

Apprentissage non supervisé de structures de graphes

Analyse de réseaux d'interactions pour comprendre l'organisation d'écosystème

Sophie Donnet¹ et François Massol²

¹ MIA Paris-Saclay,



² Center for Infection and Immunity of Lille, CNRS



1. A propos de la collaboration
2. De l'intérêt des réseaux en écologie
3. Unsupervised clustering of nodes by probabilistic models
 - Matrix representation
 - Descriptive statistics
 - Probabilistic models for bipartite networks
 - Variational inference
4. Application

A propos de la collaboration

De l'intérêt des réseaux en écologie

Unsupervised clustering of nodes by probabilistic models

Application



Sophie Donnet

- Chercheuse en modélisation et apprentissage statistique
- Spécialisée en modèles probabilistes pour les réseaux d'interactions



François Massol

- Chercheur en écologie
- Axes de recherche : écologie spatiale, écologie des communautés et écologie évolutive

Intérêt commun

Réseaux d'interactions

- **MIRES**: Méthodes Interdisciplinaires pour les Réseaux d'Échanges de Semences
 - Effet des échanges de graines sur la biodiversité cultivée
- **GDR RESODIV** :
 - Etude des réseaux de circulation d'objets biologiques (plantes et animaux), et des savoirs et savoir-faire qui leurs sont associés, dans les agricultures des pays des Nords et des Suds.
- **ANR Econet**
 - Advanced statistical modelling of ecological networks
- **ANR NGB**
 - Biosurveillance Next-Gen des changements dans la structure et le fonctionnement des écosystèmes

- Coencadrement de la thèse de Tam Le Minh (avec S. Robin) sur la comparaison de réseaux écologiques
- Dispense de formations à destination d'étudiants et chercheurs sur l'analyse de réseaux
- Copublication

A propos de la collaboration

De l'intérêt des réseaux en écologie

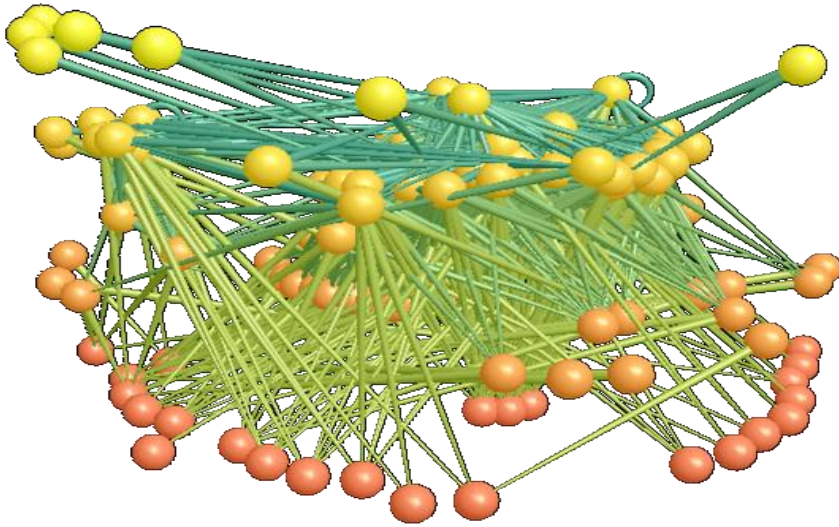
Unsupervised clustering of nodes by probabilistic models

Application

Question principale

Comment étudier les objets réseaux issus de l'écologie et de l'agronomie ?

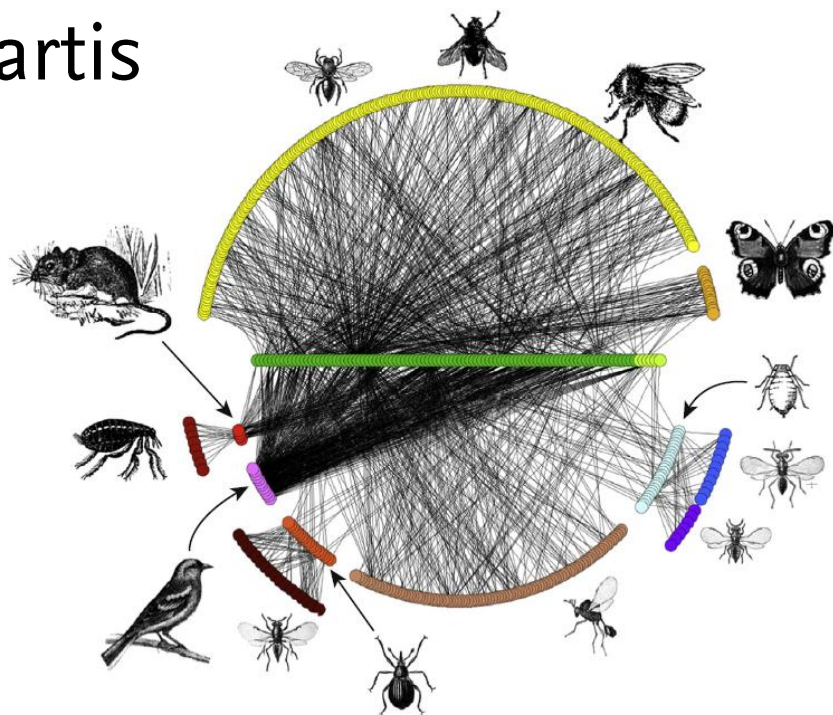
Réseaux écologiques



- Réseaux d'interactions, de cooccurrences, etc. entre espèces
- Réseaux de sites connectés dans l'espace
- De plus en plus de jeux de données (metabarcoding)

Réseaux écologiques

- Comprendre la dynamique des réseaux nécessite de comprendre leur structure
- Beaucoup de réseaux bipartis (2 niveaux d'espèces) ou multipartis (plus) = sans interactions intra-niveau



Réseaux plantes-pollinisateurs

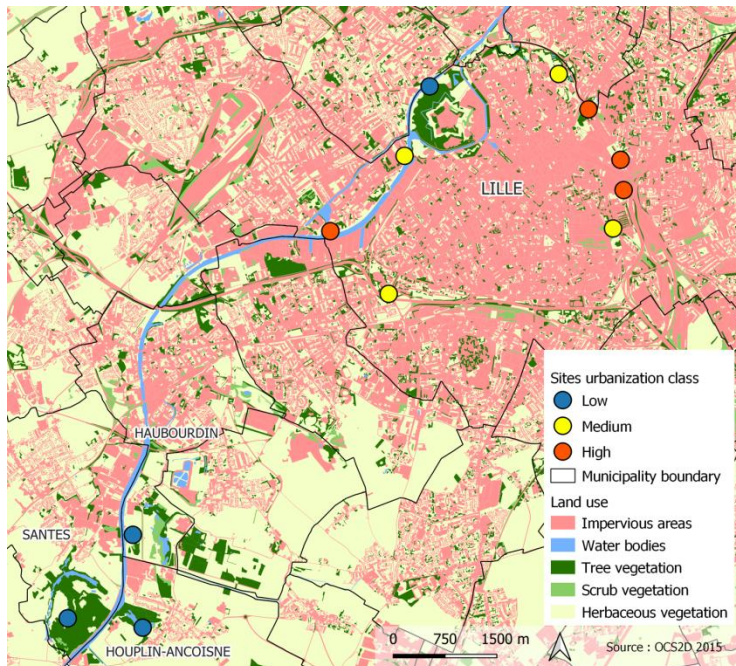


photo : N. de Manincor

- A la base du service de pollinisation (~ 87,5 % des angiospermes pollinisés par animaux, Ollerton et al. 2011)
- Dépendances réciproques entre plantes et pollinisateurs
⇒ potentielles extinctions en cascade
- Diversité & complexité des réseaux liées à leur stabilité

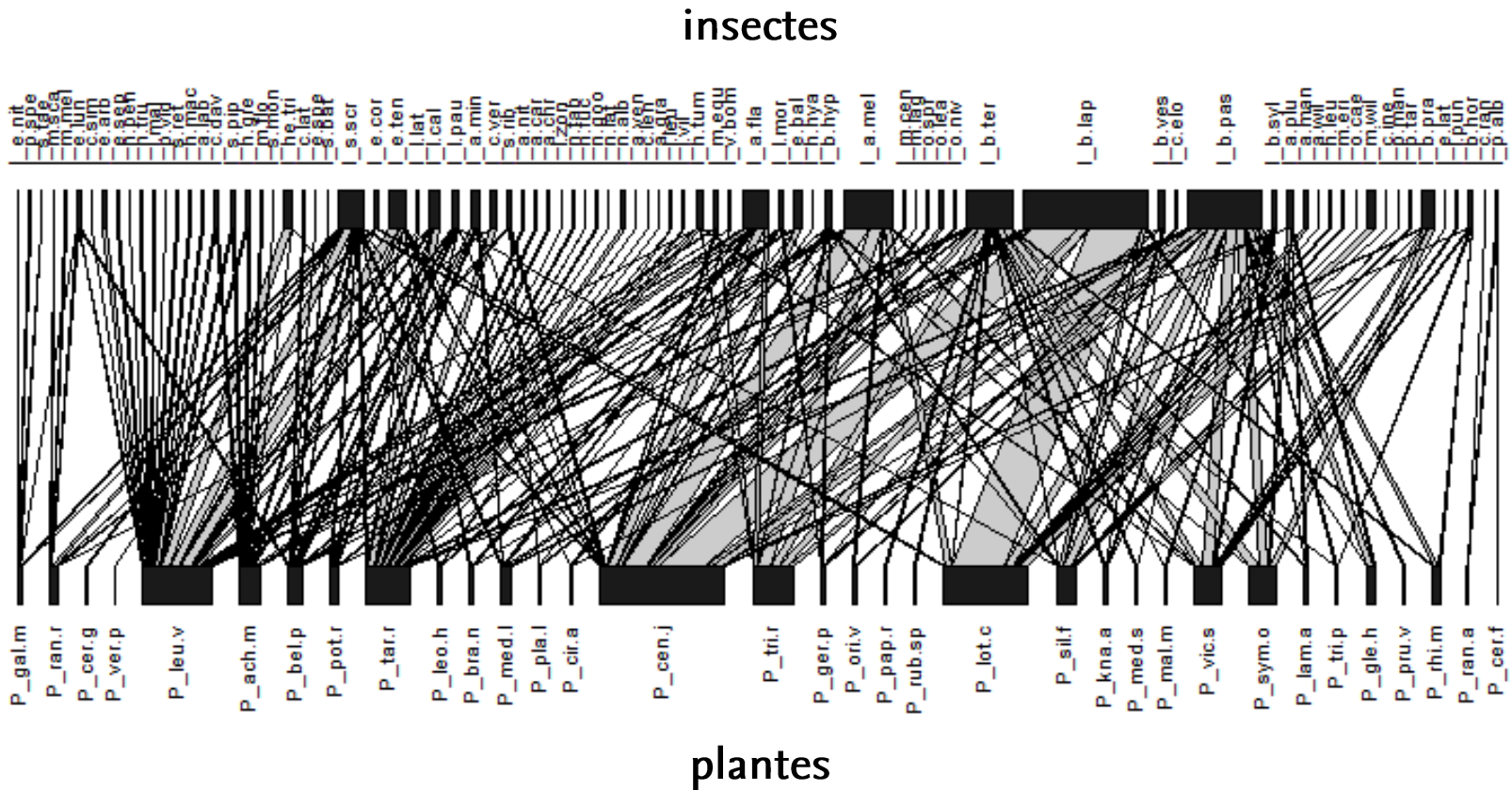
Jeu de données plantes-pollinisateurs

- Travail de post-doctorat d'Alessandro Fisogni
- Etude des réseaux plantes-pollinisateurs en ville
- 12 sites = 3 catégories d'urbanisation x 4 « répliqués » (sites semés par la MEL)



Jeu de données plantes-pollinisateurs

- Données SIG sur environnements
- Couverture florale



faible urbanisation, tous sites et dates confondus

Jeu de données plantes-pollinisateurs

- Premier travail publié sur la phénologie (dates de floraison et dates de sortie des insectes)



Research

Urbanization drives an early spring for plants but not for pollinators

Alessandro Fisogni, Nina Hautekèete, Yves Piquot, Marion Brun, Cédric Vanappelghem, Denis Michez and François Massol

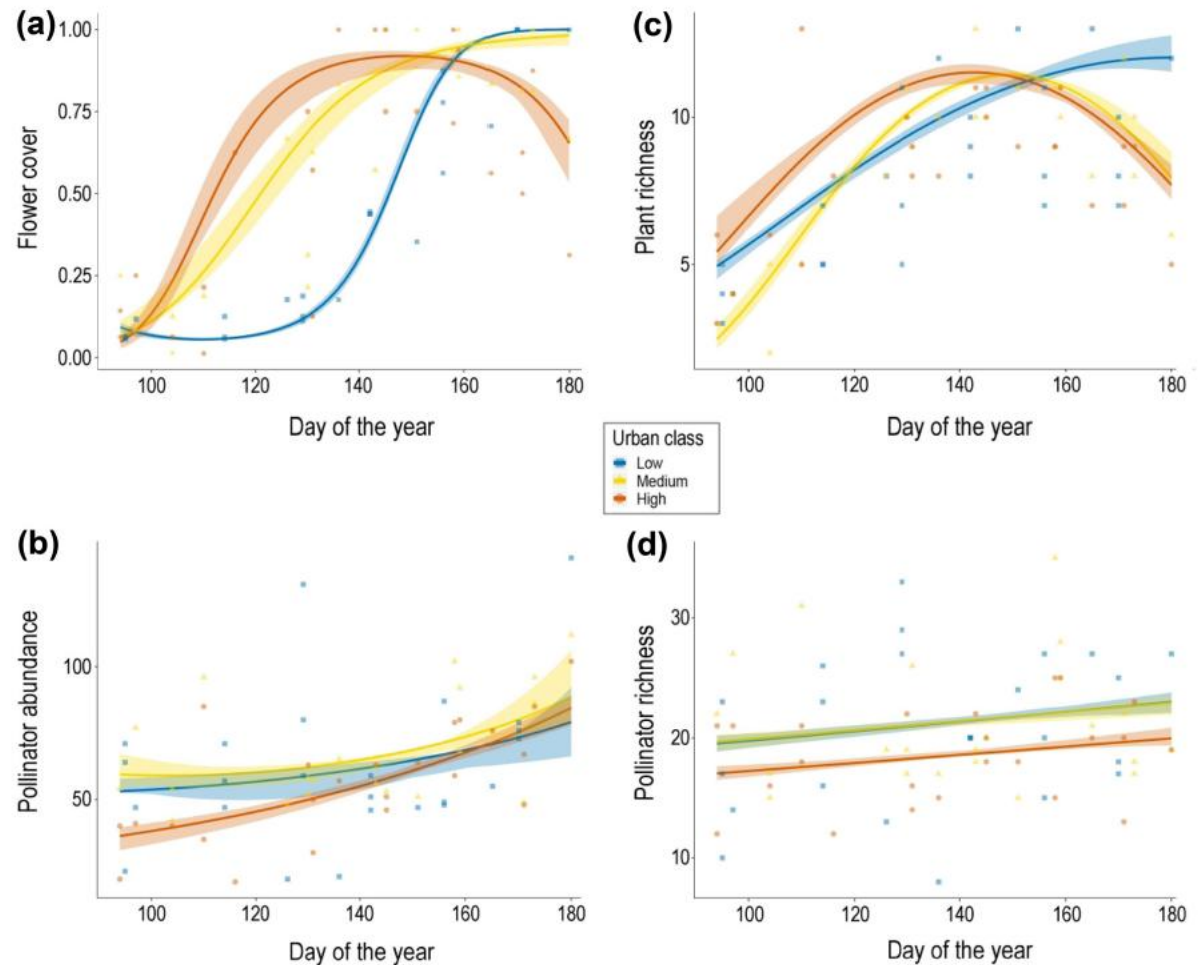
- Données réseaux = en révision...
- Données accessibles en ligne :

<https://zenodo.org/record/5570297>

Jeu de données plantes-pollinisateurs

Les plantes sont plus précoces en ville, mais pas les pollinisateurs

⇒ potentiel problème pour le réseau PP ?



A propos de la collaboration

De l'intérêt des réseaux en écologie

Unsupervised clustering of nodes by probabilistic models

Application

Objectif

- Highlight heterogeneity in the data (here nodes)
- Identify groups of nodes
 - Grouping nodes showing the same connection behavior
 - Each group represents a significantly different understanding from the other groups

A propos de la collaboration

De l'intérêt des réseaux en écologie

Unsupervised clustering of nodes by probabilistic models

Matrix representation

Descriptive statistics

Probabilistic models for bipartite networks

Variational inference

Application

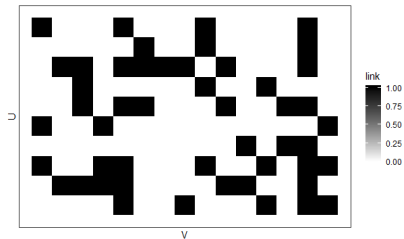
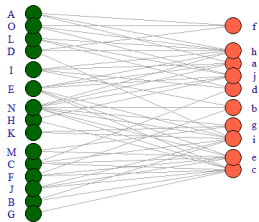
Matrix Representation

- Bipartite networks can be represented by an **incidence matrix** or **bi-adjacency matrix**
- For $i \in \text{🌱}, j \in \text{🐝}$,

$$Y_{ij} = \begin{cases} 1 & \text{if there is an interaction between } i \text{ and } j \\ 0 & \text{otherwise.} \end{cases}$$

- Rectangular matrix
- In most cases : $Y_{ij} \in \{0, 1\}$. However, sometimes $Y_{ij} \in \mathbb{R}$, **weighted bipartite graph**
- Directed bipartite graph : not classical. Proposition $Y_{ij} \in \{-1, 0, 1\}$

Matrix Representation



A propos de la collaboration

De l'intérêt des réseaux en écologie

Unsupervised clustering of nodes by probabilistic models

Matrix representation

Descriptive statistics

Probabilistic models for bipartite networks

Variational inference

Application

Aim : give a short description of the network, give a hint about its structure, look for heterogeneity in the connections

- Many metrics supplied for simple networks
- Have been extended to bipartite networks

R-packages

| Name | Usage |
|--------------|---|
| Networksis | Tool to simulate bipartite networks |
| enaR | Provides algorithms for the analysis of ecological networks |
| Netpredictor | Prediction of missing links in any given bipartite network |
| biGRAPH | Extension of the igraph library for bipartite graphs |
| bipartite | Visualising Bipartite Networks and Calculating Some (Ecological) Indices |

Degree

$$\begin{aligned} \deg(u) &= \sum_{v \in \mathcal{N}} (u \leftrightarrow v), & \deg(v) &= \sum_{u \in \mathcal{N}} (u \leftrightarrow v) \\ \deg_i &= \sum_{j=1}^n Y_{ij} & \deg_j &= \sum_{i=1}^n Y_{ij} \end{aligned}$$

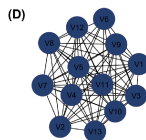
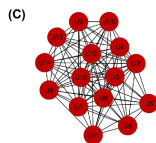
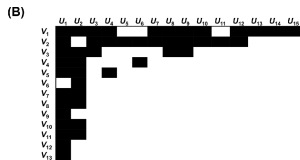
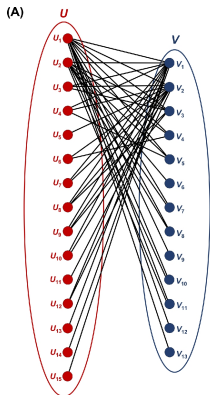
- Nodes with high degree are **hubs**
- Nodes with null degree are **isolated**
- If edges are oriented : in- and out- degrees can be computed.

Property on the network

Definition

- Important property in ecology
- Defined as a pattern of interactions in which specialists (e.g. pollinators that visit few plant species) interact with plants that are visited by generalists.
- Mathematically, looking for a reordering of rows and columns such that Y is nested

Nestedness



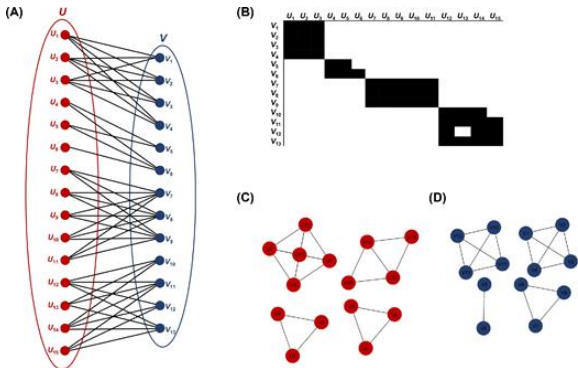
[PKP⁺18]

Modularity

Property on the network

Definition

Existence of clusters (blocks, module, communities) where nodes are much more connected than with other clusters



A propos de la collaboration

De l'intérêt des réseaux en écologie

Unsupervised clustering of nodes by probabilistic models

Matrix representation

Descriptive statistics

Probabilistic models for bipartite networks

Variational inference

Application

A first probabilistic model

- **Context:** our incidence matrix Y is the realization of a stochastic process.
- **Aim:** Propose a stochastic process is able to mimic heterogeneity in the connections.

Naive model





$$\forall (i, j) \in \text{🌱} \times \text{🐝}, \quad Y_{ij} \sim \text{Bern}(p)$$

- Homogeneity of the connections
- No hubs, no community, no nestedness

Latent Block Model

- **Aim** : introduce heterogeneity in the connections
- **Tool** : introduce blocks of nodes gathering entities that interact roughly similarly in the network

Latent Block Model with equations : latent variables i

- Each group of nodes ( and ) is divided into **blocks / clusters**
- K_{seedling} number of blocks in  and L_{bee} number of blocks in 
- For any $i \in \{1, \dots, n_{\text{seedling}}\}$, let Z_i^{seedling} be such that

$$Z_i^{\text{seedling}} = k \quad \text{if } \text{seedling } i \text{ belongs to cluster } k$$

- For any $j \in \{1, \dots, n_{\text{bee}}\}$, let W_j^{bee} be such that

$$W_j^{\text{bee}} = \ell \quad \text{if entity } j \text{ of group } \text{bee} \text{ belongs to cluster } \ell$$

Latent Block Model with equations : latent variables ii

Random latent variables

$(Z_i^{\text{🌱}})_{i=1\dots n}$ and $(W_j^{\text{🐝}})_{j=1\dots n}$ independent random variables, such that,

$$\mathbb{P}(Z_i^{\text{🌱}} = k) = \pi_k^{\text{🌱}},$$

$$\mathbb{P}(W_j^{\text{🐝}} = \ell) = \pi_\ell^{\text{🐝}}$$

with

$$\sum_{k=1}^{K^{\text{🌱}}} \pi_k^{\text{🌱}} = 1 \quad \text{and} \quad \sum_{\ell=1}^{L^{\text{🐝}}} \pi_\ell^{\text{🐝}} = 1$$

Latent Block Model with equations : connection probability

Conditionally to the latent variables...

$$C = (Z_{\text{plant}}, W_{\text{bee}}) = \{Z_i^{\text{plant}}, i = 1 \dots n_{\text{plant}}, W_j^{\text{bee}}, j = 1 \dots n_{\text{bee}}\} :$$

$$\mathbb{P}(Y_{ij} = 1 | Z_i^{\text{plant}} = k, W_j^{\text{bee}} = \ell) = \alpha_{k\ell}.$$

Other emission distributions

- Previous model adapted to 0-1 network
- If Y_{ij} is a count

$$Y_{ij} | Z_i^{\text{plant}} = k, W_j^{\text{bee}} = \ell \sim \mathcal{P}(\alpha_{k\ell})$$

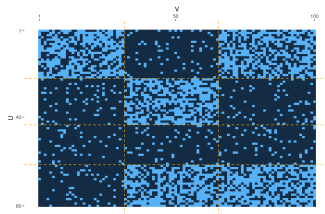
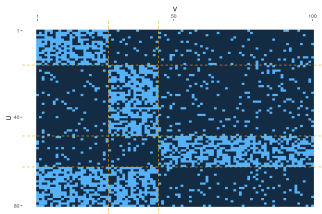
- If $Y_{ij} \in \mathbb{R}$

$$Y_{ij} | Z_i^{\text{plant}} = k, W_j^{\text{bee}} = \ell \sim \mathcal{N}(\alpha_{k\ell}, \sigma_{k\ell})$$

[GN08]

A very flexible model i

LBM able to generate **communities**...

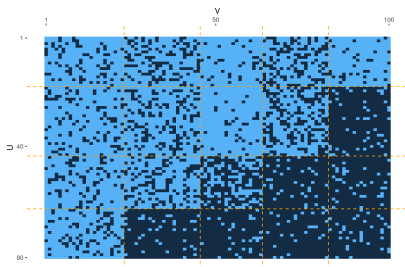


$$\alpha = \begin{pmatrix} 0.60 & 0.09 & 0.09 \\ 0.09 & 0.60 & 0.09 \\ 0.09 & 0.09 & 0.60 \\ 0.60 & 0.60 & 0.09 \end{pmatrix}$$

$$\alpha = \begin{pmatrix} 0.60 & 0.09 & 0.60 \\ 0.09 & 0.60 & 0.09 \\ 0.09 & 0.09 & 0.09 \\ 0.09 & 0.60 & 0.60 \end{pmatrix}$$

A very flexible model ii

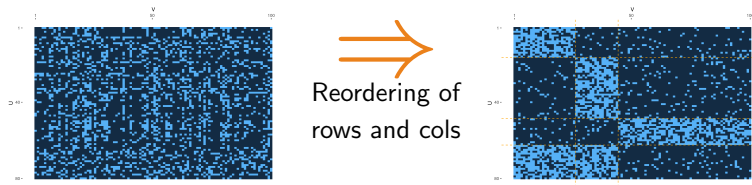
... or nested networks



$$\alpha = \begin{pmatrix} 0.80 & 0.70 & 0.90 & 0.60 & 0.90 \\ 0.80 & 0.70 & 0.90 & 0.60 & 0.09 \\ 0.80 & 0.70 & 0.40 & 0.09 & 0.09 \\ 0.80 & 0.09 & 0.09 & 0.09 & 0.09 \end{pmatrix}$$

Inference for LBM

Aim : From an incidence matrix, discovering the clusters



Remarks

- Looking for the blocks such that, under the assumption that my data come from the LBM model, the observed data Y is most probable (= most likely to occur)
- No specific prior structure
- Entities gathered because they have similar behavior in the network

A propos de la collaboration

De l'intérêt des réseaux en écologie

Unsupervised clustering of nodes by probabilistic models

Matrix representation

Descriptive statistics

Probabilistic models for bipartite networks

Variational inference

Application

Maximum likelihood inference

$$\begin{aligned}\hat{\theta} &= \operatorname{argmax}_{\theta} p(Y; \theta) = \operatorname{argmax}_{\theta} \sum_C p_{\theta}(Y, C) \\ &= \operatorname{argmax}_{\theta} \sum_{Z^{\text{🌱}}, W^{\text{🐝}}} p_{\theta}(Y, Z^{\text{🌱}}, W^{\text{🐝}})\end{aligned}$$

- Complete likelihood: $p_{\theta}(Y, Z^{\text{🌱}}, W^{\text{🐝}})$ easy to compute
- Likelihood $p(Y; \theta)$: integration over all the possible row and column clusterings $(Z^{\text{🌱}}, W^{\text{🐝}})$ (sum of $K^{\text{🌱}n} \times L^{\text{🐝}n}$ terms)
- Latent variables \Rightarrow **Expectation-Maximization**
 - Requires to evaluate $p(Z^{\text{🌱}}, W^{\text{🐝}} | Y; \theta)$
 - No independence in this distribution
 - Complicated distribution

Variational Inference

- Use a variational version of the Expectation-Maximization algorithm [DPR08, BKM17, MRV10]
- Penalized criterion to select the numbers of blocks $K_{\text{🌱}}$ and $L_{\text{🐝}}$.
- R-package `sbm` + `blockmodels`

▶ Skip details

Variational inference

Principle [WJ08, BKM17].

- Choose a divergence measure $D(q \parallel p)$
- Choose a class of distributions \mathcal{Q}
- Maximize w.r.t. θ and $q \in \mathcal{Q}$ the lower bound

$$J(Y; \theta, q) = \log p_\theta(Y) - D(q(C) \parallel p_\theta(C \mid Y)) \leq \log p_\theta(Y)$$

Popular choice for SBMs. [GN08, DPR08, Leg16, MM15]

- $D = KL$:

$$\begin{aligned} J(Y; \theta, q) &= \log p_\theta(Y) - KL(q(C) \parallel p_\theta(C \mid Y)) \\ &= \mathbb{E}_q(\log p_\theta(Y, C)) + \mathcal{H}(q(C)) \end{aligned}$$

- q factorizable: $\mathcal{Q} = \left\{ q(C) : q(C) = \prod_i q_i(Z_i) \prod_j q_j(W_j) \right\}$

$$\tau_{ik} = \mathbb{P}_q(Z_i = k)$$

→ mean field approximation

Algorithm

At iteration (t) , given $(\theta^{(t-1)}, q_{\tau^{(t-1)}})$,

- **Step 1** Maximization w.r.t. τ

$$\begin{aligned}\tau^{(t)} &= \arg \max_{\tau \in \mathcal{T}} J(Y; \theta^{(t-1)}, q_{\tau}) \\ &= \arg \max_{\tau \in \mathcal{T}} \mathbb{E}_{q_{\tau}} [\log p_{\theta^{(t-1)}}(Y, C)] + \mathcal{H}(q_{\tau}(C)) \\ &= \arg \min_{\tau \in \mathcal{T}} \text{KL}[q_{\tau}, p(\cdot | Y; \theta^{(t-1)})]\end{aligned}$$

- **Step 2** Maximization w.r.t. θ

$$\begin{aligned}\theta^{(t)} &= \arg \max_{\theta} J(Y; \theta, q_{\tau^{(t)}}) \\ &= \arg \max_{\theta} \mathbb{E}_{q_{\tau^{(t)}}} [\log p_{\theta}(Y, C)]\end{aligned}$$

Outputs

For a given $(K_{\text{🌱}}, L_{\text{🐝}})$

- Parameter estimates $\hat{\theta}$: $\hat{\alpha}$ and $\hat{\pi}_{\text{🌱}}, \hat{\pi}_{\text{🐝}}$
- The best clustering of nodes

$$\hat{C} = \underset{C}{\operatorname{argmax}} p(C|Y; \hat{\theta}, \mathcal{M}) \approx \underset{C}{\operatorname{argmax}} q_{\hat{\tau}}(C|Y; \hat{\theta}, \mathcal{M})$$

Choice of the numbers of blocks

ICL : penalized criterion that relies on the complete likelihood

For any model \mathcal{M} defined by $(K_{\text{🌱}}, L_{\text{🐝}})$

$$ICL(\mathcal{M}) = \log p_{\theta}(Y, \hat{C}; \mathcal{M}) - \frac{1}{2} pen_{\mathcal{M}}$$

with

$$pen_{\mathcal{M}} = \underbrace{(K_{\text{🌱}} - 1) \log n_{\text{🌱}} + (L_{\text{🐝}} - 1) \log n_{\text{🐝}}}_{\text{Clustering}} + \underbrace{K_{\text{🌱}} L_{\text{🐝}} \log(n_{\text{🌱}} n_{\text{🐝}})}_{\text{Connection}}$$

A propos de la collaboration

De l'intérêt des réseaux en écologie

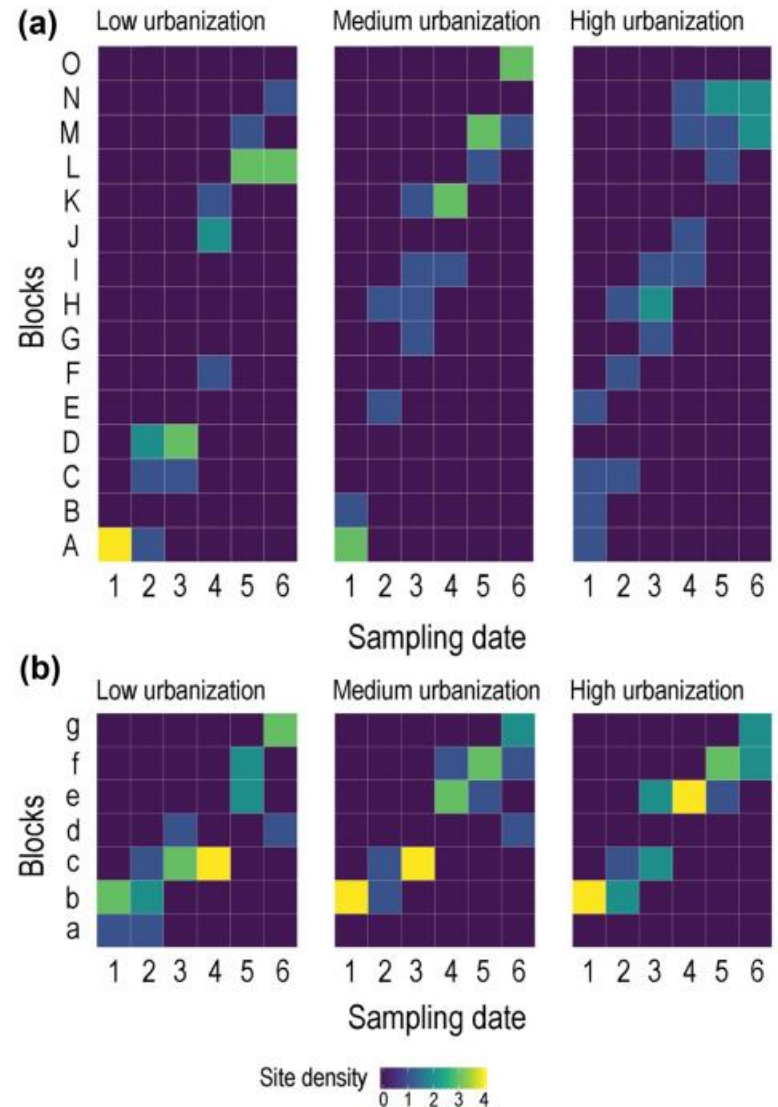
Unsupervised clustering of nodes by probabilistic models

Application

Applications

Mise en évidence du décalage phénologique en utilisant un clustering des (sites x dates) x espèces

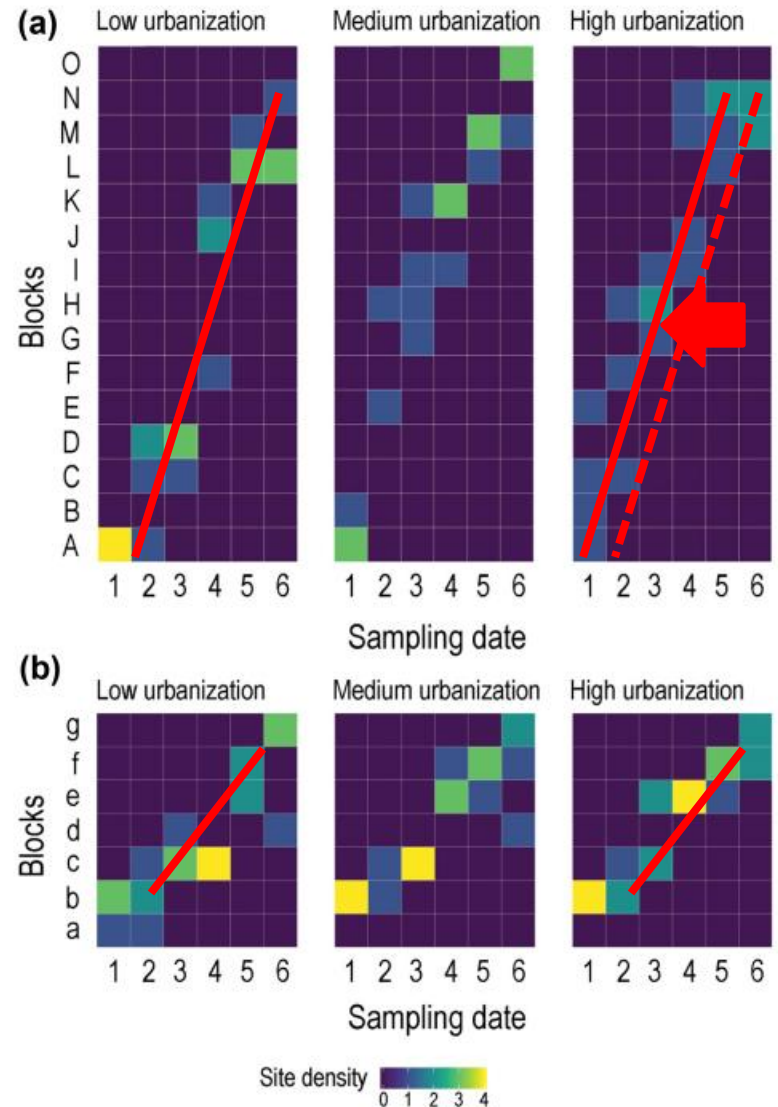
Représentation des densités de site par date et par bloc



Applications

Mise en évidence du décalage phénologique en utilisant un clustering des (sites x dates) x espèces

Représentation des densités de site par date et par bloc



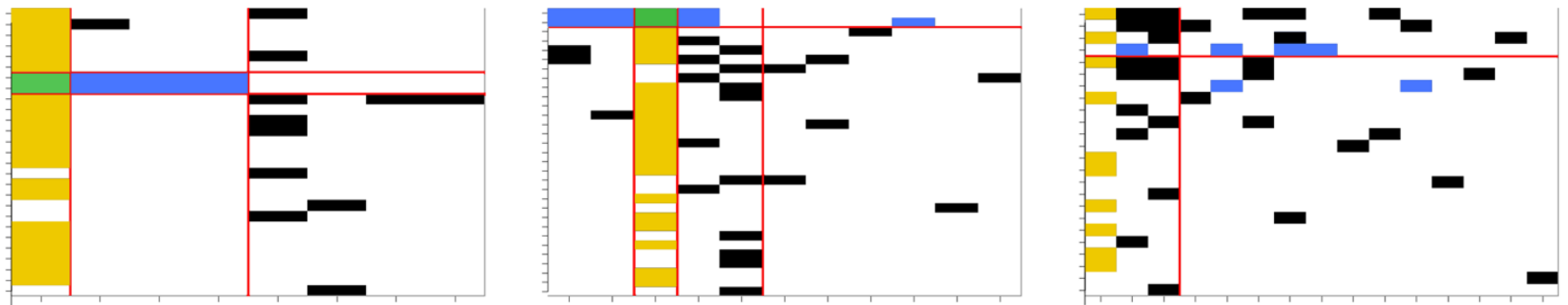
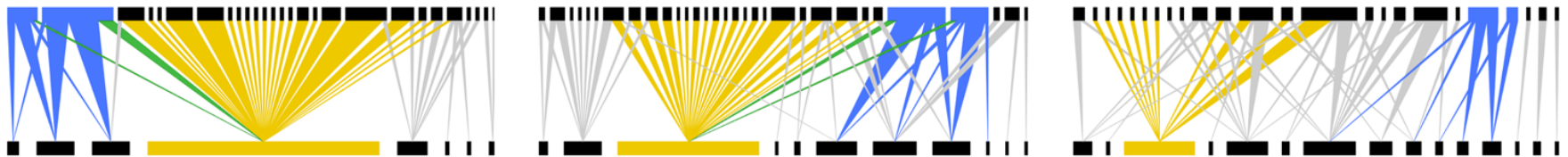
Applications

Mois d'avril seulement (plus clair)

Low urbanization

Medium urbanization

High urbanization



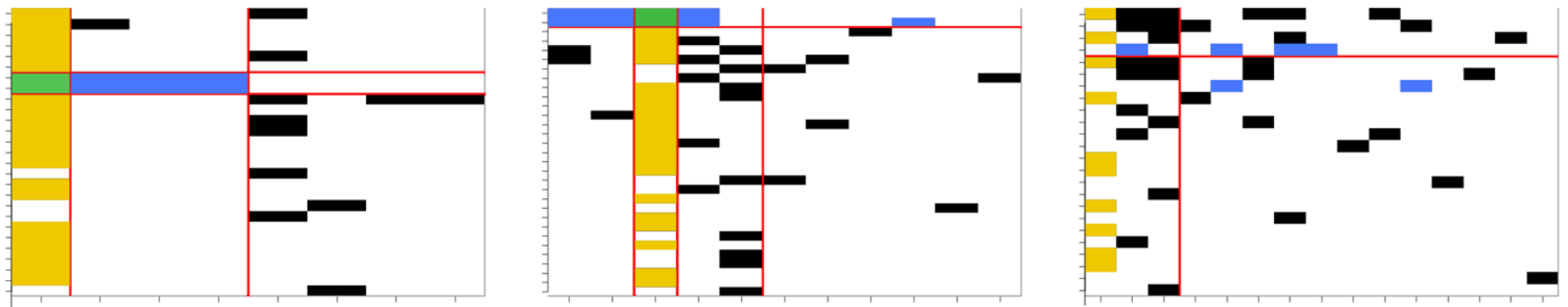
jaune = *Taraxacum* section *ruderalia*

bleu = *Anthophora plumipes* & *Bombus pascuorum*

Applications

Mois d'avril seulement (plus clair)

Permet de mettre en évidence une réorganisation des blocs de plantes et d'insectes entre types de site



jaune = *Taraxacum* section *ruderalia*

bleu = *Anthophora plumipes* & *Bombus pascuorum*

Conclusions

- Sur la saison entière, les espèces super-généralistes forment un bloc par espèce
- La phénologie semble un bon indicateur des blocs
- Les blocs sont liés aux degrés
- Réorganisation des blocs entre types de site

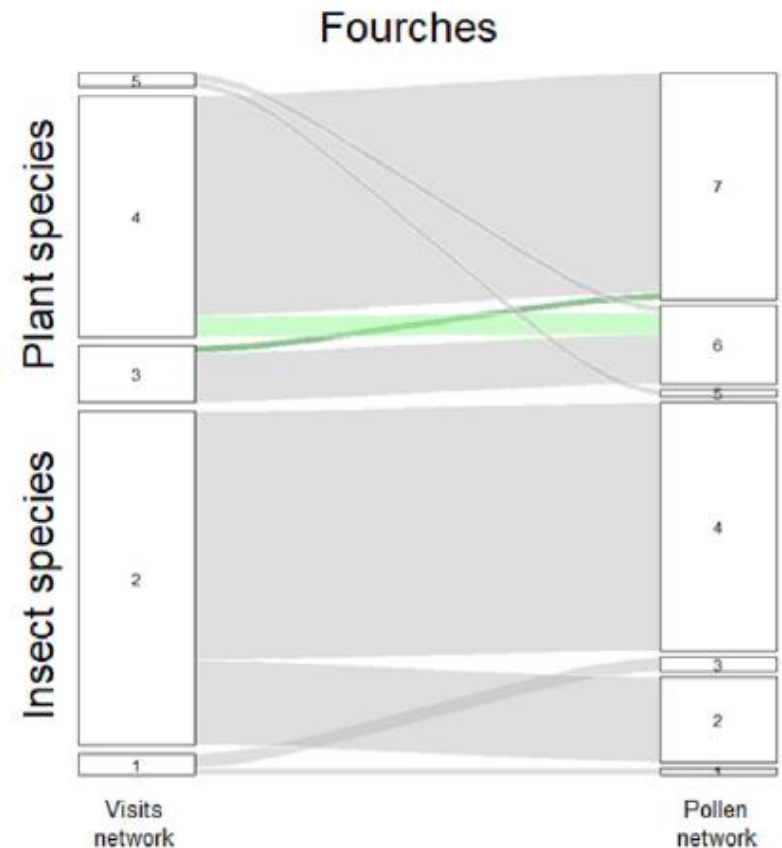
Perspectives

1. Variation et comparaisons des réseaux

Utilisation des LBM pour comparer la congruence générale des réseaux

Comparaison des blocs obtenus par LBM pour le même « réseau », observé de deux manières différentes (observation des insectes sur fleurs vs. pollens sur insectes)

De Manincor et al. 2020



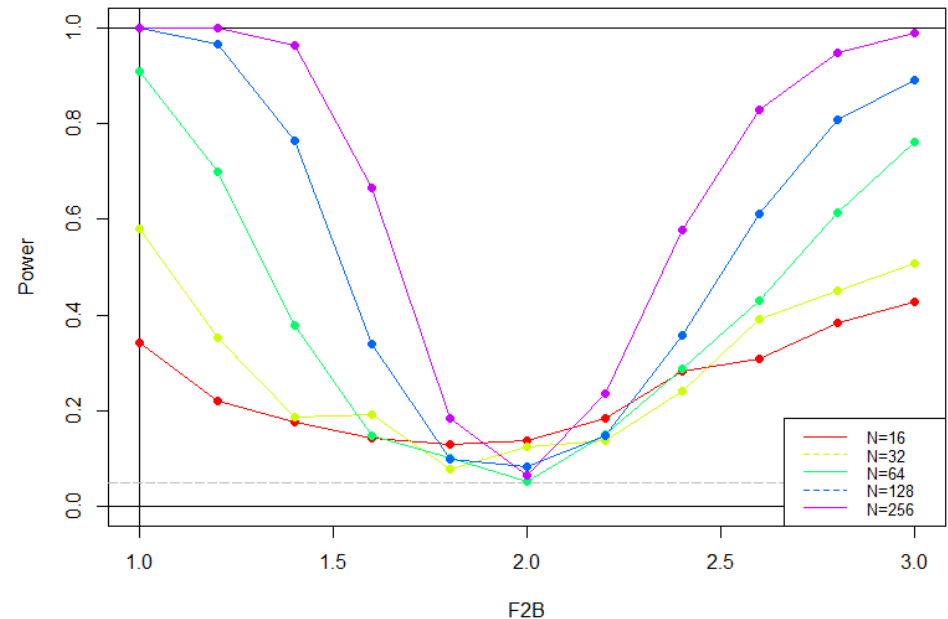
Perspectives

1. Variation et comparaisons des réseaux

Utilisation des U-statistiques pour mettre en évidence les différences de propriétés des réseaux

Puissance d'un test de comparaison (F2) entre réseaux fondé sur le modèle WBEDD (thèse de Tâm Le Minh)

Dubart et al. 2022



Perspectives

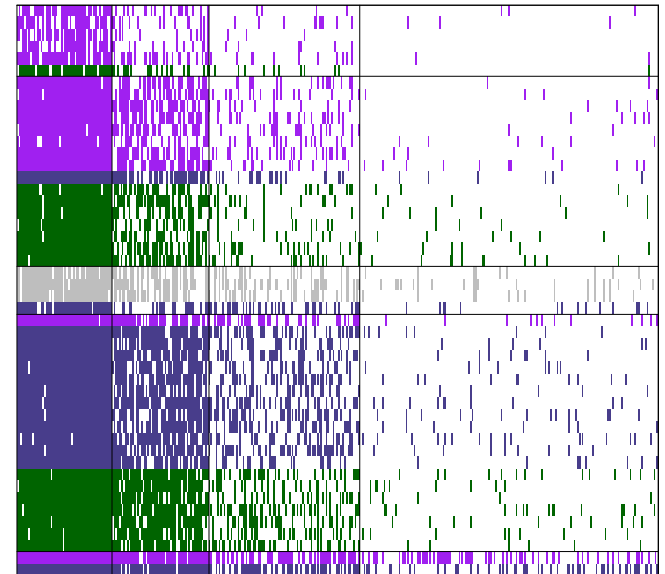
2. LBM sur réseaux plus grands

Problème potentiel = taille des réseaux, en particulier gros réseaux (ex. microbiote)

- solution 1 : travailler à l'échelle de la famille

LBM sur réseaux de microbiote (espèces de moules x OTU obtenues sur séquençage 16S)
couleurs = espèces de *Mytilus*
colonnes = familles de bactéries

Ben Cheikh & Massol, in prep.



Perspectives

3. Metabarcoding et reconstruction de réseaux

Echantillonnage coûteux en temps, dépendant d'une expertise taxonomique rare

⇒ passage au metabarcoding ?

Pour les espèces microbiennes, pas moyen de faire autrement...

Perspectives

3. Metabarcoding et reconstruction de réseaux

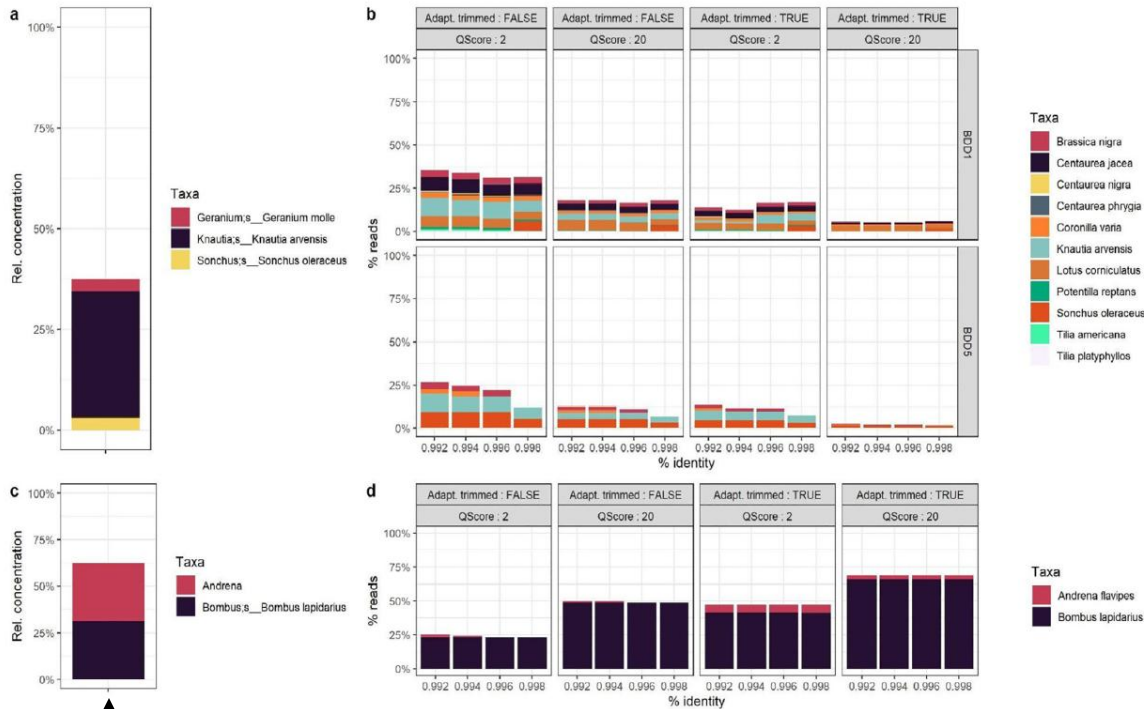
Difficultés multiples :

- abondances des espèces à partir d'ADN ?
- taux d'erreur d'assignation ?
- taux d'espèces non perçues ?

- quels pipelines bioinformatiques ?
- quelles séquences utiliser ?
- quel coût ?

Perspectives

3. Metabarcoding et reconstruction de réseaux



Tests sur communautés
« mocks » pour retrouver
plantes et pollinisateurs selon
différentes méthodes bioinfo

Dubart et al. 2022

ce qu'on retrouve après bioinfo...

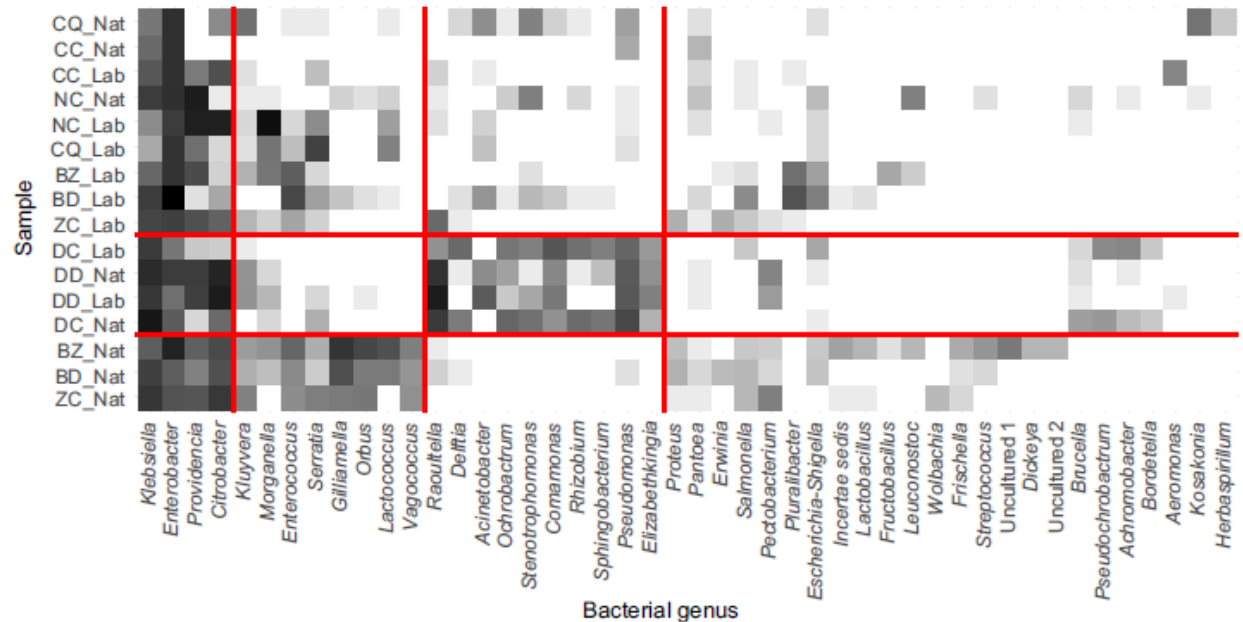
ce qu'on met dans le tube

Perspectives





3. Metabarcoding et reconstruction de réseaux

LBM entre espèces de mouches, élevées au labo ou échantillonnées sur le terrain et leur microbiote

Ravné et al. 2022



References

-  D. M. Blei, A. Kucukelbir, and J. D. McAuliffe.
Variational inference: A review for statisticians.
Journal of the American Statistical Association, 112(518):859–877, 2017.
-  J.-J. Daudin, F. Picard, and S. Robin.
A mixture model for random graphs.
Stat. Comput., 18(2):173–83, 2008.
-  Gérard Govaert and Mohamed Nadif.
Block clustering with bernoulli mixture models: Comparison of different approaches.
Comput. Stat. Data Anal., 52(6):3233–3245, February 2008.
-  Jean-Benoist Leger.
Blockmodels: A R-package for estimating in latent block model and stochastic block model, with various probability functions, with or without covariates.
Technical report, arXiv:1602.07587, 2016.

Remerciements

A. Fisogni, N. de Manincor, M. Dubart, T. Le Minh,
S.-C. Chabert-Liddell

ARSENIC, NGB & EcoNet ANR consortiums

Groupe MIRES, GDR Resodiv

