.
4. K. Faust and J. Raes, "Microbial interactions : from networks to models," *Nature Reviews Microbiology*, vol. 10, pp. 538–550, Aug. 2012.
5. K. Faust, J. F. Sathirapongsasuti, J. Izard, N. Segata, D. Gevers, J. Raes, and C. Huttenhower, "Microbial Co-occurrence Relationships in the Human Microbiome," *PLOS Computational Biology*, vol. 8, p. e1002606, July 2012.
6. J. Friedman and E. J. Alm, "Inferring Correlation Networks from Genomic Survey Data," *PLOS Computational Biology*, vol. 8, p. e1002687, Sept. 2012.
7. Y. Ban, L. An, and H. Jiang, "Investigating microbial co-occurrence patterns based on metagenomic compositional data," *Bioinformatics*, vol. 31, pp. 3322–3329, Oct. 2015.
8. J. Aitchison, "A new approach to null correlations of proportions," *Journal of the International Association for Mathematical Geology*, vol. 13, pp. 175–189, Apr. 1981.
9. L. C. Xia, D. Ai, J. Cram, J. A. Fuhrman, and F. Sun, "Efficient statistical significance approximation for local similarity analysis of high-throughput time series data," *Bioinformatics*, vol. 29, pp. 230–237, Jan. 2013.
10. Y. Deng, Y.-H. Jiang, Y. Yang, Z. He, F. Luo, and J. Zhou, "Molecular ecological network analyses," *BMC Bioinformatics*, vol. 13, p. 113, 2012.
11. Z. D. Kurtz, C. L. Müller, E. R. Miraldi, D. R. Littman, M. J. Blaser, and R. A. Bonneau, "Sparse and Compositionally Robust Inference of Microbial Ecological Networks," *PLOS Computational Biology*, vol. 11, p. e1004226, May 2015.
12. S. Weiss, W. Van Treuren, et al., "Correlation detection strategies in microbial data sets vary widely in sensitivity and precision," *The ISME Journal*, vol. 10, no. 7, pp. 1669–1681, 2016.
13. M. Layeghifard, D. M. Hwang, and D. S. Guttman, "Disentangling Interactions in the Microbiome: A Network Perspective" *Trends in Microbiology*, 2017, 25, 217-228