

Master 2 Statistique et M2A
Sorbonne Université
2019-2020

Inférence statistique de graphes

Anna Ben-Hamou
anna.ben-hamou@upmc.fr

Table des matières

CHAPITRE 1. Généralités sur les graphes	3
1. Premières définitions	3
2. Marches aléatoires	4
3. Conductances et inégalités de Cheeger	9
4. Clustering de réseaux	12
4.1. Le clustering spectral	12
4.2. Le clustering hiérarchique	13
4.3. L’intermédiation des arêtes	14
4.4. La modularité de Newman	15
CHAPITRE 2. Echantillonnage et estimation	16
1. Estimation de moyennes	16
2. Algorithme de Metropolis–Hastings	19
3. Estimation de la taille du graphe	21
3.1. Echantillons i.i.d.	21
3.2. Estimation par marches aléatoires	25
CHAPITRE 3. Détection de communautés	27
1. Modèle d’Ising et belief propagation	29
2. Reconstruction exacte	31
3. Reconstruction faible	36
4. Détection de communautés	40
4.1. Transmission d’information sur un arbre	40
4.2. Test de la présence de communautés	41
CHAPITRE 4. Graphes géométriques aléatoires : tests et estimation de la dimension	43
1. Un modèle simple de graphe géométrique aléatoire	43
2. Détection de la géométrie	44
2.1. Le test des triangles	45
2.2. Matrices de Wishart et matrices du GOE	46
CHAPITRE 5. Archéologie des réseaux	50
1. Modèles d’arbres croissants : attachement uniforme et préférentiel	50
2. Un algorithme simple pour trouver la racine	51
3. Urnes de Polya	51
4. Performances de l’algorithme	53

5. Bornes inférieures	55
CHAPITRE 6. Modèles graphiques	57
1. Propriétés de Markov et théorème d’Hammersley–Clifford	57
2. Modèles d’arbres	62
2.1. L’algorithme de Chow et Liu	62
2.2. L’algorithme de Kruskal	63
3. Reconstruction de graphes sparses	63
4. Modèles graphiques gaussiens	67
Outils probabilistes	69
Bibliographie	70

Généralités sur les graphes

1. Premières définitions

Accordons-nous d'abord sur ce l'on appellera un graphe dans ce cours.

Définition 1.1 (Graphe, sommets, arêtes). Un graphe G est un couple (V, E) formé d'un ensemble fini V , et d'un ensemble E de parties de V de cardinal 2. Les éléments de V sont appelés les sommets, ceux de E les arêtes.

Si $\{u, v\} \in E$, on dit que les sommets u et v sont voisins, ce que l'on note $u \sim v$. Le voisinage $\mathcal{N}(u)$ d'un sommet $u \in V$ est l'ensemble de ses voisins, i.e.

$$\mathcal{N}(u) = \{v \in V, v \sim u\},$$

et le degré de u , noté $\deg(u)$, est le nombre de ses voisins, i.e. $\deg(u) = |\mathcal{N}(u)|$.

Définition 1.2 (Multigraphe). Un multigraphe G est un couple (V, E) formé d'un ensemble fini V , et d'un multiensemble fini E de parties de V de cardinal 1 ou 2. Autrement dit, dans un multigraphe, deux sommets distincts peuvent être reliés par plusieurs arêtes, et un sommet peut être relié à lui-même par une ou plusieurs boucles.

Si $G = (V, E)$ est un multigraphe et si l'on note $m_{u,v}$ est la multiplicité de l'arête $\{u, v\}$ dans E , et m_u la multiplicité de $\{u\}$ dans E , alors le degré de u est défini par

$$\deg(u) = m_u + \sum_{v \in V} m_{u,v}.$$

Définition 1.3 (Graphe dirigé). Un graphe dirigé G est un couple (V, E) formé d'un ensemble fini V , et d'un ensemble E de paires ordonnées d'éléments distincts de V . Ainsi, l'arête dirigée (u, v) peut être présente sans que l'arête opposée (v, u) le soit.

Dans un graphe dirigé, on définit le degré sortant d'un sommet u par

$$\deg^+(u) = \sum_{v \in V} \mathbb{1}_{\{(u,v) \in E\}},$$

et son degré entrant par

$$\deg^-(u) = \sum_{v \in V} \mathbb{1}_{\{(v,u) \in E\}}.$$

Proposition 1.1. Soit $G = (V, E)$ un graphe. Alors

$$\sum_{u \in V} \deg(u) = 2|E|.$$

Définition 1.4 (Matrice d'adjacence). Soit $G = (V, E)$ un graphe avec $V = \{1, \dots, n\}$. La matrice d'adjacence de G est la matrice de taille $n \times n$ dont les entrées sont données par

$$A_{i,j} = \mathbb{1}_{\{\{i,j\} \in E\}}.$$

Définition 1.5 (Chemins). Soit $G = (V, E)$ un graphe. Un chemin est une suite de sommets $\mathbf{c} = (u_0, \dots, u_k)$ telle que pour tout $j \in \llbracket 0, k-1 \rrbracket$, les sommets u_j et u_{j+1} sont voisins. La longueur d'un chemin $\mathbf{c} = (u_0, \dots, u_k)$, notée $L(\mathbf{c})$, est donnée par l'entier k (le nombre d'arêtes sur ce chemin).

On dit que deux sommets u et v communiquent s'il existe un chemin allant de u à v . Une composante connexe est un sous-ensemble $A \subset V$ tel que pour tous $u, v \in A$, u et v communiquent et il n'existe pas de sommet dans A^c qui communique avec un sommet de A . Un graphe est dit connexe s'il n'y a qu'une composante connexe. Sur un graphe connexe, on peut définir une distance naturelle, donnée par

$$d(u, v) = \min \{L(\mathbf{c}), \mathbf{c} \in \mathcal{C}_{u,v}\},$$

où $\mathcal{C}_{u,v}$ correspond à l'ensemble des chemins de u à v .

Proposition 1.2. Soit A la matrice d'adjacence d'un graphe $G = (V, E)$ avec $V = \{1, \dots, n\}$. Pour tous $i, j \in V$ et pour tout $k \in \mathbb{N}$, on a

$$A^k(i, j) = |\{\mathbf{c} \in \mathcal{C}_{i,j}, L(\mathbf{c}) = k\}|.$$

Le diamètre d'un graphe connexe $G = (V, E)$ est la plus grande distance entre deux sommets :

$$\text{diam}(G) = \max_{u,v \in V} d(u, v) = \max_{u,v \in V} \min_{\mathbf{c} \in \mathcal{C}_{u,v}} L(\mathbf{c}).$$

2. Marches aléatoires

Soit $G = (V, E)$ un graphe. La marche aléatoire simple sur G est la chaîne de Markov d'espace d'états V et de matrice de transition P donnée par

$$\forall u, v \in V, P(u, v) = \begin{cases} \frac{1}{\deg(u)} & \text{si } \{u, v\} \in E, \\ 0 & \text{sinon.} \end{cases}$$

Autrement dit, $P = D^{-1}A$ où A est la matrice d'adjacence et D est la matrice diagonale dont les coefficients diagonaux correspondent aux degrés des sommets. Ainsi, à chaque temps, la marche aléatoire se déplace vers un voisin choisi uniformément au hasard.

Définition 1.6. On appelle probabilité stationnaire de P une probabilité π sur V qui vérifie $\pi P = \pi$.

Proposition 1.3. Soit P la matrice de transition de la marche aléatoire sur un graphe $G = (V, E)$. Alors il existe une probabilité stationnaire π donnée par

$$\forall u \in V, \pi(u) = \frac{\deg(u)}{2|E|}.$$

Si G est connexe, alors cette probabilité stationnaire est unique.

Remarque 1.1. Dire que G est connexe est équivalent à dire que la matrice P est irréductible :

$$\forall x, y \in V, \exists t \in \mathbb{N}, P^t(x, y).$$

Preuve de la Proposition 1.3. Montrons que le noyau P est réversible par rapport à π , i.e. que pour tous $x, y \in V$, $\pi(x)P(x, y) = \pi(y)P(y, x)$ (on dit que P vérifie les équations d'équilibre détaillé). On a

$$(1.1) \quad \pi(x)P(x, y) = \frac{\deg(x)}{2|E|} \cdot \frac{1}{\deg(x)} \mathbb{1}_{\{x, y\} \in E} = \frac{\mathbb{1}_{\{x, y\} \in E}}{2|E|} = \pi(y)P(y, x).$$

Cela implique que π est stationnaire pour P . En effet, pour tout $y \in V$,

$$\pi P(y) = \sum_{x \in V} \pi(x)P(x, y) = \sum_{x \in V} \pi(y)P(y, x) = \pi(y).$$

Montrons maintenant que si G est connexe, π est l'unique probabilité stationnaire. En remarquant qu'une mesure stationnaire est un élément du noyau de ${}^tP - I$, et que la dimension du noyau d'une matrice est égal à la dimension de sa transposée, il suffit de montrer que $\dim \text{Ker}(P - I) = 1$. Comme P est stochastique, on sait que le vecteur $\mathbf{1}$ est dans le noyau de $P - I$ et l'on va montrer que $\text{Ker}(P - I) = \text{Vect}(\mathbf{1})$. Soit $f \in \text{Ker}(P - I)$. Comme V est fini, f atteint son maximum sur V . Soit x_0 tel que $f(x_0) = \max_{x \in V} f(x) = M$, et soit $y \in V$. Comme le graphe est connexe, il existe un chemin $(x_0, x_1, \dots, x_k = y)$ allant de x_0 à y . Supposons qu'il existe $x \in \mathcal{N}(x_0)$ tel que $f(x) < M$. Alors

$$f(x_0) = \sum_{z \in V} P(x_0, z)f(z) = \frac{1}{\deg(x_0)} \sum_{z \in \mathcal{N}(x_0)} f(z) \leq \frac{1}{\deg(x_0)} (f(x) + (\deg(x_0) - 1)M) < M,$$

ce qui est contradictoire. Ainsi, pour tout $x \in \mathcal{N}(x_0)$, $f(x) = M$. En particulier $f(x_1) = M$. En répétant le même argument de proche en proche le long du chemin de x_0 à y , on obtient $f(y) = f(x_{k-1}) = \dots = f(x_1) = M$. Ainsi f est constante et on a bien $\text{Ker}(P - I) = \text{Vect}(\mathbf{1})$. ■

Définition 1.7. Le noyau de transition P est dit ergodique si

$$\exists t \in \mathbb{N}, \forall x, y \in V, P^t(x, y).$$

Proposition 1.4. Soit P la matrice de transition de la marche aléatoire sur un graphe $G = (V, E)$ avec $|V| = n$. Alors

- (1) La matrice P possède n valeurs propres réelles $\lambda_1 \geq \dots \geq \lambda_n$ et l'espace euclidien $\ell_{\mathbb{R}}^2(\pi)$ des fonctions $f : V \rightarrow \mathbb{R}$ muni du produit scalaire $\langle f, g \rangle_{\pi} = \sum_{x \in V} f(x)g(x)\pi(x)$ possède une base orthonormée de vecteurs propres $(\psi_j)_{j=1}^n$ associés aux valeurs propres $(\lambda_j)_{j=1}^n$.

- (2) On a

$$1 = \lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_n \geq -1.$$

- (3) P est irréductible (i.e. G est connexe) si et seulement si $\lambda_2 < 1$.

(4) P est ergodique si et seulement si $\max\{\lambda_2, -\lambda_n\} < 1$.

Preuve de la Proposition 1.4.

- (1) Le premier point vient du théorème spectral appliqué à l'opérateur P , auto-adjoint pour le produit scalaire $\langle \cdot, \cdot \rangle_\pi$. En effet, comme P est réversible par rapport à π (voir (1.1)), on a, pour toutes fonctions $f, g \in \ell_{\mathbb{R}}^2(\pi)$,

$$\begin{aligned} \langle f, Pg \rangle_\pi &= \sum_{x \in V} \pi(x) f(x) P g(x) \\ &= \sum_{x, y \in V} \pi(x) P(x, y) f(x) g(y) \\ &= \sum_{x, y \in V} \pi(y) P(y, x) f(x) g(y) \\ &= \langle Pf, g \rangle_\pi. \end{aligned}$$

- (2) Comme P est stochastique, le vecteur $\mathbf{1}$ est vecteur propre pour la valeur propre 1. Montrons maintenant que toutes les valeur propre sont comprises dans $[-1, 1]$. Soit $\lambda \in \text{Sp}(P)$ et soit f un vecteur propre associé à λ . Alors

$$|\lambda| \cdot \|f\|_\infty = \|\lambda f\|_\infty = \|Pf\|_\infty = \max_{x \in V} \left| \sum_{y \in V} P(x, y) f(y) \right| \leq \|f\|_\infty,$$

et ainsi $|\lambda| \leq 1$.

- (3) On a vu dans la preuve de la Proposition 1.3 que si G est connexe, alors $\text{Ker}(P - I) = \text{Vect}(\mathbf{1})$. La valeur propre 1 a donc multiplicité 1, et $\lambda_2 < 1$. Inversement, si le graphe n'est pas connexe, soit \mathcal{C} une des composantes connexes de G et soit $f = \mathbb{1}_{\{\cdot \in \mathcal{C}\}}$. Alors

$$Pf(x) = \sum_{y \in V} P(x, y) \mathbb{1}_{\{y \in \mathcal{C}\}} = \mathbf{P}_x(X_1 \in \mathcal{C}) = f(x).$$

Ainsi la valeur propre 1 a une multiplicité strictement plus grande que 1 et $\lambda_2 = 1$.

- (4) Si P est ergodique, on dispose de $t \in \mathbb{N}$ impair et de $\varepsilon > 0$ tel que pour tous $x \in V$, $P^t(x, x) \geq \varepsilon$. Soit Q le noyau de transition défini par

$$Q = \frac{P^t - \varepsilon I}{1 - \varepsilon}.$$

Notons que si $\lambda \in \text{Sp}(P)$, alors $\frac{\lambda^t - \varepsilon}{1 - \varepsilon} \in \text{Sp}(Q)$. Donc, par le point (2), $\left| \frac{\lambda^t - \varepsilon}{1 - \varepsilon} \right| \leq 1$, c'est-à-dire $\lambda^{2t} + 2\varepsilon(1 - \lambda^t) \leq 1$. Ainsi, comme t est impair, -1 ne peut pas être valeur propre. Inversement, si $\lambda_\star = \max\{\lambda_2, -\lambda_n\} < 1$, alors en décomposant la fonction $f = \delta_y$ dans la base orthonormée (ψ_j) de vecteurs propres, on a $f = \sum_{j=1}^n \pi(y) \psi_j(y) \psi_j$ et puisque $\psi_1 = \mathbf{1}$,

$$\frac{P^t(x, y)}{\pi(y)} = 1 + \sum_{j=2}^n \lambda_j^t \psi_j(x) \psi_j(y) \geq 1 - |\lambda_\star|^t \sum_{j=2}^n |\psi_j(x) \psi_j(y)|.$$

Ainsi, dès que t est tel que $|\lambda_\star|^t \max_{x,y \in V} \sum_{j=2}^n |\psi_j(x)\psi_j(y)| < 1$, on a $P^t(x,y) > 0$ pour tous $x, y \in V$. ■

La quantité $1 - \lambda_2$ s'appelle le trou spectral de P (ou indifféremment de G). Le trou spectral absolu correspond à $1 - \lambda_\star$ avec

$$\lambda_\star = \max \{|\lambda|, \lambda \in \text{Sp}(P), \lambda \neq 1\} = \max \{\lambda_2, -\lambda_n\}.$$

Comme nous le verrons ci-dessous, ces quantités spectrales sont intimement liées aux propriétés de mélange de la marche aléatoire et aux propriétés géométriques de G (présence de « goulots d'étranglement »).

Si le graphe G est connexe (P irréductible), on sait que la chaîne possède une unique probabilité stationnaire π , proportionnelle aux degrés des sommets. Quand peut-on garantir que la chaîne va effectivement *converger* vers cette mesure π ? Plus précisément, on souhaite savoir sous quelles conditions on a

$$(1.2) \quad \forall x, y \in V, P^t(x, y) \xrightarrow[t \rightarrow +\infty]{} \pi(y).$$

La connexité de G ne suffit pas à garantir une telle convergence. Par exemple, la marche aléatoire simple sur un cycle de taille n avec n pair ne converge pas vers la probabilité uniforme. En effet, si les sommets sont étiquetés de 0 à $n - 1$, et si la marche part de $x = 0$, alors elle sera toujours sur un sommet pair aux temps pairs, et sur un sommet impair aux temps impairs. En fait, dès que le graphe est biparti (les sommets peuvent être partitionnés en deux sous-ensembles tels que toute arête a une extrémité dans chaque sous-ensemble), alors la marche simple ne peut pas converger. Pour garantir (1.2), on a besoin de plus que l'irréductibilité : l'ergodicité.

Proposition 1.5. *Si P est ergodique, alors*

$$\forall x, y \in V, P^t(x, y) \xrightarrow[t \rightarrow +\infty]{} \pi(y).$$

Remarque 1.2. Une façon simple de transformer un noyau irréductible P en un noyau ergodique est de rendre la marche *paresseuse* en considérant le noyau

$$\tilde{P} = \frac{P + I}{2}.$$

Autrement dit, on considère la marche qui à chaque temps reste sur place avec probabilité $1/2$ et effectue une transition selon P avec probabilité $1/2$.

Nous allons en fait montrer un résultat plus fort que la Proposition 1.5. Pour $x \in V$, notons $\mathcal{D}_x(t)$ la distance en variation totale entre la loi de la chaîne au temps t partie de x et la loi stationnaire, i.e.

$$\mathcal{D}_x(t) = \|P^t(x, \cdot) - \pi\|_{\text{TV}} = \max_{A \subset V} |P^t(x, A) - \pi(A)| = \sum_{y \in \Omega} (P^t(x, y) - \pi(y))_+,$$

et

$$\mathcal{D}(t) = \max_{x \in V} \mathcal{D}_x(t).$$

Théorème 1.6. *On a*

$$\frac{\lambda_\star^t}{2} \leq \mathcal{D}(t) \leq \frac{\lambda_\star^t}{2\sqrt{\pi_{\min}}},$$

où $\pi_{\min} = \min_{x \in V} \pi(x)$. En particulier, $\mathcal{D}(t)^{1/t} \xrightarrow{t \rightarrow \infty} \lambda_\star$, et, si P est ergodique, comme alors $\lambda_\star < 1$, on a $\mathcal{D}(t) \xrightarrow{t \rightarrow \infty} 0$.

Preuve du Théorème 1.6. Pour la première inégalité, soit φ un vecteur propre associé à une valeur propre $\lambda \neq 1$ de P . Comme les vecteurs propres sont orthogonaux pour le produit scalaire $\langle \cdot, \cdot \rangle_\pi$ et que $\mathbf{1}$ est vecteur propre pour la valeur propre 1, on a $\langle \varphi, \mathbf{1} \rangle_\pi = 0$. Ainsi

$$|\lambda^t \varphi(x)| = \left| \sum_{y \in V} (P^t(x, y) - \pi(y)) \varphi(y) \right| \leq 2\|\varphi\|_\infty \mathcal{D}(t).$$

Pour $x \in V$ tel que $\varphi(x) = \|\varphi\|_\infty$, on a donc $|\lambda|^t \leq 2\mathcal{D}(t)$. Et comme cela est vrai pour toute valeur propre différente de 1, $\lambda_\star^t \leq 2\mathcal{D}(t)$. Pour la deuxième inégalité, en décomposant la fonction $f = \delta_x$ dans la base orthonormée (ψ_j) de vecteurs propres, on a $f = \sum_{j=1}^n \pi(y) \psi_j(y) \psi_j(x)$ et puisque $\psi_1 = \mathbf{1}$,

$$P^t(x, y) - \pi(y) = P^t f(x) - \pi(y) = \sum_{j=2}^n \lambda_j^t \pi(y) \psi_j(y) \psi_j(x).$$

Par l'inégalité de Cauchy-Schwarz,

$$2\mathcal{D}_x(t) = \sum_{y \in V} |P^t(x, y) - \pi(y)| \leq \sqrt{\sum_{y \in V} \pi(y) \left(\sum_{j=2}^n \lambda_j^t \psi_j(y) \psi_j(x) \right)^2}.$$

Or

$$\sum_{y \in V} \pi(y) \left(\sum_{j=2}^n \lambda_j^t \psi_j(y) \psi_j(x) \right)^2 = \left\langle \sum_{j=2}^n \lambda_j^t \psi_j(x) \psi_j, \sum_{j=2}^n \lambda_j^t \psi_j(x) \psi_j \right\rangle_\pi = \sum_{j=2}^n \lambda_j^{2t} \psi_j(x)^2.$$

Et comme pour tout $j \geq 2$, $|\lambda_j| \leq \lambda_\star$, on a

$$2\mathcal{D}_x(t) \leq \lambda_\star^t \sqrt{\sum_{j=1}^n \psi_j(x)^2} = \frac{\lambda_\star^t}{\sqrt{\pi(x)}}.$$

En prenant le maximum sur $x \in V$, on obtient l'inégalité voulue. ■

Ainsi en particulier, sur tout graphe connexe, la marche aléatoire paresseuse converge vers sa probabilité stationnaire π . En pratique, on ne peut pas laisser la marche évoluer pendant un temps infini. Il est alors important non seulement de savoir que la chaîne va converger,

mais de savoir à partir de quel moment la loi de la marche sera « assez proche » de la loi stationnaire. Pour $\varepsilon \in]0, 1[$, on définit le temps de mélange (en variation totale) par

$$t_{\text{MIX}}(\varepsilon) = \min \{t \in \mathbb{N}, \mathcal{D}(t) \leq \varepsilon\} .$$

On sait alors que si $t \geq t_{\text{MIX}}(\varepsilon)$, pour tout sous-ensemble $A \subset V$, et quelque soit le point de départ $x \in V$,

$$\pi(A) - \varepsilon \leq \mathbf{P}_x(X_t \in A) \leq \pi(A) + \varepsilon .$$

Remarquons que si l'on s'intéresse à l'ordre de grandeur du temps de mélange (typiquement en fonction d'un paramètre n quantifiant la taille du graphe), alors la cible ε importe peu (pourvu qu'elle soit inférieure à $1/2$). Plus précisément, pour tout $k \geq 1$

$$(1.3) \quad t_{\text{MIX}} \left((2\varepsilon)^k \right) \leq k t_{\text{MIX}}(\varepsilon) .$$

En effet, introduisons

$$\overline{\mathcal{D}}(t) = \max_{x, y \in \Omega} \|P^t(x, \cdot) - P^t(y, \cdot)\|_{\text{TV}} .$$

Alors il est facile de voir que $\mathcal{D}(t) \leq \overline{\mathcal{D}}(t) \leq 2\mathcal{D}(t)$. La fonction $\overline{\mathcal{D}}$ est sous-multiplicative : $\overline{\mathcal{D}}(t+s) \leq \overline{\mathcal{D}}(t)\overline{\mathcal{D}}(s)$. En effet, soit $A \subset V$ et soit $B = \{z \in V, P^t(x, z) \geq P^t(y, z)\}$. En décomposant selon que la chaîne est en B ou en B^c au temps t , on a

$$\begin{aligned} P^{t+s}(x, A) - P^{t+s}(y, A) &= \sum_{z \in B} (P^t(x, z) - P^t(y, z)) P^s(z, A) - \sum_{z \in B^c} (P^t(y, z) - P^t(x, z)) P^s(z, A) \\ &\leq \|P^t(x, \cdot) - P^t(y, \cdot)\|_{\text{TV}} \max_{u, v \in V} (P^s(u, A) - P^s(v, A)) , \end{aligned}$$

où l'on a utilisé le fait que $\|P^t(x, \cdot) - P^t(y, \cdot)\|_{\text{TV}} = P^t(x, B) - P^t(y, B)$. En prenant le maximum sur $A \subset V$, on obtient

$$\|P^{t+s}(x, \cdot) - P^{t+s}(y, \cdot)\|_{\text{TV}} \leq \|P^t(x, \cdot) - P^t(y, \cdot)\|_{\text{TV}} \overline{\mathcal{D}}(s) ,$$

et en prenant le maximum sur $x, y \in V$, on a bien $\overline{\mathcal{D}}(t+s) \leq \overline{\mathcal{D}}(t)\overline{\mathcal{D}}(s)$. Ainsi

$$\mathcal{D}(k t_{\text{MIX}}(\varepsilon)) \leq \overline{\mathcal{D}}(k t_{\text{MIX}}(\varepsilon)) \leq \overline{\mathcal{D}}(t_{\text{MIX}}(\varepsilon))^k \leq (2\mathcal{D}(t_{\text{MIX}}(\varepsilon)))^k \leq (2\varepsilon)^k .$$

On peut définir d'autres temps de mélange selon la distance choisie pour mesurer l'écart entre la loi de la chaîne et la loi stationnaire. Par exemple, le temps de mélange *uniforme* est défini par

$$t_{\text{UNIF}}(\varepsilon) = \min \left\{ t \in \mathbb{N}, \max_{x, y \in V} \left| \frac{P^t(x, y)}{\pi(y)} - 1 \right| \leq \varepsilon \right\} .$$

3. Conductances et inégalités de Cheeger

Soit $S \subset V$ un sous-ensemble non-vide de V . La frontière de S , notée ∂S , est l'ensemble des arêtes dont une extrémité est dans S et l'autre dans S^c :

$$\partial S = \{\{u, v\} \in E, u \in S, v \notin S\} .$$

Le volume de S , noté $\text{Vol}(S)$, est défini comme la somme des degrés des sommets de S :

$$\text{Vol}(S) = \sum_{u \in S} \text{deg}(u) .$$

On définit alors la conductance de S comme le rapport

$$\Phi(S) = \frac{|\partial S|}{\text{Vol}(S)}.$$

La conductance de S peut aussi se voir comme la probabilité qu'une marche aléatoire initiée selon la mesure stationnaire conditionnée à S sorte de S en un pas. En effet,

$$\sum_{x \in S} \frac{\pi(x)}{\pi(S)} \sum_{y \notin S} P(x, y) = \sum_{x \in S} \frac{\deg(x)}{\text{Vol}(S)} \sum_{y \notin S} \frac{1}{\deg(x)} \mathbb{1}_{\{\{x, y\} \in E\}} = \frac{|\partial S|}{\text{Vol}(S)}.$$

La conductance minimale de G , appelée aussi constante de Cheeger, est donnée par

$$\Phi_* = \min \left\{ \Phi(S), S \subset V, S \neq \emptyset, \pi(S) \leq \frac{1}{2} \right\}.$$

Clairement, G est connexe si et seulement si $\Phi_* > 0$. Par la proposition 1.4, on a donc

$$1 - \lambda_2 > 0 \Leftrightarrow \Phi_* > 0 \Leftrightarrow G \text{ connexe}.$$

En fait, ce lien entre trou spectral et conductance minimale peut être quantifié de façon bien plus précise. Ce sont les fameuses inégalités de Cheeger.

Proposition 1.7. *Si λ_2 est la deuxième plus grande valeur propre de P , la matrice de transition sur $G = (V, E)$, et si Φ_* est la conductance minimale de G , alors*

$$\frac{\Phi_*^2}{2} \leq 1 - \lambda_2 \leq 2\Phi_*.$$

Avant de prouver ces inégalités, rappelons que la variance sous π d'une fonction $f : V \rightarrow \mathbb{R}$ peut s'écrire

$$\text{Var}_\pi(f) = \frac{1}{2} \sum_{x, y \in V} \pi(x)\pi(y) (f(y) - f(x))^2.$$

La forme de Dirichlet de f est définie par

$$\mathcal{E}(f) = \frac{1}{2} \sum_{x, y \in V} \pi(x)P(x, y) (f(y) - f(x))^2 = \langle f, (I - P)f \rangle_\pi.$$

Alors que $\text{Var}_\pi(f)$ mesure les fluctuations globales de f sur V , $\mathcal{E}(f)$ mesure les fluctuations locales de f le long des arêtes du graphe. On a

$$(1.1) \quad 1 - \lambda_2 = \min_{f \text{ non-constante}} \frac{\mathcal{E}(f)}{\text{Var}_\pi(f)}.$$

En effet, comme $\mathcal{E}(f)$ et $\text{Var}_\pi(f)$ sont toutes les deux invariante par translations, il suffit de minimiser sur les fonctions centrées (i.e. $\langle f, \mathbf{1} \rangle_\pi = 0$) et non identiquement nulles. L'égalité (1.1) équivaut alors à

$$\lambda_2 = \max_{f \perp \mathbf{1}, f \neq \mathbf{0}} \frac{\langle f, Pf \rangle_\pi}{\|f\|_\pi^2}.$$

Or, pour toute fonction $f : V \rightarrow \mathbb{R}$ avec $\langle f, \mathbf{1} \rangle_\pi = 0$ et $f \neq \mathbf{0}$, l'écriture de f dans la base orthonormée de fonctions propres (ψ_j) donne

$$f = \sum_{j=2}^n \langle f, \psi_j \rangle_\pi \psi_j.$$

Ainsi

$$\langle f, Pf \rangle_\pi = \sum_{j=2}^n \lambda_j \langle f, \psi_j \rangle_\pi^2 \leq \lambda_2 \sum_{j=2}^n \langle f, \psi_j \rangle_\pi^2 = \lambda_2 \|f\|_\pi^2,$$

avec égalité pour $f = \psi_2$.

Preuve de la Proposition 1.7. Pour l'inégalité $1 - \lambda_2 \leq 2\Phi_\star$, soit $f = \mathbb{1}_S$ avec $\emptyset \neq S \subset V$ et $\pi(S) \leq \frac{1}{2}$. On a $\text{Var}_\pi(f) = \pi(S)(1 - \pi(S))$ et

$$\mathcal{E}(f) = \sum_{x \in S} \sum_{y \notin S} \pi(x)P(x, y) = \pi(S)\Phi(S).$$

Donc, en utilisant la caractérisation (1.1) et en ne minimisant que sur les fonctions f de la forme $f = \mathbb{1}_S$ avec $S \neq \emptyset$ et $\pi(S) \leq \frac{1}{2}$,

$$1 - \lambda_2 \leq \min_{S \subset V, S \neq \emptyset, \pi(S) \leq 1/2} \frac{\Phi(S)}{1 - \pi(S)} \leq 2\Phi_\star.$$

Pour l'inégalité $\frac{\Phi_\star^2}{2} \leq 1 - \lambda_2$, soit f_2 un vecteur propre de P pour la valeur propre λ_2 tel que $\pi(f_2 > 0) \leq \frac{1}{2}$ (cela est toujours possible puisque si cela n'est pas vérifié pour un vecteur propre f_2 , ça l'est pour $-f_2$), et posons $f = \max\{f_2, 0\}$. Pour $t > 0$, soit $S_t = \{f^2 > t\}$. Comme $\pi(S_t) \leq 1/2$, on a

$$\pi(f^2 > t)\Phi_\star \leq \sum_{x \in S_t, y \notin S_t} \pi(x)P(x, y) = \sum_{x, y} \pi(x)P(x, y) \mathbb{1}_{f^2(y) \leq t < f^2(x)}.$$

En intégrant par rapport à t , on obtient

$$\|f\|_\pi^2 \Phi_\star \leq \frac{1}{2} \sum_{x, y} \pi(x)P(x, y) |f^2(x) - f^2(y)|,$$

et en utilisant l'inégalité de Cauchy-Schwarz,

$$\begin{aligned} \|f\|_\pi^4 \Phi_\star^2 &\leq \frac{1}{4} \left(\sum_{x, y} \pi(x)P(x, y)(f(x) + f(y))^2 \right) \left(\sum_{x, y} \pi(x)P(x, y)(f(x) - f(y))^2 \right) \\ &= \|f\|_\pi^4 - \langle f, Pf \rangle_\pi^2. \end{aligned}$$

Ainsi,

$$\frac{\langle f, Pf \rangle_\pi}{\|f\|_\pi^2} \leq \sqrt{1 - \Phi_\star^2} \leq \left(1 - \frac{\Phi_\star^2}{2}\right).$$

Il ne reste plus qu'à vérifier que $\langle f, Pf \rangle_\pi \geq \lambda_2 \|f\|_\pi^2$. Cela est vrai car $Pf \geq \lambda_2 f$. En effet, on a d'une part $Pf \geq 0$ car $f \geq 0$, et d'autre part, $Pf \geq Pf_2 = \lambda_2 f_2$. ■

4. Clustering de réseaux

Un graphe connexe $G = (V, E)$ étant donné, peut-on trouver la partition (S, S^c) de V qui atteint la conductance minimale Φ_* ? Plus généralement, peut-on partitionner les sommets en plusieurs sous-ensembles, appelés communautés, de telle sorte que les sommets d'un même groupe soient « mieux » connectés que les sommets de groupes différents? Ce problème (dont nous étudierons une version probabiliste au Chapitre 3) est connu sous le nom de *clustering* et est un des problèmes centraux en analyse statistique des graphes. De façon informelle, il peut s'énoncer de la façon suivante : étant donné G , comment regrouper les sommets en sous-groupes distincts de telle sorte que la densité d'arêtes soit plus élevée au sein des groupes qu'entre les groupes? Déjà, définir une notion de partition optimale pour le clustering ne va pas de soi, et de nombreux critères ont été proposés pour juger de la qualité d'une partition donnée, donnant lieu à autant d'algorithmes de clustering. Ici, nous présenterons d'abord le clustering spectral et le clustering hiérarchique. Puis nous étudierons un algorithme proposé par Girvan and Newman [12] reposant sur la notion d'intermédiarité d'une arête (*edge betweenness*). Enfin, nous introduirons la modularité de Newman [22], et présenterons un algorithme élaboré par Clauset et al. [8] pour optimiser cette modularité.

4.1. Le clustering spectral. Le principe du clustering spectral est simple : si le graphe est composé de communautés, alors il devrait exister une permutation des sommets qui rende la matrice d'adjacence quasi-diagonale par blocs. Pour des raisons de robustesse, le clustering spectral ne cherche pas à diagonaliser directement A la matrice d'adjacence de G , mais plutôt une version normalisée de celle-ci appelée matrice laplacienne. Si D correspond à la matrice diagonale dont les coefficients diagonaux $(d_i)_{i=1}^n$ correspondent aux degrés des sommets, alors on définit la matrice laplacienne de G comme

$$L = D - A.$$

On utilise souvent la matrice laplacienne normalisée définie par

$$L_{\text{norm}} = I - D^{-1/2} A D^{-1/2}.$$

On a

$$L_{\text{norm}}(i, j) = \begin{cases} 1 & \text{si } i = j, \\ -\frac{A_{i,j}}{\sqrt{d_i d_j}} & \text{si } i \neq j. \end{cases}$$

La matrice L_{norm} est symétrique. De plus, pour tout $u \in \mathbb{R}^n$, on peut facilement voir que

$${}^t u L_{\text{norm}} u = \frac{1}{2} \sum_{1 \leq i, j \leq n} A_{i,j} \left(\frac{u_i}{\sqrt{d_i}} - \frac{u_j}{\sqrt{d_j}} \right)^2.$$

En particulier, la matrice L_{norm} est semi-définie positive, et toutes ses valeurs sont positives ou nulles. On les ordonne par ordre croissant : $\mu_1 \leq \dots \leq \mu_n$. On peut aussi remarquer que 0 est valeur propre pour le vecteur propre $D^{1/2} \mathbf{1}$. Ainsi $\mu_1 = 0$. En fait, la multiplicité de la valeur propre 0 est égale au nombre de composantes connexes de G , et si $C_1, \dots, C_k \subset V$ sont les k composantes connexes de G , alors l'espace propre associé à la valeur propre 0 est engendré par $D^{1/2} \mathbf{1}_{C_1}, \dots, D^{1/2} \mathbf{1}_{C_k}$, où $\mathbf{1}_{C_j}(u) = \mathbb{1}_{u \in C_j}$ (exercice). Ainsi, si le graphe a k composantes

connexes, celles-ci se lisent exactement sur les vecteurs propres associés aux k plus petites valeurs propres de L_{norm} (qui sont alors égales à 0). L'idée du clustering spectral est que, même lorsque le graphe est connexe, les k premiers vecteurs propres (i.e. ceux associés aux k plus petites valeurs propres) devraient toujours nous renseigner sur les communautés.

On présente ici un algorithme de clustering spectral dû à Ng et al. [23]. À partir du graphe G , on forme la matrice laplacienne normalisée L_{norm} et l'on détermine ses k premiers vecteurs propres, notés $u_1, \dots, u_k \in \mathbb{R}^n$. Soit $U \in \mathcal{M}_{n,k}(\mathbb{R})$ la matrice dont la $j^{\text{ième}}$ colonne est donnée par u_j , et soit T la matrice formée à partir de U en renormalisant les entrées pour que toutes les lignes soient de norme 1 :

$$T_{i,j} = \frac{U_{i,j}}{\sqrt{\sum_{\ell=1}^k U_{i,\ell}^2}}.$$

Notons $t_i \in \mathbb{R}^k$ le vecteur unitaire donné par la $i^{\text{ième}}$ ligne de T . Pour partitionner les points $(t_i)_{i=1}^n$ en k blocs, on utilise l'algorithme des k -means : on part d'une partition $\mathbf{S}^{(0)} = (S_1^{(0)}, \dots, S_k^{(0)})$ de $\{1, \dots, n\}$ arbitraire, et à chaque temps $t \in \mathbb{N}$, si la partition $\mathbf{S}^{(t)}$ est formée, on calcule les barycentres des blocs :

$$\forall j \in \llbracket 1, k \rrbracket, \quad m_j^{(t)} = \frac{1}{|S_j^{(t)}|} \sum_{i \in S_j^{(t)}} t_i.$$

Puis on forme la partition $\mathbf{S}^{(t+1)}$ en affectant chaque point au bloc dont l'indice est celui du barycentre dont il est le plus proche, i.e.

$$S_j^{(t+1)} = \left\{ i \in [n], \forall \ell \in \llbracket 1, k \rrbracket, \|t_i - m_j^{(t)}\| \leq \|t_i - m_\ell^{(t)}\| \right\}.$$

Ng et al. [23] montrent que sous certaines conditions, pour t assez grand, la partition $\mathbf{S}^{(t)}$ permet de bien retrouver la structure de communautés de G .

On pourra se référer à Von Luxburg [24] pour un tutoriel très clair et très complet sur le clustering spectral.

4.2. Le clustering hiérarchique. Les algorithmes de clustering hiérarchiques fonctionnent de la façon suivante : on commence par attribuer un poids $\omega_{i,j}$ à toutes les paires de noeuds distincts $\{i, j\}$ (pas seulement aux arêtes), correspondant à une certaine notion de proximité entre i et j . Typiquement, on prend souvent $\omega_{i,j}$ égal au nombre minimal d'arêtes qu'il faut supprimer pour disconnecter i et j (i.e. pour que i et j se retrouvent dans deux composantes connexes distinctes). Déterminer ces poids peut se faire en temps polynomial, par exemple en utilisant le Théorème *Max-Flow Min-Cut*, et en utilisant un algorithme pour déterminer le max-flow, comme l'algorithme de *push-relabel* dont la complexité est en $O(|V|^2|E|)$ (pour une paire donnée, ce qui permet de calculer tous les poids en un temps $O(|V|^4|E|)$). Une fois ces poids déterminés, on ordonne les paires de sommets par ordre décroissant des poids. Le clustering hiérarchique s'exécute alors par fusions successives : à chaque étape, on fusionne les deux noeuds qui ont le plus grand poids parmi les paires restantes. Ainsi, partant d'un état où tous les noeuds sont séparés (n communautés formées de singletons), on arrive progressivement à l'état où tous les noeuds sont dans la même communauté. Au

cours de l’algorithme, les composantes sont proprement emboîtées et l’on peut ainsi représenter l’évolution par un arbre appelé *dendrogramme*, dans lequel une coupe à un niveau donné fournit une structure de communautés à une certaine précision (voir Figure 1).

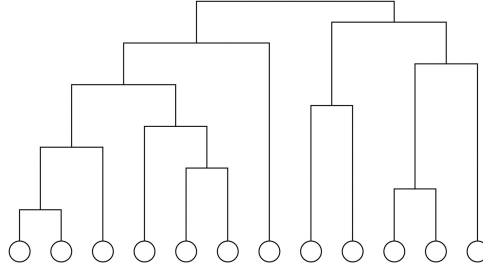


FIGURE 1. Un dendrogramme

Cet algorithme présente plusieurs pathologies. Notamment, il a tendance à séparer les sommets périphériques de la communauté à laquelle on les attribuerait naturellement. Par exemple, si un sommet u n’est attaché que par une arête au reste du réseau, on a envie de le classer dans le même groupe que son unique voisin v . Or ce n’est souvent pas ce qui se produit : la coupe minimale entre u et n’importe quel autre sommet w est égale à 1, et le sommet v n’est donc pas particulièrement privilégié.

4.3. L’intermédiarité des arêtes. Girvan and Newman [12] ont proposé une autre façon de partitionner les sommets, fondée sur la notion d’intermédiarité des arêtes : au lieu de s’intéresser en premier lieu aux arêtes susceptibles d’être intra-blocs (reliant deux sommets appartenant à la même communauté), leur idée est de s’intéresser plutôt aux arêtes susceptibles d’être des arêtes passerelles, reliant deux sommets de communautés différentes. L’algorithme ne procédera pas par fusions successives comme dans le clustering hiérarchique, mais au contraire par divisions successives.

L’intermédiarité d’une arête $e \in E$, notée $c(e)$, est définie comme le nombre de plus courts chemins entre des paires de sommets qui passent par e , i.e.

$$c(e) = \sum_{\{u,v\}} \sum_{\mathbf{c} \in \mathcal{C}_{u,v}} \mathbb{1}_{L(\mathbf{c})=d(u,v)} \mathbb{1}_{e \in \mathbf{c}}.$$

(On peut aussi définir la centralité d’un sommet $u \in V$ de façon analogue comme le nombre de plus courts chemins entre des paires de sommets qui passent par u .)

Une fois que l’on a classé les arêtes du graphe par ordre décroissant d’intermédiarité, la première étape de l’algorithme consiste à supprimer l’arête de plus grande intermédiarité. Puis, on recalcule les intermédiarités des arêtes restantes et on supprime celle de plus grande intermédiarité, et ainsi de suite jusqu’à ce qu’on arrive au graphe vide (sans arête). Chaque étape de l’algorithme fournit une structure de communautés d’une certaine précision, où les communautés correspondent aux composantes connexes.

Dans un graphe à n sommets et m arêtes, calculer l’intermédiarité de toutes les arêtes prend un temps $O(nm)$. Comme il faut recalculer les intermédiarités à chacune des $|E|$ étapes de l’algorithme, la complexité est en $O(|E|^2|V|)$.

4.4. La modularité de Newman. La notion de modularité a été introduite par Newman [22]. L'idée est qu'une bonne partition des sommets n'est finalement pas tant que ça une partition où il y a peu d'arêtes inter-blocs, mais plutôt une partition dans laquelle, entre deux blocs, il y a moins d'arêtes que ce à quoi on s'attendrait si les arêtes étaient placées au hasard. Contrairement aux deux méthodes précédentes, cette notion prend en compte la possibilité que le graphe ne possède pas de structure de communautés. Il ne s'agit pas seulement de proposer des partitions plus ou moins raffinées, mais de savoir si le graphe possède une partition pour laquelle les densités d'arêtes intra-blocs et inter-blocs sont significativement différentes des valeurs auxquelles on s'attendrait dans un graphe aléatoire.

Soit $\mathcal{P} = (V_1, \dots, V_K)$ une partition de V en K blocs. Pour $k \in \llbracket 1, K \rrbracket$, notons e_k le proportion des arêtes de G qui ont les deux extrémités dans V_k :

$$e_k = \frac{\sum_{\{u,v\} \in E} \mathbb{1}_{\{u,v \in V_k\}}}{m}.$$

On souhaite comparer e_k avec une valeur « typique ». Pour cela, on considère un modèle de graphe aléatoire, appelé *modèle de configuration*, qui permet de générer des graphes (ou plutôt des multi-graphes) dont la suite de degrés est prescrite au préalable. Initialement, on attache à chaque sommet $v \in V$ un nombre de *demi-arêtes* égal à $\deg(v)$. Puis on génère le graphe en choisissant, uniformément au hasard, un appariement sur l'ensemble des demi-arêtes. On interprète alors une paire de demi-arêtes comme une arête entre les sommets correspondants. Dans ce modèle, si l'on tire une arête uniformément au hasard, la probabilité qu'elle ait les deux extrémités dans V_k est

$$\frac{\sum_{u \in V_k} \deg(u)}{2m} \frac{\sum_{v \in V_k} \deg(v) - 1}{2m - 1} \approx \left(\frac{\sum_{u \in V_k} \deg(u)}{2m} \right)^2 = \pi(V_k)^2.$$

Cela conduit à la définition suivante pour la modularité de la partition \mathcal{P} :

$$\text{mod}(\mathcal{P}) = \sum_{k=1}^K \{e_k - \pi(V_k)^2\}.$$

Notons que $\text{mod}(\mathcal{P})$ peut être positive ou négative, avec $\text{mod}(\mathcal{P}) > 0$ indiquant une possible structure de communautés. Newman propose donc de chercher une partition avec une grande modularité, idéalement celle qui atteint la modularité maximale :

$$\text{mod}_G = \max_{\mathcal{P}} \text{mod}(\mathcal{P}),$$

où le maximum est pris sur toutes les partitions possibles de V . Malheureusement, chercher la partition optimale pour ce problème est extrêmement coûteux (le nombre de partitions d'un ensemble à n éléments est donné par le nombre de Bell B_n qui croît plus qu'exponentiellement vite en n) et il est plus raisonnable de se contenter de solutions approchées. Par exemple, Newman [21] propose un algorithme de type *greedy* où, partant de la partition où chaque sommet est un singleton, on fusionne successivement les deux blocs qui induisent la plus grande augmentation (ou la plus petite diminution) dans la modularité. On peut montrer que chaque étape peut être exécutée en un temps $O(m)$, ce qui conduit à une complexité en $O(mn)$.

Echantillonnage et estimation

Dans ce chapitre, on considère un graphe $G = (V, E)$ fixé mais initialement inconnu. Comment échantillonner sur G pour obtenir des informations? Si l'on suppose que l'on sait échantillonner selon une mesure relativement représentative, par exemple la mesure uniforme sur V ou la mesure π proportionnelle aux degrés, alors on peut appliquer des outils de statistiques classiques. Mais en général, le problème même de l'échantillonnage sur G pose problème. Des réseaux comme celui du web par exemple sont tellement grands (le web contient plus de 1,6 milliards de sites) qu'il n'est pas du tout garanti que l'on puisse tirer un nœud au hasard. En pratique, on a recours à des méthodes de type MCMC (*Markov Chains Monte Carlo*) : partant d'un sommet fixé, on lance une marche aléatoire sur le graphe, et on la laisse évoluer jusqu'à ce qu'elle mélange. Au bout d'un temps suffisant, la distribution de la marche est alors très proche de sa distribution stationnaire. Ainsi, si le temps de mélange n'est pas trop grand, on dispose d'un moyen simple et peu coûteux pour échantillonner un sommet selon π . Nous verrons que si l'on souhaite échantillonner selon une autre loi que π , il existe des moyens simples de modifier la marche pour qu'elle ait la loi stationnaire voulue. Une fois que l'on a un moyen d'échantillonner sur G se pose la question de l'inférence : que peut-on apprendre du graphe G à partir de la trajectoire d'une marche aléatoire? Et surtout, à partir de combien de temps peut-on obtenir une bonne estimation de certains paramètres du graphe?

1. Estimation de moyennes

Soit $G = (V, E)$ un graphe avec $|V| = n$ et $|E| = m$, et soit π la probabilité proportionnelle aux degrés : $\pi(u) = \frac{\deg(u)}{2m}$. Soit $f : V \rightarrow \mathbb{R}$ une fonction définie sur les sommets. On souhaite estimer la moyenne de f contre π :

$$\mathbf{E}_\pi f = \sum_{v \in V} \pi(v) f(v),$$

à partir de la trajectoire d'une marche aléatoire $(X_t)_{t \geq 0}$ sur G , partie d'un sommet arbitraire $x \in V$. Un estimateur naturel consiste à prendre la moyenne de f le long de la trajectoire :

$$\hat{f}_t = \frac{1}{t} \sum_{s=1}^{t-1} f(X_s).$$

Le théorème ergodique fournit un équivalent de la loi forte des grands nombres dans le cas d'observations markoviennes.

Theorème 2.1 (Théorème ergodique). *Soit $f : V \rightarrow \mathbb{R}$ et $(X_t)_{t \geq 0}$ une chaîne de Markov irréductible sur V , de loi stationnaire π . Alors pour tout $x \in V$,*

$$\mathbf{P}_x \left(\frac{1}{t} \sum_{s=1}^{t-1} f(X_s) \xrightarrow[t \rightarrow +\infty]{} \mathbf{E}_\pi f \right) = 1.$$

Autrement dit, si $(X_t)_{t \geq 0}$ est la marche aléatoire sur un graphe connexe $G = (V, E)$, alors, pour n'importe quel point de départ, \hat{f}_t est un estimateur consistant de $\mathbf{E}_\pi f$. Quitte à rendre la marche paresseuse pour assurer l'ergodicité, on peut aussi montrer le TCL suivant.

Theorème 2.2 (TCL pour les chaînes de Markov). *Soit $f : V \rightarrow \mathbb{R}$ et $(X_t)_{t \geq 0}$ une chaîne de Markov ergodique réversible sur V , de loi stationnaire π . Alors pour tout $x \in V$ et pour tout $\lambda \in \mathbb{R}$,*

$$\mathbf{P}_x \left(\frac{\sqrt{t}(\hat{f}_t - \mathbf{E}_\pi f)}{\sigma} \leq \lambda \right) \xrightarrow[t \rightarrow +\infty]{} \Phi(\lambda),$$

où $\sigma^2 = \text{Var}_\pi f + 2 \sum_{s=1}^{+\infty} \text{Cov}_\pi(f(X_0), f(X_s))$.

Remarque 2.1. Vérifions que la variance limite dans le TCL ci-dessous est bien finie. Quitte à recentrer, on peut supposer $\mathbf{E}_\pi f = 0$. En écrivant f dans la base orthonormée de vecteurs propres $(\psi_j)_{j=1}^n$ et en remarquant que $\langle f, \psi_1 \rangle_\pi = 0$ (car $\psi_1 = \mathbf{1}$ et f est centrée), on a

$$f = \sum_{j=2}^n \langle f, \psi_j \rangle_\pi \psi_j \quad \text{et} \quad \text{Var}_\pi f = \|f\|_\pi^2 = \sum_{j=2}^n \langle f, \psi_j \rangle_\pi^2.$$

Pour $s \geq 1$,

$$\begin{aligned} \text{Cov}_\pi(f(X_0), f(X_s)) &= \sum_{x, y \in V} \pi(x) P^s(x, y) f(x) f(y) \\ &= \langle f, P^s f \rangle_\pi \\ &= \left\langle \sum_{j=2}^n \langle f, \psi_j \rangle_\pi \psi_j, \sum_{j=2}^n \langle f, \psi_j \rangle_\pi \lambda_j^s \psi_j \right\rangle_\pi \\ &= \sum_{j=2}^n \lambda_j^s \langle f, \psi_j \rangle_\pi^2. \end{aligned}$$

Ainsi, comme pour tout $j \in \llbracket 2, n \rrbracket$, $|\lambda_j| < 1$ (car la chaîne est ergodique),

$$\begin{aligned} \sigma^2 &= \|f\|_\pi^2 + 2 \sum_{s \geq 1} \sum_{j=2}^n \lambda_j^s \langle f, \psi_j \rangle_\pi^2 \\ &= \|f\|_\pi^2 + 2 \sum_{j=2}^n \frac{\lambda_j}{1 - \lambda_j} \langle f, \psi_j \rangle_\pi^2 \\ &\leq \|f\|_\pi^2 \left(1 + \frac{2\lambda_2}{1 - \lambda_2} \right). \end{aligned}$$

Pour $(X_t)_{t \geq 0}$ la marche aléatoire paresseuse sur G connexe, l'estimateur \widehat{f}_t est donc asymptotiquement normal. Si l'on souhaite des bornes non-asymptotiques sur l'écart $|\widehat{f}_t - \mathbf{E}_\pi f|$, le problème est plus compliqué. Une solution est de laisser la marche évoluer jusqu'au temps de mélange $t_{\text{MIX}}(\varepsilon)$ est de calculer \widehat{f}_t seulement à partir de là. Le résultat suivant peut être vu comme un équivalent de l'inégalité de Bienaymé–Tchebychev pour les chaînes de Markov.

Proposition 2.3. *Soit $f : V \rightarrow \mathbb{R}$ et $(X_t)_{t \geq 0}$ une chaîne de Markov ergodique réversible sur V . Alors pour tout $x \in V$ et pour tout $t \geq 0$ et $\varepsilon, \delta > 0$,*

$$\mathbf{P}_x \left(\left| \frac{1}{t} \sum_{s=0}^{t-1} f(X_{t_{\text{MIX}}(\varepsilon)+s}) - \mathbf{E}_\pi f \right| \geq \delta \right) \leq \varepsilon + \frac{2 \text{Var}_\pi f}{\delta^2(1 - \lambda_2)t}.$$

Preuve de la Proposition 2.3. Tout d'abord, par la caractérisation de la distance en variation totale suivante :

$$\|\mu - \nu\|_{\text{TV}} = \sup_{f: V \rightarrow [0,1]} \{ \mathbf{E}_\mu f - \mathbf{E}_\nu f \},$$

on a

$$\begin{aligned} \mathbf{P}_x \left(\left| \frac{1}{t} \sum_{s=0}^{t-1} f(X_{t_{\text{MIX}}(\varepsilon)+s}) - \mathbf{E}_\pi f \right| \geq \delta \right) &= \sum_{y \in V} P^{t_{\text{MIX}}(\varepsilon)}(x, y) \mathbf{P}_y \left(\left| \frac{1}{t} \sum_{s=0}^{t-1} f(X_s) - \mathbf{E}_\pi f \right| \geq \delta \right) \\ &\leq \varepsilon + \mathbf{P}_\pi \left(\left| \frac{1}{t} \sum_{s=0}^{t-1} f(X_s) - \mathbf{E}_\pi f \right| \geq \delta \right). \end{aligned}$$

Ensuite, l'inégalité de Tchebychev donne

$$\mathbf{P}_\pi \left(\left| \frac{1}{t} \sum_{s=0}^{t-1} f(X_s) - \mathbf{E}_\pi f \right| \geq \delta \right) \leq \frac{\text{Var}_\pi \left(\sum_{s=0}^{t-1} f(X_s) \right)}{\delta^2 t^2}.$$

Quitte à considérer $f - \mathbf{E}_\pi f$, on peut supposer f centrée. On a alors

$$\begin{aligned} \text{Var}_\pi \left(\sum_{s=0}^{t-1} f(X_s) \right) &= t \text{Var}_\pi f + 2 \sum_{0 \leq i < j \leq t-1} \text{Cov}_\pi(f(X_i), f(X_j)) \\ &= t \text{Var}_\pi f + 2 \sum_{k=1}^{t-1} (t-k) \text{Cov}_\pi(f(X_0), f(X_k)) \\ &\leq t \text{Var}_\pi f + 2t \sum_{k=1}^{t-1} \text{Cov}_\pi(f(X_0), f(X_k)). \end{aligned}$$

Par la même décomposition spectrale que dans la remarque 2.1, on a $\text{Cov}_\pi(f(X_0), f(X_k)) \leq \lambda_2^k \text{Var}_\pi f$. Ainsi

$$\text{Var}_\pi \left(\sum_{s=0}^{t-1} f(X_s) \right) \leq \left(1 + \frac{2\lambda_2}{1 - \lambda_2} \right) t \text{Var}_\pi f \leq \frac{2t \text{Var}_\pi f}{1 - \lambda_2},$$

ce qui donne bien le résultat voulu. ■

Ainsi, la trajectoire de la marche aléatoire sur G permet d'estimer l'espérance d'une fonction f contre la mesure stationnaire π . Mais que peut-on dire si l'on souhaite plutôt estimer $\mathbf{E}_\mu f$, pour une mesure μ qui n'est pas forcément la mesure stationnaire ? L'algorithme de Metropolis–Hastings ci-dessous permet de transformer la marche simple en une chaîne de Markov dont μ est la loi stationnaire. Ainsi les résultats 2.1, 2.2, et 2.3 pourront s'appliquer à cette nouvelle chaîne, permettant d'estimer $\mathbf{E}_\mu f$.

2. Algorithme de Metropolis–Hastings

Les marches aléatoires fournissent un moyen d'échantillonner les sommets d'un graphe lorsque l'on n'a pas directement au graphe : si l'on attend assez longtemps (typiquement plus longtemps que le temps de mélange), la marche se trouve sur le sommet $u \in V$ avec probabilité à peu près $\pi(u) = \frac{\deg(u)}{2m}$. La plupart des grands réseaux réels sont connus pour avoir un temps de mélange relativement petit, typiquement logarithmique en la taille du graphe. On dispose donc d'un moyen simple et efficace pour échantillonner selon la loi π . Mais qu'en est-il si l'on souhaite échantillonner selon une autre loi que π ? Par exemple, on peut vouloir échantillonner les sommets selon la probabilité uniforme sur V , qui peut être très différente de π si les degrés sont très hétérogènes. L'algorithme de Metropolis–Hastings permet de modifier les transitions de la marche pour obtenir une chaîne de Markov qui possède la loi stationnaire souhaitée.

Soit P la matrice de transition de la marche aléatoire simple sur $G = (V, E)$, de loi stationnaire π . Et soit μ une autre probabilité V , selon laquelle on souhaite simuler. On définit, à partir de P , un noyau Q dont la loi stationnaire est μ , de la façon suivante : si l'état courant est $x \in V$, on génère $y \in V$ selon $P(x, \cdot)$, et l'on accepte cette transition en y avec probabilité :

$$r(x, y) = \frac{\mu(y)P(y, x)}{\mu(x)P(x, y)} \wedge 1.$$

Le rapport $r(x, y)$ est appelé rapport de Metropolis–Hastings. Le noyau de transition Q de cette nouvelle chaîne est donné par

$$(2.1) \quad Q(x, y) = \begin{cases} P(x, y)r(x, y), & \text{si } y \neq x, \\ 1 - \sum_{z \neq x} P(x, z)r(x, z), & \text{si } y = x. \end{cases}$$

Proposition 2.4. *Le noyau Q défini par (2.1) a pour loi stationnaire μ .*

Démonstration. Il suffit de vérifier que Q satisfait la condition d'équilibre ponctuel par rapport à μ , i.e. que pour tous $x, y \in V$, $\mu(x)Q(x, y) = \mu(y)Q(y, x)$. Soient x, y dans V avec $x \neq y$. Par symétrie on peut toujours supposer $\mu(y)P(y, x) \leq \mu(x)P(x, y)$, quitte à échanger les rôles de x et y (la condition d'équilibre ponctuel ne change pas si on permute x et y). Dans ce cas notons que

$$r(x, y) = \frac{\mu(y)P(y, x)}{\mu(x)P(x, y)}, \quad r(y, x) = 1.$$

Comme $x \neq y$, pour passer de x à y avec la chaîne définie par l'algorithme, il faut deux choses : générer y avec probabilité $P(x, y)$ et accepter le mouvement de x à y avec probabilité $r(x, y)$.

Ainsi

$$Q(x, y) = P(x, y)r(x, y) = \frac{\mu(y)P(y, x)}{\mu(x)}.$$

On en déduit que $\mu(x)Q(x, y) = \mu(y)P(y, x)$.

Par le même argument, on calcule $Q(y, x)$, qui vaut $P(y, x)r(y, x)$. Cette fois $r(y, x) = 1$, donc $P(y, x) = Q(y, x)$. Les deux identités précédentes mises ensemble donnent la condition d'équilibre ponctuel pour $x \neq y$, et celle-ci est immédiate si $x = y$. ■

Application. Si la loi souhaitée μ est la loi uniforme sur V , i.e. pour tout $x \in V$, $\mu(x) = \frac{1}{n}$, le rapport de Metropolis–Hastings s'écrit

$$r(x, y) = \frac{\deg(x)}{\deg(y)} \wedge 1.$$

Ainsi, lorsque la marche est en x , et que y est choisi selon P , i.e. uniformément parmi les voisins de x , la transition est acceptée avec une probabilité d'autant plus grande que le degré de y est petit. Cela se comprend bien : si l'on souhaite que la marche converge vers la probabilité uniforme, il faut favoriser les sommets de petit degré, sur lesquels π mais trop peu de poids. Cette nouvelle marche va alors permettre d'estimer des moyennes d'une fonction $f : V \rightarrow \mathbb{R}$ contre la mesure uniforme.

Échantillonnage d'importance. Une autre façon de faire, au lieu de passer par Metropolis–Hastings, est d'utiliser la méthode d'échantillonnage d'importance. Au lieu de modifier la marche, l'idée est de modifier les poids dans la moyenne empirique. Si l'on connaît les valeurs de $\pi(x)$ et $\mu(x)$ pour tout $x \in V$, alors on peut définir l'estimateur

$$\frac{1}{t} \sum_{s=0}^{t-1} \frac{\mu(X_s)}{\pi(X_s)} f(X_s),$$

qui, par le théorème ergodique, converge presque sûrement vers

$$\mathbf{E}_\pi \left[\frac{\mu(X)}{\pi(X)} f(X) \right] = \sum_{x \in V} \pi(x) \frac{\mu(x)}{\pi(x)} f(x) = \mathbf{E}_\mu f.$$

Le problème est que l'on ne connaît généralement pas les valeurs de $\pi(X_s)$ pour X_s sur la trajectoire (même si l'on suppose que l'on observe le degré des sommets visités, on ne connaît pas la constante de normalisation $2|E|$). On peut aussi ne pas connaître les valeurs de $\mu(X_s)$. Par exemple, si μ est l'uniforme sur V , on ne peut pas calculer $\mu(X_s)$ puisque l'on ne connaît pas la taille de V . Dans ce cas particulier où μ est la mesure uniforme sur V , on peut néanmoins procéder de la façon suivante : si $(X_t)_{t \geq 0}$ est la marche simple paresseuse sur G et que l'on suppose que l'on observe les degrés des sommets visités, on définit

$$\tilde{f}_t = \sum_{s=0}^{t-1} \frac{f(X_s)}{\deg(X_s)} \cdot \frac{1}{\sum_{k=0}^{t-1} \frac{1}{\deg(X_k)}}.$$

En appliquant deux fois le théorème ergodique, on a que, pour tout point de départ $x \in V$,

$$\frac{1}{t} \sum_{s=0}^{t-1} \frac{f(X_s)}{\deg(X_s)} \xrightarrow[t \rightarrow \infty]{\text{p.s.}} \sum_{x \in V} \pi(x) \frac{f(x)}{\deg(x)} = \frac{1}{2m} \sum_{x \in V} f(x),$$

et

$$\frac{1}{t \sum_{k=0}^{t-1} \frac{1}{\deg(X_k)}} \xrightarrow[t \rightarrow \infty]{\text{p.s.}} \sum_{x \in V} \pi(x) \frac{1}{\deg(x)} = \frac{n}{2m}.$$

Ainsi

$$\tilde{f}_t \xrightarrow[t \rightarrow \infty]{\text{p.s.}} \frac{1}{n} \sum_{x \in V} f(x),$$

qui correspond bien à la moyenne de f sous la loi uniforme.

3. Estimation de la taille du graphe

Dans cette question on s'intéresse à la question : peut-on, à partir d'une marche aléatoire sur G , avoir une idée de la taille de G ? Commençons d'abord par supposer que l'on a accès à des sommets distribués selon π .

3.1. Echantillons i.i.d. Soit $G = (V, E)$ un graphe, avec $|V| = n$ et $|E| = m$. Dans cette section, on considère que les observations dont on dispose pour mener l'estimation sont

$$\mathbf{X}^{(K)} = (X_1, \deg(X_1), \dots, X_K, \deg(X_K)) ,$$

où X_1, \dots, X_K sont des sommets i.i.d. distribué selon π (donnée par $\pi(u) = \frac{\deg(u)}{2m}$).

Le problème d'estimation statistique peut alors être formulé de la façon suivante : soit $\gamma(G)$ un paramètre d'intérêt du graphe, que l'on souhaite estimer (typiquement, $\gamma(G)$ sera le nombre de sommets $n = |V|$). L'objectif est de construire, à partir de l'échantillon $\mathbf{X}^{(K)}$, un estimateur $\hat{\gamma}_K$ qui soit tel que, pour tout graphe $G = (V, E)$, pour tout $\varepsilon, \delta > 0$, et pour tout $K \geq K_0(G, \varepsilon, \delta)$, on a

$$(2.1) \quad \mathbf{P} \left(\left| \frac{\hat{\gamma}_K}{\gamma(G)} - 1 \right| > \varepsilon \right) \leq \delta.$$

Autrement dit, on souhaite construire un estimateur qui approche bien le vrai paramètre, dès que l'échantillon est assez grand. Toute la difficulté se cache dans ce $K_0(G, \varepsilon, \delta)$, qui dépend non seulement des cibles de précision et d'erreur ε et δ mais surtout du graphe G lui-même, et qui correspond au moment où l'estimateur commence à fournir une bonne approximation du paramètre d'intérêt.

3.1.1. Le cas régulier. Commençons par un cas particulier relativement simple, mais qui illustre bien le problème. Supposons que le graphe $G = (V, E)$ est un graphe régulier, i.e. tous les sommets ont le même degré, disons $d \geq 1$. Dans ce cas, la mesure π est la loi uniforme sur V : pour tout $v \in V$, $\pi(v) = \frac{1}{n}$. Le paramètre que l'on souhaite estimer est $n = |V|$. Le problème revient donc à estimer la taille d'un ensemble V à partir d'un échantillon de K éléments tirés uniformément au hasard dans V . On peut déjà remarquer la chose suivante : si tous les éléments de l'échantillon sont distincts, alors il est impossible d'estimer correctement la taille. Pour pouvoir commencer à dire quelque chose, il faut qu'au moins un élément soit tiré plus d'une fois. C'est ce qu'on appelle une *collision*, et nous allons voir que l'ordre de

grandeur du temps de la première collision est en fait aussi suffisant pour avoir une bonne estimation de n . Avant cela, rappelons le principe du paradoxe des anniversaires.

Le paradoxe des anniversaires. Soit V un ensemble de taille n et (X_1, \dots, X_K) un échantillon indépendant d'éléments de V tirés uniformément au hasard. On peut se représenter V comme l'ensemble des jours de l'année, et (X_1, \dots, X_K) comme l'échantillon des dates de naissance dans une classe de K élèves, où l'on suppose que les dates de naissance sont indépendantes et uniformément distribuées sur l'année. Quelle est la probabilité que deux élèves soient nés le même jour ? La probabilité que tous les élèves aient des anniversaires différents s'écrit

$$\mathbf{P}(|\{X_1, \dots, X_K\}| = K) = \frac{n(n-1)\dots(n-K+1)}{n^K} = \exp\left\{\sum_{i=1}^{K-1} \log\left(1 - \frac{i}{n}\right)\right\}$$

Si $K > n/2$, on a $\mathbf{P}(|\{X_1, \dots, X_K\}| = K) \leq \left(\frac{3}{4}\right)^{K/2} = o(1)$. Supposons $K \leq n/2$. En utilisant que pour tout $u \in [0, 1/2]$, $-u - u^2 \leq \log(1-u) \leq -u$, on a d'une part

$$\sum_{i=1}^{K-1} \log\left(1 - \frac{i}{n}\right) \leq -\sum_{i=1}^{K-1} \frac{i}{n} = -\frac{K(K-1)}{n},$$

et d'autre part

$$\sum_{i=1}^{K-1} \log\left(1 - \frac{i}{n}\right) \geq -\sum_{i=1}^{K-1} \left(\frac{i}{n} + \frac{i^2}{n^2}\right) = -\frac{K(K-1)}{n} - \frac{K(K-1)(2K-1)}{6n^2},$$

où l'on a utilisé $\sum_{i=1}^k i^2 = \frac{k(k+1)(2k+1)}{6}$. On voit ainsi que

$$(2.2) \quad \mathbf{P}(|\{X_1, \dots, X_K\}| = K) = \begin{cases} 1 - o(1) & \text{si } K = o(\sqrt{n}), \\ e^{-a^2 + o(1)} & \text{si } K = a\sqrt{n}, \\ o(1) & \text{si } K \gg \sqrt{n}. \end{cases}$$

C'est donc lorsque la classe comporte $K \asymp \sqrt{n}$ élèves que la probabilité que deux élèves aient le même anniversaire devient non-négligeable. (Comme $\sqrt{365} \approx 19$, il est probable qu'une classe d'une vingtaine d'élèves en comporte deux nés le même jour, ce qui peut sembler surprenant, d'où le terme « paradoxe ».)

Voyons comment se servir des collisions dans l'échantillon pour construire un estimateur de n dans le cas régulier. Soit C_K le nombre de collisions observées, i.e.

$$C_K = \sum_{i < j} \mathbb{1}_{\{X_i = X_j\}}.$$

On a

$$\mathbf{E}C_K = \sum_{i < j} \mathbf{P}(X_i = X_j) = \binom{K}{2} \frac{1}{n},$$

et

$$\begin{aligned} \mathbf{E}C_K^2 &= \sum_{i < j} \mathbf{P}(X_i = X_j) + \sum_{\substack{i < j, k < \ell \\ |\{i, j, k, \ell\}|=3}} \mathbf{P}(X_i = X_j = X_k) + \sum_{\substack{i < j, k < \ell \\ |\{i, j, k, \ell\}|=4}} \mathbf{P}(X_i = X_j, X_k = X_\ell) \\ &= \binom{K}{2} \frac{1}{n} + \binom{K}{3} \frac{6}{n^2} + \binom{K}{2} \binom{K-2}{2} \frac{1}{n^2}. \end{aligned}$$

En remarquant que $\binom{K-2}{2} \leq \binom{K}{2}$, on obtient

$$\text{Var } C_K \leq \frac{K^2}{2n} + \frac{K^3}{n^2}.$$

En posant $\hat{n}_K = \frac{\binom{K}{2}}{C_K}$, on a alors, par l'inégalité de Chebyshev

$$\mathbf{P}\left(\left|\frac{n}{\hat{n}_K} - 1\right| > \varepsilon\right) = \mathbf{P}(|C_K - \mathbf{E}C_K| > \varepsilon \mathbf{E}C_K) \leq \frac{\text{Var } C_K}{\varepsilon^2 (\mathbf{E}C_K)^2} \asymp \varepsilon^{-2} \left(\frac{n}{K^2} + \frac{1}{K}\right).$$

Ainsi, dès que $K \gtrsim \frac{\sqrt{n}}{\varepsilon^2 \delta}$, on a

$$\mathbf{P}\left(\left|\frac{n}{\hat{n}_K} - 1\right| > \varepsilon\right) \leq \delta.$$

On a donc trouvé un estimateur \hat{n}_K de n qui vérifie (2.1) pour $K_0(G, \varepsilon, \delta) \asymp \frac{\sqrt{n}}{\varepsilon^2 \delta}$. Inversement, si $K < \sqrt{\log(1/\delta)n}$, alors par (2.2), $\mathbf{P}(C_K = 0) > \delta$. Or s'il n'y a aucune collision, aucun estimateur n'est capable d'approcher n universellement.

3.1.2. *Le cas général.* Voyons maintenant comment adapter ce genre de méthode au cas général où le graphe G n'est pas nécessairement régulier. On rappelle que l'on observe

$$\mathbf{X}^{(K)} = (X_1, \deg(X_1), \dots, X_K, \deg(X_K)),$$

où les sommets X_1, \dots, X_K sont i.i.d. de loi π donnée par $\pi(v) = \frac{\deg(v)}{2m}$ (avec $m = |E|$). On suppose que tous les sommets sont chargés par π , i.e. pour tous $v \in V$, $\deg(v) \geq 1$. Katzir et al. [14] ont proposé l'estimateur suivant :

$$\hat{n}_K = \frac{D_K}{C_K},$$

où

$$D_K = \sum_{i < j} \frac{\deg(X_i)}{\deg(X_j)},$$

et où $C_K = \sum_{i < j} \mathbb{1}_{\{X_i = X_j\}}$ est toujours le nombre de collisions. Remarquons que dans le cas régulier, cela correspond au même estimateur que précédemment.

Proposition 2.5. *Pour tout $\varepsilon, \delta > 0$, pour tout graphe $G = (V, E)$ avec $n = |V|$ et $m = |E|$, si*

$$K \gtrsim \frac{1}{\varepsilon^2 \delta} \max \left\{ \left(\sum \pi(u)^2 \right)^{-1/2}, \frac{m}{n} \right\},$$

alors

$$\mathbf{P}\left(\left|\frac{\hat{n}_K}{n} - 1\right| > \varepsilon\right) \leq \delta.$$

Comme on a toujours $\sum \pi(u)^2 \geq \frac{1}{n}$, il suffit d'avoir $K \gtrsim \frac{1}{\varepsilon^2 \delta} \max \left\{ \sqrt{n}, \frac{m}{n} \right\}$.

Preuve de la Proposition 2.5. On a

$$\mathbf{E}C_K = \binom{K}{2} \sum_{u \in V} \pi(u)^2,$$

et

$$\mathbf{E}D_K = \binom{K}{2} \sum_{u,v \in V} \pi(u)\pi(v) \frac{\deg(u)}{\deg(v)} = n \binom{K}{2} \sum_{u \in V} \pi(u)^2.$$

Ainsi $\frac{\mathbf{E}D_K}{\mathbf{E}C_K} = n$. Pour montrer que \hat{n}_K est bien concentré autour de n , remarquons que si $|D_K - \mathbf{E}D_K| \leq \frac{\varepsilon}{3} \mathbf{E}D_K$ et $|C_K - \mathbf{E}C_K| \leq \frac{\varepsilon}{3} \mathbf{E}C_K$, alors

$$(1 - \varepsilon)n \leq \frac{(1 - \varepsilon/3)\mathbf{E}D_K}{(1 + \varepsilon/3)\mathbf{E}C_K} \leq \frac{D_K}{C_K} \leq \frac{(1 + \varepsilon/3)\mathbf{E}D_K}{(1 - \varepsilon/3)\mathbf{E}C_K} \leq (1 + \varepsilon)n,$$

c'est-à-dire $\left| \frac{\hat{n}_K}{n} - 1 \right| \leq \varepsilon$. Il suffit donc de chercher K_0 tel que pour tout $K \geq K_0$,

$$\mathbf{P} \left(|D_K - \mathbf{E}D_K| > \frac{\varepsilon}{3} \mathbf{E}D_K \right) + \mathbf{P} \left(|C_K - \mathbf{E}C_K| > \frac{\varepsilon}{3} \mathbf{E}C_K \right) \leq \delta.$$

En utilisant l'inégalité de Chebyshev, on voit qu'il suffit de montrer que

$$(2.3) \quad \frac{\text{Var } D_K}{(\mathbf{E}D_K)^2} + \frac{\text{Var } C_K}{(\mathbf{E}C_K)^2} \leq \frac{\varepsilon^2 \delta}{9}.$$

On a d'une part

$$\mathbf{E}C_K^2 = \binom{K}{2} \sum_{u \in V} \pi(u)^2 + 6 \binom{K}{3} \sum_{u \in V} \pi(u)^3 + \binom{K}{2} \binom{K-2}{2} \left(\sum_{u \in V} \pi(u)^2 \right)^2.$$

Ainsi, en notant $\beta = \sum \pi(u)^2$,

$$\text{Var } C_K \leq \frac{K^2 \beta}{2} + K^3 \sum_{u \in V} \pi(u)^3 \leq \frac{K^2 \beta}{2} + K^3 \beta^{3/2}.$$

et

$$(2.4) \quad \frac{\text{Var } C_K}{(\mathbf{E}C_K)^2} \lesssim \frac{1}{K^2 \beta} + \frac{1}{K \sqrt{\beta}}.$$

D'autre part, en notant X, Y, Z trois variables indépendantes de loi π ,

$$\begin{aligned} \mathbf{E}D_K^2 &= \binom{K}{2} \mathbf{E} \left[\left(\frac{\deg(X)}{\deg(Y)} \right)^2 \right] \\ &\quad + 2 \binom{K}{3} \left(\mathbf{E} \left[\frac{\deg(X)}{\deg(Y)} \right] + \mathbf{E} \left[\frac{\deg(X)^2}{\deg(Y) \deg(Z)} \right] + \mathbf{E} \left[\frac{\deg(X) \deg(Y)}{\deg(Z)^2} \right] \right) \\ &\quad + \binom{K}{2} \binom{K-2}{2} \mathbf{E} \left[\frac{\deg(X)}{\deg(Y)} \right]^2. \end{aligned}$$

Comme le dernier terme de la somme ci-dessus est plus petit que $(\mathbf{E}D_K)^2$, on a

$$\begin{aligned} \text{Var } D_K &\leq \binom{K}{2} \sum_{u,v} \frac{\pi(u)^3}{\pi(v)} + 2 \binom{K}{3} \left(n \sum_u \pi(u)^2 + n^2 \sum_u \pi(u)^3 + \sum_{u,v,w} \frac{\pi(u)^2 \pi(v)^2}{\pi(w)} \right) \\ &\lesssim K^2 n m \beta^{3/2} + K^3 n \beta + K^3 n^2 \beta^{3/2} + K^3 n m \beta^2, \end{aligned}$$

où l'on a utilisé que $\sum \pi(u)^3 \leq \beta^{3/2}$ et $\pi(v) \geq \frac{1}{2m}$. Ainsi

$$(2.5) \quad \frac{\text{Var } D_K}{(\mathbf{E}D_K)^2} \lesssim \frac{m}{K^2 n \sqrt{\beta}} + \frac{1}{K n \beta} + \frac{1}{K \sqrt{\beta}} + \frac{m}{K n}.$$

En combinant (2.4) et (2.5), on voit que pour $K \gtrsim \frac{1}{\varepsilon^2 \delta} \max \left\{ \frac{1}{\sqrt{\beta}}, \frac{m}{n} \right\}$, l'inégalité (2.3) est vérifiée, et l'estimateur \hat{n}_K fournit une bonne approximation de n . ■

3.2. Estimation par marches aléatoires. Voyons maintenant comment passer au cas où l'on n'a plus directement accès à π . Soit $x \in V$ un sommet de départ fixé dans le graphe, et soient $X^{(1)}, \dots, X^{(K)}$ K marches aléatoires paresseuses indépendantes, toutes parties de $X_0^{(i)} = x$. Si le graphe G est connexe, et si l'on laisse évoluer les marches pendant un temps t assez grand, alors par le Théorème 1.6, l'échantillon des points d'arrivée $(X_t^{(1)}, \dots, X_t^{(K)})$ sera « quasiment » un échantillon i.i.d. de loi π , et l'on pourra appliquer les résultats de la section précédente pour l'estimation de la taille du graphe. Plus précisément, on montre que pour t assez grand, l'estimateur $\hat{n}_K(t)$ défini comme précédemment mais sur les points d'arrivée des marches, i.e.

$$\hat{n}_K(t) = \frac{D_K(t)}{C_K(t)} = \frac{\sum_{1 \leq i < j \leq K} \frac{\deg(X_t^{(i)})}{\deg(X_t^{(j)})}}{\sum_{1 \leq i < j \leq K} \mathbb{1}_{\{X_t^{(i)} = X_t^{(j)}\}}},$$

fournit une bonne approximation de n . Rappelons que pour $\alpha > 0$, le temps de mélange uniforme $t_{\text{UNIF}}(\alpha)$ est donné par

$$t_{\text{UNIF}}(\alpha) = \min \left\{ t \geq 0, \max_{x,y \in V} \left| \frac{\mathbf{P}_x(X_t = y)}{\pi(y)} - 1 \right| \leq \alpha \right\}.$$

Proposition 2.6. *Pour tout $\varepsilon, \delta > 0$, pour tout graphe $G = (V, E)$ connexe, pour tout point de départ $x \in V$, si $K \gtrsim \frac{1}{\varepsilon^2 \delta} \max \left\{ \sqrt{n}, \frac{m}{n} \right\}$ et si $t \gtrsim t_{\text{UNIF}}(\alpha)$ pour $\alpha \in]0, 1[$ assez petit, alors*

$$\mathbf{P}_x \left(\left| \frac{\hat{n}_K(t)}{n} - 1 \right| > \varepsilon \right) \leq \delta.$$

Preuve de la Proposition 2.6. Par définition de $t_{\text{UNIF}}(\alpha)$, si $t \geq t_{\text{UNIF}}(\alpha)$, alors pour tous $x, u \in V$,

$$1 - \alpha \leq \frac{\mathbf{P}_x(X_t = u)}{\pi(u)} \leq 1 + \alpha.$$

On a donc

$$(1 - \alpha)^2 \mathbf{E}C_K \leq \mathbf{E}_x C_K(t) = \binom{K}{2} \sum_{u \in V} \mathbf{P}_x(X_t = u)^2 \leq (1 + \alpha)^2 \mathbf{E}C_K,$$

et idem pour $\mathbf{E}_x D_K(t)$. De même

$$(1 - \alpha)^4 \mathbf{E} C_K^2 \leq \mathbf{E}_x C_K(t)^2 \leq (1 + \alpha)^4 \mathbf{E} C_K^2,$$

et idem pour $\mathbf{E}_x D_K(t)^2$. Ainsi

$$\mathrm{Var}_x C_K(t) \lesssim \mathrm{Var} C_K + \alpha (\mathbf{E} C_K)^2 \quad \text{et} \quad \mathrm{Var}_x D_K(t) \lesssim \mathrm{Var} D_K + \alpha (\mathbf{E} D_K)^2.$$

Si $\alpha \leq c\varepsilon^2\delta$ pour $c > 0$ une constante assez petite, alors il suffit d'appliquer les mêmes arguments que dans la preuve de la Proposition 2.5 et on a bien le résultat voulu. ■

Ce résultat est frustrant pour plusieurs raisons. Tout d'abord, pour construire l'estimateur $\hat{n}_{K,t}$, il faut connaître le temps de mélange t_{UNIF} de la marche G (ou au moins avoir une borne supérieure sur celui-ci). D'autre part, en lançant K marches de longueur t et en ne retenant que leurs points d'arrivée pour mener l'inférence, on néglige beaucoup d'information qui peut être utile. Ben-Hamou et al. [3] ont montré que l'on pouvait construire un estimateur plus performant (et même optimal) en utilisant toute la trajectoire des marches. Cet estimateur repose non plus sur le nombre de collisions dans l'échantillon formé par les points d'arrivée, mais sur le nombre d'intersections entre les trajectoires de deux marches aléatoires indépendantes.

Détection de communautés

Il existe énormément de variantes du problème de la détection de communautés. Ici, nous considérons une version relativement simple à deux communautés. Soit V un ensemble à $2n$ éléments, avec $n \in \mathbb{N}^*$, et soit $\mathcal{G}(2n, p, q)$ un graphe aléatoire d'ensemble de sommets V et dont les arêtes sont générées de la façon suivante : on commence par attribuer à chaque sommet un groupe d'appartenance sous la forme d'un label égal à $+1$ ou -1 en choisissant le vecteur d'appartenance $\sigma = (\sigma_v)_{v \in V}$ uniformément au hasard dans l'ensemble

$$\mathcal{X} = \left\{ \tau \in \{-1, 1\}^V, \sum_{v \in V} \tau_v = 0 \right\}.$$

Autrement dit, on choisit uniformément au hasard une partition de V en deux blocs de même taille, et l'on attribue l'étiquette $+1$ aux sommets d'un bloc, et l'étiquette -1 aux sommets de l'autre bloc. Deux sommets du même groupe sont connectés avec probabilité p , et deux sommets de groupe différent sont connectés avec probabilité q , tout cela de façon indépendante. On considérera parfois un modèle très proche, noté $\mathcal{G}'(2n, p, q)$ où l'étiquette de chaque sommet est tirée indépendamment uniformément dans $\{-1, 1\}$, i.e. le vecteur des appartenances σ est tiré uniformément dans $\{-1, 1\}^V$, et les arêtes sont générées comme dans $\mathcal{G}(2n, p, q)$. Par la concentration de la loi binomiale $\mathcal{B}(2n, 1/2)$ autour de n , les deux modèles sont très proches. Pour les deux modèles, on a, pour $g = (V, E)$ un graphe fixé,

$$(3.1) \quad \mathbf{P}(G = g \mid \sigma) = \prod_{\{u,v\} \in E} p^{\delta_{\sigma_u, \sigma_v}} q^{1 - \delta_{\sigma_u, \sigma_v}} \prod_{\{u,v\} \notin E} (1-p)^{\delta_{\sigma_u, \sigma_v}} (1-q)^{1 - \delta_{\sigma_u, \sigma_v}}.$$

On suppose $p \geq q$ (deux sommets du même groupe ont plus de chances d'être connectés que deux sommets de groupes différents). Le problème de la détection de communautés est alors de savoir si, en observant le graphe G , on est capable de retrouver le groupe des sommets (à permutation des groupes près). Il y a en fait plusieurs objectifs possibles, que l'on liste du plus au moins exigeants :

- la reconstruction exacte : on veut retrouver, avec probabilité tendant vers 1, le groupe de tous les sommets, i.e. construire, à partir de G , un vecteur $\hat{\sigma}$ tel que

$$\mathbf{P}(\hat{\sigma} = \pm \sigma) = \mathbf{P} \left(\left| \sum_{v \in V} \hat{\sigma}_v \sigma_v \right| = 2n \right) \xrightarrow{n \rightarrow \infty} 1.$$

- la reconstruction faible : on accepte de se tromper mais seulement sur un nombre négligeable de sommets, i.e. on veut construire un vecteur $\hat{\sigma}$ tel que

$$\frac{|\sum_{v \in V} \hat{\sigma}_v \sigma_v|}{2n} \xrightarrow{\mathbf{P}} 1.$$

- la détection : on veut faire strictement mieux que l'estimateur obtenu en tirant indépendamment à pile ou face le groupe de chaque sommet, i.e. on veut construire un vecteur $\hat{\sigma}$ qui soit asymptotiquement strictement positivement corrélé avec σ :

$$\mathbf{P} \left(\frac{|\sum_{v \in V} \hat{\sigma}_v \sigma_v|}{2n} \geq \varepsilon \right) \xrightarrow{n \rightarrow \infty} 1,$$

pour un certain $\varepsilon > 0$.

Observons tout d'abord que si $p = q$, alors il n'y a aucun espoir de pouvoir reconstruire les communautés : le graphe observé est un graphe d'Erdős–Rényi classique, et les sommets sont indistinguables. Intuitivement, plus l'écart $p - q$ est petit, plus le problème est difficile. Non seulement, il est possible que si $p - q$ est trop petit, alors reconstruire correctement les communautés soit impossible, mais, même s'il existe une solution, il est possible qu'elle soit extrêmement dure à trouver. Pour chacun des problèmes ci-dessous, on veut non seulement montrer qu'un tel estimateur $\hat{\sigma}$ existe, mais aussi qu'on peut le calculer en un temps raisonnable (i.e. polynomial en n). Or cela est loin d'être garanti. Pour un théoricien de l'informatique, la reconstruction de communautés fait immédiatement penser au problème de la recherche d'une bisection minimale dans un graphe, qui est connu pour être NP-complet.

Notons que, pour ce qui est de la reconstruction exacte, on dispose d'un estimateur tout indiqué pour être le meilleur estimateur possible de σ : un estimateur du maximum de vraisemblance donné par

$$\hat{\sigma} \in \arg \max_{\sigma \in \mathcal{X}} \mathbf{P}(G \mid \sigma),$$

Pour voir qu'on ne peut pas espérer faire mieux que $\hat{\sigma}$, observons que

$$\mathbf{P}(\sigma \mid G) = \frac{\frac{1}{|\mathcal{X}|} \mathbf{P}(G \mid \sigma)}{\mathbf{P}(G)},$$

où $\mathbf{P}(G) = \sum_{\eta \in \mathcal{X}} \frac{1}{|\mathcal{X}|} \mathbf{P}(G \mid \eta)$. Ainsi, maximiser $\mathbf{P}(G \mid \sigma)$, c'est maximiser $\mathbf{P}(\sigma \mid G)$. Donc pour tout autre estimateur $\hat{\eta}$,

$$\begin{aligned} \mathbf{P}(\hat{\eta} = \sigma) &= \sum_g \mathbf{P}(G = g, \sigma = \hat{\eta}(g)) = \sum_g \frac{1}{|\mathcal{X}|} \mathbf{P}(G = g \mid \sigma = \hat{\eta}(g)) \\ &\leq \sum_g \frac{1}{|\mathcal{X}|} \mathbf{P}(G = g \mid \sigma = \hat{\sigma}(g)) = \mathbf{P}(\hat{\sigma} = \sigma). \end{aligned}$$

Ainsi $\hat{\sigma}$ est l'estimateur qui a le plus de chances de retrouver la partition cachée. On pourrait alors croire le problème résolu mais il y a plusieurs problèmes :

- l'estimateur du maximum de vraisemblance peut ne pas être unique. Or s'il existe plusieurs maximiseurs de la vraisemblance, alors on ne peut pas espérer reconstruire exactement les communautés avec probabilité $1 - o(1)$: un de ces maximiseurs correspondra à la vraie partition cachée mais les autres peuvent être complètement différents ;
- si les communautés sont assez fortes, ou si le graphe est assez dense, alors alors il y a de fortes chances pour que $\hat{\sigma}$ soit effectivement égal à σ ; mais si le bruit l'emporte sur le signal, il y a alors une grande probabilité pour que $\hat{\sigma}$ ne reflète que l'effet du bruit

et non pas du signal (on parle d'*overfitting*). On a alors intérêt à prendre en compte toute la loi a posteriori $P(\sigma \mid G)$ et pas seulement son mode : cette fonction peut avoir plusieurs optima locaux qui n'ont rien à voir les un avec les autres, et σ peut ne pas être égal à l'optimum global ;

- même si cet estimateur est unique, il peut être extrêmement long à déterminer : on ne peut pas se permettre de calculer toutes les valeurs $\mathbf{P}(G \mid \sigma)$ pour σ parcourant \mathcal{X} qui est de taille $|\mathcal{X}| = \binom{2n}{n}$.

Avant de prouver certains résultats de faisabilité ou d'impossibilité concernant les différents types de reconstruction, commençons par souligner les liens entre le problème de la reconstruction de communautés et un modèle bien connu en physique statistique, le modèle d'Ising, et présentons un algorithme très performant en pratique : la propagation de croyances (*belief propagation*).

1. Modèle d'Ising et belief propagation

Étant donné un graphe $G = (V, E)$ et un paramètre $T > 0$, la probabilité sur $\sigma \in \{-1, 1\}^V$ donnée par

$$\mu(\sigma \mid G) \propto \exp \left\{ -\frac{H_E(\sigma)}{T} \right\}, \quad H_E(\sigma) = -J \sum_{\{u,v\} \in E} \delta_{\sigma_u, \sigma_v},$$

appelée distribution de Boltzmann, définit le modèle d'Ising sur G , avec température T et constante de couplage $J \in \mathbb{R}$. En physique, la quantité $H_E(\sigma)$ s'appelle le Hamiltonien, ou encore l'énergie de la configuration σ . La probabilité du modèle d'Ising met donc exponentiellement plus de poids sur les configurations de basse énergie. Lorsque $J > 0$, on parle d'interactions ferromagnétiques : la distribution de Boltzmann favorise les configurations dans lesquelles des sommets voisins ont les mêmes spins. Lorsque $J < 0$, les interactions sont dites anti-ferromagnétiques et la distribution de Boltzmann favorise les configurations où les voisins ont des spins opposés. Alors que les systèmes physiques tendent à aller vers des états à basse énergie, la température les poussent vers des états à plus haute énergie. Quand $T \rightarrow +\infty$, la distribution de Boltzmann tend vers la probabilité uniforme sur $\{-1, 1\}^V$, et quand $T \rightarrow 0$, la distribution de Boltzmann se concentre sur les états qui minimisent l'énergie (par exemple pour $J > 0$, sur les deux états complètement alignés, où les spins valent tous 1 ou -1).

On a vu que $\mathbf{P}(\sigma \mid G) \propto \mathbf{P}(G \mid \sigma)$. De plus, comme $q^{|E|}(1-q)^{\binom{n}{2}-|E|}$ est un terme constant par rapport à σ , on a

$$\begin{aligned} \mathbf{P}(\sigma \mid G) &\propto \prod_{\{u,v\} \in E} \left(\frac{p}{q}\right)^{\delta_{\sigma_u, \sigma_v}} \prod_{\{u,v\} \notin E} \left(\frac{1-p}{1-q}\right)^{\delta_{\sigma_u, \sigma_v}} \\ &\propto \exp \left\{ \log \left(\frac{p}{q}\right) \sum_{\{u,v\} \in E} \delta_{\sigma_u, \sigma_v} - \log \left(\frac{1-p}{1-q}\right) \sum_{\{u,v\} \notin E} \delta_{\sigma_u, \sigma_v} \right\} \\ &\propto \exp \left\{ -H_E(\sigma) - H_{\bar{E}}(\sigma) \right\}, \end{aligned}$$

où

$$H_E(\sigma) = -\log\left(\frac{p}{q}\right) \sum_{\{u,v\} \in E} \delta_{\sigma_u, \sigma_v} \quad \text{et} \quad H_{\bar{E}}(\sigma) = \log\left(\frac{1-q}{1-p}\right) \sum_{\{u,v\} \notin E} \delta_{\sigma_u, \sigma_v}.$$

Ainsi $\mathbf{P}(\sigma \mid G)$ peut être vue comme une distribution de Boltzmann dans un modèle d'Ising à température $T = 1$ où l'on considère à la fois des interactions à courte portée (i.e. sur les arêtes du graphe) et des interactions à longue portée (i.e. hors des arêtes du graphe). Les interactions à courte portée sont ferromagnétiques : $J_E = \log\left(\frac{p}{q}\right) > 0$ (puisque $p > q$) alors que les interactions à longue portée sont anti-ferromagnétiques : $J_{\bar{E}} = -\log\left(\frac{1-q}{1-p}\right) < 0$.

D'un point de vue bayésien, la probabilité $\mathbf{P}(\sigma \mid G)$ s'interprète comme la loi a posteriori : c'est tout ce que l'on peut espérer apprendre de σ en observant G . Même si l'on dispose de ressources de calcul illimitées, si cette distribution ne contient pas assez d'informations sur σ , on ne peut pas espérer pouvoir reconstruire les communautés. En admettant que l'on puisse la trouver en temps raisonnable, la partition la plus naturelle que l'on peut extraire de la loi $\mathbf{P}(\sigma \mid G)$ est son mode, soit l'état de plus faible énergie :

$$\hat{\sigma} = \arg \max_{\sigma} \mathbf{P}(\sigma \mid G).$$

La reconstruction exacte est possible si et seulement si $\hat{\sigma} = \sigma$ (avec probabilité $1 - o(1)$), c'est-à-dire si l'état le plus probable pour la loi a posteriori correspond à la partition cachée. Comme nous le verrons dans la prochaine section, cela ne peut se produire que si le graphe est assez dense ou si les communautés sont assez fortes. Or la plupart des systèmes physiques sont sparses, chaque atome n'interagissant qu'avec quelques voisins. Le bruit l'emporte alors sur le signal et la configuration $\hat{\sigma}$ ne sera pas exactement égale à σ . On peut alors revoir à la baisse le critère de reconstruction, et accepter de se tromper pour un petit nombre de sommets.

Plutôt que le maximum de vraisemblance global $\hat{\sigma}$, une autre façon d'estimer σ est d'attribuer à chaque sommet le spin qui a la plus grande probabilité marginale. Plus précisément, pour tout $u \in V$, et $x \in \{-1, 1\}$, on définit

$$\psi_u(x) = \mathbf{P}(\sigma_u = x \mid G) = \sum_{\sigma, \sigma_u = x} \mathbf{P}(\sigma \mid G).$$

L'estimateur $\tilde{\sigma}$ est alors donné par

$$\forall u \in V, \tilde{\sigma}_u = \arg \max_{x \in \{-1, 1\}} \psi_u(x).$$

Cet estimateur maximise l'espérance de la fraction des sommets bien étiquetés. Ainsi, c'est le meilleur estimateur pour la reconstruction faible. Si les marginales sont uniformes, alors la fraction des sommets bien étiquetés vaut 0 et même la détection est impossible¹. Comparé à $\hat{\sigma}$, l'estimateur $\tilde{\sigma}$ est mieux adapté à la reconstruction faible : quand le système est sparse et bruyé, prendre une moyenne sur toutes les lois marginales est souvent une meilleure solution que de se focaliser sur le maximum global.

1. En fait, par symétrie, les marginales proprement dites sont toujours uniformes, mais on peut rompre cette symétrie en fixant arbitrairement le label de quelques sommets.

Comment calculer les marginales ψ_u ? Une approche populaire est l'échantillonnage de Monte Carlo par chaînes de Markov (MCMC) : en flippant l'un après l'autre les spins de sommets choisis au hasard selon la distribution marginale du spin de ce sommet conditionnelle aux spins des autres sommets, on obtient une chaîne de Markov sur $\{-1, 1\}^V$, appelée la dynamique de Glauber, qui converge vers la distribution $\mathbf{P}(\sigma \mid G)$. On peut ainsi obtenir plusieurs échantillons de σ quasiment distribués selon la distribution de Boltzmann et utiliser ces échantillons pour estimer les lois marginales.

Supposons que, conditionnellement à σ_u , les spins σ_v pour $v \in \mathcal{N}(u)$ sont indépendants, c'est-à-dire que les voisins d'un sommet u ne sont corrélés que du fait de u . Notons que cette approximation est valable si le graphe G est un arbre : si tous les chemins entre les voisins de u passent par u , alors l'hypothèse d'indépendance conditionnelle est vérifiée. Si le graphe G n'est pas exactement un arbre mais est *localement arborescent*, i.e. si les cycles dans G sont relativement grands et que les corrélations décroissent assez vite avec la distance, alors on peut espérer que cette hypothèse reste une bonne approximation. Dans le cas sparse où p et q sont d'ordre $1/n$, on peut montrer que G est effectivement localement arborescent avec grande probabilité.

Faisons une autre simplification, dite de *champ moyen* : supposons que les interactions à longue portée (celles entre les non-voisins) peuvent être considérées comme un effet global de l'environnement (encore une fois, dans le cas sparse, cette approximation n'est pas aberrante). Autrement dit, on suppose que la loi $\mathbf{P}(\sigma \mid G)$ peut s'écrire

$$\mathbf{P}(\sigma \mid G) = \prod_{u \in V} \psi_u(\sigma_u) \prod_{\{u,v\} \in E} \frac{\psi_{u,v}(\sigma_u, \sigma_v)}{\psi_u(\sigma_u)\psi_v(\sigma_v)},$$

où $\psi_{u,v}$ est la loi marginale du couple (σ_u, σ_v) . L'algorithme de propagation de croyances est un algorithme de transmission de messages itératif, où à chaque itération, un sommet u envoie à un de ses voisins $v \in \mathcal{N}(u)$ un message $\psi_{u \rightarrow v}$ qui correspond à une estimée de la loi marginale de σ_u si v n'était pas voisin de u . Plus précisément, on part d'une certaine initialisation des messages $(\psi_{u \rightarrow v}^{(0)})_{u \neq v}$, et à chaque temps t , on choisit $u \in V$ uniformément au hasard, et on met à jour le message envoyé par u à chacun de ses voisins $v \in \mathcal{N}(u)$ par

$$\psi_{u \rightarrow v}^{(t+1)}(x) \propto \prod_{w \in \mathcal{N}(u) \setminus \{v\}} \left(p\psi_{w \rightarrow u}^{(t)}(x) + q\psi_{w \rightarrow u}^{(t)}(-x) \right).$$

Si tout se passe bien (et on sait que tout se passe bien sur un arbre), alors cet algorithme converge rapidement vers un point fixe qui permet de bien estimer les marginales en prenant une moyenne des messages reçus par un sommet :

$$\hat{\psi}_u(x) \propto \prod_{v \in \mathcal{N}(u)} \left(p\psi_{v \rightarrow u}^{(\infty)}(x) + q\psi_{v \rightarrow u}^{(\infty)}(-x) \right).$$

2. Reconstruction exacte

Un des premiers résultats concernant la faisabilité de la reconstruction exacte est dû à Dyer and Frieze [11] et énonce que si $p - q$ ne décroît pas trop vite vers 0, alors il est possible de retrouver exactement la partition en temps polynomial.

Proposition 3.1 (Dyer and Frieze [11]). *Soit $G \sim \mathcal{G}(2n, p, q)$ avec $p - q \geq \left(5 \frac{\log n}{n}\right)^{1/4}$. Alors la reconstruction exacte est possible en un temps $O(n^3)$.*

Remarque 3.1. Le résultat de Dyer and Frieze [11] est en fait plus précis : les auteurs montrent qu'avec probabilité $1 - o(1)$, la coupe minimale (la partition de V en deux sous-ensembles de taille égale avec le plus petit nombre d'arêtes entre les blocs) est uniquement donnée par (V_{+1}, V_{-1}) , où $V_x = \{v \in V, \sigma_v = x\}$, et que l'on peut non seulement trouver cette coupe mais *montrer* que c'est bien la coupe minimale en un temps polynomial.

Preuve de la Proposition 3.1. L'estimateur $\hat{\sigma}$ est construit de la façon suivante : soit $w \in V$ un sommet choisi arbitrairement. On fixe $\hat{\sigma}_w = 1$. Et pour tout $u \in V \setminus \{w\}$, on calcule $X_u = |V(u) \cap V(w)|$ où $V(u) = \{v \in V, \{u, v\} \in E\}$ est le voisinage de u . On classe les $(X_u)_{u \neq w}$ par ordre décroissant : $X_{(1)} \geq \dots \geq X_{(2n-1)}$. On pose $\hat{\sigma}_u = 1$ si $X_u \geq X_{(n-1)}$ et $\hat{\sigma}_u = -1$ si $X_u \leq X_{(n)}$.

Remarquons déjà que $\hat{\sigma}$ peut être construit en un temps $O(n^3)$. Montrons maintenant que $\mathbf{P}(\hat{\sigma} \neq \pm\sigma) \rightarrow 0$. Si u est un sommet de $V \setminus \{w\}$ tel que $\sigma_u = \sigma_w$, alors

$$X_u \stackrel{\mathcal{L}}{=} \sum_{k=1}^{n-2} B_k + \sum_{k=1}^n B'_k,$$

où les B_k, B'_k sont indépendantes, avec $B_k \sim \mathcal{B}(p^2)$ et $B'_k \sim \mathcal{B}(q^2)$. En particulier $\mathbf{E}[X_u \mid \sigma_u = \sigma_w] = (n-2)p^2 + nq^2$. D'autre part, si u est un sommet de $V \setminus \{w\}$ tel que $\sigma_u \neq \sigma_w$, alors

$$X_u \stackrel{\mathcal{L}}{=} \sum_{k=1}^{2n-2} B''_k,$$

où les B''_k sont indépendantes, avec $B''_k \sim \mathcal{B}(pq)$. En particulier, $\mathbf{E}[X_u \mid \sigma_u \neq \sigma_w] = (2n-2)pq$. Par une borne union et l'inégalité de Hoeffding 6.7, on obtient d'une part

$$\mathbf{P}\left(\exists u \in V \setminus \{w\}, \sigma_u = \sigma_w, X_u \leq (n-2)p^2 + nq^2 - \frac{n}{3}(p-q)^2\right) \leq ne^{-\frac{2n}{9}(p-q)^4} = o(1),$$

où l'on a utilisé $(p-q)^4 \geq 5 \frac{\log n}{n}$, et d'autre part

$$\mathbf{P}\left(\exists u \in V \setminus \{w\}, \sigma_u \neq \sigma_w, X_u \geq (2n-2)pq + \frac{n}{3}(p-q)^2\right) \leq ne^{-\frac{2n}{9}(p-q)^4} = o(1).$$

Or

$$(n-2)p^2 + nq^2 - (2n-2)pq - \frac{2n}{3}(p-q)^2 = \frac{n}{3}(p-q)^2 - 2p^2 + 2pq \gg 1.$$

Donc

$$\mathbf{P}(\exists u, v \in V \setminus \{w\}, \sigma_u = \sigma_w, \sigma_v \neq \sigma_w, X_u < X_v) = o(1).$$

Donc avec probabilité $1 - o(1)$, l'estimateur $\hat{\sigma}$ retrouve la bonne partition (au signe près). ■

Le résultat de Dyer et Frieze concerne les graphes relativement denses : le degré moyen, $(p+q)n$, doit être au moins d'ordre $n^{3/4}(\log n)^{1/4}$. Intuitivement, le problème de la reconstruction est d'autant plus difficile que le graphe est sparse (moins il y a d'arêtes) et que l'écart $p - q$ est petit. Jusqu'à quelle vitesse peut-on faire tendre l'écart $p - q$ vers 0 et toujours pouvoir retrouver la partition cachée ? Remarquons déjà que pour pouvoir retrouver le groupe de tous

les sommets, il est nécessaire qu'aucun des sommets ne soit isolé. Or il est connu que pour que le graphe d'Erdős–Renyi $\mathcal{G}(n, p)$ ne contienne pas de sommet isolé, il faut que $p \geq \frac{\log n}{n}$. De même, pour que $\mathcal{G}(2n, p, q)$ n'ait pas de points isolés, il faut que $p + q \geq \frac{\log n}{n}$. Allons plus loin et montrons que, pour que la reconstruction exacte dans $\mathcal{G}(2n, p, q)$ avec $p > q$ soit possible, il faut que p et q soient tels que, si $X \sim \mathcal{B}(n-1, p)$ et $Y \sim \mathcal{B}(n, q)$, avec X et Y indépendantes, alors

$$(3.1) \quad n\mathbf{P}(Y \geq X) \xrightarrow{n \rightarrow \infty} 0.$$

L'intuition derrière cette condition est la suivante : supposons que toutes les étiquettes ont été correctement retrouvées, sauf pour un sommet v . Conditionnellement au groupe de tous les autres sommets, comment estimer le groupe de ce sommet ? L'estimateur du maximum de vraisemblance dans ce cas consiste à attribuer à v l'étiquette présente sur la majorité de ses voisins. Mais s'il se trouve que la majorité des voisins de v ont l'étiquette opposée à celle de v , alors on se trompe. Ainsi, pour pouvoir retrouver l'étiquette de tous les sommets, il faut qu'avec probabilité $1 - o(1)$, pour tout $v \in V$, le nombre de voisins de v avec la même étiquette soit plus grand que le nombre de voisins avec l'étiquette opposée. Autrement dit, il faut que

$$\mathbf{P}\left(\bigcup_{v \in V} A_v\right) \xrightarrow{n \rightarrow \infty} 0,$$

où

$$A_v = \left\{ \sum_{u \sim v} \mathbb{1}_{\sigma_u \neq \sigma_v} \geq \sum_{u \sim v} \mathbb{1}_{\sigma_u = \sigma_v} \right\}.$$

En utilisant que $\sum_{u \sim v} \mathbb{1}_{\sigma_u = \sigma_v} \sim \mathcal{B}(n-1, p)$ et $\sum_{u \sim v} \mathbb{1}_{\sigma_u \neq \sigma_v} \sim \mathcal{B}(n, q)$, que ces deux variables sont indépendantes, et en admettant qu'une borne union donne une probabilité équivalente, on arrive à la condition (3.1). Donnons maintenant une preuve plus rigoureuse. Pour cela, commençons par le résultat intermédiaire suivant.

Lemme 3.2. *Si $n\mathbf{P}(Y \geq X)$ ne tend pas vers 0, alors il existe $\delta > 0$ tel que, pour une infinité d'entiers n ,*

$$\mathbf{P}\left(\bigcup_{v \in V} A_v\right) \geq \delta.$$

Preuve du Lemme 3.2. Supposons que $n\mathbf{P}(X \geq Y)$ ne tend pas vers 0. Posons

$$N = \sum_{v \in V} \mathbb{1}_{A_v},$$

le nombre de sommets dont la majorité des voisins a l'étiquette opposée. Par hypothèse, il existe $\varepsilon > 0$ tel que, pour une infinité d'entiers n ,

$$\mathbf{E}N = 2n\mathbf{P}(X \geq Y) \geq \varepsilon.$$

Par l'inégalité de Paley–Zygmund 6.6,

$$\mathbf{P}\left(N \geq \frac{\varepsilon}{2}\right) \geq \mathbf{P}\left(N \geq \frac{\mathbf{E}N}{2}\right) \geq \frac{\mathbf{E}[N]^2}{4\mathbf{E}[N^2]}.$$

Cherchons maintenant à majorer $\mathbf{E}[N^2]$. On a

$$\mathbf{E}[N^2] = \sum_{v \in V} \mathbf{P}(A_v) + \sum_{u \neq v} \mathbf{P}(A_u \cap A_v).$$

Pour calculer $\mathbf{P}(A_u \cap A_v)$, on décompose selon que u et v ont ou non la même étiquette et selon que u et v sont voisins ou non. On a

$$\begin{aligned} \mathbf{P}(A_u \cap A_v \cap \{u \sim v\} \cap \{\sigma_u = \sigma_v\}) &= \frac{1}{2} \cdot p \cdot \mathbf{P}(\mathcal{B}(n, q) \geq 1 + \mathcal{B}(n-2, p))^2 \\ &\leq \frac{p}{2} \mathbf{P}(X \geq Y)^2, \end{aligned}$$

et

$$\begin{aligned} \mathbf{P}(A_u \cap A_v \cap \{u \not\sim v\} \cap \{\sigma_u \neq \sigma_v\}) &= \frac{1}{2} \cdot (1-q) \cdot \mathbf{P}(\mathcal{B}(n-1, q) \geq \mathcal{B}(n-1, p))^2 \\ &\leq \frac{1-q}{2} \mathbf{P}(X \geq Y)^2. \end{aligned}$$

Les deux autres cas sont un peu moins faciles à contrôler. On a

$$\mathbf{P}(A_u \cap A_v \cap \{u \not\sim v\} \cap \{\sigma_u = \sigma_v\}) = \frac{1}{2} \cdot (1-p) \cdot \mathbf{P}(\mathcal{B}(n, q) \geq \mathcal{B}(n-2, p))^2.$$

Or

$$\mathbf{P}(Y \geq X) = p \mathbf{P}(\mathcal{B}(n, q) \geq 1 + \mathcal{B}(n-2, p)) + (1-p) \mathbf{P}(\mathcal{B}(n, q) \geq \mathcal{B}(n-2, p)).$$

Ainsi

$$\mathbf{P}(\mathcal{B}(n, q) \geq \mathcal{B}(n-2, p)) \leq \frac{\mathbf{P}(Y \geq X)}{1-p},$$

et

$$\mathbf{P}(A_u \cap A_v \cap \{u \not\sim v\} \cap \{\sigma_u = \sigma_v\}) \leq \frac{\mathbf{P}(Y \geq X)^2}{2(1-p)} \leq \frac{3\mathbf{P}(Y \geq X)^2}{2},$$

puisque'on a supposé $p \leq 2/3$. Enfin

$$\begin{aligned} \mathbf{P}(A_u \cap A_v \cap \{u \sim v\} \cap \{\sigma_u \neq \sigma_v\}) &= \frac{1}{2} \cdot q \cdot \mathbf{P}(1 + \mathcal{B}(n-1, q) \geq \mathcal{B}(n-1, p))^2 \\ &\leq \frac{q}{2} \mathbf{P}(Y+1 \geq X)^2. \end{aligned}$$

On a

$$\begin{aligned} \mathbf{P}(Y+1 \geq X) &= \sum_{k=-1}^{n-2} \mathbf{P}(X = k+1) \mathbf{P}(Y \geq k) \\ &= \mathbf{P}(X=0) + \sum_{k=0}^{n-2} \frac{p(n-2-k)}{(1-p)(k+1)} \mathbf{P}(X=k) \mathbf{P}(Y \geq k). \end{aligned}$$

Distinguons deux cas : si $p \geq \frac{12 \log n}{n}$, alors

$$\mathbf{P}(Y+1 \geq X) \leq (1-p)^{n-1} + 2n \mathbf{P}\left(X \leq \frac{p(n-1)}{4}\right) + 12 \mathbf{P}(Y \geq X).$$

Pour n assez grand, on a $(1-p)^{n-1} \leq 2n^{-10}$ et, par la concentration des variables binomiales 6.9,

$$\mathbf{P}\left(X \leq \frac{p(n-1)}{4}\right) = \mathbf{P}\left(X - \mathbf{E}X \leq -\frac{3}{4}p(n-1)\right) \leq e^{-\frac{(3/4)^2 p(n-1)}{2}} \leq n^{-3}.$$

Ainsi

$$\mathbf{P}(Y + 1 \geq X) \leq 12\mathbf{P}(Y \geq X) + O(n^{-2}),$$

et

$$\mathbf{P}(A_u \cap A_v \cap \{u \sim v\} \cap \{\sigma_u \neq \sigma_v\}) \leq 48\mathbf{P}(Y \geq X)^2 + O(n^{-2})$$

Si $p < 12\frac{\log n}{n}$, alors, comme $\mathbf{P}(X = 0) \leq \mathbf{P}(Y \geq X)$,

$$\mathbf{P}(Y + 1 \geq X) \leq (1 + 36 \log n) \mathbf{P}(Y \geq X),$$

et, comme $q < p = O\left(\frac{\log n}{n}\right)$,

$$\mathbf{P}(A_u \cap A_v \cap \{u \sim v\} \cap \{\sigma_u \neq \sigma_v\}) = O\left(\frac{(\log n)^3}{n} \mathbf{P}(Y \geq X)^2\right) = o(\mathbf{P}(Y \geq X)^2).$$

En combinant tout cela, on obtient qu'il existe une constante $C > 0$ telle que

$$\mathbf{E}[N^2] \leq 2n\mathbf{P}(Y \geq X) + C(2n)^2\mathbf{P}(Y \geq X)^2 + C = \mathbf{E}[N] + C\mathbf{E}[N]^2 + C,$$

ce qui implique qu'il existe $\delta > 0$ tel que, pour une infinité d'entiers n ,

$$\mathbf{P}(N \geq 1) \geq \delta. \quad \blacksquare$$

Preuve de la nécessité de (3.1). Par symétrie, le lemme 3.2 implique que

$$\mathbf{P}\left(\bigcup_{v \in V} A_v \cap \{\sigma_v = 1\}\right) = \mathbf{P}\left(\bigcup_{v \in V} A_v \cap \{\sigma_v = -1\}\right) \geq \frac{\delta}{2}.$$

En remarquant que les événements sont positivement corrélés, on a

$$\mathbf{P}\left(\bigcup_{u, v \in V} \{\sigma_u = 1, \sigma_v = -1\} \cap A_u \cap A_v\right) \geq \frac{\delta^2}{4}.$$

Or si le couple (G, σ) est tel qu'il existe deux tels sommets u et v , et si l'on considère la partition $\eta \in \mathcal{X}$ qui est égale à σ partout sauf qu'elle échange les labels de u et v : $\eta_u = -1$ et $\eta_v = 1$, alors, par la formule 3.1, il est facile de voir que

$$\mathbf{P}(G \mid \eta) \geq \mathbf{P}(G \mid \sigma).$$

Ainsi, un estimateur du maximum de vraisemblance peut se tromper sur ces deux sommets. On a donc montré qu'il existe $\varepsilon > 0$ tel que pour une infinité de n , pour tout estimateur $\hat{\eta}$,

$$\mathbf{P}(\hat{\eta} \neq \pm \sigma) \geq \varepsilon. \quad \blacksquare$$

On a donc montré que pour que la reconstruction exacte soit possible, il faut que $n\mathbf{P}(Y \geq X) \rightarrow 0$, où X et Y sont deux variables binomiales indépendantes, avec $X \sim \mathcal{B}(n-1, p)$ et

$Y \sim \mathcal{B}(n, q)$. On peut en fait montrer que cette condition est aussi suffisante pour pouvoir reconstruire exactement la partition cachée en temps polynomial!

Proposition 3.3 (Mossel et al. [18]). *La reconstruction exacte dans $\mathcal{G}(2n, p, q)$ est possible si et seulement si p et q sont tels que $n\mathbf{P}(Y \geq X) \rightarrow 0$, où X et Y sont deux variables binomiales indépendantes, avec $X \sim \mathcal{B}(n-1, p)$ et $Y \sim \mathcal{B}(n, q)$.*

3. Reconstruction faible

Proposition 3.4 (Mossel et al. [18]). *Le reconstruction faible dans $\mathcal{G}(2n, p, q)$ avec $p > q$ est possible si et seulement si p et q sont tels que $\mathbf{P}(Y \geq X) \rightarrow 0$, où X et Y sont deux variables binomiales indépendantes, avec $X \sim \mathcal{B}(n-1, p)$ et $Y \sim \mathcal{B}(n, q)$.*

Remarque 3.2. La condition $\mathbf{P}(Y \geq X) \rightarrow 0$ est équivalente à

$$\frac{\sqrt{n}(p-q)}{\sqrt{p+q}} \rightarrow +\infty.$$

Nous allons montrer une version plus faible de la Proposition 3.4. Plus précisément nous montrerons le résultat suivant.

Proposition 3.5. *Si p et q sont tels que*

$$(3.1) \quad \frac{\sqrt{n}(p-q)}{\sqrt{(p+q)\log n}} \rightarrow +\infty,$$

alors la reconstruction faible dans $\mathcal{G}(2n, p, q)$ est possible.

Donnons d'abord le principe général de la preuve, consistant à utiliser le deuxième vecteur propre de la matrice d'adjacence pour partitionner les sommets en deux groupes. Soit A la matrice d'adjacence de $G \sim \mathcal{G}(n, p, q)$. Quitte à réordonner les sommets, on a

$$\mathbf{E}A = \begin{pmatrix} p & \dots & p & q & \dots & q \\ \vdots & & \vdots & \vdots & & \vdots \\ p & \dots & p & q & \dots & q \\ q & \dots & q & p & \dots & p \\ \vdots & & \vdots & \vdots & & \vdots \\ q & \dots & q & p & \dots & p \end{pmatrix}$$

La matrice $\mathbf{E}A$ est de rang 2. Elle possède deux valeurs propres : $\lambda_1 = n(p+q)$ associé au vecteur propre $\mathbf{1}$, et $\lambda_2 = n(p-q)$ associé au vecteur propre $y = {}^t(1, \dots, 1, -1, \dots, -1)$. En particulier, la partition cachée se lit exactement sur le vecteur y . Si l'on arrive à montrer que la matrice A est suffisamment concentrée autour de $\mathbf{E}A$ (dans un sens que nous préciserons), alors on peut espérer qu'un vecteur propre pour la deuxième valeur propre de A aura des signes opposés sur les deux communautés.

Rappelons que si $M \in \mathcal{M}_{n,n}(\mathbb{R})$ est une matrice symétrique, alors ses valeurs propres sont réelles, et son rayon spectrale est égal à sa norme d'opérateur sur $(\mathbb{R}^n, \|\cdot\|)$ où $\|\cdot\| = \|\cdot\|_2$

est la norme euclidienne :

$$\|M\| = \max\{|\lambda|, \lambda \in \text{Sp}(M)\} = \sup_{x \in \mathbb{R}^n} \frac{\|Mx\|}{\|x\|}.$$

Si A et B sont deux matrices $n \times n$ symétriques, on notera $A \preceq B$ pour signifier que la matrice $B - A$ est semi-définie positive, autrement dit les valeurs propres de $B - A$ sont positives.

Avant de prouver la Proposition 3.5, nous allons démontrer un résultat de concentration pour les sommes de matrices indépendantes, qui est un équivalent matricielle de l'inégalité de Bernstein 6.8.

Proposition 3.6. *Soient X_1, \dots, X_N des matrices $n \times n$ symétriques indépendantes telles que $\mathbf{E}X_i = \mathbf{0}$ et $\|X_i\| \leq K$. Alors pour tout $t \geq 0$,*

$$\mathbf{P} \left(\left\| \sum_{i=1}^N X_i \right\| \geq t \right) \leq 2n \exp \left\{ -\frac{t^2}{2 \left(\sigma^2 + \frac{Kt}{3} \right)} \right\},$$

$$\text{où } \sigma^2 = \left\| \sum_{i=1}^N \mathbf{E}X_i^2 \right\|.$$

Preuve de la Proposition 3.6. Notons $S = \sum_{i=1}^N X_i$. Comme S est symétrique, ses valeurs propres sont réelles, et l'on peut les ordonner par ordre décroissant : $\lambda_1(S) \geq \dots \geq \lambda_n(S)$. On a alors $\|S\| = \max\{\lambda_1(S), -\lambda_n(S)\}$. Comme $-\lambda_n(S) = \lambda_1(-S)$, il suffit de montrer que

$$\mathbf{P}(\lambda_1(S) \geq t) \leq n \exp \left\{ -\frac{t^2}{2 \left(\sigma^2 + \frac{Kt}{3} \right)} \right\}.$$

On a, pour tout $u \geq 0$,

$$\mathbf{P}(\lambda_1(S) \geq t) \leq e^{-ut} \mathbf{E}e^{u\lambda_1(S)} = e^{-ut} \mathbf{E}\lambda_1(e^{uS}),$$

où la dernière inégalité vient du fait que $e^{u\lambda_1(S)} = \lambda_1(e^{uS})$, avec $e^{uS} = I + \sum_{k=1}^{+\infty} \frac{u^k S^k}{k!}$. Comme e^{uS} est définie positive, toutes ses valeurs propres sont positives et l'on a $\lambda_1(e^{uS}) \leq \text{tr}(e^{uS})$. À ce stade, on aimerait pouvoir dire $e^{uS} = \prod_{i=1}^N e^{uX_i}$, mais cette identité n'est pas vraie : une exponentielle de matrices ne transforme pas une somme en un produit. En revanche, on a l'inégalité suivante :

$$(3.2) \quad \mathbf{E} \text{tr}(e^{uS}) \leq \text{tr} \exp \left\{ \sum_{i=1}^N \log \mathbf{E}e^{uX_i} \right\},$$

où le logarithme d'une matrice est défini comme l'inverse de l'exponentielle : $\log(e^A) = A$. Pour montrer cette inégalité, on utilise un résultat puissant d'analyse matricielle (que nous ne montrerons pas ici), le théorème de Lieb : pour toute matrice $n \times n$ symétrique H , l'application

$$A \mapsto \text{tr} \exp(H + \log A)$$

est concave sur l'ensemble des matrices $n \times n$ symétriques définies positives. En utilisant ce théorème avec $H = \sum_{i=1}^{N-1} uX_i$ et $A = e^{uX_N}$, et en utilisant l'inégalité de Jensen conditionnellement à X_1, \dots, X_{N-1} , on obtient

$$\mathbf{E} [\operatorname{tr}(e^{uS}) \mid X_1, \dots, X_{N-1}] \leq \operatorname{tr} \exp \left\{ u \sum_{i=1}^{N-1} X_i + \log \mathbf{E} e^{uX_N} \right\}.$$

En prenant l'espérance conditionnelle sachant X_1, \dots, X_{N-2} , et en répétant le même argument, et ainsi de suite, on obtient bien l'inégalité (3.2). On a ainsi réussi à passer de la transformée génératrice des moments de S à celle des matrices X_i . Montrons maintenant que si X est une matrice $n \times n$ symétrique avec $\mathbf{E}X = \mathbf{0}$ et $\|X\| \leq K$, alors pour tout $0 \leq u < 3/K$,

$$\log \mathbf{E} e^{uX} \preceq \frac{u^2 \mathbf{E}X^2}{2(1 - \frac{uK}{3})}.$$

Tout d'abord, pour $0 \leq u < 3/K$ et pour x tel que $|x| \leq K$, on a

$$e^{ux} = 1 + ux + \sum_{k \geq 2} \frac{(ux)^k}{k!} \leq 1 + ux + \frac{u^2 x^2}{2} \sum_{k \geq 2} \left(\frac{uK}{3} \right)^{k-2} \leq 1 + ux + \frac{u^2 x^2}{2(1 - \frac{uK}{3})}.$$

Cela implique l'inégalité matricielle

$$(3.3) \quad e^{uX} \preceq I + uX + \frac{u^2 X^2}{2(1 - \frac{uK}{3})}.$$

En effet, comme X est symétrique, on peut écrire $X = U \Delta^t U$ où U est une matrice orthogonale et Δ est la matrice diagonale dont les coefficients diagonaux sont donnés par les valeurs propres (μ_i) . Pour une fonction $f : \mathbb{R} \rightarrow \mathbb{R}$, on peut alors définir $f(X)$ comme la matrice $U f(\Delta)^t U$, où $f(\Delta)$ est la matrice diagonale dont les coefficients sont donnés par les $(f(\mu_i))$. Ainsi, si X est telle que $\|X\| \leq K$, et si f et g sont deux fonctions de \mathbb{R} dans \mathbb{R} telle que pour tout x tel que $|x| \leq K$, $f(x) \leq g(x)$, alors toutes les valeurs propres de $g(X) - f(X)$ sont positives, i.e. $f(X) \preceq g(X)$. En prenant l'espérance dans (3.3) et en se rappelant que $\mathbf{E}X = \mathbf{0}$, on a

$$\mathbf{E} e^{uX} \preceq I + \frac{u^2 \mathbf{E}X^2}{2(1 - \frac{uK}{3})}.$$

On laisse en exercice (pas trivial du tout) de montrer que si $\mathbf{0} \prec A \preceq B$, alors $\log A \preceq \log B$. Puis, en utilisant $\log(1+z) \leq z$ pour tout $z \geq 0$ (et les mêmes arguments que ci-dessus pour passer à une inégalité matricielle), on obtient bien

$$\log \mathbf{E} e^{uX} \preceq \frac{u^2 \mathbf{E}X^2}{2(1 - \frac{uK}{3})}.$$

En revenant à (3.2), on a donc

$$\mathbf{E} \operatorname{tr}(e^{uS}) \leq \operatorname{tr} \exp \left\{ \frac{u^2 \sum_{i=1}^N \mathbf{E}X_i^2}{2(1 - \frac{uK}{3})} \right\} \leq n \exp \left\{ \frac{\sigma^2 u^2}{2(1 - \frac{uK}{3})} \right\},$$

et

$$\mathbf{P}(\lambda_1(S) \geq t) \leq n e^{-ut} \exp \left\{ \frac{\sigma^2 u^2}{2(1 - \frac{uK}{3})} \right\} = n \exp \left\{ -ut + \frac{\sigma^2 u^2}{2(1 - \frac{uK}{3})} \right\}.$$

En optimisant sur $0 \leq u < 3/K$, on voit que le membre de droit est minimal pour $u = \frac{t}{\sigma^2 + Kt/3}$, ce qui donne

$$\mathbf{P}(\lambda_1(S) \geq t) \leq n \exp \left\{ -\frac{t^2}{2(\sigma^2 + \frac{Kt}{3})} \right\}.$$

■

Preuve de la Proposition 3.5. Soit A la matrice d'adjacence de $\mathcal{G}(n, p, q)$. On peut écrire

$$A = \sum_{i < j} X_{i,j},$$

où $X_{i,j}$ est la matrice dont tous les coefficients sont nuls sauf les coefficients en (i, j) et en (j, i) qui sont égaux à 1 si $i \sim j$, et à 0 sinon. On a

$$\sigma^2 = \left\| \sum_{i < j} \mathbf{E}[(X_{i,j} - \mathbf{E}X_{i,j})^2] \right\| \leq n(p+q), \quad \text{et} \quad \|X_{i,j} - \mathbf{E}X_{i,j}\| \leq 2.$$

Ainsi, par la Proposition 3.6, pour tout $t \geq 0$,

$$\mathbf{P}(\|A - \mathbf{E}A\| \geq t) \leq 4n \exp \left\{ -\frac{t^2}{2(\sigma^2 + \frac{2t}{3})} \right\}.$$

Ainsi, pour tout $\varepsilon > 0$,

$$\mathbf{P} \left(\frac{\|A - \mathbf{E}A\|}{n(p-q)} \geq \varepsilon \right) \leq 4n \exp \left\{ -\frac{\varepsilon^2 n(p-q)^2}{2(p+q + \frac{2\varepsilon}{3}(p-q))} \right\} = o(1),$$

puisque, par l'hypothèse (3.1), $\frac{n(p-q)^2}{p+q} \gg \log n$. Autrement dit

$$\frac{\|A - \mathbf{E}A\|}{n(p-q)} \xrightarrow{\mathbf{P}} 0.$$

Considérons le vecteur x défini par

$$\hat{x} = \arg \max \{ {}^t x A x, x \perp \mathbf{1}, \|x\| = 1 \}.$$

On peut décomposer \hat{x} en

$$\hat{x} = \alpha y + \beta z,$$

où $y = \frac{1}{\sqrt{2n}}(1, \dots, 1, -1, \dots, -1) = \frac{\sigma}{\sqrt{2n}}$, z est un vecteur de norme 1 orthogonal à $\mathbf{1}$ et à y , et où $\alpha^2 + \beta^2 = 1$. On a

$$\begin{aligned} {}^t \hat{x} A \hat{x} &= {}^t \hat{x}(\mathbf{E}A)\hat{x} + {}^t \hat{x}(A - \mathbf{E}A)\hat{x} \\ &= \alpha^2 {}^t y(\mathbf{E}A)y + {}^t \hat{x}(A - \mathbf{E}A)\hat{x} \\ &= (\alpha^2 + o_{\mathbf{P}}(1))n(p-q), \end{aligned}$$

où l'on a utilisé que $|{}^t \hat{x}(A - \mathbf{E}A)\hat{x}| \leq \|A - \mathbf{E}A\| = o_{\mathbf{P}}(n(p-q))$. Cela implique que $|\alpha| = 1 - o_{\mathbf{P}}(1)$. Ainsi, si l'on considère l'estimateur $\hat{\sigma}$ donné par

$$\forall v \in V, \hat{\sigma}_v = \text{sign}(\hat{x}_v),$$

alors on a

$$\frac{|\sum_{v \in V} \hat{\sigma}_v \sigma_v|}{2n} \xrightarrow{\mathbf{P}} 1.$$

En effet, supposons que $\alpha \geq 0$ (le cas $\alpha \leq 0$ se traite de la même façon). On a

$$\begin{aligned} \sum_{v \in V} \mathbb{1}_{\{\hat{\sigma}_v \neq \sigma_v\}} &= \sum_{v \in V} \mathbb{1}_{\{\text{sign}(\alpha \sigma_v + \beta \sqrt{2n} z_v) \neq \sigma_v\}} \\ &\leq \sum_{v \in V} \mathbb{1}_{\{|\beta \sqrt{2n} z_v| > \frac{\alpha}{2}\}} \\ &\leq \sum_{v \in V} \mathbb{1}_{\{|z_v|^2 > \frac{\alpha^2}{8\beta^2 n}\}} \\ &\leq \frac{8\beta^2 n}{\alpha^2} = o_{\mathbf{P}}(n), \end{aligned}$$

où l'on a utilisé le fait que $\sum_{v \in V} |z_v|^2 = 1$ et que $\beta^2 = 1 - \alpha^2 = o_{\mathbf{P}}(1)$. ■

4. Détection de communautés

Dans cette section, on s'intéresse au modèle $G \sim \mathcal{G}'(n, p, q)$ et l'on se place dans le cas « sparse », où $p = \frac{a}{n}$ et $q = \frac{b}{n}$ avec $a, b > 0$ des constantes fixées. Dans ce régime, il est clair que l'objectif de reconstruction exacte peut être abandonné. La question est plutôt celle de la détection : peut-on trouver un estimateur $\hat{\sigma}$ qui soit positivement corrélé avec σ , i.e. algorithme qui fasse strictement mieux que de prédire à pile-ou-face le label de chaque sommet ? Au cours des dix dernières années, de nombreux travaux de recherche ont été consacrés à cette question difficile, pour aboutir à la preuve d'une conjecture fascinante qui avait été faite par des physiciens [9] : la détection est possible en un temps polynomial si et seulement si

$$(3.1) \quad \frac{(a-b)^2}{2} > a+b.$$

Mossel et al. [19] ont d'abord montré le résultat d'impossibilité : si $\frac{(a-b)^2}{2} \leq a+b$, alors la détection est impossible. Puis, Mossel et al. [20] et Massoulié [16] ont indépendamment la réciproque : si $\frac{(a-b)^2}{2} > a+b$, alors on peut trouver un algorithme polynomial pour détecter les communautés.

Nous ne donnerons pas ici la preuve de ces résultats difficiles. En revanche, nous tenterons d'expliquer l'apparition du seuil (3.1), appelé seuil de Kesten–Stigum. Pour cela, nous considérerons un problème proche mais relativement plus simple à étudier, celui de la transmission de messages sur les arbres, et verrons comment les deux modèles sont liés.

4.1. Transmission d'information sur un arbre. Soit T un arbre enraciné en ρ et $\varepsilon > 0$. Le processus de transmission (*broadcast process*) sur T est défini de la façon suivante. La racine ρ reçoit une information initiale sous forme d'un bit aléatoire $\sigma_\rho \sim \mathcal{B}(1/2)$. Ce bit est alors propagé dans l'arbre, avec erreurs : chaque sommet reçoit indépendamment le bit de son parent avec probabilité $1 - \varepsilon$ et le bit opposé avec probabilité ε . De façon équivalente, pour $\eta = 2\varepsilon$, le bit du parent est transmis sans erreur à son enfant avec probabilité $1 - \eta$, et, avec probabilité η , le bit est rejoué à pile-ou-face. Si l'on note T_n l'ensemble des sommets à

distance n de la racine, on peut se poser la question suivante : si l'on connaît uniquement la valeurs des bits sur T_n , peut-on retrouver le bit en ρ avec probabilité de succès strictement supérieure à $1/2$? (Notons que l'on peut toujours atteindre une probabilité de succès de $1/2$ en tirant à pile-ou-face).

Supposons que chaque sommet de T a le même nombre d d'enfants (la racine est de degré d , et tous les autres sommets sont de degré $d + 1$). Dans ce cas $|T_n| = d^n$. On dit qu'un sommet $u \in T_n$ est fidèle si, tout le long du chemin de ρ à u , le bit a toujours été transmis sans erreur. La probabilité qu'un sommet soit fidèle est $(1 - \eta)^n$ et il y a ainsi en moyenne $d^n(1 - \eta)^n$ sommets fidèles en T_n . Remarquons aussi que le bit des sommets infidèles est distribué selon une Bernoulli $\mathcal{B}(1/2)$. Si l'on omet les corrélations en prétendant que ces bits sont indépendants, alors le théorème central-limite nous dit que, parmi ces sommets infidèles, il y a environ une moitié de sommets de chaque bit, avec des fluctuations d'ordre $\sqrt{d^n}$. Si l'on choisit d'estimer le bit en ρ par le bit majoritaire sur T_n , il faut, pour garantir une probabilité de succès supérieure à $1/2$, qu'il y ait un nombre suffisamment grand de sommets fidèles pour pouvoir surpasser les fluctuations des sommets infidèles. Plus précisément, il faut que

$$d^n(1 - \eta)^n \gg \sqrt{d^n},$$

soit $d(1 - \eta)^2 > 1$.

L'analogie avec $G \sim \mathcal{G}'(n, \frac{a}{n}, \frac{b}{n})$ peut se voir en interprétant d comme le degré moyen dans G , i.e. $d = \frac{a+b}{2}$, et ε comme la probabilité que $\sigma_u \neq \sigma_v$ sachant que u et v sont voisins dans G , i.e.

$$\varepsilon = \mathbf{P}(\sigma_u \neq \sigma_v \mid u \sim v) = \frac{\mathbf{P}(\sigma_u \neq \sigma_v)\mathbf{P}(u \sim v \mid \sigma_u \neq \sigma_v)}{\mathbf{P}(u \sim v)} = \frac{\frac{1}{2} \cdot \frac{b}{n}}{\frac{1}{2} \cdot \frac{a}{n} + \frac{1}{2} \cdot \frac{b}{n}} = \frac{b}{a + b}.$$

Avec ces identifications, le seuil $d(1 - \eta)^2 > 1$ se réécrit $(a - b)^2 > 2(a + b)$, qui correspond bien au seuil (3.1).

4.2. Test de la présence de communautés. Soit $a > b > 0$ deux constantes fixées. Notons P la loi de d'un graphe aléatoire d'Erdős-Rényi $\mathcal{G}(n, \frac{a+b}{2n})$ et Q la loi $\mathcal{G}'(n, \frac{a}{n}, \frac{b}{n})$. On souhaite tester

$$H_0 : G \sim P \quad \text{contre} \quad H_1 : G \sim Q,$$

à partir d'une seule réalisation de G . Notons \mathcal{X} l'ensemble des graphes à n sommets. On définit le risque d'un test $T : \mathcal{X} \rightarrow \{0, 1\}$ par

$$\mathbf{R}(T) = \frac{1}{2} (\mathbf{P}_0(T(G) = 1) + \mathbf{P}_1(T(G) = 0)).$$

Il s'agit d'un problème de test de type bayésien : on tire à pile-ou-face pour choisir entre P et Q , puis on tire G selon la loi choisie. Le problème est alors de retrouver la loi qui a généré G à partir de la seule observation de G . Le risque $\mathbf{R}(T)$ correspond alors précisément à la

probabilité que le test T se trompe.

$$\begin{aligned} \mathbf{R}(T) &= \frac{1}{2} \sum_{g \in \mathcal{X}} (P(g) \mathbb{1}_{T(g)=1} + Q(g) \mathbb{1}_{T(g)=0}) \\ &= \frac{1}{2} \left\{ 1 + \sum_{g \in \mathcal{X}} (P(g) - Q(g)) \mathbb{1}_{T(g)=1} \right\} \\ &= \frac{1}{2} (1 - \{\mathbf{E}_Q T - \mathbf{E}_P T\}) \\ &\geq \frac{1}{2} (1 - \|P - Q\|_{\text{TV}}), \end{aligned}$$

où $\|P - Q\|_{\text{TV}}$ est la distance en variation totale entre P et Q donnée par

$$\|P - Q\|_{\text{TV}} = \max_{A \subset \mathcal{X}} |P(A) - Q(A)| = \max_{f: \|f\|_{\infty} \leq 1} |\mathbf{E}_P f - \mathbf{E}_Q f| = \frac{1}{2} \sum_{g \in \mathcal{X}} |P(g) - Q(g)|.$$

Inversement, on vérifie aisément que la borne inférieure $\frac{1}{2} (1 - \|P - Q\|_{\text{TV}})$ est atteinte par le test T^* du maximum de vraisemblance, défini par

$$T^*(g) = \mathbb{1}_{\{P(g) \leq Q(g)\}}.$$

Ainsi, le risque minimal est

$$\mathbf{R}_* = \mathbf{R}(T^*) = \frac{1}{2} (1 - \|P - Q\|_{\text{TV}}).$$

En particulier, on a toujours $\mathbf{R}_* \leq 1/2$: le test qui consiste à répondre 0 ou 1 avec probabilité 1/2, quel que soit G , se trompe avec probabilité 1/2. Peut-on faire strictement mieux que 1/2? La réponse est oui, dès que $a > b$, par comparaison du nombre de triangles dans chacun des deux modèles. En effet, le nombre de triangles dans $\mathcal{G}(n, \frac{a+b}{2n})$ a pour espérance $\binom{n}{3} \left(\frac{a+b}{2n}\right)^3 \sim \frac{(a+b)^3}{48}$, et l'on peut montrer que ce nombre suit approximativement (au sens où la distance en variation totale tend vers 0) une loi de Poisson avec la même espérance. Dans $\mathcal{G}'(n, \frac{a}{n}, \frac{b}{n})$, le nombre de triangles a pour espérance

$$\binom{n}{3} \left(\frac{a^3}{4n^3} + \frac{3ab^2}{4n^3} \right) \sim \frac{a^3}{24} + \frac{ab^2}{8} = \frac{(a+b)^3}{48} \left(1 + \left(\frac{a-b}{a+b} \right)^3 \right),$$

et suit approximativement une loi de Poisson avec cette espérance. Ainsi

$$\liminf_{n \rightarrow +\infty} \|P - Q\|_{\text{TV}} \geq \left\| \mathcal{P} \left(\frac{(a+b)^3}{48} \right) - \mathcal{P} \left(\frac{a^3}{24} + \frac{ab^2}{8} \right) \right\|_{\text{TV}} > 0,$$

et $\limsup_{n \rightarrow +\infty} \mathbf{R}_* < \frac{1}{2}$.

Une question bien plus difficile est de savoir sous quelles conditions sur a et b on peut garantir que $\lim_{n \rightarrow \infty} \mathbf{R}_* = 0$ (voir [19]).

Graphes géométriques aléatoires : tests et estimation de la dimension

Dans de nombreux réseaux réels, il est raisonnable de supposer que deux individus ont plus de chances d'être reliés par une arête s'ils sont « proches » dans un certain espace métrique. Par exemple, dans un réseau social, on peut penser que deux personnes ont plus de chances d'être amies si elles sont géographiquement proches, ou de façon plus générale, si elles sont proches dans l'espace social (prenant en compte non seulement la position géographique mais la position dans l'échelle sociale, le niveau d'éducation, etc). Les graphes aléatoires géométriques permettent de modéliser cela en supposant l'existence d'une géométrie sous-jacente, mais non-observée. La question est alors de savoir si l'on peut détecter la présence de cette géométrie (problème de test), et si oui, si l'on peut estimer certains aspects de cette géométrie, par exemple la dimension de l'espace sous-jacent.

1. Un modèle simple de graphe géométrique aléatoire

Dans un graphe géométrique, chaque sommet est associé à un point dans un espace métrique, et une arête est présente entre deux sommets si la distance entre les points correspondants est inférieure à un seuil fixé. Ici, on considère le cas où l'espace métrique sous-jacent est la sphère euclidienne

$$\mathbb{S}^{d-1} = \left\{ x \in \mathbb{R}^d, \|x\| = 1 \right\},$$

où $d \in \mathbb{N}^*$ et $\|\cdot\|$ est la norme euclidienne. Le graphe aléatoire géométrique $G \sim \mathcal{G}(n, p, d)$ est généré de la façon suivante : n points X_1, \dots, X_n sont tirés indépendamment et uniformément sur la sphère \mathbb{S}^{d-1} , et pour $1 \leq i \neq j \leq n$, les sommets i et j sont reliés si et seulement si

$$\langle X_i, X_j \rangle \geq t_{p,d},$$

où $t_{p,d} \in [-1, 1]$ est tel que $\mathbf{P}(\langle X_i, X_j \rangle \geq t_{p,d}) = p$. Par exemple, si $p = 1/2$, $t_{1/2,d} = 0$. De façon équivalente, les sommets i et j sont connectés si et seulement si

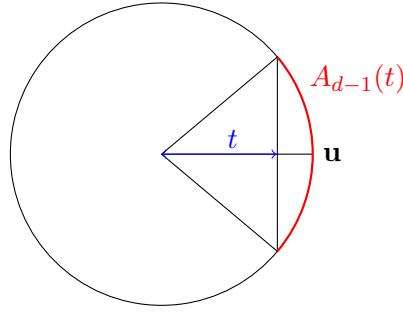
$$\|X_i - X_j\| \leq \sqrt{2(1 - t_{p,d})}.$$

Pour comprendre la façon dont $t_{p,d}$ dépend de p , remarquons que

$$p = \mathbf{P}(\langle X_i, \mathbf{u} \rangle \geq t_{p,d}),$$

où $\mathbf{u} \in \mathbb{S}^{d-1}$ est le vecteur $(1, 0, \dots, 0)$. Si l'on note μ_{d-1} la mesure uniforme sur \mathbb{S}^{d-1} , et $A_{d-1}(t) = \{x \in \mathbb{S}^{d-1}, \langle x, \mathbf{u} \rangle \geq t\}$, alors

$$p = \mu_{d-1}(A_{d-1}(t_{p,d})).$$



2. Détection de la géométrie

Si l'on souhaite tester la présence d'une structure géométrique sous-jacente, le graphe le plus naturel avec lequel comparer $\mathcal{G}(n, p, d)$ est $\mathcal{G}(n, p)$, le graphe d'Erdős–Renyi, dans lequel deux sommets sont indépendamment connectés avec probabilité p . On peut alors formuler le problème comme un simple test d'hypothèses :

$$H_0 : G \sim \mathcal{G}(n, p) \quad \text{contre} \quad H_1 : G \sim \mathcal{G}(n, p, d).$$

Comme dans la Section 4.2, en notant \mathcal{X} l'ensemble des graphes à n sommets, on définit le risque d'un test $T : \mathcal{X} \rightarrow \{0, 1\}$ par

$$\mathbf{R}(T) = \frac{1}{2} (\mathbf{P}_0(T(G) = 1) + \mathbf{P}_1(T(G) = 0)).$$

Le risque minimal $\mathbf{R}_\star = \inf_T \mathbf{R}(T)$ est alors donné par

$$\mathbf{R}_\star = \frac{1}{2} (1 - \|\mathbf{P}_0(G \in \cdot) - \mathbf{P}_1(G \in \cdot)\|_{\text{TV}}),$$

et le test optimal atteignant ce risque est le test du maximum de vraisemblance

$$T^\star(g) = \mathbb{1}_{\{\mathbf{P}_0(G=g) \leq \mathbf{P}_1(G=g)\}}.$$

Ainsi, comprendre le risque optimal revient à comprendre la distance en variation totale entre les deux lois. Si $\mathcal{G}(n, p)$ et $\mathcal{G}(n, p, d)$ sont asymptotiquement indistinguables en variation totale, au sens où $\|\mathbf{P}_0(G \in \cdot) - \mathbf{P}_1(G \in \cdot)\|_{\text{TV}} \rightarrow 0$ quand $n \rightarrow +\infty$, alors $\mathbf{R}_\star \rightarrow 1/2$: aucun test ne peut faire mieux que celui qui tire à pile-ou-face pour décider. Inversement, si l'on peut parfaitement distinguer les deux lois au sens où $\|\mathbf{P}_0(G \in \cdot) - \mathbf{P}_1(G \in \cdot)\|_{\text{TV}} \rightarrow 1$ quand $n \rightarrow +\infty$, alors $\mathbf{R}_\star \rightarrow 0$: avec probabilité qui tend vers 1, le test optimal ne se trompe pas. Remarquons que dans ce cas, la question n'est pas complètement réglée, car le test du maximum de vraisemblance peut être extrêmement long à déterminer (car $\mathbf{P}_1(G = g)$ est une intégrale compliquée). Il s'agit alors de savoir si l'on peut construire un test calculable en temps polynomial et dont le risque tende toujours vers 0 quand $n \rightarrow +\infty$.

Intuitivement, plus la dimension de la sphère sous-jacente est grande, plus il est difficile de la détecter, tandis qu'une géométrie de petite dimension aura un fort effet sur le graphe et sera potentiellement détectable. Jusqu'à quelle dimension peut-on espérer pouvoir détecter la géométrie ? Le théorème ci-dessous, dû à Bubeck et al. [5], montre que la dimension critique est $d \sim n^3$.

Theorème 4.1 ([5]). *Soit \mathbf{R}_\star le risque optimal du problème de test de $\mathcal{G}(n, p)$ contre $\mathcal{G}(n, p, d)$, avec $p \in]0, 1[$ fixé. Alors on a*

$$\lim_{n \rightarrow +\infty} \mathbf{R}_\star = \begin{cases} 0 & \text{si } d \ll n^3, \\ \frac{1}{2} & \text{si } d \gg n^3. \end{cases}$$

Pour la première moitié de ce théorème ($\mathbf{R}_\star \rightarrow 0$ quand $d \ll n^3$), il suffit d'exhiber un test dont le risque tende vers 0. Nous verrons en Section 2.1 que compter le nombre de triangles (ce qui se fait rapidement) dans le graphe permet de bien distinguer les deux lois, pourvu que $d \ll n^3$. La deuxième partie du théorème vient d'une borne supérieure sur la distance en variation totale entre les lois des matrices d'adjacence dans les deux hypothèses, établie en Section 2.2.

2.1. Le test des triangles. Intuitivement, on s'attend à ce que le nombre de triangles soit plus élevé dans un graphe géométrique que dans un graphe d'Erdős–Renyi. En effet, dans $\mathcal{G}(n, p)$, la probabilité que trois points forme un triangle est égale à p^3 , alors que dans $\mathcal{G}(n, p, d)$, sachant que deux sommets u et v sont connectés, la probabilité qu'un troisième sommet w soit connecté à la fois à u et à v est plus grande que p^2 , puisque si w est proche de u dans la sphère sous-jacente, alors, comme u est proche de v , il est probable que w soit proche de v aussi.

En vérité, comparer simplement le nombre de triangles ne permet pas d'aller jusqu'au seuil $d \ll n^3$ (la variance est trop grande). Pour réduire la variance, l'astuce est de considérer comme statistique de test une version signée du nombre de triangles. Plus précisément, soit

$$\Delta = \sum_{1 \leq i < j < k \leq n} (A_{i,j} - p)(A_{j,k} - p)(A_{k,i} - p),$$

où A est la matrice d'adjacence de G . Sous H_0 , on a $\mathbf{E}_0 \Delta = 0$, et

$$\begin{aligned} \text{Var}_0(\Delta) &= \binom{n}{3} \mathbf{E} [(A_{i,j} - p)^2 (A_{j,k} - p)^2 (A_{k,i} - p)^2] \\ &= \binom{n}{3} p^3 (1 - p)^3. \end{aligned}$$

D'autre part, sous H_1 ,

$$\begin{aligned} \mathbf{E}_1 \Delta &= \binom{n}{3} (\mathbf{P}_1(A_{1,2}A_{2,3}A_{3,1} = 1) - 3p\mathbf{P}_1(A_{1,2}A_{1,3} = 1) + 3p^2\mathbf{P}_1(A_{1,2} = 1) - p^3) \\ &= \binom{n}{3} (\mathbf{P}_1(A_{1,2}A_{2,3}A_{3,1} = 1) - p^3) \end{aligned}$$

où l'on a utilisé $\mathbf{P}_1(A_{1,2} = 1) = p$ et $\mathbf{P}_1(A_{1,2}A_{1,3} = 1) = p^2$ (les événements $\langle X_1, X_2 \rangle \geq t_{p,d}$ et $\langle X_1, X_3 \rangle \geq t_{p,d}$ sont indépendants). Notons pour alléger les notations $t = t_{p,d}$. Par invariance par rotation, on a

$$\begin{aligned} \mathbf{P}_1(A_{1,2}A_{2,3}A_{1,3} = 1) &= \mathbf{P}(\langle X_1, X_2 \rangle \geq t, \langle X_2, X_3 \rangle \geq t, \langle X_1, X_3 \rangle \geq t) \\ &= \mathbf{P}(\langle X, \mathbf{u} \rangle \geq t, \langle Y, \mathbf{u} \rangle \geq t, \langle X, Y \rangle \geq t), \end{aligned}$$

où X et Y sont deux indépendants uniformes sur \mathbb{S}^{d-1} et où $\mathbf{u} = (1, 0, \dots, 0)$. La preuve est un peu technique mais l'on peut montrer (voir [5]) que pour tout $p \in]0, 1[$, il existe c_p tel que

$$\mathbf{P}(\langle X, \mathbf{u} \rangle \geq t, \langle Y, \mathbf{u} \rangle \geq t, \langle X, Y \rangle \geq t) \geq p^3 + \frac{c_p}{\sqrt{d}},$$

ce qui implique que

$$\mathbf{E}_1 \Delta \geq \frac{c_p \binom{n}{3}}{\sqrt{d}}.$$

D'autre part, on a (voir encore [5])

$$\text{Var}_1(\Delta) \leq n^3 + \frac{3n^4}{d}.$$

Considérons le test T donné par

$$T(g) = \mathbb{1} \left\{ \Delta(g) \geq \frac{c_p \binom{n}{3}}{2\sqrt{d}} \right\}.$$

Par l'inégalité de Chebyshev, on a, lorsque $d \ll n^3$,

$$\mathbf{P}_0(T(G) = 1) = \mathbf{P}_0 \left(\Delta \geq \frac{c_p \binom{n}{3}}{2\sqrt{d}} \right) \leq \frac{4d \text{Var}_0(\Delta)}{c_p^2 \binom{n}{3}} \leq \frac{c'_p d}{n^3} = o(1),$$

et

$$\begin{aligned} \mathbf{P}_1(T(G) = 0) &= \mathbf{P}_1 \left(\Delta < \frac{c_p \binom{n}{3}}{2\sqrt{d}} \right) = \mathbf{P}_0 \left(\Delta - \mathbf{E}_0 \Delta < -\frac{c_p \binom{n}{3}}{2\sqrt{d}} \right) \\ &\leq \frac{4d \text{Var}_1(\Delta)}{c_p^2 \binom{n}{3}} \leq \frac{c''_p (d+n)}{n^3} = o(1). \end{aligned}$$

Ainsi $\mathbf{R}(T) \rightarrow 0$ et $\mathbf{R}_\star \rightarrow 0$.

2.2. Matrices de Wishart et matrices du GOE. Pour montrer que $\mathbf{R}_\star \rightarrow 1$ (ou de façon équivalente que $\|\mathbf{P}_0(G \in \cdot) - \mathbf{P}_1(G \in \cdot)\|_{\text{TV}} \rightarrow 0$) quand $d \gg n^3$, l'idée est de remarquer que les matrices d'adjacence de $\mathcal{G}(n, p, d)$ et $\mathcal{G}(n, p)$ peuvent s'écrire comme des fonctions de matrices aléatoires dont la loi est bien étudiée. Tout d'abord, remarquons qu'un vecteur X uniformément distribué sur la sphère \mathbb{S}^{d-1} peut s'écrire comme

$$X = \frac{Y}{\|Y\|},$$

où Y est un vecteur gaussien standard $\mathcal{N}_d(0, I_d)$. Une matrice de Wishart W est une matrice aléatoire $n \times n$ dont les entrées sont données par

$$W_{i,j} = \langle Y_i, Y_j \rangle,$$

où Y_1, \dots, Y_n sont des vecteurs gaussiens $\mathcal{N}_d(0, I_d)$ indépendants. Ainsi, comme $W_{i,i} = \|Y_i\|^2$, la matrice d'adjacence A de $\mathcal{G}(n, p, d)$ est donnée par

$$A_{i,j} = \begin{cases} 1 & \text{si } \frac{W_{i,j}}{\sqrt{W_{i,i}W_{j,j}}} \geq t_{p,d} \text{ et } i \neq j, \\ 0 & \text{sinon.} \end{cases}$$

Notons φ la fonction qui envoie W sur A . La matrice d'adjacence E du graphe $\mathcal{G}(n, p)$ peut elle se voir comme fonction d'une matrice du GOE (*Gaussian Orthogonal Ensemble*), c'est-à-dire d'une matrice symétrique \widetilde{M} dont les entrées diagonales sont i.i.d de loi $\mathcal{N}(0, 2)$ et dont les entrées au-dessus de la diagonale sont i.i.d. $\mathcal{N}(0, 1)$, avec les entrées sur et au-dessus de la diagonale toutes indépendantes. Alors la matrice d'adjacence E de $\mathcal{G}(n, p)$ peut s'écrire

$$E_{i,j} = \begin{cases} 1 & \text{si } \widetilde{M}_{i,j} \geq \overline{\Phi}^{-1}(p) \text{ et } i \neq j, \\ 0 & \text{sinon,} \end{cases}$$

où $\overline{\Phi}^{-1}(p)$ est le quantile d'ordre $1 - p$ de la loi $\mathcal{N}(0, 1)$. Comme E ne dépend que des entrées non-diagonales de \widetilde{M} , on peut écrire

$$E_{i,j} = \begin{cases} 1 & \text{si } M_{i,j} \geq \sqrt{d}\overline{\Phi}^{-1}(p) \text{ et } i \neq j, \\ 0 & \text{sinon,} \end{cases}$$

où $M = \sqrt{d}\widetilde{M} + dI_n$. Notons ψ la fonction qui envoie M sur E . On a

$$\begin{aligned} \|\mathbf{P}_0(G \in \cdot) - \mathbf{P}_1(G \in \cdot)\|_{\text{TV}} &= \|A - E\|_{\text{TV}} = \|\varphi(W) - \psi(M)\|_{\text{TV}} \\ &\leq \|\varphi(W) - \varphi(M)\|_{\text{TV}} + \|\varphi(M) - \psi(M)\|_{\text{TV}} \\ &\leq \|W - M\|_{\text{TV}} + \|\varphi(M) - \psi(M)\|_{\text{TV}}. \end{aligned}$$

Les fonctions φ et ψ sont très proches, si bien que le deuxième terme ci-dessus est petit. En fait, dans le cas $p = 1/2$, on a même $\varphi = \psi$. En effet dans ce cas, on a $t_{1/2,d} = \overline{\Phi}^{-1}(1/2) = 0$, et

$$A_{i,j} = \begin{cases} 1 & \text{si } W_{i,j} \geq 0 \text{ et } i \neq j, \\ 0 & \text{sinon.} \end{cases} \quad \text{et} \quad E_{i,j} = \begin{cases} 1 & \text{si } M_{i,j} \geq 0 \text{ et } i \neq j, \\ 0 & \text{sinon,} \end{cases}$$

Admettons que $\|\varphi(M) - \psi(M)\|_{\text{TV}} \rightarrow 0$ quand $d \gg n^3$ (voir [5] pour la preuve), et montrons que, quand $d \gg n^3$ alors une matrice de Wishart devient indistinguable d'une matrice du GOE, au sens où $\|W - M\|_{\text{TV}} \rightarrow 0$.

La densité d'une matrice du GOE contre la mesure de Lebesgue sur l'ensemble des matrices symétriques, notée λ , est donnée par

$$A \mapsto \frac{1}{(2\pi)^{\frac{n(n+1)}{4}} 2^{\frac{n}{2}}} \exp\left\{-\frac{\text{tr}(A^2)}{4}\right\},$$

et ainsi la densité de M contre λ est donnée par

$$f_M(A) = \frac{1}{(2\pi d)^{\frac{n(n+1)}{4}} 2^{\frac{n}{2}}} \exp\left\{-\frac{1}{4d} \text{tr}((A - dI_n)^2)\right\}.$$

D'autre part, pour $d \geq n$, la densité de W contre λ est donnée par

$$f_W(A) = \frac{(\det(A))^{\frac{d-n-1}{2}}}{2^{\frac{dn}{2}} \pi^{\frac{n(n-1)}{4}} \prod_{i=1}^n \Gamma\left(\frac{d+1-i}{2}\right)} \exp\left\{-\frac{\text{tr}(A)}{2}\right\} \mathbb{1}_{\{A \succeq 0\}},$$

où $A \succcurlyeq 0$ signifie que A est semi-définie positive. On a

$$\begin{aligned} \|W - M\|_{\text{TV}} &= \int (f_W(A) - f_M(A))_+ d\lambda(A) = \int \left(\frac{f_W(A)}{f_M(A)} - 1 \right)_+ f_M(A) d\lambda(A) \\ &= \mathbf{E} \left[\left(\frac{f_W(M)}{f_M(M)} - 1 \right)_+ \right]. \end{aligned}$$

Pour montrer que cette espérance tend vers 0, il suffit de montrer que

$$\frac{f_W(M)}{f_M(M)} \xrightarrow{\mathbf{P}} 1.$$

Tout d'abord, énonçons un résultat sur le spectre des matrices aléatoires du GOE, qui vaut en fait plus généralement pour les matrices de Wigner, voir par exemple [1].

Théorème 4.2 (Théorème de Wigner). *Soit \widetilde{M} une matrice $n \times n$ du GOE. Soient $\lambda_1 \leq \dots \leq \lambda_n$ les valeurs propres (réelles) de \widetilde{M} et*

$$\mu_n = \frac{1}{n} \sum_{i=1}^n \delta_{\lambda_i/\sqrt{n}}$$

la mesure spectrale empirique de $\frac{1}{\sqrt{n}}\widetilde{M}$. Alors μ_n converge faiblement en probabilité vers la loi du demi-cercle, de densité

$$\sigma(x) = \frac{1}{2\pi} \sqrt{4 - x^2} \mathbb{1}_{\{|x| \leq 2\}}.$$

Plus précisément, pour toute fonction f continue bornée sur \mathbb{R} et pour tout $\varepsilon > 0$, on a

$$\mathbf{P} \left(\left| \frac{1}{n} \sum_{i=1}^n f \left(\frac{\lambda_i}{\sqrt{n}} \right) - \int_{\mathbb{R}} f(x) \sigma(x) dx \right| > \varepsilon \right) \rightarrow 0.$$

De plus,

$$\mathbf{P} (\forall i \in \llbracket 1, n \rrbracket, \lambda_i \in [-2\sqrt{n}, 2\sqrt{n}]) \xrightarrow{n \rightarrow \infty} 1,$$

et pour tout $k \in \mathbb{N}^*$, il existe $\sigma_k^2 > 0$ tel que

$$n^{-k/2} \sum_{i=1}^n \left(\lambda_i^k - \mathbf{E} \lambda_i^k \right) \xrightarrow{\mathcal{L}} \mathcal{N}(0, \sigma_k^2).$$

Notons tout d'abord que, par le Théorème 4.2, avec probabilité $1 - o(1)$, toutes les valeurs propres λ_i de \widetilde{M} sont plus grandes que $-2\sqrt{n}$, et ainsi toutes les valeurs propres de $M = \sqrt{d}\widetilde{M} + dI_n$ sont plus grandes que $d - 2\sqrt{dn}$. Comme $d \gg n^3$, avec probabilité $1 - o(1)$ toutes les valeurs propres de M sont positives et la matrice M est définie positive. En utilisant que $\det(M) = \prod_{i=1}^n (\sqrt{d}\lambda_i + d)$ et que $\text{tr}(M) = \sum_{i=1}^n (\sqrt{d}\lambda_i + d)$, on a

$$\begin{aligned} \log \left(\frac{f_W(M)}{f_M(M)} \right) &= \frac{1}{2} \sum_{i=1}^n \left\{ (d - n - 1) \log \left(d + \sqrt{d}\lambda_i \right) - d - \sqrt{d}\lambda_i + \frac{\lambda_i^2}{2} \right\} - \sum_{i=1}^n \log \Gamma \left(\frac{d + 1 - i}{2} \right) \\ &\quad + \left(\frac{n(n+3)}{4} - \frac{dn}{2} \right) \log(2) + \frac{n}{2} \log(\pi) + \frac{n(n+1)}{4} \log(d) + o_{\mathbf{P}}(1). \end{aligned}$$

En utilisant la formule de Stirling $\Gamma(z) = (z - \frac{1}{2}) \log z - z + \frac{\log(2\pi)}{2} + O(\frac{1}{z})$, en développant les log et en utilisant $\frac{n^3}{d} = o(1)$, on obtient

$$\begin{aligned} \log \left(\frac{f_W(M)}{f_M(M)} \right) &= \frac{1}{2} \sum_{i=1}^n \left\{ (d - n - 1) \log \left(d + \sqrt{d} \lambda_i \right) - \sqrt{d} \lambda_i + \frac{\lambda_i^2}{2} \right\} \\ &\quad - \sum_{i=1}^n \frac{d-i}{2} \log(d+1-i) + \frac{n(n+1)}{4} \log(d) - \frac{n(n-1)}{4} + o_{\mathbf{P}}(1) \\ &= \frac{1}{2} \sum_{i=1}^n \left\{ (d - n - 1) \log \left(1 + \frac{\lambda_i}{\sqrt{d}} \right) - \sqrt{d} \lambda_i + \frac{\lambda_i^2}{2} \right\} \\ &\quad - \sum_{i=1}^n \frac{d-i}{2} \log \left(1 - \frac{i-1}{d} \right) - \frac{n(n-1)}{4} + o_{\mathbf{P}}(1) \\ &= \frac{1}{2} \sum_{i=1}^n \left\{ (d - n - 1) \log \left(1 + \frac{\lambda_i}{\sqrt{d}} \right) - \sqrt{d} \lambda_i + \frac{\lambda_i^2}{2} \right\} + o_{\mathbf{P}}(1). \end{aligned}$$

Pour conclure, un développement de Taylor du log jusqu'à l'ordre 3, et le fait qu'avec probabilité $1 - o(1)$, toutes les valeurs propres λ_i vérifient $|\lambda_i| \leq 2\sqrt{n}$ donnent

$$\log \left(\frac{f_W(M)}{f_M(M)} \right) = \frac{1}{6\sqrt{d}} \sum_{i=1}^n \lambda_i^3 - \frac{n+1}{2\sqrt{d}} \sum_{i=1}^n \lambda_i + o_{\mathbf{P}}(1).$$

Or, par le Théorème 4.2, les variables aléatoires

$$\frac{1}{\sqrt{n}} \sum_{i=1}^n \lambda_i \quad \text{et} \quad \frac{1}{n^{3/2}} \sum_{i=1}^n \lambda_i^3$$

convergent toutes les deux vers des variables gaussiennes. Ainsi, en utilisant encore que $\frac{n^3}{d} = o(1)$, on obtient bien

$$\log \left(\frac{f_W(M)}{f_M(M)} \right) = o_{\mathbf{P}}(1),$$

ce qui implique $\|W - M\|_{\text{TV}} = o(1)$.

Archéologie des réseaux

Dans les chapitres précédents, les graphes considérés étaient toujours statiques dans le temps. Cependant la plupart des réseaux réels sont en évolution constante, et il est important d'avoir des modèles pour refléter cet aspect dynamique des réseaux. Ici, nous considérerons des modèles de graphes aléatoires dynamiques dits croissants : initialement, le graphe ne contient qu'un seul sommet, et, à chaque instant, une nouvelle arête arrive dans le graphe en s'accrochant aléatoirement à un des sommets déjà présents, selon une certaine loi. En particulier, à chaque temps, le graphe obtenu est un arbre. La question est alors de savoir si, en observant le graphe à un certain temps n , on est capable de déterminer où se trouve le sommet initial, ou au moins de trouver un petit sous-ensemble de sommets tel que le sommet initial y est contenu avec grande probabilité. Les résultats présentés dans ce chapitre ont été montrés par Bubeck et al. [6].

1. Modèles d'arbres croissants : attachement uniforme et préférentiel

Soit $(T_n)_{n \geq 1}$ une suite d'arbres aléatoires construite de la façon suivante :

- au temps $n = 1$, l'arbre T_1 est formé d'un seul sommet isolé, étiqueté 1 ;
- pour $n \geq 1$, conditionnellement à T_n , l'arbre T_{n+1} est formé à partir de T_n par l'ajout d'un nouveau sommet étiqueté $n + 1$ et d'une nouvelle arête $\{n + 1, v\}$ où v est un sommet de T_n tiré au hasard selon une certaine loi de probabilité.

Nous nous concentrerons ici sur deux lois d'attachement naturelles : l'attachement uniforme et l'attachement préférentiel. Dans le modèle d'attachement uniforme, la nouvelle arête s'attache à un sommet de T_n choisi uniformément au hasard, i.e. pour $i \in \{1, \dots, n\}$,

$$\mathbf{P}_{\text{unif}}(v = i \mid T_n) = \frac{1}{n}.$$

Dans le modèle d'attachement préférentiel, la nouvelle arête s'attache à un sommet de T_n choisi avec probabilité proportionnelle à son degré, i.e. pour $i \in \{1, \dots, n\}$,

$$\mathbf{P}_{\text{pref}}(v = i \mid T_n) = \frac{\text{deg}_{T_n}(i)}{2(n-1)},$$

où $\text{deg}_{T_n}(i)$ est le degré de i dans T_n (l'arbre T_n a n sommets donc n_1 arêtes et la somme des degrés est égale à $2(n-1)$). En fait, pour l'attachement préférentiel, le modèle n'est pas bien défini au temps 1 puisque la racine a degré 0. On considère alors qu'on commence la dynamique au temps 2, avec T_2 l'arbre formé par deux sommets relié par une arête.

Si l'on observe le graphe T_n au temps n sans les étiquettes des sommets, est-on capable de retrouver la racine, c'est-à-dire le sommet 1 de départ (ou bien un des deux sommets de départ dans le cas de l'attachement préférentiel) ? Plus précisément, pour une probabilité d'erreur

$\varepsilon \in]0, 1[$ fixée, le but est de trouver dans T_n un sous-ensemble de sommets A de cardinal aussi petit possible, et qui soit tel que la racine appartienne à cet ensemble avec probabilité supérieure à $1 - \varepsilon$. Le résultat surprenant que nous allons montrer est que, dans chacun des deux modèles d'attachement, il est possible de construire un tel sous-ensemble de cardinal d'ordre constant, i.e. dont la taille dépend de la cible ε mais pas de n . Remarquons que cela n'est pas garanti. Par exemple, une idée naïve serait de prendre les K sommets de plus grand degré (on peut en effet remarquer que plus un sommet est ancien, plus son degré est grand en espérance). Mais pour s'assurer que la racine soit dans cet ensemble avec probabilité proche de 1, on peut montrer qu'il faut prendre K de l'ordre de $\log n$. Il va donc falloir être plus fin que cela. L'idée va être de chercher les sommets qui sont les plus « centraux » dans l'arbre, en un sens précisé ci-dessous. L'ingrédient principal de la preuve sera une analogie entre l'évolution de la taille des sous-arbres autour d'un sommet donné et l'évolution des couleurs dans une urne de Polya.

2. Un algorithme simple pour trouver la racine

Pour un arbre $T = (V, E)$, et $u \in V$, le graphe $T \setminus \{u\}$ (obtenu en retirant le sommet u ainsi que les arêtes adjacentes à u) possède $\deg(u)$ composantes connexes. Pour $v \in \mathcal{N}(u)$, notons $T_v^{(u)}$ la composante connexe de $T \setminus \{u\}$ qui contient v . Introduisons la fonction $\psi_T : V(T) \rightarrow \mathbb{N}$ définie par

$$\psi_T(u) = \max_{v \in \mathcal{N}(u)} |T_v^{(u)}|.$$

En d'autres termes, $\psi_T(u)$ donne la taille du plus grand sous-arbre autour de u . Pour $K \in \llbracket 1, |V| \rrbracket$, soit $A_K(T)$ qui contient les K sommets avec la plus petite valeur de ψ_T . On va montrer que pour T_n généré selon le modèle à attachement uniforme ou préférentiel, et pour K assez grand mais ne dépendant que de ε , l'ensemble $A_K(T_n)$ contient la racine avec probabilité supérieure à $1 - \varepsilon$. Avant d'énoncer ces résultats, rappelons la dynamique d'une urne de Polya, et voyons comment les modèles étudiés font intervenir cette dynamique.

3. Urnes de Polya

Dans une urne de Polya classique, l'urne contient initialement $b \geq 1$ boules blanches et $r \geq 1$ boules rouges. À chaque temps, on tire uniformément au hasard une boule dans l'urne et on la remplace avec une autre boule de la même couleur. Notons X_n le nombre de boules rouges dans l'urne après le $n^{\text{ième}}$ tirage (initialement $X_0 = r$), et $x_n = \frac{X_n}{n+b+r}$ la proportion des boules rouges dans l'urne. Comme à l'instant $n+1$, on tire une boule rouge avec probabilité x_n , on a

$$\mathbf{E}[x_{n+1} \mid x_n] = \frac{(n+b+r)x_n + x_n}{n+b+r+1} = x_n.$$

Ainsi, la suite $(x_n)_{n \geq 0}$ est une martingale. Comme cette suite est bornée ($x_n \in [0, 1]$), il s'ensuit qu'elle converge presque sûrement vers une variable aléatoire x . Quelle est la loi de cette variable limite? Remarquons que pour $k \in \llbracket 0, n \rrbracket$, la probabilité de tirer d'abord k boules rouges puis $n-k$ boules blanches s'écrit

$$\frac{r}{b+r} \cdot \frac{r+1}{b+r+1} \cdots \frac{r+k-1}{b+r+k-1} \cdot \frac{b}{b+r+k} \cdots \frac{b+n-k-1}{b+r+n-1}$$

et que cette probabilité ne dépend en fait pas de l'ordre dans lequel on tire les boules, mais seulement du nombre de boules rouges tirées. Ainsi on a

$$\mathbf{P}(X_n = r + k) = \binom{n}{k} \frac{r(r+1)\dots(r+k-1) \cdot b(b+1)\dots(n+n-k-1)}{(b+r)(b+r+1)\dots(b+r+n-1)}.$$

On dit que $X_n - r$ suit la loi Beta-Binomiale de paramètres n, r, b . Une autre façon d'obtenir cette loi est de d'abord tirer une valeur p dans $[0, 1]$ selon la loi Beta(r, b) de densité

$$z \mapsto \frac{\Gamma(b+r)}{\Gamma(b)\Gamma(r)} z^{b-1}(1-z)^{r-1} \mathbb{1}_{[0,1]}(z),$$

puis, conditionnellement à p , de tirer une binomiale de paramètres n et p . Conditionnellement à p , la loi des grands nombres nous dit que $\frac{X_n - b}{n}$ converge presque sûrement vers p , et comme $x_n = \frac{X_n - b}{n} + o(1)$, il en est de même de x_n . Ainsi la suite (x_n) converge presque sûrement vers une variable limite de loi Beta(r, b).

On peut généraliser l'urne de Polya classique en considérant un nombre K de couleurs et un nombre d de boules ajoutées après chaque tirage : initialement, l'urne contient a_1 boules de la couleur 1, a_2 boules de la couleur 2, ... , a_K boules de la couleur K . Notons $a = a_1 + \dots + a_K$ le nombre total de boules initialement présentes dans l'urne. À chaque temps, on tire uniformément une boule dans l'urne et on la remplace avec en plus d boules de la même couleur. On note $X_n(k)$ le nombre de boules de couleur k présentes dans l'urne après le $n^{\text{ième}}$ tirage, et $x_n(k) = \frac{X_n(k)}{n+a}$ la proportion de boules de couleur k . Une généralisation de la loi Beta en dimension supérieure est la loi de Dirichlet : pour $K \geq 2$, la loi de Dirichlet de paramètres $\alpha_1, \dots, \alpha_K > 0$, notée $\text{Dir}(\alpha_1, \dots, \alpha_K)$ est la loi à valeurs dans le simplexe

$$\mathcal{S}_K = \left\{ z = (z_1, \dots, z_K) \in (\mathbb{R}_+^*)^K, \sum_{i=1}^K z_i = 1 \right\}$$

dont la densité est donnée par

$$z = (z_1, \dots, z_K) \mapsto \frac{\Gamma\left(\sum_{i=1}^K \alpha_i\right)}{\prod_{i=1}^K \Gamma(\alpha_i)} \prod_{i=1}^K z_i^{\alpha_i-1} \mathbb{1}_{\mathcal{S}_K}(z).$$

On a la convergence suivante :

$$(x_n(1), \dots, x_n(K)) \xrightarrow{\text{p.s.}} x \quad \text{où} \quad x \sim \text{Dir}\left(\frac{a_1}{d}, \dots, \frac{a_K}{d}\right).$$

Voyons en quoi les modèles d'arbres étudiés possèdent une dynamique très proche de celle d'une urne de Polya. Soit T_n un arbre aléatoire généré selon le modèle à attachement uniforme et notons $T_{1,k}, \dots, T_{k,k}$ les k sous-arbres obtenus à partir de T_n en enlevant toutes les arêtes entre les sommets $\{1, \dots, k\}$. À partir du temps k , les nouveaux sommets viennent s'attacher au sous-arbre $T_{i,k}$ avec probabilité proportionnelle au nombre de sommets dans $T_{i,k}$ ainsi le vecteur $\left(\frac{|T_{1,k}|}{n}, \dots, \frac{|T_{k,k}|}{n}\right)$ peut se voir comme le vecteur des proportions des couleurs dans une urne de Polya à k couleur avec $d = 1$, et contenant initialement une boule de chaque couleur.

Et l'on a

$$\left(\frac{|T_{1,k}|}{n}, \dots, \frac{|T_{1,k}|}{n} \right) \xrightarrow{\text{p.s.}} \text{Dir}(1, \dots, 1).$$

Dans le modèle à attachement préférentiel, c'est plutôt la somme des degrés dans les sous-arbres qui évolue selon une urne de Polya. En effet, à partir de temps k , les nouvelles arêtes viennent s'attacher dans $T_{i,k}$ avec probabilité proportionnelle à la somme des degrés dans ce sous-arbre, et cette somme augmente alors de 2. En remarquant que

$$\sum_{u \in T_{i,k}} \deg_{T_n}(u) = \deg_{T_k}(i) + 2(|T_{i,k}| - 1),$$

et que $\deg_{T_k}(i)$ ne peut pas être plus grand que $k - 1$, on a

$$\left(\frac{\sum_{u \in T_{i,k}} \deg_{T_n}(u)}{2(n-1)}, \dots, \frac{\sum_{u \in T_{k,k}} \deg_{T_n}(u)}{2(n-1)} \right) - \left(\frac{|T_{1,k}|}{n}, \dots, \frac{|T_{1,k}|}{n} \right) \xrightarrow{\text{p.s.}} (0, \dots, 0),$$

et ainsi

$$\left(\frac{|T_{1,k}|}{n}, \dots, \frac{|T_{1,k}|}{n} \right) \xrightarrow{\text{p.s.}} \text{Dir} \left(\frac{\deg_{T_k}(1)}{2}, \dots, \frac{\deg_{T_k}(k)}{2} \right).$$

4. Performances de l'algorithme

Proposition 5.1. *Soit $\varepsilon \in]0, 1/2[$.*

(1) *Pour $K \geq \frac{30 \log(1/\varepsilon)}{\varepsilon}$, on a*

$$\liminf_{n \rightarrow +\infty} \mathbf{P}_{\text{unif}}(1 \in A_K(T_n)) \geq 1 - \varepsilon.$$

(2) *Pour $K \geq \frac{C \log^2(1/\varepsilon)}{\varepsilon^4}$ avec $C > 0$ une constante universelle, on a*

$$\liminf_{n \rightarrow +\infty} \mathbf{P}_{\text{pref}}(1 \in A_K(T_n)) \geq 1 - \varepsilon.$$

Remarque 5.1. Pour l'attachement uniforme, Bubeck et al. [6] construisent un algorithme plus performant que celui présenté ici. Ils montrent que l'on peut trouver un sous-ensemble $B_K(T_n)$ de taille K tel que, si

$$(5.1) \quad K \geq a \exp \left(\frac{b \log(1/\varepsilon)}{\log \log(1/\varepsilon)} \right),$$

avec $a, b > 0$ des constantes universelles, alors $\liminf_{n \rightarrow +\infty} \mathbf{P}_{\text{unif}}(1 \in B_K(T_n)) \geq 1 - \varepsilon$.

Preuve de la Proposition 5.1. (1) Commençons par montrer le résultat concernant l'attachement uniforme. Pour alléger les notations, notons ψ au lieu de ψ_{T_n} . Observons que

$$\begin{aligned} \mathbf{P}(1 \notin A_K(T_n)) &\leq \mathbf{P}(\exists i > K, \psi(i) \leq \psi(1)) \\ &\leq \mathbf{P}(\psi(1) \geq (1 - \varepsilon/4)n) + \mathbf{P}(\exists i > K, \psi(i) \leq (1 - \varepsilon/4)n). \end{aligned}$$

Clairement $\psi(1) \leq \max\{|T_{1,2}|, |T_{2,2}|\}$. Or les variables $|T_{1,2}|$ et $|T_{2,2}|$ sont i.i.d. et, par les résultats de la section précédente,

$$\frac{|T_{1,2}|}{n} \xrightarrow{\mathcal{L}} \text{Beta}(1, 1) = \text{Unif}[0, 1].$$

Ainsi

$$\limsup_{n \rightarrow +\infty} \mathbf{P}(\psi(1) \geq (1 - \varepsilon/4)n) \leq 2 \lim_{n \rightarrow +\infty} \mathbf{P}\left(\frac{|T_{1,2}|}{n} \geq 1 - \varepsilon/4\right) = \frac{\varepsilon}{2}.$$

D'autre part, pour tout $i > K$,

$$\psi(i) \geq \min_{1 \leq k \leq K} \sum_{\substack{j=1 \\ j \neq k}}^K |T_{j,K}|.$$

Or, comme $\left(\frac{|T_{1,K}|}{n}, \dots, \frac{|T_{K,K}|}{n}\right)$ converge vers la loi $\text{Dir}(1, \dots, 1)$, la variable $\frac{1}{n} \sum_{\substack{j=1 \\ j \neq k}}^K |T_{j,K}|$ converge vers la loi $\text{Beta}(K-1, 1)$. Ainsi

$$\begin{aligned} \limsup_{n \rightarrow +\infty} \mathbf{P}(\exists i > K, \psi(i) \leq (1 - \varepsilon/4)n) &\leq \lim_{n \rightarrow +\infty} \mathbf{P}\left(\exists k \in \llbracket 1, K \rrbracket, \frac{1}{n} \sum_{\substack{j=1 \\ j \neq k}}^K |T_{j,K}| \leq 1 - \varepsilon/4\right) \\ &\leq K(K-1) \int_0^{1-\varepsilon/4} z^{K-2} dz = K(1 - \varepsilon/4)^{K-1}. \end{aligned}$$

On a donc

$$\limsup_{n \rightarrow +\infty} \mathbf{P}(1 \notin A_K(T_n)) \leq \frac{\varepsilon}{2} + K(1 - \varepsilon/4)^{K-1}.$$

On peut vérifier que pour $0 < \varepsilon < 2$ et $K \geq \frac{30 \log(1/\varepsilon)}{\varepsilon}$, cette quantité est inférieure à ε .

(2) Considérons maintenant le modèle à attachement préférentiel. Comme précédemment, observons que

$$\begin{aligned} \mathbf{P}(1 \notin A_K(T_n)) &\leq \mathbf{P}(\exists i > K, \psi(i) \leq \psi(1)) \\ &\leq \mathbf{P}(\psi(1) \geq (1 - \varepsilon^2/4)n) + \mathbf{P}(\exists i > K, \psi(i) \leq (1 - \varepsilon^2/4)n). \end{aligned}$$

En utilisant que $\psi(1) \leq \max\{|T_{1,2}|, |T_{2,2}|\}$, que les variables $|T_{1,2}|$ et $|T_{2,2}|$ sont i.i.d. et que

$$\frac{|T_{1,2}|}{n} \xrightarrow{\mathcal{L}} \text{Beta}(1/2, 1/2),$$

on a

$$\limsup_{n \rightarrow +\infty} \mathbf{P}(\psi(1) \geq (1 - \varepsilon^2/4)n) \leq 2 \lim_{n \rightarrow +\infty} \mathbf{P}\left(\frac{|T_{1,2}|}{n} \geq 1 - \varepsilon^2/4\right) = \frac{2}{\pi} \arcsin(\varepsilon/2) \leq \frac{\varepsilon}{2}.$$

On a utilisé que la loi $\text{Beta}(1/2, 1/2)$, aussi appelée loi de l'arcsinus, a pour densité $z \mapsto \frac{1}{\pi\sqrt{z(1-z)}}$ et donc pour fonction de répartition

$$\int_0^x \frac{1}{\pi\sqrt{z(1-z)}} dz = \int_0^{\sqrt{x}} \frac{2}{\pi\sqrt{1-t^2}} dt = \frac{2}{\pi} \arcsin(\sqrt{x}).$$

D'autre part, pour tout $i > K$,

$$\psi(i) \geq \min_{1 \leq k \leq K} \sum_{\substack{j=1 \\ j \neq k}}^K |T_{j,K}|,$$

et pour $1 \leq k \leq K$, la variable $\sum_{\substack{j=1 \\ j \neq k}}^K |T_{j,K}|$ est stochastiquement plus grande que $\sum_{j=2}^K |T_{j,K}|$.

Or comme $\left(\frac{|T_{1,K}|}{n}, \dots, \frac{|T_{K,K}|}{n}\right)$ converge vers la loi Dir $\left(\frac{\deg_{T_K}(1)}{2}, \dots, \frac{\deg_{T_K}(K)}{2}\right)$, la variable $\frac{1}{n} \sum_{j=2}^K |T_{j,K}|$ converge vers la loi

$$\text{Beta} \left(\frac{1}{2} \sum_{j=2}^K \deg_{T_K}(j), \frac{\deg_{T_K}(1)}{2} \right) = \text{Beta} \left(K - 1 - \frac{\deg_{T_K}(1)}{2}, \frac{\deg_{T_K}(1)}{2} \right).$$

Ainsi

$$\begin{aligned} \limsup_{n \rightarrow +\infty} \mathbf{P} \left(\exists i > K, \psi(i) \leq (1 - \varepsilon^2/4)n \right) &\leq \lim_{n \rightarrow +\infty} \mathbf{P} \left(\exists k \in \llbracket 1, K \rrbracket, \frac{1}{n} \sum_{\substack{j=1 \\ j \neq k}}^K |T_{j,K}| \leq 1 - \varepsilon^2/4 \right) \\ &\leq K \mathbf{P} \left(\text{Beta} \left(K - 1 - \frac{\deg_{T_K}(1)}{2}, \frac{\deg_{T_K}(1)}{2} \right) \leq 1 - \varepsilon^2/4 \right) \end{aligned}$$

On peut montrer que pour $K \geq C \frac{\log^2(1/\varepsilon)}{\varepsilon}$ avec $C > 0$ assez grand, ce terme est inférieur à $\varepsilon/2$. ■

5. Bornes inférieures

Nous avons vu que, dans chacun des deux modèles, il est possible de construire un sous-ensemble de confiance pour la racine, dont la taille ne dépend que de ε . Cette dépendance en ε est-elle optimale? Peut-on obtenir des bornes inférieures pour la taille d'un ensemble de confiance? D'autre part, dans lequel des deux modèles est-il plus difficile de retrouver la racine? Intuitivement, on pourrait se dire que le phénomène de renforcement des degrés dans l'attachement préférentiel devrait aider à retrouver la racine plus facilement. Or c'est l'inverse qui se produit : dans l'attachement préférentiel, s'il se trouve qu'au début du processus, un sommet autre que la racine commence à avoir un degré relativement grand, alors cet écart sera renforcé et la racine ne pourra plus être retrouvée grâce à son degré. Dans cette section, nous allons montrer une borne inférieure dans l'attachement préférentiel, bien plus grande que la borne supérieure que l'on a établi pour l'attachement uniforme.

Proposition 5.2. *Soit T_n un arbre généré par le modèle à attachement préférentiel. Il existe $c > 0$ tel que pour tout $\varepsilon \in]0, 1[$, tout sous-ensemble $A(T_n)$ de sommets vérifiant $\liminf_{n \rightarrow +\infty} \mathbf{P}_{\text{pref}}(1 \in A(T_n)) \geq 1 - \varepsilon$ doit vérifier $|A(T_n)| \geq \frac{c}{\varepsilon}$.*

Remarque 5.2. Par la proposition 5.2, combinée à la borne supérieure (5.1) pour l'attachement uniforme, on voit que le problème de retrouver la racine est exponentiellement plus dur dans le modèle à attachement préférentiel que dans le modèle à attachement uniforme.

Preuve de la Proposition 5.2. Soit $\varepsilon \in]0, 1[$. Nous allons montrer que pour un certain $n \in \mathbb{N}$ (qui pourra dépendre de ε), la procédure optimale à K sommets se trompe avec probabilité supérieure à ε si $K < \frac{c}{\varepsilon}$. Comme la probabilité d'erreur de la procédure optimale est croissante en n (on peut coupler les procédures optimales sur T_n et T_{n+1} de telle sorte que si la procédure

se trompe sur T_n , elle se trompe sur T_{n+1}), cela suffit pour avoir le résultat sur la limite inférieure.

L'idée est de montrer que pour un certain n , avec probabilité supérieure à 2ε , la racine est une feuille attachée au sommet 2, et qu'il y a $2c/\varepsilon$ feuilles attachées au sommet 2. Ainsi, tout sous-ensemble qui contient moins de c/ε sommets se trompe avec probabilité supérieure à ε .

La probabilité que la racine soit une feuille de T_n est égale à

$$\prod_{k=2}^n \left(1 - \frac{1}{2(k-1)}\right) = \exp \left\{ \sum_{k=1}^{n-1} \log \left(1 - \frac{1}{2k}\right) \right\} \geq \exp \left\{ - \sum_{k=1}^{n-1} \left(\frac{1}{2k} + \frac{1}{4k^2} \right) \right\},$$

où l'on a utilisé que pour $x \in]0, 1/2]$, $\log(1-x) \geq -x - x^2$. Comme $\sum_{k=1}^{n-1} \frac{1}{k} \leq 1 + \log(n)$ et $\sum_{k=1}^{n-1} \frac{1}{k^2} \leq \frac{\pi^2}{6}$, cette probabilité est supérieure à

$$\frac{e^{-\frac{1}{2} - \frac{\pi^2}{24}}}{\sqrt{n}} \geq \frac{2}{5\sqrt{n}},$$

Ainsi pour $n = \lfloor \frac{1}{100\varepsilon^2} \rfloor$, la probabilité que la racine soit une feuille dans T_n est supérieure à 4ε . Maintenant, conditionnellement à l'événement que 1 est une feuille dans T_n , le nombre d'autres feuilles attachées au sommet 2 dans T_n a la même loi que le nombre de feuilles attachées à la racine dans T_{n-1} . Montrons que ce nombre est de l'ordre de \sqrt{n} . Notons D_n le degré de la racine dans T_n est F_n le nombre de feuilles attachées à la racine. Conditionnellement à T_n , on a

$$D_{n+1} - D_n = \begin{cases} 1 & \text{avec probabilité } \frac{D_n}{2(n-1)}, \\ 0 & \text{avec probabilité } 1 - \frac{D_n}{2(n-1)}. \end{cases}$$

Par récurrence, on obtient

$$\mathbf{E}D_n \asymp \sqrt{n} \quad \text{et} \quad \mathbf{E}D_n^2 \asymp n.$$

De même, conditionnellement à T_n , on a

$$F_{n+1} - F_n = \begin{cases} 1 & \text{avec probabilité } \frac{D_n}{2(n-1)}, \\ -1 & \text{avec probabilité } \frac{F_n}{2(n-1)}, \\ 0 & \text{avec probabilité } 1 - \frac{D_n + F_n}{2(n-1)}. \end{cases}$$

Par récurrence et en utilisant les ordres de grandeur obtenus pour $\mathbf{E}D_n$ et $\mathbf{E}D_n^2$, on obtient

$$\mathbf{E}F_n \asymp \sqrt{n} \quad \text{et} \quad \mathbf{E}F_n^2 \asymp n.$$

Ainsi, par l'inégalité de Paley–Zygmund 6.6, on voit qu'on peut trouver $\delta > 0$ assez petit tel que

$$\mathbf{P}(F_n \geq \delta\sqrt{n}) \geq \frac{1}{2}.$$

Il existe donc $c > 0$ tel que, pour $n = \lfloor \frac{1}{100\varepsilon^2} \rfloor$, avec probabilité supérieure à 2ε la racine est une feuille et il y a plus de $2c/\varepsilon$ autres feuilles attachées au sommet 2. Ainsi toute procédure qui ne retient que c/ε sommets manque la racine avec probabilité supérieure à ε . ■

Modèles graphiques

Soit $X = (X_1, \dots, X_p)$ un vecteur aléatoire de loi P inconnue. Comment apprendre la loi P à partir de n réalisations indépendantes $X^{(1)}, \dots, X^{(n)}$ de X ? Dans le régime où p est petit par rapport à n , les méthodes statistiques classiques, comme le maximum de vraisemblance, permettent généralement de construire des estimateurs efficaces de P . Mais quand la dimension p est grande par rapport au nombre d'observations, ces méthodes s'avèrent souvent très insatisfaisantes et bien trop coûteuse. Dans de nombreuses situations, on peut cependant espérer que, même si la dimension p est grande, la structure de dépendance du vecteur est en fait relativement creuse. Comprendre cette structure de dépendance devient un enjeu crucial. Par exemple, dans le cas extrême où P est en fait une loi produit $P = Q^{\otimes p}$, le problème devient unidimensionnel et il s'agit simplement d'apprendre la loi Q à partir de $p \times n$ observations. De façon plus général, si la loi P se factorise en plusieurs « petits » blocs indépendants, il est important de le savoir. Les modèles graphiques ont été introduits pour modéliser cette structure de dépendance des lois multivariées. Grosso modo, à une distribution p -dimensionnelle, on associe un graphe G à p sommets, où l'absence d'arête entre i et j signifie que, conditionnellement à $(X_k)_{k \neq i, j}$, les variables X_i et X_j sont indépendantes. La question devient alors celle d'apprendre le graphe G à partir des observations $X^{(1)}, \dots, X^{(n)}$.

Bibliographie :

- Lauritzen [15];
- Drton and Maathuis [10].

1. Propriétés de Markov et théorème d'Hammersley–Clifford

Soit $G = (V, E)$ un graphe (simple et non-dirigé) avec $|V| = p \geq 1$, et soit $X = (X_v)_{v \in V}$ un vecteur aléatoire indexé par V , à valeurs dans un espace mesurable $\mathcal{X} = \prod_{v \in V} \mathcal{X}_v$. On suppose que la loi P de X possède une densité f par rapport à une mesure produit $\mu = \bigotimes_{v \in V} \mu_v$ sur \mathcal{X} (typiquement les mesures μ_v seront soit toutes égales à la mesure de Lebesgue sur \mathbb{R} ou toutes égales à la mesure de comptage sur un ensemble dénombrable). Pour $A \subset V$, on note $X_A = (X_u)_{u \in A}$. On dit que X a la propriété de Markov par paires (notée (P)) si pour toute paire de sommets $u \neq v$ avec $\{u, v\} \notin E$,

$$X_u \perp X_v \mid X_{V \setminus \{u, v\}}.$$

On dit que X a la propriété de Markov locale (notée (L)) si pour tout $v \in V$,

$$X_v \perp X_{V \setminus \mathcal{N}^+(v)} \mid X_{\mathcal{N}(v)},$$

où $\mathcal{N}(v)$ est le voisinage de v et $\mathcal{N}^+(v) = \mathcal{N}(v) \cup \{v\}$. Enfin, on dit que X a la propriété de Markov globale (notée (G)) si pour tout triplet (A, B, S) de sous-ensembles disjoints de V tels que S sépare A et B (i.e. tout chemin allant d'un sommet de A à un sommet de B contient un sommet de S),

$$X_A \perp X_B \mid X_S.$$

On a toujours

$$(G) \Rightarrow (L) \Rightarrow (P).$$

En effet, pour $(G) \Rightarrow (L)$, il suffit de remarquer que $\mathcal{N}(v)$ sépare $\{v\}$ et $V \setminus (\{v\} \cup \mathcal{N}(v))$. Pour $(L) \Rightarrow (P)$, énonçons d'abord deux propriétés de l'indépendance conditionnelle : soit (X, Y, Z) un triplet de variables aléatoires, alors

- (i) si $X \perp Y \mid Z$ et si h est une fonction mesurable, alors $h(X) \perp Y \mid Z$;
- (ii) si $X \perp Y \mid Z$ et si h est une fonction mesurable, alors $X \perp Y \mid (Z, h(X))$.

Supposons (L) et soit $\{u, v\} \notin E$. En particulier $u \in V \setminus \mathcal{N}^+(v)$, et

$$V \setminus \{u, v\} = \mathcal{N}(v) \cup ((V \setminus \mathcal{N}^+(v)) \setminus \{u\}).$$

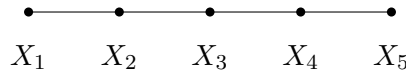
Par hypothèse $X_v \perp X_{V \setminus \mathcal{N}^+(v)} \mid X_{\mathcal{N}(v)}$. Ainsi, la propriété (ii) donne

$$X_v \perp X_{V \setminus \mathcal{N}^+(v)} \mid X_{V \setminus \{u, v\}},$$

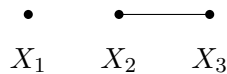
et par la propriété (i),

$$X_v \perp X_u \mid X_{V \setminus \{u, v\}}.$$

Exemple 6.1. Soit $X = (X_1, X_2, X_3, X_4, X_5)$ un vecteur de variable de Bernoulli défini de la façon suivante : X_1 et X_5 sont indépendantes et toutes les deux de loi $\mathcal{B}(1/2)$, $X_2 = X_1$, $X_4 = X_5$ et $X_3 = X_2 X_4$. On peut montrer (exercice) que X vérifie (L) mais pas (G) pour le graphe



Exemple 6.2. Soit $X = (X_1, X_2, X_3)$ un vecteur de variable de Bernoulli avec X_1 de loi $\mathcal{B}(1/2)$, et $X_2 = X_3 = X_1$. On peut montrer (exercice) que X vérifie (P) mais pas (L) pour le graphe



L'indépendance conditionnelle et les propriétés de Markov sont intimement liées à la factorisation de la distribution P . Une clique de G est un sous-graphe complet maximal, c'est-à-dire un sous-ensemble $A \subset V$ tel que pour tous $u, v \in A$, u et v sont voisins et il n'existe pas de sommet dans A^c qui soit voisin de tous les sommets de A . Notons \mathcal{C} l'ensemble des

cliques de G . On dit que la loi P a la propriété de factorisation (notée (F)) selon le graphe G si la densité f de P par rapport à μ peut s'écrire

$$(6.1) \quad f(x) = \prod_{C \in \mathcal{C}} \psi_C(x_C),$$

pour des fonctions positives $\psi_C : \prod_{u \in C} \mathcal{X}_u \rightarrow \mathbb{R}_+$ (pour que cette formulation ait un sens, il faut spécifier un ordre sur les sommets de C . Si V est identifié à $\{1, \dots, p\}$, on peut par exemple choisir l'ordre croissant).

Proposition 6.1. *Pour tout graphe $G = (V, E)$ et toute distribution P sur \mathcal{X} , on a*

$$(F) \Rightarrow (G).$$

Preuve de la Proposition 6.1. Soit (A, B, S) un triplet de sous-ensembles disjoints de V tels que S sépare A et B . Notons \tilde{A} l'ensemble des sommets qui appartiennent à l'une des composantes connexes du graphe induit sur $V \setminus S$ qui contiennent un élément de A et $\tilde{B} = V \setminus (\tilde{A} \cup S)$. Comme A et B sont séparés par S , toute clique de G est soit contenue dans $\tilde{A} \cup S$ soit contenue dans $\tilde{B} \cup S$. Si \mathcal{C}_A est l'ensemble des cliques contenues dans $\tilde{A} \cup S$, alors on a

$$f(x) = \prod_{C \in \mathcal{C}_A} \psi_C(x_C) \prod_{C \in \mathcal{C} \setminus \mathcal{C}_A} \psi_C(x_C) = f_1(x_{\tilde{A} \cup S}) f_2(x_{\tilde{B} \cup S}).$$

Cela montre que $X_{\tilde{A}} \perp X_{\tilde{B}} \mid X_S$, et donc, par la propriété (i), que $X_A \perp X_B \mid X_S$. ■

Lorsque la densité f de P est continue et strictement positive sur \mathcal{X} , le théorème de Clifford–Hammersley [13] énonce que (P) implique (F) et qu'il y a donc équivalence entre toutes ces propriétés.

Théorème 6.2 (Théorème de Clifford–Hammersley [13]). *Soit P une distribution de densité f continue et strictement positive par rapport à une mesure produit μ sur \mathcal{X} . Alors P satisfait la propriété (P) de Markov par paires par rapport au graphe G si et seulement si P a la propriété de factorisation (F) selon G .*

Avant de prouver ce théorème, rappelons la formule d'inversion de Möbius.

Lemme 6.3 (Formule d'inversion de Möbius). *Soient Ψ et Φ des fonctions définies sur l'ensemble des parties d'un ensemble fini V , à valeurs dans un groupe abélien. Alors les deux énoncés suivants sont équivalents :*

$$(i) \text{ pour tout } A \subset V, \Psi(A) = \sum_{B \subset A} \Phi(B) ;$$

$$(ii) \text{ pour tout } A \subset V, \Phi(A) = \sum_{B \subset A} (-1)^{|A \setminus B|} \Psi(B).$$

Preuve du Lemme 6.3. Montrons (i) \Rightarrow (ii). Soit $A \subset V$. On a

$$\begin{aligned} \sum_{B \subset A} (-1)^{|A \setminus B|} \Psi(B) &= \sum_{B \subset A} (-1)^{|A \setminus B|} \sum_{C \subset B} \Phi(C) \\ &= \sum_{C \subset A} \Phi(C) \left\{ \sum_{B, C \subset B \subset A} (-1)^{|A \setminus B|} \right\} \\ &= \sum_{C \subset A} \Phi(C) \left\{ \sum_{H \subset A \setminus C} (-1)^{|H|} \right\} \end{aligned}$$

Or, comme tout ensemble non-vidé a le même nombre de parties paires que de parties impaires,

$$\sum_{H \subset A \setminus C} (-1)^{|H|} = \begin{cases} 1 & \text{si } A \setminus C = \emptyset, \\ 0 & \text{sinon .} \end{cases}$$

Ainsi $\sum_{B \subset A} (-1)^{|A \setminus B|} \Psi(B) = \Phi(A)$. La réciproque se montre de la même manière. \blacksquare

Preuve du théorème 6.2. On a déjà vu que $(F) \Rightarrow (G) \Rightarrow (P)$, il suffit donc de montrer que $(P) \Rightarrow (F)$. Comme f est strictement positive, on peut prendre son logarithme et l'équation (6.1) à établir se ré-écrit

$$\log f(x) = \sum_{A \subset V} \phi_A(x),$$

avec $\phi_A = 0$ dès que A n'est pas un sous-graphe complet de G . Soit x^* un élément fixé de \mathcal{X} et, pour tout $A \subset V$, posons

$$H_A(x) = \log f(x_A, x_{A^c}^*),$$

où $(x_A, x_{A^c}^*)$ est l'élément de \mathcal{X} qui coïncide avec x sur les coordonnées $i \in A$, et avec x^* sur les coordonnées $i \notin A$. Posons aussi

$$\phi_A(x) = \sum_{B \subset A} (-1)^{|A \setminus B|} H_B(x).$$

Les fonctions H_A et ϕ_A ne dépendent de x que par x_A . Par la formule d'inversion de Möbius 6.3, on a

$$\log f(x) = H_V(x) = \sum_{A \subset V} \phi_A(x).$$

Le théorème sera démontré si l'on montre que $\phi_A = 0$ dès que le sous-graphe induit par A n'est pas complet. Supposons qu'il existe $u, v \in A$ tels que $\{u, v\} \notin A$ et notons $C = A \setminus \{u, v\}$. Alors on a

$$(6.2) \quad \phi_A(x) = \sum_{B \subset C} (-1)^{|C \setminus B|} (H_B(x) - H_{B \cup \{u\}}(x) - H_{B \cup \{v\}}(x) + H_{B \cup \{u, v\}}(x)).$$

En utilisant la propriété (P) de Markov par paires,

$$\begin{aligned}
H_B(x) - H_{B \cup \{u\}}(x) &= \log \frac{f(x_B, x_u^*, x_v^*, x_{V \setminus (B \cup \{u, v\})}^*)}{f(x_B, x_u, x_v^*, x_{V \setminus (B \cup \{u, v\})}^*)} \\
&= \log \frac{f(x_B, x_{V \setminus (B \cup \{u, v\})}^*) f(x_u^* \mid x_B, x_{V \setminus (B \cup \{u, v\})}^*) f(x_v^* \mid x_B, x_{V \setminus (B \cup \{u, v\})}^*)}{f(x_B, x_{V \setminus (B \cup \{u, v\})}^*) f(x_u \mid x_B, x_{V \setminus (B \cup \{u, v\})}^*) f(x_v^* \mid x_B, x_{V \setminus (B \cup \{u, v\})}^*)} \\
&= \log \frac{f(x_u^* \mid x_B, x_{V \setminus (B \cup \{u, v\})}^*)}{f(x_u \mid x_B, x_{V \setminus (B \cup \{u, v\})}^*)}.
\end{aligned}$$

De même,

$$\begin{aligned}
H_{B \cup \{v\}}(x) - H_{B \cup \{u, v\}}(x) &= \log \frac{f(x_B, x_u^*, x_v, x_{V \setminus (B \cup \{u, v\})}^*)}{f(x_B, x_u, x_v, x_{V \setminus (B \cup \{u, v\})}^*)} \\
&= \log \frac{f(x_v \mid x_B, x_{V \setminus (B \cup \{u, v\})}^*)}{f(x_v \mid x_B, x_{V \setminus (B \cup \{u, v\})}^*)}.
\end{aligned}$$

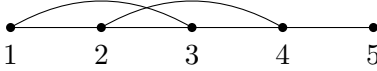
Ainsi tous les termes de la somme (6.2) sont nuls et on a bien $\phi_A(x) = 0$. \blacksquare

Définition 6.1. Soit $G = (V, E)$ un graphe. On appelle modèle graphique, ou champ aléatoire markovien sur G , un vecteur aléatoire $X = (X_u)_{u \in V}$ à valeurs dans $\mathcal{X} = \prod_{u \in V} \mathcal{X}_u$, possédant une densité strictement positive sur \mathcal{X} , et vérifiant, de façon équivalente la propriété (P), (L), (G) ou (F) par rapport à G .

Exemple 6.3 (Modèle graphique gaussien). Soit X un vecteur gaussien en dimension $p = 5$, d'espérance nulle et de matrice de covariance Σ . En notant $K = \Sigma^{-1}$, la densité s'écrit

$$f(x) = \frac{1}{\sqrt{(2\pi)^5 \det(\Sigma)}} \exp \left\{ -\frac{1}{2} \sum_{u, v} K_{u, v} x_u x_v \right\}.$$

La distribution factorise selon le graphe ci-dessous si et seulement si $K_{1,4} = K_{1,5} = K_{2,5} = K_{3,5} = 0$. Plus généralement, le modèle gaussien associé à un graphe $G = (V, E)$ comprend toutes les distributions gaussiennes avec une matrice de covariance inverse définie positive $K \in \mathcal{M}_{|V|, |V|}(\mathbb{R})$ telle que $K_{u, v} = 0$ dès que $\{u, v\} \notin E$.



Exemple 6.4 (Modèle d'Ising). Le modèle d'Ising sur G est la loi de probabilité sur $\sigma : V \rightarrow \{-1, 1\}$ donnée par

$$\forall \sigma \in \{-1, 1\}^V, \quad \mu(\sigma) \propto \exp \left\{ - \sum_{\{u, v\} \in E} \theta_{u, v} \sigma_u \sigma_v \right\},$$

où les constantes $\theta_{u,v} \in \mathbb{R}$ sont appelées constantes de couplage et corresponde à la force de l'interaction entre u et v . Si $\theta_{u,v} > 0$, l'interaction est dite ferromagnétique. Si $\theta_{u,v} < 0$, l'interaction est dite anti-ferromagnétique.

On observe n réalisations de la variable X , i.i.d. de loi P , supposée de densité f sur \mathcal{X} , continue, strictement positive, et vérifiant la propriété (F) (de façon équivalente, (P), (L) et (G)) par rapport à un graphe $G = (V, E)$, avec $|V| = p$. Estimer le graphe G à partir des observations $X^{(1)}, \dots, X^{(n)}$ est en général un problème très difficile. Néanmoins, dans certains cas particuliers, on sait construire des algorithmes efficaces pour estimer le graphe sous-jacent, par exemple lorsque le graphe G est un arbre.

2. Modèles d'arbres

Soit $X = (X_v)_{v \in V}$ un vecteur aléatoire de \mathcal{X} dont la densité (par rapport à une mesure produit μ), notée f_T , est continue, strictement positive, et se factorise par rapport à un arbre $T = (V, E)$. On peut alors écrire

$$(6.1) \quad \forall x \in \mathcal{X}, \quad f_T(x) = \prod_{\{u,v\} \in E} \frac{f_{uv}(x_u, x_v)}{f_u(x_u)f_v(x_v)} \prod_{v \in V} f_v(x_v),$$

où f_v est la densité marginale de X_v et f_{uv} est la densité marginale de (X_u, X_v) . Supposons pour simplifier que les variables X_u sont discrètes, à valeurs dans un ensemble dénombrable \mathcal{X}_u (avec μ la mesure de comptage sur \mathcal{X} et f strictement positive sur \mathcal{X}). On observe n réalisations $X^{(1)}, \dots, X^{(n)}$ i.i.d. de X .

2.1. L'algorithme de Chow et Liu. L'algorithme proposé par Chow and Liu [7] pour estimer T (et en fait la densité f elle-même) fonctionne de la façon suivante. On définit les fréquences empiriques \hat{f}_u et \hat{f}_{uv} : pour $x_u \in \mathcal{X}_u$ et $x_v \in \mathcal{X}_v$,

$$\hat{f}_u(x_u) = \frac{1}{n} \sum_{i=1}^n \mathbb{1}_{\{X_u^{(i)}=x_u\}} \quad \text{et} \quad \hat{f}_{uv}(x_u, x_v) = \frac{1}{n} \sum_{i=1}^n \mathbb{1}_{\{X_u^{(i)}=x_u, X_v^{(i)}=x_v\}}.$$

Par « plug-in » dans (6.1), on pose, pour $x \in \mathcal{X}$,

$$\hat{f}_T(x) = \prod_{\{u,v\} \in E} \frac{\hat{f}_{uv}(x_u, x_v)}{\hat{f}_u(x_u)\hat{f}_v(x_v)} \prod_{v \in V} \hat{f}_v(x_v).$$

La log-vraisemblance en T de l'échantillon $X^{(1)}, \dots, X^{(n)}$ est alors définie par

$$\ell_T(X^{(1)}, \dots, X^{(n)}) = \sum_{i=1}^n \log \hat{f}_T(X^{(i)}).$$

En inversant les sommes, on obtient

$$\frac{1}{n} \ell_T(X^{(1)}, \dots, X^{(n)}) = \sum_{\{u,v\} \in E} \sum_{\substack{x_u \in \mathcal{X}_u \\ x_v \in \mathcal{X}_v}} \hat{f}_{uv}(x_u, x_v) \log \frac{\hat{f}_{uv}(x_u, x_v)}{\hat{f}_u(x_u)\hat{f}_v(x_v)} + \sum_{v \in V} \sum_{x_v \in \mathcal{X}_v} \hat{f}_v(x_v) \log \hat{f}_v(x_v).$$

En remarquant que la deuxième somme (celle sur V) ne dépend pas de l'arbre T considéré, on voit que maximiser $\ell_T(X^{(1)}, \dots, X^{(n)})$ en T revient à maximiser

$$W(T) = \sum_{\{u,v\} \in E} I(\hat{f}_{uv}),$$

où

$$I(\hat{f}_{uv}) = \sum_{\substack{x_u \in \mathcal{X}_u \\ x_v \in \mathcal{X}_v}} \hat{f}_{uv}(x_u, x_v) \log \frac{\hat{f}_{uv}(x_u, x_v)}{\hat{f}_u(x_u) \hat{f}_v(x_v)}$$

est l'information mutuelle de la loi empirique \hat{f}_{uv} . Ainsi, trouver $\hat{T} \in \arg \max_T W(T)$ revient à trouver un arbre couvrant de V de poids maximal (le poids d'un arbre étant défini par les sommes des poids de ses arêtes), lorsque le poids de chaque arête $\{u, v\}$ est donné par $I(\hat{f}_{uv})$. Or il existe plusieurs algorithmes simples pour trouver un tel arbre, par exemple l'algorithme de Kruskal.

2.2. L'algorithme de Kruskal. Soit $G = (V, E)$ un graphe connexe. À chaque arête $e \in E$ est associé un poids $w(e) \geq 0$. Un arbre couvrant T de G est un arbre dont toutes les arêtes sont dans E et qui contient tous les sommets de V . L'algorithme de Kruskal fournit une méthode pour trouver rapidement un arbre couvrant maximal, i.e. un arbre couvrant T qui maximise

$$W(T) = \sum_{e \in T} w(e).$$

La procédure est la suivante : on commence par classer les arêtes par ordre de poids décroissant $w(e_1) \geq \dots \geq w(e_{|E|})$. L'arbre couvrant T^* est alors construit progressivement : pour $k = 1, \dots, |E|$, on inclut l'arête e_k pourvu que cette arête ne crée pas de cycle (si elle crée un cycle, on passe à l'arête e_{k+1}), jusqu'à ce que tous les sommets soient couverts par T^* . La complexité de cet algorithme correspond essentiellement à la première étape de tri des arêtes, soit $O(|E| \log |E|)$.

Chow and Liu [7] ont montré que l'arbre \hat{T} ainsi obtenu étant un estimateur consistant de T , au sens où

$$\mathbf{P}(\hat{T} = T) \rightarrow 1 \text{ quand } n \rightarrow +\infty.$$

3. Reconstruction de graphes sparses

Dans le cas où le graphe G à reconstruire n'est pas un arbre, le problème est bien plus difficile. On peut néanmoins s'en sortir si l'on se place dans le cas sparse où tous les degrés sont majorés par une constante.

Soit $\mathcal{G}_{d,p}$ l'ensemble des graphes $G = (V, E)$ avec $|V| = p$ et $\max_{v \in V} \deg(v) \leq d$. On se place dans le cadre où $\mathcal{X} = \prod_{u \in V} \mathcal{X}_u$ est ensemble fini, avec $\max_{u \in V} |\mathcal{X}_u| \leq K$, si bien que $|\mathcal{X}| \leq K^p$. Étant donné un échantillon $X^{(1)}, \dots, X^{(n)}$ i.i.d. distribué selon un champ markovien sur $G \in \mathcal{G}_{d,p}$ à valeurs dans \mathcal{X} , l'objectif est de construire un estimateur \hat{G} de G tel que $\mathbf{P}(\hat{G} = G)$ soit le plus grand possible. Plus précisément, la question est de savoir à partir de quelles valeurs de n (en fonction de p) on peut espérer avoir $\mathbf{P}(\hat{G} = G) \rightarrow 1$ quand $p \rightarrow \infty$.

Commençons par énoncer une borne inférieure sur le nombre d'observations nécessaires à la reconstruction.

Proposition 6.4 (Bresler et al. [4]). *Soit G un graphe aléatoire tiré uniformément dans $\mathcal{G}_{d,p}$. Alors il existe une constante $c > 0$ telle que, si $n \leq \frac{cd}{\log K} \log p$, tout estimateur \widehat{G} vérifie $\mathbf{P}(\widehat{G} = G) = o(1)$.*

Preuve de la Proposition 6.4. Remarquons déjà que l'estimateur \widehat{G}_\star qui a le plus de chance de réussir est l'estimateur du maximum a posteriori

$$\widehat{G}_\star(\mathbf{X}) = \arg \max_{g \in \mathcal{G}_{d,p}} \mathbf{P}(G = g \mid \mathbf{X}),$$

où $\mathbf{X} = (X^{(1)}, \dots, X^{(n)})$. Dans cette définition, si le maximum est atteint par plusieurs graphes, on en choisit un selon une règle arbitraire. Ainsi \widehat{G}_\star est une fonction déterministe de \mathbf{X} . En particulier, si $S \subset \mathcal{G}_{d,p}$ est l'ensemble des valeurs possibles pour $\widehat{G}_\star(\mathbf{X})$, alors $|S| \leq K^{np}$ (puisque \mathbf{X} prend au plus K^{np} valeurs). En remarquant que pour $g \notin S$, $\mathbf{P}(\widehat{G}_\star = g \mid G = g) = 0$, on a

$$\begin{aligned} \mathbf{P}(\widehat{G}_\star = G) &= \sum_{g \in S} \mathbf{P}(G = g) \mathbf{P}(\widehat{G}_\star = g \mid G = g) \\ &\leq \sum_{g \in S} \mathbf{P}(G = g) = \frac{|S|}{|\mathcal{G}_{d,p}|} \leq \frac{K^{np}}{|\mathcal{G}_{d,p}|}. \end{aligned}$$

Montrons que pour d fixé, $\log |\mathcal{G}_{d,p}| \gtrsim dp \log p$. Commençons par minorer la taille de $\mathcal{G}_{d,p+2}$ en fonction de celle de $\mathcal{G}_{d,p}$. Soit g un graphe à p sommets de degré au plus d auquel on souhaite ajouter deux sommets a et b . On commence par choisir d sommets dans g , notés u_1, \dots, u_d , et on place une arête entre a et u_j , en étiquetant cette arête avec le label a_j . Il y a $\binom{p}{d} d!$ façons de le faire. Si u_j était déjà de degré d , on sait qu'il existe au moins un voisin v de u_j qui n'est pas voisin de a (puisque, à part u_j , a n'a que $d-1$ autres voisins). On retire alors l'arête entre v et u_j pour en faire une arête entre v et b , et l'on étiquette cette arête avec le label b_j . En enlevant les étiquettes sur les arêtes de a et b (et il y a au plus $d!^2$ étiquetages possibles), on obtient un graphe à $p+2$ sommets de degrés au plus d . Ainsi

$$|\mathcal{G}_{d,p+2}| \geq |\mathcal{G}_{d,p}| \binom{p}{d} d! \frac{1}{d!^2} = |\mathcal{G}_{d,p}| \binom{p}{d} \frac{1}{d!}.$$

Pour p pair, on obtient par récurrence

$$|\mathcal{G}_{d,p}| \geq \prod_{i=1}^{p/2} \binom{p-2i}{d} \frac{1}{d!} \geq \left(\binom{p/2}{d} \frac{1}{d!} \right)^{p/4},$$

et pour p impair, on peut juste dire $|\mathcal{G}_{d,p}| \geq |\mathcal{G}_{d,p-1}|$. En passant au logarithme, on a bien

$$\log |\mathcal{G}_{d,p}| = \Omega(dp \log p).$$

Ainsi,

$$\mathbf{P}(\widehat{G}_\star = G) \leq K^{np} e^{-\Omega(dp \log p)} = e^{p(\log(K)n - \Omega(d \log p))},$$

qui tend vers 0 avec p dès que $n \leq \frac{c}{\log K} d \log p$ pour une constante $c > 0$ assez petite. ■

Pour la borne supérieure, on se place dans le cas $\mathcal{X} = \{-1, 1\}^{|V|}$ ($K = 2$) mais la preuve s'étend facilement au cas d'alphabets finis quelconques. Pour $A \subset V$, on rappelle que x_A est le vecteur restreint aux coordonnées de A : $x_A = (x_u)_{u \in A}$. Pour $u \in A$, on note $x_A^{(u)}$ le vecteur égal à x_A sauf sur la coordonnée u où x_u est remplacé par $-x_u$. De plus, pour $v \in V$, on note $\mathcal{N}^+(v) = \mathcal{N}(v) \cup \{v\}$.

Proposition 6.5 (Bresler et al. [4]). *Supposons qu'il existe $\varepsilon, \delta > 0$ tels que pour tout $v \in V$, pour tout $u \in \mathcal{N}(v)$, et pour tout $A \subset V \setminus \mathcal{N}^+(v)$ avec $|A| \leq d$, il existe $x_{\mathcal{N}(v)} \in \{-1, 1\}^{|\mathcal{N}(v)|}$ et $x_A \in \{-1, 1\}^{|A|}$ tels que*

$$(6.1) \quad \left| \mathbf{P}(X_v = 1 \mid X_{\mathcal{N}(v)} = x_{\mathcal{N}(v)}) - \mathbf{P}(X_v = 1 \mid X_{\mathcal{N}(v)} = x_{\mathcal{N}(v)}^{(u)}) \right| > \varepsilon,$$

et

$$(6.2) \quad \min \left\{ \mathbf{P}(X_{\mathcal{N}(v)} = x_{\mathcal{N}(v)}, X_A = x_A), \mathbf{P}(X_{\mathcal{N}(v)} = x_{\mathcal{N}(v)}^{(u)}, X_A = x_A) \right\} > \delta.$$

Alors pour $C = C(\varepsilon, \delta) > 0$, si $n \geq Cd \log p$, il existe un estimateur \widehat{G} , calculable en temps polynomial, tel que $\mathbf{P}(\widehat{G} = G) \rightarrow 1$ quand $p \rightarrow +\infty$.

Preuve de la Proposition 6.5. Notons $\widehat{\mathbf{P}}_n$ la mesure de probabilité donnée par les fréquences empiriques observées dans $\mathbf{X} = (X^{(1)}, \dots, X^{(n)})$, i.e. pour tout $x \in \{-1, 1\}^p$,

$$\widehat{\mathbf{P}}_n(X = x) = \frac{1}{n} \sum_{i=1}^n \mathbb{1}_{\{X^{(i)}=x\}}.$$

L'estimateur \widehat{G} est construit de la façon suivante : pour tout $v \in V$, et pour tout $U \subset V \setminus \{v\}$ avec $|U| \leq d$ (l'ensemble U est à voir comme un potentiel candidat pour être le voisinage de v), on calcule la quantité

$$f(v, U) = \min_{A, u \in U} \max_{x_A, x_U} \left| \widehat{\mathbf{P}}_n(X_v = 1 \mid X_A = x_A, X_U = x_U) - \widehat{\mathbf{P}}_n(X_v = 1 \mid X_A = x_A, X_U = x_U^{(u)}) \right|,$$

où le minimum est pris sur les sous-ensembles $A \subset V \setminus (\{v\} \cup U)$ avec $|A| \leq d$ et $u \in U$, et où, pour de tels A et u , le maximum est pris sur les éléments x_A, x_U vérifiant

$$\widehat{\mathbf{P}}_n(X_A = x_A, X_U = x_U) > \frac{\delta}{2},$$

et

$$\widehat{\mathbf{P}}_n(X_A = x_A, X_U = x_U^{(u)}) > \frac{\delta}{2}.$$

Pour tout $v \in V$, l'algorithme choisit le plus grand sous-ensemble $U \subset V \setminus \{v\}$ avec $|U| \leq d$ tel que $f(v, U) > \varepsilon/2$. Avant de passer à l'analyse de l'algorithme, montrons que pour n assez grand, la mesure empirique $\widehat{\mathbf{P}}_n$ est proche de \mathbf{P} . Par une borne union et l'inégalité de Hoeffding, on a

$$\begin{aligned} & \mathbf{P} \left(\exists W \subset V, |W| \leq 2d+1, x_W \in \{-1, 1\}^{|W|}, \left| \widehat{\mathbf{P}}_n(X_W = x_W) - \mathbf{P}(X_W = x_W) \right| > \gamma \right) \\ & \leq \sum_{k=1}^{2d+1} \binom{p}{k} 2^{k+1} e^{-2\gamma^2 n} \leq (2d+1)(2p)^{2d+2} e^{-2\gamma^2 n}. \end{aligned}$$

Or $(2d+1)(2p)^{2d+2}e^{-2\gamma^2 n} \xrightarrow[p \rightarrow \infty]{} 0$ dès que $n \geq \frac{Cd}{\gamma^2} \log p$ pour une constante $C > 0$ assez grande.

Cela implique en particulier que pour de tels n , avec probabilité $1 - o(1)$, pour tout $v \in V$, pour tout $B \subset V \setminus \{v\}$ avec $|B| \leq 2d$ et pour tout $x_B \in \{-1, 1\}^{|B|}$ tel que

$$\widehat{\mathbf{P}}_n(X_B = x_B) > \frac{\delta}{2},$$

on a

$$\begin{aligned} & \left| \widehat{\mathbf{P}}_n(X_v = 1 \mid X_B = x_B) - \mathbf{P}(X_v = 1 \mid X_B = x_B) \right| \\ &= \left| \frac{\widehat{\mathbf{P}}_n(X_v = 1, X_B = x_B)}{\widehat{\mathbf{P}}_n(X_B = x_B)} - \frac{\mathbf{P}(X_v = 1, X_B = x_B)}{\mathbf{P}(X_B = x_B)} \right| \\ &\leq \left| \frac{\widehat{\mathbf{P}}_n(X_v = 1, X_B = x_B)}{\widehat{\mathbf{P}}_n(X_B = x_B)} - \frac{\mathbf{P}(X_v = 1, X_B = x_B)}{\widehat{\mathbf{P}}_n(X_B = x_B)} \right| + \left| \frac{1}{\widehat{\mathbf{P}}_n(X_B = x_B)} - \frac{1}{\mathbf{P}(X_B = x_B)} \right| \\ &< \frac{\gamma}{\delta/2} + \frac{\gamma}{\delta/2(\delta/2 - \gamma)} < \frac{\varepsilon}{4} \end{aligned}$$

pour $\gamma = \frac{\varepsilon\delta^2}{48}$. Plaçons-nous maintenant sur cet événement (pour tout $v \in V$, pour tout $B \subset V \setminus \{v\}$...) et montrons qu'alors l'algorithme trouve le bon graphe. Soit $v \in V$ et $U \subset V \setminus \{v\}$ avec $|U| \leq d$. Supposons $U \not\subset \mathcal{N}(v)$. Alors $U \setminus \mathcal{N}(v) \neq \emptyset$. Par la propriété de Markov locale, on a, pour tous $u \in U \setminus \mathcal{N}(v)$ et $x_{\mathcal{N}(v) \cup U}$,

$$\begin{aligned} \mathbf{P}(X_v = 1 \mid X_{\mathcal{N}(v) \cup U} = x_{\mathcal{N}(v) \cup U}) &= \mathbf{P}(X_v = 1 \mid X_{\mathcal{N}(v)} = x_{\mathcal{N}(v)}) \\ &= \mathbf{P}(X_v = 1 \mid X_{\mathcal{N}(v) \cup U} = x_{\mathcal{N}(v) \cup U}^{(u)}). \end{aligned}$$

Si $x_{\mathcal{N}(v) \cup U}$ est tel que $\widehat{\mathbf{P}}_n(X_{\mathcal{N}(v) \cup U} = x_{\mathcal{N}(v) \cup U}) > \frac{\delta}{2}$ et $\widehat{\mathbf{P}}_n(X_{\mathcal{N}(v) \cup U} = x_{\mathcal{N}(v) \cup U}^{(u)}) > \frac{\delta}{2}$, on a alors

$$\left| \widehat{\mathbf{P}}_n(X_v = 1 \mid X_{\mathcal{N}(v) \cup U} = x_{\mathcal{N}(v) \cup U}) - \widehat{\mathbf{P}}_n(X_v = 1 \mid X_{\mathcal{N}(v) \cup U} = x_{\mathcal{N}(v) \cup U}^{(u)}) \right| < \frac{\varepsilon}{2}.$$

On a donc montré que sur un événement donc la probabilité tend vers 1 dès que $n \geq \frac{Cd}{\varepsilon\delta^2} \log p$, pour tout $v \in V$ et $U \subset V \setminus \{v\}$ avec $|U| \leq d$, $f(v, U) \leq \varepsilon/2$. Maintenant si $U = \mathcal{N}(v)$, alors, sur ce même événement et en utilisant les hypothèses (6.1) et (6.2) et la propriété de Markov locale, pour $u \in \mathcal{N}(v)$ et $A \subset V \setminus \mathcal{N}^+(v)$, il existe $x_{\mathcal{N}(v) \cup A}$ tel que

$$\begin{aligned} & \left| \widehat{\mathbf{P}}_n(X_v = 1 \mid X_{\mathcal{N}(v) \cup A} = x_{\mathcal{N}(v) \cup A}) - \widehat{\mathbf{P}}_n(X_v = 1 \mid X_{\mathcal{N}(v) \cup A} = x_{\mathcal{N}(v) \cup A}^{(u)}) \right| \\ &> \left| \mathbf{P}(X_v = 1 \mid X_{\mathcal{N}(v) \cup A} = x_{\mathcal{N}(v) \cup A}) - \mathbf{P}(X_v = 1 \mid X_{\mathcal{N}(v) \cup A} = x_{\mathcal{N}(v) \cup A}^{(u)}) \right| - \frac{\varepsilon}{2} \\ &= \left| \mathbf{P}(X_v = 1 \mid X_{\mathcal{N}(v)} = x_{\mathcal{N}(v)}) - \mathbf{P}(X_v = 1 \mid X_{\mathcal{N}(v)} = x_{\mathcal{N}(v)}^{(u)}) \right| - \frac{\varepsilon}{2} \\ &> \frac{\varepsilon}{2}. \end{aligned}$$

Ainsi avec probabilité $1 - o(1)$, pour tout $v \in V$, $\mathcal{N}(v)$ est bien le plus grand sous-ensemble U de taille au plus d tel que $f(v, U) > \varepsilon/2$. ■

4. Modèles graphiques gaussiens

Un autre cas particulier où le problème prend une forme un peu plus simple est le cas où $X \sim \mathcal{N}_p(0, \Sigma)$ est un vecteur gaussien en dimension p . Comme vu dans l'exemple 6.3, dans ce cas, estimer le graphe G sous-jacent revient à déterminer les coordonnées (u, v) pour lesquelles $K_{u,v} > 0$, où $K = \Sigma^{-1}$ (la matrice K est parfois appelée matrice de concentration ou de précision). Soient $X^{(1)}, \dots, X^{(n)}$ des vecteurs i.i.d. de loi $\mathcal{N}_p(0, \Sigma)$, avec Σ définie positive. Rappelons que la densité des $X^{(i)}$ est donnée par

$$\forall x \in \mathbb{R}^p, f(x) = \frac{1}{\sqrt{(2\pi)^p \det(\Sigma)}} \exp \left\{ -\frac{1}{2} {}^t x \Sigma^{-1} x \right\}.$$

Soit S la matrice $p \times p$ donnée par

$$S_{u,v} = \frac{1}{n} \sum_{i=1}^n X_u^{(i)} X_v^{(i)}.$$

À constantes près, la log-vraisemblance s'écrit

$$-\frac{n}{2} \log \det(\Sigma) - \frac{1}{2} \sum_{i=1}^n {}^t X^{(i)} \Sigma^{-1} X^{(i)} = \frac{n}{2} \{ \log \det(K) - \text{tr}(SK) \}$$

Ainsi, chercher le maximum de vraisemblance revient à maximiser

$$\ell(K) = \log \det(K) - \text{tr}(SK)$$

sous la contrainte que K est symétrique définie positive. Lorsque $n \geq p$, on peut montrer que la matrice S est presque sûrement définie positive et qu'alors le maximum de vraisemblance existe et est unique presque sûrement. Mais quand $p > n$, la matrice S est singulière et il se peut que la vraisemblance n'admette pas de maximum (voir le paragraphe ci-dessous sur la distribution de Wishart). D'autre part, même lorsque le maximum de vraisemblance \widehat{K} est bien défini, si le graphe $\widehat{G} = (V, \widehat{E})$ est construit en posant $\{u, v\} \notin \widehat{E}$ dès que $\widehat{K}_{u,v} = 0$, alors on obtient généralement un graphe excessivement dense, car il est très peu probable qu'une entrée de \widehat{K} soit précisément égale à 0. Une solution proposée par Meinshausen et al. [17] et étudiée par Yuan and Lin [25] et Banerjee et al. [2] consiste à ajouter à la log-vraisemblance une pénalité de type ℓ_1 de façon à favoriser les solutions « sparses ». Plus précisément, on estime K par

$$\widehat{K} = \arg \min \{ -\log \det(K) + \text{tr}(SK) + \lambda \|K\|_1, K \text{ symétrique définie positive} \},$$

où $\|K\|_1 = \sum_{u \neq v} |K_{u,v}|$ et $\lambda > 0$. Notons que l'introduction d'une pénalité rend le problème strictement convexe, le minimum est donc toujours uniquement atteint, que la matrice S soit ou non de rang plein. Cette pénalisation est analogue à la méthode du Lasso en régression linéaire.

La distribution de Wishart. Soient $X^{(1)}, \dots, X^{(n)}$ des vecteurs (colonnes) i.i.d. de loi $\mathcal{N}_p(0, \Sigma)$ avec Σ définie positive, et soit \mathbf{X} la matrice $n \times p$ dont la $i^{\text{ième}}$ ligne correspond à

${}^tX^{(i)}$. La loi de la matrice aléatoire

$$W = {}^t\mathbf{X}\mathbf{X} = \sum_{i=1}^n X^{(i)} {}^tX^{(i)}$$

est appelée loi de Wishart. On note $W \sim \mathcal{W}_p(n, \Sigma)$. Comme $\text{rg}(X^{(i)} {}^tX^{(i)}) = 1$, on a toujours $\text{rg}(W) \leq n$. Ainsi, si $n < p$, la matrice W est singulière. Inversement, si $n \geq p$, la matrice W est inversible avec probabilité 1. En effet, on peut écrire $X^{(i)} = \Sigma^{-1/2}Z^{(i)}$, où $Z^{(i)} \sim \mathcal{N}_p(0, I_p)$. Ainsi

$$W = \Sigma^{1/2} \left(\sum_{i=1}^n Z^{(i)} {}^tZ^{(i)} \right) \Sigma^{1/2} = \Sigma^{1/2} {}^t\mathbf{Z}\mathbf{Z}\Sigma^{1/2},$$

où \mathbf{Z} est la matrice $n \times p$ dont la $i^{\text{ième}}$ ligne correspond à ${}^tZ^{(i)}$. Montrons que si $n \geq p$, alors le rang de \mathbf{Z} est égal à p . Notons que les colonnes de \mathbf{Z} , notées Z_1, \dots, Z_p , sont des vecteurs i.i.d. de loi $\mathcal{N}_n(0, I_n)$. Ainsi

$$\begin{aligned} \mathbf{P}(\text{rg}(\mathbf{Z}) < p) &= \mathbf{P}(\exists k \in \llbracket 1, p \rrbracket, Z_k \in \text{Vect}(Z_j, j \neq k)) \\ &\leq \sum_{k=1}^p \mathbf{P}(Z_k \in \text{Vect}(Z_j, j \neq k)). \end{aligned}$$

Le sous-espace $\text{Vect}(Z_j, j \neq k)$ est de dimension inférieure ou égale à $p - 1$, donc strictement inférieure à n . Or si $Z \sim \mathcal{N}_n(0, I_n)$ et si M est un sous-espace de \mathbb{R}^n de dimension strictement inférieure à n , alors $\mathbf{P}(Z \in M) = 0$. Ainsi $\text{rg}(\mathbf{Z}) = p$ avec probabilité 1, et, comme $\Sigma^{1/2}$ est aussi de rang p , il en est de même de W . On a donc montré que W est de rang p avec probabilité 1 si et seulement si $n \geq p$.

Outils probabilistes

Proposition 6.6 (Inégalité de Paley–Zygmund). *Soit X une variable aléatoire positive telle que $\mathbf{E}X^2 < \infty$. Alors pour tout $\delta \in]0, 1[$,*

$$\mathbf{P}(X \geq \delta \mathbf{E}X) \geq (1 - \delta)^2 \frac{\mathbf{E}[X]^2}{\mathbf{E}[X^2]}.$$

Proposition 6.7 (Inégalité de Hoeffding). *Soient X_1, \dots, X_n des variables aléatoires indépendantes avec $a_i \leq X_i \leq b_i$. Alors pour tout $t \geq 0$,*

$$\mathbf{P}\left(\left|\sum_{i=1}^n (X_i - \mathbf{E}X_i)\right| \geq t\right) \leq \exp\left\{-\frac{2t^2}{\sum_{i=1}^n (b_i - a_i)^2}\right\}.$$

Proposition 6.8 (Inégalité de Bennett). *Soient X_1, \dots, X_n des variables aléatoires indépendantes avec $X_i \leq c$. Notons $Z = \sum_{i=1}^n X_i$ et $v = \sum_{i=1}^n \mathbf{E}X_i^2$. Alors pour tout $\lambda > 0$,*

$$\log \mathbf{E}e^{\lambda(Z - \mathbf{E}Z)} \leq \frac{v}{c^2} \phi(c\lambda),$$

où $\phi(x) = e^x - x - 1$. En particulier, cela donne une inégalité de type Bernstein à droite : pour tout $t \geq 0$,

$$\mathbf{P}(Z - \mathbf{E}Z \geq t) \leq \exp\left\{-\frac{t^2}{2\left(v + \frac{ct}{3}\right)}\right\}.$$

Proposition 6.9 (Concentration de la binomiale). *Soit $X \sim \text{Bin}(n, p)$ une variable binomiale. Alors pour tout $t \geq 0$,*

$$\mathbf{P}(X - \mathbf{E}X \geq t) \leq \exp\left\{-\frac{t^2}{2(np + t/3)}\right\},$$

et

$$\mathbf{P}(X - \mathbf{E}X \leq -t) \leq \exp\left\{-\frac{t^2}{2np}\right\}.$$

Bibliographie

- [1] G. W. Anderson, A. Guionnet, and O. Zeitouni. *An introduction to random matrices*, volume 118. Cambridge university press, 2010.
- [2] O. Banerjee, L. E. Ghaoui, and A. d'Aspremont. Model selection through sparse maximum likelihood estimation for multivariate gaussian or binary data. *Journal of Machine learning research*, 9(Mar) :485–516, 2008.
- [3] A. Ben-Hamou, R. I. Oliveira, and Y. Peres. Estimating graph parameters with random walks. *Mathematical Statistics and Learning*, 1(3/4) :375–399, 2018.
- [4] G. Bresler, E. Mossel, and A. Sly. Reconstruction of markov random fields from samples : Some observations and algorithms. *SIAM Journal on Computing*, 42(2) :563–578, 2013.
- [5] S. Bubeck, J. Ding, R. Eldan, and M. Z. Rácz. Testing for high-dimensional geometry in random graphs. *Random Structures & Algorithms*, 49(3) :503–532, 2016.
- [6] S. Bubeck, L. Devroye, and G. Lugosi. Finding adam in random growing trees. *Random Structures & Algorithms*, 50(2) :158–172, 2017.
- [7] C. Chow and C. Liu. Approximating discrete probability distributions with dependence trees. *IEEE transactions on Information Theory*, 14(3) :462–467, 1968.
- [8] A. Clauset, M. E. Newman, and C. Moore. Finding community structure in very large networks. *Physical review E*, 70(6) :066111, 2004.
- [9] A. Decelle, F. Krzakala, C. Moore, and L. Zdeborová. Asymptotic analysis of the stochastic block model for modular networks and its algorithmic applications. *Physical Review E*, 84(6) :066106, 2011.
- [10] M. Drton and M. H. Maathuis. Structure learning in graphical modeling. *Annual Review of Statistics and Its Application*, 4 :365–393, 2017.
- [11] M. E. Dyer and A. M. Frieze. The solution of some random np-hard problems in polynomial expected time. *Journal of Algorithms*, 10(4) :451–489, 1989.
- [12] M. Girvan and M. E. Newman. Community structure in social and biological networks. *Proceedings of the national academy of sciences*, 99(12) :7821–7826, 2002.
- [13] J. M. Hammersley and P. Clifford. Markov fields on finite graphs and lattices. *Unpublished manuscript*, 46, 1971.
- [14] L. Katzir, E. Liberty, O. Somekh, and I. A. Cosma. Estimating sizes of social networks via biased sampling. *Internet Mathematics*, 10(3-4) :335–359, 2014.
- [15] S. L. Lauritzen. *Graphical models*, volume 17. Clarendon Press, 1996.
- [16] L. Massoulié. Community detection thresholds and the weak ramanujan property. In *Proceedings of the forty-sixth annual ACM symposium on Theory of computing*, pages 694–703. ACM, 2014.

- [17] N. Meinshausen, P. Bühlmann, et al. High-dimensional graphs and variable selection with the lasso. *The annals of statistics*, 34(3) :1436–1462, 2006.
- [18] E. Mossel, J. Neeman, and A. Sly. Consistency thresholds for the planted bisection model. In *Proceedings of the forty-seventh annual ACM symposium on Theory of computing*, pages 69–75. ACM, 2015.
- [19] E. Mossel, J. Neeman, and A. Sly. Reconstruction and estimation in the planted partition model. *Probability Theory and Related Fields*, 162(3-4) :431–461, 2015.
- [20] E. Mossel, J. Neeman, and A. Sly. A proof of the block model threshold conjecture. *Combinatorica*, 38(3) :665–708, 2018.
- [21] M. E. Newman. Fast algorithm for detecting community structure in networks. *Physical review E*, 69(6) :066133, 2004.
- [22] M. E. Newman. Modularity and community structure in networks. *Proceedings of the national academy of sciences*, 103(23) :8577–8582, 2006.
- [23] A. Y. Ng, M. I. Jordan, and Y. Weiss. On spectral clustering : Analysis and an algorithm. In *Advances in neural information processing systems*, pages 849–856, 2002.
- [24] U. Von Luxburg. A tutorial on spectral clustering. *Statistics and computing*, 17(4) : 395–416, 2007.
- [25] M. Yuan and Y. Lin. Model selection and estimation in the gaussian graphical model. *Biometrika*, 94(1) :19–35, 2007.