

Détection de structures et *clustering* dans des réseaux bipartites

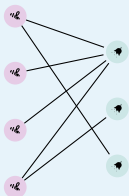
Soutenance de Master MathSV

Louis Lacoste, encadré par Pierre Barbillon et Sophie Donnet
Laboratoire MIA Paris-Saclay

29 août 2024

Contexte écologique

- Nombreux réseaux disponibles pour interactions similaires.
- Suivi biodiversité, robustesse et risque d'effondrement . . .



$$\begin{pmatrix} 1 & 0 & 1 \\ 1 & 0 & 0 \\ 1 & 0 & 0 \\ 1 & 1 & 0 \end{pmatrix}$$

Figure 1 – Exemple
d'un réseau
plantes-pollinisateurs

Matrice
d'adjacence
associée

Contexte écologique

- Nombreux réseaux disponibles pour interactions similaires.
- Suivi biodiversité, robustesse et risque d'effondrement ...

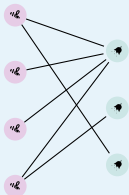


Figure 1 – Exemple d'un réseau plantes-pollinisateurs

$$\begin{pmatrix} 1 & 0 & 1 \\ 1 & 0 & 0 \\ 1 & 0 & 0 \\ 1 & 1 & 0 \end{pmatrix}$$

Matrice d'adjacence associée

Contexte mathématique

Pour un unique réseau : variables latentes, *embedding*, ...

Motivations pour proposer des méthodes adaptées aux collections de réseaux :

- Espèces différentes, rôles analogues.
- Transfert d'informations grands vers petits réseaux.
- Regrouper les réseaux selon leur similarité (*clustering* de réseaux).

Latent Block Model (LBM¹)

Proposé par Govaert et Nadif, 2005.

Pour

- $Q_1 = |\{\bullet, \bullet, \bullet\}|$ blocs fixés en ligne
- $Q_2 = |\{\bullet, \bullet, \bullet\}|$ blocs fixés en colonne

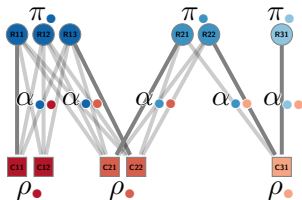


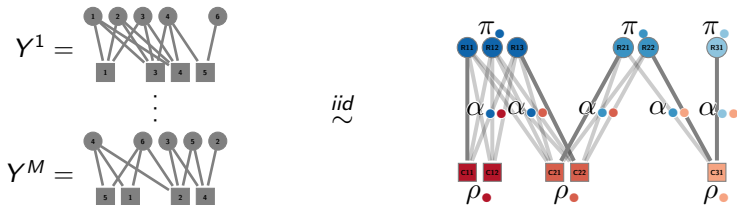
Figure 2 – Exemple de LBM¹

Paramètres

- $\pi_{\bullet} = \mathbb{P}(Z_i = \bullet)$ en ligne et
- $\rho_{\bullet} = \mathbb{P}(W_j = \bullet)$ en colonne
- $\alpha_{\bullet, \bullet} = \mathbb{P}(X_{ij} = 1 | Z_i = \bullet, W_j = \bullet)$

1. Que j'appellerai par la suite BiSBM

Collections bipartites

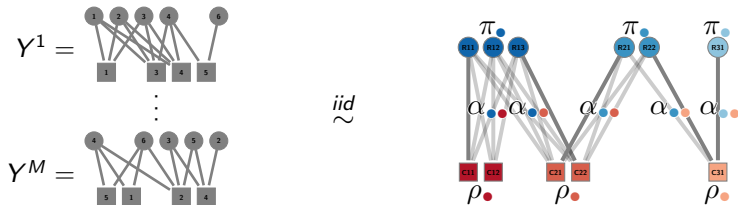


- $Q_1 = |\{\bullet, \bullet, \bullet\}|$ blocs fixés en ligne
- $Q_2 = |\{\bullet, \bullet, \bullet\}|$ blocs fixés en colonne

Paramètres

- $\pi_{\bullet} = \mathbb{P}(Z_i = \bullet)$ en ligne et $\rho_{\bullet} = \mathbb{P}(W_j = \bullet)$ en colonne
- $\alpha_{\bullet, \bullet} = \mathbb{P}(X_{ij} = 1 | Z_i = \bullet, W_j = \bullet)$

Différents modèles

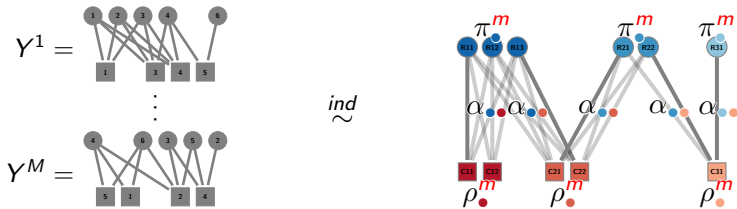


iid-colBiSBM

$$\boldsymbol{\pi} = (\pi_1, \dots, \pi_{Q_1}) \text{ et } \boldsymbol{\rho} = (\rho_1, \dots, \rho_{Q_2})$$

Dans tous les modèles la structure de connectivité (α) est supposée identique au sein de la collection.

Différents modèles



$\pi\rho$ -colBiSBM

$\pi = ((\pi_1^m, \dots, \pi_{Q_1}^m))_{m=1, \dots, M}$ et $\rho = ((\rho_1^m, \dots, \rho_{Q_2}^m))_{m=1, \dots, M}$
 avec $\forall q, m \in \llbracket 1, Q_1 \rrbracket \times \llbracket 1, M \rrbracket, \pi_q^m \in [0, 1]$ et $\forall r, m \in \llbracket 1, Q_2 \rrbracket \times \llbracket 1, M \rrbracket, \rho_r^m \in [0, 1]$

Dans tous les modèles la structure de connectivité (α) est supposée identique au sein de la collection.

Estimation des paramètres

En adaptant Chabert-Liddell et al., 2024 qui se base sur la méthode proposée par Daudin et al., 2008 utilisant l'algorithme *Variational EM*.

$$\ell(\mathbf{X}; \theta) \geq \sum_{m=1}^M \left(Q^m(\theta \mid \theta^{(t)}) + \mathcal{H}(\mathcal{R}_{\mathbf{X}^m, \theta^{(t)}}(\mathbf{Z}^m, \mathbf{W}^m)) \right) =: J(\tau; \theta)$$

où $Q^m(\theta \mid \theta^{(t)}) = \mathbb{E}_{\mathbf{Z}^m, \mathbf{W}^m \sim \mathcal{R}_{\mathbf{X}^m, \theta^{(t)}}(\cdot)} [\log p(\mathbf{X}^m, \mathbf{Z}^m, \mathbf{W}^m \mid \theta)]$

Approximation variationnelle

$\mathcal{R}_{\mathbf{X}^m, \theta^{(t)}}(\mathbf{Z}^m, \mathbf{W}^m) = P(\mathbf{Z}^m \mid \mathbf{X}^m, \theta^{(t)})P(\mathbf{W}^m \mid \mathbf{X}^m, \theta^{(t)})$, c'est à dire avoir une indépendance lignes, colonnes.

Formule développée de l'EM variationnel

$$\begin{aligned}
 \ell(\mathbf{X}; \theta) \geq & \sum_{m=1}^M \left(\sum_{i=1}^{n_1^m} \sum_{j=1}^{n_2^m} \sum_{q \in \mathcal{Q}_{1,m}} \sum_{r \in \mathcal{Q}_{2,m}} \tau_{i,q}^{1,m} \tau_{j,r}^{2,m} \log f(X_{ij}^m; \alpha_{qr}) \right. \\
 & + \sum_{i=1}^{n_1^m} \sum_{q \in \mathcal{Q}_{1,m}} \tau_{i,q}^{1,m} \log \pi_q^m + \sum_{j=1}^{n_2^m} \sum_{r \in \mathcal{Q}_{2,m}} \tau_{j,r}^{2,m} \log \rho_r^m \\
 & \left. - \sum_{i=1}^{n_1} \tau_{i,q}^{1,m} \log \tau_{i,q}^{1,m} - \sum_{j=1}^{n_2} \tau_{j,r}^{2,m} \log \tau_{j,r}^{2,m} \right) =: J(\tau; \theta),
 \end{aligned}$$

où ℓ désigne la log vraisemblance.

Approximation variationnelle

$$\tau_{iq}^{1,m} = P_{\mathcal{R}_m}(Z_{iq}^m = 1 | X_{i\bullet}^m) \text{ et } \tau_{jr}^{2,m} = P_{\mathcal{R}_m}(W_{jr}^m = 1 | X_{\bullet j}^m)$$

Étape *Variational Expectation*

$$\hat{\tau}^{(t+1)} = \arg \max_{\tau} \mathcal{J}(\tau, \hat{\theta}^{(t)})$$

$$\begin{cases} \hat{\tau}_{iq}^{1,m} \propto \hat{\pi}_q^{m(t)} \prod_{j=1}^{n_2^m} \prod_{r \in \mathcal{Q}_2^m} f(X_{ij}^m; \hat{\alpha}_{qr}^{(t)}) \hat{\tau}_{jr}^{2,m(t+1)} & \forall i = 1, \dots, n_1^m, q \in \mathcal{Q}_1^m \\ \hat{\tau}_{jr}^{2,m} \propto \hat{\rho}_r^{m(t)} \prod_{i=1}^{n_1^m} \prod_{q \in \mathcal{Q}_1^m} f(X_{ij}^m; \hat{\alpha}_{qr}^{(t)}) \hat{\tau}_{iq}^{1,m(t+1)} & \forall j = 1, \dots, n_2^m, r \in \mathcal{Q}_2^m \end{cases}$$

2. Initialisation des $\hat{\tau}$ avec un *spectral clustering* sur les réseaux.

Étape *Maximization*

$$\hat{\theta}^{(t+1)} = \arg \max_{\theta} \mathcal{J}(\hat{\tau}^{(t+1)}, \theta)$$

Paramètres de connectivité

$$\hat{\alpha}_{qr} = \frac{\sum_{m=1}^M \sum_{i=1}^{n_1^m} \sum_{j=1}^{n_2^m} \tau_{iq}^{1,m} \tau_{jr}^{2,m} X_{ij}^m}{\sum_{m=1}^M \sum_{i=1}^{n_1^m} \sum_{j=1}^{n_2^m} \tau_{iq}^{1,m} \tau_{jr}^{2,m}}$$

Proportions pour *iid*

$$\hat{\pi}_q = \frac{\sum_{m=1}^M \sum_{i=1}^{n_1^m} \tau_{iq}^{1,m}}{\sum_{m=1}^M n_1^m}$$

$$\hat{\rho}_r = \frac{\sum_{m=1}^M \sum_{j=1}^{n_2^m} \tau_{jr}^{2,m}}{\sum_{m=1}^M n_2^m}$$

Étape *Maximization*

$$\hat{\theta}^{(t+1)} = \arg \max_{\theta} \mathcal{J}(\hat{\tau}^{(t+1)}, \theta)$$

Paramètres de connectivité

$$\hat{\alpha}_{qr} = \frac{\sum_{m=1}^M \sum_{i=1}^{n_1^m} \sum_{j=1}^{n_2^m} \tau_{iq}^{1,m} \tau_{jr}^{2,m} X_{ij}^m}{\sum_{m=1}^M \sum_{i=1}^{n_1^m} \sum_{j=1}^{n_2^m} \tau_{iq}^{1,m} \tau_{jr}^{2,m}}$$

Proportions pour $\pi\rho$

$$\hat{\pi}^m_q = \frac{\sum_{i=1}^{n_1^m} \tau_{iq}^{1,m}}{n_1^m}$$

$$\hat{\rho}^m_r = \frac{\sum_{j=1}^{n_2^m} \tau_{jr}^{2,m}}{n_2^m}$$

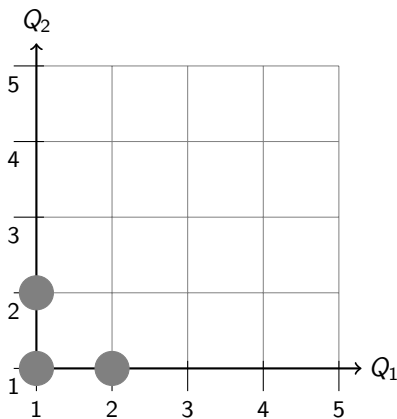
Problème de choix de (Q_1, Q_2)

L'estimation de paramètres se fait à Q_1, Q_2 blocs fixés, il faut donc déterminer les “meilleures” coordonnées. Nous maximisons un critère, le *Bayesian Information Criterion - Like* (BIC-L), de vraisemblance pénalisée en adaptant les formules de [Chabert-Liddell et al., 2024](#).

Problèmes de l'exploration

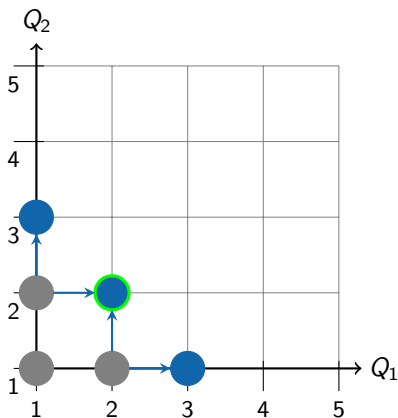
- Exploration de l'espace \mathbb{N}^2 coûteux, besoin d'une stratégie.
- Sensibilité aux initialisations et à l'aléatoire.

Choix de (Q_1, Q_2) - Approche gloutonne



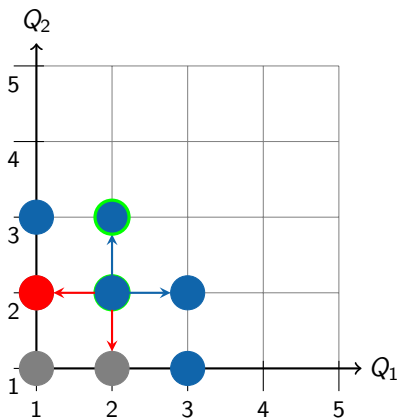
- Modèle initialisé :
●

Choix de (Q_1, Q_2) - Approche gloutonne



- Modèle initialisé :
●
- Modèle après *split* :
●
- Modèle maximisant le critère :
○

Choix de (Q_1, Q_2) - Approche gloutonne



- Modèle initialisé :
●
- Modèle après *split* :
●
- Modèle maximisant le critère :
○
- Modèle après *merge* :
●

Choix de (Q_1, Q_2) - Fenêtre glissante

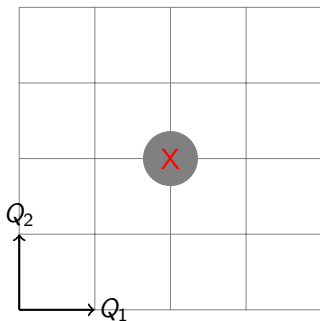


Figure 3 – Fenêtre glissante

Choix de (Q_1, Q_2) - Fenêtre glissante

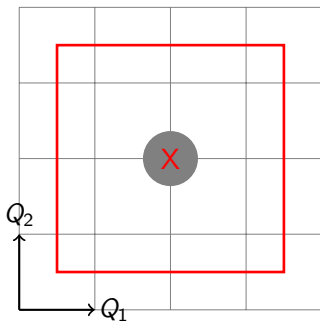


Figure 3 – Fenêtre glissante

Choix de (Q_1, Q_2) - Fenêtre glissante

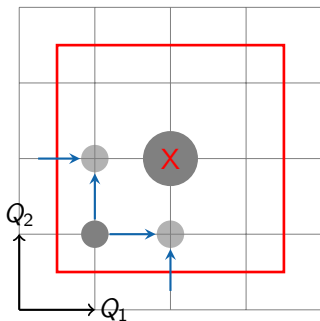


Figure 3 – Fenêtre glissante

Initialisation du modèle si nécessaire

Choix de (Q_1, Q_2) - Fenêtre glissante

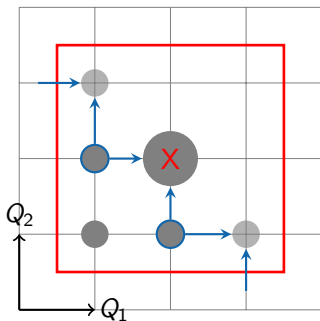


Figure 3 – Fenêtre glissante

Choix de (Q_1, Q_2) - Fenêtre glissante

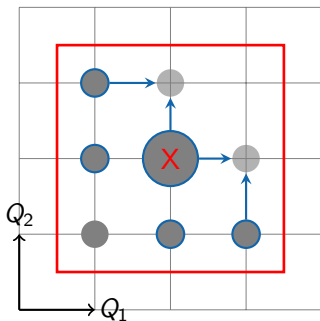


Figure 3 – Fenêtre glissante

Choix de (Q_1, Q_2) - Fenêtre glissante

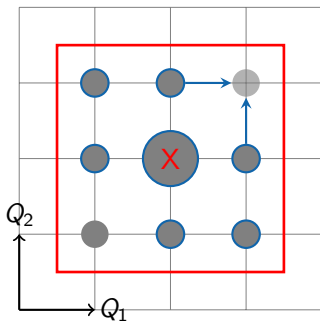


Figure 3 – Fenêtre glissante

Choix de (Q_1, Q_2) - Fenêtre glissante

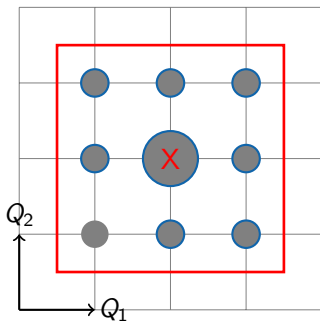


Figure 3 – Fenêtre glissante

Choix de (Q_1, Q_2) - Fenêtre glissante

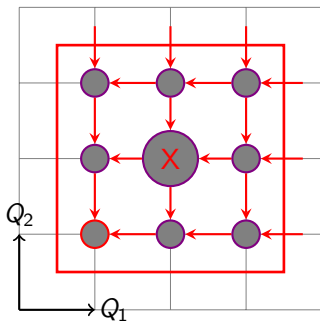


Figure 3 – Fenêtre glissante

Choix de (Q_1, Q_2) - Fenêtre glissante

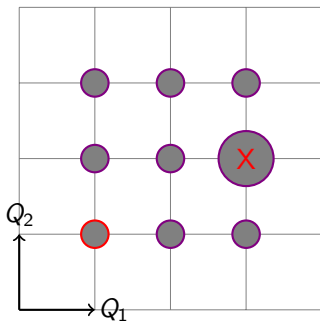


Figure 3 – Fenêtre glissante

Localisation du nouveau mode

Choix de (Q_1, Q_2) - Fenêtre glissante

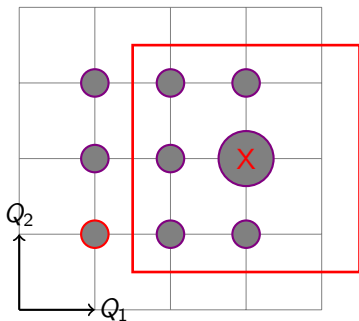
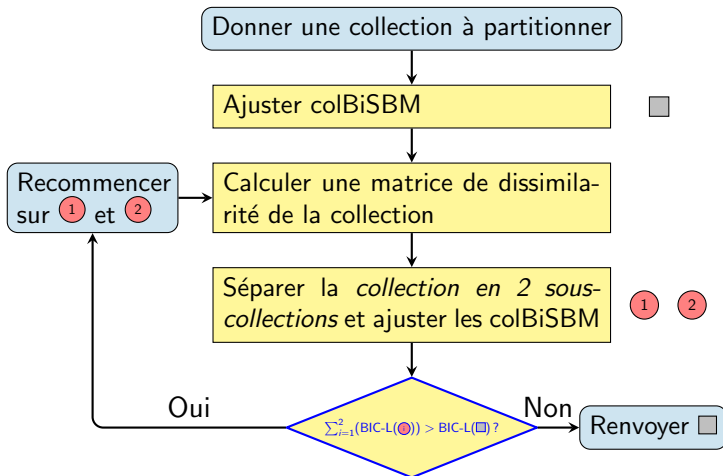


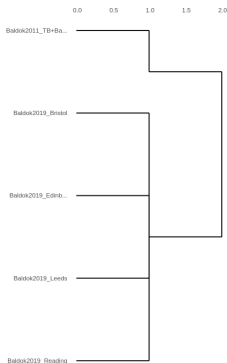
Figure 3 – Fenêtre glissante

Déplacement sur le nouveau mode puis itération

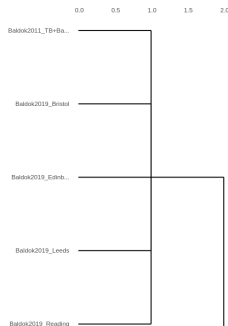
Clustering de réseaux



Application à Baldock et al., 2011, 2019 I



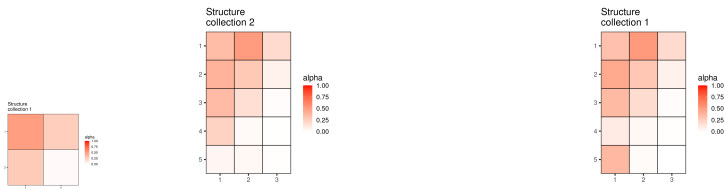
(a) Modèle iid ,
séparent réseau africain et réseaux anglais



(b) Modèle $\pi\rho$,
fusionnent réseaux africain et anglais

Figure 4 – Partitionnement des réseaux de Baldock et al., 2011, 2019

Application à Baldock et al., 2011, 2019 II



(a) Modèle *iid*,
séparet réseau africain et réseaux anglais

(b) Modèle $\pi\rho$,
fusionnent réseaux africain et anglais

Figure 5 – Structures détectées pour les réseaux de Baldock et al., 2011, 2019

Conclusion et perspectives

Capacités

- 4 modèles dont 3 qui ont une flexibilité sur au moins une des dimensions (adaptabilité aux données).
- Détecter structures classiques et moins classique de façon agnostique.
- Partitionner un ensemble de réseaux selon leurs structures.

Perspectives

- Investiguer stabilité face à l'aléatoire et aux *optima* locaux.
- Preuve d'identifiabilité du modèle $\pi\rho$.
-

Package et applications

- Intégration au package `co1SBM` et publication CRAN
- Intégrer possibilité d'un critère supplémentaire pour le clustering
- Appliquer clustering données de [Pichon et al., 2024](#) ; [Doré et al., 2021](#)
-

Merci pour votre attention !

Bibliographie I

- Govaert, G., & Nadif, M. (2005). An EM Algorithm for the Block Mixture Model. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 27(4), 643-647.
<https://doi.org/10.1109/TPAMI.2005.69>
- Chabert-Liddell, S.-C., Barbillon, P., & Donnet, S. (2024). Learning Common Structures in a Collection of Networks. An Application to Food Webs. *The Annals of Applied Statistics*, 18(2), 1213-1235.
<https://doi.org/10.1214/23-AOAS1831>
- Daudin, J.-J., Picard, F., & Robin, S. (2008). A mixture model for random graphs. *Stat Comput*, 18(2), 173-183.
<https://doi.org/10.1007/s11222-007-9046-7>
- Baldock, K. C. R., Memmott, J., Ruiz-Guajardo, J. C., Roze, D., & Stone, G. N. (2011). Daily temporal structure in African savanna flower visitation networks and consequences for network sampling. *Ecology*, 92(3), 687-698. <https://doi.org/10.1890/10-1110.1>

Bibliographie II

- Baldock, K. C. R., Goddard, M. A., Hicks, D. M., Kunin, W. E., Mitschunas, N., Morse, H., Osgathorpe, L. M., Potts, S. G., Robertson, K. M., Scott, A. V., Staniczenko, P. P. A., Stone, G. N., Vaughan, I. P., & Memmott, J. (2019). A systems approach reveals urban pollinator hotspots and conservation opportunities. *Nat Ecol Evol*, 3(3), 363-373. <https://doi.org/10.1038/s41559-018-0769-y>
- Pichon, B., Le Goff, R., Morlon, H., & Perez-Lamarque, B. (2024). Telling mutualistic and antagonistic ecological networks apart by learning their multiscale structure. *Methods in Ecology and Evolution*, 15(6), 1113-1128. <https://doi.org/10.1111/2041-210X.14328>
- Doré, M., Fontaine, C., & Thébault, E. (2021). Relative effects of anthropogenic pressures, climate, and sampling design on the structure of pollination networks at the global scale. *Global Change Biology*, 27(6), 1266-1280. <https://doi.org/10.1111/gcb.15474>

Bibliographie des annexes I