



# **De la détection de la structure communautaire des réseaux complexes**

**Mémoire**

**Jean-Gabriel Young**

**Maîtrise en physique**  
Maître ès sciences (M.Sc.)

Québec, Canada

© Jean-Gabriel Young, 2014



# Résumé

Une description précise de la structure mésoscopique des réseaux complexes permet de modéliser efficacement les processus de propagation sur réseaux ainsi que la croissance de ces derniers. Cette structure s'exprime en termes d'une décomposition en communautés (ou groupes denses de noeuds structurellement rapprochés), de sorte que l'identification d'une organisation optimale constitue un problème de décision de classe NP-difficile. Plusieurs algorithmes récents permettent d'obtenir des solutions approchées dans un temps polynomial. Toutefois, la nature ad hoc de ces algorithmes rend difficile l'évaluation de leurs qualités et de leurs faiblesses.

Cet ouvrage fait état d'un formalisme analytique unifiant la théorie de la détection communautaire via une description matricielle. Dans un premier temps, on démontre qu'une grande classe d'algorithmes de détection est équivalente à un problème d'optimisation dont la solution approximative peut être obtenue par la décomposition spectrale de matrices de coût, fonction de la structure du réseau à l'étude. Ces développements établissent un cadre théorique permettant l'étude rigoureuse d'algorithmes ad hoc, et mènent également à un algorithme de détection original, basé sur des principes fondamentaux d'optimisation continue. Dans un deuxième temps, il est démontré par le biais de la théorie des matrices aléatoires que, pour une classe particulière de réseaux, la structure communautaire (ici a priori connue) laisse une empreinte aisément identifiable dans le spectre des matrices de coût associées. Ces deux points de vue complémentaires, optimisation spectrale et théorie des matrices aléatoires, donnent accès à de nouvelles observations importantes qu'une simple étude numérique ne peut expliquer, tel l'apparition d'une limite de détection intrinsèque.

Ces développements analytiques restent toutefois confinés à des modèles de réseaux simples. Pour des problèmes plus complexes, une approche numérique est préconisée. On introduit donc une méthode heuristique de détection permettant d'améliorer les performances de tout algorithme imparfait. Dans la perspective de calibrer cette méthode, on présente également un processus de croissance local polyvalent qui produit des réseaux réalistes possédant une structure communautaire connue.



# Abstract

A precise description of the mesoscopic structure of complex systems is necessary to improve models of the dynamical processes *on* and *of* networks. However, knowledge of this structure comes at great cost, since finding an optimal decomposition in communities is a problem that belongs to the NP hard complexity class. Multiple recent algorithms yield approximate solutions in polynomial time. Most of these algorithms are collections of ad hoc methods, such that only extensive numerical studies lead to insightful comparisons.

In this thesis, we present the basis of a unified theory of community detection, which builds upon recent advances of the spectral theory of complex networks. First, we prove that a large class of detection algorithm is equivalent to an optimization problem that can be solved approximately through the spectral decomposition of a very general cost matrix. Within this framework, otherwise ad hoc algorithms can be studied analytically and rigorously. This point of view also leads to a new, original and first-principled spectral *detection* algorithm. Second, using random matrix theory, we generalize existing results and prove that the spectrum of a class of modular networks contains valuable information on their mesoscopic structure. These complementary approaches, spectral optimization and random matrix theory, give powerful insights into the spectral theory of complex networks, and their relevance to community structure.

These analytical results are unfortunately not yet generalizable to arbitrary networks. For complex cases, we prefer a purely numerical approach. Hence, we introduce a heuristic method that drastically improves the efficiency of existing, imperfect algorithms. To test this method, we also present a local growth process that produces realistic modular networks with known community structure. These networks can then be used as versatile benchmarks.



# Table des matières

<b>Résumé</b>	<b>iii</b>
<b>Abstract</b>	<b>v</b>
<b>Table des matières</b>	<b>vii</b>
<b>Liste des tableaux</b>	<b>ix</b>
<b>Liste des figures</b>	<b>xi</b>
<b>Avant-propos</b>	<b>xv</b>
<b>Liste des abréviations</b>	<b>xvii</b>
<b>Notation</b>	<b>xix</b>
<b>Glossaire</b>	<b>xxi</b>
<b>Liste des contributions</b>	<b>xxiii</b>
<b>Introduction</b>	<b>1</b>
<b>1 Introduction aux réseaux complexes et à leur structure communautaire</b>	<b>5</b>
1.1 Notions élémentaires . . . . .	5
1.2 Structure mésoscopique : réseaux modulaires . . . . .	16
<b>2 Dévoilement des communautés cachées via la détection en cascade sur réseaux</b>	<b>33</b>
2.1 Avant-propos . . . . .	34
2.2 Résumé . . . . .	34
2.3 Abstract . . . . .	34
2.4 Introduction . . . . .	35
2.5 Resolution limit due to shadowing . . . . .	35
2.6 Cascading detection . . . . .	39
2.7 Results and discussion . . . . .	40
2.8 Conclusion and perspectives . . . . .	43
2.9 Appendix : Description and properties of networks, and numerical results . . . . .	44
<b>3 Détection et optimisation spectrale</b>	<b>45</b>
3.1 Formalisme de partitionnement : motivations . . . . .	45

3.2	Formulation matricielle du problème de partitionnement . . . . .	46
3.3	Optimisation spectrale continue . . . . .	53
3.4	Algorithmes de détection spectral . . . . .	59
<b>4</b>	<b>Limite intrinsèque des méthodes de détection : théorie des matrices aléatoires</b>	<b>71</b>
4.1	Modèle stochastique par blocs . . . . .	71
4.2	Densité spectrale de la matrice de déviation . . . . .	77
4.3	Spectre de la matrice de modularité . . . . .	86
4.4	Limite de détectabilité dans le modèle stochastique par blocs . . . . .	89
<b>5</b>	<b>Attachement préférentiel structurel des communautés</b>	<b>93</b>
5.1	Avant-propos . . . . .	94
5.2	Résumé . . . . .	94
5.3	Abstract . . . . .	94
5.4	Introduction . . . . .	95
5.5	Considerations on the structure of communities . . . . .	96
5.6	Link creation mechanism . . . . .	97
5.7	Improved growth model . . . . .	99
5.8	Details of the community structure . . . . .	101
5.9	Relation to Dunbar's number . . . . .	103
5.10	Conclusion . . . . .	104
	<b>Conclusion</b>	<b>107</b>
<b>A</b>	<b>Fonction génératrice de probabilités</b>	<b>109</b>
<b>B</b>	<b>Compendium de preuves et de définitions</b>	<b>111</b>
B.1	Induction d'une organisations de noeuds ou de liens . . . . .	111
B.2	Graphe adjoint . . . . .	112
B.3	Produit de Hadamard et trace . . . . .	114
B.4	Dérivées d'une transformation de la PGF binomiale . . . . .	115
B.5	Formule de Stieltjes-Perron . . . . .	116
<b>C</b>	<b>Vecteurs simplexes réguliers</b>	<b>119</b>
C.1	$n$ -simplexes . . . . .	119
C.2	Méthode numérique de construction des vecteurs simplexes réguliers . . . . .	120
C.3	$n$ -simplexes et figure de sommet de l'hypercube $(n + 1)$ -dimensionnel . . . . .	121
C.4	Rotation dans un espace de dimension arbitraire . . . . .	121
C.5	Projection de la figure de sommet d'un hypercube $(n + 1)$ -dimensionnel . . . . .	125
	<b>Bibliographie</b>	<b>129</b>



# Liste des tableaux

1.1	Organisations de noeuds et de liens pouvant être reliées par induction . . . . .	22
1.2	Algorithmes de détection utilisés pour l'étude de l'arXiv circa 2005 . . . . .	28
1.3	Exposants MLE des distributions des quantités mésoscopiques de arXiv 2005 . . . . .	30
1.4	Modularité et coupe des partitions de l'arXiv . . . . .	31
2.1	Description and properties of real networks used to test the cascading detection algorithm . . . . .	44
2.2	Summary of the results of the cascading detection algorithm . . . . .	44



# Liste des figures

1.1	Réseau dorsal d'Internet circa 2006 . . . . .	7
1.2	Classes de réseaux et matrices d'adjacences associées . . . . .	8
1.3	Propriétés structurelles locales . . . . .	11
1.4	Similarité des voisinages . . . . .	13
1.5	Redondance dans un réseau technologique réel et dans ses équivalents aléatoires . . . . .	14
1.6	Redondance dans un réseau social réel et dans ses équivalents aléatoires . . . . .	16
1.7	Hiérarchie des types d'organisation . . . . .	19
1.8	Taille de l'espace des configurations des divers types d'organisations communautaires . . . . .	20
1.9	Distribution des propriétés mésoscopiques de arXiv cond-mat 2005 . . . . .	29
2.1	Shadowing effect for the CPA . . . . .	37
2.2	Calculation of the similarity between two links . . . . .	37
2.3	Shadowing effect for the LCA . . . . .	38
2.4	Fraction of remaining assignable links for real networks using the cascading approach . . . . .	40
2.5	Distribution of the size of the detected communities . . . . .	41
2.6	Sample of the communities detected with the cascading approach on the Words network . . . . .	42
3.1	Vecteurs simplexes réguliers de dimension 1,2 et 3 . . . . .	50
3.2	Simplexe en tant que projection d'une figure de sommet . . . . .	51
3.3	Paramétrisation de la matrice de partition . . . . .	55
3.4	Solutions relaxées versus vecteurs de base paramétrés . . . . .	60
3.5	Comparaison des solutions continues et discrètes . . . . .	62
3.6	Spectres expérimentaux de plusieurs matrices de coût . . . . .	64
3.7	Solution relaxée versus solution inversée : sélection de modèle . . . . .	65
3.8	Application : Zachary Karate Club . . . . .	67
3.9	Application : Réseau en anneau et convergence du modèle . . . . .	68
3.10	Application : Transition dans la détectabilité de la structure communautaire . . . . .	69
4.1	Modèle stochastique par blocs . . . . .	72
4.2	Procédure de construction des graphes combinatoires . . . . .	79
4.3	Graphes maximaux pour un nombre d'arêtes pair et impair . . . . .	81
4.4	Graphe maximal : cas général . . . . .	81
4.5	Graphes doublés en tant qu'arbres planaires enracinés . . . . .	82
4.6	Densité spectrale de la matrice de déviation . . . . .	85
4.7	Relation entre le spectre des matrices de déviation et de modularité. . . . .	88
4.8	Plus grande valeur propre de la matrice de modularité du SBM à deux blocs . . . . .	90
4.9	Limite de détectabilité du SBM à deux blocs . . . . .	92

- 5.1 Example of unacceptable community structures . . . . . 97
- 5.2 Internal degree distribution in the SPA\* model . . . . . 98
- 5.3 Reproduction of a real networks with the SPA\* model . . . . . 102
- 5.4 Structural correlations . . . . . 102
- 5.5 Dunbar’s number in the SPA\* model . . . . . 103
  
- B.1 Graphe et graphe adjoint en tant que projection d’un graphe bipartie . . . . . 113
  
- C.1 Vecteurs simplexes réguliers de basse dimension . . . . . 120
- C.2 Figures de sommet d’hypercubes de basse dimension . . . . . 122
- C.3 Orientation d’une figure de sommet en dimension 3 . . . . . 123

“A good question is, of course, the key by which infinite answers can be educed.”

---

-Isaac ASIMOV, *Foundation's Edge*  
(1982)



# Avant-propos

En entreprenant l'écriture de ce mémoire, j'ai voulu regrouper plusieurs visions différentes d'une même question dans un seul ouvrage. Même si ce genre de projet ne peut jamais être parfaitement achevé, j'ose espérer que les travaux présentés ici parviendront à mettre en lumière les liens fondamentaux qui existent entre les diverses facettes du problème de la détection communautaire.

Bien que je sois le seul signataire de ce mémoire, je suis redevable à plusieurs collaborateurs, collègues, et amis, à qui je dois mon inspiration, les connaissances et la motivation m'ayant permis de mener ce projet à bien. J'aimerais profiter de cet avant-propos pour les mentionner et les remercier explicitement.

D'entrée de jeu, je voudrais remercier mon directeur, Louis J. Dubé, pour m'avoir offert une place au sein de son équipe, à trois reprises, en tant qu'étudiant d'été, à la maîtrise, et à nouveau au doctorat. C'est à lui que je dois l'existence de ce mémoire, grâce à ses subtiles (et moins subtiles !) poussées dans la bonne direction, ses conseils, sa passion et son ouverture d'esprit envers un projet de recherche en constante évolution.

Évidemment, ce mémoire n'aurait pas non plus vu le jour sans la contribution de l'équipe Dynamica au complet. Plus spécifiquement, j'aimerais remercier les collègues de la section réseau, Antoine Allard et Laurent Hébert-Dufresne, auxquels je dois tant de bons conseils et de bons moments, ainsi que notre dernier impétrant, Edward Laurence, avec qui "ça va brasser", comme disait l'autre. Je voudrais aussi saluer les gars de lumière, Denis Gagnon, Jean-Luc Déziel, Joey Dumont, avec qui j'ai partagé idées et rigolades, autour de l'éternel café (et autres breuvages). Je vous suis redevable de votre point de vue unique sur mes idées et de votre expertise de la chose numérique. Je souhaiterais aussi souligner la contribution de ceux qui nous ont quittés pour de nouveaux horizons ; Pierre-André Noël, Guillaume Painchaud et Vincent Marceau, pour les conseils très nombreux de l'un, très réalistes de l'autre, et les discussions enrichissantes avec tous.

J'aimerais faire une mention bien particulière aux professeurs Alain Hertz et Patrick Desrosiers, qui ont accepté d'examiner ce mémoire de maîtrise. Je n'aurais pu espérer avoir l'avis de meilleurs examinateurs, puisque la tangente méta-heuristique de certains des mes travaux est due à une visite du Prof. Alain Hertz en 2012 et que j'ai été introduit à la théorie des matrices aléatoires sous la tutelle du Prof. Patrick Desrosiers.

Je voudrais aussi saluer mes collègues du Baccalauréat, qui feront à jamais partie de la famille

élargie. Un remerciement tout spécial à Ludovic, sans qui je n'aurais jamais développé la curiosité qui m'a mené à poursuivre mes études, ainsi qu'à Sophie pour nos éternelles discussions sur les aléas de la vie d'étudiant. Je voudrais aussi saluer Alex Côté, qui a promis de lire ce mémoire !

Je tiens à remercier tout spécialement ceux que je côtoie dans la vie de tous les jours. Mes parents, Louise et Ronald, mon frère Raphaël et ma soeur Sarah, pour leurs encouragements, leur intérêt et leur soutien constant. Les amis, qui m'aident à me rappeler de ne pas prendre la vie trop au sérieux. Finalement, je voudrais remercier Laurence du fond du coeur, pour avoir partagé les bonheurs, les angoisses et les rêves qui ont parsemés les deux dernières années.

Jean-Gabriel Young  
11 août 2014



# Liste des abréviations

## Abréviations françaises

CG	Côté gauche
CD	Côté droit
ER	Graphe Erdős-Rényi
TMA	Théorie des matrices aléatoires

## Abréviations anglaises

CM	Configuration Model
CPA	Clique percolation algorithm
LCA	Link clustering algorithm
MGF	Moment generating function
PGF	Probability generating function
PGP	Pretty-good-privacy
SBM	Stochastic block model
SPA	Structural preferential attachment
SVD	Singular value decomposition



# Notation

## Notation générale

$\mathbf{A}$	Matrice.
$\mathbf{A}^T$	Matrice transposée.
$a_{ij}$	Élément $(i, j)$ de la matrice $\mathbf{A}$ .
$a_{ij}^T$	Élément $(i, j)$ de la matrice transposée $\mathbf{A}^T$ .
$[\mathbf{A}]_{i,j}$	Élément $(i, j)$ de la matrice $\mathbf{A}$ .

## Vecteurs

Par soucis de clarté, les produits mixtes entre matrices et vecteurs sont représentés explicitement en termes de produit matriciel. Les matrices  $1 \times 1$  sont considérées comme des scalaires <sup>1</sup>.

$\mathbf{v}$	Matrice <i>colonne</i> associée au vecteur $\vec{v}$ , i.e. $[\mathbf{v}]_{i1} \equiv v_i$ . Agit par la droite sur une matrice.
$\mathbf{u}^T$	Matrice <i>rangée</i> associée au vecteur $\vec{u}$ , i.e. $[\mathbf{u}^T]_{1i} \equiv u_i$ . Agit par la gauche sur une matrice.
$\mathbf{u}^T \mathbf{v}$	Produit scalaire $\vec{u} \cdot \vec{v}$ .
$\mathbf{uv}^T$	Produit tensoriel $\vec{u} \otimes \vec{v}$ .

## Matrices spéciales

L'indice est omis lorsqu'il n'y a pas d'ambiguïté.

$\mathbf{1}_n$	Matrice colonne $n \times 1$ dont tous les éléments sont égaux à 1.
$\mathbf{I}_{n \times n}$	Matrice identité de dimension $n \times n$ .
$\mathbf{0}_{n \times n}$	Matrice nulle de dimension $n \times n$ .
$\mathbf{e}_j^{(n)}$	Vecteur orthogonal unitaire de $\mathbb{R}^n$ .

---

1. Une attention particulière est portée au respect des dimensions matricielles, car cette équivalence ne tient généralement pas.



# Glossaire

## Ensembles

$\mathcal{E}$	Ensemble des liens du graphe.	p. 6
$n_+(i)$	Voisinage du noeud $i$ , incluant $i$ .	p. 13
$\Omega_x$	Collection de communautés de noeuds ( $x = n$ ) ou de liens ( $x = \ell$ ).	p. 17
$\Psi^{(x)}$	Communauté (ensemble) de noeuds ( $x = n$ ) ou de liens ( $x = \ell$ ).	p. 17
$\mathcal{V}$	Ensemble des noeuds du graphe.	p. 6

## Matrices et vecteurs

$A$	Matrice d'adjacence.	p. 6
$B$	Matrice de modularité.	p. 27
$C$	Matrice symétrique de coût.	p. 47
$D$	Matrice diagonale étirant les vecteurs de base de $P$ .	p. 54
$k$	Vecteur des degrés.	p. 9
$\Lambda$	Matrice des valeurs propres de $C$ .	p. 58
$m$	Vecteur des appartenances.	p. 25
$P$ et $P_o$	Matrice d'incidence communautaire (représentation orthogonale).	p. 24
$P_s$	Matrice d'incidence communautaire (représentation simplexe).	p. 51
$\tilde{P}$	Matrice d'incidence communautaire (base orthogonale transformée).	p. 54
$p$	Matrice de projection $\mathbb{R}^n \mapsto \mathbb{R}^{n-1}$ .	p. 50
$R$	Matrice d'incidence communautaire relaxée (rangées $\in \mathbb{R}^g$ ).	p. 56
$R_X(z_j)$	Matrice résolvante de la matrice $X$	p. 77
$S$	Matrice de vecteurs simplex (1 vecteur par colonne).	p. 50
$s$	Vecteur des tailles.	p. 25
$X$	Matrice aléatoire de déviation.	p. 72
$Z$	Matrice d'incidence des liens (binaire).	p. 112

## Autres

$\{^n_k\}$	Nombre de Stirling de seconde espèce.	p. 20
$B_n$	$n^{\text{ème}}$ nombre de Bell.	p. 21
$C_x$	Couverture complète ( $x = c$ ) ou incomplète ( $x = i$ ).	p. 18
$C^{(t)}$	Ratio de transitivité.	p. 11
$C_\ell^{(a)}$	Coefficient d'agrégation du noeud $\ell$ .	p. 12
$C$	Fonction de coût (ou critère de sélection).	p. 22
$c(\Psi_\alpha, \Psi_\beta)$	Chevauchement absolu des communautés $\Psi_\alpha$ et $\Psi_\beta$	p. 26
$c_n(\Psi_\alpha, \Psi_\beta)$	Chevauchement normalisé des communautés $\Psi_\alpha$ et $\Psi_\beta$	p. 26
$D(\mathbf{u}, \mathbf{v})$	Distance séparant deux points $\mathbb{R}^n$	p. 60
$\varepsilon$	Tolérance sur l'analyse de Procrustes (nombre de rangées réassignées).	p. 62
$G(\mathcal{E}, \mathcal{V})$	Graphe formé du couple de liens $\mathcal{E}$ et de noeuds $\mathcal{V}$ .	p. 6
$g$	Nombre de communautés.	p. 17
$k_i$	Degré (du noeud $i$ ).	p. 9
$L_\alpha$	Nombre de liens internes dans la communauté de noeuds $\alpha$ .	p. 26
$M$	Nombre de liens.	p. 6
$m_i$	Nombre d'appartenance du noeud $i$ .	p. 25
$N$	Nombre de noeuds.	p. 6
$\mathcal{P}_x$	Partition complète ( $x = c$ ) ou incomplète ( $x = i$ ).	p. 18
$Q$	Modularité.	p. 27
$\rho(\lambda)$	Densité spectrale d'une matrice.	p. 10
$\rho$	Densité du réseau.	p. 10
$\rho_\alpha$	Densité de la communauté de noeuds $\alpha$ .	p. 27
$\tilde{\rho}_\alpha$	Densité excédentaire de la communauté de noeuds $\alpha$ .	p. 27
$R$	Taille de la coupure d'une partition.	p. 27
$r$	Nombre d'orientations initiales dans l'analyse de Procrustes.	p. 62
$\sigma_i$	Indices des communautés du noeud $i$ .	p. 47
$S(i, j)$	Similarité du voisinage des noeuds $i$ et $j$ .	p. 13
$s_\alpha$	Taille de la communauté $\Psi_\alpha$ .	p. 25

# Liste des contributions

## Articles

- L. Hébert-Dufresne, A. Allard, **J.-G. Young** & L. J. Dubé, *Global efficiency of local immunization on complex networks*, Nature Scientific Reports, 3 :2171 (2013).
- L. Hébert-Dufresne, A. Allard, **J.-G. Young** & L. J. Dubé, *Percolation on random networks with arbitrary k-core structure*, Phys. Rev. E 88 (2013), 062820.
- A. Allard, L. Hébert-Dufresne, **J.-G. Young** & L. J. Dubé, *Coexistence of phases and the observability of random graphs*, Phys. Rev. E 89 (2014), 022801.
- L. Hébert-Dufresne, A. Allard, **J.-G. Young** & L. J. Dubé, *Universal growth constraints of human systems*, (soumis).
- **J.-G. Young**, L. Hébert-Dufresne, A. Allard & L. J. Dubé, *Structural preferential attachment of community structure and its relation to Dunbar's number*, (en préparation, voir chapitre 5).
- L. Hébert-Dufresne, E. Laurence, A. Allard, **J.-G. Young** & L. J. Dubé, *Complex networks as an emerging property of hierarchical preferential attachment*, (en préparation).
- **J.-G. Young**, A. Allard, L. Hébert-Dufresne & L. J. Dubé, *Unveiling hidden communities through cascading detection on network structures*, (en préparation, voir chapitre 2).

## Conférences

- **J.-G. Young**, L. Hébert-Dufresne, A. Allard & L. J. Dubé, *Structural preferential attachment of community structure and its relation to Dunbar's number*, NetSci 2014, Berkeley, États-Unis.
- **J.-G. Young**, L. Hébert-Dufresne, A. Allard & L. J. Dubé, *Local and global solutions to community detection : when resolution matters*, NetSci 2013, Copenhague, Danemark.
- **J.-G. Young**, L. Hébert-Dufresne, A. Allard<sup>2</sup> & L. J. Dubé, *Unveiling hidden communities through cascading detection on network structures*, COMPLEX 2012, Sante-Fe, États-Unis.

---

2. Présentateur





# Introduction

## Réductionnisme et phénomènes émergents

L'approche réductionniste a longtemps été le principal moteur de la science moderne, et mène encore à des triomphes scientifiques de nos jours. Cette vision du progrès scientifique divise essentiellement les avancées en deux catégories : la recherche fondamentale et la recherche extensive [7]. Dans cette conception, la recherche fondamentale est celle qui est dédiée à repousser les limites de notre compréhension des lois de base de l'univers, alors que la recherche dite extensive consiste essentiellement à appliquer les nouvelles découvertes de la science fondamentale aux problèmes inexplicables des autres domaines de la science. Le réductionnisme admet donc implicitement une hiérarchie du savoir, ce qui se traduit par la division de la recherche en deux directions distinctes ; une vers une description de plus en plus fondamentale de l'univers et une vers une description des phénomènes concrets de plus haut niveau. Évidemment, cette conception extrême du réductionnisme ne survit pas au contact avec la recherche réelle, puisqu'une théorie des phénomènes de haut niveau exprimée dans le langage de la physique quantique n'est ni possible, ni utile. Une vision plus modérée consiste donc à relier les niveaux de descriptions par paires, à décrire des phénomènes compliqués à partir des théories s'appliquant aux éléments plus simples prenant part à ces phénomènes [56], e.g. la température d'un gaz par des collisions atomiques, la résistance du SIDA aux vaccins par la mutation de l'agent infectieux. On dira donc que le réductionnisme approche la science du bas vers le haut, du fondamental vers le concret.

Bien que très efficace, l'approche réductionniste ne permet pas d'appréhender l'ensemble du monde qui nous entoure [11]. En effet, il est facile d'identifier des phénomènes où une compréhension des parties ne permet pas de décrire directement l'ensemble, où l'ensemble est plus grand que la somme des parties. On pensera par exemple au paramagnétisme, à la conscience, aux fluctuations des marchés boursiers, etc. Ces phénomènes dits *émergents* ont mené à l'apparition d'une approche opposée au réductionnisme, dite holiste, où un système est considéré comme un tout indivisible. Cette approche brille lorsqu'elle est appliquée à des systèmes formés d'une multitude de sous-éléments, appelés *systèmes complexes* [99]. Holisme et systèmes complexes vont en fait de pair, puisque la diversité et la *complexité* de ces systèmes demandent une vision plus globale, ce qui garantit simultanément l'échec d'une description purement microscopique.

## Hiérarchie et structure des systèmes complexes

Il a toujours été difficile de définir précisément le concept de système complexe, un fait encore plus avéré lors des balbutiements de ce domaine d'étude. On pensera par exemple aux premières définitions qui plaçaient systématiquement la physique statistique du côté de la science de la complexité, à cause du *nombre* impressionnant de particules contribuant à l'émergence de phénomènes de masses. Il est rapidement devenu apparent que la *quantité* de sous éléments n'est pas un bon indicateur de la complexité d'un système. Après tout, un circuit de  $10^9$  résistances parallèles peut être à toutes fins pratiques considéré comme une seule résistance d'impédance quasi nulle.

C'est la hiérarchie, des éléments du système et non du savoir, qui a d'abord été dégagée comme propriété fondamentale des systèmes complexes [99]. Selon cette définition, tout système ayant une hiérarchie plus ou moins raffinée est complexe, que ce soit un gaz formé de particules elles-mêmes divisibles en sous parties sous-divisibles, ou une société civile formée de groupes maintes fois sous-divisée en petites cellules sociales. Manifestement, cette définition sous-tend un spectre de systèmes plus ou moins complexes à l'intérieur même de la science de la complexité ; la profondeur de la hiérarchie est alors un bon indicateur de complexité. Prenons à nouveau comme exemple un gaz : celui-ci est intuitivement considéré moins complexe qu'une (hypothétique) organisation sociale possédant le même nombre d'acteurs. Au premier abord, on pourrait penser que cette disparité est entièrement due à la position de ces systèmes dans la hiérarchie des domaines du savoir, que la description théorique des gaz est plus directe car ils sont plus rapprochés conceptuellement des lois naturelles régissant notre univers. Il s'agit toutefois d'une illusion réductionniste, puisqu'une description tout aussi simple des systèmes sociaux est fournie par l'approche holiste [111]. La différence cruciale entre ces deux systèmes réside plutôt dans le fait que le gaz est constitué de particules interagissant plus ou moins aléatoirement, alors qu'une organisation sociale est basée sur un enchevêtrement inextricable d'individus.

Par conséquent, si la hiérarchie définit les systèmes complexes, c'est la structure des interactions des sous-éléments du système qui les organise. Cette structure peut en fait être vue comme une *conséquence* de la hiérarchie, une observation très importante dans le contexte de la recherche, où la hiérarchie est souvent implicite et où seule la structure des interactions est connue. Cette dernière devient alors la propriété permettant de situer un système dans l'échelle de la complexité. En effet, si on reprend l'exemple des gaz, on constate que des interactions aléatoires, de proche en proche, suffisent à donner une description satisfaisante des phénomènes de masse (e.g. modèle de Ising, turbulence des fluides), alors que cette même hypothèse mène à une caricature de la réalité lorsqu'elle est appliquée aux systèmes sociaux (e.g. épidémiologie sur une population mélangée uniformément). Dans le deuxième cas, on doit explicitement tenir compte de la hiérarchie de la structure (structure non homogène).

## Naissance de la science des réseaux complexes

C'est en poussant ce raisonnement plus loin qu'il est devenu clair que *l'essence même* d'un système complexe est essentiellement capturée par sa structure [11]. Il n'est donc plus nécessaire de connaître la nature d'un système complexe pour tirer des conclusions sur sa fonction, son organisation, etc. Pour répondre à ce besoin de représenter la structure d'un système de façon abstraite, la science de la complexité s'est appropriée le concept de graphe, un objet mathématique représentant les constituants d'un système comme des noeuds et leurs interactions comme des liens, ignorant ainsi toute information superflue à l'exception de la structure des interactions.

C'est dans ce contexte que le concept réseau complexe a émergé et évolué pour finalement devenir un sujet d'étude à part entière. D'abord étudié en tant qu'objet mathématique ayant des propriétés régies par des théorèmes stricts, les graphes sont ainsi peu à peu devenus un simple moyen de représenter la structure hautement variable des systèmes réels, du moins dans le cadre de la science des réseaux. Le traitement déterministe des graphes a ainsi fait place à un traitement statistique, la structure exacte n'étant plus dans la mire de la recherche.

Du point de vue théorique, ce sont initialement des ensembles statistiques de réseaux qui ont été étudiés [35, 65]. Cette approche a permis d'isoler l'effet d'une série de propriétés structurelles sur, par exemple, la capacité des réseaux à supporter une propagation d'information régie par une dynamique donnée [70]. Parallèlement, grâce à l'explosion du domaine de l'acquisition de données, une multitude de réseaux réels sont devenus accessibles, permettant ainsi de comparer les résultats théoriques aux comportements de systèmes réels [79]. Bien que ces premiers résultats théoriques aient permis d'établir des méthodes très importantes qui sont en elles-mêmes d'intérêt [3], ils n'ont pas bien survécu à la comparaison avec les comportements réels [47].

## Structure communautaire

Parmi les causes de ces disparités, la description théorique trop simpliste de la structure a rapidement été identifiée comme étant la grande coupable. En effet, la hiérarchie, pourtant si centrale dans la définition initiale des systèmes complexes, avait initialement été laissée de côté. Elle se traduit par une densité d'interaction plus importante pour les éléments du système qui font partie des mêmes sous-groupes. Les premiers modèles supposaient une organisation beaucoup trop uniforme. Cette organisation a subséquentement été popularisée sous le nom de *structure communautaire*, un nom tiré de l'observation de ce type de structure dans les réseaux d'interactions sociales.

On peut attribuer l'omission initiale de la structure communautaire au flou entourant sa définition ainsi qu'au défi technique que présente l'extraction de cette structure. Même si le premier obstacle persiste, un certain consensus a été établi quant à la définition *approximative* de la structure communautaire : il s'agit d'une décomposition d'un réseau complexe en sous-ensembles d'éléments densément connectés (par rapport à la densité moyenne du réseau), avec ou sans recouvrement entre les sous-

ensembles [76]. Cette définition large permet d’englober plusieurs définitions *exactes* de la structure communautaire dans un même cadre conceptuel. Le deuxième obstacle est quant à lui intrinsèquement relié au concept de structure communautaire. En effet, la division d’un graphe en sous-ensemble des noeuds est typiquement un problème de décision difficile (au sens strict), de sorte que l’identification de la structure communautaire est un problème qui doit être abordé avec des méthodes approximatives ingénieuses.

## Vers une compréhension unifiée

Lorsque le problème de la *détection* de la structure communautaire a été relié à la question connexe et abondamment étudiée du partitionnement de graphes, la quantité de nouveaux algorithmes de *détection communautaire* a explosée [38]. Dans un contexte où autant la définition de la structure recherchée que la méthode employée pour l’extraire sont laissées libres, il est rapidement devenu difficile, voire impossible, de comparer ces algorithmes. La variété des mécanismes de détection, souvent introduits de façon *ad hoc*, a compliqué cette tâche. Jusqu’à tout récemment, les seuls outils permettant cette comparaison étaient ainsi des bancs d’essai numériques incluant des biais difficiles à identifier.

En réponse à ce problème important, le présent mémoire établit les bases d’outils heuristiques, analytiques et numériques offrant une compréhension *unifiée* de la théorie de la détection communautaire (l’objectif étant de traiter les différents algorithmes dans un même cadre conceptuel). L’ouvrage est divisé comme suit. Au chapitre 1, on établit formellement tous les concepts nécessaires à l’étude de la structure communautaire des réseaux complexes. On introduit ensuite une méthode méta-heuristique permettant d’améliorer les performances de tout algorithme de détection imparfait au chapitre 2. Au chapitre 3, les bases d’un formalisme analytique unifié de détection communautaire sont établies. Une approche matricielle y est préconisée. Le calcul du spectre des matrices introduites au chapitre 3 est effectué au chapitre 4, et confirme la validité du formalisme pour une classe idéalisée de réseaux complexes. Au chapitre 5, un processus de croissance local et réaliste est présenté puis intégré à un processus de croissance déjà existant, ce qui pose les bases d’un banc d’essai numérique extrêmement polyvalents et exempt de biais. Cette approche numérique vient compléter les résultats analytiques et heuristiques des chapitres précédents. En conclusion, les grandes lignes des développements futurs nécessaires pour pousser ces outils à maturité sont établies.

# Chapitre 1

## Introduction aux réseaux complexes et à leur structure communautaire

L'étude de la structure communautaire des réseaux complexes est une question relativement récente et éloignée du champs d'application traditionnel de la physique théorique. Un physicien typique sera familier avec les *méthodes mathématiques* employées dans cet ouvrage, mais généralement étranger aux problèmes auxquels elles sont appliquées. Ce premier chapitre se veut donc une introduction aux concepts de réseaux complexes (Sec. 1.1) et de leur organisation communautaire (Sec. 1.2).

### 1.1 Notions élémentaires

La science des réseaux s'intéresse à l'étude de la structure des systèmes complexes [13, 66, 75, 102]. Plus spécifiquement, elle s'intéresse à la *structure des interactions* des *constituants* d'un système complexe. Dans le cadre de cette science, ces deux concepts sont abstraits conjointement à l'aide de graphes [15], ou réseaux<sup>1</sup>. Le processus d'abstraction permet de traiter des systèmes de natures variées – e.g. le monde académique [81, 83], des organisations sociales [2, 23, 44, 61] ou conceptuelles [81], des systèmes technologiques [12, 17, 42, 44, 55, 82, 91, 110] – à l'aide d'un même ensemble d'outils. Le passage vers la représentation en graphes est un travail de modélisation en soi, influencé tant par la nature du système à l'étude que les choix du modélisateur.

Afin de bien illustrer les concepts de *constituants* d'un système et d'*interactions*, considérons par exemple la modélisation du système technologique qu'est Internet<sup>2</sup>. Plusieurs niveaux de description s'offrent au modélisateur, ce qui détermine la *nature des constituants* du système. En première approximation, on pourrait décrire ce système au niveau de la dorsale<sup>3</sup>. On considèrera alors que les serveurs à haut débit connectés internationalement sont les constituants de base, ce qui nous donnera

---

1. Nomenclature utilisée ci-après de façon interchangeable.

2. Système informatique sur lequel est bâti le web [75].

3. De l'anglais, *Internet backbone* : Infrastructures à longue distance et haut débit d'Internet (Office québécois de la langue française).

accès à une bonne description du système au niveau global. Si une représentation au niveau des systèmes autonomes est désirée, il faut alors considérer que les serveurs des fournisseurs d'accès internet sont aussi des éléments de base du système. Si l'objectif est d'analyser le réseaux de distribution d'information pair à pair, il faut alors repousser la précision de la description et inclure les clients des fournisseurs Internet comme éléments de base, et ainsi de suite.

La *nature des interactions* dépend quant à elle du niveau de description sélectionné. Par exemple, pour une représentation d'Internet au niveau de la dorsale, on peut décréter que deux serveurs sont en interaction s'ils sont reliés par un câble de fibre optique. Ce choix est évident puisque Internet prend directement la forme d'un réseau d'interaction. L'exercice est toutefois plus subtil et subjectif dans le cas de systèmes complexes différents. Par exemple, dans le cadre de l'analyse d'un réseau de contacts sociaux, les interactions peuvent être définies comme des contacts physiques, des relations d'amitiés, des relations professionnelles, etc.

### 1.1.1 Graphes et représentations

L'étape suivante consiste à représenter les éléments de ce système et les interactions entre ces éléments à l'aide de graphes.

**Définition 1.** Un graphe  $G = (\mathcal{V}, \mathcal{E})$  est un ensemble de  $N$  noeuds  $\mathcal{V}$  (de l'anglais vertices) et de  $M$  liens  $\mathcal{E}$  (edges) connectant les noeuds entre eux. On réfère aux noeuds à l'aide d'indice latin  $i, j, k, \dots$ . On réfère aux liens à l'aide de la notation  $e_{ij}$ , qui dénote un lien connectant les noeuds  $i$  et  $j$ .

Dans l'exemple d'Internet modélisé à l'échelle de la dorsale, chaque serveur est représenté par un noeud, tandis que les liaisons par fibres optiques sont représentées à l'aide de liens entre ces noeuds (Fig. 1.1).

La structure d'un graphe est ensuite exprimée mathématiquement à l'aide de matrices d'interactions, dont la *matrice d'adjacence* est l'exemple le plus répandu.

**Définition 2.** La matrice d'adjacence  $\mathbf{A}$  est une matrice  $N \times N$  dont les colonnes et rangées sont associées aux noeuds du graphe  $G = (\mathcal{V}, \mathcal{E})$ . L'élément<sup>4</sup>  $a_{ij}$  de cette matrice donne le poids de l'interaction entre les noeuds  $i$  et  $j$ .

Les interactions des constituants des systèmes complexes sont classifiées à l'aide de deux critères ; elles sont *pondérées* ou non, ainsi qu'*orientées* ou non. Sur la base de ces propriétés, on distingue quatre classes de réseaux et donc quatre classes de matrices d'adjacences (Fig. 1.2).

Un réseau est non dirigé lorsque les noeuds sont reliés par des liens uniques et sans poids. Les relations entre les éléments du système sont alors binaires ( $\exists$  ou  $\nexists$ ) de sorte que  $a_{ij} = 0, 1$  uniquement. Au contraire, un réseau est dit pondéré lorsqu'une force est associée à chaque interaction. Notons qu'on peut réinterpréter les réseaux où des liens multiples existent comme des réseaux pondérés [71].

---

4.  $a_{ij}$  est un élément de matrice donnant un poids, alors que  $e_{ij}$  fait référence à un objet mathématique (un lien).

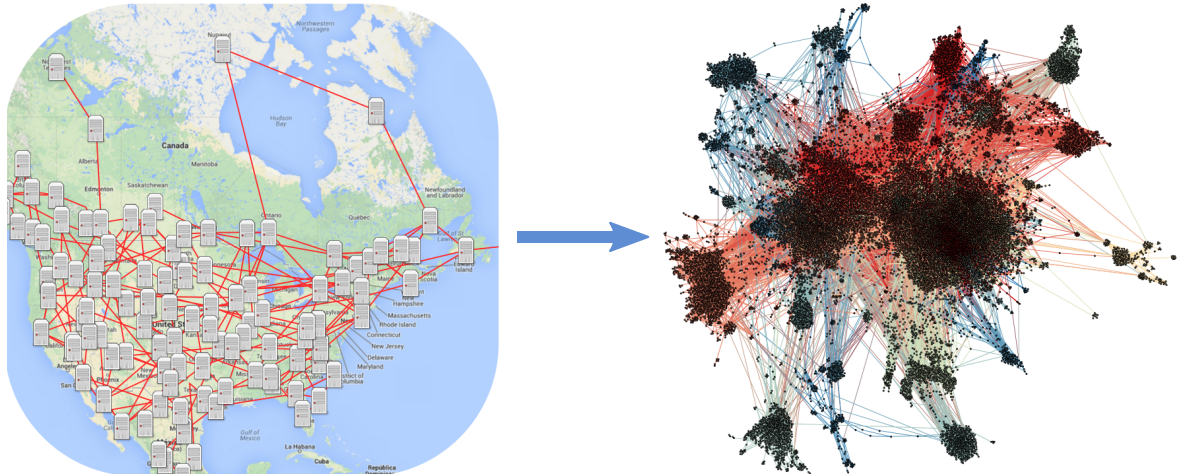


FIGURE 1.1 – **Réseau dorsal d’Internet circa 2006.** (gauche) Représentation schématique d’une partie du système complexe qu’est la dorsale d’Internet. (droite) Abstraction de la structure de système à l’aide d’un réseau complexe non dirigé et non pondéré. Chaque noeud correspond à un serveur, et chaque lien à la présence d’une connexion entre ces serveurs, tel que recensé dans les tables BGP de *archive.routeviews.org* [44]. La coloration indique des ensembles de noeuds (serveurs) structurellement rapprochés.

Ainsi, après normalisation, les éléments de la matrice d’adjacence des graphes pondérés et des graphes à liens multiples sont tous dans l’intervalle  $[0, 1]$ .

Lorsque les interactions sont réciproques (e.g. le serveur  $A$  est connecté au serveur  $B$ , et vice-versa), le réseau associé est dit non orienté. Ceci se traduit par une matrice d’adjacence symétrique, i.e.  $a_{ij} = a_{ji} \forall i, j$ . Lorsque les interactions sont dirigées (e.g. une protéine  $A$  active une protéine  $B$ , alors que l’inverse est faux), le réseaux associé est dit orienté. La directionnalité des interactions est représentée en introduisant une asymétrie dans la matrice d’adjacence, i.e.  $a_{ij} \neq a_{ji}$ . Par convention [75], l’élément  $a_{ij}$  réfère à un lien *quittant* le noeud  $j$  et étant dirigé *vers* le noeud  $i$ .

En outre, un réseau peut compter des boucles – des liens connectant un noeud à lui-même – qui se traduisent par une diagonale non nulle pour la matrice d’adjacence  $\mathbf{A}$ . On parle de *graphes simples* lorsque les liens sont non dirigés, non pondérés, et que les boucles sont interdites. Dans le cadre de cet ouvrage, les réseaux seront simples sauf indication contraire. L’extension au cas pondéré et avec boucle est la plupart du temps triviale, tandis que l’extension au cas dirigé demande un effort analytique important, car cette modification se traduit par une brisure de la symétrie de la matrice  $\mathbf{A}$ .

### 1.1.2 Propriétés structurelles versus description complète

Les graphes permettent de représenter la structure d’un ensemble extrêmement varié de systèmes. Cette structure peut être exprimée sans ambiguïté à l’aide de matrices d’adjacence, qui contiennent toute l’information structurelle du réseau. Toutefois, il est souvent plus utile de décrire un graphe à



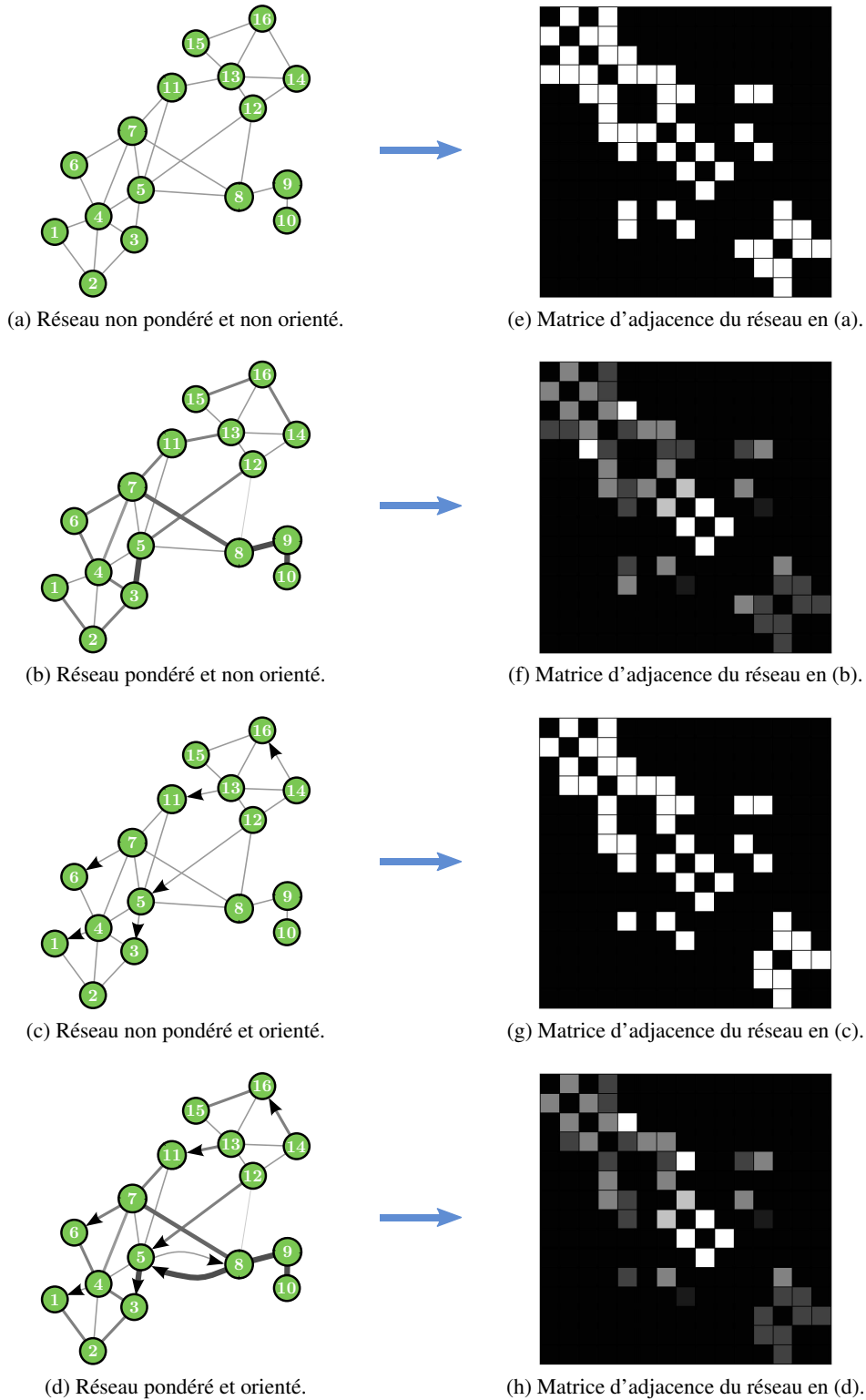


FIGURE 1.2 – **Classes de réseaux et matrices d'adjacences associées.** Dans le cas des réseaux dirigés, un lien est considéré réciproque si les flèches sont omises. Les matrices d'adjacence  $\mathbf{A}$  sont de dimension  $16 \times 16$  et l'intervalle  $a_{ij} \in [0, 1]$  est représenté par un dégradé de noir vers blanc.



l'aide d'une série de propriétés robustes qui permettent simplement de qualifier sa structure *effective*.<sup>5</sup> En fait, si on considère que les diverses propriétés – plutôt que la matrice d'adjacence – définissent un réseau, on ouvre la porte à une approche théorique très puissante [4, 47, 80], soit l'étude d'*ensembles* de réseaux dont certaines propriétés sont fixées et où le reste de la structure est laissé au hasard<sup>6</sup>. Formellement, on étudie alors des ensembles de matrices d'adjacence reliées par une relation d'équivalence.

On divise ces propriétés en trois groupes : microscopique, mésoscopique et macroscopique, sur la base de la quantité d'information nécessaire à leur calcul. Les propriétés microscopiques (locales) sont spécifiées par le *voisinage immédiat* des noeuds, où le voisinage  $n(i)$  du noeud  $i$  est constitué de tous les noeuds directement connecté à ce dernier. Les propriétés macroscopiques (globales) demandent une connaissance *complète* du réseau, alors que les propriétés mésoscopiques (intermédiaires) se situent – comme le nom l'indique – entre les deux niveaux de description.

Cet ouvrage portant principalement sur le niveau mésoscopique d'organisation, on commence par donner quelques exemples de propriétés utiles des niveaux micro et macroscopique, avant de dédier une section entière à la structure mésoscopiques (Sec. 1.2).

### 1.1.3 Mesures locales et globales : une liste non exhaustive

La séquence des degrés d'un réseau complexe est sans doute sa propriété locale la plus importante. Elle joue un rôle central tant pour la description de la structure des réseaux réels [12, 65, 75] que pour l'étude de la propagation d'information sur réseaux complexes [46, 70].

**Définition 3.** Le degré  $k_i$  est le nombre de voisin du noeud  $i$ . En termes de la matrice d'adjacence, ce degré est donné par la somme

$$k_i := \sum_{j=1}^N a_{ij} \equiv \sum_{j=1}^N a_{ji}, \quad (1.1.1)$$

où l'équivalence tient uniquement pour les réseaux non dirigés ( $\mathbf{A}$  symétrique). Voir la figure 1.3 pour un exemple visuel du degré.

Le vecteur  $\mathbf{k}$  (i.e. la séquence) contenant tous les degrés du réseau est donc obtenu en appliquant un vecteur de sommation (cf. glossaire des notations) sur la matrice d'adjacence

$$\mathbf{k} = \mathbf{A}\mathbf{1}. \quad (1.1.2)$$

5. Le parallèle avec la physique statistique est évident : l'énergie moyenne d'un gaz est beaucoup plus parlante qu'une liste exhaustive des positions et vitesses de chacune des particules.

6. Encore une fois, l'analogie avec la physique statistique est directe. Dans l'approche canonique, l'énergie moyenne et les niveaux d'énergies accessibles sont fixes tandis que le peuplement des états est déterminé au hasard (suivant un facteur de Boltzmann).

La logique impose que le nombre *total* de degrés soit pair, afin que chaque *lien* soit complet (un lien est constitué de deux degrés : un pour chaque extrémité), ce qui se traduit mathématiquement par

$$\sum_{i=1}^N k_i = 2M, \quad (1.1.3)$$

avec  $M$  la cardinalité de  $\mathcal{E}$  (nombre de liens). Alternativement, en notation matricielle, le nombre total de degrés est donné par  $\mathbf{1}^T \mathbf{k} = \mathbf{1}^T \mathbf{A} \mathbf{1} = 2M$ , tel que le degré moyen est

$$\langle k \rangle = \frac{2M}{N} = \frac{1}{N} \mathbf{1}^T \mathbf{A} \mathbf{1}. \quad (1.1.4)$$

Puisqu'il est nécessaire de recourir à la matrice d'adjacence complète pour calculer  $M$  et  $\langle k \rangle$ , on peut qualifier ces mesures de macroscopiques, bien qu'une évaluation approximative de  $\langle k \rangle$  ne demande qu'une connaissance limitée du réseau. La nature macroscopique du degré *moyen* est bien illustrée par le fait que la *densité* du réseau – une autre quantité macroscopique – peut être obtenue à l'aide de la même information.

**Définition 4.** La densité  $\rho$  d'un réseau est le rapport du nombre de liens existants sur le nombre de liens possibles. Un réseau de  $N$  noeuds peut compter jusqu'à  $\binom{N}{2}$  liens, tel que  $\rho$  est donné par

$$\rho = \frac{M}{\binom{N}{2}} = \frac{2M}{N(N-1)} = \frac{1}{N(N-1)} \mathbf{1}^T \mathbf{A} \mathbf{1} \simeq \frac{1}{N^2} \mathbf{1}^T \mathbf{A} \mathbf{1}, \quad (1.1.5)$$

où la dernière égalité tient rigoureusement dans la limite  $N \rightarrow \infty$ .

Par construction,  $\rho$  est borné à l'intervalle  $[0, 1]$ . Un réseau est dit creux<sup>7</sup> lorsque la densité est faible ( $\rho \sim 0$ ), et dense dans le cas inverse<sup>8</sup> ( $\rho \sim 1$ ).

On pourrait s'attendre à ce que les réseaux réels, qui sont généralement creux (voir tableau 2.1, par exemple) aient une *structure quasiment en arbre*<sup>9</sup> et donc que seuls quelques parcours différents existent entre toute paire de noeuds. En pratique, les liens des réseaux réels sont toutefois redondants et ce malgré la faible densité de ces derniers. Formellement, cette redondance des liens peut être exprimée en termes de *corrélations* [69, 107, 110] : le fait que deux noeuds soient connectés influence positivement la probabilité que les noeuds de leurs voisinages soient également connectés. Cette corrélation se traduit par une grande abondance de *triangles* dans les réseaux réels.

**Définition 5.** Un triangle est un groupe de trois noeuds complètement connectés à l'aide de trois liens.

---

7. En anglais : *sparse*.

8. Certains auteurs [67, 75] définissent les graphes creux comme ceux pour lesquels  $\rho \rightarrow 0$  lorsque  $N \rightarrow \infty$ . Cette dichotomie requiert toutefois que la limite  $N \rightarrow \infty$  puisse être prise, e.g. que le réseau soit tiré d'un ensemble statistique permettant le calcul analytique de  $\rho$ , ou que des instances de tailles différentes du même réseau soient connues.

9. Arbre : Graphe connexe sans cycle et lien multiple, i.e. graphe connexe de densité minimale [15].

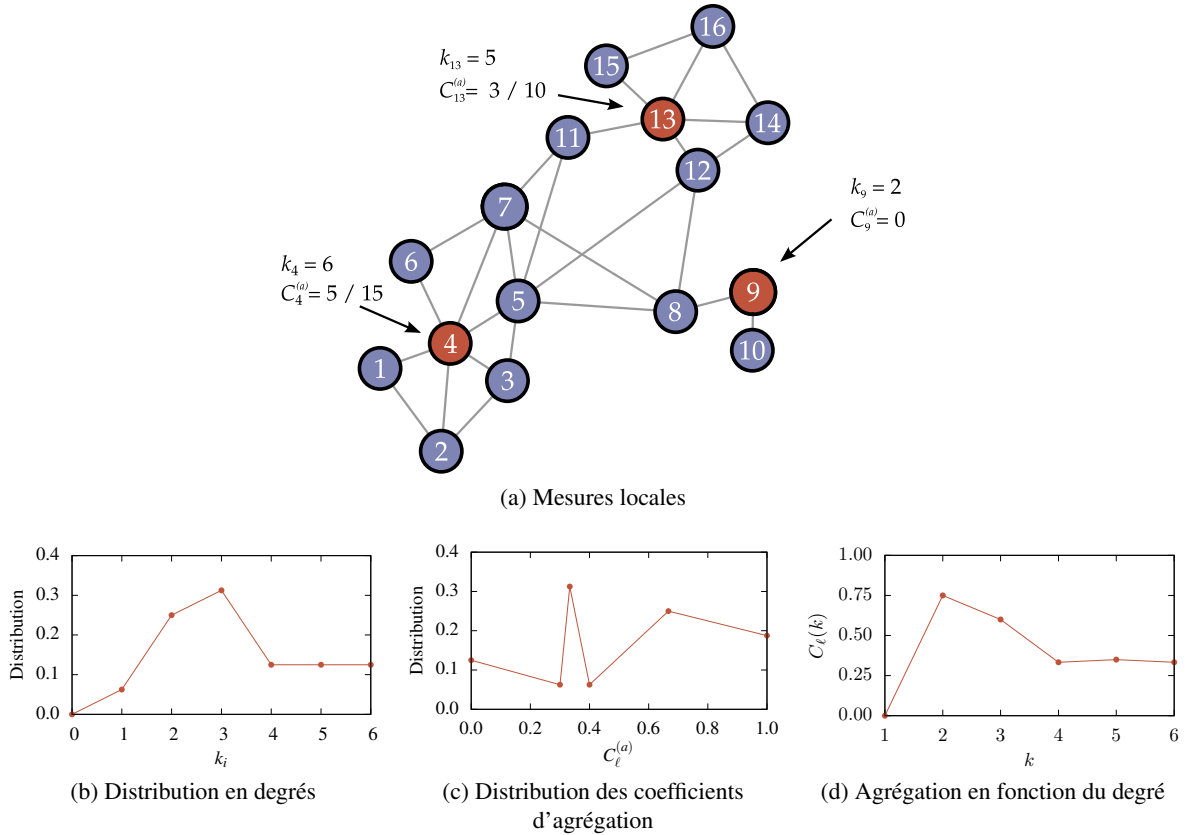


FIGURE 1.3 – **Propriétés structurelles locales.** (a) Réseau artificiel non dirigé et non dirigé de  $N = 16$  noeuds et de  $M = 27$  liens introduit à la Fig. 1.2. Ce réseau possède une densité  $\rho = 0.225$  et un ratio de transitivité  $C^{(t)} = 0.40741$ . (b) Distribution en degrés de moyenne  $\langle k \rangle = 3.3750$ . (c) Distribution des coefficients d'agrégation de moyenne  $\langle C^{(a)} \rangle = 0.50208$ . (d) Relation entre le degré des noeuds et leur coefficient d'agrégation.

La corrélation entre les liens peut donc être calculée de façon indirecte par le biais du *coefficient d'agrégation*, qui mesure la densité de triangles d'un réseau. Il existe deux définitions différentes de ce coefficient dans la littérature [50, 109], qu'on donne ici sans discrimination, car elles sont toutes deux utiles dans des contextes différents.

**Définition 6.** *Le ratio de transitivité d'un réseau  $C^{(t)}$  est le rapport du nombre de triangles  $n_{\Delta}$  sur le nombre de triades (triangles dont un lien est manquant)  $n_t$*

$$C^{(t)} = 3 \times \frac{n_{\Delta}}{n_t} = \frac{3 \times \left[ \frac{1}{6} \text{Tr}(\mathbf{A}^3) \right]}{\frac{1}{2} \left[ \mathbf{1}^T \mathbf{A}^2 \mathbf{1} - \text{Tr}(\mathbf{A}^2) \right]}, \quad (1.1.6)$$

où le facteur 3 est dû au fait qu'un triangle est constitué de trois triades, et où on utilise la propriété  $[\mathbf{A}^{\ell}]_{ij} \equiv$  nombre de parcours dirigé de longueur  $\ell$  allant du noeud  $j$  au noeud  $i$  [75].

Le numérateur de l'eq. (1.1.6) compte le nombre de parcours de longueur 3 partant et revenant au même noeud (d'où l'utilisation de l'opérateur trace). La division par 6 apparaît car chaque triangle

introduit 6 parcours dirigés distincts. La même logique est utilisée pour construire le dénominateur, à la différence que les triades forment des parcours quittant un noeud et arrivant à un noeud différent (parcours de deux liens), tel qu'on doit sommer tous les éléments de  $\mathbf{A}^2$  sauf ceux sur la diagonale.

**Définition 7.** *Le coefficient d'agrégation du noeud  $\ell$  est le nombre de paires de voisins connectés de  $\ell$  normalisée par le nombre de paires de voisins de  $\ell$*

$$C_\ell^{(a)} = \frac{\sum_{i \neq j} (a_{\ell i} a_{\ell j}) a_{ij}}{\sum_{i \neq j} (a_{\ell i} a_{\ell j})} = \frac{\sum_{i,j=1}^N (a_{\ell i} a_{\ell j}) a_{ij}}{\sum_{i,j=1}^N (a_{\ell i} a_{\ell j}) - \sum_{i=1}^N a_{\ell i}^2}. \quad (1.1.7)$$

On utilise le fait que les entrées de la matrice  $\mathbf{A}$  sont binaires, tel que  $(a_{\ell i} a_{\ell j}) a_{ij}$  est égale à un s.s.i le noeud  $\ell$  est connecté aux noeuds  $i$  et  $j$  et que ces deux noeuds sont également connectés. Voir la figure 1.3 pour un exemple visuel du coefficient d'agrégation.

**Définition 8.** *Le coefficient d'agrégation du réseau est la moyenne des coefficients locaux*

$$\langle C^{(a)} \rangle = \frac{1}{N} \sum_{\ell=1}^N C_\ell^{(a)}. \quad (1.1.8)$$

Le ratio de transitivité (1.1.6) présente un fort intérêt théorique car il est déterminé par une expression matricielle pouvant être calculée dès que les traces des puissances de la matrice d'adjacence sont connues (on pourrait par exemple utiliser les résultats des références [68, 117]). En pratique, le coefficient d'agrégation (1.1.7) est toutefois la mesure de redondance la plus utilisée [75]. En effet, la moyenne  $\langle C^{(a)} \rangle$  donne directement la probabilité qu'une triade soit complétée, i.e. la force de la corrélation induite par l'existence d'un lien pour un réseau donné. De plus, le coefficient d'agrégation est propre aux noeuds plutôt qu'à l'entièreté du réseau, ce qui donne accès à beaucoup d'information, sous la forme de relations fonctionnelles, e.g. la relation entre les coefficients  $C_\ell^{(a)}$  et le degré des noeuds [88].

Bien que les deux précédentes mesures d'agrégation quantifient la redondance du voisinage d'un noeud via un décompte de triangles, il ne s'agit pas du tout de la seule approche possible [75]. On peut par exemple s'intéresser à la *similarité* du voisinage de deux noeuds [1].

**Définition 9.** *La similarité du voisinage de deux noeuds  $i$  et  $j$  est définie comme la taille du voisinage commun de  $i$  et  $j$  normalisée par la taille de leur voisinage total :*

$$S(i, j) = \frac{|n_+(i) \cap n_+(j)|}{|n_+(i) \cup n_+(j)|} \equiv \frac{\sum_{l=1}^N a_{il} a_{jl} + 2a_{ij} + \delta_{ij}}{\sum_{l=1}^N (a_{il} + a_{jl}) - \sum_{l=1}^N a_{il} a_{jl} + 2(1 - a_{ij}) - \delta_{ij}}, \quad (1.1.9)$$

où  $n_+(i)$  est l'ensemble formé par l'union du noeud  $i$  et de son voisinage, et où  $|X|$  dénote la cardinalité de l'ensemble  $X$ . Notons que l'expression de gauche est un indice de Jaccard. Au numérateur, la première somme donne le nombre de voisins communs, tandis que les termes additionnels tiennent compte des cas spéciaux (le cas où le lien direct entre  $i$  et  $j$  existe et le cas où  $i = j$ , respectivement). Au dénominateur, la première somme compte tous les voisins des noeuds  $i$  et  $j$  sans discrimination,

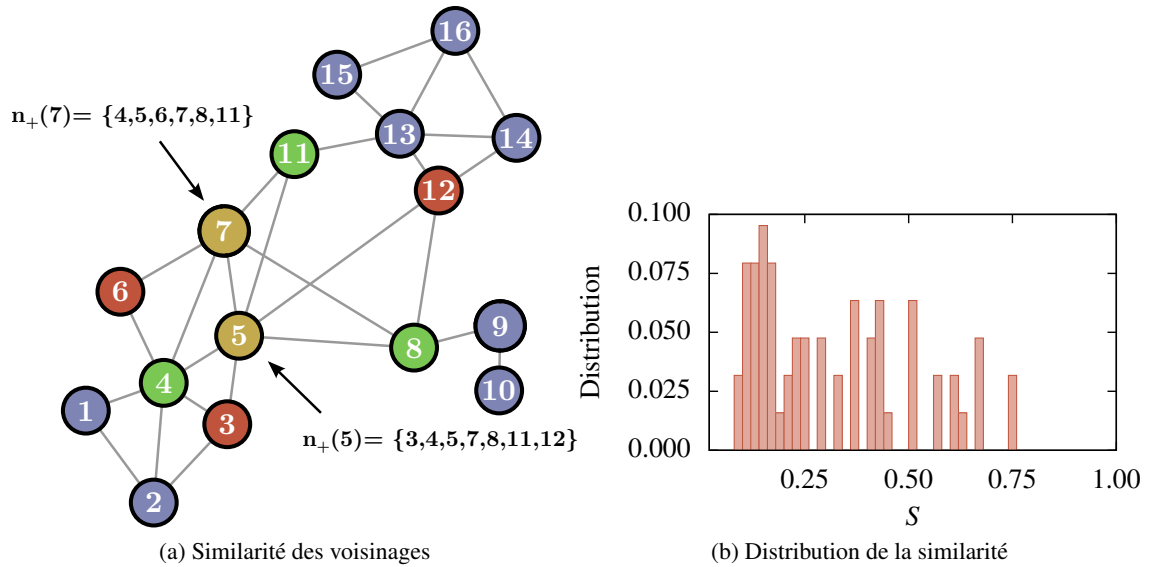


FIGURE 1.4 – **Similarité des voisinages.** (a) Exemple visuel du calcul de la similarité, pour un réseau artificiel. Le voisinage commun des noeud 5 et 7 est coloré en vert (et jaune), alors que les noeuds faisant partie uniquement du voisinage du noeud 5 ou 7 sont colorés en rouge. Ici, l'union est de taille 8 et l'intersection de taille 5, tel que  $S(5,7) = 5/8$ . Par extension, les paires de liens  $S(e_{4,5}, e_{4,7}), S(e_{11,5}, e_{11,7}), S(e_{8,5}, e_{8,7})$  partagent cette similarité. On note que ces paires de liens sont celles qui relient les noeuds 5 et 7 (en jaune) aux noeuds de l'intersection de leurs voisinages (en vert). (b) Distribution de la similarité pour toutes paires de noeuds partageant au moins un voisin commun.

alors que la deuxième somme retire les termes inclus deux fois. Encore une fois les termes additionnels tiennent compte des cas spéciaux (lien  $e_{ij} \exists$  et le cas  $i = j$ ). Voir la figure 1.4 pour un exemple visuel du calcul de la similarité.

Cette mesure de similarité, pourtant propre aux noeuds, est généralement utilisée pour quantifier la similarité de *paires de liens* [1]. Ainsi, on définit la similarité de deux liens  $e_{ki}, e_{kj}$  partageant un noeud commun  $k$  (appelé clé de voûte, de l'anglais *keystone* [1]), comme la similarité des voisinages  $n_+(i), n_+(j)$ , i.e.

$$S(e_{ki}, e_{kj}) := S(i, j) \quad (1.1.10)$$

#### 1.1.4 Un point de comparaison : réseaux aléatoires

Afin d'apprécier à quel point la redondance des liens et des noeuds quantifiée à travers  $\{C_\ell^{(a)}\}$ ,  $C^{(t)}$  et  $\{S(e_{ki}, e_{kj})\}$  est une caractéristique importante et surprenante des réseaux réels, on compare ces derniers à des *réseaux aléatoires équivalents*. Un réseau aléatoire équivalent est obtenu en brassant les liens d'un réseau donné, tout en conservant une série de propriétés fixes. Les procédures de brassage produisent donc des *ensembles* de réseaux sans structure, à l'exception de celle qui est explicitement imposée par contrainte. Dans cette section, on considèrera deux ensembles classiques, à savoir un ensemble où le nombre de liens est conservé, et un ensemble où les degrés des noeuds sont conservés.

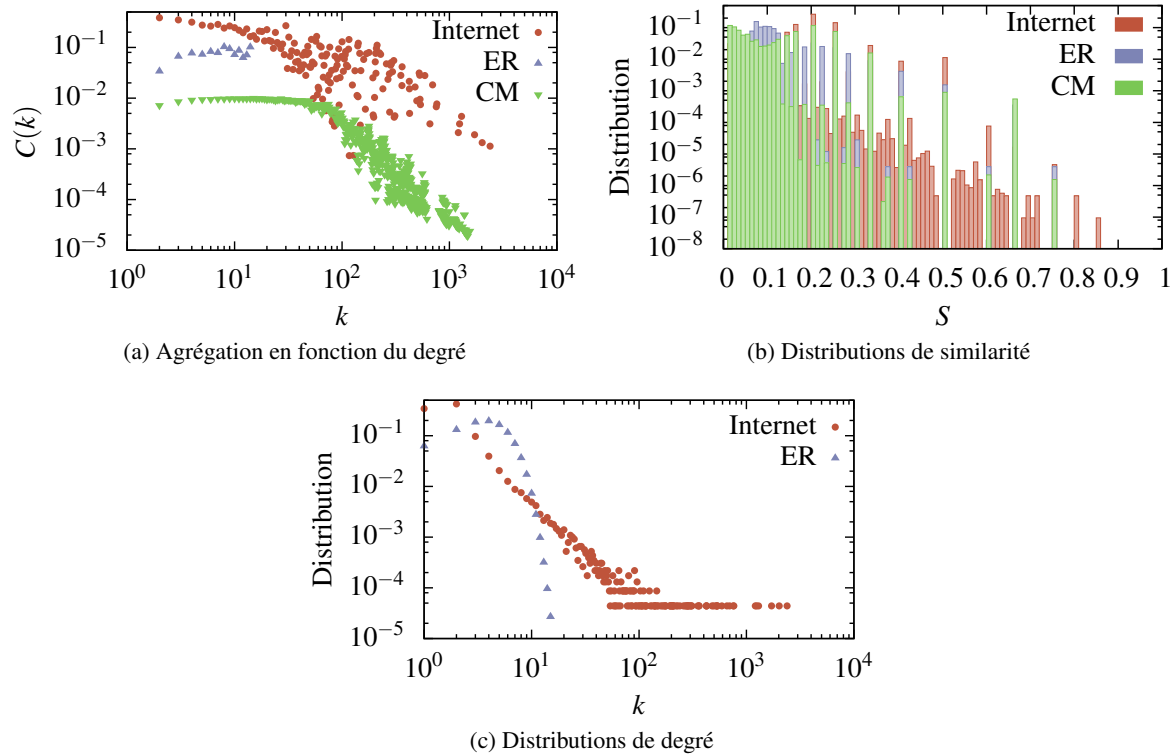


FIGURE 1.5 – **Redondance dans un réseau technologique réel et dans ses équivalents aléatoires.** On compare le réseau de la dorsale d’Internet (rouge) aux réseaux Erdős-Rényi de degré moyen identique (bleu) et aux réseaux de configuration ayant la même séquence degré (vert). Les résultats sont moyennés sur 25 réalisations de chaque réseau aléatoire. (a) Relation moyenne  $C(k)$  entre le coefficient d’agrégation local  $C_\ell^{(a)}$  et le degré  $k_\ell, \forall \ell$ . (b) Distribution de la similarité pour toutes paires de noeuds partageant au moins un voisin commun. (c) Distribution des degrés de la dorsale d’Internet et du réseau ER équivalent.

Pour effectuer un brassage aléatoire avec *conservation de liens*, on déconnecte d’abord tous les liens du réseau. Ensuite, pour chaque lien supprimé, on connecte une paire de noeuds choisie de façon uniforme parmi l’ensemble  $\mathcal{V}$ . Formellement, cette procédure est équivalente à la construction de graphes Erdős-Rényi (ER) [35], ou graphes aléatoires. Plus spécifiquement, il s’agit de la variante qu’on pourrait dire *micro-canonique*, dans laquelle le nombre de noeuds et de liens est exactement fixé<sup>10</sup>. Un graphe aléatoire équivalent peut donc être obtenu plus directement en tirant un réseau au hasard dans l’ensemble des graphes ER de  $N$  noeuds et  $M$  liens.

Le brassage aléatoire avec conservation des degrés est un processus plus subtil. La technique consiste à couper les liens en deux, créant ainsi des demi-liens. On reconnecte ensuite les demi-liens, au hasard, ce qui produit un graphe qui conserve la séquence de degrés initial  $\mathbf{k}$ , sans posséder la redondance locale. Cette seconde procédure est équivalente<sup>11</sup> à la construction d’un graphe tiré du

10. Dans la variante *canonique*, l’ensemble des noeuds  $\mathcal{V}$  est gardé fixe, mais les liens sont placés avec une *probabilité* donnée. On fixe ainsi la moyenne du nombre de liens plutôt que la quantité exacte.

11. L’équivalence tient strictement dans la limite  $N \rightarrow \infty$ . En taille finie, des boucles et des liens parallèles apparaissent

modèle de configurations (CM), i.e. un graphe possédant une séquence de degrés fixe [65, 79].

Les graphes produits à l’aide des techniques introduites dans les précédents paragraphes étant intrinsèquement *aléatoires*, il devient important de considérer un grand ensemble de graphes et de conclure sur les propriétés moyennes, plutôt que sur les propriétés d’une réalisation donnée. On effectue ici cet exercice, en construisant des graphes aléatoires équivalents au réseau de la dorsale d’Internet, et en obtenant les comportements moyens de toutes les propriétés structurelles introduites à la Sec. 1.1.3 pour ces ensembles.

Dans le cas du brassage aléatoire (ER), on observe que la distribution en degrés n’est pas conservée (Fig. 1.5c) : la nouvelle distribution est binomiale plutôt qu’en loi de puissance. Le degré moyen reste toutefois le même, par construction. La distribution en degrés est évidemment conservée dans le cas du modèle de configuration (CM).

Dans les deux cas, l’agrégation locale, quoique présente, est beaucoup plus faible ( $\langle C^{(a)} \rangle_{\text{ER}} \approx 0.00013$  et  $\langle C^{(a)} \rangle_{\text{CM}} \approx 0.043$ ) que dans le réseau original ( $\langle C^{(a)} \rangle_{\text{réel}} \approx 0.23$ ), et ce pour tous les degrés (Fig. 1.5a). L’importante décroissance de la fonction  $C(k)$  pour  $k \gg 1$  est simplement due au fait qu’il devient difficile de maintenir un haut coefficient d’agrégation pour les noeuds ayant beaucoup de voisins. En effet, il faudrait que le réseau soit essentiellement complètement connecté pour que leur coefficient d’agrégation ne chute pas. Les similarités entre les comportements de  $C_{\text{CM}}(k)$  et  $C(k)$  sont donc dues à la nature même du coefficient d’agrégation plutôt qu’à des ressemblances structurelles. Notons que le ratio de transitivité, qui est analytiquement<sup>12</sup>  $\langle C^{(t)} \rangle \approx 1.83 \cdot 10^{-4}$  pour les réseaux sans corrélation (i.e. CM et ER), est aussi plus élevé dans le réseau original (de l’ordre de 0.001 [75, p.200]). L’absence de corrélation structurelle dans les réseaux aléatoires simples se manifeste aussi à travers une distribution de similarité concentrée vers le bas (Fig. 1.5b).

On notera qu’Internet est un réseau réel dont la structure est particulièrement peu redondante, ce qui est souvent observé pour les réseaux technologiques en général [75, § 8.6]. Les liens sont beaucoup plus corrélés dans les réseaux sociaux [78], comme en témoigne la figure 1.6, où le même exercice de comparaison a été effectué pour le réseau de collaboration scientifique construit à partir des articles de la section “Condensed matter” du site *arXiv.org* [81]. Le coefficient d’agrégation moyen du réseau réel est  $\langle C^{(a)} \rangle \approx 0.63$ , alors qu’on a  $\langle C^{(a)} \rangle_{\text{ER}} \approx 0.00024$  et  $\langle C^{(a)} \rangle_{\text{CM}} \approx 0.0013$  pour les versions aléatoires, soit des résultats inférieurs par 3 et 2 ordres de grandeur, contre 3 et 1 ordres de grandeur dans le cas d’Internet.

---

avec une probabilité  $\sim 1/N$ . On introduit généralement un léger biais en brassant les liens fautifs à nouveau.

12. Un noeud a en moyenne  $\langle k \rangle$  voisins. En absence de corrélations structurelles, la probabilité que deux des voisins d’un noeud donné soit mutuellement voisins est  $\langle k \rangle / N$ , car les voisins de ces noeuds sont choisis au hasard parmi l’ensemble des  $N$  noeuds du réseau. Ceci donne directement, par définition, le coefficient de transitivité.

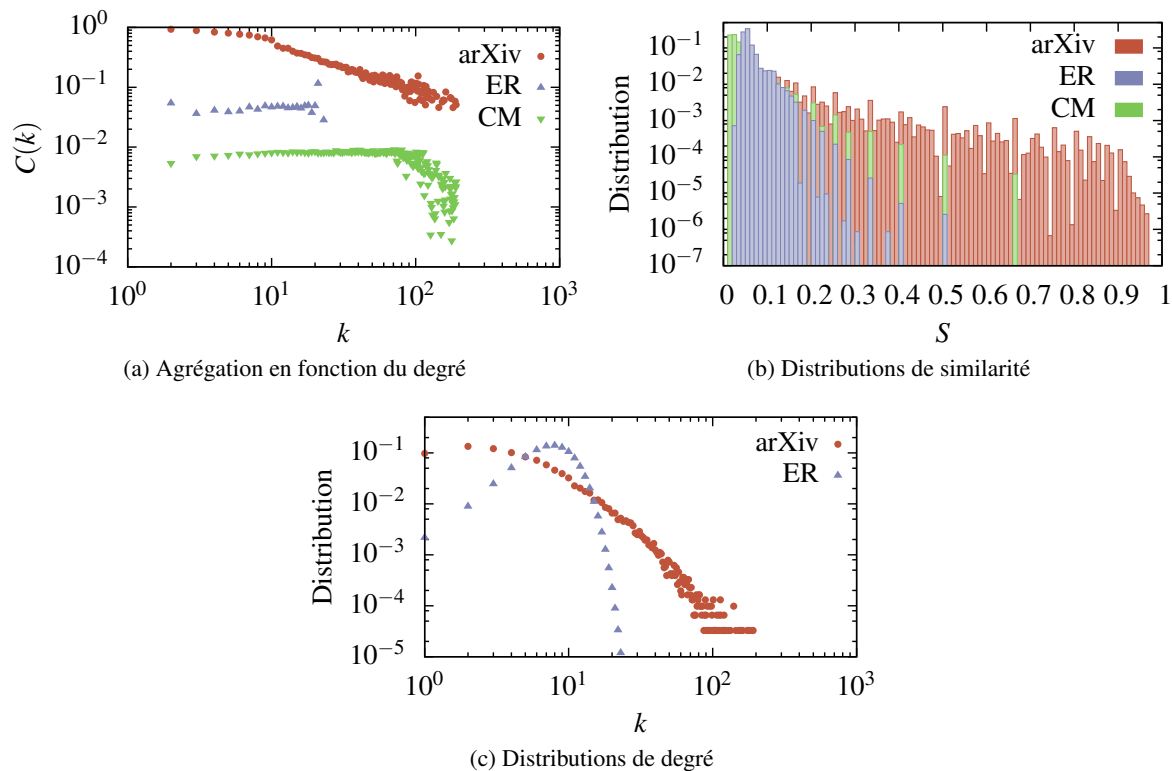


FIGURE 1.6 – **Redondance dans un réseau social réel et dans ses équivalents aléatoires.** On compare le réseau *arXiv cond-mat circa 2005* [81] (rouge) aux réseaux Erdős-Rényi de degré moyen identique (bleu) et aux réseaux de configuration ayant la même séquence degré (vert). Les résultats sont moyennés sur 25 réalisations de chaque réseau aléatoire. (a) Relation moyenne  $C(k)$  entre le coefficient d’agrégation local  $C_\ell^{(a)}$  et le degré  $k_\ell, \forall \ell$ . (b) Distribution de la similarité pour toutes paires de noeuds partageant au moins un voisin commun. (c) Distribution des degrés d’arXiv et du réseau ER équivalent.

## 1.2 Structure mésoscopique : réseaux modulaires

Comme le montre la Sec. 1.1.4, l’agrégation locale des liens et des noeuds est une caractéristique importante des réseaux réels. Ce point de vue purement local peut toutefois être trompeur. En effet, pour avoir une meilleure compréhension de la structure des réseaux complexes, il convient d’élargir notre champ de vision passablement, et d’adopter un point de vue à l’échelle *mésoscopique* : « au milieu, médian »<sup>13</sup>. Ainsi, au lieu de se restreindre au voisinage immédiat des noeuds, on considèrera leur voisinage plus ou moins rapproché, sans toutefois s’attarder au réseau complet. Il devient alors apparent que l’agrégation est en fait la signature d’un autre phénomène, soit le regroupement de noeuds en *communautés* [40, 72].

Ce choix curieux de nomenclature est dû aux origines sociologiques du concept de communauté. En effet, dès leur introduction en sociologie, les graphes ont menés à la formalisation du concept du

13. Petit Robert, Édition 1993



groupe social [18]. De multiples définitions de la communauté sociale ont émergé de cette formalisation [87], généralement basées sur l'idée implicite qu'une communauté est un ensemble de noeuds ou de liens densément connectés<sup>14</sup>. C'est ce passage d'une définition humaine (groupes d'individus s'identifiant comme amis, collègues, famille) à une définition mathématique (groupes densément connectés) qui a ultimement permis d'établir que la communauté est un élément d'intérêt du niveau mésoscopique de tout réseau, non pas uniquement des réseaux sociaux. Ainsi, pour des raisons historiques, le terme "communauté" est généralement utilisé en science des réseaux, même si on parle parfois de "groupes" (*cluster*) ou de "modules", des concepts analogues<sup>15</sup> introduits et étudiés dans d'autres contextes (e.g. recherche opérationnel ou informatique) [21].

### 1.2.1 Types d'organisations communautaires

De la même façon qu'il existe plusieurs sortes de réseaux complexes (dirigés, pondérés, etc.), plusieurs types d'organisations communautaires sont possibles. On distingue les types d'organisations communautaires sur la base de deux critères. La première distinction porte sur l'identité des éléments d'une communauté, alors que la deuxième et plus importante distinction est empruntée de la théorie des ensembles, et a trait à la multiplicité de l'appartenance des éléments d'une communauté.

#### Communautés de noeuds et de liens

On se rappellera qu'un graphe est un *couple*  $G = (\mathcal{V}, \mathcal{E})$  entre un ensemble de noeuds  $\mathcal{V}$  et des liens  $\mathcal{E}$ . Ceci suggère donc déjà deux types d'organisations communautaires de base, i.e. une organisation en groupes de noeuds [38] et une organisation en groupes de liens [1].

**Définition 10.** Une *décomposition en communautés de noeuds* du graphe  $G = (\mathcal{V}, \mathcal{E})$  est une collection  $\Omega$  de  $g$  sous-ensembles de  $\mathcal{V}$  telle que chaque noeud appartient à au moins un élément de  $\Omega$ . Formellement,  $\Omega = \{\Psi_\alpha\}$  avec  $\Psi_\alpha \subset \mathcal{V}$ , vérifiant la propriété suivante :

$$\forall v \in \mathcal{V}, \exists \alpha \text{ tel que } v \in \Psi_\alpha \quad (1.2.1)$$

Une définition similaire de la *décomposition en communautés de liens* est obtenue en substituant  $\mathcal{V}$  par  $\mathcal{E}$ .

La décomposition en communautés de noeuds est celle qui est la plus étudiée (voir les références de [38]), car ces derniers sont généralement les éléments qui sont au coeur de l'exercice de modélisation. Par exemple, décrire le réseau de la dorsale d'Internet (cf. Sec. 1.1) comme un ensemble de *groupes* de serveurs interconnectés est certainement plus intuitif que de le décrire à l'aide de groupes de connexions à fibre optiques. Toutefois, une description en termes de liens est plus pertinente dans certains contextes [1]. On pense par exemple aux groupes sociaux où le lien, plutôt que l'individu,

14. Définition valide pour les réseaux non dirigés pondérés ou non, qui sont l'objet de ce mémoire.

15. Le *data clustering* consiste à diviser des propriétés en groupes pertinents, à partir d'une matrice exprimant la similarité entre ces propriétés. Or, cette matrice peut être vue comme la matrice d'adjacence pondérée d'un graphe. Un *cluster* (groupe) de propriétés similaires correspond donc à une communauté dans le graphe.

définit la communauté. En effet, un individu peut assumer plusieurs rôles dans une société (ami, parent, collègue), de sorte qu'un groupe social n'est pas défini strictement par ses membres (noeuds) mais plutôt par les liens entre ces membres. L'existence de ce groupe se manifeste alors au niveau du réseau social correspondant par un regroupement dense de *liens*.

### Partitions et couvertures

Afin de d'établir la deuxième distinction, considérons un réseau de  $N$  noeuds et  $M$  liens décomposé en communautés (telles que définies à la Déf. 10). On peut alors se retrouver face à deux situations.

**Définition 11.** Une décomposition en communautés (de liens ou de noeuds)  $\Omega = \{\Psi_1, \Psi_2, \dots, \Psi_g\}$  du graphe  $G = (\mathcal{V}, \mathcal{E})$  est une partition  $\mathcal{P}$  si  $\Omega$  remplit les conditions

$$|\Psi_\alpha \cap \Psi_\beta| = 0 \quad \forall \alpha, \beta \quad (1.2.2a)$$

$$\Psi_1 \cup \Psi_2 \cup \dots \cup \Psi_g = \begin{cases} \mathcal{V} & \text{Communautés de noeuds} \\ \mathcal{E} & \text{Communautés de liens} \end{cases} \quad (1.2.2b)$$

Autrement dit, une partition est une décomposition où tout élément appartient obligatoirement à une seule et unique communauté. Dans la littérature de détection communautaire, les partitions incomplètes (i.e. qui ne satisfont pas la condition 1.2.2b) sont quand même désignées sous le nom de partition [38, § A.1]. Il s'agit d'une simple technicité, puisque cette disparité peut être réconciliée formellement en considérant que les éléments *non assignés*, sont implicitement assignés à des communautés distinctes d'un seul élément. Le critère central définissant la partition est donc la condition (1.2.2a). Un relâchement de ce critère mène à un deuxième type d'organisation communautaire.

**Définition 12.** Une décomposition en communautés (de liens ou de noeuds)  $\Omega$  du graphe  $G = (\mathcal{V}, \mathcal{E})$  est une couverture  $\mathcal{C}$  si  $\Omega$  remplit les conditions

$$|\Psi_\alpha \cap \Psi_\beta| \geq 0 \quad \forall \alpha, \beta \quad (1.2.3a)$$

$$\Psi_\alpha \neq \Psi_\beta \quad \forall \alpha, \beta \quad (1.2.3b)$$

$$\Psi_1 \cup \Psi_2 \cup \dots \cup \Psi_g = \begin{cases} \mathcal{V} & \text{Communautés de noeuds} \\ \mathcal{E} & \text{Communautés de liens} \end{cases} \quad (1.2.3c)$$

Dans le cas de la couverture, un élément *peut* appartenir à plus d'une communauté. À nouveau, la deuxième condition (1.2.3c) n'est pas très stricte, puisqu'elle peut être satisfaite facilement en ajoutant des communautés artificielles de cardinalité 1.

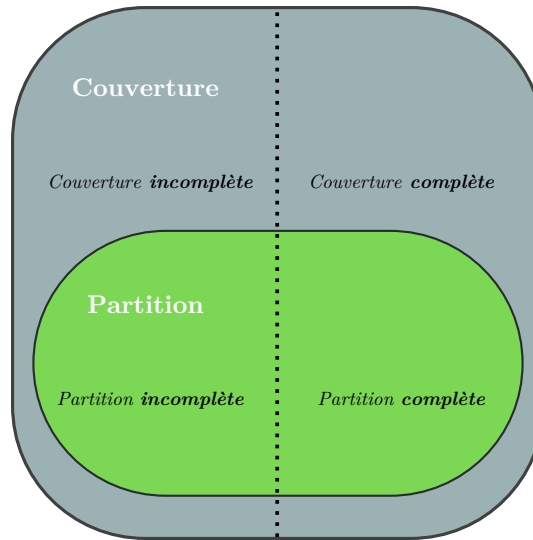


FIGURE 1.7 – **Hierarchie des types d'organisation.** Visualisation des relations formelles (1.2.4). La couverture est le cas le plus général de décomposition en communautés.

Formellement, les différentes organisations communautaires respectent la hiérarchie suivante :

$$\mathcal{P}_c \cap \mathcal{P}_i = \emptyset \quad (1.2.4a)$$

$$\mathcal{C}_c \cap \mathcal{C}_i = \emptyset \quad (1.2.4b)$$

$$\mathcal{P} \subset \mathcal{C} \quad (1.2.4c)$$

$$\mathcal{P}_c \subset \mathcal{C}_c \quad (1.2.4d)$$

$$\mathcal{P}_i \subset \mathcal{C}_i \quad (1.2.4e)$$

où  $\mathcal{P} := \mathcal{P}_c \cup \mathcal{P}_i$  et  $\mathcal{C} := \mathcal{C}_c \cup \mathcal{C}_i$  sont respectivement une partition et une couverture, et où l'indice indique si elles sont complètes ou incomplètes. Les relations (1.2.4a)-(1.2.4b) indiquent que l'incomplétude et la complétude sont des concepts exclusifs. La relation (1.2.4c) formalise le fait qu'une partition est un cas spécial de couverture, sous les modalités spécifiées par les relations (1.2.4a)-(1.2.4b). Cette hiérarchie peut être représentée simplement à l'aide d'un diagramme de Venn, un exercice qui est effectué à la figure 1.7.

Il serait tentant de sauter à la conclusion qu'une description en termes de couverture est la meilleure, de par sa polyvalence. Cependant, deux obstacles majeurs limitent l'utilisation de la couverture. D'une part, des contraintes de modélisation imposent parfois que les éléments du réseau n'appartiennent qu'à un seul groupe. Il n'est alors pas nécessaire d'utiliser une description en termes de couverture, car il faudra de toute façon contraindre les descriptions possibles au sous-ensemble des partitions. D'autre part, pour sélectionner une bonne division en groupes (qu'on traitera plus en détail à la Sec. 1.2.2) il faut explorer – du moins efficacement – l'espace des configurations possibles, qui

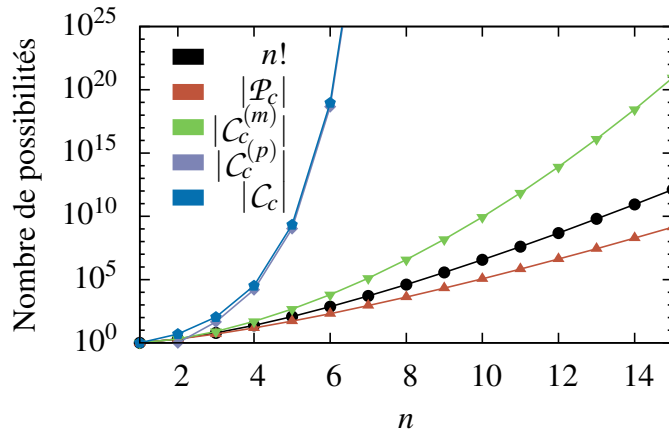


FIGURE 1.8 – **Taille de l'espace des configurations des divers types d'organisations communautaires.** Sur une échelle semi-logarithmique, le nombre de couvertures (pentagones bleus) et de couvertures propres (losanges mauves) explose rapidement (numériquement l'infini pour  $n > 10$ ). Cette explosion est tellement rapide que le facteur exponentiel  $2^{2^n}$  apparaissant dans le calcul du nombre de couvertures propres devient rapidement négligeable. Par contraste, le nombre de partitions possibles (triangles rouges) croît plus lentement que la factorielle du nombre d'éléments à organiser (cercles noirs). Les couvertures minimales (triangles inversés verts) se situent entre les deux extrêmes, et sont donc un bon compromis lorsqu'un chevauchement est désiré.

est énorme dans le cas de la couverture. En effet, le nombre de couvertures complètes

$$|C_c| = \frac{1}{2} \sum_{j=0}^n (-1)^j \binom{n}{j} 2^{2^{n-j}}, \quad (1.2.5a)$$

grandit comme une somme de termes exponentiels [63]. Même si on limite les configurations considérées aux couvertures propres  $C_c^{(p)}$  (couverture ne pouvant inclure l'ensemble complet comme sous-ensemble) et aux couvertures minimales  $C_c^{(m)}$  (couverture où toutes les communautés sont essentielles à la complétude), l'espace des configurations croît comme [63, 64]

$$|C_c^{(p)}| = |C_c| - \frac{1}{4} 2^{2^n} \quad (1.2.5b)$$

$$|C_c^{(m)}| = \sum_{j=1}^n \frac{1}{j!} \sum_{k=j}^{\min(n, 2^{j-1})} k! \begin{Bmatrix} n \\ k \end{Bmatrix} \binom{2^j - j - 1}{k - j}. \quad (1.2.5c)$$

À titre de comparaison, le nombre de partitions (complètes) d'un ensemble de  $n$  éléments distinguables en  $k$  sous-ensembles non vides est donné par un nombre de Stirling de seconde espèce [100, p.81]

$$\begin{Bmatrix} n \\ k \end{Bmatrix} = \frac{1}{k!} \sum_{j=0}^k (-1)^{k-j} \binom{k}{j} j^n, \quad (1.2.6)$$

de sorte que le nombre de partitions possibles d'un réseau constitué de  $n$  éléments (noeuds/liens)

en un nombre arbitraire de communautés (de noeuds/liens) est

$$\sum_{k=1}^n \left\{ \begin{matrix} n \\ k \end{matrix} \right\} =: B_n, \quad (1.2.7)$$

où  $B_n$  est le  $n^{\text{ème}}$  nombre de Bell [100, p.82]. Les nombres de Bell, pourtant énormes, grandissent beaucoup plus lentement que le nombre de couvertures possibles (Fig 1.8). En pratique, puisque les communautés sont intuitivement définies comme des groupes d'éléments contigus (chaque élément doit être rejoignable à partir de tous les éléments de la communauté), une grande partie des couvertures et partitions comptées en (1.2.5)-(1.2.7) doivent être rejetées. Les eqs. (1.2.5)-(1.2.7) sont donc des bornes supérieures pour les quantités réelles, qui sont plus petites et fonction de la structure de chaque réseau. Les approximations (1.2.5)-(1.2.7) donnent toutefois une bonne idée des ordres de grandeur relatifs pour chacun des cas.

### Classes d'organisation et induction

La combinaison des concepts de couverture, partition, communauté de liens et communauté de noeuds permettent de définir quatre grandes classes d'organisations communautaires, soit la partition de noeuds, la partition de liens, la couverture de noeuds et la couverture de liens.

Une relation simple permet de passer d'une organisation de noeuds à une organisation de liens [1, 36] et vice-versa. Celle-ci est basée sur l'observation que les ensembles  $\mathcal{V}$ ,  $\mathcal{E}$  sont couplés par le réseau, de sorte qu'une organisation communautaire pour l'un de ces ensembles implique une organisation pour l'autre ensemble.

**Définition 13.** *Soit les noeuds  $v_i, v_k$  connectés via le lien  $e_{ik}$  et une organisation (partition/couverture) de liens  $\Omega_\ell$ . Cette décomposition en communautés de liens induit une décomposition en communautés de noeuds  $\Omega_n$  défini par*

$$v_i, v_k \in \Psi_\alpha^{(n)} \quad \text{si} \quad e_{ik} \in \Psi_\alpha^{(\ell)} \quad (1.2.8)$$

**Définition 14.** *Soit une organisation de noeuds  $\Omega_n$ . Cette décomposition en communautés de noeuds induit une décomposition en communautés de liens  $\Omega_\ell$  défini par*

$$e_{ik} \in \Psi_\alpha^{(\ell)} \quad \text{si} \quad v_i \in \Psi_\alpha^{(n)}, v_k \in \Psi_\alpha^{(n)} \quad (1.2.9)$$

Dans la premier cas, les noeuds héritent en quelque sorte de l'assignation communautaire des liens qui les relie, alors que dans la deuxième, c'est l'inverse qui se produit.

Le processus d'induction est valide autant dans le cas d'une partition que d'une couverture, mais mène à des organisations différentes, selon le cas. Des relations formelles et complètes sont difficiles à établir car plusieurs cas pathologiques doivent être considérés. On se limite donc au cas réaliste et

	$\mathcal{P}_c^{(n)}$	$\mathcal{P}_i^{(n)}$	$\mathcal{C}_c^{(n)}$	$\mathcal{C}_i^{(n)}$
$\mathcal{P}_c^{(\ell)}$	×	×	✓*	×
$\mathcal{P}_i^{(\ell)}$	✓	✓	✓*	✓*
$\mathcal{C}_c^{(\ell)}$	×	×	✓	×
$\mathcal{C}_i^{(\ell)}$	✓*	✓*	✓	✓

TABLE 1.1 – **Organisation de noeuds et de liens pouvant être reliées par induction.** Ce tableau recense les types d’organisations de noeuds ( $n$ ) pouvant être induites par une organisation de liens ( $\ell$ ) donnée, et vice-versa. L’exercice est effectué pour toutes les combinaisons de partitions  $\mathcal{P}$  ou couvertures  $\mathcal{C}$ , complètes ( $c$ ) ou incomplètes ( $i$ ). Une étoile (\*) indique que seul un sous-ensemble restreint peut être induit (détails en Annexes Sec. B.1).

non trivial d’un réseau constitué d’une composante connexe (tout noeud est rejoignable) divisée en plus d’une communauté contiguë (tout élément est rejoignable à l’intérieur d’une communauté) d’au moins deux noeuds ou d’au moins un lien. Dans ce cas particulier, l’induction est bijective, de sorte qu’il devient possible d’établir la liste des organisations  $\Omega_\ell$  pouvant être induite par une organisation  $\Omega_n$  donnée, et vice-versa. Ces relations sont présentées dans le tableau 1.1 (la démonstration est reléguée à l’annexe Sec. B.1).

## 1.2.2 Algorithmes de détection

Une structure communautaire pertinente, qu’elle soit une partition ou une couverture, de noeuds et ou de liens, n’est pas aussi triviale à identifier que le degré moyen ou le ratio de transitivité d’un réseau. Alors que ces informations découlent simplement et directement de la matrice d’adjacence (Sec. 1.1.3), l’organisation communautaire est donnée par la solution d’un problème de de partition de la classe de complexité NP-difficile [38, p.83, §3.1].

L’identification de l’organisation communautaire est effectuée à l’aide d’*algorithmes de détection communautaire*, dont l’objectif est de visiter un sous-ensemble restreint de l’espace des configurations possibles (cf. Fig. 1.8) afin d’identifier la décomposition en communautés qui optimise un critère de sélection (ou fonction de coût)  $C$  donné. Cette définition large comporte donc deux parties, soit un mécanisme d’exploration et un critère de sélection.

### Critère de sélection

L’idée sous-jacente à tout algorithme de détection communautaire est d’identifier des sous-ensembles de noeuds ou de liens similaires structurellement (densément connectés, plus rapprochés que la moyenne, etc.). Par conséquent, bien que le critère de sélection exacte  $C$  varie énormément d’un algorithme à l’autre, les optima de  $C$  sont en général associés à des divisions en groupes où les noeuds/liens similaires (selon une mesure donnée) sont assignés aux mêmes communautés. La latitude sur le choix de critère permet d’adapter la structure détectée au problème de modélisation qui motive l’utilisation d’un algorithme de détection.

Par exemple, regrouper des liens ayant une haute similarité structurelle (eq. 1.1.10) mène à une structure communautaire qui permet de bien modéliser la propagation sur réseau [49] ou de quantifier l'importance des éléments du réseau en termes de leur potentiel épidémiologique [46]. Dans quelques cas rares [51], d'autres critères (e.g. densité) mènent à l'identification de groupes fonctionnels [95], i.e. des sous-ensembles de noeuds ou de liens dont la fonction (e.g. biologique, technologique) est révélée par la structure mésoscopique du réseau. D'autres choix de critères (e.g. recherche de sous-graphes ER) permettent d'extraire une information communautaire qui peut ensuite être utilisée pour reproduire les propriétés observables des réseaux complexes réels (e.g. distributions de degrés, d'agré-gations) à l'aide de modèles génératifs de faible dimension (peu de paramètres) [44, 45, 96].

Ici, le message important est qu'il n'existe pas de *critère universel*, ce qui n'est pas surprenant lorsqu'on réalise qu'il existe  $2^N$  réseaux non dirigés et non dirigés différents de taille  $N$  : il serait pour le moins inattendu qu'un seul critère permette d'extraire une structure pertinente, pour tout problème et pour chacun de ces réseaux (pour tout  $N$ ).

### Mécanisme exploratoire

Un algorithme de détection n'est pas complet sans mécanisme de sélection efficace. En effet, l'énumération et l'évaluation de  $C$  pour toutes les solutions possibles est à proscrire, même dans des graphes de petite taille (Fig. 1.8). On doit donc recourir à des algorithmes approximatifs permettant d'obtenir des solutions approchées dans un temps polynomial déterministe [38, p.83]. Le développement d'algorithmes approximatifs est un domaine de recherche très actif, de sorte qu'il existe aujourd'hui une foule d'algorithmes s'attaquant au même problème par des angles différents (cf. références de [38]).

Historiquement, les premiers algorithmes approximatifs ont été développés pour solutionner le problème de la partition de noeuds. Il convient de souligner une distinction subtile entre algorithme de partitionnement et algorithme de détection communautaire menant à une partition de noeuds : un algorithme de détection doit procéder sans *a priori* [76], i.e. le nombre de communautés et la taille de ces dernières ne sont pas des informations qui sont connues à l'avance, alors que cette information est *imposée* dans le cas des algorithmes de partition.

Les algorithmes de *détection* de partition de noeuds, qui doivent implicitement ou explicitement solutionner un problème de sélection de modèle [115, § 5], sont donc apparus plus tard [38, §4-5]. Ils ont rapidement été suivis d'algorithmes permettant d'identifier des organisations communautaires avec chevauchement (couverture, traditionnellement d'intersections relativement petites) [114].

Sans prétendre effectuer une revue complète des types de mécanismes d'exploration (voir [38, 75, 114] pour plus de détails), on peut distinguer les grandes classes d'algorithmes suivantes :

1. Algorithmes méta-heuristique purs [104] : l'espace des configurations est exploré à l'aide de mouvements explicites (e.g. le noeud  $i$  passe de la communauté  $\Psi_\alpha^{(n)}$  à  $\Psi_\beta^{(n)}$ ), selon une heuris-

tique donnée (e.g. recuit simulé [103], colonie de fourmis [43]).

2. Algorithmes d'agrégation : partant d'un ensemble de communautés initiales choisies systématiquement (e.g. [1, 81]) ou de façon aléatoire (e.g. [59]), on ajoute des éléments aux communautés lorsqu'ils répondent à un critère donné (e.g. vraisemblance [59, 60], propagation de  $k$ -clique [81], similarité de liens [1], distance dans un espace de noeuds [62]).
3. Algorithmes de division : le réseau est séparé en plusieurs parties en retirant des éléments (noeuds / liens) selon une mesure de leur non-importance (e.g. plus petite centralité [40], ou minimisation du nombre de liens inter-modules [54]).
4. Algorithmes dynamiques : une dynamique est imposée sur le réseau (e.g. marcheur aléatoire [94], modèle de spins [92]) et les états stables ou asymptotiques sont utilisés pour identifier les communautés.
5. Algorithmes d'inférence statistique : on suppose que le réseau est généré par un membre d'une famille de modèles génératifs. Le modèle est sélectionné en comparant la structure du réseau aux membres de cette famille (e.g. par estimation du maximum de vraisemblance [10, 115]).

Toutes ces classes d'algorithmes peuvent être adaptées à n'importe quel type d'organisation communautaire (partition, couverture, noeuds, liens), à l'exception des algorithmes de divisions, qui mènent uniquement à des partitions (parfois incomplètes).

### 1.2.3 Propriétés de l'organisation communautaire

Tout comme la structure d'un réseau peut être quantifiée à partir de la matrice d'adjacence  $\mathbf{A}$  (Sec. 1.1.3), la structure *communautaire* d'un réseau peut être quantifiée à l'aide d'une *matrice d'incidence communautaire*  $\mathbf{P}$ . La matrice  $\mathbf{P}$  est propre à chaque réalisation d'un algorithme (constitué d'un couple mécanisme-critère) et à chaque réseau.

**Définition 15.** Soit une organisation communautaire  $\Omega_x = \{\Psi_1^{(x)}, \dots, \Psi_g^{(x)}\}$  (où  $x = n$  pour des communautés de noeuds et  $x = \ell$  pour des communautés de liens) comportant  $|\Omega_x| = g$  groupes. La matrice d'incidence communautaire  $\mathbf{P}$  est une matrice binaire  $N \times g$  (communautés de noeuds) ou  $M \times g$  (communautés de liens). L'élément  $p_{i\alpha}$  de  $\mathbf{P}$  est égale à 1 si l'élément  $i$  appartient à la communauté  $\Psi_\alpha^{(x)}$ , et 0 autrement.

Dans cette section, on effectue un tour d'horizon limité des propriétés de l'organisation communautaire qui seront pertinentes pour le reste du présent mémoire. Ces propriétés sont divisées en deux grandes catégories, à savoir les propriétés qui dépendent uniquement de l'organisation mésoscopique  $\mathbf{P}$ , et les propriétés *mixtes* qui dépendent autant de la décomposition  $\mathbf{P}$  que de la structure du réseau.

#### Propriétés strictement communautaires

Prise seule, la matrice  $\mathbf{P}$  contient toute l'information qui est relative uniquement aux communautés, comme leur taille.



**Définition 16.** La taille d'une communauté  $\Psi_\alpha^{(x)}$  est la cardinalité de l'ensemble  $\Psi_\alpha^{(x)}$  des éléments de la communauté, i.e.  $|\Psi_\alpha^{(x)}| = \sum_i p_{i\alpha}$ . Le vecteur des tailles  $\mathbf{s}$  est obtenu de la matrice d'incidence communautaire via le produit

$$\mathbf{s}^T = \mathbf{1}^T \mathbf{P} \implies \mathbf{s} = \mathbf{P}^T \mathbf{1} \quad (1.2.10)$$

Afin d'alléger la notation, on utilisera de façon interchangeable l'expression matricielle  $[\mathbf{s}]_\alpha$  et l'expression scalaire  $s_\alpha$  pour la taille de la communauté  $\Psi_\alpha$ .

**Définition 17.** Le nombre d'appartenance d'un élément est le nombre de communautés auquel il participe. Le vecteur d'appartenance  $\mathbf{m}$  est obtenu de la matrice d'incidence communautaire via le produit

$$\mathbf{m} = \mathbf{P} \mathbf{1} \quad (1.2.11)$$

Afin d'alléger la notation, on utilisera de façon interchangeable l'expression matricielle  $[\mathbf{m}]_i$  et l'expression scalaire  $m_i$  pour le nombre d'appartenance du noeud  $i$ .

On notera que si  $\Omega_x$  est une partition complète, la somme des tailles doit être égale à la taille du réseau

$$\mathbf{1}^T \mathbf{s} = \mathbf{1}^T \mathbf{P}^T \mathbf{1} = \begin{cases} N & \text{Communautés de noeuds} \\ M & \text{Communautés de liens} \end{cases} \quad (1.2.12)$$

et tous les éléments du vecteur  $\mathbf{m}$  doivent être exactement égaux à 1 :

$$\mathbf{m} = \mathbf{P} \mathbf{1} = \mathbf{1}. \quad (1.2.13)$$

Ces propriétés sont perdues dans le cas d'une couverture, car les rangées de  $\mathbf{P}$  contiennent alors plus d'un élément.

Comme dans le cas du vecteur des degrés  $\mathbf{k}$  (e.g. Fig. 1.6), les séquences de tailles  $\mathbf{s}$  et d'appartenances  $\mathbf{m}$  peuvent être utilisées pour construire des *distributions normalisées* de taille et d'appartenance. Pour les réseaux réels, ces distributions sont généralement en loi de puissance  $p_x \propto x^{-\gamma}$  [44, 58], ce qui a des implications importantes pour la détection communautaire : de très petites communautés côtoient des communautés de grande taille, de sorte que les algorithmes de détection doivent pouvoir opérer à plusieurs ordres de grandeur, afin d'extraire la structure communautaire complète.

Une autre information importante contenue strictement dans la matrice  $\mathbf{P}$  a trait à la force du chevauchement des communautés.

**Définition 18.** Le chevauchement absolu de deux communautés est défini comme la taille de l'intersection des ensembles des éléments de ces communautés

$$c(\Psi_\alpha^{(x)}, \Psi_\beta^{(x)}) := |\Psi_\alpha^{(x)} \cap \Psi_\beta^{(x)}| = \sum_i p_{i\alpha} p_{i\beta} = [\mathbf{P}^T \mathbf{P}]_{\alpha\beta}, \quad (1.2.14)$$

Le chevauchement moyen est donné par l'expression matricielle

$$\begin{aligned} \langle c(\Psi_\alpha^{(x)}, \Psi_\beta^{(x)}) \rangle &= \frac{1}{g(g-1)} \sum_{\alpha \neq \beta} c(\Psi_\alpha^{(x)}, \Psi_\beta^{(x)}) \\ &= \frac{1}{g(g-1)} \left( \sum_{\alpha, \beta} - \sum_{\alpha = \beta} \right) c(\Psi_\alpha^{(x)}, \Psi_\beta^{(x)}) = \frac{1}{g(g-1)} \mathbf{1}^\top \mathbf{P}^\top \mathbf{P} \mathbf{1} - \frac{1}{g(g-1)} \text{Tr}(\mathbf{P}^\top \mathbf{P}). \end{aligned} \quad (1.2.15)$$

**Définition 19.** Le chevauchement normalisé de deux communautés est quantifié via l'indice de Jaccard (définition 9) des ensembles des éléments de ces communautés

$$c_n(\Psi_\alpha^{(x)}, \Psi_\beta^{(x)}) := \frac{|\Psi_\alpha^{(x)} \cap \Psi_\beta^{(x)}|}{|\Psi_\alpha^{(x)} \cup \Psi_\beta^{(x)}|} = \frac{\sum_i p_{i\alpha} p_{i\beta}}{\sum_i (p_{i\alpha} + p_{i\beta} - p_{i\alpha} p_{i\beta})} = \frac{[\mathbf{P}^\top \mathbf{P}]_{\alpha\beta}}{[\mathbf{s}]_\alpha + [\mathbf{s}]_\beta - [\mathbf{P}^\top \mathbf{P}]_{\alpha\beta}}. \quad (1.2.16)$$

Par construction, une partition possède un chevauchement moyen nul, puisque l'intersection des communautés est définie comme nulle. Dans le cas des couvertures, une borne supérieure subtile apparaît, à cause de la nature même du problème de l'identification communautaire. En effet, l'information communautaire est extraite de la structure du réseau, tel qu'un recouvrement parfait n'a pas de raisons d'être, puisqu'aucune information structurelle ne peut le justifier [1, SOM, p.10].

### Propriétés mixtes

Certaines propriétés peuvent être qualifiées de *mixtes*, car elles dépendent de la structure  $\mathbf{A}$  et de la décomposition communautaire  $\mathbf{P}$ . Ici, elles seront définies pour des communautés de noeuds, mais des propriétés analogues peuvent être calculées pour les communautés de liens, en obtenant d'abord l'organisation induite  $\Omega_\ell \rightarrow \Omega_n$ , telle que définie à la Sec. 1.2.1.

**Définition 20.** Le nombre de liens internes de la communauté  $\Psi_\alpha^{(n)}$  est le nombre de liens strictement contenu dans le sous-graphe correspondant à  $\Psi_\alpha^{(n)}$ .

$$L_\alpha = \frac{1}{2} \sum_{i \neq j} a_{ij} p_{i\alpha} p_{j\alpha} = \frac{1}{2} [\mathbf{P}^\top \mathbf{A} \mathbf{P}]_{\alpha\alpha}, \quad (1.2.17)$$

tel que le nombre total de liens internes est

$$L = \sum_{\alpha=1}^g L_\alpha = \frac{1}{2} \text{Tr}(\mathbf{P}^\top \mathbf{A} \mathbf{P}). \quad (1.2.18)$$

Il faut noter que le même lien peut potentiellement être compté plus qu'une fois dans le calcul de  $L$  si  $\Omega_n$  est une couverture.  $L$  perd donc son sens absolu dans ce cas. Toutefois, la mesure locale  $L_\alpha$  reste pertinente dans tous les cas (partition, couverture), car on peut l'utiliser pour calculer la *densité d'une communauté*.

**Définition 21.** La densité  $\rho_\alpha$  d'une communauté de  $[\mathbf{s}]_\alpha$  noeuds  $\Psi_\alpha^{(n)}$  est le rapport du nombre de liens existant dans cette communauté, sur le nombre de liens possibles

$$\rho_\alpha = \frac{L_\alpha}{\binom{[\mathbf{s}]_\alpha}{2}} = \frac{[\mathbf{P}^\top \mathbf{A} \mathbf{P}]_{\alpha\alpha}}{[\mathbf{s}]_\alpha([\mathbf{s}]_\alpha - 1)}. \quad (1.2.19)$$

**Définition 22.** La densité excédentaire  $\tilde{\rho}_\alpha$  d'une communauté de noeuds  $\Psi_\alpha^{(n)}$  est le rapport du nombre de liens existant dans cette communauté, sur le nombre de liens possibles, une fois les liens nécessaires à la connectivité retirés. La structure en arbre ( $[\mathbf{s}]_\alpha - 1$  liens) étant la structure minimale garantissant la connectivité d'un groupe de noeuds,  $\tilde{\rho}_\alpha$  est donné par

$$\tilde{\rho}_\alpha = \frac{L_\alpha - ([\mathbf{s}]_\alpha - 1)}{\binom{[\mathbf{s}]_\alpha}{2} - ([\mathbf{s}]_\alpha - 1)} = \frac{[\mathbf{P}^\top \mathbf{A} \mathbf{P}]_{\alpha\alpha} - 2([\mathbf{s}]_\alpha - 1)}{([\mathbf{s}]_\alpha - 1)([\mathbf{s}]_\alpha - 2)}. \quad (1.2.20)$$

Pour des fins de moyennage et de normalisation, la densité des communautés de taille  $[\mathbf{s}]_\alpha = 1$  ainsi que la densité excédentaire des communautés de taille  $[\mathbf{s}]_\alpha \leq 2$  sont définies comme nulles [1].

### Propriétés des partitions

Lorsque la division en communautés est une partition, deux propriétés additionnelles sont traditionnellement définies [75, § 11], soit la taille de la coupure et la modularité.

**Définition 23.** La taille de la coupure d'une partition est définie comme le nombre de liens qui ne sont pas des liens internes

$$R = M - L \quad (1.2.21)$$

**Définition 24.** La modularité  $Q$  d'une partition mesure la différence entre le nombre de liens internes réalisés et le nombres de liens internes attendus, selon un modèle nul donné.

$$Q_{MN} = \sum_\alpha \sum_{ij} (a_{ij} - \langle a_{ij} \rangle_{MN}) p_{i\alpha} p_{j\alpha} \quad (1.2.22)$$

Le modèle de configurations (Sec. 1.1.4) est l'hypothèse nulle la plus répandue<sup>16</sup>. Dans un réseau de taille fini construit à l'aide du modèle de configurations, le nombre moyen de liens entre les deux noeuds  $i$  et  $j$  est  $k_i k_j / 2M$ , soit le produit de leurs degrés, normé par le nombre total de degrés, tel que

$$Q_{CM} = \sum_\alpha \sum_{ij} \left( a_{ij} - \frac{k_i k_j}{2M} \right) p_{i\alpha} p_{j\alpha} \quad (1.2.23)$$

est la définition la plus courante (à un facteur multiplicatif près [74]).

<sup>16</sup>. Plusieurs autres choix ont été explorés et donnent parfois de meilleurs résultats [106, § II.A]. Ici, on utilise le CM pour des raisons historiques.

Acronyme	Algorithme	Référence
CCPA	Percolation de cliques en cascade	[Chap.2]
GCE	Algorithme glouton d'expansion de cliques	[60]
GreedyMod	Algorithme glouton de modularité	[25]
InfoMap	Flot d'information	[94]
LG3	Modularité sur graphe adjoint (B.2)	[36]
Louvain $\emptyset$	Algorithme de Louvain sans modèle nul	[16]
Louvain CM	Algorithme de Louvain avec modèle nul standard	[16]
Louvain ER	Algorithme de Louvain avec modèle nul Erdős-Rényi	[16]
OSLOM	Inférence statistique locale	[59]

TABLE 1.2 – **Algorithmes de détection utilisés pour l'étude de l'arXiv circa 2005.** Liste des algorithmes utilisés pour produire la Fig. 1.9. CCPA, OSLOM, et GCE produisent des couvertures de noeuds, LG3 produit une partition de liens, tandis que Louvain  $\emptyset$ /CM/ER, InfoMap et GreedyMod mènent à des partitions de noeuds. On notera que toutes les variantes de l'algorithme de Louvain utilisent le même mécanisme d'exploration de l'espace des configurations, et que GreedyMod et Louvain CM ont le même critère de sélection (modularité).

Ces deux mesures sont *a priori* complètement différentes. Une similarité frappante peut toutefois être mise en lumière par une réécriture de  $R$  en termes de la matrice Laplacienne  $\mathbf{L}$  (matrice diagonale des degrés moins la matrice d'adjacence) et une réécriture de  $Q$  en termes de la matrice de modularité  $\mathbf{B}$  (matrice d'adjacence moins la matrice du produit des degrés  $\mathbf{k}\mathbf{k}^T/M$ ). En effet, utilisant le fait que  $\sum_{\alpha} \sum_j \delta_{ij} p_{i\alpha} p_{j\alpha} = \sum_{\alpha} p_{i\alpha}^2 = 1$  dans le cas d'une partition, la coupure  $R$  peut être réécrite comme

$$R = \frac{1}{2} \sum_{i,j} a_{ij} - \frac{1}{2} \sum_{\alpha} \sum_{i,j} a_{ij} p_{i\alpha} p_{j\alpha} = \frac{1}{2} \sum_{\alpha} \sum_{i,j} (k_i \delta_{ij} - a_{ij}) p_{i\alpha} p_{j\alpha} = \frac{1}{2} \sum_{\alpha} \sum_{i,j} l_{ij} p_{i\alpha} p_{j\alpha} = \frac{1}{2} \text{Tr}(\mathbf{P}^T \mathbf{L} \mathbf{P}), \quad (1.2.24)$$

tandis que la modularité prend la forme

$$Q_{CM} = \sum_{\alpha} \sum_{i,j} \left( a_{ij} - \frac{k_i k_j}{2M} \right) p_{i\alpha} p_{j\alpha} = \sum_{\alpha} \sum_{i,j} b_{ij} p_{i\alpha} p_{j\alpha} = \text{Tr}(\mathbf{P}^T \mathbf{B} \mathbf{P}). \quad (1.2.25)$$

Ces formes sont fonctionnellement identiques.

### Propriétés mésoscopiques des réseaux réels

Expérimentalement, on observe que les réseaux réels d'une même catégorie conceptuelle (e.g. réseaux sociaux) se comportent de façon similaire au niveau mésoscopique, et ce peu importe l'algorithme de détection utilisé pour extraire  $\mathbf{P}$  [75]. Afin d'illustrer le type de comportement observé, on calcule ici la distribution de ces quantités pour un réseau déjà visité à la Sec. 1.1.4, à savoir l'*arXiv* de la matière condensée circa 2005 [81]. Les distributions résultantes sont rassemblées dans la figure 1.9, pour 9 algorithmes différents.

En premier lieu, on note que tous les algorithmes semblent produire une organisation communautaire où les tailles  $\{s_k\}$  et appartenances  $\{m_i\}$  sont distribuées en loi de puissance  $p_x \propto x^{-\gamma}$  (tableau 1.3). On extrait les exposants  $\gamma$  de la distribution expérimentale en estimant le maximum

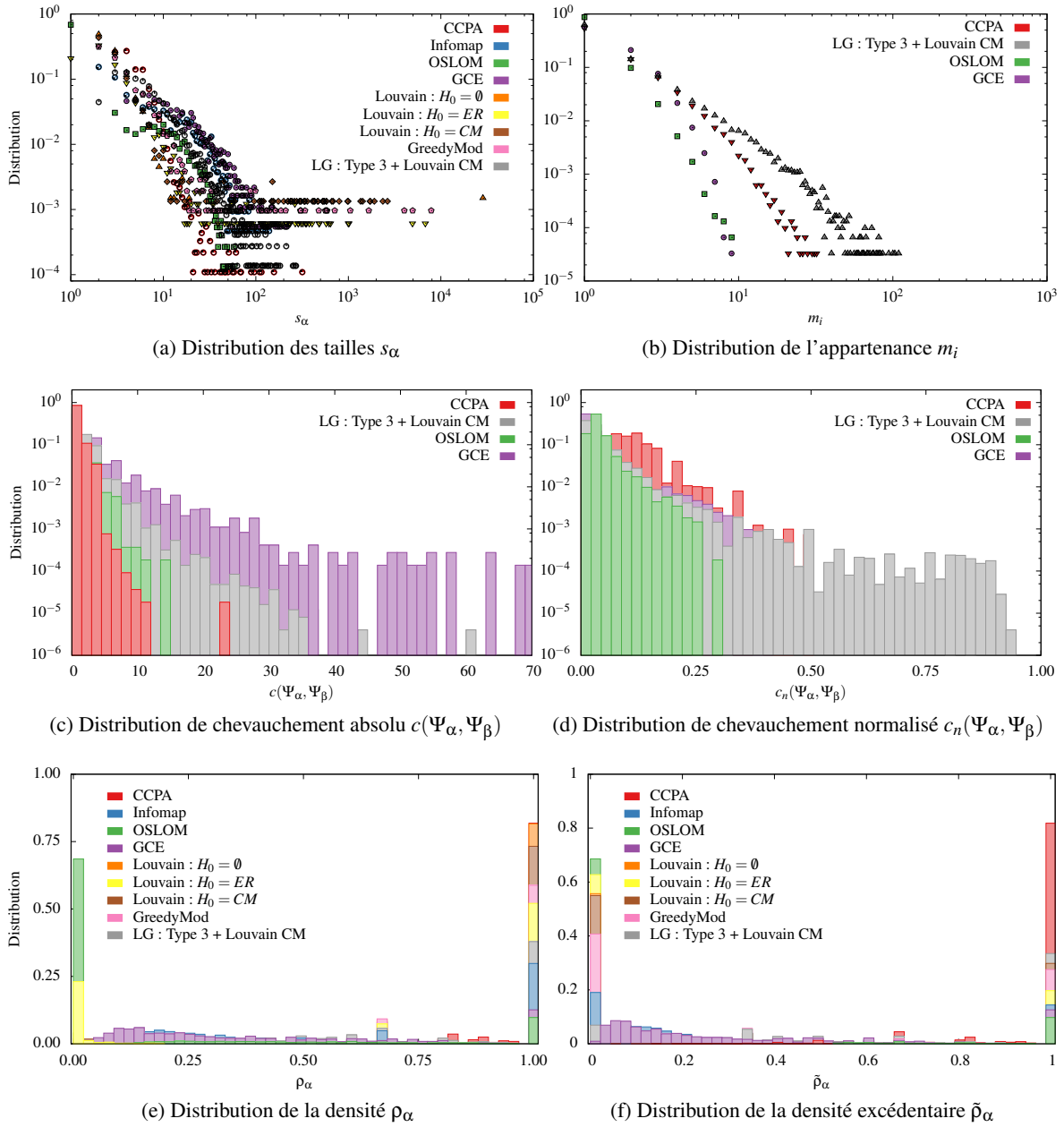


FIGURE 1.9 – **Distribution des diverses propriétés mésoscopiques de l'arXiv circa 2005.** Afin d'unifier la comparaison, on a utilisé l'organisation  $C_n$  induite par la partition de liens qui est naturellement produite par LG3, avant de calculer toutes les quantités pour des organisation de noeuds. Les algorithmes produisant des partitions de noeuds ont été omis en (b)-(c)-(d) car l'appartenance est trivialement  $m_i = 1 \forall i$ , et le chevauchement  $c(\Psi_\alpha, \Psi_\beta), c_n(\Psi_\alpha, \Psi_\beta) = 0 \forall \alpha, \beta$ .

Algorithme	$\gamma$	$s_{\min}$	$\mathcal{D}$
CCPA	4.31	6	0.018
GCE	3.63	57	0.047
GreedyMod	2.13	3	0.047
InfoMap	4.34	41	0.035
LG3	3.23	43	0.027
Louvain $\emptyset$	3.38	3	0.094
Louvain CM	2.23	297	0.048
Louvain ER	2.10	2	0.027
OSLOM	5.92	27	0.037

(a) Distribution de taille  $\{s_k\}$ 

Algorithme	$\gamma$	$m_{\min}$	$\mathcal{D}$
CCPA	4.68	11	0.022
GCE	15.53	7	0.018
LG3	2.10	1	0.014
OSLOM	3.49	1	0.012

(b) Distribution d'appartenance  $\{m_i\}$ 

TABLE 1.3 – **Exposants MLE des distributions en loi de puissance de la taille et de l'appartenance pour arXiv cond-mat 2005.** Couple  $(\gamma, x_{\min})$  le plus vraisemblable, en supposant que la taille des communautés (a) et l'appartenance des noeuds (b) sont distribuées selon la loi de puissance  $p(x) \propto x^{-\gamma}$  pour  $x > x_{\min}$ . La statistique de Kolmogorov-Smirnov  $\mathcal{D}$  apparaît dans la colonne de droite. Elle quantifie la distance entre l'hypothèse et les données empiriques [86, p.736] : plus  $\mathcal{D}$  est proche de 0, plus l'hypothèse correspond aux observations.

de vraisemblance de l'hypothèse en loi de puissance [26]. Cette hypothèse est invraisemblable dans quelques cas rares, e.g. GCE, où un exposant  $\gamma \approx 15$  est inféré à partir de 3 points<sup>17</sup>. L'hypothèse de la loi de puissance, ici confirmée pour 8 algorithmes sur 9, est en fait généralisable à *la grande majorité* des réseaux et des algorithmes (e.g. référence [44]), de sorte que la plupart des modèles tentant de reproduire la structure des réseaux réels incluent cette propriété implicitement ou explicitement [58, 96].

En deuxième lieu, on observe que le chevauchement est généralement limité à une petite partie des communautés, ou de façon équivalente que  $c_n(\Psi_\alpha, \Psi_\beta)$  est *beaucoup* plus souvent près de 0 que de 1 (Fig. 1.9.d, échelle semi-logarithmique). À la définition 18, on avait déjà postulé l'existence d'une telle borne supérieure pour  $c_n(\Psi_\alpha, \Psi_\beta)$ , relevant de la nature même du problème de détection : il est impossible de discerner deux communautés identiques structurellement (même noeuds/liens), tel qu'un chevauchement parfait ne peut être obtenu sur la base du réseau seul (de la méta-information est nécessaire, e.g. groupes de recherche pour l'arXiv). Naturellement, le chevauchement est trivialement nul dans le cas des algorithmes détectant des partitions de noeuds.

En troisième lieu, on observe que la majorité des communautés sont soit extrêmement denses, soit extrêmement creuses (figures 1.9.e-f.). Ce comportement trahit un défaut important dû aux mécanismes exploratoires employés : les algorithmes trop gloutons tendent à identifier des communautés triviales de noeuds très densément connectées, avant de placer les noeuds restant dans des commu-

17. On peut conclure que l'hypothèse en loi de puissance est invraisemblable car un exposant plausible se situe généralement dans l'intervalle  $\gamma \in [2, 5]$ . La borne inférieure est due au fait qu'une loi de puissance ne possède pas de moyenne finie pour  $\gamma < 2$ , tandis que la borne supérieure est due au fait qu'on ne peut pas déterminer de façon statistiquement significative qu'un exposant est grand (e.g.  $\gamma \gg 1$ ) pour des systèmes de petites tailles [26].

Algorithme	$Q_{CM}$	$R$
GreedyMod	0.65	21276
InfoMap	0.65	44284
Louvain $\emptyset$	0.04	0
Louvain CM	0.72	29015
Louvain ER	0.63	20370

TABLE 1.4 – **Modularité et coupe des partitions de l’arXiv**. Modularité (1.2.25) et coupe (1.2.24) des partitions de noeuds de arXiv cond-mat circa 2005. Ce réseau est constitué de  $M = 125959$  liens, ce qui correspondrait à la coupe maximale. La coupe  $R = 0$  nulle obtenue par l’algorithme Louvain  $\emptyset$  indique que cet algorithme n’a identifié que les composantes déconnectées du réseau.

nautés “fourre-tout”, qui sont forcément de faible densité.

En dernier lieu, soulignons que les distributions présentées à la figure 1.9 nous permettent d’illustrer les conséquences du choix de mécanisme exploratoire et de critère de sélection  $C$ .

Les effets du critère de sélection peuvent être observés en comparant les distributions associées à l’algorithme de Louvain (figures 1.9.a,e-f), qui est essentiellement un mécanisme d’exploration rapide, applicable à plusieurs critères  $C$ <sup>18</sup>. Sans surprise, même si le mécanisme d’exploration est identique dans tous les cas, la *nature* des communautés identifiées diffère : la variante sans modèle nul extrait une énorme communauté de faible densité et quelques communautés très denses, la variante CM mène à plusieurs communautés de haute densité *excédentaire*, et la variante ER se situe à mi-chemin entre ces deux extrêmes.

Les conséquences du choix de mécanisme d’exploration peuvent aussi être observées, en comparant les distributions associées à l’algorithme glouton d’optimisation de la modularité (GreedyMod) et l’algorithme Louvain CM, qui utilisent un même critère de sélection  $Q_{CM}$  et des mécanismes d’exploration différents (figures 1.9.a,e-f). Les partitions identifiées comme optimales ne sont manifestement pas identiques dans les deux cas, car autant les distributions de tailles que de densité diffèrent. On peut d’ailleurs confirmer la disparité en se référant à la modularité  $Q_{CM}$  de ces partitions (tableau 1.4), qui nous confirme que GreedyMod identifie une partition de qualité objectivement inférieure.

#### 1.2.4 Problèmes ouverts

Le bref tour d’horizon effectué dans ce chapitre introductif permet à peine d’entrevoir la complexité et la richesse du problème de la détection communautaire. En effet, malgré la taille toujours grandissante du corpus scientifique traitant de la structure mésoscopique des réseaux complexes, plusieurs problèmes ouverts subsistent. Ce mémoire identifie et offre des pistes de solutions à trois de ces problèmes, à l’aide de trois approches complémentaires : heuristique, analytique et numérique.

Dans un premier temps, on s’attaque au problème de la résolution. Tel que souligné à la Sec. 1.2.3,

18. Ce mécanisme d’exploration permet d’optimiser des variantes arbitraires de la modularité (1.2.22), i.e. d’optimiser  $Q_{MN}$  pour tout choix de modèle nul  $\langle a_{ij} \rangle_{MN}$ .

la structure communautaire des réseaux réels est libre d'échelle, ce qui cause problème aux algorithmes incapables d'identifier une structure communautaire distribuée sur plusieurs ordres de grandeur [39, 53, 106]. Une solution heuristique est proposée au chapitre 2.

On traite ensuite du problème de l'unification. En effet, l'abondance de mécanismes d'exploration et de critères de sélection (Sec. 1.2.2) rend la comparaison d'algorithmes de détection difficile. Ce problème est exploré analytiquement aux chapitres 3-4, où les bases d'une théorie unifiée de la détection communautaire sont établies.

Finalement, au chapitre 5, on pose les bases d'un modèle stochastique permettant de reproduire la structure effective *et* la structure communautaire des réseaux réels. Ce modèle possède peu de paramètres et permet de construire des réseaux dans des régimes structurels qualitativement très différents. Quoique le travail de comparaison ne soit pas effectué dans le cadre de cet ouvrage, ce modèle se veut un *banc d'essai numérique* pour les algorithmes de détection. Il s'agit donc d'un premier pas vers une nouvelle solution numérique au problème de l'évaluation quantitative et qualitative des algorithmes de détection.



## **Chapitre 2**

# **Dévoilement des communautés cachées via la détection en cascade sur réseaux**

**Unveiling Hidden Communities Through Cascading Detection on Network Structures**

Jean-Gabriel Young, Antoine Allard

Laurent Hébert-Dufresne et Louis J. Dubé

Département de Physique, de Génie Physique, et d'Optique,  
Université Laval, Québec, Québec, Canada G1V A06.

En préparation

## 2.1 Avant-propos

Dans ce premier chapitre de recherche, on identifie un type de problème de résolution jusqu'alors inconnu. On établit ensuite que plusieurs algorithmes de détection communautaire existants sont victimes de ce problème. En réponse, un méta-algorithme<sup>1</sup> de détection est proposé et testé sur deux algorithmes de détection populaires. On montre que leur performance est drastiquement améliorée par le méta-algorithme, et ce pour plusieurs (8) réseaux réels.

Il s'agit d'une version préliminaire d'un article initialement préparé pour la conférence COMPLEX 2012 (Santa-Fe, New-Mexico).

## 2.2 Résumé

Les algorithmes de détection communautaire visent à assigner les noeuds et les liens d'un réseau complexe à des communautés significatives (e.g. groupes, modules fonctionnels). Leur développement a mené à une meilleure compréhension de l'organisation des réseaux complexes. Ici, nous avançons que la plupart des algorithmes de détection identifient correctement les communautés prédominantes des réseaux de grandes tailles, mais que ce succès n'est pas reproduit à toutes les échelles. Par conséquent, ces algorithmes laissent une portion importante des noeuds et liens mal ou non assignés. Nous montrons que les communautés denses sont à l'origine de ce problème, car elles agissent comme un écran qui cachent les petites communautés. Nous proposons un méta-algorithme de détection en cascade polyvalent qui élimine ce problème. Nous illustrons notre approche en l'appliquant à deux algorithmes populaires, que nous utilisons ensuite pour identifier la structure communautaire de plusieurs réseaux réels. Les algorithmes modifiés nous permettent d'identifier des communautés jusqu'alors ignorées par les algorithmes originaux, ce qui mène à une réduction importante de la fraction de liens non-assignés.

## 2.3 Abstract

Community detection is the process of assigning nodes and links in significant communities (e.g. clusters, functional modules) and its development has led to a better understanding of complex networks. When applied to sizable networks, we argue that most detection algorithms correctly identify prominent communities, but fail to do so across multiple scales. As a result, a significant fraction of the network is left uncharted. We show that this problem stems from larger or denser communities overshadowing smaller or sparser ones, and that this effect accounts for most of the undetected communities and unassigned links. We propose a generic cascading approach to community detection that circumvents the problem. Using real network datasets with two widely used community detection algorithms, we show how a cascading procedure allows for the detection of the missing communities and results in a significant increase in the fraction of assigned links.

---

1. Un algorithme mettant en relation des algorithmes déjà existants.

## 2.4 Introduction

Over the course of the last decade, network science has attracted an ever growing interest since it provides important insights on a large class of interacting complex systems. One of the features that has drawn much attention is the structure of interactions highlighted by the network representation. Indeed, it has become increasingly clear that global structural patterns emerge in most real networks [40]. One such pattern, where links and nodes are aggregated into larger groups, is called the community structure of a network.

While the exact definition of communities is still not agreed upon [38], the general consensus is that these groups should be denser than the rest of the network. The notion that communities form some sort of independent units (families, friend circles, coworkers, protein complexes, etc.) within the network is thus embedded in that broader definition. It follows that communities represent functional modules, and that understanding their layout as well as their organization on a global level is crucial to a fuller understanding of the system under scrutiny [49, 44].

By developing techniques to extract this organization, one assumes that communities are encoded in the way nodes are interconnected, and that their structure may be recovered from limited, incomplete topological information. Various algorithms and models have been proposed to tackle the problem, each featuring a different definition of the community structure while sharing the same general objective. Although these tools have been used with success in several different contexts [38, 1, 81], a number of shortcomings are still to be addressed.

In this chapter, we show how to improve existing algorithms independent of the procedure or the definitions they use. More precisely, we first demonstrate that present algorithms tend to overlook small communities found in the neighborhood of larger, denser ones. Then, we propose and develop a *cascading* approach to community detection that greatly enhance their performance.

## 2.5 Resolution limit due to shadowing

It is known that a resolution limit exists for a large class of community detection algorithms that rely on the optimization of a quality function (e.g., modularity [40]) over non-overlapping partitions of the network [39]. Indeed, it appears that the size of the smallest detectable community is related to the size of the network. This leads to counterintuitive cases where clearly separated clusters of nodes are considered as one larger community because they are too small to be resolved by the detecting algorithm. A possible solution could be to conduct a second analysis on all detected communities in order to verify that no smaller modules can be identified.

However, the optimal covering of a network should include overlapping communities, as they capture the multiplicity of functions that a node might fulfill since nodes can then be shared between many communities. We argue that a different resolution limit, due to an effect that we refer to as *shadowing*, arises in detection algorithms that :

1. allow for *overlapping communities*;
2. rely on some *global resolution parameter*.

Shadowing typically occurs when large/dense communities act as screens preventing the detection of smaller/sparser adjacent communities. To illustrate this phenomenon, we study two families of detection algorithms based on two different paradigms of community structure, namely nodes and links communities.

### 2.5.1 Clique percolation algorithm

The clique percolation algorithm (CPA) [81] defines communities as maximal *k-clique* percolation chains, where a *k-clique* is a fully connected subgraphs of *k* nodes, and where a percolation chain is a group of cliques that can be reached from one adjacent<sup>2</sup> *k-clique* to another [33]. The complete community structure is obtained by detecting every maximal percolation chains for a given value of *k*.

It is noteworthy that the definition of a community in this context is consistent with the general description of communities outlined in Sec. 2.4. Indeed, *k-clique* percolation chains are dense by definition, and a sparser neighboring region is required to stop a *k-clique* percolation chain, ensuring that communities are denser than their surroundings. We expect shadowing since both conditions listed in Sec. 2.5 are met :

1. because percolation chains—communities—consist of *k-cliques* sharing  $k - 1$  nodes, overlapping communities occur whenever two cliques share less than  $k - 1$  nodes ;
2. the size of the cliques, *k*, acts as a global resolution parameter.

Let us explain this last point. In principle, low values of *k* lead to a more flexible detection of communities as a smaller clique size allows a wider range of configurations. However, low values of *k* often yield an excessively coarse-grained community structure of the network since percolation chains may grow almost unhindered and include a significant fraction of the nodes. In contrast, large values of *k* may leave most of the network uncharted as only large and dense clusters of nodes are then detected as communities. An *optimal value*<sup>3</sup> corresponding to a compromise between these two extreme outcomes must therefore be chosen. As this value of *k* attempts to balance these two unwanted effects for the entire network as a whole, a shadowing effect is expected to arise causing the algorithm to overlook smaller communities, or to merge them with larger ones. See Fig. 2.1 for an illustration of this effect.

### 2.5.2 Link clustering algorithm

The link clustering algorithm (LCA) [1] aggregates links—and hence the nodes they connect—into communities based on the similarity of their respective neighborhood. Denoting  $e_{ab}$  the link

2. Two *k-cliques* are said to be adjacent if they share  $k - 1$  nodes.

3. For the purpose of this study, we use the lowest value of *k* such that no extensive community is detected. As suggested in [81], the largest community is considered extensive if it contains about twice as many nodes as the second largest community. In networks with large unbreakable cliques (Internet and Protein networks), a lower ratio is used (0.25 and 0.31, respectively).

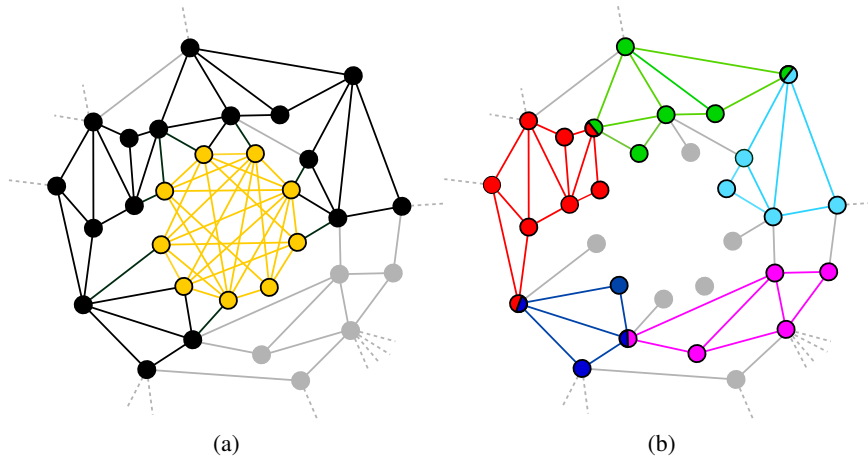


FIGURE 2.1 – **Shadowing effect for the CPA.** (a) The yellow region is the sole detectable community with  $k = 4$  or  $5$ , while its union with the black region corresponds to the community detected with  $k = 3$ . This pathological example illustrates the two undesirable extreme effects mentioned in the main text : either most of the network is detected as a single community, or only large and dense clusters are detected. No optimal value of  $k$  can be found in this case. (b) The structure of this subgraph nevertheless suggests that it could be decomposed in a dense community in the middle, surrounded by smaller communities. We see that if the links involved in the dense community (detected with  $k = 4$  or  $5$ ) were removed, a second iteration of the algorithm with  $k = 3$  would lead to the detection of several smaller communities that were overshadowed by the larger one. This illustrates the essence of the cascading detection method discussed in Sec. 2.6.

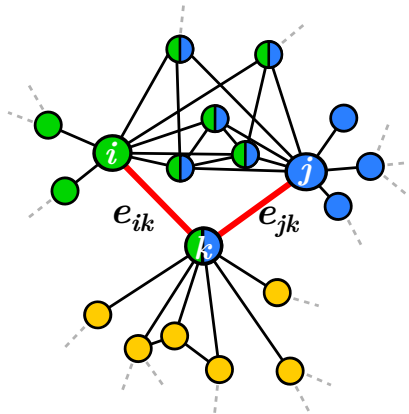


FIGURE 2.2 – **Calculation of the similarity between two links.** The sets  $n_+(i)$  and  $n_+(j)$  are respectively colored in green and blue. From Eq. (2.5.1), we have that  $S(e_{ik}, e_{jk}) = 6/13$ . Note that apart from nodes  $i$  and  $j$ , the neighboring nodes of the keystone  $k$  (colored in yellow) are not considered in this calculation of  $S(e_{ik}, e_{jk})$ .

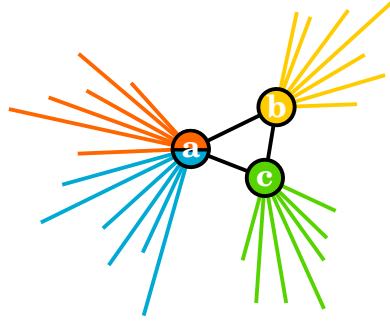


FIGURE 2.3 – **Shadowing effect for the LCA.** The pairwise unions of the three sets  $n_+(a)$ ,  $n_+(b)$  and  $n_+(c)$  contain considerably more elements than their corresponding intersections since nodes  $a$ ,  $b$  and  $c$  all have high degrees. According to Eq. 2.5.1, this implies that  $e_{ab}$ ,  $e_{bc}$  and  $e_{ac}$  share lower similarities—namely  $S(e_{ac}, e_{bc}) = S(e_{ab}, e_{bc}) = 3/22$  and  $S(e_{ab}, e_{ac}) = 3/17$ —than if the triangle had been isolated (see Sec. 2.6). It is therefore likely that these three links will be left unassigned.

between nodes  $a$  and  $b$ , the similarity of two adjacent links  $e_{ik}$  and  $e_{jk}$  (attached to a same node  $k$  called the *keystone*) is quantified through the Jaccard index

$$S(e_{ik}, e_{jk}) = \frac{|n_+(i) \cap n_+(j)|}{|n_+(i) \cup n_+(j)|}, \quad (2.5.1)$$

where  $n_+(q)$  is the set of node  $q$  and its neighbors, and  $|n_+(q)|$  is the cardinality of the set. Figure 2.2 illustrates the calculation of  $S(e_{ik}, e_{jk})$ . Once the similarity has been calculated for all adjacent pair of links, communities are built by iteratively aggregating adjacent links whose similarity exceeds a given threshold  $S_c$ . We refer to links that are left after this process (i.e., communities consisting of one single link) as *unassigned*.

Again, a shadowing effect is expected, since the two aforementioned conditions are fulfilled :

1. because communities are built by aggregating links, this algorithm naturally allows communities to overlap (to share nodes) since a node can belong to as many communities as its degree (number of links it is attached to) can allow ;
2. the similarity threshold  $S_c$  acts as a global resolution parameter as it dictates whether two links belong to the same community or not.

To elucidate the global aspect of  $S_c$ , we describe how its value is chosen (as proposed in [1]). Let us first define the density  $\rho_j$  of community  $j$  as

$$\rho_j = \frac{d_j - (n_j - 1)}{\frac{n_j(n_j - 1)}{2} - (n_j - 1)}, \quad (2.5.2)$$

where  $d_j$  and  $n_j$  are the number of links and nodes in community  $j$ , respectively. Considering that a community of  $n$  nodes must at least include  $n - 1$  links,  $\rho_j$  computes the fraction of potential “excess links” that are present in the community. The similarity threshold  $S_c$  is then chosen such that it maximizes the overall density of the communities

$$\rho(S_c) = \frac{1}{D} \sum_{j \in \mathcal{C}(S_c)} d_j \rho_j \quad (2.5.3)$$

where  $C(S_c)$  is the set of communities detected for a given  $S_c$ ,  $D$  is the total number of links assigned to communities of more than one link (i.e.,  $d_j > 1$ ). Note that  $\rho(S_c)$  is typically a well-behaved function of  $S_c$  and normally displays a single maximal plateau [1]. The value of  $S_c$  corresponding to this plateau is then selected as it leads, on average, to the denser set of communities, hence its global nature.

Following an analysis similar to Sec. 2.5.1, we expect small communities to be left undetected as they are eclipsed by larger and denser ones. This is mainly due to the use of a resolution parameter ( $S_c$ ) that cannot be adjusted locally. For instance, links in a small community could exhibit vanishing similarities because some of the associated nodes are hubs (nodes of high degree). This is especially true in the vicinity of large and dense clusters whose nodes are typically of high degree (see Fig. 2.3 for an illustration).

## 2.6 Cascading detection

Figures 2.1 and 2.3 suggest that the inability to detect small or sparse communities in the vicinity of larger or denser ones—the shadowing effect—could be circumvented by removing these structures from the networks. We formalize this idea and propose a *cascading* approach to community detection that proceeds as follows :

1. identify large or dense communities—by tuning the resolution parameter accordingly—using a given community detection algorithm ;
2. remove the internal links<sup>4</sup> of the communities identified in step 1 ;
3. repeat until no new significant communities are found.

The first iteration of this algorithm detects the communities that are normally targeted by detection algorithms, thus ensuring that the cascading approach retains the main features of the “canonical” community structure. After removal of links involved in the detected communities, a new iteration of the detection algorithm is then performed on a sparser network in which previously hidden communities are now apparent. This process is repeated until a final and more thorough covering of the network into overlapping communities is eventually obtained.

For example, in the case of the CPA, a high value of  $k$  (which leads to the traditional community structure) is selected for the first iteration of the algorithm. The network then becomes significantly sparser since all cliques of size  $k' > k$  are destroyed by the removal of internal links in step 2. Subsequent iterations of the detection algorithm can thus be conducted at lower  $k$ , unveiling finer structures, as the pathways formed by dense cluster are no longer available. The process naturally comes to a halt at  $k = 3$ , since  $k = 2$  only detects the disjoint components of the network. In the case of the LCA, the detection is stopped *before* the partition density reaches zero, for  $\rho(S_c) \simeq 0$  only detects chains of links (the keystone ensures a non-vanishing similarity), which in general are not classified as significant communities.

---

4. Internal links are defined as links that join two nodes belonging to the same community.

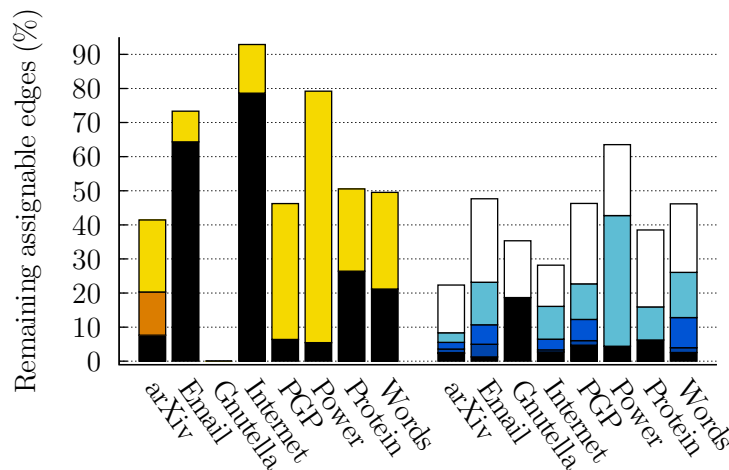


FIGURE 2.4 – **Fraction of remaining assignable links for real networks using the cascading approach.** (*left*) The number of unassigned links after one iteration of the CPA—corresponding to a typical use—is shown in yellow, and the final state is shown in black. Whenever more than 2 iterations were performed, the intermediate results are shown in orange. For the Gnutella network, the optimal value was  $k = 3$  at the first iteration, leading to an immediate complete detection of the community structure. (*right*) Results of a canonical use of the LCA are shown in white and shades of blue correspond to subsequent iterations. The final state is again shown in black. Note that all results are normalized to the number of assignable links in the original network. For the CPA, this corresponds to the number of links that belong to at least one 3-clique. For the LCA, a link is considered assignable if at least one of the two nodes it joins have a degree greater than one. Numerical results are summarized in Table 2.2 (Appendix 2.9).

It is worth mentioning that conducting this repeated analysis does not increase the computational cost significantly, since the cascading algorithm scales exactly like the community detection algorithm used at each iteration, and since the number of iterations that can be carried is small (typically less than 10). Moreover, the size of the networks (number of links and nodes) effectively decreases after each iteration, further reducing the cost.

## 2.7 Results and discussion

To investigate the efficiency and the behavior of the cascading detection, we have applied our approach to 8 network datasets : arXiv cond-mat circa 2004 [hereafter : *arXiv*] [81], university Rovira i Virgili email exchanges [*Email*] [42], Gnutella peer-to-peer data [*Gnutella*] [91], internet autonomous systems [*Internet*] [44], Pretty-Good-Privacy data exchange [*PGP*] [17], Western States Power Grid [*Power*] [110], Protein-protein interactions [*Protein*] [81] and word associations [*Words*] [81]. Their properties are summarized in Table 2.1 (Appendix 2.9).

First and foremost, our results show that cascading detection *always* improves the thoroughness of the community structure detection. Figure 2.4 shows that while a traditional use of the algorithms yields partitions with high fractions of unassigned links, the cascading approach leads to community



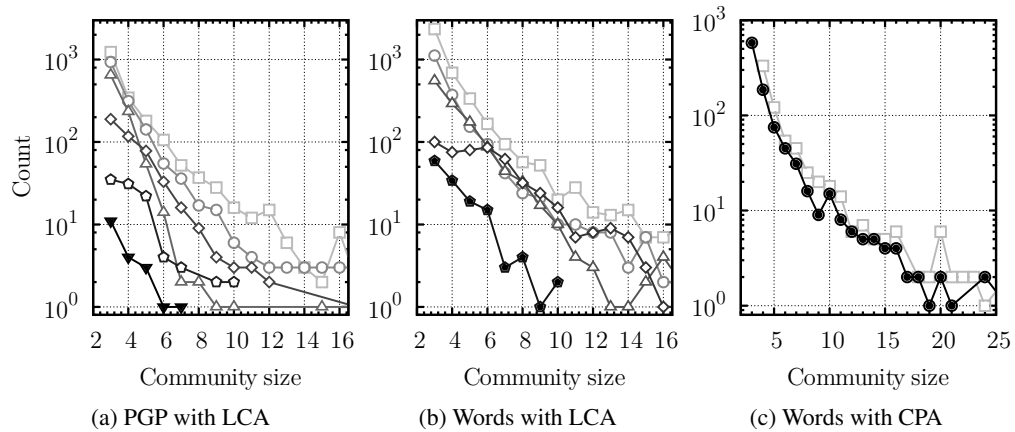


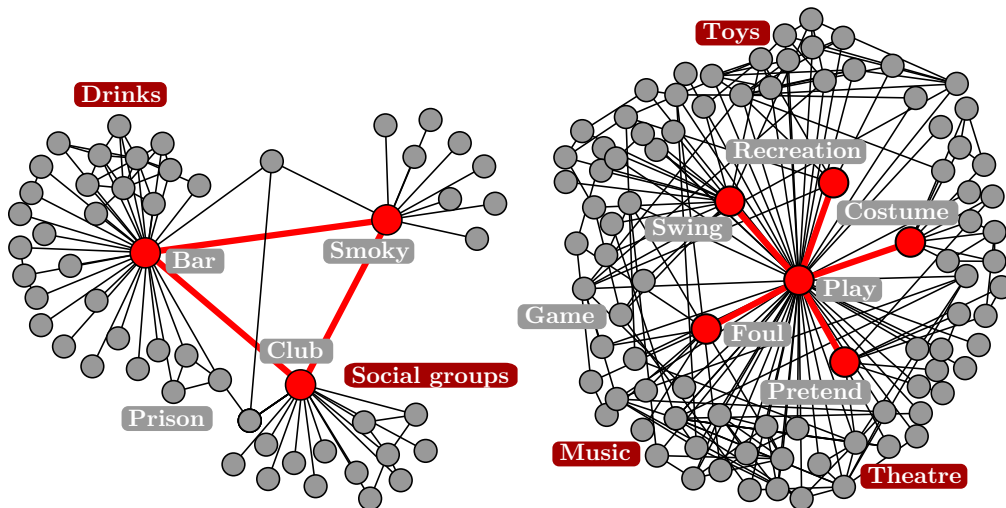
FIGURE 2.5 – **Distribution of the size of the detected communities.** Distribution of the size of the detected communities (in terms of nodes) at each iteration of the cascading approach for (a) the PGP network and (b) Words network using the LCA, and for (c) the Words network using the CPA. The distributions obtained after the first iteration are shown using light gray square markers, and subsequent iterations (whenever required) are respectively marked by circles, triangles, rhombuses, pentagons and inverted triangles. Filled black markers are used for the last iteration. Interestingly, the size of the detected communities roughly follows the same distribution at each iteration. Although this is not a direct proof, it suggests that the communities unveiled through cascading are similar to the ones detected by a “traditional” use of a community detection algorithm. In other words, these communities are significant and are not simple artifacts of the cascading approach.

structures where this fraction is significantly reduced. More precisely, the percentage of remaining assignable<sup>5</sup> links drops from 54.1% to 26.3% on average in the case of CPA, and from 41.0% to 5.3% in the case of LCA. Note how cascading detection is more efficient when applied to the LCA. This is due to the fact that the effective network gets increasingly sparse with each iteration, and that link clustering works equally well on sparse and dense networks, whereas clique percolation requires a high level of clustering to yield any results.

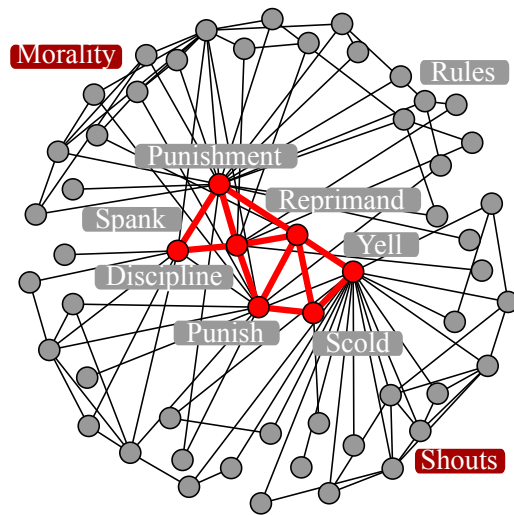
Figure 2.5 confirms that as the cascading detection proceeds, smaller—and previously masked—communities are detected, regardless of the algorithm used. For instance, Fig. 2.5 clearly shows how a significant number of 3-cliques are overlooked by “traditional” use of the CPA. However, large communities are also found after many iterations, suggesting that the shadowing effect is not restricted to small communities.

Visual inspection of the detected communities not only verifies the quality of the hidden communities, but also confirms our intuition of the shadowing effect. For instance, Fig. 2.6 (a) shows a triangle detected at the third iteration (out of five) of the LCA on the Words network. This structure was missed during the initial detection due to the high degree of its three nodes, as speculated in Fig.

5. Links that are not part of any triangles cannot be assigned to a community by the CPA since they cannot take part in any  $k$ -clique, whereas the LCA can potentially assign every link to a community, since isolated links were removed from the datasets.



(a) Triangle detected with LCA at the third iteration. (b) Star detected with LCA at the third iteration.



(c) Dense community detected with the CPA at the second iteration.

FIGURE 2.6 – **Sample of the communities detected with the cascading approach on the Words network.** The detected communities are shown (red) as well as their neighboring nodes (grey). Red and grey labels identify respectively semantic fields and individual words. .

2.3. Similarly, although  $k = 4$  was initially chosen (according to the criterion discussed in Sec. 2.5.1) for the CPA on the Words network, a second iteration using  $k = 3$  has permitted the detection of other significant communities such as the one shown in Fig. 2.6 (c).

More complex structures and correlations are also brought to light using this approach. Figure 2.6 (b) presents a star of high-degree nodes detected at the third iteration of the LCA on the word association network. None of these nodes are directly connected to each other, but they share many neighbors. Hence, once the main communities were removed—here semantic fields related to toys, theatre and music—the shadow was lifted such that this correlated, but unconnected structure could be detected. Whether this particular structure should be defined as a relevant community is up for debate. Keeping in mind that there are no consensus on the definition of a proper *community* in complex networks, the role of algorithms, and consequently of the cascading method, is to infer plausible significant structures.

Internal link removal is destructive in the sense that information about shadowed communities is lost in the process, as some of the internal links are shared by more than one community [1]. Leaving these links untouched would certainly enhance the quality of the detected communities while further reducing the uncharted portion of the network. Nevertheless, without using refined algorithm and by only resorting to our simple idea, we obtain surprisingly good results. This suggest that shadowing is not necessarily due to the density of the prominent communities but rather to the stiffness of the resolution parameter. In essence, by using a cascading approach, we allow this parameter to vary artificially from a region of the network to the other, as the algorithm is effectively applied to a new network – partially retaining the structure of the original network – at each iteration. A once rigid global parameter can now flexibly adapt to small changes in the topology of the network to better reveal subtle structures.

## 2.8 Conclusion and perspectives

In conclusion, we have managed to significantly reduce the uncharted portion of a network by assigning an important fraction of seemingly random links to relevant communities. This significant improvement in community detection will help shrink the gap between analytical models and their real network counterparts. The difficult problem of accurately modeling the dynamical properties of real networks might be better tackled if one includes complex community structure through comprehensive distributions or solved motifs [3, 52], two applications for which a reliable and complete partition is fundamental.

Moreover, this work opens the way to more subtle cascading approach, as envisioned at the end of the previous section. For instance, we could build an extreme version of the algorithm where communities are detected one by one. Such an approach would enable a perfect adaptation of the resolution parameter to the situation at hand. And while it would certainly come to significant computational cost, it could lead to the mapping of the community detection problem unto simpler problems. If we accept

to detect communities one at a time, the detection of the most significant ones can be done through well optimized methods, such as modularity optimization [74], which would otherwise be incapable of detecting overlapping communities. Finally, perhaps the most significant observation that emerges from our work could be simply stated : since community structure occurs at all scales, global partitioning of overlapping communities must be done sequentially, cascading through the organizational layers of the network.

## 2.9 Appendix : Description and properties of networks, and numerical results

TABLE 2.1 – Description and properties of real networks used in this study.

Network	$N^a$	$M^b$	$\langle k \rangle^c$	$C^d$	Ref.
arXiv	30 561	125 959	8.24	0.63	[81]
Email	1 134	5 143	9.07	0.46	[42]
Gnutella	36 682	88 328	4.82	0.01	[91]
Internet	22 963	48 436	4.22	0.23	[44]
PGP	10 680	24 316	4.55	0.27	[17]
Power	4 941	6 594	2.67	0.08	[110]
Protein	2 640	6 379	4.83	0.21	[81]
Words	7 207	31 784	8.82	0.15	[81]

a. Number of nodes.

b. Number of links (undirected projection with no multi-edges and removed self-loops).

c. Average degree, with  $\langle k \rangle = 2M/N$ .

d. Average clustering coefficient, i.e. average fraction of existing edges between neighbors of a node.

TABLE 2.2 – Summary of the results presented in Fig. 2.4.

Network	$CPA_n^a$	$CPA_f^b$	# iterations <sup>c</sup>	$LCA_n^d$	$LCA_f^e$	# iterations <sup>f</sup>
arXiv	41.4	7.6	3	22.3	2.3	7
Email	73.4	64.5	2	47.7	1.3	5
Gnutella	0	0	1	35.3	18.7	2
Internet	92.9	78.6	2	28.2	2.7	5
PGP	46.2	6.4	2	46.3	4.5	6
Power	79.2	5.5	2	63.6	4.4	3
Protein	50.5	26.4	2	38.5	6.2	3
Words	49.5	21.2	2	46.1	2.6	5

a. Percentage of remaining assignable links for a *normal* use of the CPA.

b. Percentage of remaining assignable links after the cascading approach is applied, for the CPA.

c. Number of applications of the cascading algorithm before the final state is reached, for the CPA.

d. Percentage of remaining assignable links for a *normal* use of the LCA.

e. Percentage of remaining assignable links after the cascading approach is applied, for the LCA.

f. Number of applications of the cascading algorithm before the final state is reached, for the LCA.

## Chapitre 3

# Détection et optimisation spectrale

Une méthode rapide, simple et efficace permettant d'améliorer plusieurs algorithmes de détection a été introduite au chapitre précédent. Toutefois, ce chapitre illustre aussi un problème important de la recherche sur le problème de la détection communautaire, soit la prolifération de méthodes de détection *ad hoc*, qui semblent bien fonctionner en principes mais qui ne sont pas accompagnées d'un cadre théorique solide pouvant *démontrer* leur bon fonctionnement.

Dans ce chapitre, on établira le premier volet d'une théorie unifiée de la détection communautaire. Construisant sur le découplage entre mécanisme d'exploration et critère de sélection introduit dans la Sec. 1.2.2, on considèrera formellement que le problème de la détection communautaire est solutionné lorsqu'on réussit à identifier une division en communautés  $\Omega_x$  optimale au sens d'un critère de sélection  $C$  – ou fonction objective – donné [1, 57, 59, 77, 90, 93, 105]. Afin d'obtenir une théorie analytique, on ré-exprimera le problème de la détection communautaire à l'aide d'un formalisme matriciel, basé sur les matrices  $\mathbf{A}$  et  $\mathbf{P}$  introduites aux Sec. 1.1-1.2. Conséquemment, dans cette théorie, les critères de sélection  $C$  seront fonction de  $\mathbf{A}$  et  $\mathbf{P}$ , tandis que le partitionnement spectral de noeuds sera un choix naturel de mécanisme exploratoire aux fins d'un traitement analytique.

### 3.1 Formalisme de partitionnement : motivations

Il pourrait sembler curieux de baser une théorie complète de la détection communautaire sur le partitionnement de noeuds, car ce type de décomposition communautaire ne peut pas exprimer de prime abord tous les types d'organisations mésoscopiques (cf. Sec. 1.2.1). Une étude plus rapprochée de la question révèle toutefois qu'il s'agit d'un choix judicieux, car les limitations du partitionnement de noeuds ne sont qu'apparentes. En effet, à l'aide de généralisations successives, il est possible d'étendre les types d'organisations communautaires pouvant être décrites comme des partitions à pratiquement tous les types de décompositions présentées à la Sec. 1.2.1.

Tout d'abord, on peut unifier les problèmes de partitionnement de liens et de noeuds à l'aide du graphe adjoint  $G' = (\mathcal{V}', \mathcal{E}')$  [36]. Il s'agit d'un graphe résultant d'une bijection [112] qui associe les

liens du graphe original  $G = (\mathcal{V}, \mathcal{E})$  à des noeuds dans  $G'$  (cf. annexe B.2). Partitionner les noeuds de  $G'$  selon un critère  $C'$  revient à partitionner les liens de  $G$  selon un critère  $C$ , de sorte qu'un algorithme de partitionnement de noeuds applicable à un critère de sélection arbitraire est suffisant pour solutionner les problèmes de partitionnement  $\mathcal{P}^{(\ell)}$  et  $\mathcal{P}^{(n)}$ . On peut ensuite obtenir n'importe quel type d'organisation communautaire de noeuds (e.g.  $C_c^{(n)}$ ), puisqu'une partition complète de liens  $\mathcal{P}^{(\ell)}$  est reliée à tout les types d'organisation de noeuds par induction (cf. Tableau. 1.1) Ainsi, le partitionnement de noeuds permet d'extraire tous les types d'organisations communautaires, pour des réseaux non dirigés et non dirigés ordinaires.

De plus, il a été récemment suggéré que le problème de l'identification d'une partition optimale pour des réseaux biparties et/ou directionnels peut être ramené à un problème de partitionnement sur réseaux pondérés ordinaires, à l'aide d'une transformation<sup>1</sup> qui symétrise la matrice d'adjacence [29]. Cette observation étend la portée des algorithmes de partitionnement de noeuds à tout type de réseau, car le partitionnement de réseaux pondérés est normalement une extension triviale<sup>2</sup> du partitionnement ordinaire [90].

### 3.2 Formulation matricielle du problème de partitionnement

La première étape menant à une théorie unifiée de la détection communautaire consiste à formuler le problème de partitionnement de noeuds de façon matricielle. Ici, on adoptera une formulation similaire à celle employée dans les références [5, 6, 74, 90, 118] dans un contexte de partitionnement pur (où la quantité et la taille des communautés sont connues). Cette formulation comporte deux parties distinctes, qui reflètent la distinction entre mécanisme d'exploration et critère de sélection introduite à la Sec. 1.2.2. Plus précisément, on cherchera à écrire les fonctions objectives  $C$  donnant la qualité d'une partition comme des fonctions matricielles de la matrice d'incidence communautaire  $\mathbf{P}$  et de la matrice d'adjacence  $\mathbf{A}$  (cf. chapitre 1). Ce faisant, le mécanisme d'exploration de l'espace des configurations est complètement abstrait et se manifeste uniquement par le choix de matrices  $\mathbf{P}$ , alors que le critère de sélection, ou fonction objective, est donné par une expression matricielle  $C(\mathbf{A}, \mathbf{P})$ .

Deux points de vue différents existent lorsqu'on traite du problème de partitionnement. Dans la première approche, la matrice d'adjacence est considérée comme une donnée du problème. On définit alors des fonctions objectives  $C$  à partir de cette matrice, e.g. la modularité, et on cherche la partition qui optimise cette fonction. Dans la deuxième approche, on construit une matrice d'adjacence particulièrement adaptée au problème de partitionnement (e.g. par apprentissage automatique, Ref. [9, p.6]) directement à partir du système complexe à l'étude, avant d'appliquer une méthode de partitionnement quelconque (e.g. Ref. [14, 22]). Ce deuxième point de vue est courant dans le domaine de l'agrégation de données, alors que le premier point de vue est généralement celui qui est adopté en science des

---

1. La transformation en question est :  $\mathbf{A}' = \begin{pmatrix} 0 & 1 \\ 0 & 0 \end{pmatrix} \otimes \mathbf{A} + \begin{pmatrix} 0 & 0 \\ 1 & 0 \end{pmatrix} \otimes \mathbf{A}^T$ .

2. Il suffit que le critère de sélection  $C$  tiennent compte du poids des noeuds.

réseaux. Afin d'éviter toute confusion, on se concentrera sur le premier point de vue, ne serait-ce que pour éviter les complications liées à l'utilisation d'un algorithme d'apprentissage.

### 3.2.1 Fonction objective

Dans le cadre de cet ouvrage, on considèrera des fonctions objectives particulières, soit des fonctions *additives*.

**Définition 25.** Une fonction objective est dite additive lorsqu'elle satisfait les conditions

$$C = \sum_{\alpha} C_{\alpha} \quad (3.2.1a)$$

$$C_{\alpha} = \sum_{i \in \Psi_{\alpha}} \sum_{j \in \Psi_{\alpha}} c_{ij} \quad (3.2.1b)$$

où  $C_{\alpha}$  est la fonction objective donnant la qualité du sous-graphe induit par la communauté de noeuds  $\Psi_{\alpha}^{(n)}$ , et où  $c_{ij}$  est un facteur donnant le coût associé au fait de placer les noeuds  $i$  et  $j$  dans la même communauté.

Cette catégorie de fonction exclut par exemple des fonctions du type  $C_{\alpha} = \prod_{i,j \in \Psi_{\alpha}} c_{ij}$  ou  $C = \prod_{\alpha} C_{\alpha}$ . On note au passage qu'on utilise ici une définition plus stricte de l'additivité que la plupart des références (e.g. Ref. [106]), qui ne requiert pas la condition (3.2.1b). Cette restriction est toutefois une condition nécessaire aux développements qui suivent.

Afin de pouvoir traiter le problème du partitionnement de façon formelle, on doit aussi définir une application reliant les noeuds à leur communauté.

**Définition 26.** Pour une partition de noeuds complète  $\mathcal{P}^{(n)}$ , l'application

$$\sigma_i : i \mapsto \Psi_{\alpha} \quad (3.2.2)$$

associe l'indice d'une et une seule communauté  $\Psi_{\alpha}$  à chaque noeud  $i$  du réseau.

La correspondance est unique car on traite ici d'une partition, de sorte que  $\sigma_i$  signifie littéralement "la communauté du noeud  $i$ ". En combinant les définitions (3.2.1b)-(3.2.2), la qualité  $C_{\alpha}$  de la communauté  $\alpha$  est donnée par

$$C_{\alpha} = \sum_{i,j} c_{ij} \delta_{\sigma_i, \alpha} \delta_{\sigma_j, \alpha}, \quad (3.2.3)$$

de sorte que la fonction objective totale d'une partition  $\Omega_n$  est

$$C = \sum_{\alpha} C_{\alpha} = \sum_{\alpha} \sum_{i,j} c_{ij} \delta_{\sigma_i, \alpha} \delta_{\sigma_j, \alpha} = \sum_{i,j} c_{ij} \delta_{\sigma_i, \sigma_j}, \quad (3.2.4)$$

où on a réorganisé la somme de la dernière égalité. L'utilisation de la quantité doublement indicée  $c_{ij}$  suggère qu'on pourra écrire  $C$  comme une fonction d'une matrice de coût  $\mathbf{C} = \{c_{ij}\}_{i,j=1,\dots,N}$ .

### 3.2.2 Matrice de partition

Le passage vers une formulation matricielle sera donc complet si on représente  $\delta_{\sigma_i\sigma_j}$  sous la forme d'une expression matricielle. Dans le cas d'une partition, la matrice d'incidence communautaire  $\mathbf{P}$  introduite à la définition 15 mène rapidement à une formulation correcte, car toutes les paires de rangées de  $\mathbf{P}$  satisfont l'expression

$$\delta_{\sigma_i\sigma_j} = \sum_{\alpha} p_{i\alpha} p_{j\alpha}, \quad (3.2.5)$$

tel que la fonction objective 3.2.4 peut directement être réécrite comme

$$\mathbf{C} = \sum_{i,j} c_{ij} \delta_{\sigma_i\sigma_j} = \sum_{i,j} c_{ij} \sum_{\alpha} p_{i\alpha} p_{j\alpha} = \text{Tr}(\mathbf{P}^T \mathbf{C} \mathbf{P}), \quad (3.2.6)$$

ou  $\mathbf{C}$  est une matrice symétrique de coût.

Ce résultat d'apparence simple cache toutefois une interprétation beaucoup plus puissante de la forme matricielle de la fonction objective, qui est révélée en adoptant une approche légèrement différente.

L'argument va comme suit : on considère que chaque communauté  $\alpha$  est représentée par un objet mathématique distinct, tiré d'un ensemble de  $g$  objets (lorsqu'il y a  $g$  communautés). On doit sélectionner ces objets de sorte qu'un produit entre deux objets nous permette d'exprimer le delta de Kronecker  $\delta_{\alpha\beta}$  sur les communautés (Eq. 3.2.4). L'espace vectoriel  $\mathbb{R}^g$  et le produit scalaire  $\mathbf{u}^T \mathbf{v}$  conviennent parfaitement, car il suffit alors d'associer  $g$  vecteurs orthonormaux aux communautés (e.g. pour  $g = 2$ ,  $\mathbf{e}_1^{(2)} = [1, 0]^T$  est associée à  $\alpha = 1$ ,  $\mathbf{e}_2^{(2)} = [0, 1]^T$  à  $\alpha = 2$ ) pour obtenir

$$\delta_{\sigma_i\sigma_j} = (\mathbf{e}_{\sigma_i}^{(g)})^T \mathbf{e}_{\sigma_j}^{(g)} = \begin{cases} 1 & i, j \text{ dans la même communauté} \\ 0 & \text{sinon} \end{cases}. \quad (3.2.7)$$

L'intérêt de cette formulation réside dans le nouvel éclairage qu'elle apporte sur la matrice d'incidence  $\mathbf{P}$  : cette dernière peut être vue comme une somme de produits tensoriels de vecteurs orthonormaux

$$\mathbf{P} := \begin{bmatrix} (\mathbf{e}_{\sigma_1}^{(g)})^T \\ (\mathbf{e}_{\sigma_2}^{(g)})^T \\ \vdots \\ (\mathbf{e}_{\sigma_N}^{(g)})^T \end{bmatrix} \implies \mathbf{P} = \sum_{i=1}^N \mathbf{e}_i^{(N)} \otimes (\mathbf{e}_{\sigma_i}^{(g)})^T \quad (3.2.8)$$

où  $\{\mathbf{e}_i^{(N)}\}_{i=1,\dots,N}$  est une base orthonormale dans  $\mathbb{R}^N$ , (similaire à  $\{\mathbf{e}_{\alpha}^{(g)}\}_{\alpha=1,\dots,g}$ , de plus haute dimension). Afin d'alléger la notation, on utilisera  $\mathbf{e}'_i$  pour dénoter les vecteurs dans  $\mathbb{R}^N$ , et  $\mathbf{e}_{\alpha}$  pour les vecteurs dans  $\mathbb{R}^g$ .

Cette représentation tensorielle a pour effet d'explicitier le calcul des deux propriétés fonamen-



tales de la partition, soit l'appartenance des noeuds et la taille des communautés (cf. Sec. 1.2.3), car

$$\begin{aligned}
 \mathbf{P}\mathbf{P}^T &= \left( \sum_{i=1}^N \mathbf{e}'_i \otimes \mathbf{e}_{\sigma_i}^T \right) \left( \sum_{j=1}^N \mathbf{e}'_j \otimes \mathbf{e}_{\sigma_j}^T \right)^T & \mathbf{P}^T\mathbf{P} &= \left( \sum_{i=1}^N \mathbf{e}'_i \otimes \mathbf{e}_{\sigma_i}^T \right)^T \left( \sum_{j=1}^N \mathbf{e}'_j \otimes \mathbf{e}_{\sigma_j}^T \right) \\
 &= \left( \sum_{i=1}^N \mathbf{e}'_i \otimes \mathbf{e}_{\sigma_i}^T \right) \left( \sum_{j=1}^N (\mathbf{e}'_j)^T \otimes \mathbf{e}_{\sigma_j} \right) & &= \left( \sum_{i=1}^N (\mathbf{e}'_i)^T \otimes \mathbf{e}_{\sigma_i} \right) \left( \sum_{j=1}^N \mathbf{e}'_j \otimes \mathbf{e}_{\sigma_j}^T \right) \\
 &= \sum_{i,j=1}^N \mathbf{e}'_i (\mathbf{e}'_j)^T \otimes \mathbf{e}_{\sigma_i}^T \mathbf{e}_{\sigma_j} & &= \sum_{i,j=1}^N (\mathbf{e}'_j)^T \mathbf{e}'_i \otimes \mathbf{e}_{\sigma_i} \mathbf{e}_{\sigma_j}^T \\
 &= \begin{bmatrix} \mathbf{e}_{\sigma_1}^T \mathbf{e}_{\sigma_1} & \mathbf{e}_{\sigma_1}^T \mathbf{e}_{\sigma_2} & \dots & \mathbf{e}_{\sigma_1}^T \mathbf{e}_{\sigma_N} \\ \mathbf{e}_{\sigma_2}^T \mathbf{e}_{\sigma_1} & \mathbf{e}_{\sigma_2}^T \mathbf{e}_{\sigma_2} & \dots & \mathbf{e}_{\sigma_2}^T \mathbf{e}_{\sigma_N} \\ \vdots & \vdots & \ddots & \vdots \\ \mathbf{e}_{\sigma_N}^T \mathbf{e}_{\sigma_1} & \mathbf{e}_{\sigma_N}^T \mathbf{e}_{\sigma_2} & \dots & \mathbf{e}_{\sigma_N}^T \mathbf{e}_{\sigma_N} \end{bmatrix} & &= \begin{bmatrix} s_1 & 0 & \dots & 0 \\ 0 & s_2 & \dots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \dots & s_g \end{bmatrix} \\
 &= \begin{bmatrix} 1 & \delta_{\sigma_1\sigma_2} & \dots & \delta_{\sigma_1\sigma_N} \\ \delta_{\sigma_2\sigma_1} & 1 & \dots & \delta_{\sigma_2\sigma_N} \\ \vdots & \vdots & \ddots & \vdots \\ \delta_{\sigma_N\sigma_1} & \delta_{\sigma_N\sigma_2} & \dots & 1 \end{bmatrix} & &= \text{diag}(\mathbf{s}) & (3.2.9)
 \end{aligned}$$

Autrement dit, le produit  $\mathbf{P}\mathbf{P}^T$  donne une matrice  $N \times N$  de deltas de Kronecker sur les communautés, alors que  $\mathbf{P}^T\mathbf{P}$  donne une matrice  $g \times g$  diagonale de tailles.

### 3.2.3 Représentation alternative

La représentation vectorielle de la matrice d'incidence permet aussi de relier l'expression 3.2.6 à d'autres formulations déjà existantes dans la littérature (e.g. Ref. [90]). En effet, bien que les vecteurs orthogonaux de  $\mathbb{R}^g$  soient conceptuellement simples à utiliser, ce choix n'est certainement pas unique. Tout ensemble de  $g$  vecteurs satisfaisant une relation de pseudo-orthogonalité sous un produit scalaire donné est un choix valable. Les vecteurs simplexes réguliers en sont un exemple populaire.

Essentiellement, un  $n$ -simplexe régulier est une généralisation du concept de triangle à l'espace  $\mathbb{R}^n$ ,  $\forall n$  (cf. annexe C). Les  $(n+1)$  vecteurs simplexes réguliers  $\{\mathbf{t}_\alpha\}_{\alpha=1,\dots,n+1}$  associés sont les vecteurs reliant le centre du  $n$ -simplexe régulier à ses  $n+1$  sommets (Fig. 3.1). Parmi les propriétés intéressantes de ces vecteurs, on a en particulier que le produit scalaire de deux vecteurs  $\mathbf{t}_\alpha$  différents donne toujours une quantité scalaire dépendante de la dimension de l'espace vectoriel dans lequel ils résident

$$\mathbf{t}_\alpha^T \mathbf{t}_\beta = -\frac{1}{n+1} \quad (3.2.10a)$$

$$\mathbf{t}_\alpha^T \mathbf{t}_\alpha = 1 - \frac{1}{n+1} \quad (3.2.10b)$$

(pour le choix de normalisation 3.2.10b). Ce résultat peut être résumé simplement sous la forme

$$\mathbf{t}_\alpha^T \mathbf{t}_\beta = \delta_{\alpha\beta} - \frac{1}{n+1}. \quad (3.2.11)$$

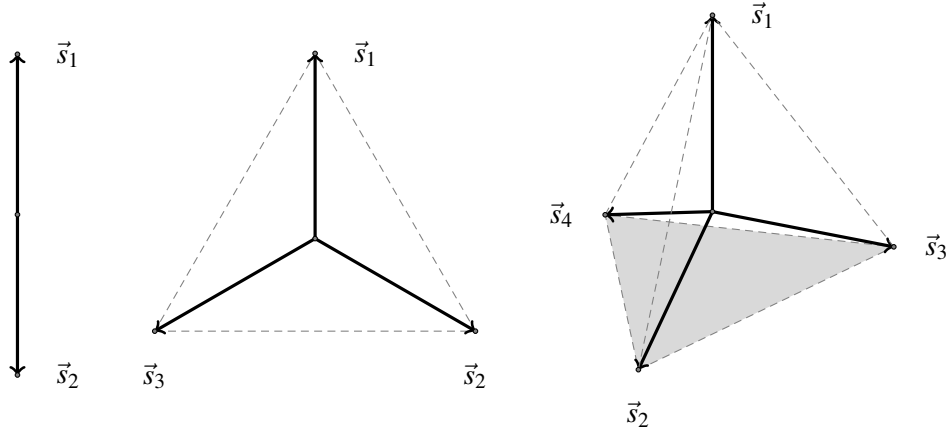


FIGURE 3.1 – **Vecteurs simplexes réguliers  $n = 1$ ,  $n = 2$  et  $n = 3$ .** En dimension  $\mathbb{R}^n$  pour  $n = 1, 2, 3$  (gauche, centre, droite), les simplexes réguliers sont simplement une droite finie de dimension fini, un triangle équilatéral et une pyramide à base triangulaire. Les vecteurs simplexes  $\vec{s}_\alpha^{(n)}$  partent de l'origine et pointent vers les coins de ces objets géométriques.

La présence d'un delta de Kronecker montre que les vecteurs simplexes sont aussi un choix valable de représentation d'une partition.

### Relation entre les vecteurs simplexes réguliers et les vecteurs orthonormaux

En fait on peut démontrer (cf. calcul détaillé dans l'annexe C) que la représentation en termes de  $g$  vecteurs simplexes réguliers dans  $\mathbb{R}^{g-1}$  est simplement la projection de la représentation en termes des  $g$  vecteurs orthonormaux de  $\mathbb{R}^g$  dans  $\mathbb{R}^{g-1}$ . Essentiellement, si on place les  $g$  vecteurs simplexes  $\{\mathbf{t}_\alpha\}_{\alpha=1,\dots,g}$  côte à côte dans une matrice  $(g-1) \times g$  de vecteurs simplexes  $\mathbf{S}$ , et les  $g$  vecteurs orthonormaux  $\{\mathbf{e}_\alpha\}_{\alpha=1,\dots,g}$  côte à côte dans une matrice  $g \times g$  de vecteurs orthonormaux  $\mathbf{I}_g$  (il s'agit en fait d'une matrice identité), on peut prouver que les deux bases sont reliées par

$$\mathbf{S} = -\mathbf{p} \left( \prod_{j=1}^{g-1} \mathbf{r}_{j,\theta_j}^{(g)} \right) \mathbf{I}_g \equiv -\mathbf{p} \mathbf{R} \mathbf{I}_g \quad (3.2.12)$$

où  $\mathbf{R}$  est une matrice orthogonale quelconque,  $\mathbf{p}$  est la matrice  $(g-1) \times g$  de projection  $\mathbb{R}^g \rightarrow \mathbb{R}^{(g-1)}$

$$\mathbf{p} = \begin{bmatrix} 0 & 1 & 0 & \dots & 0 \\ 0 & 0 & 1 & \dots & 0 \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & 0 & \dots & 1 \end{bmatrix}, \quad (3.2.13)$$

et où les  $\{\mathbf{r}_{j,\theta_j}^{(g)}\}$  sont des matrices de rotation en dimension  $g$  dans le plan  $j$ , avec des angles de rotation donnés par la relation récursive

$$\theta_j = \operatorname{acot} \left( (-1)^{g-j} \prod_{i=j+1}^{g-1} \sin \theta_i \right) \quad \theta_{g-1} = \operatorname{acot}(-1) = -\frac{\pi}{4}. \quad (3.2.14)$$

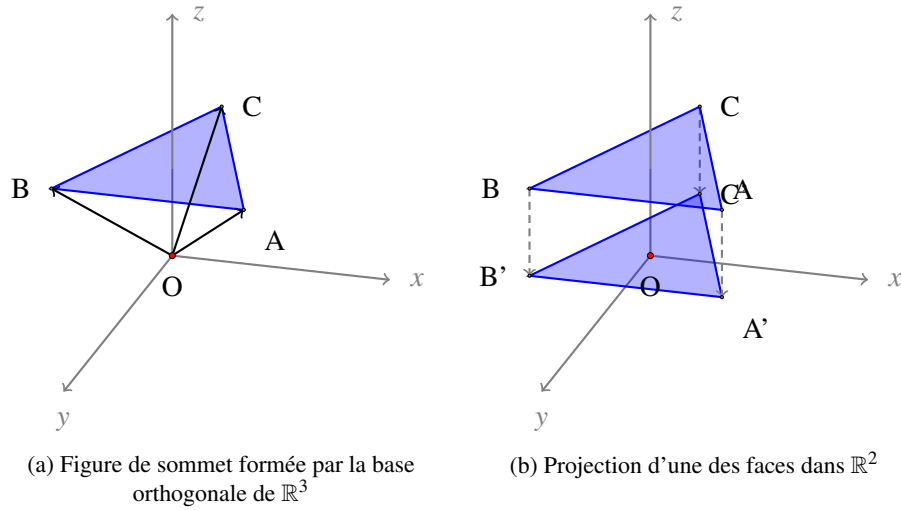


FIGURE 3.2 – **Simplexe en tant que projection d'une figure de sommet.** Le simplexe associé à  $\mathbb{R}^2$  (triangle équilatéral) correspond à une des faces de la figure de sommet formée par la base orthogonale de  $\mathbb{R}^3$ . La projection de cette face dans  $\mathbb{R}^2$  permet donc d'obtenir les vecteurs simplexes réguliers à partir des vecteurs orthogonaux. Cette projection est décomposable en deux opérations : la première opération, représentée par  $\mathbf{R}$  (Eq. 3.2.12), consiste à orienter la figure de sommet correctement par rapport au plan de projection. La deuxième partie, représentée par  $\mathbf{p}$  (Eq. 3.2.13), supprime la composante superflue (ici  $z$ ) des vecteurs  $\vec{OS}_1, \vec{OS}_2, \dots, \vec{OS}_g$  (ici  $S_1 = A, S_2 = B, S_g = C$ ).

Armé de l'expression (3.2.12), on peut passer d'une représentation de la matrice de partition à l'autre via

$$\mathbf{P}_s = \mathbf{P}_o \mathbf{S}^T \iff \begin{bmatrix} \mathbf{t}_{\sigma_1}^T \\ \mathbf{t}_{\sigma_2}^T \\ \vdots \\ \mathbf{t}_{\sigma_N}^T \end{bmatrix} = \begin{bmatrix} (\mathbf{e}_{\sigma_1}^{(g)})^T \\ (\mathbf{e}_{\sigma_2}^{(g)})^T \\ \vdots \\ (\mathbf{e}_{\sigma_N}^{(g)})^T \end{bmatrix} \begin{bmatrix} \mathbf{t}_1^T \\ \vdots \\ \mathbf{t}_g^T \end{bmatrix}, \quad (3.2.15)$$

où les indices  $s$  et  $o$  dénotent la représentation en termes de vecteurs simplexes et de vecteurs orthogonaux, respectivement. La matrice  $\mathbf{S}$  n'étant pas carrée, cette transformation n'est pas inversible, ce qui est ultimement dû à la présence d'une projection  $\mathbf{p}$ . Ainsi, un formalisme exprimé en termes des vecteurs orthogonaux est plus général qu'un formalisme exprimé en termes de vecteurs simplexes réguliers, car la transformation  $\mathbf{P}_o \rightarrow \mathbf{P}_s$  est permise, alors que l'inverse est faux.

Par définition (Eqs. 3.2.10), les deux produits de la matrice de vecteurs simplexes réguliers sont

$$\mathbf{S}^T \mathbf{S} = (\mathbf{pR})^T (\mathbf{pR}) = \mathbf{I}_{g \times g} - \frac{1}{g} \mathbf{1}_g \mathbf{1}_g^T \quad (3.2.16a)$$

$$\mathbf{S} \mathbf{S}^T = (\mathbf{pR}) (\mathbf{pR})^T = \mathbf{p} \mathbf{p}^T = \mathbf{I}_{(g-1) \times (g-1)}, \quad (3.2.16b)$$

de sorte que les propriétés équivalentes aux propriétés (3.2.9) dans le cas simplexes sont

$$\begin{aligned}
 \mathbf{P}_s \mathbf{P}_s^T &= \mathbf{P}_o \mathbf{S}^T \mathbf{S} \mathbf{P}_o^T & \mathbf{P}_s^T \mathbf{P}_s &= \mathbf{S} \mathbf{P}_o^T \mathbf{P}_o \mathbf{S}^T \\
 &= \mathbf{P}_o \left( \mathbf{I}_{g \times g} - \frac{1}{g} \mathbf{1}_g \mathbf{1}_g^T \right) \mathbf{P}_o^T & &= \mathbf{S} \text{diag}(\mathbf{s}) \mathbf{S}^T \\
 &= \mathbf{P}_o \mathbf{P}_o^T - \frac{1}{g} \mathbf{1}_N \mathbf{1}_N^T \\
 &= \begin{bmatrix} 1 - \frac{1}{g} & \delta_{\sigma_1 \sigma_2} - \frac{1}{g} & \cdots & \delta_{\sigma_1 \sigma_N} - \frac{1}{g} \\ \delta_{\sigma_2 \sigma_1} - \frac{1}{g} & 1 - \frac{1}{g} & \cdots & \delta_{\sigma_2 \sigma_N} - \frac{1}{g} \\ \vdots & \vdots & \ddots & \vdots \\ \delta_{\sigma_N \sigma_1} - \frac{1}{g} & \delta_{\sigma_N \sigma_2} - \frac{1}{g} & \cdots & 1 - \frac{1}{g} \end{bmatrix}
 \end{aligned} \tag{3.2.17}$$

Ici on a utilisé le fait que chaque noeud n'est assigné qu'à une seule communauté, de sorte que  $\mathbf{P}_o \mathbf{1}_g = \mathbf{1}_N$ .

Le produit matriciel  $\mathbf{P}_s \mathbf{P}_s^T$  conserve une interprétation simple : il s'agit d'une matrice d'assignation communautaire à laquelle on a soustrait une matrice de constantes. Le produit  $\mathbf{P}_s^T \mathbf{P}_s$  perd quant à lui toute forme d'interprétation directe, sauf dans le cas particulier où toutes les tailles sont égales (i.e.  $s_\alpha = N/g \forall \alpha$ ). En effet, dans ce cas, on obtient un multiple de la matrice identité de dimension  $g - 1$  :

$$\mathbf{P}_s^T \mathbf{P}_s = \frac{N}{g} \mathbf{S} \mathbf{S}^T = \frac{N}{g} \mathbf{I}_{(g-1) \times (g-1)}. \tag{3.2.18}$$

Dans les autres cas, seule la trace garde une interprétation directe puisque

$$\begin{aligned}
 \text{Tr}(\mathbf{P}_s^T \mathbf{P}_s) &= \text{Tr}(\mathbf{S} \text{diag}(\mathbf{s}) \mathbf{S}^T) = \text{Tr}(\text{diag}(\mathbf{s}) \mathbf{S}^T \mathbf{S}) \\
 &= \text{Tr}(\text{diag}(\mathbf{s})) - \frac{1}{g} \text{Tr}(\mathbf{1}_g^T \text{diag}(\mathbf{s}) \mathbf{1}_g) = N - \frac{N}{g},
 \end{aligned} \tag{3.2.19}$$

ce qui implique  $N - \text{Tr}(\mathbf{P}_s^T \mathbf{P}_s) = \langle s_\alpha \rangle$ .

### Utilisation des vecteurs simplexes réguliers

Les vecteurs simplexes réguliers nous permettent ultimement de représenter la partition en communautés de façon moins intuitive, mais à l'aide d'une dimension de moins, ce qui peut s'avérer crucial lorsqu'on traite des matrices de grande dimension. Cette économie de dimension est d'autant plus intéressante qu'on peut démontrer facilement que le problème d'optimisation est identique dans les deux représentations. En effet, si on substitue la relation (3.2.15) dans l'équivalent simplexe de la formulation matricielle de la fonction objective (3.2.6), on obtient la correspondance

$$\begin{aligned}
 C_s &:= \text{Tr}(\mathbf{P}_s^T \mathbf{C} \mathbf{P}_s) = \text{Tr}(\mathbf{S} \mathbf{P}_o^T \mathbf{C} \mathbf{P}_o \mathbf{S}^T) \\
 &= \text{Tr} \left( \left( \mathbf{I}_{g \times g} - \frac{1}{g} \mathbf{1}_g \mathbf{1}_g^T \right) \mathbf{P}_o^T \mathbf{C} \mathbf{P}_o \right) = \text{Tr}(\mathbf{P}_o^T \mathbf{C} \mathbf{P}_o) - \frac{1}{g} \text{Tr}(\mathbf{1}_g^T \mathbf{P}_o^T \mathbf{C} \mathbf{P}_o \mathbf{1}_g)
 \end{aligned} \tag{3.2.20}$$

entre les deux variantes de la fonction objective. Puisque chaque noeud n'est attribué qu'à un seul groupe (i.e.  $\mathbf{P}_o \mathbf{1}_g = \mathbf{1}_N$ ), on peut simplifier cette relation

$$C_s = \text{Tr}(\mathbf{P}_o^T \mathbf{C} \mathbf{P}_o) - \frac{1}{g} \text{Tr}(\mathbf{1}_N^T \mathbf{C} \mathbf{1}_N) = C_o - \frac{1}{g} \mathbf{1}_N^T \mathbf{C} \mathbf{1}_N \quad (3.2.21)$$

ce qui démontre que les deux formalismes ne diffèrent que par un facteur additif indépendant du choix de partition  $\mathbf{P}_o, \mathbf{P}_s$ . On peut ignorer ce facteur puisqu'on s'intéresse aux optima *relatifs*. De plus, sous la condition additionnelle que la fonction de coût somme à zéro pour chaque noeud (ce qui est fréquemment le cas, e.g. modularité, matrice Laplacienne Sec. 1.2.3), la correspondance est parfaite car  $\mathbf{C} \mathbf{1}_N = \mathbf{0}_N$ , i.e.

$$\text{Tr}(\mathbf{P}_s^T \mathbf{C} \mathbf{P}_s) = \text{Tr}(\mathbf{P}_o^T \mathbf{C} \mathbf{P}_o). \quad (3.2.22)$$

Malgré les économies associées au formalisme simplexe, on utilisera le formalisme orthogonal pour la suite, dans l'a perspective de conserver la généralité des résultats.

### 3.3 Optimisation spectrale continue

La réécriture de la fonction objective sous la forme matricielle (3.2.6) (ou de façon équivalente sous la forme 3.2.20) ne change pas la difficulté du problème de partition. En effet, pour optimiser la fonction objective par énumération, il faut évaluer les  $B_n$  choix de matrices  $\mathbf{P}_o$  ( $\mathbf{P}_s$ ) possibles (cf. Sec. 1.2.1). Même le problème extrêmement simplifié de la division d'un réseau de  $N$  noeuds en  $g$  groupes de tailles  $s_1, s_2, \dots, s_g$  reste NP-difficile [30, GT25, p.16], avec son espace de solutions de taille  $\binom{N}{s_1 \ s_2 \ \dots \ s_g}$ .

On cherche donc un mécanisme d'exploration permettant d'obtenir des solutions approchées dans un temps polynomial déterministe. Des méta-heuristiques pourraient être utilisées pour obtenir cette solution rapidement (cf. Sec. 1.2.2). Toutefois, par définition, le traitement analytique des solutions résultantes est limité par le processus utilisé pour y parvenir, ce qui ne cadre pas avec notre objectif avoué de développer une théorie unifiée de la détection communautaire. On cherchera donc à obtenir une solution spectrale approximée via une optimisation continue, analytique. Cette démarche est similaire à celles introduites par Newman [74, 73] et Riolo [90], mais diffère par sa généralité (applicable à toute matrice de coût  $\mathbf{C}$ ), par sa formulation (vecteurs orthogonaux plutôt que simplexes) et par sa philosophie (détection communautaire versus partition).

#### 3.3.1 Liberté sur le choix de représentation

La Sec. 3.2.3 laisse entrevoir que le choix de représentation n'est pas unique. On peut en fait utiliser cette observation à notre avantage, en écrivant la fonction objective sous une forme plus générale possédant des paramètres libres pouvant être ajustés au problème. Plus spécifiquement, on peut

ré-exprimer la matrice de partition de vecteurs orthogonaux comme

$$\mathbf{P}_o := \tilde{\mathbf{P}}\mathbf{D}\mathbf{Q}^T + \mathbf{1}_N\mathbf{a}^T \quad (3.3.1a)$$

$$\mathbf{P}_o^T := \mathbf{Q}\mathbf{D}\tilde{\mathbf{P}}^T + \mathbf{a}\mathbf{1}_N^T, \quad (3.3.1b)$$

où  $\tilde{\mathbf{P}}$  est une matrice dont les rangées sont tirées d'un ensemble de  $g$  vecteurs arbitraires de  $\mathbb{R}^g$ ,  $\mathbf{D}$  est une matrice diagonale,  $\mathbf{Q}$  est une matrice orthogonale, et où  $\mathbf{a}$  est un vecteur  $g \times 1$  quelconque. Ces matrices ont pour effet :

$\mathbf{D}$  : d'étirer les vecteurs de  $\tilde{\mathbf{P}}$  selon les  $g$  directions de  $\mathbb{R}^g$ , indépendamment ;

$\mathbf{Q}^T$  : d'effectuer une rotation / réflexion quelconque dans  $\mathbb{R}^g$  des vecteurs  $\tilde{\mathbf{P}}\mathbf{D}$  ;

$\mathbf{1}_N\mathbf{a}^T$  : d'appliquer une translation uniforme aux vecteurs  $\tilde{\mathbf{P}}\mathbf{D}\mathbf{Q}^T$ .

Tout ensemble de  $g$  vecteurs  $\tilde{\mathbf{P}}$  pouvant être transformés en vecteurs orthogonaux à l'aide de l'opération (3.3.1) devient ainsi un choix valable de représentation matricielle de la partition. La fonction objective associée prend la forme

$$\begin{aligned} C &= \text{Tr}(\mathbf{P}_o^T \mathbf{C} \mathbf{P}_o) = \text{Tr}\left(\left(\mathbf{Q}\mathbf{D}\tilde{\mathbf{P}}^T + \mathbf{a}\mathbf{1}_N^T\right)\mathbf{C}\left(\tilde{\mathbf{P}}\mathbf{D}\mathbf{Q}^T + \mathbf{1}_N\mathbf{a}^T\right)\right) \\ &= \text{Tr}\left(\mathbf{Q}\mathbf{D}\tilde{\mathbf{P}}^T \mathbf{C} \tilde{\mathbf{P}}\mathbf{D}\mathbf{Q}^T\right) + 2\text{Tr}\left(\mathbf{Q}\mathbf{D}\tilde{\mathbf{P}}^T \mathbf{C} \mathbf{1}_N\mathbf{a}^T\right) + \text{Tr}\left(\mathbf{a}\mathbf{1}_N^T \mathbf{C} \mathbf{1}_N\mathbf{a}^T\right) \\ &= \text{Tr}\left(\tilde{\mathbf{P}}^T \mathbf{C} \tilde{\mathbf{P}}\mathbf{D}^2\right) + 2\text{Tr}\left(\mathbf{Q}\mathbf{D}\tilde{\mathbf{P}}^T \mathbf{C} \mathbf{1}_N\mathbf{a}^T\right) + |\mathbf{a}|^2 \mathbf{1}_N^T \mathbf{C} \mathbf{1}_N, \end{aligned} \quad (3.3.2)$$

où on a utilisé l'invariance de la trace sous permutation cyclique et sous transposition, ainsi que l'orthogonalité de  $\mathbf{Q}$ .

Le but de la paramétrisation (3.3.1) est de retirer la division arbitraire entre l'information relative à la partition (assignation des noeuds) et la méta-information (taille des groupes). On se rappellera (cf. Eq. 3.2.9 et Def. 16, Sec. 1.2) que les produits  $\mathbf{P}_o^T \mathbf{1}_N$  et  $\mathbf{P}_o^T \mathbf{P}_o$  donnent le vecteur des tailles  $\mathbf{s}$  et la matrice diagonale des tailles  $\text{diag}(\mathbf{s})$ , respectivement. Or, l'existence de ces identités est simplement due au fait qu'on peut effectuer un comptage à partir des rangées de  $\mathbf{P}_o$  : les vecteurs correspondant à ces rangées ne contiennent pas de méta-information à proprement parler.

La paramétrisation (3.3.1) permet de faire passer l'information relative à la taille des groupes de la matrice de partition (comme dans le cas orthogonal) vers les *vecteurs de partition*. Pour ce faire, on demande explicitement que la matrice de partition transformée  $\tilde{\mathbf{P}}$  ne contienne *plus* de méta-information, de sorte que celle-ci se retrouvera automatiquement dans le paramètre libre  $\mathbf{a}$  :

$$\tilde{\mathbf{P}}^T \mathbf{1}_N = \mathbf{D}^{-1} \mathbf{Q}^T \mathbf{1}_g \quad (3.3.3a)$$

$$\tilde{\mathbf{P}}^T \tilde{\mathbf{P}} = \mathbf{I}_{g \times g}, \quad (3.3.3b)$$

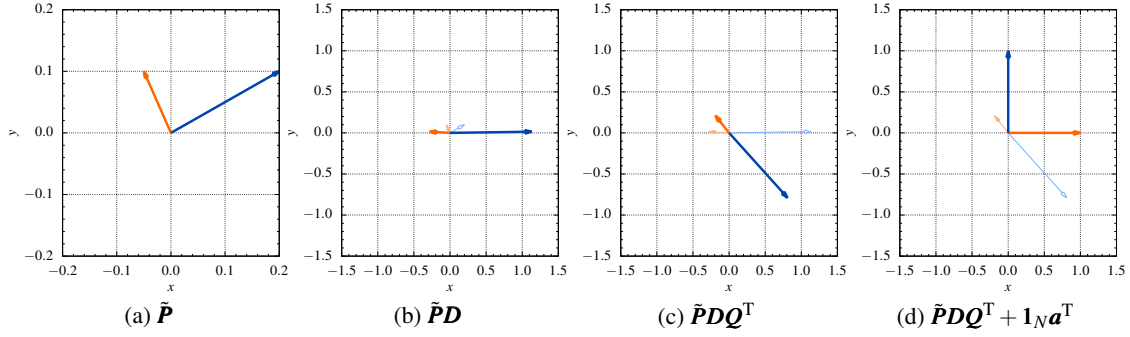


FIGURE 3.3 – **Paramétrisation de la matrice de partition.** Dans cet exemple, on étudie l’effet de la paramétrisation (3.3.1) sur  $\tilde{\mathbf{P}}$ , pour une division en deux groupes de tailles  $\mathbf{s} = [20, 80]^T$ . À chaque étape, on a laissé les vecteurs de l’étape précédente en filigrane afin de faciliter la comparaison. (a) Vecteurs de partition contenant de l’information sur la taille des groupes. (b) Vecteurs de partition étirés par un facteur  $\mathbf{D} = \mathbf{\Delta}^{1/2}$ , avec  $\mathbf{\Delta}$  les valeurs propres de  $\mathbf{M} = \text{diag}(\mathbf{s}) - \mathbf{s}\mathbf{s}^T/N + \mathbf{1}_g\mathbf{1}_g^T/N$ . (c) Vecteurs de partitions étirés et tournés dans la plan selon une rotation  $\mathbf{Q}^T$ , où  $\mathbf{Q}$  est la matrice de vecteurs propres de  $\mathbf{M}$ . (d) Vecteurs orthogonaux résultants de la transformation complète, dont la dernière étape est une translation uniforme  $(\mathbf{s} + \mathbf{1}_g)/N$ .

Les équations (3.3.3) peuvent être utilisées pour calculer les paramètres libres  $\mathbf{Q}, \mathbf{D}, \mathbf{a}$  comme suit : on prend l’inverse de la paramétrisation (3.3.1)

$$\tilde{\mathbf{P}} = (\mathbf{P}_o - \mathbf{1}_N \mathbf{a}^T) \mathbf{Q} \mathbf{D}^{-1} \quad (3.3.4a)$$

$$\tilde{\mathbf{P}}^T = \mathbf{D}^{-1} \mathbf{Q}^T (\mathbf{P}_o^T - \mathbf{a} \mathbf{1}_N^T) \quad (3.3.4b)$$

et on substitue  $\tilde{\mathbf{P}}, \tilde{\mathbf{P}}^T$  dans les Eqs. (3.3.3). Un réarrangement mène alors à

$$\mathbf{a} = \frac{1}{N} \mathbf{P}_o^T \mathbf{1}_N - \frac{1}{N} \mathbf{1}_g \quad (3.3.5a)$$

$$\mathbf{D}^{-1} \mathbf{Q}^T [\mathbf{P}_o^T \mathbf{P}_o - (\mathbf{a} \mathbf{1}_N^T \mathbf{P}_o + \mathbf{P}_o^T \mathbf{1}_N \mathbf{a}^T) + N \mathbf{a} \mathbf{a}^T] \mathbf{Q} \mathbf{D}^{-1} = \mathbf{I}_{g \times g}. \quad (3.3.5b)$$

Ce réarrangement suffit à établir que  $\mathbf{a}$  est uniquement fonction du vecteur de tailles  $\mathbf{s}$

$$\mathbf{a} = \frac{1}{N} \mathbf{P}_o^T \mathbf{1}_N - \frac{1}{N} \mathbf{1}_g = \frac{1}{N} (\mathbf{s} - \mathbf{1}_g), \quad (3.3.6)$$

un résultat qu’on peut ensuite utiliser pour montrer que  $\mathbf{M} := [\mathbf{P}_o^T \mathbf{P}_o - (\mathbf{a} \mathbf{1}_N^T \mathbf{P}_o + \mathbf{P}_o^T \mathbf{1}_N \mathbf{a}^T) + N \mathbf{a} \mathbf{a}^T]$  est une matrice symétrique fonction de ce même vecteur

$$\begin{aligned} \mathbf{M} &= \mathbf{P}_o^T \mathbf{P}_o - (\mathbf{a} \mathbf{1}_N^T \mathbf{P}_o + \mathbf{P}_o^T \mathbf{1}_N \mathbf{a}^T) + N \mathbf{a} \mathbf{a}^T \\ &= \text{diag}(\mathbf{s}) - \mathbf{a} \mathbf{s}^T - \mathbf{s} \mathbf{a}^T + N \mathbf{a} \mathbf{a}^T \\ &= \text{diag}(\mathbf{s}) - \frac{1}{N} \mathbf{s} \mathbf{s}^T + \frac{1}{N} \mathbf{1}_g \mathbf{1}_g^T. \end{aligned} \quad (3.3.7)$$

Cette matrice étant symétrique, on peut la réécrire sous la forme diagonalisée  $\mathbf{U} \mathbf{\Delta} \mathbf{U}^T$  (avec  $\mathbf{U}$  une matrice orthogonale et  $\mathbf{\Delta}$  une matrice diagonale) en solutionnant le problème aux valeurs propres

$$\mathbf{M} \mathbf{U} = \mathbf{U} \mathbf{\Delta}. \quad (3.3.8)$$

La relation (3.3.3b) est alors vérifiée en choisissant  $\mathbf{Q} = \mathbf{U}$  et  $\mathbf{D} = \mathbf{\Delta}^{1/2}$ , chose qu'une substitution directe confirme rapidement :

$$\tilde{\mathbf{P}}^T \tilde{\mathbf{P}} = \mathbf{D}^{-1} \mathbf{Q}^T \mathbf{M} \mathbf{Q} \mathbf{D}^{-1} = \mathbf{D}^{-1} \mathbf{Q}^T (\mathbf{U} \mathbf{\Delta} \mathbf{U}^T) \mathbf{Q} \mathbf{D}^{-1} = \mathbf{\Delta}^{-1/2} \mathbf{U}^T \mathbf{U} \mathbf{\Delta} \mathbf{U}^T \mathbf{U} \mathbf{\Delta}^{-1/2} = \mathbf{I}_{g \times g}. \quad (3.3.9)$$

### 3.3.2 Optimisation approximative

Malgré tout le travail effectué pour exprimer la fonction objective sous une forme matricielle, l'optimisation de l'Eq. (3.3.2) reste difficile. Paradoxalement, ce problème devient beaucoup plus facile si on étend l'espace de solution, en sélectionnant les rangées de  $\tilde{\mathbf{P}}$  dans un ensemble *continu* de vecteurs plutôt qu'un ensemble discret de  $g$  vecteurs. En effet, le problème se prête alors à des techniques d'optimisation continues et analytiques. On considère donc la relaxation  $\tilde{\mathbf{P}} \rightarrow \mathbf{R}$ , où  $\mathbf{R}$  est une matrice dont les rangées sont des vecteurs quelconques de  $\mathbb{R}^g$ . La fonction objective associée

$$\text{Tr}(\mathbf{R}^T \mathbf{C} \mathbf{R} \mathbf{D}^2) + 2\text{Tr}(\mathbf{Q} \mathbf{D} \mathbf{R}^T \mathbf{C} \mathbf{1}_N \mathbf{a}^T) + |\mathbf{a}|^2 \mathbf{1}_N^T \mathbf{C} \mathbf{1}_N, \quad (3.3.10)$$

peut alors être optimisée *exactement*. Toutefois, cette fonction objective diffère de la fonction objective initiale  $\text{Tr}(\mathbf{P}_o^T \mathbf{C} \mathbf{P}_o)$ . En effet, notre mécanisme exploratoire comporte une approximation, sous la forme d'un passage du discret vers le continu.

Les optima de cette fonction objective approximative se situent aux points où *toutes* les dérivées  $\delta_{r_{i\alpha}}$  sont égales à 0, i.e. aux solutions de

$$\frac{\partial}{\partial \mathbf{R}} [\text{Tr}(\mathbf{R}^T \mathbf{C} \mathbf{R} \mathbf{D}^2) + 2\text{Tr}(\mathbf{Q} \mathbf{D} \mathbf{R}^T \mathbf{C} \mathbf{1}_N \mathbf{a}^T) + |\mathbf{a}|^2 \mathbf{1}_N^T \mathbf{C} \mathbf{1}_N] = \mathbf{0}_{N \times g}, \quad (3.3.11)$$

où  $\partial/\partial \mathbf{X}$  est compris comme

$$\frac{\partial \sigma}{\partial \mathbf{X}} = \begin{bmatrix} \frac{\partial \sigma}{\partial x_{11}} & \frac{\partial \sigma}{\partial x_{12}} & \cdots & \frac{\partial \sigma}{\partial x_{1n}} \\ \frac{\partial \sigma}{\partial x_{21}} & \frac{\partial \sigma}{\partial x_{22}} & \cdots & \frac{\partial \sigma}{\partial x_{2n}} \\ \frac{\partial \sigma}{\partial x_{31}} & \frac{\partial \sigma}{\partial x_{32}} & \cdots & \frac{\partial \sigma}{\partial x_{3n}} \\ \vdots & \vdots & \ddots & \vdots \\ \frac{\partial \sigma}{\partial x_{m1}} & \frac{\partial \sigma}{\partial x_{m2}} & \cdots & \frac{\partial \sigma}{\partial x_{mn}} \end{bmatrix}, \quad (3.3.12)$$

avec  $\sigma$  un scalaire quelconque. Avant de procéder au calcul de l'Eq. (3.3.11), on rappelle au lecteur les identités suivantes [85, §2] :

$$\text{Matrice } \mathbf{X} \text{ sans structure} \quad \frac{\partial a}{\partial \mathbf{X}} = \mathbf{O}_{m \times n} \quad (3.3.13a)$$

$$\frac{\partial}{\partial \mathbf{X}} \text{Tr}(\mathbf{X}^T \mathbf{A}) = \mathbf{A} \quad (3.3.13b)$$

$$\frac{\partial}{\partial \mathbf{X}} \text{Tr}(\mathbf{X}^T \mathbf{A} \mathbf{X} \mathbf{B}) = \mathbf{A} \mathbf{X} \mathbf{B} + \mathbf{A}^T \mathbf{X} \mathbf{B}^T \quad (3.3.13c)$$

$$\frac{\partial}{\partial \mathbf{X}} \text{Tr}(\mathbf{X} \mathbf{A} \mathbf{X}^T) = \mathbf{X} (\mathbf{A} + \mathbf{A}^T) \quad (3.3.13d)$$

$$\text{Matrice } \mathbf{X} \text{ symétrique} \quad \frac{\partial}{\partial \mathbf{X}} \text{Tr}(\mathbf{A} \mathbf{X}) = \mathbf{A} + \mathbf{A}^T - (\mathbf{A} \circ \mathbf{I}) \quad (3.3.13e)$$

$$\text{Matrice } \mathbf{X} \text{ diagonale} \quad \frac{\partial}{\partial \mathbf{X}} \text{Tr}(\mathbf{A} \mathbf{X}) = \mathbf{A} \circ \mathbf{I} \quad (3.3.13f)$$



où  $\mathbf{X}$  est une matrice  $m \times n$ , où  $a$  (scalaire),  $\mathbf{A}$  et  $\mathbf{B}$  (matrices) ne sont pas des fonctions de  $\mathbf{X}$ , où  $\circ$  dénote le produit de Hadamard défini par  $[A \circ B]_{ij} = a_{ij}b_{ij}$ , et où les dimensions de  $\mathbf{A}$  et  $\mathbf{B}$  sont telles que les matrices résultantes sont carrées.

Appliquant les identités (3.3.13a)-(3.3.13c), on obtient directement

$$\begin{aligned} \mathbf{C}\mathbf{R}\mathbf{D}^2 + \mathbf{C}^T\mathbf{R}(\mathbf{D}^2)^T + 2\mathbf{C}\mathbf{1}_N\mathbf{a}^T\mathbf{Q}\mathbf{D} &= \mathbf{0}_{N \times g} \\ \mathbf{C}\mathbf{R}\mathbf{D}^2 + \mathbf{C}\mathbf{1}_N\mathbf{a}^T\mathbf{Q}\mathbf{D} &= \mathbf{0}_{N \times g}, \end{aligned} \quad (3.3.14)$$

où la simplification est due au fait que la matrice de coût  $\mathbf{C}$  est symétrique par définition du problème de partitionnement<sup>3</sup>, et que  $\mathbf{D}$  est diagonale, de telle façon que  $(\mathbf{D}^2)^T = \mathbf{D}^2$ . Par inspection, on obtient que la forme

$$\mathbf{R} = -\mathbf{1}_N\mathbf{a}^T\mathbf{Q}\mathbf{D}^{-1} \quad (3.3.15)$$

solutionne le problème d'optimisation.

Ce résultat est problématique, et ce pour trois raisons.

1. Toutes les rangées de  $\mathbf{R}$  sont identiques, i.e. l'optimisation continue assigne tous les noeuds *exactement* au même groupe, peu importe le choix de  $\mathbf{D}, \mathbf{Q}, \mathbf{a}$ .
2.  $\mathbf{R}$  est obtenue indépendamment de  $\mathbf{C}$ , i.e. le choix de matrice de coût n'a pas d'incidence sur la partition.
3. La fonction objective relaxée est trivialement nulle lorsque  $\mathbf{R}$  est réinjectée dans (3.3.10).

Il s'agit en fait d'un problème récurant de tout problème d'optimisation continue possédant une solution triviale. Ces solutions peuvent habituellement être évitées en imposant des contraintes *via* les multiplicateurs de Lagrange.

### Optimisation approximative contrainte

Dans le cas à l'étude, une contrainte s'offrent naturellement à nous, puisqu'on avait demandé précédemment que

$$\tilde{\mathbf{P}}^T\tilde{\mathbf{P}} = \mathbf{I}_{g \times g}, \quad (3.3.16)$$

lorsqu'on a paramétré  $\mathbf{P}_o$ . On impose donc que la matrice de partition transformée relaxée  $\mathbf{R}$  respecte la relation équivalente

$$\mathbf{R}^T\mathbf{R} = \mathbf{I}_{g \times g}, \quad (3.3.17)$$

i.e. que la matrice  $\mathbf{R}$  soit formée de  $g$  vecteurs orthonormaux contenant de l'information sur la taille des groupes.

---

3. Dans le cadre d'une partition de noeuds, une matrice de coût asymétrique n'a pas de sens. Le coût associé à l'inclusion des noeuds  $i$  et  $j$  est soit compté deux fois (via  $c_{ij}$  et  $c_{ji}$ ) soit ignoré. Toute asymétrie est donc superflue et peut être éliminée en répartissant le coût également, via la transformation  $\mathbf{C}' = (\mathbf{C} + \mathbf{C}^T)/2$ , par exemple.

En fait, l'Eq. (3.3.17) exprime  $\binom{g}{2}$  contraintes simultanément, soit les  $g(g+1)/2$  contraintes sur les produits scalaires entre les colonnes de  $\mathbf{R}$ . Conséquemment, on a besoin de  $\binom{g}{2}$  multiplicateurs de Lagrange, qu'on place dans une matrice symétrique  $\mathbf{\Lambda}$  de dimension  $g \times g$ . La contrainte complète multipliée par les multiplicateurs de Lagrange est alors donnée par

$$\frac{1}{2} \mathbf{1}_g^T [(\mathbf{\Lambda} + \mathbf{\Lambda} \circ \mathbf{I}_{g \times g}) \circ (\mathbf{R}^T \mathbf{R} - \mathbf{I}_{g \times g})] \mathbf{1}_g = \underbrace{\lambda_{11} \sum_i (r_{i1} r_{i1} - 1) + \lambda_{12} \sum_i (r_{i1} r_{i2} - 0) + \dots}_{g(g-1)/2 \text{ termes}} \quad (3.3.18)$$

L'expression  $(\mathbf{\Lambda} + \mathbf{\Lambda} \circ \mathbf{I})$  et le facteur  $1/2$  permettent de tenir compte de la symétrie de  $\mathbf{\Lambda}$  explicitement. Une identité simple (cf. démonstration dans l'annexe B.3) et la linéarité du produit de Hadamard permettent de ré-exprimer la contrainte sous forme de trace

$$\frac{1}{2} \text{Tr}(\mathbf{R}(\mathbf{\Lambda} + \mathbf{\Lambda} \circ \mathbf{I}_{g \times g})\mathbf{R}^T) - \text{Tr}(\mathbf{\Lambda}) \equiv \frac{1}{2} \mathbf{1}_g^T [(\mathbf{\Lambda} + \mathbf{\Lambda} \circ \mathbf{I}_{g \times g}) \circ (\mathbf{R}^T \mathbf{R} - \mathbf{I}_{g \times g})] \mathbf{1}_g \quad (3.3.19)$$

tel que la fonction complète à optimiser est finalement donnée par (Eq. 3.3.10)

$$\text{Tr}(\mathbf{R}^T \mathbf{C} \mathbf{R} \mathbf{D}^2) + 2 \text{Tr}(\mathbf{Q} \mathbf{D} \mathbf{R}^T \mathbf{C} \mathbf{1}_N \mathbf{a}^T) + |\mathbf{a}|^2 \mathbf{1}_N^T \mathbf{C} \mathbf{1}_N - \frac{1}{2} \text{Tr}(\mathbf{R}(\mathbf{\Lambda} + \mathbf{\Lambda} \circ \mathbf{I}_{g \times g})\mathbf{R}^T) + \text{Tr}(\mathbf{\Lambda}). \quad (3.3.20)$$

Sous cette forme contrainte, on a que  $\mathbf{R}$  et  $\mathbf{\Lambda}$  doivent solutionner

$$\frac{\partial}{\partial \mathbf{R}} \left[ \text{Tr}(\mathbf{R}^T \mathbf{C} \mathbf{R} \mathbf{D}^2) + 2 \text{Tr}(\mathbf{Q} \mathbf{D} \mathbf{R}^T \mathbf{C} \mathbf{1}_N \mathbf{a}^T) + |\mathbf{a}|^2 \mathbf{1}_N^T \mathbf{C} \mathbf{1}_N - \frac{1}{2} \text{Tr}(\mathbf{R}(\mathbf{\Lambda} + \mathbf{\Lambda} \circ \mathbf{I})\mathbf{R}^T) + \text{Tr}(\mathbf{\Lambda}) \right] = \mathbf{0}_{N \times g} \quad (3.3.21a)$$

$$\frac{\partial}{\partial \mathbf{\Lambda}} \left[ \text{Tr}(\mathbf{R}^T \mathbf{C} \mathbf{R} \mathbf{D}^2) + 2 \text{Tr}(\mathbf{Q} \mathbf{D} \mathbf{R}^T \mathbf{C} \mathbf{1}_N \mathbf{a}^T) + |\mathbf{a}|^2 \mathbf{1}_N^T \mathbf{C} \mathbf{1}_N - \frac{1}{2} \text{Tr}(\mathbf{R}(\mathbf{\Lambda} + \mathbf{\Lambda} \circ \mathbf{I})\mathbf{R}^T) + \text{Tr}(\mathbf{\Lambda}) \right] = \mathbf{0}_{g \times g} \quad (3.3.21b)$$

La solution de l'équation (3.3.21b) redonne trivialement la contrainte (3.3.17) (on utilise les identités 3.3.13e-3.3.13f), alors que l'Eq. (3.3.21a) impose que le couple  $(\mathbf{R}, \mathbf{\Lambda})$  solutionne

$$2\mathbf{C} \mathbf{R} \mathbf{D}^2 + 2\mathbf{C} \mathbf{1}_N \mathbf{a}^T \mathbf{Q} \mathbf{D} = \frac{1}{2} \mathbf{R}(\mathbf{\Lambda} + \mathbf{\Lambda}^T + 2\mathbf{\Lambda} \circ \mathbf{I}_{g \times g}). \quad (3.3.22)$$

Pour progresser plus loin, on cherche à convertir cette équation en problème aux valeurs propres, ce qui est possible uniquement si  $\mathbf{C} \mathbf{1}_N \mathbf{a}^T \mathbf{Q} \mathbf{D} = \mathbf{0}_{N \times g}$  et si  $\mathbf{\Lambda}$  est diagonale (i.e. seuls les  $g$  multiplicateurs de Lagrange sur la diagonale de  $\mathbf{\Lambda}$  sont non nuls). Le premier constat nous restreint à des cas où les rangées de la matrice de coût  $\mathbf{C}$  somment à 0 (tel que  $\mathbf{C} \mathbf{1}_N = \mathbf{0}_N$ ), tandis que le deuxième constat impose qu'on *choisisse*  $\mathbf{\Lambda}$  comme étant diagonale *a priori*.

Pour un tel choix de  $(\mathbf{C}, \mathbf{\Lambda})$  la démonstration menant à l'Eq (3.3.21b) reste valide<sup>4</sup>, et on obtient finalement que le couple  $(\mathbf{R}, \mathbf{\Lambda})$  doit solutionner le problème aux valeurs propres

$$\mathbf{C} \mathbf{R} = \mathbf{R} \mathbf{\Lambda} \mathbf{D}^{-2}. \quad (3.3.23)$$

4. Il suffit d'utiliser l'identité (3.3.13f) partout où on aurait utilisé (3.3.13e).

On notera que *la norme* des  $g$  vecteurs formant les colonnes de  $\mathbf{R}$  est encore imposée par  $\mathbf{\Lambda}$ , mais que la valeur des produits scalaires mixtes entre colonnes ne l'est plus. Autrement dit, on impose que la base  $\mathbf{R}$  soit normalisable, sans imposer qu'elle soit orthogonale<sup>5</sup>. On notera aussi que la contrainte (3.3.17) sera toujours vérifiée, puisque  $\mathbf{R}$  est constituée des vecteurs propres de  $\mathbf{C}$ , une matrice symétrique réelle.

### Solution optimale

L'équation aux valeurs propres donne la *forme* que les solutions  $(\mathbf{R}, \mathbf{\Lambda})$  doivent respecter, sans indiquer comment *sélectionner* ces solutions. En effet, une matrice  $\mathbf{C}$  symétrique possédera  $n \leq N$  vecteurs propres  $\{\mathbf{r}_i\}_{i=1,\dots,n}$  associés à  $\{\lambda_i\}_{i=1,\dots,n}$  valeurs propres réelles et non dégénérées, ce qui implique qu'on devra identifier l'optimum recherché parmi  $\binom{n}{g} \sim n^g$  solutions valides. Pour ce faire, on recourt à la fonction objective<sup>6</sup>

$$C = \text{Tr}(\mathbf{R}^T \mathbf{C} \mathbf{R} \mathbf{D}^2) \quad (3.3.24)$$

dans laquelle on substitue les solutions de (3.3.23) afin d'obtenir la qualité des optima :

$$C = \text{Tr}(\mathbf{R}^T \mathbf{C} \mathbf{R} \mathbf{D}^2) = \text{Tr}(\mathbf{R}^T \mathbf{R} \mathbf{\Lambda}) = \text{Tr}(\mathbf{\Lambda}) = \sum_{j=1}^g \lambda_j [\mathbf{D}^2]_{jj}, \quad (3.3.25)$$

pour une sélection donnée de valeurs propres  $\{\lambda_j\}_{j=1,\dots,g} \subseteq \{\lambda_k\}_{k=1,\dots,n}$ . La matrice  $\mathbf{D}$  étant une matrice positive avec des éléments placés en ordre croissant (par convention), l'Eq. (3.3.25) nous fournit une méthode exacte pour sélectionner les vecteurs propres constituant  $\mathbf{R}$  :

La matrice de partition relaxée  $\mathbf{R}$  maximisant (minimisant) la fonction objective (3.3.24) est obtenue en concaténant les  $g$  vecteurs propres associés aux  $g$  plus grandes (plus petites) valeurs propres de  $\mathbf{C}$ . L'ordre de ces vecteurs doit suivre l'ordre croissant des normes des valeurs propres associées.

## 3.4 Algorithmes de détection spectral

Dans la section précédente, on a introduit une méthode pour partitionner approximativement les noeuds d'un réseau  $G = (\mathcal{V}, \mathcal{E})$  en  $g$  groupes de tailles  $\mathbf{s}$  selon une fonction de coût additive  $C$  de somme nulle, représentée par  $\mathbf{C}$ . On doit maintenant étudier le problème de l'inversion de la relaxation  $\mathbf{R} \rightarrow \tilde{\mathbf{P}}$  et de la sélection d'un modèle, i.e. du choix de  $\mathbf{s}$ . Une fois qu'on aura établi la solution de ces deux derniers problèmes, on aura un algorithme de détection spectral complet en main.

### 3.4.1 Inversion de la relaxation

Les problèmes d'inversion et de sélection de modèles sont intrinsèquement reliés, puisque un choix de  $\mathbf{s}$  implique un choix de base pour  $\tilde{\mathbf{P}}$ , ce qui influence évidemment le passage  $\mathbf{R} \rightarrow \tilde{\mathbf{P}}$ . Afin de

5. Cette contrainte est de toutes manières superflue, car  $\mathbf{C}$  est symétrique par construction, de sorte que ses vecteurs propres forment une base de orthogonale dans  $\mathbb{R}^N$  par définition.

6. On a retiré les termes mixtes, car on se restreint au cas  $\mathbf{C} \mathbf{1}_N = \mathbf{0}_N$ .

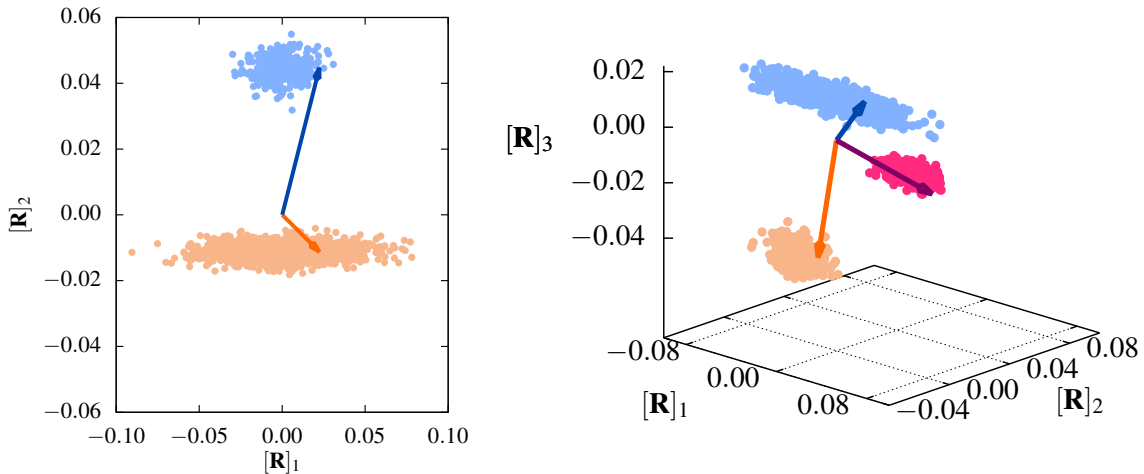


FIGURE 3.4 – **Solutions relaxées versus vecteurs de base paramétrés.** On a optimisé une matrice de coût (modularité, Sec. 24) sur des réseaux artificiels comptant des communautés de tailles (gauche)  $\mathbf{s} = [400, 1600]^T$  (droite)  $\mathbf{s} = [400, 1300, 300]^T$ . Les nuages de points représentent les extrémités des vecteurs correspondant aux rangées de la matrice  $\mathbf{R}$  optimale (solution *continue* optimale). Les vecteurs de références optimaux sont connus puisqu'on connaît la structure naturelle du réseau. Ils sont illustrés ici à titre de référence.

clarifier le développement de la présente section, on supposera pour l'instant qu'on connaît le modèle  $\mathbf{s}$  approprié pour  $\mathbf{C}$ , i.e. la taille naturelle des communautés du réseau pour notre choix de  $\mathbf{C}$ .

On se rappellera que la relaxation  $\tilde{\mathbf{P}} \rightarrow \mathbf{R}$  a pour effet de multiplier le nombre de rangées distinctes de  $\mathbf{R}$ , faisant passer ce nombre de  $g$  à potentiellement  $N \gg g$  rangées différentes. Ainsi, les rangées de la matrice  $\mathbf{R}$  optimale définissent des nuages de points dans  $\mathbb{R}^g$ , plutôt que  $g$  points dans ce même espace (cf. Fig. 3.4). Le processus d'inversion consiste à associer chaque rangée de  $\mathbf{R}$  au vecteur de base le plus près dans  $\mathbb{R}^g$ , selon une distance  $D(\mathbf{u}, \mathbf{v})$  quelconque. Pour

$$D(\mathbf{u}, \mathbf{v}) = [D_e(\mathbf{u}, \mathbf{v})]^2 = \sum_{\alpha=1}^g (u_\alpha - v_\alpha)^2, \quad (3.4.1)$$

le carré de la distance euclidienne, on a que la distance *totale*  $D_{\text{tot}}$  entre les vecteurs approximatifs  $\{\mathbf{r}_i^T\}_{i=1, \dots, N}$  (rangées de  $\mathbf{R}$ ) et les vecteurs de référence  $\{\tilde{\mathbf{p}}_i\}_{i=1, \dots, N}$  (rangées de  $\tilde{\mathbf{P}}$ ) est

$$\begin{aligned} D_{\text{tot}} &= \sum_{i=1}^N [D_e(\mathbf{r}_i, \tilde{\mathbf{p}}_i)]^2 = \sum_{i=1}^N \sum_{\alpha=1}^g (r_{i\alpha} - \tilde{p}_{i\alpha})^2 = \sum_{i=1}^N \sum_{\alpha=1}^g [\mathbf{R} - \tilde{\mathbf{P}}]_{i\alpha}^2 = \sum_{i=1}^N \sum_{\alpha=1}^g [\mathbf{R} - \tilde{\mathbf{P}}]_{i\alpha} [\mathbf{R}^T - \tilde{\mathbf{P}}^T]_{\alpha i} \\ &= \text{Tr} \left( (\mathbf{R} - \tilde{\mathbf{P}}) (\mathbf{R}^T - \tilde{\mathbf{P}}^T) \right) \\ &= \text{Tr} (\mathbf{R} \mathbf{R}^T) - 2 \text{Tr} (\mathbf{R} \tilde{\mathbf{P}}^T) + \text{Tr} (\tilde{\mathbf{P}} \tilde{\mathbf{P}}^T) \\ &= 2 \left\{ g - \text{Tr} (\mathbf{R} \tilde{\mathbf{P}}^T) \right\}. \end{aligned} \quad (3.4.2)$$

La solution discrète approximativement optimale est donc donnée par la matrice  $\tilde{\mathbf{P}}$  sujette à la paramétrisation (3.3.1) minimisant la trace (3.4.2). En principe, cette trace *pourrait* être minimisée en considérant chaque rangées de  $\mathbf{r}_i$  individuellement en la comparant aux  $g$  choix possibles de  $\{\tilde{\mathbf{p}}_\alpha\}_{\alpha=1, \dots, g}$ .

Cette comparaison exhaustive est une opération de complexité  $O(Ng)$  dans le pire des cas, ce qui ne change pas l'ordre de complexité de la méthode complète, puisque la solution du problème aux valeurs propres (3.3.23) est toujours plus complexe (même la méthode Lanczos, qui est uniquement applicable aux matrices creuses est  $O(Ng^2)$ ) [28][90, p.9].

En pratique, puisque la solution  $\mathbf{R}$  est définie à une réflexion prêt (le signe des vecteurs propres est libre), on fait face à potentiellement  $2^g N$  comparaisons, un constat qui est d'autant plus décourageant si on vise à appliquer notre algorithme spectral à répétition. Heureusement une méthode rapide, empruntée à l'analyse de Procrustes [90, p.8], permet d'identifier une matrice  $\tilde{\mathbf{P}}$  minimisant  $D_{\text{tot}}$  dans un temps  $O(g^4 N) \ll O(2^g N)$  pour  $g \gg 1$ . L'idée est de faire converger une matrice initiale non optimale  $\tilde{\mathbf{P}}_0$  vers une matrice  $\tilde{\mathbf{P}}_f$  minimisant (3.4.2), à l'aide d'une série de réflexions / rotations  $\{\mathbf{G}_i\}$  (des matrices orthogonales de dimension  $g \times g$ ).

La matrice initiale  $\tilde{\mathbf{P}}_0$  est construite en associant chaque rangée de  $\mathbf{R}$  au vecteur de référence  $\{\tilde{\mathbf{p}}_\alpha\}_{\alpha=1,\dots,g}$  le plus proche, selon le carré de la distance euclidienne (on arrondit  $\mathbf{R}$ ). Afin d'éviter toute forme de biais, il convient d'abord d'effectuer une réflexion aléatoire de ces vecteurs de référence, à l'aide d'une matrice diagonale  $\mathbf{G}_0$  d'éléments  $[\mathbf{G}_0]_{\alpha\beta} = \pm\delta_{\alpha\beta}$ .

En général, ce premier essai n'est pas le meilleur choix. On considère donc une réflexion / rotation  $\mathbf{G}_{i+1}\tilde{\mathbf{P}}_i$  du premier essai ( $i = 0$ ). La distance totale  $D_{\text{tot}}$  entre cette nouvelle matrice inconnue et la solution relaxée  $\mathbf{R}$  est alors donnée par

$$D_{\text{tot}} = \text{Tr}\left((\mathbf{R} - \tilde{\mathbf{P}}_i \mathbf{G}_{i+1})(\mathbf{R}^T - \mathbf{G}_{i+1}^T \tilde{\mathbf{P}}_i^T)\right) = \text{Tr}(\mathbf{R}^T \mathbf{R}) - 2\text{Tr}\left(\mathbf{G}_{i+1}^T \tilde{\mathbf{P}}_i^T \mathbf{R}\right) + \text{Tr}\left(\tilde{\mathbf{P}}_i^T \tilde{\mathbf{P}}_i\right). \quad (3.4.3)$$

La transformation orthogonale  $\mathbf{G}_{i+1}$  qui minimise  $D_{\text{tot}}$  est celle qui maximise  $\text{Tr}\left(\mathbf{G}_{i+1}^T \tilde{\mathbf{P}}_i^T \mathbf{R}\right)$ , car les autres termes ne dépendent pas de  $\mathbf{G}_{i+1}$ . On peut calculer cette transformation analytiquement à l'aide de la décomposition en valeurs singulières de  $\tilde{\mathbf{P}}_i^T \mathbf{R}$  :

$$\tilde{\mathbf{P}}_i^T \mathbf{R} = \mathbf{X}_i \boldsymbol{\Sigma}_i \mathbf{Y}_i^T, \quad (3.4.4)$$

où  $\mathbf{X}_i$  et  $\mathbf{Y}_i^T$  sont des matrices orthogonales et  $\boldsymbol{\Sigma}_i$  est une diagonale de valeurs singulières. L'observation clé est que  $\mathbf{Y}_i^T \mathbf{G}_i^T \mathbf{X}_i$  est orthogonale, puisqu'elle est donnée par un produit de matrices orthogonales. Ceci nous permet de borner la trace

$$\text{Tr}\left(\mathbf{G}_i^T \tilde{\mathbf{P}}_i^T \mathbf{R}\right) = \text{Tr}\left(\mathbf{G}_i^T \mathbf{X}_i \boldsymbol{\Sigma}_i \mathbf{Y}_i^T\right) = \text{Tr}\left(\mathbf{Y}_i^T \mathbf{G}_i^T \mathbf{X}_i \boldsymbol{\Sigma}_i\right) \leq \text{Tr}(\boldsymbol{\Sigma}_i), \quad (3.4.5)$$

car les éléments d'une matrice orthogonale sont tous égaux ou inférieurs à 1. En particulier, le maximum est atteint lorsque  $\mathbf{Y}_i^T \mathbf{G}_i^T \mathbf{X}_i = \mathbf{I}$ , i.e. lorsque

$$\mathbf{G}_{i+1} = \mathbf{X}_i \mathbf{Y}_i^T. \quad (3.4.6)$$

Cette dernière équation fournit  $\mathbf{G}_{i+1}$ , la réflexion / rotation à appliquer aux vecteurs de référence, comme une fonction de l'itération  $i$  précédente. On peut alors construire la matrice de partition actualisée  $\tilde{\mathbf{P}}_{i+1}$  en arrondissant les rangées de  $\mathbf{R}$  vers les nouveaux vecteurs de référence, avant de recommencer la démarche. On considère qu'on a convergé lorsque  $\tilde{\mathbf{P}}_{i+1} = \tilde{\mathbf{P}}_i$ , ou lorsque seulement  $\leq \varepsilon$

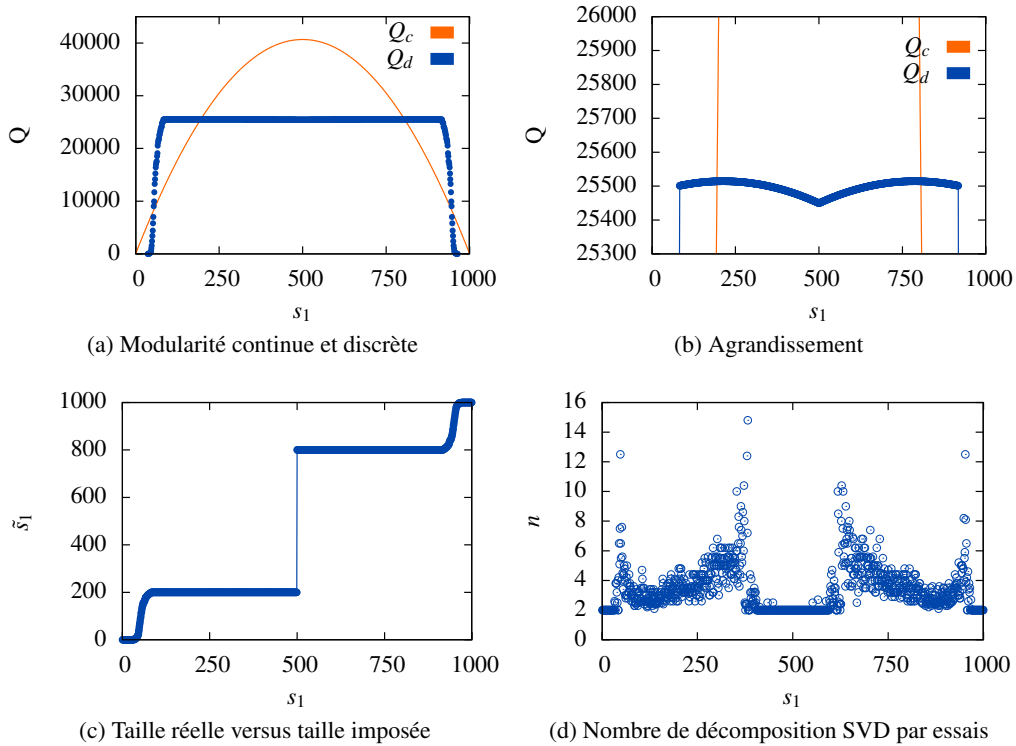


FIGURE 3.5 – **Comparaison des solutions continues et discrètes.** On maximise la modularité  $Q$  sur un réseau jouet de 2 blocs de 200 et 800 noeuds. (a) Modularité  $Q_c$  de la solution continue et modularité  $Q_d$  de la solution inversée, en fonction de la taille imposée à la communauté #1. (b) Agrandissement de la sous-figure a. Les maxima de la solution discrète apparaissent aux points où le modèle  $\mathbf{s}$  correspond à la division naturelle du réseau. (c) Taille  $\tilde{s}_1$  du groupe obtenu après inversion de la solution relaxée, en fonction de la taille imposée  $s_1$ . On aurait obtenu une droite si la procédure d'inversion avait conservé parfaitement le modèle  $\mathbf{s}$ . (d) Nombre  $n$  d'itérations de la décomposition SVD nécessaires pour converger vers une matrice de partition optimale  $\tilde{\mathbf{P}}_f$ , en fonction de la taille imposée  $s_1$ .

rangées (paramètre libre) sont réassignées, dans les cas plus difficiles. Quelques applications suffisent généralement pour obtenir  $\tilde{\mathbf{P}}_f$ . Il faut noter que (a) cette méthode n'est pas exacte, car seul un très petit sous-ensemble des  $2^s$  configurations est visité (b) le choix initial aléatoire  $\mathbf{G}_0$  détermine entièrement la solution. Il est donc suggéré de répéter la procédure d'inversion pour  $r$  choix d'orientation initiale  $\mathbf{G}_0$  (paramètre libre) et de sélectionner celle qui mène à la plus petite distance totale  $D_{\text{tot}}$ , afin de d'augmenter les chances d'identifier la solution optimale.

### 3.4.2 Sélection de modèle

La méthode présentée dans ce chapitre ne permet *pas* de sélectionner le modèle  $\mathbf{s}$  *analytiquement* pour une matrice de coût et un réseau arbitraire. En effet, la somme (3.3.25) pourrait théoriquement nous donner accès au modèle optimale, car il s'agit de la seule équation qui fait intervenir la paramétrisation (3.3.1) (à travers la matrice  $\mathbf{D}$ ) et la fonction objective (à travers les valeurs de propres

de  $\mathbf{C}$ ). Or, cette équation étant uniquement valide pour le problème relaxé, on ne peut pas s'attendre à ce que le choix de modèle  $\mathbf{s}$  menant à l'optimum absolu de la fonction objective continue corresponde aussi à l'optimum du problème original. Notre intuition est d'ailleurs rapidement confirmée lorsqu'on compare la solution continue du problème de partitionnement à la solution discrète obtenue après inversion de cette-dernière, dans un cas trivial (Fig. 3.5). On constate que l'optimum continu correspond toujours au modèle  $\mathbf{s}$  plaçant les noeuds dans des communautés de taille identique<sup>7</sup>, alors que l'optimum du problème exact est donné par un modèle collant aux communautés naturelles du réseau.

Afin de transformer notre algorithme de partitionnement spectral en algorithme de *détection* spectral, on doit absolument trouver une méthode algorithmique simple permettant d'identifier le modèle optimal pour un réseau dont la structure mésoscopique naturelle est inconnue. En théorie, on pourrait identifier ce modèle en appliquant l'algorithme complet à tous les modèles  $\mathbf{s}$  (pour  $g = 1, \dots, N$ ) possibles. Cette solution est toutefois à proscrire, car le nombre de modèles croît comme  $p(N)$ , le nombre de partition de l'entier naturel  $N$ , i.e. comme

$$p(N) \sim \frac{1}{4n\sqrt{3}} \exp\left(\pi\sqrt{\frac{2N}{3}}\right) \quad N \rightarrow \infty. \quad (3.4.7)$$

Notre solution alternative repose sur quelques observations simples :

- Le nombre de communautés naturelles est généralement révélé par la présence d'un saut dans le spectre de  $\mathbf{C}$  (voir Ref. [97] et Fig. 3.6, où plusieurs matrices  $\mathbf{C}$  sont étudiées pour des réseaux réels et synthétiques, Ref [67, 76, 84, 117] et Chap. 4, où l'existence d'un tel saut est démontrée pour des réseaux explicitement modulaires).
- Un mauvais choix de modèle  $\mathbf{s}$  mène rarement à la solution optimale, car les rangées de  $\tilde{\mathbf{P}}$  ne sont alors pas parfaitement adaptées à la configuration géométrique de la solution relaxée. Le processus d'inversion est toutefois robuste, de sorte que la partition résultante contient généralement une quantité d'information non négligeable sur le modèle optimal (cf. Fig 3.7, colonne de droite et légende). Cette propriété, qui est vue comme une faiblesse dans le cadre du partitionnement (où on cherche à *imposer* une division selon  $\mathbf{s}$ ) sera ici utilisée à notre avantage.

Ces observations peuvent être combinées pour obtenir une méthode de sélection heuristique rapide :

1. Le plus important espacement dans le spectre dicte le nombre de groupe  $g_s$  à imposer.
2. Un modèle exploratoire  $\mathbf{s}'$  (par exemple  $g_s$  groupes de taille identique) nous permet alors d'identifier la taille réelle des groupes  $\mathbf{s}$  par application itérative de l'algorithme de partitionnement. Ces applications répétées ne sont pas coûteuses, car l'étape la plus complexe (solution du problème aux valeurs propres) n'est effectuée qu'une seule fois.
3.  $g_s$  est uniquement un repère, puisque le spectre d'un réseau réel n'est pas aussi "propre" que celui d'un réseau stochastique par blocs. On essaye donc plusieurs  $g_s$  autour du choix initial.

7. Conjecture : il s'agit d'un résultat général.

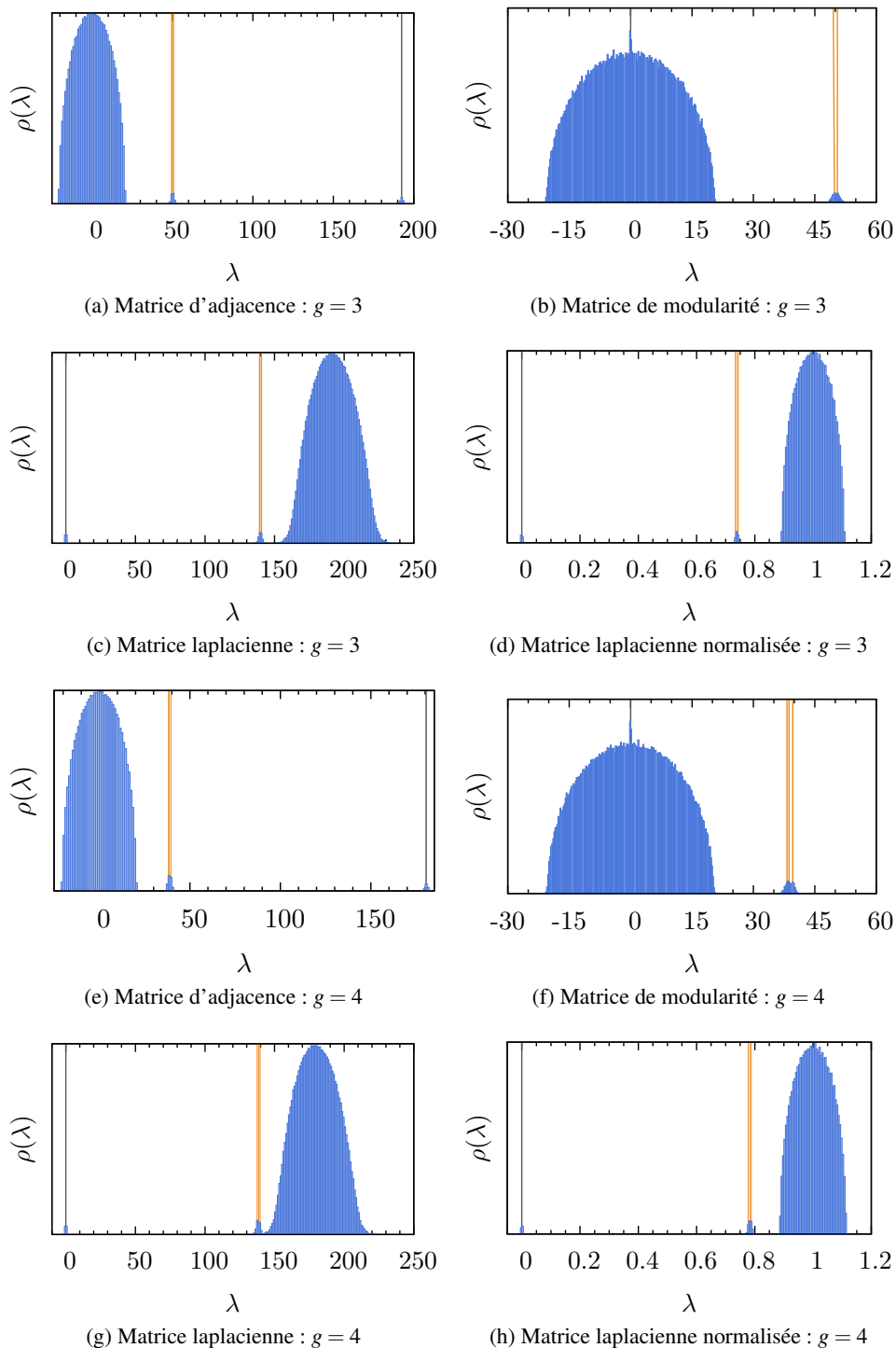


FIGURE 3.6 – **Spectres expérimentaux de plusieurs matrices de coût.** Spectres moyens de quelques matrices de coût [75] pour des réseaux de  $N = 480$  noeuds et de  $g = 3, 4$  communautés générés à l'aide du modèle stochastique par blocs [40]. Les valeurs propres isolées associées à la structure communautaire sont mise en évidence à l'aide de droites verticales oranges, tandis que les valeurs propres isolées qui sont des artefacts du choix de matrice sont illustrées à l'aide de droites verticales grises (e.g.  $L$  la matrice Laplacienne, possède toujours une valeur propre nulle.)



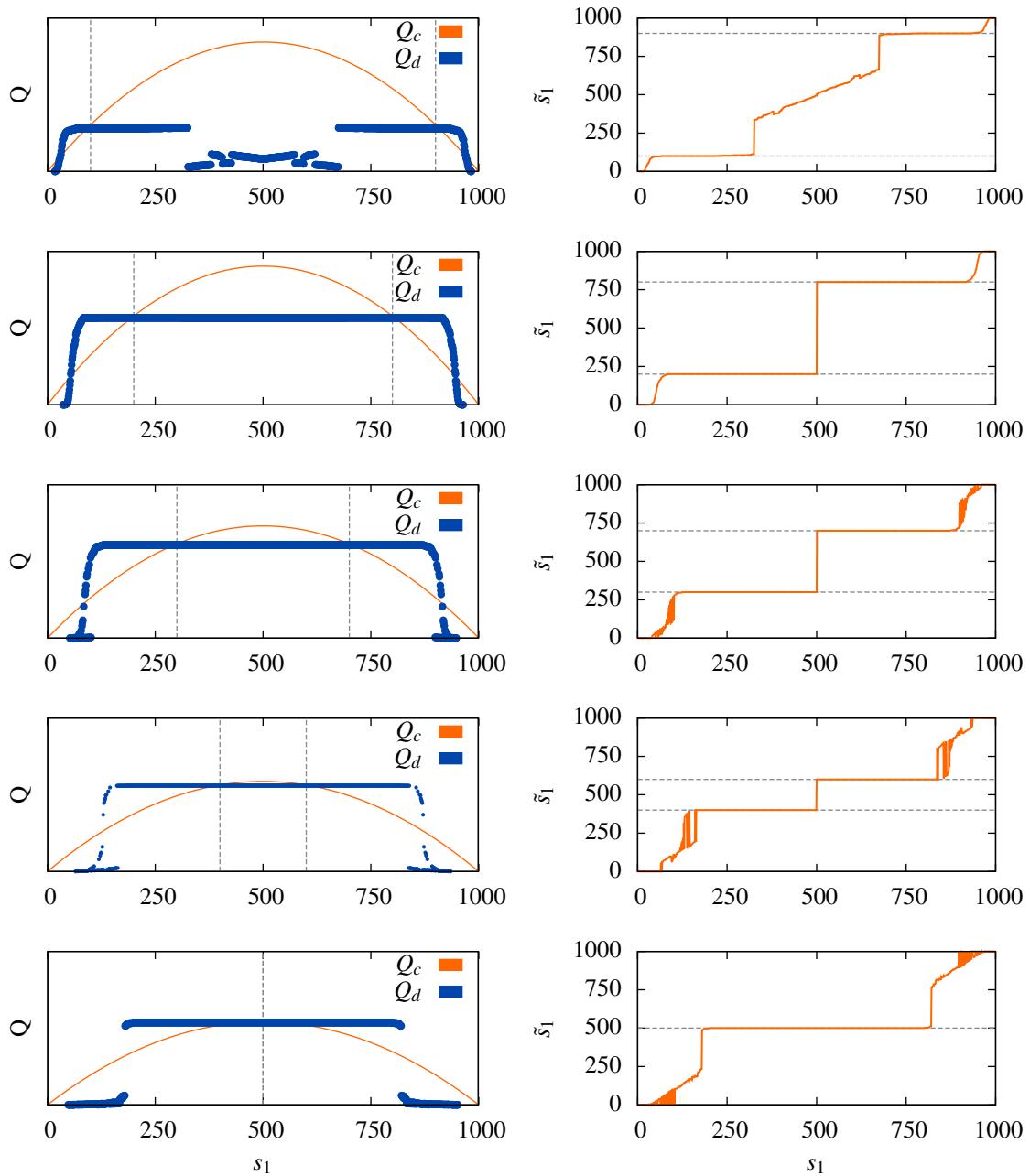


FIGURE 3.7 – **Analyse des solutions relaxées et solutions inversées.** On optimise la modularité pour des réseaux jouets de 1000 noeuds divisés en 2 communautés naturelles de tailles (de haut en bas) :  $\mathbf{s} = [100, 900]^T, [200, 800]^T, \dots, [500, 500]^T$ . (gauche) Modularité  $Q_c$  des solutions continues et modularité  $Q_d$  des solutions inversées, en fonction de la taille imposée à la communauté #1. L'optimum de  $Q_d$  est indiqué par une droite verticale. (droite) Taille  $\tilde{s}_1$  du groupe obtenu après inversion de la solution relaxée, en fonction de la taille imposée  $s_1$ . La solution inversée retombe sur la structure naturelle du réseau (en pointillé), malgré un mauvais choix initial de modèle, et ce sur un grand intervalle de taille.

### 3.4.3 Algorithme complet

Cette méthode de sélection de modèle complète notre algorithme spectral de détection. Bien que la preuve de la validité de l'algorithme soit relativement complexe, l'algorithme lui-même n'est pas très compliqué. En effet on peut le résumer en 9 étapes, soit :

Pour une fonction de coût additive de somme nulle  $C$ , représentée par une matrice de coût symétrique  $\mathbf{C}$  satisfaisant la contrainte  $\mathbf{C}\mathbf{1}_N$  ;

1. Calculer les  $\ell$  plus grandes (plus petites) valeurs propres  $\{\lambda_j\}_{j=1,\dots,\ell}$  de  $\mathbf{C}$  ainsi que les vecteurs propres associés, puis déterminer le nombre  $g_s$  de valeurs propres isolées par le plus grand espacement dans le spectre partiel  $\{\lambda_j\}_{j=1,\dots,\ell}$ .
2. Pour  $g$  près de  $g_s$ ,
  - a) Définir  $\mathbf{R}$ , la matrice des  $g$  plus grands (plus petits) vecteurs propre de  $\mathbf{C}$ .
  - b) Calculer la paramétrisation des vecteurs de référence (Sec. 3.3.1) correspondant à un modèle  $\mathbf{s}'$  de  $g$  groupes de tailles identiques.
  - c) Tant que  $\mathbf{s}'$  n'a pas convergé ( $\mathbf{s}' \neq \mathbf{s}''$ ) à  $\varepsilon'$  noeuds près, définir  $\mathbf{s}' = \mathbf{s}''$  puis :
    - i. Obtenir la solution discrète  $\tilde{\mathbf{P}}_f$  la plus proche de  $\mathbf{R}$  par analyse de Procrustes (Sec. 3.4.1).
    - ii. Extraire le nouveau modèle  $\mathbf{s}''$  de  $\tilde{\mathbf{P}}_f$ .
  - d) Évaluer la fonction objective discrète  $C(g) = \text{Tr}(\tilde{\mathbf{P}}_f^T \mathbf{C} \tilde{\mathbf{P}}_f \mathbf{D}^2)$  et définir  $\tilde{\mathbf{P}} = \tilde{\mathbf{P}}_f$  s.s.i.  $C(g)$  est plus grande (plus petite) que  $C(g') \forall$  les  $g'$  essayés jusqu'à maintenant.
3.  $\tilde{\mathbf{P}}$  est la partition naturelle des noeuds qui maximise (minimise)  $C$ .

### 3.4.4 Application

Comme nous l'avons déjà mentionné, l'objectif de ce chapitre n'est pas d'introduire un algorithme heuristique rapide, mais plutôt d'introduire un formalisme nous permettant de traiter analytiquement la détection communautaire. Tout de même, il convient de démontrer que l'algorithme présenté dans la Sec. 3.4.3 est fonctionnel. Pour ce faire, on valide notre algorithme en l'appliquant à trois réseaux complètement différents, soit un réseau mythique que tout bon algorithme se *doit* de partitionner correctement (le Zachary Karate Club [113], Fig. 3.8), un réseau en anneau dont le modèle  $\mathbf{s}$  n'est pas facile à identifier [38] (Fig. 3.9), et un réseau synthétique exhibant une transition de phase entre un régime possédant une structure mésoscopique et un régime aléatoire [68, 90] (Fig. 3.10). Le premier réseau sert simplement à démontrer que l'algorithme est fonctionnel, alors que les deux derniers réseaux sont choisis spécifiquement pour démontrer deux propriétés importantes de notre algorithme, soit le fait qu'il converge vers le modèle optimal (réseau en anneau) et le fait que les communautés identifiées sont statistiquement significatives (réseau jouet).

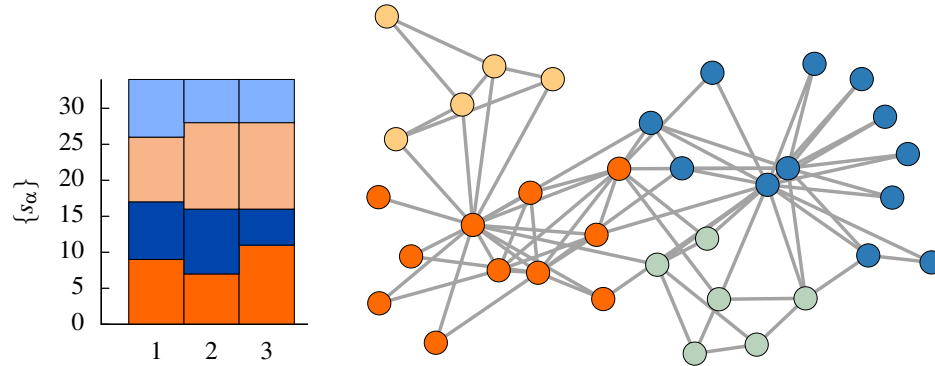


FIGURE 3.8 – **Application : Zachary Karate Club.** On optimise la modularité  $Q$  [74] avec une tolérance nulle sur la convergence du modèle  $\mathbf{s}'$ ,  $r = 5$  essais d'orientation initiale  $\mathbf{G}_0$ , et une tolérance  $\varepsilon$  nulle sur la convergence de l'analyse de Procrustes. (gauche) Évolution du modèle  $\mathbf{s}'$  en fonction de l'itération. (droite) Partition optimale détectée par notre algorithme.

### Zachary Karate Club

Le test classique de tout algorithme de détection est le Zachary Karate Club [113]. Il s'agit d'un réseau social tiré de l'étude d'un club de Karaté de 34 membres, où les liens indiquent des relations sociales en dehors du club. Lors de l'étude, une querelle a menée le groupe à se séparer en deux factions. Dans le cadre de la détection communautaire, les deux factions ont historiquement été considérées comme des groupes à identifier à partir de la structure du réseau. Avec l'avancement des méthodes de partition, il est cependant devenu apparent qu'une meilleure division en 4 groupes existe (du moins, pour plusieurs fonctions de coûts, e.g. la modularité  $Q_{CM}$  [74]). Lorsqu'on applique notre algorithme spectral avec une matrice de coût  $\mathbf{C} = \mathbf{B}$  (matrice de modularité, donnant la modularité classique à une constante multiplicative  $4M$  près), on identifie correctement la partition en 4 groupes comme étant optimale (Fig. 3.8). Ce résultat illustre bien le fait que fonction et structure ne vont pas toujours de pairs [51].

### Réseau en anneau

Le réseau en anneau est un graphe formé de  $c$  cliques complètement connectées de  $n$  noeuds (Fig. 3.9) On sait que la modularité<sup>8</sup> de la division naturelle en clique est donnée par [39]

$$Q_s = 1 - \frac{2}{n(n-1)+2} - \frac{1}{c}. \quad (3.4.8)$$

Or, si on calcule la modularité d'une division en communautés dans laquelle les cliques adjacentes sont groupées deux par deux,

$$Q_d = 1 - \frac{1}{n(n-1)+2} - \frac{2}{c}, \quad (3.4.9)$$

8. Cette version de la modularité diffère de la notre par un facteur multiplicatif  $4M$  inclut par convention chez certains auteurs [74].

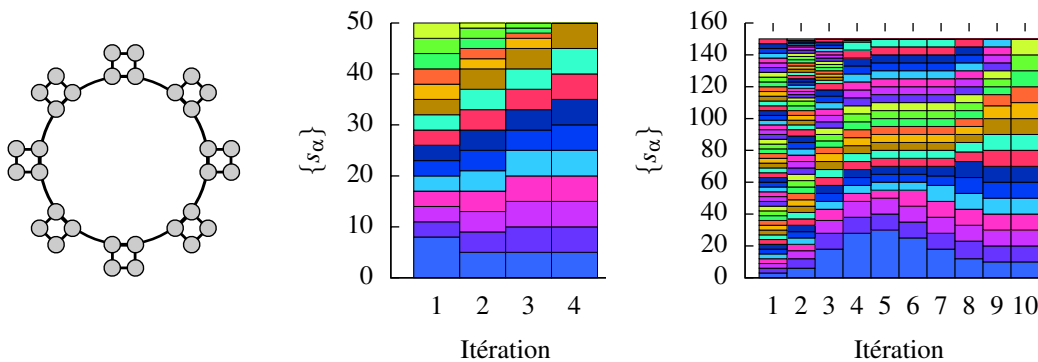


FIGURE 3.9 – **Application : Réseau en anneau et convergence du modèle.** (gauche) Réseau en anneau de  $c = 8$  cliques de taille  $n = 4$ . (centre) Évolution du modèle pour un réseau de  $c = 10$  cliques de  $n = 5$  noeud. La partition optimale de modularité  $Q = 0.809$  et de  $g = c$  communautés est correctement identifiée (droite) Évolution du modèle pour un réseau de  $c = 30$  cliques de  $n = 5$  noeud. La partition optimale de modularité  $Q = 0.888$  et de  $g = c/2$  communautés est correctement identifiée. Dans tous les cas, on a appliqué l’algorithme avec une tolérance nulle sur la convergence du modèle  $s'$ ,  $r = 20$  essais d’orientation initiale  $\mathbf{G}_0$ , et une tolérance  $\varepsilon$  de 5 rangées sur la convergence de l’analyse de Procrustes.

on se rend compte que  $Q_d$  surpasse  $Q_s$  lorsque

$$n(n-1) + 2 > c. \quad (3.4.10)$$

Autrement dit, la partition groupant les cliques deux à deux devient la division optimale lorsque la condition (3.4.10) est vérifiée. Cette inégalité implique une limite de résolution<sup>9</sup>, puisque la partition en communautés qui nous semble la plus logique devient sous-optimale.

À la figure 3.9, on vérifie que notre algorithme identifie le bon modèle lorsqu’il est appliqué à un réseau en anneau, avec la matrice de coût  $\mathbf{B}$  (matrice de modularité). Dans tous les cas, le modèle  $s'$  converge vers la partition de modularité maximale, plutôt que vers la division ‘logique’.

### Transition de détectabilité

Un troisième type de réseau artificiel permet d’étudier la validité des communautés qui sont détectées. Ces réseaux sont générés par le modèle stochastique par blocs. On analysera ce type de réseau en détails dans le prochain chapitre. Pour l’instant, il suffit de savoir que ces réseaux peuvent être paramétrés de façon à imposer qu’une fraction fixe des liens soit interne aux communautés. Le problème de détection devient progressivement plus difficile lorsque cette fraction diminue, car les communautés sont de plus en plus similaires structurellement. En fait, il existe toujours une fraction au-delà de laquelle les communautés deviennent moins denses que le reste du réseau. Ce point marque le début d’un régime où tout algorithme de détection *doit* échouer, puisque aucune information structurelle ne peut justifier la division en communautés imposée au modèle. Dans ce régime, on s’attend à ce que notre algorithme ne fasse pas mieux qu’un algorithme de division purement aléatoire.

9. À ne pas confondre avec le *shadowing* introduit au chapitre 2.

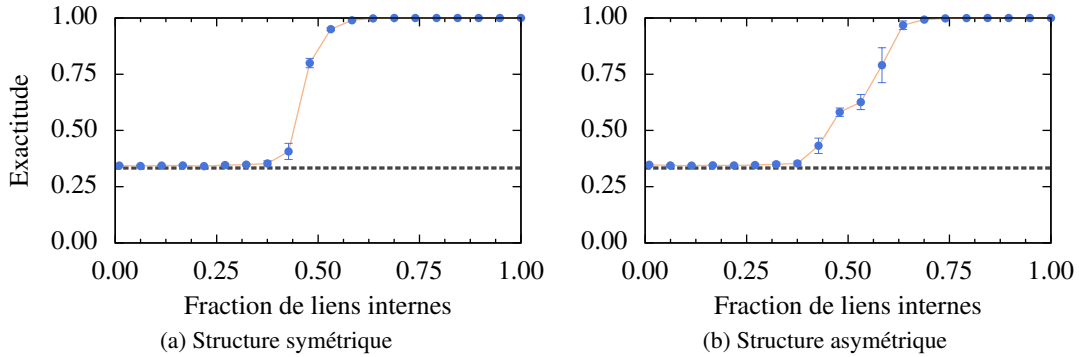


FIGURE 3.10 – **Application : Transition dans la détectabilité de la structure communautaire** Fraction des noeuds correctement classifiés par l’algorithme de détection spectral, pour un modèle par blocs et une matrice de coût  $\mathbf{B}$  (modularité, modèle nul de configuration). Les réseaux sont de taille  $N = 1800$ , et les noeuds sont placés dans des communautés de taille (a)  $\mathbf{s} = [600, 600, 600]^T$  (b)  $\mathbf{s} = [900, 600, 300]^T$ . Chaque point est obtenu en moyennant sur 50 réalisations du réseau, pour  $r = 10$  choix d’orientation initiale  $\mathbf{G}_0$ , une tolérance de  $\varepsilon = 5$  rangées sur la convergence de l’analyse de Procrustes et une tolérance de  $\varepsilon' = 10$  noeuds réassignés sur la convergence du modèle  $\mathbf{s}''$ . Le résultat qui serait obtenu pour une classification aléatoire est indiqué à l’aide d’une droite pointillée.

Afin de tester notre algorithme, on construit plusieurs instances de ce modèle, comportant toutes  $N = 1800$  noeuds, qui sont placés dans des communautés de taille  $\mathbf{s} = [600, 600, 600]^T$  (cas symétrique) et  $\mathbf{s} = [900, 600, 300]^T$  (cas asymétrique). On demande que la densité  $\rho_\alpha$  des trois communautés soit identique pour chaque réalisation, et on sélectionne ces densités de façon à maintenir le degré moyen égale à  $\langle k \rangle = 40$ . Cette précaution évite les potentiels biais liés au degré des noeuds [39, 106]. On optimise ensuite la modularité normalement, avant de comptabiliser le nombre de noeuds correctement classifié (mesure grossière de l’exactitude). Encore une fois, notre méthode réussit le test, et ne détecte *pas* de communautés structurellement non significatives (voir Fig. 3.10).

Ce constat n’est toutefois pas analytique, car il s’agit d’une observation a posteriori sur les *résultats* de notre algorithme de détection (qui est lui motivé analytiquement). Le prochain chapitre est donc dédié à une *preuve* de l’existence de cette transition de détectabilité dans le modèle stochastique par blocs. Cette analyse plus poussée nous permettra d’ailleurs d’établir de nouveaux résultats surprenants.



## Chapitre 4

# Limite intrinsèque des méthodes de détection : théorie des matrices aléatoires

Au chapitre précédent, on a démontré qu’il est possible d’extraire la division en communautés de noeuds en optimisant une fonction de coût matricielle  $\mathbf{C}$ , grâce à l’étude du spectre de  $\mathbf{C}$ . Plus spécifiquement, on a démontré qu’une partition en  $g$  groupes est donnée par une transformation des  $g$  vecteurs propres associés aux valeurs propres extrêmes de  $\mathbf{C}$ .

Dans ce chapitre, on étudie le spectre d’une matrice de coût importante [38, p.88], soit la matrice de modularité  $\mathbf{B}$ , à l’aide de la théorie des matrices aléatoires (TMA). Ce point de vue, complémentaire de celui adopté au Chap. 3, nous permet de compléter notre théorie spectrale de la détection communautaire. En effet, on démontrera que  $\mathbf{B}$  possède des valeurs propres isolées pour un réseau  $G = (\mathcal{V}, \mathcal{E})$  ayant une structure modulaire idéalisée, alors que son spectre s’apparente à celui d’une matrice purement aléatoire lorsque la structure communautaire disparaît. Ce faisant, on montre que bien que l’optimisation spectrale soit toujours applicable, le résultat est indistinguable d’une assignation aléatoire lorsque le réseau ne possède pas de structure modulaire.

### 4.1 Modèle stochastique par blocs

Puisqu’on veut obtenir des résultats analytiques exacts, on se penchera sur des réseaux simples (déjà utilisés comme banc d’essais au chapitre 3), soient les réseaux générés par le modèle stochastique par blocs (SBM). Ce modèle [27, 40] génère des réseaux complexes paramétrés par le nombre de communautés de noeuds  $g$ , la taille en noeuds  $\{s_\alpha\}_{\alpha=1,\dots,g}$  de chacun de ces groupes, ainsi que l’ensemble de probabilités  $\{p_{\alpha\beta}\}_{\alpha,\beta=1,\dots,g}$  d’existence d’un lien entre les noeuds des groupes  $\Psi_\alpha$  et  $\Psi_\beta$  (Fig. 4.1). Les boucles seront permises afin de simplifier le traitement analytique<sup>1</sup>. Les réseaux générés par le SBM sont évidemment très éloignés des systèmes réels. Leur intérêt ne réside donc pas

---

1. Cette modification a un effet négligeable sur les résultats finaux, car la présence de boucles introduit plusieurs termes asymptotiquement ( $N \rightarrow \infty$ ) nuls.

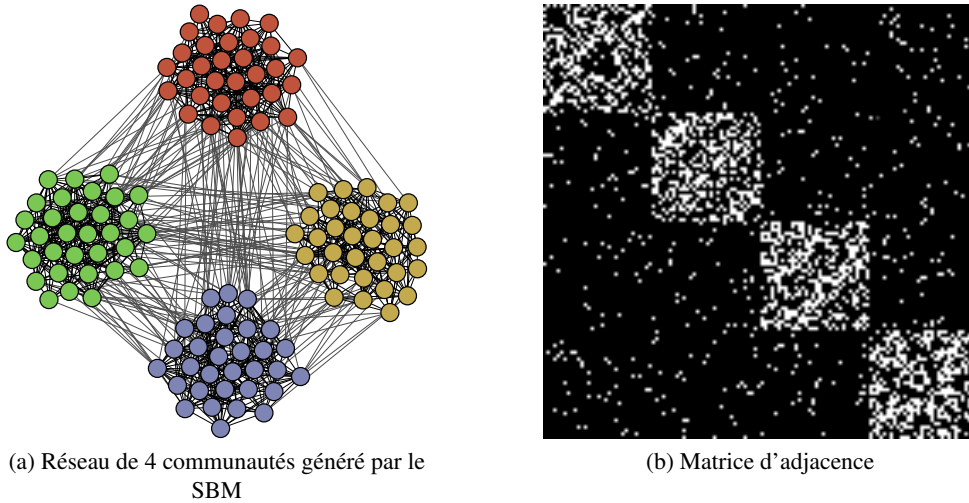


FIGURE 4.1 – **Réalisation du modèle stochastique par blocs.** Le SBM génère un réseau aléatoire de  $g$  communautés à partir d'un vecteur  $N \times 1$  de tailles  $\mathbf{s}$  et d'une matrice  $g \times g$  d'éléments  $\{p_{\alpha,\beta}\}_{\alpha,\beta=1,\dots,g}$ , où  $p_{\alpha\beta}$  est la probabilité qu'un noeud de la communauté  $\alpha$  soit connecté à un noeud de la communauté  $\beta$ . Pour l'exemple qui est illustré ici, on a utilisé un vecteur de taille  $\mathbf{s} = [24, 24, 24, 24]^T$  ainsi que deux probabilités différentes, soit une probabilité de connexion *interne*  $p = (1/3)$  et une probabilité de connexion *externe*  $q = (3/100)$ , ce qui peut être résumé sous la forme  $p_{\alpha\beta} = (1/3)\delta_{\alpha\beta} + (3/100)(1 - \delta_{\alpha\beta})$ . Le réseau résultant correspond en fait à un cas classique de réseau d'essais pour les algorithmes de détection, soit le réseau Girvan-Newman [40].

dans leur similarité avec les systèmes analysés sur une base régulière, mais plutôt dans leur extrême simplicité, qui nous permettra de prouver des résultats fondamentaux analytiquement.

#### 4.1.1 Propriétés pertinentes du modèle stochastique par blocs

Trois propriétés du SBM seront importantes pour les calculs effectués dans ce chapitre. Plus précisément, on aura besoin d'une forme générale pour les matrices d'adjacence  $\mathbf{A}_{\text{SBM}}$  et de modularité  $\mathbf{B}_{\text{SBM}}$ , ainsi que des distributions en degrés  $\mathbf{k}_\alpha$  des noeuds de chacune des communautés  $\{\Psi_\alpha\}_{\alpha=1,\dots,g}$ .

##### Matrice d'adjacence

On peut toujours écrire la matrice d'adjacence d'un réseau tiré d'un ensemble aléatoire sous la forme

$$\mathbf{A} = \langle \mathbf{A} \rangle + \mathbf{X}, \quad (4.1.1)$$

où  $\mathbf{X}$  est une matrice  $N \times N$  contenant la déviation par rapport à la matrice d'adjacence moyenne  $\langle \mathbf{A} \rangle$ . En particulier, dans le cas d'un réseau non dirigé avec des liens non corrélés, l'élément  $\bar{a}_{ij} = \bar{a}_{ji}$  de la matrice  $\langle \mathbf{A} \rangle$  correspond à la probabilité d'existence du lien  $e_{ij}$ , tandis que la matrice  $\mathbf{X}$  est une matrice symétrique de variables aléatoires indépendantes de moyenne nulle.



Dans le cas du SBM, la matrice  $\langle \mathbf{A} \rangle$  définit directement et complètement le modèle. Si on arrange les noeuds par communautés (i.e. les noeuds  $i = 1, \dots, s_1$  sont dans la communauté  $\Psi_1$ , les noeuds  $i = s_1 + 1, \dots, s_1 + s_2$  dans la communauté  $\Psi_2$ , etc.),  $\langle \mathbf{A}_{\text{SBM}} \rangle(\mathbf{s}, \{p_{\alpha\beta}\})$  prend la forme d'une matrice par blocs

$$\langle \mathbf{A} \rangle_{\text{SBM}} = \begin{bmatrix} p_{11} \mathbf{1}_{s_1} \mathbf{1}_{s_1}^T & p_{12} \mathbf{1}_{s_1} \mathbf{1}_{s_2}^T & \cdots & p_{1g} \mathbf{1}_{s_1} \mathbf{1}_{s_g}^T \\ p_{21} \mathbf{1}_{s_2} \mathbf{1}_{s_1}^T & p_{22} \mathbf{1}_{s_2} \mathbf{1}_{s_2}^T & \cdots & p_{2g} \mathbf{1}_{s_2} \mathbf{1}_{s_g}^T \\ \vdots & \vdots & \ddots & \vdots \\ p_{g1} \mathbf{1}_{s_g} \mathbf{1}_{s_1}^T & p_{g2} \mathbf{1}_{s_g} \mathbf{1}_{s_2}^T & \cdots & p_{gg} \mathbf{1}_{s_g} \mathbf{1}_{s_g}^T \end{bmatrix}. \quad (4.1.2)$$

### Distributions en degrés et moments

La fonction génératrice de probabilité (PGF) générant la distribution en degrés des noeuds de la communauté  $\Psi_\alpha$  est construite par arguments successifs (voir annexe A pour une introduction aux PGFs). Tout d'abord, on a que le polynôme

$$(1 - p_{\alpha\beta} + xp_{\alpha\beta}) \quad (4.1.3)$$

génère la distribution des résultats d'une épreuve de Bernoulli avec probabilité de succès  $p_{\alpha\beta}$ . Dans le cas du SBM, on peut interpréter (4.1.3) comme étant la PGF qui génère la distribution du nombre de liens entre *un* noeud de la communauté  $\Psi_\alpha$  et *un* noeud de la communauté  $\Psi_\beta$ . La  $s_\beta^{\text{ème}}$  puissance de (4.1.3)

$$(1 - p_{\alpha\beta} + xp_{\alpha\beta})^{s_\beta} \quad (4.1.4)$$

génère alors la somme du résultat de  $s_\beta$  épreuves de Bernoulli, i.e. une distribution binomiale (cf. propriété A.0.6 des PGF). Cette PGF peut être vue comme celle qui génère la partie de la distribution en degrés des noeuds de  $\Psi_\alpha$  due à leurs liens avec les noeuds de  $\Psi_\beta$ . En prenant le produit sur toutes les communautés, on obtient la PGF complète recherchée

$$g_\alpha(x) = \prod_{\beta=1}^g (1 - p_{\alpha\beta} + xp_{\alpha\beta})^{s_\beta}. \quad (4.1.5)$$

Afin de calculer les moments  $\langle k^n \rangle_\alpha$  de la distribution générée par  $g_\alpha(x)$ , on considère la fonction génératrice des moments (cf. eq A.0.8) associée, soit

$$m_\alpha(t) = \prod_{\beta=1}^g (1 - p_{\alpha\beta} + e^t p_{\alpha\beta})^{s_\beta} \quad \langle k^n \rangle_\alpha = \left. \frac{\partial^n m_\alpha(t)}{\partial t^n} \right|_{t=0}. \quad (4.1.6)$$

Il sera plus facile d'obtenir une expression analytique pour  $\langle k^n \rangle_\alpha$  en ré-exprimant  $m_\alpha(t)$  en termes de la quantité

$$\Delta_{\alpha\beta}^{(\ell)}(t) := \frac{s_\beta!}{(s_\beta - \ell)!} p_{\alpha\beta}^\ell e^{\ell t} (1 - p_{\alpha\beta} + e^t p_{\alpha\beta})^{s_\beta - \ell}, \quad (4.1.7)$$

sous la forme

$$m_\alpha(t) = \prod_{\beta=1}^g \Delta_{\alpha\beta}^{(0)}(t). \quad (4.1.8)$$

En effet, on a alors

$$\langle k^n \rangle_\alpha = \left. \frac{\partial^n m_\alpha(t)}{\partial t^n} \right|_{t=0} = \sum_{\mathcal{N}} \binom{n}{n_1, \dots, n_g} \prod_{\beta=1}^g \frac{\partial^{n_\beta}}{\partial t^{n_\beta}} \Delta_{\alpha\beta}^{(0)}(t) \Big|_{t=0} \quad (4.1.9)$$

où  $\mathcal{N}$  est l'ensemble des entiers  $n_1, n_2, \dots, n_g$  satisfaisant  $\sum_\ell n_\ell = n$  et où les dérivées sont données par (voir annexe B.4 pour la calcul complet)

$$\frac{\partial^{n_\beta}}{\partial t^{n_\beta}} \Delta_{\alpha\beta}^{(0)}(t) = \sum_{\ell=1}^{n_\beta} \left\{ \begin{matrix} n_\beta \\ \ell \end{matrix} \right\} \Delta_{\alpha\beta}^{(\ell)}(t). \quad (4.1.10)$$

Lorsqu'on l'évalue à  $t = 0$ , le résultat de la dérivée (4.1.10) devient

$$\theta_{\alpha\beta}^{(n_\beta)} := \left. \frac{\partial^{n_\beta}}{\partial t^{n_\beta}} \Delta_{\alpha\beta}^{(0)}(t) \right|_{t=0} = \begin{cases} \sum_{\ell=1}^{n_\beta} \left\{ \begin{matrix} n_\beta \\ \ell \end{matrix} \right\} \Delta_{\alpha\beta}^{(\ell)}(0) & = \sum_{\ell=1}^{n_\beta} \left\{ \begin{matrix} n_\beta \\ \ell \end{matrix} \right\} \frac{s_\beta!}{(s_\beta - \ell)!} p_{\alpha\beta}^\ell & n_\beta > 0 \\ \Delta_{\alpha\beta}^{(0)}(0) & = 1, & n_\beta = 0, \end{cases} \quad (4.1.11)$$

de sorte que les moments de la distribution en degrés des noeuds du bloc  $\alpha$  peuvent être “simplement” exprimés comme

$$\langle k^n \rangle_\alpha = \left. \frac{\partial^n m_\alpha(t)}{\partial t^n} \right|_{t=0} = \sum_{\mathcal{N}} \binom{n}{n_1, \dots, n_g} \prod_{\beta=1}^g \theta_{\alpha\beta}^{(n_\beta)}. \quad (4.1.12)$$

À titre d'exemple, le développement explicite de (4.1.12) pour les deux premiers moments d'un SBM arbitraire de 2 blocs donne

$$\begin{aligned} \langle k \rangle_1 &= \theta_{11}^{(1)} \theta_{12}^{(0)} + \theta_{11}^{(0)} \theta_{12}^{(1)} \\ &= s_1 p_{11} + s_2 p_{12} \end{aligned} \quad (4.1.13a)$$

$$\begin{aligned} \langle k \rangle_2 &= \theta_{21}^{(1)} \theta_{22}^{(0)} + \theta_{21}^{(0)} \theta_{22}^{(1)} \\ &= s_1 p_{21} + s_2 p_{22} \end{aligned} \quad (4.1.13b)$$

$$\begin{aligned} \langle k^2 \rangle_1 &= \theta_{11}^{(2)} \theta_{12}^{(0)} + \theta_{11}^{(0)} \theta_{12}^{(2)} + 2\theta_{11}^{(1)} \theta_{12}^{(1)} \\ &= s_1 p_{11} [1 + (s_1 - 1)p_{11}] + s_2 p_{12} [1 + (s_2 - 1)p_{12}] + 2s_1 s_2 p_{11} p_{12} \end{aligned} \quad (4.1.13c)$$

$$\begin{aligned} \langle k^2 \rangle_2 &= \theta_{21}^{(2)} \theta_{22}^{(0)} + \theta_{21}^{(0)} \theta_{22}^{(2)} + 2\theta_{21}^{(1)} \theta_{22}^{(1)} \\ &= s_1 p_{21} [1 + (s_1 - 1)p_{21}] + s_2 p_{22} [1 + (s_2 - 1)p_{22}] + 2s_1 s_2 p_{21} p_{22} \end{aligned} \quad (4.1.13d)$$

L'intérêt du résultat (4.1.12) réside surtout dans le fait qu'il permet de vérifier que les degrés des noeuds du groupe  $\Psi_\alpha$  d'un SBM arbitraire sont concentrés autour de la moyenne lorsqu'on se place dans la limite thermodynamique  $s_\beta \rightarrow \infty$ ,  $\forall \beta \neq \alpha$ . En effet, on peut montrer que la largeur relative des

distributions en degrés tend vers 0 dans cette limite :

$$\begin{aligned}
 \frac{\sqrt{\langle k^2 \rangle_\alpha - \langle k \rangle_\alpha^2}}{\langle k \rangle_\alpha} &= \frac{\sqrt{\sum_{\mathcal{N}_2} \binom{n}{n_1, \dots, n_g} \prod_{\beta=1}^g \theta_{\alpha\beta}^{(n_\beta)} - \left[ \sum_{\mathcal{N}_1} \binom{n}{n_1, \dots, n_g} \prod_{\beta=1}^g \theta_{\alpha\beta}^{(n_\beta)} \right]^2}}{\sum_{\mathcal{N}_1} \binom{n}{n_1, \dots, n_g} \prod_{\beta=1}^g \theta_{\alpha\beta}^{(n_\beta)}} \\
 &= \frac{\sqrt{\sum_{\beta_1=1}^g \theta_{\alpha\beta_1}^{(2)} + \sum_{\beta_1 \neq \beta_2} \theta_{\alpha\beta_1}^{(1)} \theta_{\alpha\beta_2}^{(1)} - \left[ \left( \sum_{\beta_1=1}^g \theta_{\alpha\beta_1}^{(1)} \right) \left( \sum_{\beta_2=1}^g \theta_{\alpha\beta_2}^{(1)} \right) \right]}}{\sum_{\beta_1=1}^g \theta_{\alpha\beta_1}^{(1)}} \\
 &= \frac{\sqrt{\sum_{\beta_1=1}^g \left[ \theta_{\alpha\beta_1}^{(2)} - \left( \theta_{\alpha\beta_1}^{(1)} \right)^2 \right]}}{\sum_{\beta_1=1}^g \theta_{\alpha\beta_1}^{(1)}} = \frac{\sqrt{\sum_{\beta_1=1}^g s_{\beta_1} p_{\alpha\beta_1} (1 - p_{\alpha\beta_1})}}{\sum_{\beta_1=1}^g s_{\beta_1} p_{\alpha\beta_1}} \sim \frac{1}{\sqrt{s_g}} \rightarrow 0, \quad (4.1.14)
 \end{aligned}$$

où on a ré-exprimé les sommes sur les coefficients multinomiaux en énumérant manuellement toutes les combinaisons possibles d'entiers  $\mathcal{N}_1$ ,  $\mathcal{N}_2$ , et où le résultat asymptotique  $\sim s_g^{-1/2}$  suppose un arrangement par taille des groupes  $s_1 < s_2 < \dots < s_g$ .

### Matrice de modularité

On se rappellera que la matrice de modularité  $\mathbf{B}$  est reliée à la matrice d'adjacence par une transformation simple (Eq. 1.2.25)

$$\mathbf{B} = \mathbf{A} - \mathbf{A}_{H_0}, \quad (4.1.15)$$

où  $\mathbf{A}_{H_0}$  est une matrice dont les éléments  $\{a_{ij}^{(H_0)}\}$  donnent le nombre de liens attendu entre les noeuds  $(i, j)$ , en supposant que  $\mathbf{A}$  est générée par un modèle nul  $H_0$ . En substituant  $\mathbf{A}$  par la paramétrisation (4.1.1), la matrice de modularité peut aussi être exprimée comme une fonction de la matrice de déviation aléatoire  $\mathbf{X}$ , sous la forme

$$\mathbf{B} = \mathbf{A} - \mathbf{A}_{H_0} = \langle \mathbf{A} \rangle - \mathbf{A}_{H_0} + \mathbf{X}. \quad (4.1.16)$$

Dans le cadre de cet ouvrage, on se limite à un seul choix de modèle nul, soit le modèle de configuration (CM). Ce modèle (cf. Sec. 1.2.3) prédit que l'élément  $a_{ij}^{(\text{CM})}$  de la matrice nulle est seulement une fonction du degré des noeuds  $i, j$

$$a_{ij}^{(\text{CM})} = \frac{k_i k_j}{2M}, \quad (4.1.17)$$

où  $m$  est le nombre total de liens. Puisque les méthodes de la TMA utilisées ici sont uniquement valides dans la limite  $N \rightarrow \infty$ , on peut remplacer le modèle nul exact par un modèle nul moyen

$$\mathbf{A}_{\text{CM}} \rightarrow \langle \mathbf{A}_{\text{CM}} \rangle \quad (4.1.18)$$

qui tient dans la limite asymptotique, tel que

$$a_{ij}^{(\text{CM})} \rightarrow \frac{\langle k_i \rangle \langle k_j \rangle}{2M}. \quad (4.1.19)$$

Pour le cas d'un réseau modélisé par le SBM, la matrice du modèle nul prend la forme par bloc

$$\mathbf{A}_{\text{CM}} = \frac{1}{2M} \begin{bmatrix} \langle k_1 \rangle \langle k_1 \rangle \mathbf{1}_{s_1} \mathbf{1}_{s_1}^T & \langle k_1 \rangle \langle k_2 \rangle \mathbf{1}_{s_1} \mathbf{1}_{s_2}^T & \dots & \langle k_1 \rangle \langle k_g \rangle \mathbf{1}_{s_1} \mathbf{1}_{s_g}^T \\ \langle k_2 \rangle \langle k_1 \rangle \mathbf{1}_{s_2} \mathbf{1}_{s_1}^T & \langle k_2 \rangle \langle k_2 \rangle \mathbf{1}_{s_2} \mathbf{1}_{s_2}^T & \dots & \langle k_2 \rangle \langle k_g \rangle \mathbf{1}_{s_2} \mathbf{1}_{s_g}^T \\ \vdots & \vdots & \ddots & \vdots \\ \langle k_g \rangle \langle k_1 \rangle \mathbf{1}_{s_g} \mathbf{1}_{s_1}^T & \langle k_g \rangle \langle k_2 \rangle \mathbf{1}_{s_g} \mathbf{1}_{s_2}^T & \dots & \langle k_g \rangle \langle k_g \rangle \mathbf{1}_{s_g} \mathbf{1}_{s_g}^T \end{bmatrix}, \quad (4.1.20)$$

tandis que la matrice de modularité complète prend la forme par bloc perturbée

$$\begin{bmatrix} (p_{11} + \frac{1}{2M} \langle k_1 \rangle \langle k_1 \rangle) \mathbf{1}_{s_1} \mathbf{1}_{s_1}^T & (p_{12} + \frac{1}{2M} \langle k_1 \rangle \langle k_2 \rangle) \mathbf{1}_{s_1} \mathbf{1}_{s_2}^T & \dots & (p_{1g} + \frac{1}{2M} \langle k_1 \rangle \langle k_g \rangle) \mathbf{1}_{s_1} \mathbf{1}_{s_g}^T \\ (p_{21} + \frac{1}{2M} \langle k_2 \rangle \langle k_1 \rangle) \mathbf{1}_{s_2} \mathbf{1}_{s_1}^T & (p_{22} + \frac{1}{2M} \langle k_2 \rangle \langle k_2 \rangle) \mathbf{1}_{s_2} \mathbf{1}_{s_2}^T & \dots & (p_{2g} + \frac{1}{2M} \langle k_2 \rangle \langle k_g \rangle) \mathbf{1}_{s_2} \mathbf{1}_{s_g}^T \\ \vdots & \vdots & \ddots & \vdots \\ (p_{g1} + \frac{1}{2M} \langle k_g \rangle \langle k_1 \rangle) \mathbf{1}_{s_g} \mathbf{1}_{s_1}^T & (p_{g2} + \frac{1}{2M} \langle k_g \rangle \langle k_2 \rangle) \mathbf{1}_{s_g} \mathbf{1}_{s_2}^T & \dots & (p_{gg} + \frac{1}{2M} \langle k_g \rangle \langle k_g \rangle) \mathbf{1}_{s_g} \mathbf{1}_{s_g}^T \end{bmatrix} + \mathbf{X}. \quad (4.1.21)$$

### Cas spécial : 2 blocs identiques

Bien que l'argument qui est présenté dans les sections 4.2-4.3 soit applicable à un SBM général, on se concentrera ici sur un cas simple afin de bien illustrer les résultats importants pouvant être obtenus à l'aide de la TMA (sans s'encombrer de détails inutiles). Plus spécifiquement, on étudiera le cas d'un SBM de deux communautés de taille égale, connectés par deux probabilités différentes : une pour les liens intra-communautaires, et une pour les liens extra-communautaires, i.e.

$$g = 2 \quad s_1 = s_2 = N/2 \quad p_{11} = p_{22} := p \quad p_{12} = p_{21} = q. \quad (4.1.22)$$

De l'Eq. 4.1.2, on tire la matrice moyenne  $\langle \mathbf{A} \rangle$  de ce modèle

$$\langle \mathbf{A} \rangle = \begin{bmatrix} p & p & \dots & q & q \\ p & p & \dots & q & q \\ \vdots & \vdots & \ddots & \vdots & \vdots \\ q & q & \dots & p & p \\ q & q & \dots & p & p \end{bmatrix}, \quad (4.1.23)$$

qui peut être formulée de façon compacte comme

$$\langle \mathbf{A} \rangle = \frac{1}{2}(p+q)\mathbf{1}_N \mathbf{1}_N^T + \frac{1}{2}(p-q)\mathbf{u}_N \mathbf{u}_N^T \quad (4.1.24)$$

à l'aide de la matrice  $\mathbf{1}_N$  habituelle et d'une nouvelle matrice  $N \times 1$

$$\mathbf{u}_N^T := \underbrace{[1 \dots 1]}_{\times N/2} \underbrace{[-1 \dots -1]}_{\times N/2}. \quad (4.1.25)$$

Les équations (4.1.13) prédisent des degrés moyens identiques pour tous les noeuds

$$\langle k \rangle_1 = s_1 p_{11} + s_2 p_{12} = \frac{N}{2}(p+q) \quad (4.1.26a)$$

$$\langle k \rangle_2 = s_1 p_{21} + s_2 p_{22} = \frac{N}{2}(p+q), \quad (4.1.26b)$$

de sorte que le modèle nul de ce SBM simple correspond en fait à un réseau ER

$$\mathbf{A}_{\text{CM}} = \frac{(p+q)}{2} \mathbf{1}_N \mathbf{1}_N^T, \quad (4.1.27)$$

car  $\mathbf{A}_{H_0}$  est un multiple d'une matrice uniforme. Finalement, l'équation (4.1.16) qui relie  $\mathbf{B}$  à  $\mathbf{A}$  et la forme compacte (4.1.24) de  $\mathbf{A}$  nous permettent d'écrire  $\mathbf{B}$  comme

$$\mathbf{B} = \frac{1}{2}(p+q)\mathbf{1}_N\mathbf{1}_N^T + \frac{1}{2}(p-q)\mathbf{u}_N\mathbf{u}_N^T - \frac{1}{2}(p+q)\mathbf{1}_N\mathbf{1}_N^T + \mathbf{X} = \frac{1}{2}(p-q)\mathbf{u}_N\mathbf{u}_N^T + \mathbf{X}. \quad (4.1.28)$$

Cette équation est particulièrement importante car elle relie les matrices de modularité  $\mathbf{B}$  et de déviation  $\mathbf{X}$  d'une façon qui nous permet également de relier leurs *spectres*. Sans rentrer dans le détail immédiatement (on relègue cet exercice à la Sec. 4.3), on mentionnera simplement que cette équation nous permettra d'utiliser le spectre de  $\mathbf{X}$  – calculable à l'aide la TMA – pour borner la distribution des valeurs propres  $\{z_j\}_{j=1,\dots,N}$  de  $\mathbf{B}$ . L'étape suivante consiste donc à calculer la densité spectral  $\rho(\lambda)$  de  $\mathbf{X}$ .

## 4.2 Densité spectrale de la matrice de déviation

Le calcul de la densité spectrale de la matrice de déviation est un exercice classique en TMA [37, p.33 & p.289], car  $\mathbf{X}$  est en fait une variante de la matrice de Wigner [8]. Dans cette section, on effectue le calcul à l'aide d'une approche non orthodoxe, basée purement sur des arguments combinatoires et topologiques. Il s'agit d'une version explicite des arguments partiels présentés dans la Ref. [68].

Pour arriver à nos fins, on se basera sur l'identité de Stieltjes-Perron (démontrée dans l'annexe B.5), qui relie la densité spectrale  $\rho(\lambda)$  d'une matrice symétrique réelle à la partie imaginaire de la trace de sa matrice résolvante  $R_{\mathbf{X}}(z)$  par

$$\rho(\lambda) = \frac{1}{\pi} \lim_{\varepsilon \rightarrow 0^+} \text{Im} \left\{ \left\langle \text{Tr}(R_{\mathbf{X}}(\lambda + i\varepsilon)) \right\rangle \right\}, \quad (4.2.1)$$

où la matrice résolvante est définie comme

$$R_{\mathbf{X}}(z) := (\mathbf{X} - z\mathbf{I})^{-1} = \sum_{j=1}^N \frac{\boldsymbol{\phi}_j \boldsymbol{\phi}_j^T}{(\lambda_j - z)} \quad z \in \mathbb{C} \setminus \mathbb{R}, \quad (4.2.2)$$

avec  $\{z_j\}_{j=1,\dots,N}$  et  $\{\boldsymbol{\phi}_j\}_{j=1,\dots,N}$ , les valeurs propres et vecteurs propres de  $\mathbf{X}$ . Dans cette section, le calcul du CD de l'équation (4.2.1) sera effectué en deux étapes. Dans un premier temps (Sec. 4.2.1), on utilisera des arguments issus de la théorie des graphes pour obtenir une expression asymptotique au premier ordre pour la moyenne de la trace apparaissant dans l'Eq. (4.2.1). Des techniques standards d'analyse pourront ensuite être appliquées (Sec. 4.2.2) pour obtenir  $\rho(\lambda)$ .

### 4.2.1 Trace de la matrice résolvante

La série de Neumann de la matrice résolvante  $R_{\mathbf{X}}$  existe et converge car  $\mathbf{X}$  est un opérateur borné<sup>2</sup>. On peut donc écrire  $R_{\mathbf{X}}$  comme

$$R_{\mathbf{X}}(\lambda + i\varepsilon) = (\mathbf{X} - (\lambda + i\varepsilon)\mathbf{I})^{-1} = - \sum_{k=0}^{\infty} \frac{\mathbf{X}^k}{(\lambda + i\varepsilon)^{k+1}}, \quad (4.2.3)$$

afin d'exprimer la moyenne de la trace de  $R_{\mathbf{X}}$  sous la forme d'une série de puissance

$$\langle \text{Tr}(R_{\mathbf{X}}(\lambda + i\varepsilon)) \rangle = - \sum_{k=0}^{\infty} \frac{\langle \text{Tr}(\mathbf{X}^k) \rangle}{(\lambda + i\varepsilon)^{k+1}}. \quad (4.2.4)$$

La moyenne de la trace d'une matrice arbitraire à la  $k^{\text{ème}}$  puissance est de la forme

$$\langle \text{Tr}(\mathbf{X}^k) \rangle = \sum_{\mathcal{N}} \langle x_{i_1 i_2} x_{i_2 i_3} \dots x_{i_k i_1} \rangle \quad (4.2.5)$$

où  $\mathcal{N}$  est l'ensemble des combinaisons d'indices possibles  $(i_1, i_2, \dots, i_k)$ , où  $i_n$  peut prendre les valeurs  $i_n = 1, 2, \dots, N \forall n$ . Les indices répétés sont permis, ce qui implique qu'on aura toujours  $v \leq k$  indices distincts pour une matrice à la  $k^{\text{ème}}$  puissance.

Cet intervalle de sommation n'étant pas intuitif, on utilise une représentation alternative permettant d'énumérer efficacement les éléments de  $\mathcal{N}$ , et donc d'obtenir les termes de l'Eq. (4.2.5) ainsi que leurs poids respectifs (voir Ref. [19, 37, 119] pour des méthodes similaires). On utilise des graphes dirigés dont les noeuds correspondent à des indices *distincts*  $i_a, i_b, \dots, i_y, i_v$  et dont les arêtes  $i_a \rightarrow i_b, \dots, i_y \rightarrow i_v$  représentent les éléments  $x_{i_a i_b}, \dots, x_{i_y i_v}$  de la somme (4.2.5) (cf. Fig. 4.2).

L'organisation cyclique des indices de  $\mathcal{N}$  impose une contrainte importante sur ces graphes : les graphes permis doivent toujours contenir un cycle d'Euler partant de chaque noeud.

Une procédure de construction simple permet d'obtenir tous les graphes qui respectent cette contrainte (illustrée à la figure 4.2).

1. Pour une puissance  $k$  donnée, on commence l'énumération en traçant un graphe cyclique dirigé de  $k$  noeuds et  $k$  arêtes.
2. Des termes contenant des indices répétés (et donc un nombre moins élevé  $v$  d'indices libres) sont représentés en fusionnant les noeuds existants. Les arêtes sont conservées et transformées en boucles.
3. On énumère toutes les séquences possibles de fusion de noeuds qui mènent à des graphes topologiquement distincts.

Le graphe initial contient évidemment les parcours d'Euler imposés par la contrainte. Puisque les arêtes sont préservées lors de la fusion des noeuds, les parcours déjà existant ne sont pas détruits. En

2. Cette condition est nécessaire, car on peut arriver à la série de Neumann en réarrangeant la définition de la résolvante sous la forme  $\mathbf{R} = -\frac{1}{\lambda}\mathbf{I} + \frac{1}{\lambda}\mathbf{X}\mathbf{R} = \mathbf{I}$ , puis en substituant  $\mathbf{R}$  à l'infini. Cette procédure converge uniquement pour  $\mathbf{X}$  borné.

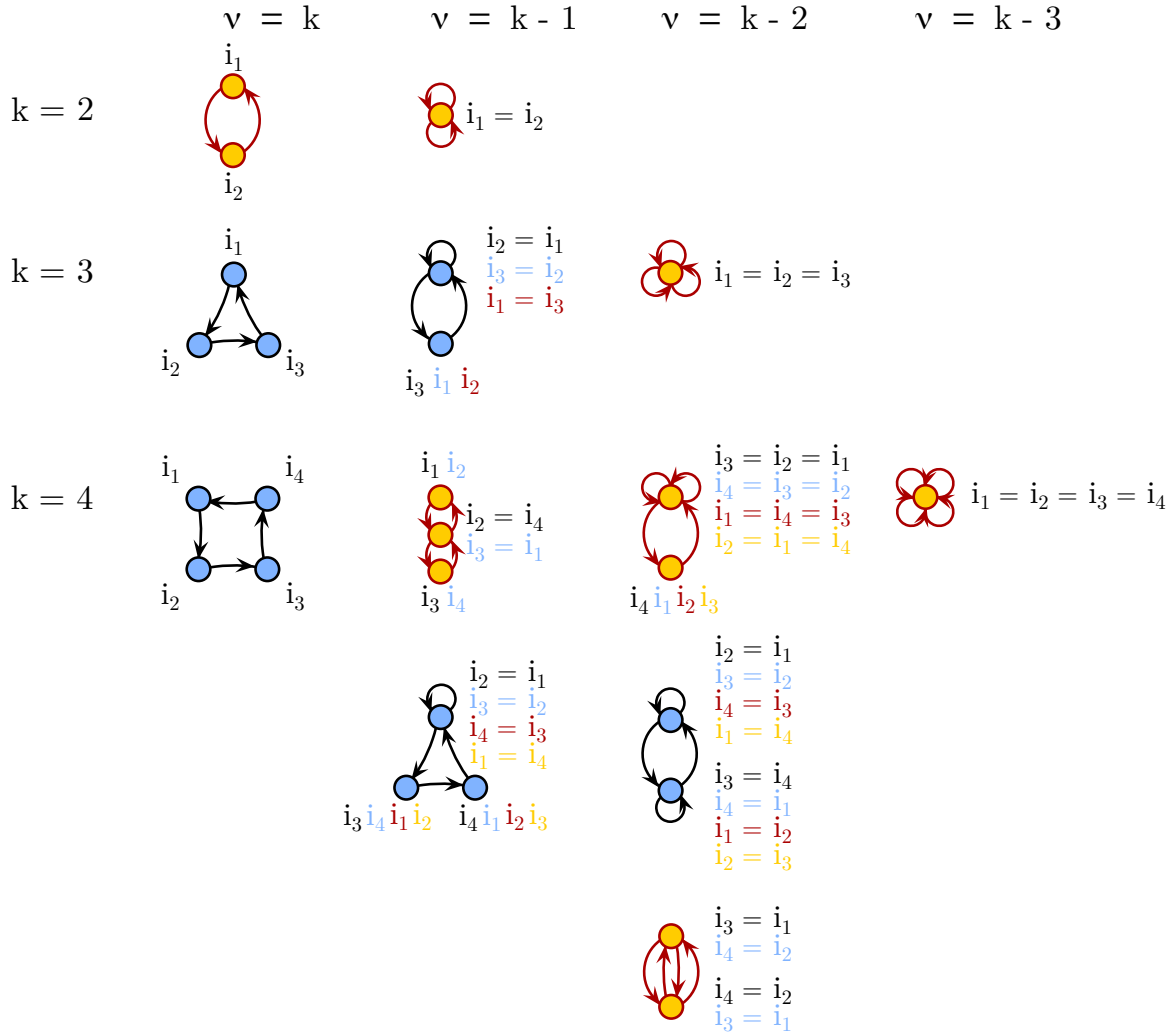


FIGURE 4.2 – **Procédure de construction des graphes combinatoires.** Les graphes initiaux sont placés à gauche. Toutes les combinaisons topologiquement distinctes sont considérées, pour l'intervalle  $v = 1, \dots, k$ . Les graphes illustrés ici représentent les termes de la forme (pour  $\mathbf{X}$  symétrique)

( $k = 2$  arêtes) :  $\langle x_{i_1 i_2}^2 \rangle, \langle x_{i_1 i_1}^2 \rangle$

( $k = 3$  arêtes) :  $\langle x_{i_1 i_2} \rangle \cdot \langle x_{i_2 i_3} \rangle \cdot \langle x_{i_3 i_1} \rangle; \langle x_{i_1 i_1} \rangle \cdot \langle x_{i_1 i_3}^2 \rangle; [\text{dégénéré 3 fois}] \langle x_{i_1 i_1}^3 \rangle;$

( $k = 4$  arêtes) :  $\langle x_{i_1 i_2} \rangle \cdot \langle x_{i_2 i_3} \rangle \cdot \langle x_{i_3 i_4} \rangle \cdot \langle x_{i_4 i_1} \rangle; \langle x_{i_1 i_2}^2 \rangle \cdot \langle x_{i_2 i_3}^2 \rangle; [\text{dégénéré 2 fois}] \langle x_{i_1 i_1} \rangle \cdot \langle x_{i_1 i_3} \rangle \cdot \langle x_{i_3 i_4} \rangle \cdot \langle x_{i_4 i_1} \rangle; [\text{dégénéré 4 fois}] \langle x_{i_1 i_1}^2 \rangle \cdot \langle x_{i_1 i_4}^2 \rangle; [\text{dégénéré 4 fois}] \langle x_{i_1 i_1} \rangle \cdot \langle x_{i_1 i_3}^2 \rangle \cdot \langle x_{i_3 i_3}^2 \rangle; [\text{dégénéré 4 fois}] \langle x_{i_1 i_2}^4 \rangle; [\text{dégénéré 2 fois}] \langle x_{i_1 i_1}^4 \rangle.$  Les termes nuls sont colorés en bleu.

suivant la procédure susmentionnée, on obtient donc une représentation de toutes les combinaisons d'indices distincts possibles.

Cette procédure d'énumération est laborieuse . Cependant, dans le cas à l'étude, une grande partie des graphes peut être éliminée. En effet, les éléments de  $\mathbf{X}$  étant indépendants et de moyenne nulle, une arête simple se traduirait par un élément  $\langle x_{ijik} \rangle$  isolé dans l'équation (4.2.5) et donc par un terme nul. On obtient une représentation des termes non nuls en ne conservant que les graphes dont les arêtes et boucles sont *au moins* doublés.

On veut maintenant calculer le *poids* de chacun de ces graphes, i.e. leur contribution à la somme (4.2.5). Cette contribution peut être séparée en trois parties :

1. un facteur d'échelle  $N^v$  comptant le nombre de substitutions possibles pour un arrangement de  $v$  indices donnés ;
2. un facteur multiplicatif  $M_k(v; \mathbf{X})$  dû aux moments  $\langle x_{jk}^m \rangle$  des éléments de la matrice de déviation ;
3. un facteur combinatoire  $a_k(v)$  comptant le nombre de graphes de  $v$  noeuds topologiquement distincts ayant le même poids.

Le poids d'un graphe donné est donc égale à

$$N^v M_k(v; \mathbf{X}) a_k(v). \quad (4.2.6)$$

À titre d'exemple, on effectue le calcul explicite pour le premier graphe  $k = 4, v = k - 2 = 2$  de la figure 4.2. Le premier facteur est égale à  $N^2$  car les deux noeuds / indices  $(i_1 = i_2 = i_3), (i_4)$  peuvent prendre  $N$  valeurs différentes. Le deuxième facteur est égale à la *valeur moyenne* de  $\langle x_{i_1 i_1}^2 \rangle \langle x_{i_1 i_2}^2 \rangle$ , qu'on notera ici  $M_k(\mu; \mathbf{X}) = \langle x_{i_1 i_1}^2 \rangle \langle x_{i_1 i_2}^2 \rangle = m_2^{(d)} m_2^{(o)}$  (où  $(d)$  réfère à un élément sur la diagonale et  $(o)$  aux autres éléments). Le dernier facteur combinatoire  $a_k(v)$  est difficile à obtenir dans le cas général, mais est égale à 4 dans ce cas particulier, comme le montre l'énumération exhaustive effectué à la figure 4.2. Le graphe représente donc une contribution  $\sim 4N^2 m_2^{(d)} m_2^{(o)}$  à la trace (4.2.5).

Seul le calcul de  $a_k(v)$  pose problème en général. On simplifie donc cette étape en s'intéressant seulement à une expression asymptotique au premier ordre en  $N$ . Pour un  $k$  donné, le terme dominant de l'équation (4.2.5) est celui pour lequel le facteur  $N^v$  est le plus grand, i.e. celui correspondant au graphe contenant un nombre maximal de noeuds  $v$ .

Dans le cas où  $k$  est pair, *sous la contrainte que toutes les arêtes doivent être au moins doubles*, ce graphe maximal est celui formé par un arbre d'arêtes doubles, sans boucles, de  $v = k/2 + 1$  noeuds (voir Fig. 4.3-4.4). L'avantage de ces graphes doublés est qu'ils sont équivalents à des *arbres planaires enracinés* (cf. Fig. 4.5). Or, il est connu que le nombre d'arbres planaires enracinés non isomorphes de  $k/2 + 1$  noeuds est donné par le nombre de Catalan  $C_{k/2}$  [101, E. 6.19 (e)]. Puisque notre procédure de construction énumère justement ces arbres non isomorphes, on a directement que

$$a_k(k/2 + 1) = C_{k/2}. \quad (4.2.7)$$



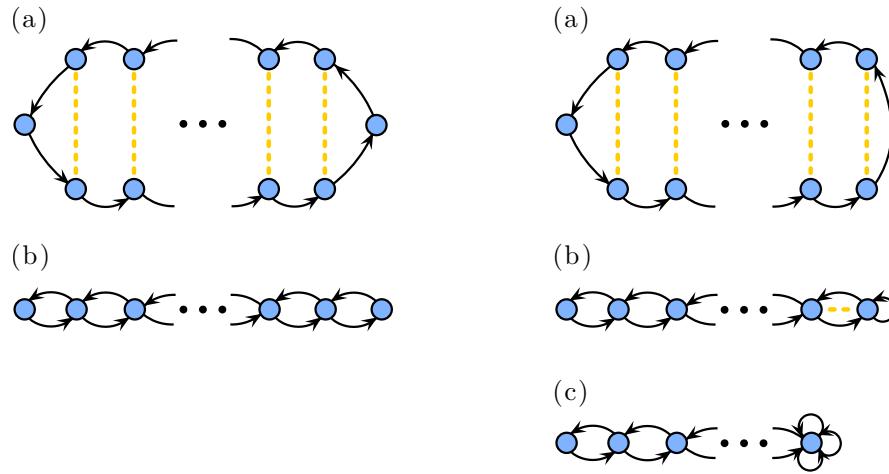


FIGURE 4.3 – **Graphes maximaux pour un nombre d’arêtes pair et impair.** (gauche) Pour un nombre initial de noeuds pair, on obtient un graphe maximal dont toutes les arêtes sont doublées en plaçant les noeuds selon la configuration (a), et en fusionnant les noeuds par paires (lignes oranges pointillées). Le graphe résultant contient  $v = k/2 + 1$  noeuds. (droite) Pour un nombre initial de noeuds impair, cette procédure laisse une boucle orpheline dans la chaîne d’arêtes doubles. On doit donc fusionner le noeud contenant une boucle orpheline à nouveau, ce qui introduit une boucle triple. Le graphe résultant contient  $[k/2]$  noeuds, où  $[n]$  dénote la partie entière de  $n$ .

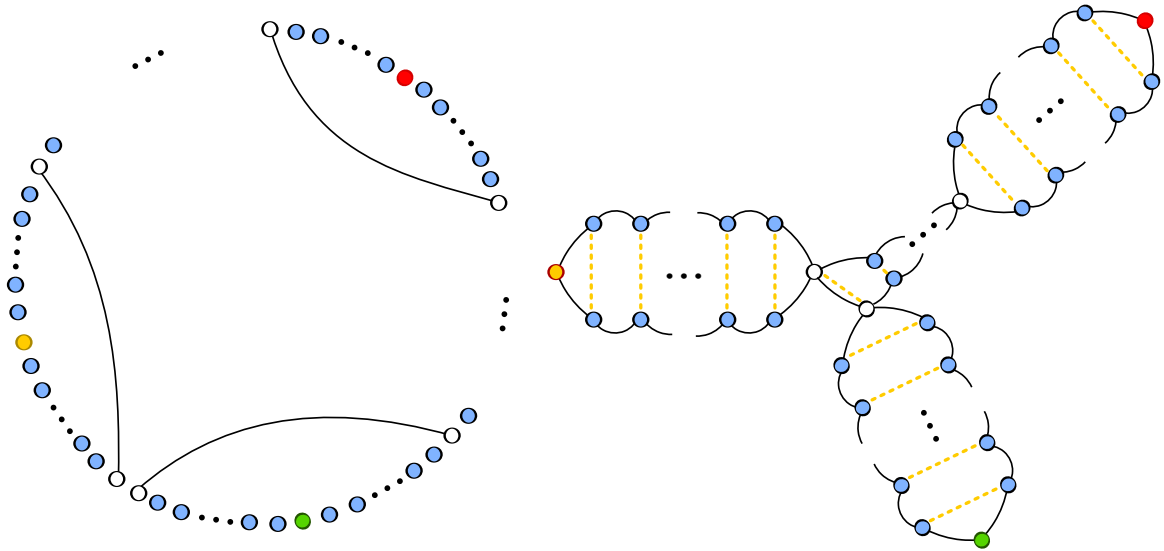


FIGURE 4.4 – **Graphe maximal : cas général.** La procédure de la figure 4.3 permet de comprendre intuitivement comment obtenir *une* des classes de graphe maximal. Mais en fait, pour un nombre initial de noeuds pair, tous les graphes avec  $\mu = 2, \dots, k/2$  pointes sont accessibles. Le nombre de noeuds est toujours égal à  $k/2 + 1$  dans ces cas, car chaque pointe additionnelle est contrebalancée par une fusion supplémentaire.

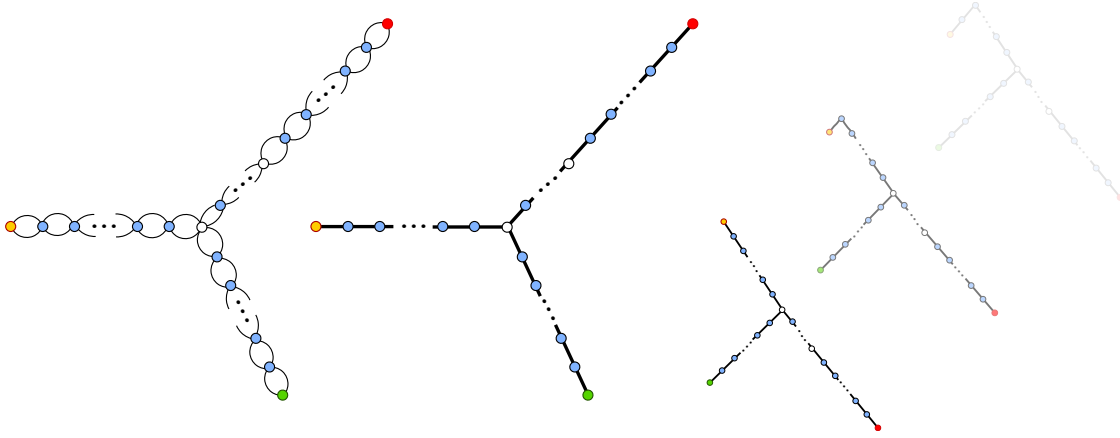


FIGURE 4.5 – **Graphes doublés en tant qu'arbres planaires enracinés.** Les arbres planaires sont obtenus en remplaçant les arêtes doublés (gauche) des graphes obtenus *via* la procédure de construction (Figs. 4.2-4.4) par des arêtes uniques (centre). Tous les noeuds peuvent être sélectionnés comme racine de l'arbre planaire (droite). En les choisissant un à un, on obtient l'ensemble des arbres planaires *enracinés* à  $\mu$  pointes. L'ensemble complet des arbres planaires enracinés non isomorphe de  $\nu$  noeuds est couvert en considérant tous les noeuds des arbres à  $\mu = 2, \dots, k/2$  pointes.

Toujours pour  $k$  pair, on a que le facteur multiplicatif  $M_k(k/2 + 1; \mathbf{X})$  apparaissant dans le calcul du poids (4.2.6) ne dépend que de la moyenne  $m_2$  des deuxièmes moments  $\langle x_{jk}^2 \rangle$  des éléments de la matrice de déviation, i.e.

$$m_2 := \frac{1}{N^2} \sum_{jk} \langle x_{jk}^2 \rangle \quad (4.2.8)$$

En effet, dans ce cas, les graphes maximaux sont formés de  $k/2$  arêtes doubles et représentent donc des termes où la quantité  $m_2$  apparaît  $k/2$  fois, tel que

$$M_k(k/2 + 1; \mathbf{X}) = m_2^{k/2}. \quad (4.2.9)$$

Dans le cas d'un nombre impair d'arêtes, les graphes maximaux contiennent  $[k/2]$  noeuds, où les parenthèses carrées dénotent la partie entière d'une quantité. Ces graphes prennent aussi la forme d'arbres enracinés non isomorphes, mais un facteur combinatoire additionnel  $f([k/2])$  doit être introduit pour tenir compte de la boucle triple qui doit être placée sur une des pointes de l'arbre (cf. Fig. 4.3). Sans effectuer le calcul explicite, on sait au moins que la forme de la contribution combinatoire complète est

$$a_k([k/2]) = f([k/2]) C_{[k/2]-1}, \quad (4.2.10)$$

tandis que la contribution des moments est

$$M_k \left( \left[ \frac{k}{2} \right]; \mathbf{X} \right) = m_2^{[k/2]-1} m_3, \quad (4.2.11)$$

où  $m_3 := (\sum \langle x_{jk}^3 \rangle) / N^2$  est la moyenne des troisièmes moments, et où la puissance 3 est due au fait que trois boucles apparaissent.

En regroupant tous les résultats obtenus jusqu'à présent, on a au premier ordre,

$$\langle \text{Tr}(\mathbf{X}^k) \rangle \sim \begin{cases} C_{k/2} \cdot N^{k/2+1} \cdot m_2^{k/2} & k \text{ pairs} \\ f([k/2]) C_{[k/2]-1} \cdot N^{[k/2]} \cdot m_2^{[k/2]-1} m_3^3 & k \text{ impairs} \end{cases}, \quad (4.2.12)$$

tel que les termes  $k$  impairs sont négligeables.

#### 4.2.2 Spectre de la matrice de déviation

On se rappellera que la densité spectrale  $\rho(\lambda)$  de la matrice  $\mathbf{X}$  est donnée par (Eq. 4.2.1)

$$\rho(\lambda) = \frac{1}{\pi} \lim_{\varepsilon \rightarrow 0^+} \text{Im} \left\{ \left\langle \text{Tr}(R_{\mathbf{X}}(\lambda + i\varepsilon)) \right\rangle \right\} \quad (4.2.13)$$

où l'expression  $\langle \text{Tr}(R_{\mathbf{X}}(\lambda + i\varepsilon)) \rangle$  peut être exprimée comme un développement en série en puissances sur les moments la trace de la matrice  $\mathbf{X}$ . Notre expression approximée pour  $\langle \text{Tr}(\mathbf{X}^k) \rangle$ , nous permet de sommer cette série (CD de l'Eq. 4.2.4). En effet, en ne conservant que les termes pairs, on obtient :

$$\langle \text{Tr}(R_{\mathbf{X}}(\lambda + i\varepsilon)) \rangle \approx - \sum_{k=0}^{\infty} \frac{\langle \text{Tr}(\mathbf{X}^{2k}) \rangle}{(\lambda + i\varepsilon)^{2k+1}} = - \sum_{k=0}^{\infty} \frac{N^{k+1} C_k m_2^k}{(\lambda + i\varepsilon)^{2k+1}} \quad (4.2.14)$$

Cette somme peut être vue comme une série de puissance générant les nombres de Catalan. Or, la fonction génératrice des nombres de Catalan est facilement accessible. En effet, partant de la relation de récurrence

$$C_{k+1} = \sum_{j=0}^k C_j C_{k-j}, \quad (4.2.15)$$

qui définit les nombres de Catalan, on multiplie chaque côté par  $x^k$  et on somme sur les  $k$  :

$$\sum_{k=0}^{\infty} C_{k+1} x^k = \sum_{k=0}^{\infty} \sum_{j=0}^k C_j C_{k-j} x^k. \quad (4.2.16)$$

Le CD est un produit de Cauchy de deux séries de puissances en  $x$ , d'où

$$\sum_{k=0}^{\infty} C_{k+1} x^k = \left( \sum_{k=0}^{\infty} C_k x^k \right) \left( \sum_{j=0}^{\infty} C_j x^j \right) := [C(x)]^2, \quad (4.2.17)$$

avec  $C(x) := \sum_k C_k x^k$ . Le CG peut aussi être exprimé en termes de  $C(x)$

$$\sum_{k=0}^{\infty} C_{k+1} x^k = \sum_{k=1}^{\infty} C_k x^{k-1} = \frac{1}{x} \sum_{k=1}^{\infty} C_k x^k = \frac{1}{x} (C(x) - C_0) \quad (4.2.18)$$

et puisque  $C_0 \equiv 1$ ,  $C(x)$  est donnée par les solutions de

$$xC^2(x) - C(x) + 1 = 0 \quad (4.2.19)$$

i.e  $C(x) = (1 \pm \sqrt{1-4x})/2x$ . On sélectionne directement la racine négative car la racine positive est de la forme  $1/0$  à  $x = 0$ .

Utilisant ce résultat, l'équation (4.2.14) devient

$$\left\langle \text{Tr}(R_{\mathbf{X}}(\lambda + i\varepsilon)) \right\rangle = -\frac{N}{(\lambda + i\varepsilon)} C(Nm_2/(\lambda + i\varepsilon)^2) = \frac{-(\lambda + i\varepsilon) + \sqrt{(\lambda + i\varepsilon)^2 - 4Nm_2}}{2m_2} \quad (4.2.20)$$

Avant de prendre la partie imaginaire, on exprime l'argument de la racine sous sa forme polaire

$$re^{i\varphi + i2n\pi} = r[\cos(\varphi) + i\sin(\varphi)] = (\lambda + i\varepsilon)^2 - 4Nm_2 \quad (4.2.21a)$$

$$r^2 = \lambda^4 - 8\lambda^2Nm_2 - 2\lambda^2\varepsilon^2 + 16N^2m_2^2 + 8Nm_2\varepsilon^2 + \varepsilon^4 + 4\varepsilon^2\lambda^2 \quad (4.2.21b)$$

$$\varphi = \tan^{-1} \left( \frac{2\varepsilon\lambda}{\lambda^2 - 4Nm_2 - \varepsilon^2} \right) + \pi, \quad (4.2.21c)$$

où la phase de  $\pi$  est introduite pour sélectionner la bonne branche dans le cas  $\lambda^2 < 4Nm_2 - \varepsilon^2$ . La partie imaginaire de la trace est ainsi

$$\text{Im} \left\{ \left\langle \text{Tr}(R_{\mathbf{X}}(\lambda + i\varepsilon)) \right\rangle \right\} = \frac{\varepsilon}{2m_2} r^{1/2} \sin(\varphi/2 + n\pi). \quad (4.2.22)$$

Dans la limite  $\varepsilon \rightarrow 0^+$ , on a

$$\lim_{\varepsilon \rightarrow 0^+} r^2 = \lambda^4 - 8\lambda^2Nm_2 + 16N^2m_2^2 = (\lambda^2 - 4Nm_2)^2 \implies r = \pm(4Nm_2 - \lambda^2) \quad (4.2.23a)$$

$$\lim_{\varepsilon \rightarrow 0^+} \varphi = \pi, \quad (4.2.23b)$$

où la racine positive/négative de  $r$  concorde avec  $n$  pair/impair. Le calcul complet du CD de l'Eq. (4.2.1) mène donc à

$$\rho(\lambda) = \lim_{\varepsilon \rightarrow 0^+} \text{Im} \left\{ \frac{1}{\pi} \left\langle \text{Tr}(R_{\mathbf{X}}(\lambda + i\varepsilon)) \right\rangle \right\} = \pm \frac{\sqrt{4Nm_2 - \lambda^2}}{2m_2\pi} \sin(\pi/2 + n\pi) = \frac{\sqrt{4Nm_2 - \lambda^2}}{2m_2\pi}, \quad (4.2.24)$$

la densité spectrale de la matrice  $\mathbf{X}$ , entièrement contenue dans l'intervalle  $\lambda \in [-2\sqrt{Nm_2}, 2\sqrt{Nm_2}]$ , à cause du choix de phase pour  $\varphi$ . Il s'agit d'une forme modifiée de la loi du demi-cercle de Wigner.

### Application au SBM à deux blocs

Les démarches effectuées jusqu'à présent sont très générales, car le spectre de la matrice de déviation ne dépend que très faiblement du modèle de réseau utilisé : en fait, le modèle intervient uniquement dans le calcul de  $m_2$ .

Pour un modèle par blocs parfaitement général, les moments de  $\mathbf{X}$  sont donnés par

$$\langle x_{jk} \rangle = \langle a_{jk} \rangle - \langle a_{jk} \rangle = 0 \quad (4.2.25a)$$

$$\langle x_{jk}^2 \rangle = \langle a_{jk}^2 - 2a_{jk} \langle a_{jk} \rangle + \langle a_{jk} \rangle^2 \rangle = \langle a_{jk}^2 \rangle - \langle a_{jk} \rangle^2 = p_{\sigma_j \sigma_k} (1 - p_{\sigma_j \sigma_k}) \quad (4.2.25b)$$

...

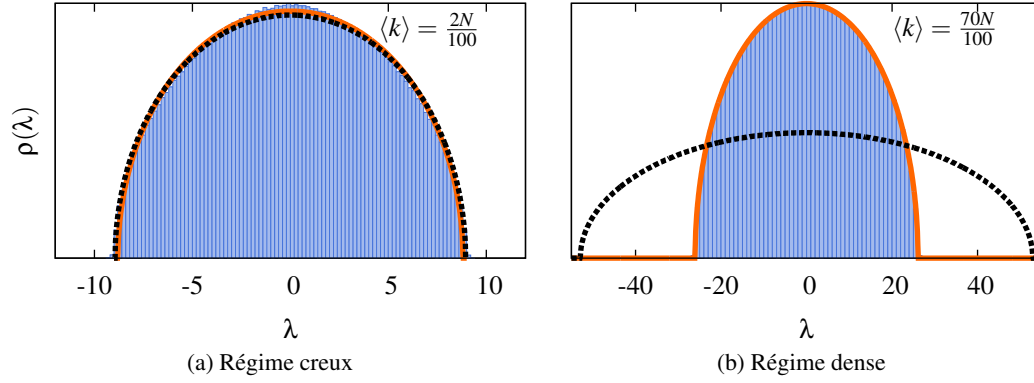


FIGURE 4.6 – **Densité spectrale de la matrice de déviation pour le SBM à deux blocs.** Comparaison du spectre empirique (histogramme bleu), de la densité analytique générale avec  $m_2 = q(1 - q) + p(1 - p)/2$  (orange) et de la prédiction de Nadakuditi et Newman avec  $m_2 = N(p + q)/2$  (pointillé noir) dans le régime creux (a) et dense (b) du SBM à deux blocs. Les résultats sont obtenus en calculant le spectre moyen d'un réseau de  $N = 1000$  noeuds sur 500 réalisations du SBM. Dans le régime creux, la densité analytique et empirique diffèrent autour de  $\lambda \sim 0$ , autant pour le formalisme binomial que le formalisme de Poisson. L'erreur est due à l'omission des termes d'ordres supérieurs dans le développement de la trace (4.2.4).

où  $\sigma_i$  est l'application (déjà introduite à la définition 26 du Chap. 3)

$$\sigma_i : i \mapsto \Psi_\alpha \quad (4.2.26)$$

qui associe l'indice  $\alpha$  d'une et une seule communauté  $\Psi_\alpha$  à chaque noeud  $i$  du réseau. La *moyenne* des deuxièmes moments  $m_2$  est donc égale à

$$m_2 = \frac{1}{N^2} \sum_{j,k=1}^g p_{\sigma_j \sigma_k} (1 - p_{\sigma_j \sigma_k}) = \sum_{\alpha, \beta=1}^g \frac{s_\alpha s_\beta}{N^2} p_{\alpha\beta} (1 - p_{\alpha\beta}). \quad (4.2.27)$$

Pour le cas spécial de SBM à l'étude,  $m_2$  prend la forme simple

$$m_2 = \frac{\binom{N}{2} \binom{N}{2}}{N^2} [2p(1 - p) + 2q(1 - q)] = \frac{q(1 - q) + p(1 - p)}{2}. \quad (4.2.28)$$

Ce qui mène ultimement une densité spectrale de la forme

$$\rho(\lambda) = \frac{\sqrt{2N[q(1 - q) + p(1 - p)] - \lambda^2}}{\pi[q(1 - q) + p(1 - p)]} \quad (4.2.29)$$

pour la matrice de déviation  $\mathbf{X}$  (Fig. 4.6) On notera que dans le cas particulier à l'étude,  $m_2$  peut aussi être vue comme une fonction de la variance de la distribution en degrés

$$m_2 \equiv \frac{\langle k^2 \rangle - \langle k \rangle^2}{N}. \quad (4.2.30)$$

Si on se place dans la limite où le réseau est creux, la distribution en degrés (binomiale) peut être approximée par une distribution de Poisson

$$p_k \approx \frac{\langle k \rangle^k e^{-\langle k \rangle}}{k} \quad (4.2.31)$$

dont la variance est égale à la moyenne, i.e.

$$m_2 \approx \frac{\langle k \rangle}{N} = \frac{(p+q)}{2}. \quad (4.2.32)$$

C'est précisément la limite dans laquelle Nadakuditi et Newman se placent [68], ce qui a une influence quantitative *et* qualitative sur les résultats à venir (voir 4.6). Une contribution importante de ce chapitre est donc d'introduire la correction (4.2.28) pour le régime dense du SBM à deux blocs.

### 4.3 Spectre de la matrice de modularité

On est maintenant en bonne position pour calculer le spectre de  $\mathbf{B}$ . À nouveau, ce calcul sera effectué en deux étapes. On démontrera d'abord qu'une grande partie du spectre de  $\mathbf{B}$  est bornée par le spectre de  $\mathbf{X}$  (Sec. 4.3.1), et donc que le spectre de  $\mathbf{B}$  suis essentiellement une distribution aléatoire. La position des valeurs propres non bornées sera ensuite calculée explicitement.

#### 4.3.1 Partie bornée du spectre de $\mathbf{B}$

La forme particulière de la matrice de modularité

$$\mathbf{B} = \frac{(p-q)}{2} \mathbf{u}_N \mathbf{u}_N^T + \mathbf{X} \quad (4.3.1)$$

peut être utilisée pour relier son spectre à celui de  $\mathbf{X}$ . En effet, lorsqu'on exprime le problème aux valeurs propres de  $\mathbf{B}$  en termes de cette paramétrisation

$$\mathbf{B} \Psi_j = z_j \Psi_j \quad \rightarrow \quad \left( \mathbf{X} + \frac{(p-q)}{2} \mathbf{u}_N \mathbf{u}_N^T \right) \Psi_j = z_j \Psi_j, \quad (4.3.2)$$

on peut remanier l'équation pour faire apparaître la résolvante de  $\mathbf{X}$

$$\begin{aligned} \frac{(p-q)}{2} \mathbf{u}_N \mathbf{u}_N^T \Psi_j &= (z_j \mathbf{I} - \mathbf{X}) \Psi_j \\ \left[ \mathbf{u}_N^T (z_j \mathbf{I} - \mathbf{X})^{-1} \right] \frac{(p-q)}{2} \mathbf{u}_N \mathbf{u}_N^T \Psi_j &= \left[ \mathbf{u}_N^T (z_j \mathbf{I} - \mathbf{X})^{-1} \right] (z_j \mathbf{I} - \mathbf{X}) \Psi_j \\ \mathbf{u}_N^T (z_j \mathbf{I} - \mathbf{X})^{-1} \mathbf{u}_N &= \frac{2}{(p-q)}. \end{aligned} \quad (4.3.3)$$

On se rappellera que la résolvante  $R_{\mathbf{X}}(z_j)$  s'exprime sur la base orthonormale des vecteurs propres  $\{\Phi_k\}_{k=1,\dots,N}$  et valeurs propres  $\{\lambda_k\}$  de  $\mathbf{X}$  comme (cf. Eq. 4.2.2)

$$(z_j \mathbf{I} - \mathbf{X})^{-1} = -R_{\mathbf{X}}(z_j) = - \sum_{k=1}^N \frac{\Phi_k \Phi_k^T}{(\lambda_k - z_j)} = \sum_{k=1}^N \frac{\Phi_k \Phi_k^T}{(z_j - \lambda_k)}. \quad (4.3.4)$$

Substituant cette forme dans (4.3.3), on obtient une équation reliant les deux spectres  $\{\lambda_k\}$ ,  $\{z_j\}$

$$\begin{aligned} \mathbf{u}_N^T \left( \sum_{k=1}^N \frac{\Phi_k \Phi_k^T}{(z_j - \lambda_k)} \right) \mathbf{u}_N &= \frac{2}{(p-q)} \\ \sum_{k=1}^N \frac{(\mathbf{u}_N^T \Phi_k)^2}{(z_j - \lambda_k)} &= \frac{2}{(p-q)}. \end{aligned} \quad (4.3.5)$$

Cette dernière équation ne permet pas d'accéder *exactement* aux solutions  $\{z_j\}$ , mais elle permet de les borner partiellement. En effet, on note que

1. la pente du CG de l'Eq. (4.3.5) est toujours négative

$$\frac{d}{dz_j} \left( \sum_{k=1}^N \frac{(\mathbf{u}_N^T \boldsymbol{\phi}_k)^2}{(z_j - \lambda_k)} \right) = - \sum_{k=1}^N \frac{(\mathbf{u}_N^T \boldsymbol{\phi}_k)^2}{(z_j - \lambda_k)^2} \quad (4.3.6)$$

2. il existe  $N$  discontinuités aux pôles  $\{\lambda_k\}$

$$\lim_{z_j \rightarrow \lambda_k^\pm} \sum_{k=1}^N \frac{(\mathbf{u}_N^T \boldsymbol{\phi}_k)^2}{(z_j - \lambda_k)} = \pm \infty \quad (4.3.7)$$

3. les limites convergent asymptotiquement vers 0 par le bas ou par le haut

$$\lim_{z_j \rightarrow \pm \infty} \sum_{k=1}^N \frac{(\mathbf{u}_N^T \boldsymbol{\phi}_k)^2}{(z_j - \lambda_k)} = 0^\pm. \quad (4.3.8)$$

Puisque le CD de l'Eq. (4.3.5) est positif (on suppose  $p > q$ ), ces 3 conditions déterminent uniquement les bornes sur l'ensemble  $\{z_j\}$  : pour un arrangement  $z_1 \geq z_2 \geq \dots \geq z_N, \lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_N$  on aura

$$z_1 \geq \lambda_1 \geq z_2 \geq \lambda_2 \geq \dots \geq z_n \geq \lambda_n. \quad (4.3.9)$$

Les valeurs propres de  $\mathbf{B}$  sont donc bornées par les valeurs propres de  $\mathbf{X}$ , sauf la plus grande valeur propre  $z_1$ . Cette intuition est vérifiée à la Fig. 4.7 pour une réalisation donnée de l'ensemble.

Des observations additionnelles peuvent être effectuées si on note que le CG de (4.3.5) dépend uniquement de la matrice aléatoire  $\mathbf{X}$ , alors que le CD tient uniquement compte du modèle nul. Par exemple, si on prend la limite  $p \rightarrow q$ , la droite bleue (CD de l'Eq. 4.3.5) se déplace vers l'infini positif et la courbe orange (CG de l'Eq. 4.3.5) ne change pas : les valeurs propres de  $\mathbf{B}$  (intersections des deux courbes) tendent donc à devenir indistinguables des valeurs propres de  $\mathbf{X}$  (discontinuités du CG de l'Eq. 4.3.5). Ce résultat est attendu, puisque la matrice de modularité est égale à la matrice de déviation  $\mathbf{X}$  lorsque  $p = q$  (Eq. 4.3.1). À l'inverse, lorsque la différence  $(p - q)$  tend vers 1, la valeur propre non bornée se distingue de plus en plus du reste du spectre, car le dernier point de rencontre des deux courbes se déplace vers la droite. Finalement, puisque la densité spectrale de  $\mathbf{X}$  (i.e. densité de discontinuité de la courbe orange) varie en  $\sqrt{N}$  (Eq. 4.2.24), les deux distributions de valeurs propres sont asymptotiquement identiques pour  $N \rightarrow \infty$  (à l'exception de  $z_1$ ).

### 4.3.2 Plus grande valeur propre

Tel qu'indiqué par l'Eq. (4.3.9), la plus grande valeur propre  $z_1$  de la matrice  $\mathbf{B}$  n'est pas bornée par les valeurs propres de  $\mathbf{X}$ . On cherche donc une expression analytique pour cette valeur.

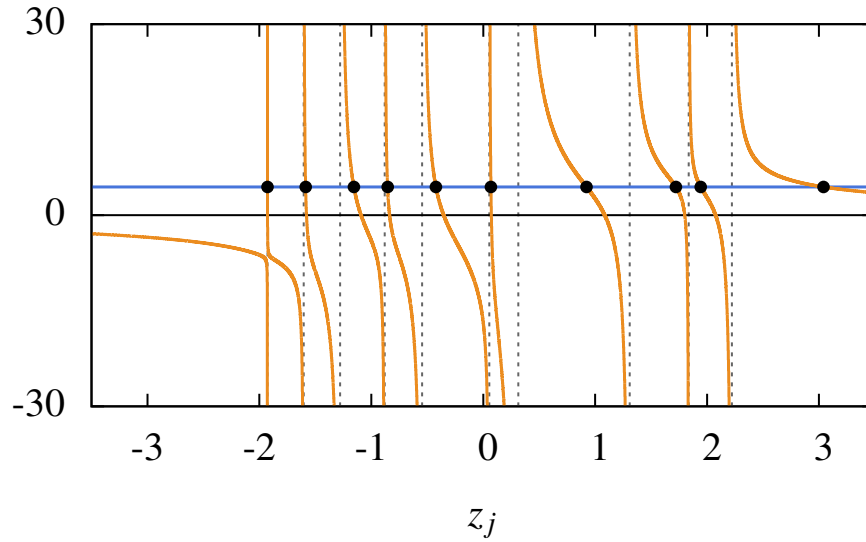


FIGURE 4.7 – **Relation entre le spectre des matrices de déviation et de modularité.** Pour une réalisation donnée du SBM à deux blocs, avec les paramètres  $N = 10$ ,  $p = 0.75$ ,  $q = 0.30$ , on représente (courbe orange) la valeur du CG de l'Eq. 4.3.5 en fonction de  $z_j$ , (droite bleu horizontale) la valeur du CD de l'Eq. 4.3.5 (il s'agit d'une constante). Les points noirs indiquent les valeurs propres  $z_j$  expérimentales de  $\mathbf{B}$ , et apparaissent aux intersections des courbes bleu et orange. Les droites grises indiquent la position des valeurs propres  $\lambda_k$  expérimentales de  $\mathbf{X}$ , et apparaissent aux discontinuités du CG de l'Eq. 4.3.5.

Afin de progresser, on doit d'abord calculer la somme  $\sum_{k=1}^N (\mathbf{u}_N^T \boldsymbol{\phi}_k)^2$  explicitement

$$\begin{aligned} \sum_{k=1}^N (\mathbf{u}_N^T \boldsymbol{\phi}_k)^2 &= \sum_{k=1}^N \left( \sum_{i=1}^{N/2} \phi_i^{(k)} - \sum_{i=N/2+1}^N \phi_i^{(k)} \right)^2 \\ &= \left\{ \sum_{i=1}^{N/2} \sum_{j=1}^{N/2} \sum_{k=1}^N \phi_i^{(k)} \phi_j^{(k)} + \sum_{i=N/2+1}^N \sum_{j=N/2+1}^N \sum_{k=1}^N \phi_i^{(k)} \phi_j^{(k)} - 2 \sum_{i=1}^{N/2} \sum_{j=N/2+1}^N \sum_{k=1}^N \phi_i^{(k)} \phi_j^{(k)} \right\}. \end{aligned} \quad (4.3.10)$$

en termes de la matrice des vecteurs propres  $\boldsymbol{\Phi}$ , on a

$$\sum_{k=1}^N (\mathbf{u}_N^T \boldsymbol{\phi}_k)^2 = \left\{ \sum_{i=1}^{N/2} \sum_{j=1}^{N/2} [\boldsymbol{\Phi} \boldsymbol{\Phi}^T]_{ij} + \sum_{i=N/2+1}^N \sum_{j=N/2+1}^N [\boldsymbol{\Phi} \boldsymbol{\Phi}^T]_{ij} - 2 \sum_{i=1}^{N/2} \sum_{j=N/2+1}^N [\boldsymbol{\Phi} \boldsymbol{\Phi}^T]_{ij} \right\}. \quad (4.3.11)$$

Pour  $\mathbf{X}$  symétrique, la matrice de vecteurs propres est orthogonale et satisfait l'identité  $\boldsymbol{\Phi} \boldsymbol{\Phi}^T = \mathbf{I}$ , de sorte que

$$\sum_{k=1}^N (\mathbf{u}_N^T \boldsymbol{\phi}_k)^2 \equiv \sum_i^N [\boldsymbol{\Phi} \boldsymbol{\Phi}^T]_{ii} = \text{Tr}(\boldsymbol{\Phi} \boldsymbol{\Phi}^T) = N, \quad (4.3.12)$$

car tous les termes  $[\boldsymbol{\Phi} \boldsymbol{\Phi}^T]_{ij}$  sont nuls sauf les  $N$  termes diagonaux. La quantité  $\langle (\mathbf{u}_N^T \boldsymbol{\phi}_k)^2 \rangle$  est donc égale à 1 pour une moyenne sur l'ensemble. Cette dernière observation permet d'obtenir une relation



pour  $z_1$ . En effet, la moyenne sur l'ensemble de l'Eq. (4.3.5) est

$$\begin{aligned} \left\langle \sum_{k=1}^N \frac{(\mathbf{u}_N^T \boldsymbol{\phi}_k)^2}{(z_j - \lambda_k)} \right\rangle &= \left\langle \frac{2}{(p-q)} \right\rangle \\ \left\langle \sum_{k=1}^N \frac{1}{(z_j - \lambda_k)} \right\rangle &= \frac{2}{(p-q)} \end{aligned} \quad (4.3.13)$$

où on a utilisé le résultat préliminaire pour  $\langle (\mathbf{u}_N^T \boldsymbol{\phi}_k)^2 \rangle$ . Or, on peut réécrire la somme apparaissant dans le CG en fonction de la *trace* de la résolvante en utilisant sa décomposition spectrale (4.3.4)

$$-\text{Tr}(\mathbf{R}_X(z_j)) = \text{Tr} \left( \sum_{k=1}^N \frac{\boldsymbol{\phi}_k \boldsymbol{\phi}_k^T}{z_j - \lambda_k} \right) = \sum_{k=1}^N \frac{1}{z_j - \lambda_k}, \quad (4.3.14)$$

de sorte que la relation

$$-\langle \text{Tr}(\mathbf{R}_X(z_j)) \rangle = \frac{2}{(p-q)} \quad (4.3.15)$$

tient aussi. On a déjà établi précédemment dans la Sec.4.2.2 (cf. Eq. 4.2.20) que

$$\left\langle \text{Tr}(\mathbf{R}_X(\lambda + i\varepsilon)) \right\rangle = \frac{-(\lambda + i\varepsilon) + \sqrt{(\lambda + i\varepsilon)^2 - 4Nm_2}}{2m_2}. \quad (4.3.16)$$

Cette fois, on prend la limite  $\varepsilon \rightarrow 0^+$  avant de prendre la partie imaginaire, ce qui nous mène à

$$\lim_{\varepsilon \rightarrow 0^+} \left\langle \text{Tr}(\mathbf{R}_X(\lambda)) \right\rangle = \frac{-\lambda + \sqrt{\lambda^2 - 4Nm_2}}{2m_2}. \quad (4.3.17)$$

Ce résultat nous permet de déterminer le CG de (4.3.15)

$$-\langle \text{Tr}(\mathbf{R}_X(z_j)) \rangle = \frac{z_j - \sqrt{z_j^2 - 4Nm_2}}{2m_2}, \quad (4.3.18)$$

ce qui nous mène finalement à

$$\frac{z_j - \sqrt{z_j^2 - 4Nm_2}}{2m_2} = \frac{2}{(p-q)} \implies z_j = \frac{1}{2}N(p-q) + 2\frac{m_2}{(p-q)} \quad (4.3.19)$$

Cette expression n'est valide que pour  $z_j \geq 2\sqrt{Nm_2}$  car les valeurs propres  $\{z_j\}$  doivent être réelles ( $\mathbf{B}$  est symétrique). Puisque les  $N-1$  valeurs propres  $\{z_j\}_{j=2,\dots,N}$  sont par définition plus petites ou égales à  $2\sqrt{Nm_2}$ , l'Eq. (4.3.19) ne peut décrire que la valeur propre  $z_1$ . On compare notre prédiction aux résultats empiriques à la figure 4.8.

## 4.4 Limite de détectabilité dans le modèle stochastique par blocs

La conclusion principale du chapitre 3 était que les  $g$  valeurs propres extrêmes d'une matrice de coût appropriée permettent de détecter la structure communautaire naturelle d'un réseau. La matrice

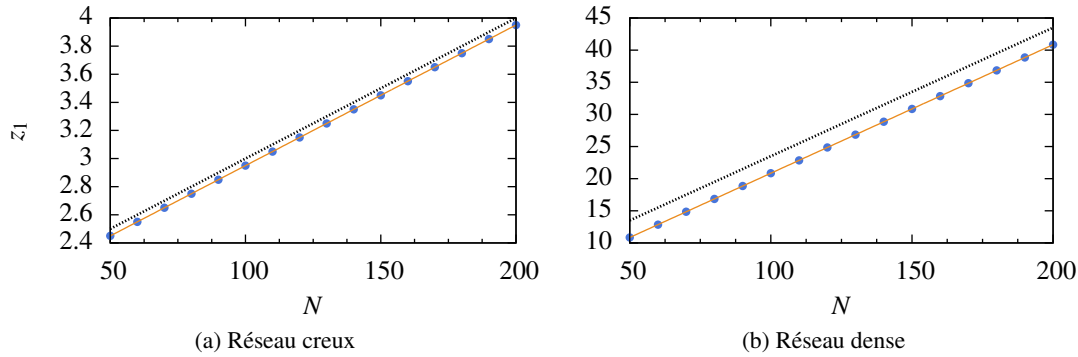


FIGURE 4.8 – **Plus grande valeur propre de la matrice de modularité du SBM à deux blocs.** Comparaison de la valeur moyenne empirique de  $z_1$  (points bleus), de notre prédiction (droite orange) et de la prédiction de Nadakuditi et Newman [67] (droite pointillée noire) dans (a) le régime creux et dans (b) le régime dense.

Les résultats sont obtenus en prenant la valeur moyenne de  $z_1$  sur 500 réalisations du SBM. On utilise les paramètres  $\langle k \rangle = 0.02N$  avec  $p = 0.03$ ,  $q = 0.01$  pour se placer dans un régime approximativement creux, et les paramètres  $\langle k \rangle = 0.7N$  avec  $p = 0.9$ ,  $q = 0.4$  pour se placer dans le régime dense.

de modularité  $\mathbf{B}$  est une de ces matrices. Dans le cas de la modularité, on tente de *maximiser* la fonction de coût, de sorte que la structure communautaire est donnée par les vecteurs propres de  $\mathbf{B}$  associés à ses plus grandes valeurs propres.

Dans le SBM, cette structure communautaire est imposée. En effet, on peut considérer que les  $g$  blocs de noeuds qui ont une densité de liens  $p_{\alpha\beta}$  plus importante que la moyenne  $\langle p_{\alpha\beta} \rangle$  forment  $g$  communautés naturelles. Un bon algorithme de détection devrait donc identifier ces  $g$  blocs correctement. Dans le cas d’un algorithme spectral, la détection parfaite d’une structure naturelle en  $g$  groupes sera uniquement possible si on peut tirer  $(g - 1)$  vecteurs propres contenant de l’information de la matrice de coût<sup>3</sup>.

S’il n’existe pas  $(g - 1)$  valeurs propres détachées, la détection spectrale de la structure “naturelle” sera impossible. En effet, une généralisation récente [84] du cas spécial présenté ici démontre que le spectre de la matrice de modularité d’un SBM est toujours divisé en deux parties, soit une bande continue bornée par la densité spectrale d’une matrice purement aléatoire, et une série de valeurs propres isolées. Les valeurs propres de la bande continue étant associées à une matrice aléatoire, les vecteurs propres correspondant sont eux-même aléatoires, i.e. ils ne contiennent *aucune* information sur la structure modulaire du réseau. Ainsi, la détection spectrale à l’aide d’une matrice de coût  $\mathbf{B}$  devient *strictement* impossible si moins de  $(g - 1)$  valeurs propres sont détachées de la bande continue, pour un SBM de  $g$  communautés naturelles. Le corollaire de cette observation est que si  $(g - 1)$  valeurs propres sont détachées de la bande continue, la détection spectrale est *strictement* possible.

3. Seuls  $(g - 1)$  vecteurs propres sont nécessaires, puisqu’on a montré (Sec. 3.2.3) que la formulation simplexe du problème de détection, qui requiert un vecteur propre de moins, est équivalente. Cette étrange compression d’information est possible puisque le fait qu’un noeud n’appartient pas à un des  $(g - 1)$  groupes nous informe implicitement sur son appartenance au  $g^{\text{ème}}$  groupe.

On peu donc affirmer qu'une limite de détectabilité existe pour un SBM de  $g$  blocs, et qu'elle se situe au point où la  $(g - 1)$ ème plus grande valeur propre rejoint le spectre.

### Limite de détection dans un SBM de 2 blocs identiques

Pour le cas extrêmement simple de SBM étudié dans ce chapitre, on pourrait s'attendre à ce que la valeur propre isolée  $z_1$  rejoigne la bande continue lorsque  $p = q$ , i.e. lorsqu'il n'existe pas de structure communautaire (le réseau est alors indistinguable d'un graphe ER de densité  $p$ ). Or, les points de l'espace  $(p, q)$  où la valeur propre isolée  $z_1 = \frac{1}{2}N(p - q) + 2\frac{m_2}{(p - q)}$  rencontre la limite de la bande aléatoire  $2\sqrt{Nm_2}$  sont plutôt donnés par

$$p = \frac{Nq + 1}{N + 2} \pm \frac{\sqrt{1 + 4q(1 - q)(N + 1)}}{N + 2} \Leftrightarrow q = \frac{Np + 1}{N + 2} \pm \frac{\sqrt{1 + 4p(1 - p)(N + 1)}}{N + 2}, \quad (4.4.1)$$

les solutions de

$$\frac{1}{2}N(p - q) + 2\frac{m_2}{(p - q)} = 2\sqrt{Nm_2} \Leftrightarrow N(p - q) = \sqrt{2N[p(1 - p) + q(1 - q)]} \quad (4.4.2)$$

pour  $N$  fixe et  $m_2$  donné par le calcul exact en régime dense.

Encore une fois, ces résultats diffèrent de ceux obtenus dans la Ref. [68], qui prédisent une limite au point

$$N(p - q) = \sqrt{2N(p + q)}. \quad (4.4.3)$$

Notre correction est particulièrement pertinente dans le contexte de la limite de détection, car elle restaure une symétrie naturelle qui est perdue dans le calcul en régime creux (Fig. 4.9). En effet, on devrait s'attendre à ce que la difficulté du problème de détection soit symétrique autour de l'axe  $q = 1 - p$ , puisqu'on peut toujours partitionner le réseau complémentaire d'un réseau donné. Autrement dit, un réseau très dense devrait être aussi facile à partitionner qu'un réseau très creux, car on peut toujours construire une matrice d'adjacence  $\mathbf{A}'$  duale

$$a'_{ij} = 0 \quad \text{ssi} \quad a_{ij} = 1 \qquad a'_{ij} = 1 \quad \text{ssi} \quad a_{ij} = 0 \quad (4.4.4)$$

et appliquer l'algorithme de détection à ce nouveau réseau.

### Modèle, universalité et limite de détection

Bien qu'on ait adopté une approche spectrale pour démontrer l'existence d'une limite de détection, celle-ci est universelle. *Tout* algorithme de détection *doit* échouer pour  $(p, q, N)$  en dessous de la limite, car la structure communautaire n'est alors pas statistiquement significative. C'est d'ailleurs pourquoi l'existence de la limite a pu être démontrée à l'aide d'une approche complètement différente [32] (calcul en régime creux). Cette assertion est justifiée par le fait que la meilleure façon d'extraire de l'information d'un réseau consiste à ajuster les paramètres du modèle du réseau à une de ses réalisations. Or, dans le cas à l'étude, *on sait* que les réseaux sont produits par un SBM de deux blocs à 2 paramètres. Ainsi, s'il n'est pas possible d'extraire la structure du réseau en connaissant explicitement

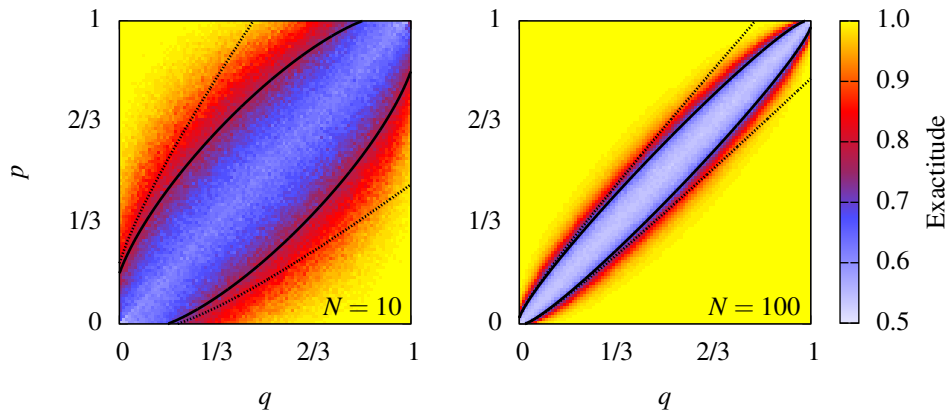


FIGURE 4.9 – **Limite de détectabilité du SBM à deux blocs.** Comparaison de l’exactitude empirique de l’algorithme de détection spectral du chapitre 3 (dégradé), de la prédiction générale de la limite de détection (ligne pleine noire) et de la prédiction de Nadakuditi et Newman (ligne pointillée noire). La différence entre les deux calculs est importante dans le régime dense  $p, q \sim 1$ . Les résultats sont obtenus en moyennant l’exactitude de l’algorithme de détection Sec. 3.4.3 sur 50 réalisations du SBM, pour les  $10^6$  combinaisons de probabilités  $p = 0.01, 0.02, \dots, 1$  et  $q = 0.01, 0.02, \dots, 1$ . Pour  $p < q$ , on détecte la structure mésoscopique du réseau à l’aide des *plus petites* valeurs propres de  $\mathbf{B}$  (le problème est symétrique, cf. Fig. 4.7). Dans ce cas, la structure mésoscopique prend la forme d’une structure coeur-périphérie [75] plutôt que d’une structure modulaire. La symétrie parfaite autour de l’axe  $p = q$  démontre bien que ces problèmes de détection sont équivalents.

le modèle, aucun algorithme ne réussira à obtenir des communautés statistiquement significatives (i.e. autrement que par hasard).

La portée de nos résultats reste toutefois limitée, car peu de systèmes réels sont correctement modélisés par un modèle par blocs. Lorsque les méthodes analytiques échouent, il convient de se tourner vers des méthodes numériques. C’est le passage qu’on initie au prochain et dernier chapitre.

## Chapitre 5

# Attachement préférentiel structurel des communautés

**Structural preferential attachment of community structure and its relation to Dunbar's number**

Jean-Gabriel Young, Laurent Hébert-Dufresne

Antoine Allard et Louis J. Dubé

Département de Physique, de Génie Physique, et d'Optique,  
Université Laval, Québec, Québec, Canada G1V A06.

En préparation

## 5.1 Avant-propos

Ce dernier chapitre de recherche originale est en fait une version préliminaire d'un nouvel article. On y introduit une extension locale du modèle d'attachement préférentiel structurel (SPA) [44], sous la forme d'un processus de croissance stochastique reproduisant la structure *interne* des communautés des réseaux complexes réels. Le modèle SPA offrait déjà une excellente description des réseaux réels au niveau mésoscopique, de sorte que le modèle complet (SPA\*) fournit une bonne compréhension microscopique *et* mésoscopique des réseaux complexes.

La raison principale ayant menée à ces nouveaux développements est la nécessité d'obtenir un nouveau standard de détection, i.e. une classe de réseaux réalistes, modulaires et paramétrés, dont la structure communautaire est connue. Le modèle SPA\* est un candidat idéal, puisqu'il satisfait tous ces critères.

## 5.2 Résumé

On introduit un mécanisme de croissance intuitif afin de décrire l'évolution interne des communautés des réseaux sociaux. Notre modèle est basé sur l'hypothèse simple que la trame sociale d'une communauté donnée est produite par la juxtaposition de deux comportements différents pour les individus, soient la création de liens à l'intérieur du groupe et l'introduction de nouveaux membres. Des comportements complexes émergent des échelles temporelles opposées que ce processus impliquent pour l'activité des individus et de leurs groupes sociaux. En particulier, on montre que le modèle reproduit des comportements observés dans les réseaux sociaux réels, comme l'apparition d'une limite sur le nombre de relations sociales pouvant être maintenues par un individu moyen dans un groupe social donné (nombre de Dunbar). En intégrant ce mécanisme de croissance à un processus d'attachement préférentiel structurel (SPA) récemment introduit, on réussit à reproduire la structure communautaire, *et la distribution en degrés* des réseaux réels avec une précision inégalée jusqu'à maintenant. Ce faisant, on souligne également l'existence d'une nouvelle propriété intéressante de la structure communautaire des réseaux sociaux : une corrélation importante existe entre le nombre de communautés auxquelles appartient un noeud, et la densité de ces dernières.

## 5.3 Abstract

We introduce an intuitive mechanism to describe the evolution of the internal structure of communities within social networks. Our idea is based on the simple assumption that each individual can, for every social group to which it belongs, develop connections and introduce new members. Complex behaviors emerge from opposing time scales for the activities of individuals and for the sum of individuals gathered in groups. We show how the resulting model reproduces behaviors observed in real social networks and in the anthropological theory known as Dunbar's number, i.e. the empirical observation of a maximal number of ties which an average individuals can sustain within a social group.

In fact, using this growth mechanism within a recently introduced structural preferential attachment (SPA) model we reproduce the community structure *and the degree distribution* of actual complex networks with an unprecedented accuracy. In so doing, we highlight an interesting property of the community structure of social networks, i.e. a strong correlation between the number of communities to which a node belong and the density of the communities to which its belong.

## 5.4 Introduction

Networks are at the center of most quantitative analysis of social systems [109]. They encode the links between different individuals within a mathematical construct that allows quantitative studies of the role of individuals in social networks through multiple metrics, and the analysis of correlations between links [20, 24, 75, 108]. These correlations have received particular attention since links tend to be clustered (meaning the friend of my friend is also my friend) in tightly connected groups [110]. These groups inform us not only about the current structure of a network, but also help predict missing or future links [24].

In a recent study [44, 45], we have presented evidence supporting that networks (as a set of linked individuals) can be interpreted as projections of higher structural levels, such as communities or social groups. From this community-centric point of view, groups do not emerge because links get clustered ; instead, links emerge because individuals join communities. Based on this idea, we showed that a simple growth process based on preferential attachment [12] at the level of communities is sufficient to explain many universal properties of complex networks [44].

According to this *structural preferential attachment model* (SPA) and previous studies in community detection [89], we expect community structure to occur at all scales, with community sizes (individuals per community) and memberships (communities per individual) both following heavy-tailed distribution. These distributions are the signature of preferential growth processes [12, 48, 98]. SPA models this property by considering that all growth events consist of an individual joining a community. The event marks the birth of a new individual with probability  $q$ , and the creation of a new community with probability  $p$ . When an existing individual or community is involved (with complementary probabilities  $1 - q$  and  $1 - p$  respectively), it is chosen *preferentially* to its past activity : an individual with  $x$  memberships or a community of size  $x$  are  $x$  times more likely to be chosen than a node/community with 1 membership/member. By considering basic units of size  $s$  (size of new communities), we distinguish between link-based ( $s = 2$ ) and node-based (or individual-based,  $s = 1$ ) systems.

This model was shown to be able to reproduce the community structure of complex networks [44], but also, to some extent, features of their degree distribution and of their self-similar properties [45]. However, it then remained an open question as to *exactly* how links are created within these communities. We now fill this gap by considering not only the growth of communities, but also the creation of links therein. Our contribution is rooted in a simple idea : nodes recruit new members

and make new links within their communities at two independent and constant rates. In a stochastic framework, this general and simple hypothesis is shown to yield a wide spectrum of possible internal organization (or equivalently : *internal* degree distribution) that is solely parametrized by the ratio of the two aforementioned rates and the size of the community of interest.

To confirm our results, we investigate real complex networks through community detection. This method highlights the fact that, while our ability to detect community based on the links of a network is getting better as new algorithms constantly outperform each other [38, 76, 114], it is still unclear whether we should expect social groups to be homogeneous or heterogeneous, sparse or dense. In that regard, our model not only sheds light on the finer details of community structure, but also provides insights on the community detection itself.

The outline of the chapter is as follows. In Sec. 5.5, we present general theoretical considerations on the nature of community structure to root our model in realistic assumptions. In Sec. 5.6, we introduce and characterize an internal link creation mechanism. In Sec. 5.7 we incorporate this mechanism to a previously proposed preferential attachment scheme at the level of communities (SPA). The resulting improved model (SPA\*) is investigated empirically in Sec. 5.8 through the reproduction of the community structure of real networks. Finally, in Sec. 5.9, we link SPA\* to the anthropological theory known as Dunbar’s number, and offer an alternative explanation to this social phenomenon solely based on network structure.

## 5.5 Considerations on the structure of communities

Even though the field of community detection on networks is blooming, there are currently no consensus on exactly what defines a community [38, 51]. Consequently, the features that should be included in models of network growth with community structure remain unclear. However, some common properties are shared by almost all community detection algorithms, such that one could hope to design a model flexible enough to be useful in association with most algorithms. Figure 5.1 presents the features that tend *not* to appear, across multiple algorithms, and that should be rejected based on intuitive properties of social groups.

First, Fig. 5.1a presents a community formed through a Poissonian scheme (or Bernoulli trial) where each pair of nodes are potentially linked with a fixed probability. This scheme does not reject the possibility of having nodes of degree zero. Yet, intuitively, a community that includes individuals sharing no links with the other members does not make much sense. Using the structure of the network as its sole source of information, no algorithm could potentially assign a node to a group without shared connections. Our first consideration thus implies community connectedness : each member of a community should be reachable from any other member.

Second, Fig. 5.1b presents highly heterogeneous communities. Community structure, under its simplest form, is defined in opposition to random organization, implying correlations between one’s



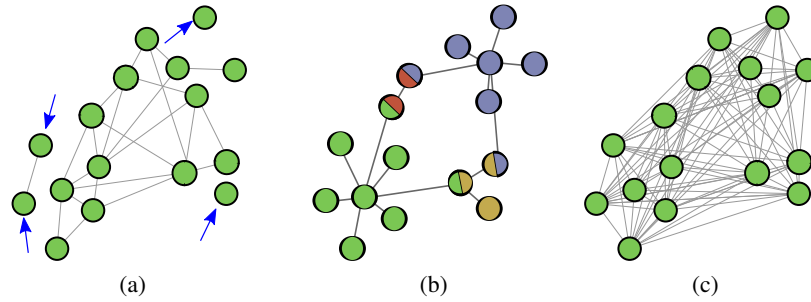


FIGURE 5.1 – **Examples of unacceptable community structures.** (a) Disconnected communities cannot be meaningfully detected on the basis of the network alone. Connectedness is explicitly built within our proposed growth mechanism, such that disconnected structure is never observed. Although (b) heterogeneous communities and (c) large cliques are seldom encountered in social networks, independent of the algorithm used to recover their community structure, our mechanism is flexible enough to allow such extreme behaviors, need it be. See Sec. 5.6 for details.

friends, and the friends of his friends. Community can be heterogeneous, but they must at least be distinguishable from a corresponding rewired (or random) network. This second consideration requires a certain redundancy between the neighborhoods of members of a given community, and consequently, that small community should be denser than their random counterparts.

Third, Fig. 5.1c presents a large, fully connected community. While these communities respect our first two considerations. They are obviously far too rigid to be of practical use. For instance, in a network of commercial ties where we would seek to identify different fields and companies, we would never require a given individual to be connected with every single other employee of his company. Obviously, while a density of 1 might be expected for smaller groups, it is hardly ever achieved in larger social structure (companies, universities, cities). This third consideration merely relaxes the second one : large communities can be sparse even while featuring redundancy in the neighborhoods of their members.

These three considerations and the SPA model discussed in the introduction offer sufficient insight to propose a simple and intuitive, yet accurate, mechanism that describes how social groups tend to grow.

## 5.6 Link creation mechanism

A preferential attachment scheme at the level of communities implies that a community has a growth rate (or recruitment rate) proportional to its current size. This can be interpreted as if each current member of a community introduces a new member at a given rate. This preferential attachment mechanism is not only intuitive, but does in fact reproduce properties of real networks [44]. It is then natural to consider a similar logic for the creation of new links between members of a community.

Our model is simply stated. Each member of a community of size  $n$  recruits a new member at

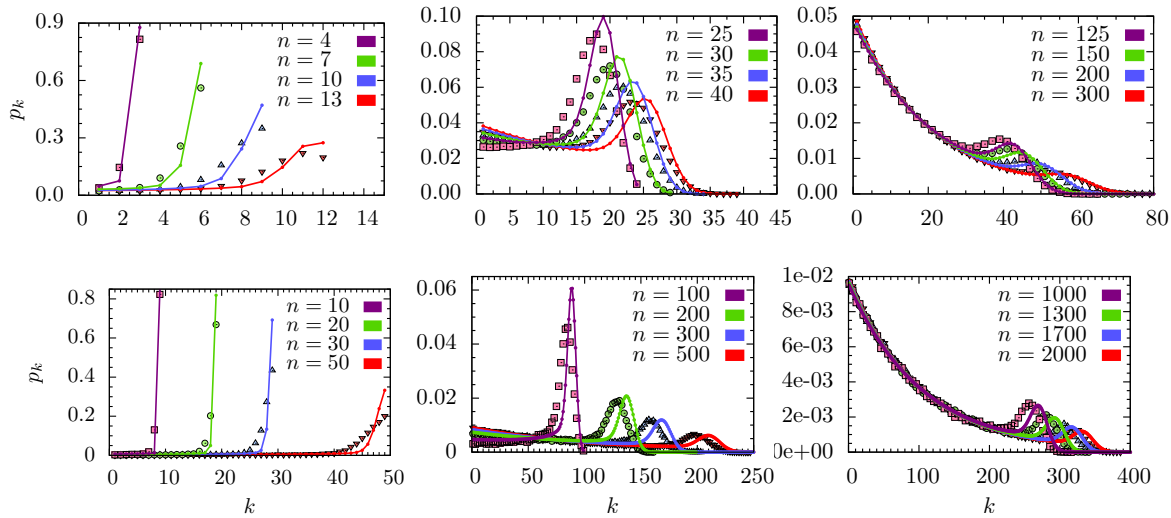


FIGURE 5.2 – **Internal degree distribution in the SPA\* model.** Internal degree distribution  $\{p_k\}$  for various community sizes with  $r = \rho_\ell/\rho_n = 9$  (top) and  $r = 49$  (bottom) in 3 different regimes (left to right). The results of Eq. 5.6.1 (lines) are compared against the average of  $10^4$  Monte-Carlo simulations (closed squares, circles and triangles). Small (left) and medium (center) size communities are highly homogeneous, while degree distributions of larger communities (right) are heavily skewed : with a mode appearing at  $k = 1$  for an average approaching  $\langle k \rangle_n = 20$  (top) and 100 (bottom) in the limit  $n \rightarrow \infty$ . Note that while the same qualitative behavior is observed for both choices of  $r$ , the transition from dense to sparse communities occurs on very different scales.

rate  $\rho_n$ , such that the growth rate  $dn/dt$  of a community of size  $n$  is proportional to  $n\rho_n$ . The new member is initially connected only to the individual who introduced it within the group (its degree, i.e. number of links, equals 1 within this community). This ensures connectedness (consideration #1). Each member also creates a new link at a rate  $\rho_\ell$  until its degree equals  $n - 1$  (such that it is connected to every other member). Bringing in as few assumptions as possible, we postulate that the receiving individual is chosen uniformly from existing (and yet unconnected) members. Consequently, a single member can gain degree faster than rate  $\rho_\ell$  if other members make the effort of creating the link, thus helping smaller communities maintain a high density (consideration #2).

The number  $n_k(t)$  of individuals with degree  $k$  within an average community of size  $n$  that was initially a fully connected clique of size  $s$  can be followed through continuous time  $t$  by a master equation :

$$\begin{aligned} \frac{d}{dt}n_k(t) = & n\rho_n\delta_{k,1} + \rho_\ell (n_{k-1}\bar{\delta}_{k,n} - n_k\bar{\delta}_{k,n-1}) \\ & + \rho_n(n_{k-1} - n_k) + \left( \rho_\ell \sum_{k'=1}^{n-2} n_{k'} \right) \frac{(n-k)n_{k-1} - (n-1-k)n_k}{\sum_{k'=1}^{n-1} (n-1-k')n_{k'}} \end{aligned} \quad (5.6.1)$$

where  $\delta_{i,j}$  is the Kronecker delta and  $\bar{\delta}_{i,j} = 1 - \delta_{i,j}$ . The first term accounts for the arrival of new members (with  $k = 1$ ) at rate  $n\rho_n$ . The second term is due to the creation of links, which brings an individual of degree  $k - 1$  to degree  $k$  [positive for  $n_k(t)$ ], and individual of degree  $k$  to degree  $k + 1$

[negative for  $n_k(t)$ ]. The third term is due to the receiving end of links created when a new individual joins the community; while the last term accounts for the link creation between existing members. The parenthesis gives the creation rate, while the ratio yields the probabilities of it affecting a node of degree  $k - 1$  [positive for  $n_k(t)$ ] or of degree  $k$  [negative for  $n_k(t)$ ]. This complete description of the average state of a community is validated in Fig. 5.2.

A simpler point of view can be adopted to gain further insights into the relation between the internal degree of an average individual (node)  $\langle k \rangle_n$  and the size  $n(t) = \sum_{k'} n_{k'}(t)$  of the community. Since the average internal degree  $\langle k \rangle_n$  is directly related to the average number of links  $L(t)$  within a community of size  $n(t)$  at time  $t$ , we obtain a mean-field equation for the latter as it is easier to follow analytically. Assuming a uniform and uncorrelated distribution of links among individuals, we can write

$$\frac{dL}{dt} = n\rho_n + n\rho_\ell \left[ 1 - \left( \frac{L}{L_{\max}(n)} \right)^{n-1} \right], \quad (5.6.2)$$

since any given individual will create a link at rate  $\rho_\ell$  if he has currently less than  $n - 1$  links. Here  $L_{\max}(n)$  simply equals  $n(n - 1)/2$  (the case where every link exists). A straightforward transformation yields  $L$  as a function of the average size at time  $t$  (using  $dn/dt = n\rho_n$ ):

$$\frac{dL}{dn} = \frac{dL}{dt} \frac{dt}{dn} = 1 + r \left[ 1 - \left( \frac{L}{L_{\max}(n)} \right)^{n-1} \right], \quad (5.6.3)$$

where we defined the ratio  $r := \rho_\ell/\rho_n$ . While the degree distribution is neither uniform nor uncorrelated in the complete model (see Fig. 5.2), we will see that our simplification is robust enough, and that Eq. (5.6.3) accurately reproduces the average density of communities (see Sec. 5.7-5.9).

A simple analysis of Eq. (5.6.3) highlights an interesting feature of this model. For large sizes  $n$ , the ratio  $L/L_{\max}(n)$  goes to zero, such that a maximal link creation rate

$$\frac{dL}{dn} \simeq 1 + r \quad (5.6.4)$$

is attained. Hence, the intensive quantity  $L/n \rightarrow (1 + r)$  converges toward a constant that depends on the parametrization of the model alone. Considering that one link equals two stubs (or degree), the asymptotic average degree is directly related to the parameter  $r$  through :

$$\langle k \rangle_\infty = \frac{2L}{n} = 2(1 + r). \quad (5.6.5)$$

This indicates a maximal average number of connections in social group (consideration #3), which will be investigated in details in Sec.5.9.

## 5.7 Improved growth model

We now link this local community growth model to the complete SPA model [44]. In our link creation mechanism, time  $t$  is continuous and independent of any event. We know however that the

size  $n(t)$  of a community appearing at continuous time  $t_0$  (of initial size  $n(t_0) = s$ ) will follow an exponential growth

$$\frac{dn(t)}{dt} = n(t)\rho_n \quad (5.7.1)$$

In pure SPA, time  $\tilde{T}$  is measured in number of events which can take one of two forms : the birth of a new community — which happens with probability  $p$  — or the growth of an existing one — with complementary probability  $1 - p$ . Without loss of generality, let us define a rescaled discrete time scale  $T$  in which a fraction  $\varepsilon$  of time steps lead to SPA events, such that  $\varepsilon T = \tilde{T}$ . The size  $n(T)$  of a community appearing at discrete time  $T_0$  (of initial size  $n(T_0) = s$ ) is then governed by the master equation

$$\begin{aligned} n(T+1) &= n(T) + \varepsilon(1-p) \frac{n(T)}{\varepsilon T [1 + p(s-1)]} \\ &= n(T) \left[ 1 + \frac{1}{T} \alpha(p; s) \right] \end{aligned} \quad (5.7.2)$$

where we have defined  $\alpha(p; s) := (1-p)/[1 + p(s-1)]$ , and where  $\varepsilon T [1 + p(s-1)]$  is the sum of all sizes at time  $T$ . In the limit of large sizes, (5.7.2) can be rewritten as

$$\frac{dn(T)}{dT} = \frac{n(T)}{T} \alpha(p; s). \quad (5.7.3)$$

Combining both time derivatives then yields the relation between time scales

$$\frac{dt}{dT} = \frac{dt}{dn(t)} \frac{dn(T)}{dT} = \frac{1}{n(t)\rho_n} \frac{n(T)}{T} \alpha(p; s) = \frac{\alpha(p; s)}{\rho_n T} \quad (5.7.4)$$

Let us define  $\tilde{n}(t)$  as the effective size of a community of size  $n(t)$  : the number of members allowed to create links (i.e. those with less than  $n(t) - 1$  links within the community). Then, in the local model, the number of links  $L(t)$  in a community of effective size  $\tilde{n}(t)$  grows at a rate

$$\frac{dL(t)}{dt} = \rho_\ell \tilde{n}(t) \quad (5.7.5)$$

such that links are introduced in a community at rate

$$\frac{dL(T)}{dT} = \frac{dL(t)}{dt} \frac{dt}{dT} = \rho_\ell \tilde{n}(t) \frac{\alpha(p; s)}{\rho_n T} = r\varepsilon(1-p) \frac{\tilde{n}(T)}{\sum n(T)} \quad (5.7.6)$$

under complete SPA growth,

This last result can be interpreted in two ways. Straightforwardly, at each time step ( $dT$ ) of the SPA process, a new link is created between existing members of a given community of effective size  $\tilde{n}$  at rate  $r\varepsilon(1-p)\tilde{n}/\sum n(T)$ . Or equivalently, at each time step ( $dT$ ) of the SPA process, a new link creation is attempted by a randomly chosen member (among all members of all communities) at rate  $r\varepsilon(1-p)$ .

Note that the ratio  $\tilde{n}/\sum n$ , is not equal to one when summed over all communities as  $\sum n(T)$  is the sum of actual (and not effective) sizes. However, if we choose a membership in a uniform random

fashion among all existing memberships, and create the link only if this particular node has at least one connection to be made in this particular community, the rate  $r(1-p)\tilde{n}/T$  will be effectively respected.

The above mathematical construct allows direct incorporation of the link creation mechanism within the SPA model through the following algorithm. Starting with  $s_0$  disjoint and fully connected communities of size  $s$ , at each time step

- **I. (Community growth)** a new community of size  $s$  is created with probability  $p\varepsilon$ , or an existing one (chosen preferentially to its size) grows with probability  $(1-p)\varepsilon$  ;
- **I.a** if a community birth event has been selected, one of the  $s$  involved individual is a new one with probability  $q$  or an existing one (chosen preferentially to its current memberships) with complementary probability  $1-q$ . The other  $s-1$  individuals are chosen preferentially ;
- **I.b** if a community growth event has occurred, the involved individual is a new one with probability  $q$  or an existing one (chosen preferentially) with probability  $1-q$ . This individual then randomly selects an existing member of the community (uniformly) and creates a link ;
- **II. (Link creation)** a new link is created with probability  $r(1-p)\varepsilon$  between an individual chosen preferentially to its memberships and a random available stub within one of its communities (chosen uniformly).

In this process, the parameter  $p$  controls the characteristic scaling of community sizes,  $q$  controls the scaling of node memberships, and  $r$  the density of communities, while  $\varepsilon$  is only set to ensure that all probabilities are smaller than one. In practice, whenever  $(1-p)r > 1$ , we set  $\varepsilon = [r(1-p)]^{-1}$  such that link creation events occur with probability 1 at each step, hence removing unnecessary computation.

Finally, note that this mapping is only approximate, since our analytical expressions for sizes  $n(T)$  and  $n(t)$  are obtained from mean-field arguments. Experiment shows that it holds nicely in general, but that it fails in extreme regime (e.g.  $p, q \sim 0$ ). In such extreme cases, we must forgo shortcuts and perform full simulations, for each communities.

## 5.8 Details of the community structure

We can now use our improved SPA model (coined SPA\*) to reproduce complex social networks to an unprecedented degree : by modelling both the structure at the level of communities, the structure within communities, and the global degree distribution. Our model uses three parameters controlling the importance of modularity ( $p$ ), connectivity between communities ( $q$ ), and connectivity within communities ( $r$ ). The first two are set by fitting the memberships ( $q$ ) and size ( $p$ ) distributions of communities as detected by diverse algorithms. The  $r$  parameter is set by fitting the asymptotic behavior of the observed  $\langle k \rangle_\infty$  to Eq. 5.6.5. A good approximation for the degree distributions of both the complete network and individual communities then *naturally emerge*. In Fig. 5.3, we model a network of co-authorship in the arXiv before 2006 [81] (which acts as a good proxy for social networks). The success of this endeavor shows that SPA\* can reproduce the statistical properties of social networks.

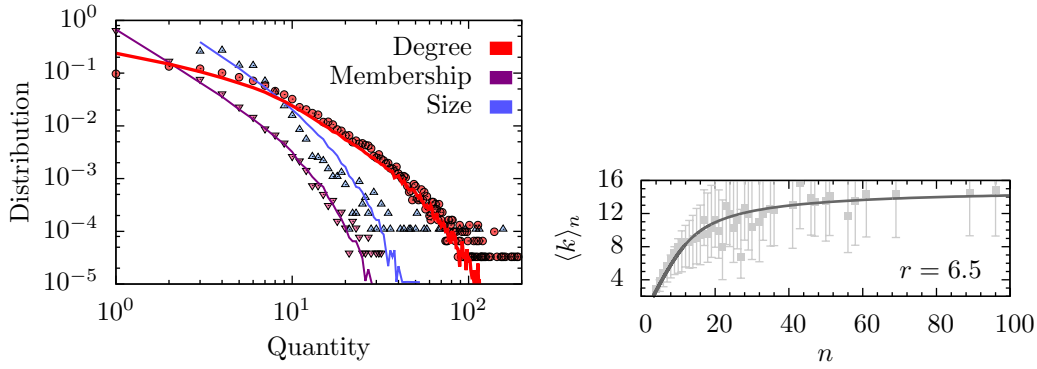


FIGURE 5.3 – **Reproduction of a real networks with the SPA\* model.** Empirical (dots) and predicted (lines) quantities related to the community structure of the arXiv co-authorships network circa 2005 [81], as detected by a cascading approach to clique percolation [116, cf. chapter 2] (CCPA). (left) Community memberships (inverted purple triangle), community size (blue triangle) and degree (red circle) distributions. Predicted results are obtained by averaging over 100 simulations of the SPA\* process, with parameters  $p = 0.57, q = 0.28, s = 1, s_0 = 400, r = 6.5$ . The membership and size distributions are directly fitted to the empirical data using  $p, q$ . The degree distribution is *inferred* from the fit on the right. (right) Mean internal degree  $\langle k \rangle_n$  as a function of community size  $n$ . Predicted results are obtained by integrating eq. 5.6.3 numerically.

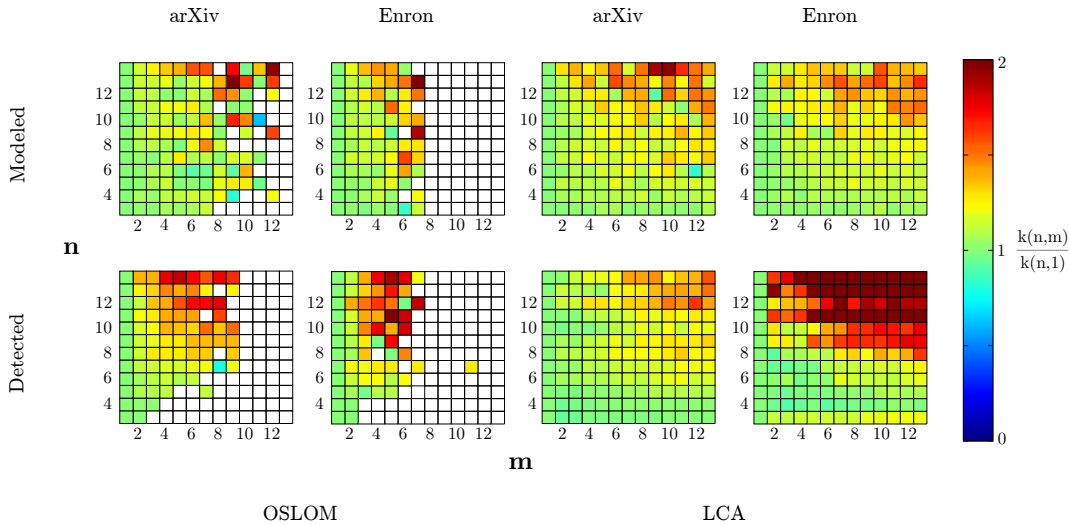


FIGURE 5.4 – **Structural correlations.** Correlations between node memberships  $m$ , community sizes  $n$  and average internal degree  $k(n, m)$  measured in relation to  $k(n, 1)$ . We investigate the the Enron email exchanges [55] and arXiv co-authorships [81] networks. Community structures are (top) produced by SPA\* and (bottom) detected by an algorithm. The empirical and modeled community structures use the result of an (left) order statistics local optimization algorithm, and a (right) link clustering algorithm. Blank squares denote missing data. Note that, without the addition of our link creation mechanism, the classic SPA model does not include any correlations ( $k(n, m) = k(n, 1)$  on average for all  $n$  and  $m$ ) even when considering a given density function of community sizes.

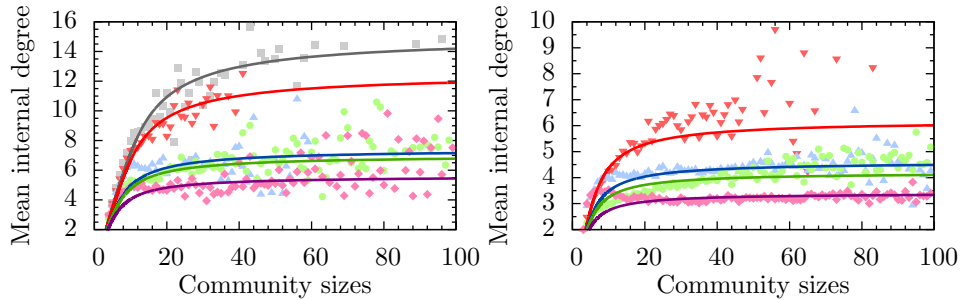


FIGURE 5.5 – **Dunbar’s number in the SPA\* model.** Average number of connections for a given individual within a group of size  $n$ . (left) arXiv [81] (right) MathSciNet [83]. The community detection algorithms are : (gray squares) CCPA [116], (red inverted triangles) LCA [1], (blue triangles) OSLOM [59], (green circles) a greedy modularity optimization of line-graphs algorithm [36], and (magenta lozenge) a greedy clique expansion method [60] (GCE). Analytical predictions of Eq. 5.6.3 are shown in solid lines, tuned for the organization represented by each algorithm. All values of  $r$  were set using Eq. 5.6.5 and are, from top to bottom : (arXiv) 6.5, 4.7, 2.7, 2.5 and 1.8. (MathSci) 2.5, 1.3, 1.1 and 0.7.

Furthermore, results shown in Fig. 5.4 investigate correlations between the organization within communities and the overarching community structure. To do so, we apply an order statistics local optimization method [59] (OSLOM), and a link clustering algorithm [1] (LCA) to the arXiv network, as well as a network of internal email communications within the Enron firm [55]. We find that nodes active in the community structure (high number  $m$  of memberships) tend to also be active within communities (high average internal degree  $\langle k \rangle_n$  in communities of size  $n$ ). This effect is reproduced by SPA\* through age-memberships and age-degree correlations. While the available data does not tell whether these correlations are indeed age-related, it is natural to assume that authors or employees who have been active for a longer time within the archive or the company, tend to have both more groups and more relations within them. Such correlations are seldom considered within growth model, but naturally emerge in SPA\* from our link creation mechanism.

## 5.9 Relation to Dunbar’s number

Previous results and those presented in Fig. 5.5 illustrate how different community detection algorithms possess qualitatively equivalent features even if they use different definitions of communities and consequently investigate different organizations. More precisely, results of Fig. 5.5 highlight how there exist two different behaviors for the average number of links per individual in relation to the size of a social group. For low sizes  $n$ , the mean degree  $\langle k \rangle_n$  essentially scales linearly with the community size as *everybody knows each other within small groups* (e.g. family or close friends). For larger sizes  $n$ ,  $\langle k \rangle_n$  reaches a plateau, Eq. (5.6.5), where a typical individual will not gain new connections when the potential number of connections is increased. So while there is no maximal community size per se, there is a *maximal number of connections that an average individual might possess within a given group* (e.g. large companies or online communities). This effectively implies a maximal community



size as connection density will decrease as size is increased until the community barely qualifies as a group.

Interestingly, this behavior of individual activity  $\langle k \rangle_n$  with respect to community size  $n$  has been previously observed in studies related to an anthropological theory known as Dunbar's number [34]. This theory is based on the observed relation between neocortical sizes in primates and the characteristic sizes of their social groups. Its interpretation usually involves information constraints related to the quality of interpersonal connections and our ability to maintain such relationships. While the importance of neocortical sizes is surely disputable [31], the fact remains that empirical evidence supports the existence of an upper bound in the absolute number of active relationships for a given individual [41]. Similarly, our empirical results indicate that an upper bound may exist in the number of connections one individual can maintain within a given social group.

In our intuitive model, this upper bound naturally emerges and is solely dependent on the  $r$  parameter. This parameter can be interpreted as the ratio between the involvement of an individual in a community, in the sense of bonding with other members, and its contribution to the growth rate of the community. For low  $r$  or large community sizes, the rate of change in the population is higher than an individual's involvement, such that the maximal degree stagnates. Whereas, for high  $r$  and small communities, the individual is able to follow population changes and hence create relationships with most of its members. Thus, different types of social groups will feature different  $r$  and consequently different values of "Dunbar's number".

In this interpretation, the upper bound for individual degree in social groups is due to the fact that connections and introduction of new members have linear requirements for individuals, whereas this implies that groups grow exponentially. Other mathematical models exist to describe Dunbar's number (e.g. [41]), usually based on arguments of priority and/or time resources. However, our model is based on the *observed* community structure of real world networks and consequently, explains Dunbar's number in terms of its two base units, individuals and groups, and the ratio of their respective characteristic rates.

## 5.10 Conclusion

In this chapter, we have presented a model of community growth, SPA\*, which describes both the observed community structure of the network and the internal organization of these communities. Our model follows an intuitive mechanism based on the role of members in recruiting new individuals and creating new links within their communities : a given individual introduces a new member in a community at one rate and creates a new link with existing members at another rate. This simple growth process yields a scale independent community structure, and a plateau of individual involvement in communities (i.e. a maximal degree). These behaviors were shown to reflect both actual data obtained through community detection on complex networks and previous studies. In particular, providing an intuitive explanation for the emergence of the well-known Dunbar's number.



The integration of this mechanism within the existing structural preferential attachment model grants a complete, simple, and intuitive process for generating networks with community structure and tunable density of connections. Future work will investigate the use of this complete model as a benchmark for community detection algorithm. Investigating, for instance, how Dunbar's number implies a limit in detectability for communities of very large sizes as they converge to null density.



# Conclusion et perspectives

Tour à tour, des méthodes heuristique, analytique et numérique ont apportées des points de vue complémentaires du problème de la détection de la structure communautaire des réseaux complexes.

On s'est d'abord attaqué pragmatiquement à ce problème, en adoptant une approche purement heuristique (chapitre 2). Cette approche a permis d'améliorer plusieurs algorithmes de détection communautaire, en apportant une simple correction à la manière d'appliquer ces algorithmes. Les effets importants de cette modification mineure ont non seulement aidé à mettre la complexité du problème de détection en lumière, mais aussi permis de souligner des difficultés inhérentes à la question de la comparaison d'algorithmes.

C'est cette comparaison, essentielle à une bonne compréhension des méthodes de détection, qui a été le fil conducteur des chapitres subséquents. Ainsi, on a d'abord cherché à unifier le problème de la détection sous le chapeau de la théorie spectrale de l'optimisation (chapitre 3). Ce nouveau point de vue a permis de poser les bases d'un formalisme général, construit à partir de principes fondamentaux d'optimisation continue. Cette formulation originale en termes de vecteurs orthogonaux a mené à un algorithme de détection pouvant optimiser une très grande classe de fonctions objectives. De plus, la généralité de l'approche a permis d'établir une relation à sens unique (non inversible) avec une approche bien établie (simplexe).

Fort d'un cadre théorique ouvrant la voie à une comparaison analytique des forces et faiblesses d'une grande classe d'algorithmes, on s'est penché sur les implications de la structure modulaire des réseaux pour cette approche (chapitre 4). Ce faisant, on a réussi à généraliser un résultat important, en démontrant qu'une limite de détection naturelle moins restrictive qu'attendue existe pour les réseaux modulaires denses.

Lorsque la comparaison théorique des algorithmes de détection est devenue impossible, on s'est tourné vers une comparaison numérique. Pour ce faire, on a introduit un nouveau mécanisme de croissance local, qui a par la suite été intégré à un modèle de croissance global, afin d'obtenir un processus stochastique générant des réseaux de structure modulaire variée.

Bien entendu, les développements qui ont été présentés dans cet ouvrage ne couvrent pas complètement le problème. Chaque approche pourra être poussée plus avant dans de futurs travaux.

Sur la plan heuristique, il sera intéressant d'explorer de nouvelles méthodes de détection en cascade. En effet, dans cet ouvrage, une approche radicale a été préconisée, puisqu'on a décidé de simplement retirer les liens problématiques du réseau à chaque itération du méta-algorithme de détection. Des méthodes plus subtiles, qui conservent l'information structurelle à chaque itération, mèneront sûrement à des résultats encore plus probants.

Sur la plan numérique, il faudra compléter le passage du mécanisme de croissance SPA\* vers un banc d'essais pour algorithmes de détection. En effet, bien que les bases du modèle aient été complètement établies dans le cadre de cet ouvrage, il reste encore un dernier effort à accomplir. L'idée sera de complètement caractériser l'espace des paramètres  $[p, q, r]$  (analytiquement lorsque possible, numériquement sinon), afin d'identifier les régions *qualitativement* différentes de l'espace des phases. On pourra alors valider un algorithme de détection en quantifiant sa capacité à récupérer la structure communautaire naturelle de réseaux SPA\* tirés de chacune de ces régions. Cet exercice de validation mènera à une compréhension beaucoup plus approfondie des limitations de chaque algorithme, puisqu'on prédit déjà que l'espace  $[p, q, r]$  de SPA\* sera plus riche et varié que celui du SBM, où des autres bancs d'essais standards du moment, comme les réseaux LFR [58]. Il sera d'ailleurs pertinent de combiner les approches numérique et heuristique, afin de valider le méta-algorithme de détection plus rigoureusement.

L'approche spectrale est toutefois celle qui est la plus prometteuse. Plusieurs directions naturelles de recherche devront être explorées.

Du côté de l'optimisation, il faudra absolument étendre la classe des fonctions de coût pouvant être traitées à l'aide de la méthode. L'inclusion des matrices de coût de somme non nulle (e.g. matrice d'adjacence) sera probablement une extension simple. Pour ce faire, il conviendra d'explorer les choix de paramétrisation plus en profondeur, puisque c'est l'apparition des paramètres libres dans l'équation d'optimisation qui a initialement motivé l'exclusion de ces matrices. Idéalement on cherchera une paramétrisation enracinée dans des principes fondamentaux, afin que le formalisme complet soit réellement général. L'extension à des fonctions de coût non additives sera quant à elle certainement plus difficile (voire impossible), ce qui en fait une question d'autant plus intéressante.

Du côté de la théorie des matrices aléatoires, il faudra évidemment établir le pont entre le cas spécial de SBM à deux blocs et un SBM général. Bien que les spectres des matrices de modularité et d'adjacence soient d'ores et déjà calculables à l'aide d'une méthode alternative [84], une généralisation de l'approche utilisée dans cet ouvrage restera pertinente, puisqu'elle est directement exprimée en termes de vecteurs simplexes<sup>1</sup>, et formulable en termes de vecteurs orthogonaux. Ces calculs permettront de maintenir un lien explicite entre l'optimisation spectrale et la théorie des matrices aléatoires.

---

1. La matrice  $\mathbf{u}$  intervenant dans la paramétrisation de la matrice de modularité est en fait une matrice de vecteurs simplexes dans  $\mathbb{R}^1$  !

## Annexe A

# Fonction génératrice de probabilités

Une fonction génératrice de probabilités (PGF)  $A(x)$  est une série de puissance formelle construite à partir des probabilités  $a_n$  formant une distribution de probabilités discrète  $\{a_n\}$

$$A(x) = a_0 + a_1x + a_2x^2 + a_3x^3 + \dots \quad (\text{A.0.1})$$

Autrement dit, le  $n^{\text{ième}}$  coefficient de la série de puissance de  $A(x)$  est le  $n^{\text{ième}}$  coefficient de la distribution discrète. Les propriétés suivantes sont particulièrement importantes :

**1. Coefficients :** On dit que la distribution  $\{a_n\}$  est générée par  $A(x)$  puisqu'on peut extraire les coefficients de la distribution en appliquant un opérateur dérivée récursivement sur  $A(x)$

$$a_n = \frac{1}{n!} \left. \frac{d^n A(x)}{dx^n} \right|_{x=0} \quad (\text{A.0.2})$$

**2. Normalisation :** La distribution  $\{a_n\}$  est normalisée si la fonction génératrice associée est égale à 1 lorsque évaluée à  $x = 1$

$$A(x)|_{x=1} = \sum_{n=0}^{\infty} a_n = 1 \quad (\text{A.0.3})$$

**3. Moment :** Le  $k^{\text{ième}}$  moment de la distribution  $\{a_n\}$  est donné par l'opérateur

$$\langle n^k \rangle = \left[ \left( x \frac{d}{dx} \right)^k A(x) \right]_{x=1}. \quad (\text{A.0.4})$$

Intuitivement, chaque application de  $x \frac{d}{dx}$  extrait une puissance de  $n$

$$x \frac{d}{dx} \sum_{n=0}^{\infty} a_n x^n = \sum_{n=0}^{\infty} n a_n x^n \quad (\text{A.0.5})$$

tel que l'évaluation à  $x = 1$  mène effectivement aux moments de  $\{a_n\}$ . Le cas spécial  $\langle n \rangle = A'(x)|_{x=1}$  donne la moyenne .

**4. Puissances d'une fonction génératrice :** La  $m^{\text{ième}}$  puissance de la fonction génératrice  $A(x)$  génère la distribution de probabilités de la somme de  $m$  réalisations indépendantes d'une expérience dont les résultats sont distribués selon  $\{a_n\}$ . Cette propriété est plus facile à visualiser à l'aide d'un développement explicite. Par exemple, pour  $m = 2$ , on a

$$[A(x)]^2 = a_0a_0 + (a_0a_1 + a_1a_0)x + (a_0a_2 + 2a_1a_1 + a_2a_0)x^2 + \dots \quad (\text{A.0.6})$$

i.e. le coefficient de  $x^n$  prend en compte toutes les combinaisons menant à un résultat total de  $n$ .

**5. Fonction génératrice de moment :** La fonction génératrice de moments (MGF)  $M(t)$  d'une distribution  $\{a_n\}$  de moments  $\langle n \rangle_{\{a_n\}}, \langle n^2 \rangle_{\{a_n\}}, \dots$  bornés est définie comme la série de puissance formelle

$$M(t) = 1 + \langle n \rangle t + \frac{1}{2!} \langle n^2 \rangle t^2 + \frac{1}{3!} \langle n^3 \rangle t^3 + \dots \equiv \langle e^{nt} \rangle. \quad (\text{A.0.7})$$

On peut récupérer le  $m^{\text{ème}}$  moment de  $\{a_n\}$  en calculant directement la  $m^{\text{ème}}$  dérivée de  $M(t)$  puis en évaluant le résultat à  $t = 0$ . La PGF et la MGF d'une distribution sont reliées analytiquement, car la PGF  $A(x)$  peut aussi être vue comme la moyenne  $A(x) \equiv \langle x^n \rangle_{\{a_n\}}$ , telle que

$$A(e^t) = \langle e^{tn} \rangle = M(t) \quad (\text{A.0.8})$$

## Annexe B

# Compendium de preuves et de définitions

### B.1 Induction d'une organisations de noeuds ou de liens

Soit une composante connexe divisée en communautés contiguës d'au moins deux noeuds, ou d'au moins un lien. Soit les organisations communautaires  $\Omega_x = \{\Psi_1^{(x)}, \Psi_2^{(x)}, \dots, \Psi_g^{(x)}\}$  de noeuds ( $x = n$ ) ou de liens ( $x = \ell$ ) du graphe  $G(\mathcal{V}, \mathcal{E})$ , reliées par induction  $\Omega_n \rightarrow \Omega_\ell$  et  $\Omega_\ell \rightarrow \Omega_n$  (cf. définitions 13-14, Sec. 1.2). Pour ce graphe, l'induction est bijective.

On établit que les relations suivantes sont possibles (une étoile (\*) indique que seul un sous-ensemble restreint peut être induit) :

	$\mathcal{P}_c^{(n)}$	$\mathcal{P}_i^{(n)}$	$\mathcal{C}_c^{(n)}$	$\mathcal{C}_i^{(n)}$
$\mathcal{P}_c^{(\ell)}$	×	×	✓*	×
$\mathcal{P}_i^{(\ell)}$	✓	✓	✓*	✓*
$\mathcal{C}_c^{(\ell)}$	×	×	✓	×
$\mathcal{C}_i^{(\ell)}$	✓*	✓*	✓	✓

Procédant par rangées, on a :

- $\mathcal{P}_c^{(\ell)}$  : Une partition complète de liens induit toujours une couverture complète de noeuds. Il s'agit d'une couverture, parce que tous les liens sont assignés à une communauté tel qu'au moins un noeud se trouvera à l'intersection de deux communautés, ou plus. Cette couverture est forcément complète car chaque noeud est relié à au moins un lien par construction, et que ces derniers sont tous assignés. Le degré d'un noeud place toutefois une borne supérieure sur le nombre de communautés auquel il peut participer, tel qu'une couverture de noeuds arbitraire ne peut être obtenue à partir d'une partition incomplète de liens.
- $\mathcal{P}_i^{(\ell)}$  : Une partition incomplète de liens peut induire n'importe quel type d'organisation de noeuds.
  - Afin de visualiser le type d'induction qui mène à une partition complète de noeuds, on adopte le point de vue inverse et on considère une division d'un graphe en communa-

tés contiguës de noeuds. Deux cas distincts seront rencontrés : (a) un lien est borné par deux noeuds dans la même communauté et hérite de l'assignation communautaire de ces noeuds, (b) un lien relie des noeuds assignés à des communautés différentes et ne peut pas être classifié. La partition de liens induite est donc une partition incomplète de liens.

- Toujours avec un point de vue inverse, si on considère une partition de noeuds incomplète, le cas (b) défini ci-haut sera rencontré plus souvent. De plus, la possibilité qu'un lien soit borné par au moins un noeud non assigné s'ajoute. L'organisation induite est donc à nouveau une partition de liens incomplète.
- Les couvertures (in)complètes sont des cas particuliers de partitions (in)complètes, et sont donc des possibilités par définition.
- Des couvertures au sens fort (intersection plus grande que 0) sont aussi possibles, car un noeud connecté par deux liens assignés à des communautés différentes héritera de deux assignations communautaires (comme dans le cas de la partition complète de liens).

En pratique, une partition de liens incomplète mène généralement à une couverture incomplète, au sens fort.

- $C_c^{(\ell)}$  : Les mêmes considérations que dans le cas  $\mathcal{P}_c^{(\ell)}$  s'appliquent ici, à la différence qu'une couverture arbitraire de noeuds peut maintenant être obtenue, car le degré n'impose plus de borne supérieure.
- $C_i^{(\ell)}$  : Les mêmes considérations que dans le cas  $\mathcal{P}_i^{(\ell)}$  s'appliquent ici à la différence qu'une couverture arbitraire de noeuds peut maintenant être obtenue, car le degré n'impose plus de borne supérieure. Il convient de noter qu'une partition de noeuds est induite uniquement dans le cas trivial où la couverture de liens est en fait une partition de liens. En effet, si un lien est assigné à plus d'une communauté, les noeuds qui le bornent seront automatiquement associés à plus d'une communauté.

## B.2 Graphe adjoint

Le graphe adjoint  $G' = (\mathcal{V}', \mathcal{E}')$  est un graphe résultant d'une bijection<sup>1</sup> qui associe les liens du graphe original  $G = (\mathcal{V}, \mathcal{E})$  à des noeuds dans  $G'$ . La matrice d'adjacence  $\mathbf{A}'$  du graphe adjoint  $G'$  est évidemment reliée analytiquement à la matrice d'adjacence  $\mathbf{A}$  du graphe original  $G$ . Ce passage n'est toutefois pas direct, dans le sens où les matrices d'adjacence  $\mathbf{A}$  et  $\mathbf{A}'$  peuvent être vues comme des projections d'une matrice plus générale, soit la matrice d'incidence *des liens*  $\mathbf{Z}$ .

**Définition 27.** *La matrice d'incidence  $\mathbf{Z}$  est une matrice  $N \times M$  binaire décrivant la relation entre l'ensemble de noeuds  $\mathcal{V}$  et de liens  $\mathcal{E}$  de  $G$ . Adoptant la convention que les liens sont numérotés à l'aide d'indices grecs (plutôt qu'avec les indices des noeuds qu'ils connectent), on a que l'élément  $b_{i\alpha}$  prend la valeur 1 si le lien  $\alpha$  est connecté au noeud  $i$  dans  $G$ .*

---

1. Sauf dans le cas particulier d'une étoile de 4 noeuds ou d'un triangle déconnecté [36, p.4].



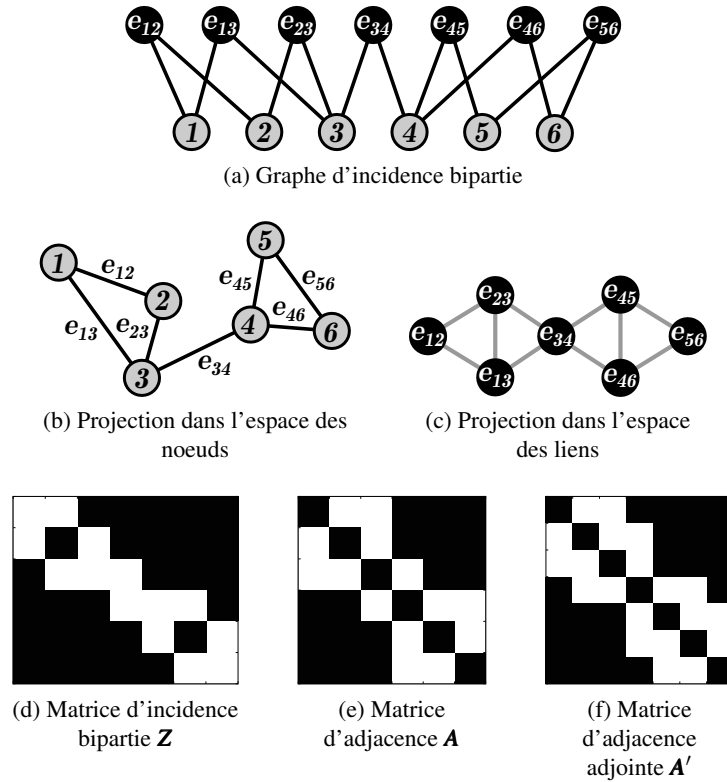


FIGURE B.1 – **Graphe et graphe adjoint en tant que projection d'un graphe bipartie.** Pour les représentations en graphe (a)-(c), les noeuds sont colorés en gris alors que les liens sont représentés à l'aide de noeuds noirs. Pour les représentations matricielles (d)-(f), une case noire indique l'entrée binaire 0, alors qu'une case blanche indique l'entrée binaire 1. (a) et (d) Représentation la plus générale du graphe  $G = (\mathcal{V}, \mathcal{E})$ , sous la forme de couples entre noeuds et liens. (b) et (e) Projection vers l'espace habituel des noeuds à l'aide de l'eq. (B.2.1a). (c) et (f) Projection vers l'espace adjoint des liens à l'aide de l'eq. (B.2.1b).

Cette matrice peut donc être vue comme la matrice d'adjacence d'un réseau bipartite qui met en relation d'un côté les noeuds et de l'autre côté les liens de  $G$ .

On a alors les deux projections suivantes (graphe non dirigé, non dirigé, sans boucle)

$$a_{ij} = \sum_{\alpha=1}^M z_{i\alpha} z_{j\alpha} (1 - \delta_{ij}) \quad \Longrightarrow \quad \mathbf{A} = \mathbf{Z}\mathbf{Z}^T - \text{diag}(\mathbf{k}) \quad (\text{B.2.1a})$$

$$a'_{\alpha\beta} = \sum_{i=1}^N z_{i\alpha} z_{i\beta} (1 - \delta_{\alpha\beta}) \quad \Longrightarrow \quad \mathbf{A}' = \mathbf{Z}^T \mathbf{Z} - 2\mathbf{I}_m, \quad (\text{B.2.1b})$$

où  $\text{diag}(\mathbf{k})$  est la matrice  $N \times N$  diagonale des degrés. Les projections (B.2.1a)-(B.2.1b) sont illustrées à la figure B.1 pour un réseau jouet arbitraire.

Le lecteur intéressé est invité à consulter la référence [36], où d'autres transformations plus complexes incluant des boucles et des poids sont aussi été proposées.

### B.3 Produit de Hadamard et trace

**Théorème 1.** Pour  $\mathbf{A}, \mathbf{B}$ , des matrices colonnes  $n \times 1$ , l'égalité

$$\mathbf{1}_n^T (\mathbf{A} \circ \mathbf{B}) \equiv \text{Tr} (\mathbf{A} \mathbf{B}^T),$$

tient, où  $\circ$  est le produit de Hadamard défini par  $[\mathbf{A} \circ \mathbf{B}]_{ij} = a_{ij} b_{ij}$ .

**Preuve :** Par définition du produit d'Hadamard,

$$[\mathbf{A} \circ \mathbf{B}]_{i1} = a_{i1} b_{i1}. \quad (\text{B.3.1})$$

Puisque  $\mathbf{A}, \mathbf{B}$  sont des matrices colonnes, on peut réorganiser la somme

$$\mathbf{1}_n^T (\mathbf{A} \circ \mathbf{B}) = \sum_{i=1}^n a_{i1} b_{i1} \quad (\text{B.3.2})$$

trivialement sous la forme d'une trace

$$\sum_{i=1}^n a_{i1} b_{i1} = \sum_{i=1}^n a_{i1} b_{1i} = \text{Tr} (\mathbf{A} \mathbf{B}^T) \quad (\text{B.3.3})$$

**Théorème 2.** Pour  $\mathbf{A}, \mathbf{B}$ , des matrices  $m \times n$  et  $\mathbf{X}$ , une matrice  $n \times n$ , l'égalité

$$\mathbf{1}_n^T (\mathbf{X} \circ (\mathbf{A}^T \mathbf{B})) \mathbf{1}_n \equiv \text{Tr} (\mathbf{A} \mathbf{X} \mathbf{B}^T),$$

tient, où  $\circ$  est le produit de Hadamard défini par  $[\mathbf{A} \circ \mathbf{B}]_{ij} = a_{ij} b_{ij}$ .

**Preuve :** Par définition du produit d'Hadamard,

$$[\mathbf{X} \circ (\mathbf{A}^T \mathbf{B})]_{ij} = x_{ij} [\mathbf{A}^T \mathbf{B}]_{ij} = x_{ij} \sum_{k=1}^m a_{ki} b_{kj}. \quad (\text{B.3.4})$$

Ainsi, la somme de toutes les entrées est donnée de façon équivalente par une triple sommation

$$\mathbf{1}_n^T (\mathbf{X} \circ (\mathbf{A}^T \mathbf{B})) \mathbf{1}_n = \sum_{i=1}^n \sum_{j=1}^n [\mathbf{X} \circ (\mathbf{A}^T \mathbf{B})]_{ij} = \sum_{i=1}^n \sum_{j=1}^n x_{ij} \sum_{k=1}^m a_{ki} b_{kj} \quad (\text{B.3.5})$$

qui peut être réorganiser sous la forme d'une trace

$$\sum_{i=1}^n \sum_{j=1}^n x_{ij} \sum_{k=1}^m a_{ki} b_{kj} = \sum_{i,j,k} [\mathbf{A}]_{ki} [\mathbf{X}]_{ij} [\mathbf{B}]_{kj} = \sum_k [\mathbf{A} \mathbf{X} \mathbf{B}^T]_{kk} = \text{Tr} (\mathbf{A} \mathbf{X} \mathbf{B}^T). \quad (\text{B.3.6})$$

## B.4 Dérivées d'une transformation de la PGF binomiale

**Théorème 3.** La  $n_\beta$ ème dérivée de la fonction

$$\Delta_{\alpha\beta}^{(\ell)}(t) := \frac{s_\beta!}{(s_\beta - \ell)!} p_{\alpha\beta}^\ell e^{\ell t} (1 - p_{\alpha\beta} + e^t p_{\alpha\beta})^{s_\beta - \ell}, \quad (\text{B.4.1})$$

est donnée par

$$\frac{\partial^{n_\beta}}{\partial t^{n_\beta}} \Delta_{\alpha\beta}^{(0)}(t) = \sum_{\ell=1}^{n_\beta} \left\{ \begin{matrix} n_\beta \\ \ell \end{matrix} \right\} \Delta_{\alpha\beta}^{(\ell)}(t) \quad (\text{B.4.2})$$

**Preuve :** On procède par induction. On veut donc démontrer que

$$\frac{\partial^{n_\beta+1}}{\partial t^{n_\beta+1}} \Delta_{\alpha\beta}^{(0)}(t) = \sum_{\ell=1}^{n_\beta+1} \left\{ \begin{matrix} n_\beta+1 \\ \ell \end{matrix} \right\} \Delta_{\alpha\beta}^{(\ell)}(t) \quad (\text{B.4.3})$$

à partir de (B.4.2). Dans un premier temps, on remarque que la dérivée  $\frac{\partial}{\partial t} \Delta_{\alpha\beta}^{(\ell)}(t)$  satisfait la pseudo relation de récurrence

$$\begin{aligned} \frac{\partial}{\partial t} \Delta_{\alpha\beta}^{(\ell)}(t) &= \frac{s_\beta!(s_\beta - \ell)}{(s_\beta - \ell)!} p_{\alpha\beta}^{\ell+1} e^{(\ell+1)t} (1 - p_{\alpha\beta} + e^t p_{\alpha\beta})^{s_\beta - (\ell+1)} + \frac{s_\beta! \ell}{(s_\beta - \ell)!} p_{\alpha\beta}^\ell e^{\ell t} (1 - p_{\alpha\beta} + e^t p_{\alpha\beta})^{s_\beta - \ell} \\ &= \Delta_{\alpha\beta}^{(\ell+1)}(t) + \ell \Delta_{\alpha\beta}^{(\ell)}(t), \end{aligned} \quad (\text{B.4.4})$$

à condition que  $\ell < s_\beta$ . Dérivant l'eq. (B.4.2) à nouveau, on obtient

$$\begin{aligned} \frac{\partial^{n_\beta+1}}{\partial t^{n_\beta+1}} \Delta_{\alpha\beta}^{(0)}(t) &= \sum_{\ell=1}^{n_\beta} \left\{ \begin{matrix} n_\beta \\ \ell \end{matrix} \right\} (\Delta_{\alpha\beta}^{(\ell+1)}(t) + \ell \Delta_{\alpha\beta}^{(\ell)}(t)) \\ &= \sum_{\ell=1}^{n_\beta} \left\{ \begin{matrix} n_\beta \\ \ell \end{matrix} \right\} \Delta_{\alpha\beta}^{(\ell+1)}(t) + \sum_{\ell=1}^{n_\beta} \left\{ \begin{matrix} n_\beta+1 \\ \ell \end{matrix} \right\} \Delta_{\alpha\beta}^{(\ell)}(t) - \sum_{\ell=1}^{n_\beta} \left\{ \begin{matrix} n_\beta \\ \ell-1 \end{matrix} \right\} \Delta_{\alpha\beta}^{(\ell)}(t), \end{aligned} \quad (\text{B.4.5})$$

où on a appliqué la relation (B.4.4) à la première égalité et utilisé la propriété

$$\ell \left\{ \begin{matrix} n_\beta \\ \ell \end{matrix} \right\} = \left\{ \begin{matrix} n_\beta+1 \\ \ell \end{matrix} \right\} - \left\{ \begin{matrix} n_\beta \\ \ell-1 \end{matrix} \right\} \quad (\text{B.4.6})$$

des nombres de Stirling de la deuxième espèce [100, p.82]. La dernière somme de (B.4.5) peut être manipulée algébriquement pour obtenir une forme similaire à la première somme

$$\sum_{\ell=1}^{n_\beta} \left\{ \begin{matrix} n_\beta \\ \ell-1 \end{matrix} \right\} \Delta_{\alpha\beta}^{(\ell)}(t) = \sum_{\ell=0}^{n_\beta-1} \left\{ \begin{matrix} n_\beta \\ \ell \end{matrix} \right\} \Delta_{\alpha\beta}^{(\ell+1)}(t) = \sum_{\ell=1}^{n_\beta} \left\{ \begin{matrix} n_\beta \\ \ell \end{matrix} \right\} \Delta_{\alpha\beta}^{(\ell+1)}(t) + \left\{ \begin{matrix} n_\beta \\ 0 \end{matrix} \right\} \Delta_{\alpha\beta}^{(1)}(t) - \left\{ \begin{matrix} n_\beta \\ n_\beta \end{matrix} \right\} \Delta_{\alpha\beta}^{(n_\beta+1)}(t). \quad (\text{B.4.7})$$

Le deuxième terme de (B.4.7) est nul (car  $\left\{ \begin{matrix} n_\beta \\ 0 \end{matrix} \right\} \equiv 0$ ), tandis que le troisième terme peut être exprimé comme

$$\left\{ \begin{matrix} n_\beta \\ n_\beta \end{matrix} \right\} \Delta_{\alpha\beta}^{(n_\beta+1)}(t) = \Delta_{\alpha\beta}^{(n_\beta+1)}(t) = \left\{ \begin{matrix} n_\beta+1 \\ n_\beta+1 \end{matrix} \right\} \Delta_{\alpha\beta}^{(n_\beta+1)}(t). \quad (\text{B.4.8})$$

La somme (B.4.7) est donc égale à

$$\sum_{\ell=1}^{n_\beta} \left\{ \begin{matrix} n_\beta \\ \ell-1 \end{matrix} \right\} \Delta_{\alpha\beta}^{(\ell)}(t) = \sum_{\ell=1}^{n_\beta} \left\{ \begin{matrix} n_\beta \\ \ell \end{matrix} \right\} \Delta_{\alpha\beta}^{(\ell+1)}(t) + \left\{ \begin{matrix} n_\beta+1 \\ n_\beta+1 \end{matrix} \right\} \Delta_{\alpha\beta}^{(n_\beta+1)}(t). \quad (\text{B.4.9})$$

On substitue cette expression dans l'expression originale (B.4.5), qui se simplifie alors à

$$\frac{\partial^{n_\beta+1}}{\partial t^{n_\beta+1}} \Delta_{\alpha\beta}^{(0)}(t) = \sum_{\ell=1}^{n_\beta+1} \left\{ \begin{matrix} n_\beta+1 \\ \ell \end{matrix} \right\} \Delta_{\alpha\beta}^{(\ell)}(t). \quad (\text{B.4.10})$$

Ce dernier résultat, couplé à la validité de la condition initiale pour  $n_\beta = 1$  (démontré par la bande en B.4.4) prouve l'eq (B.4.2) tient pour tout  $\ell$ , par induction.

## B.5 Formule de Stieltjes-Perron

**Théorème 4.** Soit  $\mathbf{X}$  une matrice symétrique réelle  $N \times N$  admettant les valeurs propres  $\{\lambda_j\}_{j=1,\dots,N}$  associées aux vecteurs propres orthonormés  $\phi_j^\top$ . La matrice résolvante de  $\mathbf{X}$

$$R_{\mathbf{X}}(z) := (\mathbf{X} - z\mathbf{I})^{-1}, \quad z \in \mathbb{C} \setminus \mathbb{R}. \quad (\text{B.5.1})$$

est reliée à la densité spectrale  $\rho(\lambda)$  de  $\mathbf{X}$  par la formule de Stieltjes-Perron

$$\rho(\lambda) = \frac{1}{\pi} \lim_{\varepsilon \rightarrow 0^+} \text{Im} \left\{ \left\langle \text{Tr}(R_{\mathbf{X}}(\lambda + i\varepsilon)) \right\rangle \right\}. \quad (\text{B.5.2})$$

**Preuve :** Les décompositions spectrales de la matrice  $\mathbf{X}$  et de l'identité  $\mathbf{I}_{N \times N}$  sur la base des vecteurs propres  $\phi_j^\top$  de  $\mathbf{X}$  sont données par

$$\mathbf{X} = \sum_{j=1}^N \lambda_j \phi_j \phi_j^\top \quad (\text{B.5.3a})$$

$$\mathbf{I} = \sum_{j=1}^N \phi_j \phi_j^\top, \quad (\text{B.5.3b})$$

tel que la décomposition spectrale de l'inverse de la matrice résolvante est

$$R_{\mathbf{X}}^{-1}(z) = \sum_{j=1}^N \phi_j \phi_j^\top (\lambda_j - z) \quad (\text{B.5.4})$$

Par inspection, on note que cette formule pour l'inverse de la résolvante détermine uniquement la décomposition de la résolvante même

$$R_{\mathbf{X}}(z) = \sum_{j=1}^N \frac{\phi_j \phi_j^\top}{(\lambda_j - z)}, \quad (\text{B.5.5})$$

ce qui peut être vérifié aisément en notant que  $\phi_j^T \phi_k = \delta_{jk}$ , de sorte que le produit matriciel  $R_{\mathbf{X}}^{-1}(z)R_{\mathbf{X}}(z)$  est effectivement égale à l'identité. La trace de la résolvante est donc donnée par

$$\mathrm{Tr}(R_{\mathbf{X}}(z)) \equiv \sum_{k=1}^N [R_{\mathbf{X}}(z)]_{kk} = \sum_{j,k=1}^N \frac{\phi_k^{(j)} \phi_k^{(j)}}{(\lambda_j - z)} = \sum_{j=1}^N (\lambda_j - z)^{-1}, \quad (\text{B.5.6})$$

où  $\phi_k^{(j)}$  dénote le  $k^{\text{ème}}$  élément du  $j^{\text{ème}}$  vecteur propre, et où on a utilisé le fait que la base est orthonormée, i.e. que  $\sum_{k=1}^N \phi_k^{(j)} \phi_k^{(j)} \equiv \phi_j^T \phi_j = 1 \forall j$ .

Considérons maintenant  $\rho(\lambda_1, \lambda_2, \dots, \lambda_N)$ , la densité de probabilité jointe des valeurs propres de  $\mathbf{X}$ . Pour une telle fonction de poids, la valeur moyenne de  $\mathrm{Tr}(R_{\mathbf{X}}(z))$  est définie comme

$$\langle \mathrm{Tr}(R_{\mathbf{X}}(z)) \rangle = \sum_{j=1}^N \int_{-\infty}^{\infty} d\lambda_1 \dots \int_{-\infty}^{\infty} d\lambda_j \dots \int_{-\infty}^{\infty} d\lambda_N \frac{\rho(\lambda_1, \dots, \lambda_j, \dots, \lambda_N)}{(\lambda_j - z)} \quad (\text{B.5.7})$$

Les intégrales sur les  $\{d\lambda_\ell\}_{\ell \neq j}$  sont triviales, puisque les valeurs propres  $\lambda_\ell$  n'apparaissent que dans la densité de probabilité, et que le passage vers une densité de probabilité marginale n'affecte pas la normalisation. Les variables  $\lambda_1, \dots, \lambda_N$  sont interchangeables, de sorte que l'intégrale (B.5.7) est équivalente à l'intégrale simple

$$\langle \mathrm{Tr}(R_{\mathbf{X}}(z)) \rangle = N \int_{-\infty}^{\infty} d\lambda_j \frac{\rho_1(\lambda_j)}{(\lambda_j - z)}, \quad (\text{B.5.8})$$

où  $\rho_1(\lambda_j)$  est la densité de probabilité marginale de la valeur propre  $\lambda_j$ , i.e. la densité spectrale  $\mathbf{X}$ .

Dans la définition de la matrice résolvante (Eq. B.5.1), la variable  $\lambda$  est considérée comme étant strictement complexe ( $z \in \mathbb{R}$  est exclu). On rend cette paramétrisation explicite en recourant à la notation  $\lambda = \mathrm{Re}\{\lambda\} + i\mathrm{Im}\{\lambda\} := \lambda + i\varepsilon$ , avec  $\varepsilon \neq 0$ . Considérons maintenant la partie imaginaire de l'eq. (B.5.8), i.e.

$$\mathrm{Im} \left\{ \langle \mathrm{Tr}(R_{\mathbf{X}}(\lambda + i\varepsilon)) \rangle \right\} = N \mathrm{Im} \left\{ \int_{-\infty}^{\infty} \frac{\rho_1(\lambda_j) d\lambda_j}{(\lambda_j - \lambda - i\varepsilon)} \right\} = N \int_{-\infty}^{\infty} \frac{\varepsilon \rho_1(\lambda_j) d\lambda_j}{(\lambda_j + \lambda)^2 + \varepsilon^2}. \quad (\text{B.5.9})$$

Dans la limite  $\varepsilon \rightarrow 0^+$ , la contrainte  $\lambda \in \mathbb{C} \setminus \mathbb{R}$  est respectée et cette intégrale devient triviale, car une distribution de Dirac apparaît. Ainsi,

$$\lim_{\varepsilon \rightarrow 0^+} \mathrm{Im} \left\{ \langle \mathrm{Tr}(R_{\mathbf{X}}(\lambda + i\varepsilon)) \rangle \right\} = N \int_{-\infty}^{\infty} \lim_{\varepsilon \rightarrow 0^+} \frac{\varepsilon \rho_1(\lambda_j) d\lambda_j}{(\lambda_j + \lambda)^2 + \varepsilon^2} = N \int_{-\infty}^{\infty} \pi \delta(\lambda_j + \lambda) \rho_1(\lambda_j) d\lambda_j = N \pi \rho_1(\lambda). \quad (\text{B.5.10})$$

Puisque la densité de probabilité marginale  $\rho_1(\lambda)$  est reliée à la densité spectrale (fonction de corrélation en un point) par

$$\rho(\lambda) = \left\langle \sum_j \delta(\lambda - \lambda_j) \right\rangle = \sum_j \langle \delta(\lambda - \lambda_j) \rangle = N \rho_1(\lambda), \quad (\text{B.5.11})$$

on obtient

$$\rho(\lambda) = \frac{1}{\pi} \lim_{\varepsilon \rightarrow 0^+} \text{Im} \left\{ \left\langle \text{Tr}(\mathcal{R}_{\mathbf{X}}(\lambda + i\varepsilon)) \right\rangle \right\}. \quad (\text{B.5.12})$$

Cette dernière équation est connue sous le nom de la formule de Stieltjes-Perron.

# Annexe C

## Vecteurs simplexes réguliers

Dans cette annexe, on effectue une preuve complète de la relation géométrique liant la base orthogonale de l'espace vectoriel  $\mathbb{R}^{n+1}$  aux vecteurs simplexes réguliers de l'espace  $\mathbb{R}^n$ .

Plus spécifiquement, on présente le concept de simplexes réguliers à la Sec. C.1, avant de survoler la méthode numérique canonique de construction de ces vecteurs (Sec. C.2). Ensuite, dans la Sec. C.3, on discute de la relation géométrique qui existe entre la figure de sommet formée par les vecteurs orthogonaux de  $\mathbb{R}^{n+1}$  et les  $n + 1$  vecteurs simplexes réguliers vivant dans l'espace  $\mathbb{R}^n$ . Une démarche analytique formalisant cette relation est finalement présentée aux sections C.4-C.5.

### C.1 $n$ -simplexes

Un  $n$ -simplexe est une généralisation du concept de triangle à l'espace  $\mathbb{R}^n$ , pour  $n > 2$ . Afin de guider cette généralisation à des dimensions où l'intuition n'est pas applicable, on définira un triangle comme l'ensemble des points qui sont contenus dans le plan formé en reliant 3 points distinct de l'espace  $\mathbb{R}^2$ . Ce triangle (ou 2-simplexe) peut être entièrement défini par les vecteurs *simplexes*  $\{\vec{s}_\alpha\}_{\alpha=1,2,3}$  qui relient l'origine aux coins du triangle. Alternativement, on peut complètement décrire ce triangle à l'aide d'un ensemble de produits scalaires. D'une part, les produits scalaires entre les différentes paires, e.g.  $\vec{s}_1 \cdot \vec{s}_2$ , de vecteurs fixent les angles relatifs entre ces vecteurs. D'autre part, les produits scalaires des vecteurs avec eux-mêmes, e.g.  $\vec{s}_1 \cdot \vec{s}_1$ , fixent la norme de ces vecteurs et donc les dimensions du triangle.

On parle de  $n$ -simplexes *réguliers* lorsque que ces objets géométriques sont une généralisation du concept de triangle *équilatérale* à un espace de dimension arbitraire. Le triangle équilatéral est définie par (a) un angle identique entre chaque paire connexe de vecteurs  $(\vec{s}_\alpha, \vec{s}_\beta)$  et (b) une norme identique pour tous les vecteurs  $\vec{s}_\alpha$ .

On choisit la normalisation telle que

$$\vec{s}_\alpha \cdot \vec{s}_\alpha = 1 - \frac{1}{n+1}. \quad (\text{C.1.1a})$$

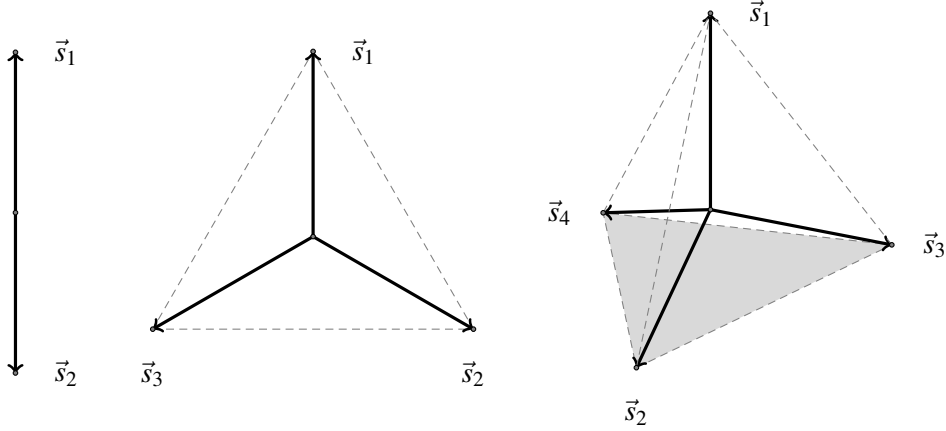


FIGURE C.1 – **Vecteurs simplexes réguliers de basse dimension.** En dimension  $\mathbb{R}^n$  pour  $n = 1, 2, 3$  (gauche, centre, droite), les simplexes réguliers sont simplement une droite finie de dimension fini, un triangle équilatéral et une pyramide à base triangulaire. Les vecteurs simplexes  $\vec{s}_\alpha^{(n)}$  partent de l'origine et pointent vers les coins de ces objet géométriques.

Sans en faire la preuve, on mentionne que ce choix de normalisation fixe les produits scalaire entre vecteurs distincts à

$$\vec{s}_\alpha \cdot \vec{s}_\beta = -\frac{1}{n+1} \quad \alpha \neq \beta. \quad (\text{C.1.1b})$$

## C.2 Méthode numérique de construction des vecteurs simplexes réguliers

Les propriétés (C.1.1a)-(C.1.1b) peuvent être utilisées pour construire explicitement l'ensemble de  $n + 1$  vecteurs décrivant un  $n$ -simplexe régulier.

La construction commence avec une matrice  $n \times (n + 1)$  non initialisée  $\mathbf{S}$ , où chaque vecteur simplexe est représenté par une colonne

$$\mathbf{S} = \begin{bmatrix} s_0^{(0)} & s_0^{(1)} & s_0^{(2)} & \dots & s_0^{(n)} \\ s_1^{(0)} & s_1^{(1)} & s_1^{(2)} & \dots & s_1^{(n)} \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ s_{n-1}^{(0)} & s_{n-1}^{(1)} & s_{n-1}^{(2)} & \dots & s_{n-1}^{(n)} \end{bmatrix}. \quad (\text{C.2.1})$$

On fixe arbitrairement le premier vecteur (première colonne de  $\mathbf{S}$ ) de la façon la plus simple : toutes ses composantes sont initialisées à 0 sauf la première, qui doit prendre la valeur  $\sqrt{1 - 1/(n+1)}$  afin de satisfaire la propriété (C.1.1a) :

$$\mathbf{S} = \begin{bmatrix} \sqrt{1 - \frac{1}{n+1}} & s_0^{(1)} & s_0^{(2)} & \dots & s_0^{(n)} \\ 0 & s_1^{(1)} & s_1^{(2)} & \dots & s_1^{(n)} \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ 0 & s_{n-1}^{(1)} & s_{n-1}^{(2)} & \dots & s_{n-1}^{(n)} \end{bmatrix}. \quad (\text{C.2.2})$$



La propriété (C.1.1b) impose alors toutes les autres composantes  $s_0^{(\alpha)}$  car seules ces dernières peuvent intervenir dans les produits scalaires  $\vec{s}_1 \cdot \vec{s}_\alpha$

$$\mathbf{s} = \begin{bmatrix} \sqrt{1 - \frac{1}{n+1}} & -\frac{1}{n}\sqrt{1 - \frac{1}{n+1}} & -\frac{1}{n}\sqrt{1 - \frac{1}{n+1}} & \cdots & -\frac{1}{n}\sqrt{1 - \frac{1}{n+1}} \\ 0 & s_1^{(1)} & s_1^{(2)} & \cdots & s_1^{(n)} \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ 0 & s_{n-1}^{(1)} & s_{n-1}^{(2)} & \cdots & s_{n-1}^{(n)} \end{bmatrix}. \quad (\text{C.2.3})$$

On fixe ensuite  $s_i^{(1)} = 0$  pour  $i = 2, 3, \dots, n - 1$ , puis on utilise la propriété (C.1.1a) pour fixer la valeur de  $s_1^{(1)}$

$$\mathbf{s} = \begin{bmatrix} \sqrt{1 - \frac{1}{n+1}} & -\frac{1}{n}\sqrt{1 - \frac{1}{n+1}} & -\frac{1}{n}\sqrt{1 - \frac{1}{n+1}} & \cdots & -\frac{1}{n}\sqrt{1 - \frac{1}{n+1}} \\ 0 & \sqrt{\frac{1}{n+1} \left[ n - \frac{1}{n} \right]} & s_1^{(2)} & \cdots & s_1^{(n)} \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & s_{n-1}^{(2)} & \cdots & s_{n-1}^{(n)} \end{bmatrix}. \quad (\text{C.2.4})$$

Cette procédure itérative peut être répétée indéfiniment afin de construire un ensemble de vecteurs simplexes de dimension arbitraire. Elle fonctionne bien numériquement et est peu sensible aux erreurs d'arrondi, mais n'offre pas de compréhension analytique complète des vecteurs simplexes.

### C.3 $n$ -simplexes et figure de sommet de l'hypercube $(n + 1)$ -dimensionnel

Il est connu que les  $n$ -simplexes réguliers forment les figures de sommet de l'hypercube  $(n + 1)$ -dimensionnel. Une figure de sommet est l'objet géométrique délimité par les  $n$  sommets adjacents à un sommet donné d'un polytope (polyèdre de dimension arbitraire)  $(n + 1)$ -dimensionnel. Lorsque le polytope est un hypercube, tous les choix de sommet de base mènent à des figures de sommet géométriquement identiques, les simplexes (Fig.C.2).

Il est donc possible d'obtenir les vecteurs simplexes analytiquement en orientant la figure de sommet de l'hypercube de dimension  $n + 1$  perpendiculairement à un des axes de l'espace  $\mathbb{R}^{n+1}$ , à l'aide d'une série de rotations. Les vecteurs simplexes sont extraits en projetant la figure de sommet dans le sous-espace  $\mathbb{R}^n$ . Une représentation graphique en  $n + 1 = 3$  dimensions est présentée à la figure C.3.

### C.4 Rotation dans un espace de dimension arbitraire

On doit calculer la rotation permettant d'orienter un simplexe régulier correctement dans un espace de dimension arbitraire. On utilise donc des observations tirée de  $\mathbb{R}^3$  afin de se guider dans notre définition de la rotation  $n$ -dimensionnelle.

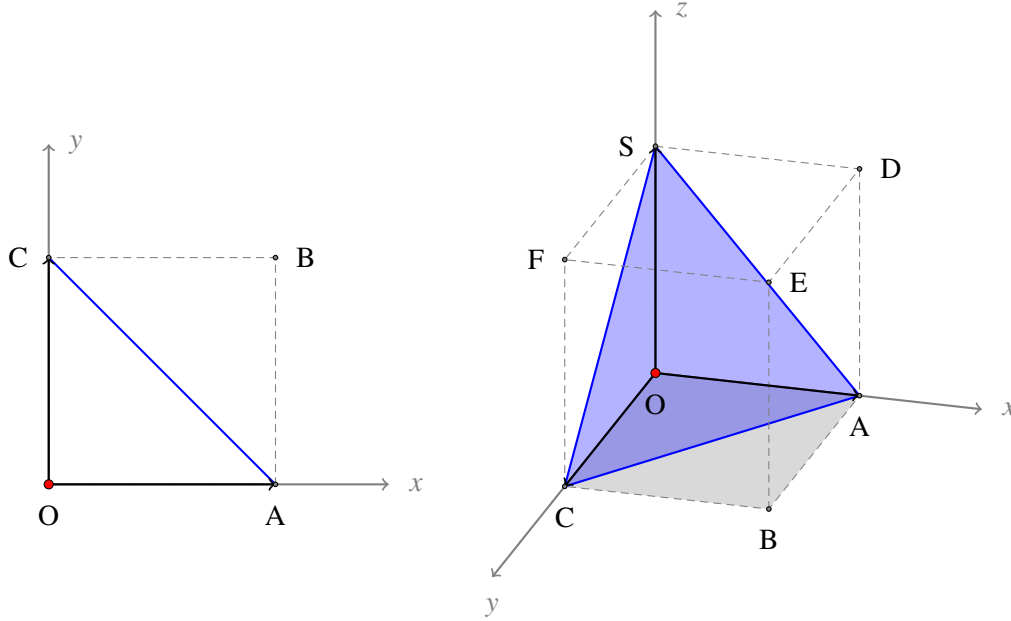


FIGURE C.2 – **Figures de sommet d'hypercube de basse dimension.** (gauche) Hypercube en  $n + 1 = 2$  dimensions (carré). La droite A-C est un simplexe de dimension  $n = 1$ . (droite) Hypercube en  $n + 1 = 3$  dimensions (cube). Le triangle équilatéral A-C-S est un simplexe de dimension  $n = 2$ .

Une rotation dans  $\mathbb{R}^3$  peut être décomposée en une série de rotations de  $\{\theta_p\}$  radians dans les plans  $p = xy, yz, zx$ . On exprime ces rotations de base via les matrices  $\mathbf{r}_{p,\theta}^{(3)}$  :

$$\mathbf{r}_{xy,\theta}^{(3)} = \begin{bmatrix} \cos \theta & -\sin \theta & 0 \\ \sin \theta & \cos \theta & 0 \\ 0 & 0 & 1 \end{bmatrix} \quad \mathbf{r}_{yz,\theta}^{(3)} = \begin{bmatrix} 1 & 0 & 0 \\ 0 & \cos \theta & -\sin \theta \\ 0 & \sin \theta & \cos \theta \end{bmatrix} \quad \mathbf{r}_{zx,\theta}^{(3)} = \begin{bmatrix} \cos \theta & 0 & \sin \theta \\ 0 & 1 & 0 \\ -\sin \theta & 0 & \cos \theta \end{bmatrix}, \quad (\text{C.4.1})$$

tel que  $\mathbf{R}(\{\theta_p\}) := \prod_p \mathbf{r}_{p,\theta_p}^{(3)}$  décrit une rotation générale<sup>1</sup>.

Les contraintes motivant ce choix de matrices de rotation de bases sont :

La rotation n'affecte pas les volumes  $\det \left[ \mathbf{r}_{p,\theta}^{(3)} \right] = 1$  (C.4.2a)

La rotation est une transformation orthogonale  $(\mathbf{r}_{p,\theta}^{(3)})^T \mathbf{r}_{p,\theta}^{(3)} = \mathbf{I} \implies (\mathbf{r}_{p,\theta}^{(3)})^T = (\mathbf{r}_{p,\theta}^{(3)})^{-1}$  (C.4.2b)

La rotation est une opération périodique  $\mathbf{r}_{p,0}^{(3)} = \mathbf{r}_{p,2\pi}^{(3)} = \mathbf{I}$  (C.4.2c)

Ne dépend pas de l'ordre pour un même plan  $[\mathbf{r}_{p,\theta}^{(3)}, \mathbf{r}_{p,\phi}^{(3)}] = \mathbf{0}$  (C.4.2d)

Dépend de l'ordre pour des plans différents  $[\mathbf{r}_{p,\theta}^{(3)}, \mathbf{r}_{p',\phi}^{(3)}] \neq \mathbf{0}$  (en général) (C.4.2e)

1. L'ordre du produit a été choisie arbitrairement ; une rotation peut être décomposée dans n'importe quelle ordre, avec des angles différents.

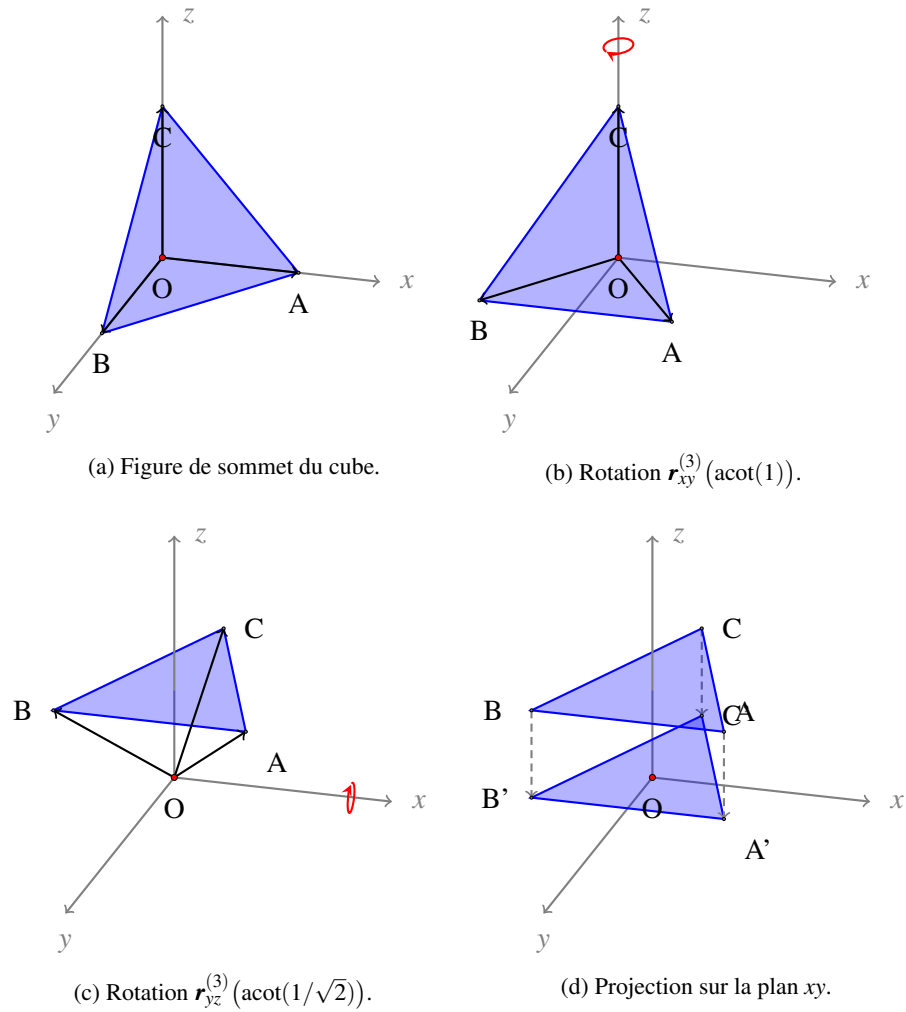


FIGURE C.3 – **Orientation d’une figure de sommet en dimension 3.** (a) Figure de sommet de l’hypercube de dimension 3 dans son orientation initiale. (b) Rotation dans le plan  $xy$ . Elle place les sommets  $A$  et  $B$  dans le même plan  $y$ . (c) Rotation dans le plan  $yz$  plaçant les sommets  $A, B$  et  $C$  dans le même plan  $z$ . (d) Projection finale vers l’espace de dimension 2.

Les matrices  $(n+1) \times (n+1)$

$$\mathbf{r}_{p,\theta}^{(n+1)} = \begin{bmatrix} 1 & 0 & \dots & 0 & 0 & \dots & 0 & 0 \\ 0 & 1 & \dots & 0 & 0 & \dots & 0 & 0 \\ \vdots & \vdots & \ddots & \vdots & \vdots & \ddots & \vdots & \vdots \\ 0 & 0 & \dots & \cos\theta & -\sin\theta & \dots & 0 & 0 \\ 0 & 0 & \dots & \sin\theta & \cos\theta & \dots & 0 & 0 \\ \vdots & \vdots & \ddots & \vdots & \vdots & \ddots & \vdots & \vdots \\ 0 & 0 & \dots & 0 & 0 & \dots & 1 & 0 \\ 0 & 0 & \dots & 0 & 0 & \dots & 0 & 1 \end{bmatrix} \leftarrow \text{rangée } p \tag{C.4.3}$$

$$\mathbf{r}_{n+1,\theta}^{(n+1)} = \begin{bmatrix} \cos\theta & 0 & \dots & 0 & 0 & \dots & 0 & \sin\theta \\ 0 & 1 & \dots & 0 & 0 & \dots & 0 & 0 \\ \vdots & \vdots & \ddots & \vdots & \vdots & \ddots & \vdots & \vdots \\ 0 & 0 & \dots & 1 & 0 & \dots & 0 & 0 \\ 0 & 0 & \dots & 0 & 1 & \dots & 0 & 0 \\ \vdots & \vdots & \ddots & \vdots & \vdots & \ddots & \vdots & \vdots \\ 0 & 0 & \dots & 0 & 0 & \dots & 1 & 0 \\ -\sin\theta & 0 & \dots & 0 & 0 & \dots & 0 & \cos\theta \end{bmatrix} \quad (\text{C.4.4})$$

sont une généralisation naturelle des matrices de rotation de base dans  $\mathbb{R}^3$ . L'intuition ne s'appliquant plus ici, on notera les plans simplement par l'indice de la première des deux composantes affectées par cette rotation.

Notons que toutes les matrices  $\mathbf{r}_{p,\theta}^{(n+1)}$  peuvent être obtenues à l'aide de permutations des colonnes et rangées de n'importe quelle matrice de l'ensemble. Ces permutations peuvent être exprimées à l'aide des matrices de permutations de base

$$\mathbf{Q}_- = \begin{bmatrix} 0 & 0 & 0 & \dots & 0 & 1 \\ 1 & 0 & 0 & \dots & 0 & 0 \\ 0 & 1 & 0 & \dots & 0 & 0 \\ \vdots & \vdots & \vdots & \ddots & \vdots & \vdots \\ 0 & 0 & 0 & \dots & 0 & 0 \\ 0 & 0 & 0 & \dots & 1 & 0 \end{bmatrix} \quad \mathbf{Q}_+ = \begin{bmatrix} 0 & 1 & 0 & \dots & 0 & 0 \\ 0 & 0 & 1 & \dots & 0 & 0 \\ 0 & 0 & 0 & \dots & 0 & 0 \\ \vdots & \vdots & \vdots & \ddots & \vdots & \vdots \\ 0 & 0 & 0 & \dots & 0 & 1 \\ 1 & 0 & 0 & \dots & 0 & 0 \end{bmatrix} \quad (\text{C.4.5})$$

qui ont, entre autres, les propriétés

$$\mathbf{Q}_+^T = \mathbf{Q}_- \quad \mathbf{Q}_- \mathbf{Q}_+ = \mathbf{I} \quad \det[\mathbf{Q}_-] = \det[\mathbf{Q}_+] = 1 \quad (\text{C.4.6})$$

En général, on peut donc exprimer  $\mathbf{r}_{p,\theta}^{(n+1)}$  comme

$$\mathbf{r}_{p,\theta}^{(n+1)} = \mathbf{Q}_-^{p-1} \mathbf{r}_{1,\theta}^{(n+1)} \mathbf{Q}_+^{p-1}. \quad (\text{C.4.7})$$

Cette paramétrisation est suffisante pour prouver que les matrices  $\mathbf{r}_{p,\theta}^{(n+1)}$  respectent les contraintes<sup>2</sup> (C.4.2).

En effet, le déterminant de  $\mathbf{r}_{1,\theta}$  est trivialement  $\cos^2\theta + \sin^2\theta = 1$ , de sorte que pour  $p, \theta$  quelconque

$$\det[\mathbf{r}_{p,\theta}] = \det[\mathbf{Q}_-^{p-1} \mathbf{r}_{1,\theta} \mathbf{Q}_+^{p-1}] = \det[\mathbf{Q}_-]^{p-1} \det[\mathbf{r}_{1,\theta}] \det[\mathbf{Q}_+]^{p-1} = (\cos^2\theta + \sin^2\theta) = 1. \quad (\text{C.4.8})$$

Similairement,  $\mathbf{r}_{1,\theta}^T \mathbf{r}_{1,\theta}$  est trivialement égale  $\mathbf{I}$ . Ainsi, pour  $p, \theta$  arbitraires

$$\mathbf{r}_{p,\theta}^T \mathbf{r}_{p,\theta} = \mathbf{Q}_-^{p-1} \mathbf{r}_{1,\theta}^T \mathbf{Q}_+^{p-1} \mathbf{Q}_-^{p-1} \mathbf{r}_{1,\theta} \mathbf{Q}_+^{p-1} = \mathbf{Q}_-^{p-1} \mathbf{r}_{1,\theta}^T \mathbf{r}_{1,\theta} \mathbf{Q}_+^{p-1} = \mathbf{Q}_-^{p-1} \mathbf{Q}_+^{p-1} = \mathbf{I} \quad (\text{C.4.9})$$

2. Afin de clarifier le traitement analytique, nous omettons la dimension  $n+1$  de la rotation dans la notation.

Les éléments diagonaux des matrices  $\mathbf{r}_{p,\theta}$  sont soit 1, soit  $\cos\theta$ . Les éléments hors diagonale sont soit 0, soit  $\sin\theta$ . Il est donc clair que la condition de périodicité est effectivement généralisée.

Les commutateurs pour  $\mathbf{r}_{1,\theta}$  sont

$$[\mathbf{r}_{1,\theta}, \mathbf{r}_{1,\varphi}] = \mathbf{0} \quad \text{trivialement} \quad (\text{C.4.10a})$$

$$[\mathbf{r}_{1,\theta}, \mathbf{r}_{2,\varphi}] = \begin{bmatrix} 0 & \sin\theta(1 - \cos\varphi) & \sin\theta\sin\varphi & 0 & \dots & 0 \\ \sin\theta(1 - \cos\varphi) & 0 & \sin\varphi(1 - \cos\theta) & 0 & \dots & 0 \\ -\sin\theta\sin\varphi & \sin\varphi(1 - \cos\theta) & 0 & 0 & \dots & 0 \\ 0 & 0 & 0 & 0 & \dots & 0 \\ \vdots & \vdots & \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & 0 & 0 & \dots & 0 \end{bmatrix} \quad (\text{C.4.10b})$$

$$[\mathbf{r}_{1,\theta}, \mathbf{r}_{p,\varphi}] = \mathbf{0} \quad \forall p > 1 \quad (\text{C.4.10c})$$

Le deuxième commutateur n'est pas nul, sauf dans les cas triviaux (rotation de  $2\pi$ ) et pour deux rotations de  $(2m - 1)\pi$  radians, où  $m \in \mathbb{N}$ .

Tous les autres commutateurs peuvent être construits à l'aide de matrices de permutations, car, pour  $p < p'$ , on a<sup>3</sup>

$$[\mathbf{r}_{p,\theta}, \mathbf{r}_{p',\varphi}] = \mathbf{Q}_-^{p-1} \mathbf{r}_{1,\theta} \mathbf{Q}_-^{p'-p} \mathbf{r}_{1,\varphi} \mathbf{Q}_+^{p'-1} - \mathbf{Q}_-^{p'-1} \mathbf{r}_{1,\varphi} \mathbf{Q}_+^{p'-p} \mathbf{r}_{1,\theta} \mathbf{Q}_+^{p-1} = \mathbf{Q}_-^{p-1} [\mathbf{r}_{1,\theta}, \mathbf{r}_{1+p'-p,\varphi}] \mathbf{Q}_+^{p-1} \quad (\text{C.4.11})$$

Ce qui complète notre vérification. On accepte donc  $\{\mathbf{r}_{p,\theta}^{(n+1)}\}_{p=1,\dots,n}$  comme *définitions* des rotations de base en  $n$ -dimensions.

## C.5 Projection de la figure de sommet d'un hypercube $(n + 1)$ -dimensionnel

Tous les éléments sont maintenant en place pour calculer la rotation nécessaire. L'objectif est de placer la figure de sommet perpendiculairement à un des axes  $a$  de  $\mathbb{R}^{n+1}$ . Par perpendiculaire à l'axe  $a$ , on veut dire que le  $a^{\text{ème}}$  élément de tous les vecteurs pointant vers les coins de cette figure prennent la même valeur. On choisit l'origine comme sommet de référence, de sorte que ces vecteurs sont les vecteurs orthogonaux  $\{\vec{e}_i\}_{i=1,\dots,n+1}$ . Afin d'avoir des vecteurs identiques à ceux obtenus à la Sec.C.2, on cherche à orienter la figure de sommet perpendiculairement à l'axe  $a = 1$ .

Initialement, la figure de sommet est représentée par un matrice identité  $(n + 1) \times (n + 1)$ , où

3. Les commutateurs inverses sont simplement obtenus à l'aide de la relation  $[\mathbf{r}_{p,\theta}, \mathbf{r}_{p',\varphi}] = -[\mathbf{r}_{p',\varphi}, \mathbf{r}_{p,\theta}]$ .

chaque colonne correspond à un vecteur orthogonal différent

$$\mathbf{S} = \begin{bmatrix} 1 & 0 & 0 & 0 & \dots & 0 \\ 0 & 1 & 0 & 0 & \dots & 0 \\ 0 & 0 & 1 & 0 & \dots & 0 \\ 0 & 0 & 0 & 1 & \dots & 0 \\ \vdots & \vdots & \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & 0 & 0 & \dots & 1 \end{bmatrix} \quad (\text{C.5.1})$$

La procédure explicite consiste à doter les 2 derniers vecteurs d'un même  $n^{\text{ème}}$  élément via une rotation dans le plan  $n$ , puis à doter les 3 derniers vecteurs d'un même  $(n-1)^{\text{ème}}$  élément à l'aide d'une rotation dans le plan  $n-1$ , etc.

Les premières étapes sont

$$\mathbf{S}' = \mathbf{r}_{n,\theta_n}^{(n+1)} = \begin{bmatrix} 1 & \dots & 0 & 0 & 0 & 0 \\ \vdots & \ddots & \vdots & \vdots & \vdots & \vdots \\ 0 & \dots & 1 & 0 & 0 & 0 \\ 0 & \dots & 0 & 1 & 0 & 0 \\ 0 & \dots & 0 & 0 & \cos \theta_n & -\sin \theta_n \\ 0 & \dots & 0 & 0 & \sin \theta_n & \cos \theta_n \end{bmatrix}$$

$$\mathbf{S}'' = \mathbf{r}_{n-1,\theta_{n-1}}^{(n+1)} \mathbf{r}_{n,\theta_n}^{(n+1)} = \begin{bmatrix} 1 & \dots & 0 & 0 & 0 & 0 \\ \vdots & \ddots & \vdots & \vdots & \vdots & \vdots \\ 0 & \dots & 1 & 0 & 0 & 0 \\ 0 & \dots & 0 & \cos \theta_{n-1} & -\cos \theta_n \sin \theta_{n-1} & \sin \theta_n \sin \theta_{n-1} \\ 0 & \dots & 0 & \sin \theta_{n-1} & \cos \theta_n \cos \theta_{n-1} & -\sin \theta_n \cos \theta_{n-1} \\ 0 & \dots & 0 & 0 & \sin \theta_n & \cos \theta_n \end{bmatrix}$$

$$\mathbf{S}''' = \mathbf{r}_{n-2, \theta_{n-2}}^{(n+1)} \mathbf{r}_{n-1, \theta_{n-1}}^{(n+1)} \mathbf{r}_{n, \theta_n}^{(n+1)} = \begin{bmatrix} 1 & \dots & 0 & 0 & 0 & 0 \\ \vdots & \ddots & \vdots & \vdots & \vdots & \vdots \\ 0 & \dots & \cos \theta_{n-2} & -\cos \theta_{n-1} \sin \theta_{n-2} & \cos \theta_n \sin \theta_{n-1} \sin \theta_{n-2} & -\sin \theta_n \sin \theta_{n-1} \sin \theta_{n-2} \\ 0 & \dots & \sin \theta_{n-2} & \cos \theta_{n-1} \cos \theta_{n-2} & -\cos \theta_n \sin \theta_{n-1} \cos \theta_{n-2} & \sin \theta_n \sin \theta_{n-1} \cos \theta_{n-2} \\ 0 & \dots & 0 & \sin \theta_{n-1} & \cos \theta_n \cos \theta_{n-1} & -\sin \theta_n \cos \theta_{n-1} \\ 0 & \dots & 0 & 0 & \sin \theta_n & \cos \theta_n \end{bmatrix}.$$

On fixe les angles  $\{\theta_p\}_{p=n, n-1, n-2}$  par inspection.

L'angle  $\theta_n$  doit solutionner  $\cos \theta_n = -\sin \theta_n$ ; on a donc  $\theta_n = -\frac{\pi}{4}$ .

À l'étape suivante, l'égalité  $\cos \theta_n = -\sin \theta_n$  est toujours satisfaite par construction, et seule l'équation

$$\cos \theta_{n-1} = \sin \theta_n \sin \theta_{n-1} \quad (\text{C.5.2})$$

présente une réelle contrainte. Ceci fixe donc l'angle  $\theta_{n-1} = \text{acot}(\sin \theta_n)$ .

Pour  $\mathbf{S}'''$ , on a à nouveau que  $\cos \theta_n = -\sin \theta_n$  et  $\cos \theta_{n-1} = \sin \theta_n \sin \theta_{n-1}$  par construction, et

$$\cos \theta_{n-2} = -\sin \theta_n \sin \theta_{n-1} \sin \theta_{n-2} \quad (\text{C.5.3})$$

est une équation non triviale, tel que  $\theta_{n-2} = \text{acot}(-\sin \theta_n \sin \theta_{n-1})$ .

Il est clair que ce type d'équations apparaît à chaque rotation additionnelle. On pose donc la forme générale

$$\theta_j = \text{acot} \left( (-1)^{n-j+1} \prod_{i=j+1}^n \sin \theta_i \right) \quad \theta_n = \text{acot}(-1) = -\frac{\pi}{4} \quad (\text{C.5.4})$$

La figure de sommet  $n$ -dimensionnelle étant correctement orientée dans  $\mathbb{R}^{n+1}$  lorsque la première rangée est identique pour toutes les colonnes, on doit appliquer  $n$  rotations. La rotation complète est donc donnée par

$$\mathbf{R}(\{\theta_p\}_{p=1, \dots, n}) = \prod_{j=1}^n \mathbf{r}_{j, \theta_j}^{(n+1)} \quad (\text{C.5.5})$$

Il suffit alors de déplacer la figure de sommet le long du dernier axe, ou simplement de la projeter dans l'espace de dimension  $n$ . Les vecteurs qui pointaient vers les coins adjacents à l'origine de l'hypercube de dimension  $n+1$  peuvent ainsi être directement interprétés comme des vecteurs pointant vers les coins du simplexe de dimension  $n$ .

Pour une matrice  $n \times (n+1)$  de projection

$$\mathbf{P} = \begin{bmatrix} 0 & 1 & 0 & \dots & 0 \\ 0 & 0 & 1 & \dots & 0 \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & 0 & \dots & 1 \end{bmatrix}, \quad (\text{C.5.6})$$

la matrice colonne des  $n + 1$  vecteurs simplexes en  $n$  dimension est donc

$$\mathbf{S} = \mathbf{p} \left( \prod_{j=1}^n \mathbf{r}_{j,\theta_j}^{(n+1)} \right) \mathbf{I}_{(n+1) \times (n+1)}. \quad (\text{C.5.7})$$

avec  $\{\theta_p\}_{p=1,\dots,n}$  défini à l'Eq. (C.5.4).

On notera que la définition naturelle de la rotation utilisée ici mène à une représentation des simplexes différant de celle obtenue à la section Sec. C.2 par un facteur global de -1. La correspondance parfaite est donc donnée par

$$\mathbf{S} = -\mathbf{p} \left( \prod_{j=1}^n \mathbf{r}_{j,\theta_j}^{(n+1)} \right). \quad (\text{C.5.8})$$

On peut vérifier *a posteriori* que les vecteurs obtenus sont bien des vecteurs simplexes ayant la bonne norme. Les propriétés (C.1.1a)-(C.1.1b) devant être respectées, on doit avoir

$$\mathbf{S}^T \mathbf{S} = \mathbf{I}_{(n+1) \times (n+1)} - \frac{1}{n+1} \mathbf{1} \mathbf{1}^T \quad (\text{C.5.9})$$

La substitution de cette hypothèse dans l'Eq. (C.5.8) mène à

$$\mathbf{S}^T \mathbf{S} = \mathbf{R}^T \mathbf{p}^T \mathbf{p} \mathbf{R} \quad (\text{C.5.10})$$

$$\mathbf{R} \mathbf{S}^T \mathbf{S} \mathbf{R}^T = \mathbf{p}^T \mathbf{p} \quad (\text{C.5.11})$$

$$\mathbf{I}_{(n+1) \times (n+1)} = \mathbf{p}^T \mathbf{p} + \frac{1}{n+1} \mathbf{R} \mathbf{1} \mathbf{1}^T \mathbf{R}^T, \quad (\text{C.5.12})$$

i.e. on doit avoir

$$\frac{1}{n+1} \mathbf{R} \mathbf{1} \mathbf{1}^T \mathbf{R}^T = \begin{bmatrix} 0 & 0 & \dots & 0 & 0 \\ 0 & 0 & \dots & 0 & 0 \\ \vdots & \vdots & \ddots & \vdots & \vdots \\ 0 & 0 & \dots & 0 & 0 \\ 0 & 0 & \dots & 0 & 1 \end{bmatrix} \implies \mathbf{R} \mathbf{1} = \sqrt{n+1} \begin{pmatrix} 0 \\ 0 \\ \vdots \\ 0 \\ 1 \end{pmatrix} \quad (\text{C.5.13})$$

pour que les vecteurs soient correctement normés et orientés. C'est effectivement le cas car (a) ce vecteur est de norme  $\sqrt{\mathbf{1}^T \mathbf{1}} = \sqrt{n+1}$ , (b) la rotation n'influence pas la norme et (c) la série d'angles (C.5.4) est sélectionnée de façon à rendre chaque élément du vecteur  $\mathbf{1}$  nul, sauf l'élément  $n + 1$ .

On notera que le produit inverse est trivialement

$$\mathbf{S} \mathbf{S}^T = \mathbf{p} \mathbf{R} \mathbf{R}^T \mathbf{p}^T = \mathbf{p} \mathbf{p}^T = \mathbf{I}_{n \times n}. \quad (\text{C.5.14})$$



# Bibliographie

- [1] AHN, Y.-Y. AND BAGROW, J.P. AND LEHMANN, S., *Link communities reveal multiscale complexity in networks*, Nature, 466 (2010), p. 761.
- [2] AHN, Y.Y. AND HAN, S. AND KWAK, H. AND MOON, S. AND JEONG, H., *Analysis of topological characteristics of huge online social networking services*, in Proceedings of the 16th international conference on World Wide Web, 2007, p. 835.
- [3] ALLARD, A. AND HÉBERT-DUFRESNE, L. AND NOËL, P.-A. AND MARCEAU, V. AND DUBÉ, L.J., *Bond percolation on a class of correlated and clustered random graphs*, J. Phys. A, 45 (2012), p. 405005.
- [4] ALLARD, A. AND NOËL, P.-A. AND DUBÉ, L.J., *Des Ponts d'Euler à la Grippe Aviaire : De l'abstraction mathématique à la réalité sociale des épidémies*, Accromath, 4 (2009), p. 24.
- [5] ALPERT, C.J. AND YAO, S.-Z., *Spectral partitioning : the more eigenvectors, the better*, in Proceedings of the 32nd annual ACM/IEEE Design Automation Conference, 1995, p. 195.
- [6] ALZATE, C. AND SUYKENS, J.A.K., *Multiway spectral clustering with out-of-sample extensions through weighted kernel PCA*, IEEE Transactions on Pattern Analysis and Machine Intelligence, 32 (2010), p. 335.
- [7] ANDERSON, P.W., *More is different*, Science, 177 (1972), p. 393.
- [8] AROUS, G.B. AND GUIONNET, A., *The Oxford Handbook on Random Matrix Theory*, Oxford University Press, 2008.
- [9] BACH, F.R. AND JORDAN, M.I., *Learning spectral clustering*, Advances in Neural Information Processing Systems, 16 (2004), p. 305.
- [10] BALL, B. AND KARRER, B. AND NEWMAN, M.E.J., *Efficient and principled method for detecting communities in networks*, Phys. Rev. E, 84 (2011), p. 36103.
- [11] BARABÁSI, A.-L., *The network takeover*, Nature Physics, 8 (2011), p. 14.
- [12] BARABÁSI, A.-L. AND ALBERT, R., *Emergence of scaling in random networks*, Science, 286 (1999), p. 509.

- [13] BARABÁSI, A.-L. AND MARTINO, M. AND PÓSFAL, M., *Network Science*, Barabási Lab, 2014.
- [14] BARNES, EARL R, *An algorithm for partitioning the nodes of graphs*, SIAM J. Alg. Disc. Disc. Meth., 3 (1982), p. 541.
- [15] BIGGS, N. AND LLOYD, E. AND WILSON, R., *Graph Theory, 1736–1936*, Oxford University Press, 1986.
- [16] BLONDEL, V.D. AND GUILLAUME, J.-L. AND LAMBIOTTE, R. AND LEFEBVRE, E., *Fast unfolding of communities in large networks*, Journal of Statistical Mechanics : Theory and Experiment, 10 (2008), p. 10008.
- [17] BOGUÑÁ, M. AND PASTOR-SATORRAS, R. AND DÍAZ-GUILERA, A. AND ARENAS, A., *Models of social networks based on social distance attachment*, Phys. Rev. E, 70 (2004), p. 56122.
- [18] BORGATTI, S.P. AND MEHRA, A. AND BRASS, D.J. AND LABIANCA, G., *Network analysis in the social sciences*, Science, 323 (2009), p. 892.
- [19] BRÉZIN, E. AND ITZYKSON, C. AND PARISI, G. AND ZUBER, J.-B., *Planar diagrams*, Communications in Mathematical Physics, 59 (1978), p. 35.
- [20] CARRINGTON, P.J. AND SCOTT, J. AND WASSERMAN, S., *Models and Methods in Social Network Analysis*, Cambridge University Press, 2005.
- [21] CATTELL, R., *A note on correlation clusters and cluster search methods*, Psychometrika, 9 (1944), p. 169.
- [22] CHAUHAN, S. AND GIRVAN, M. AND OTT, E., *Spectral properties of networks with community structure*, Phys. Rev. E, 80 (2009), p. 056114.
- [23] CHO, E. AND MYERS, S.A. AND LESKOVEC, J., *Friendship and mobility : user movement in location-based social networks*, in Proceedings of the 17th ACM SIGKDD international conference on Knowledge discovery and data mining, ACM, 2011, p. 1082.
- [24] CLAUSET, A. AND MOORE, C. AND NEWMAN, M.E.J., *Hierarchical structure and the prediction of missing links in networks*, Nature, 453 (2008), p. 98.
- [25] CLAUSET, A. AND NEWMAN, M.E.J. AND MOORE, C., *Finding community structure in very large networks*, Phys. Rev. E, 70 (2004), p. 066111.
- [26] CLAUSET, A. AND SHALIZI, C.R. AND NEWMAN, M.E.J., *Power-law distributions in empirical data*, SIAM review, 51 (2009), p. 661.
- [27] CONDON, A. AND KARP, R.M., *Algorithms for graph partitioning on the planted partition model*, Random Structures and Algorithms, 18 (2001), p. 116.
- [28] CORNELIUS, L., *An iteration method for the solution of the eigenvalue problem of linear differential and integral operators*, NIST Journal of Research, 45 (1950), p. 255.

- 
- [29] CRAMPES, M. AND PLANTIÉ, M., *A unified community detection, visualization and analysis method*, Advances in Complex Systems, 17 (2014).
- [30] CRESCENZI, P. AND KANN, V., *A compendium of NP optimization problems*, Università di Roma “La Sapienza”, 1995.
- [31] DE RUITER, J. AND WESTON, G. AND LYON, S.M., *Dunbar’s Number : Group Size and Brain Physiology in Humans Reexamined*, American Anthropologist, 113 (2011), p. 557.
- [32] DECELLE, A. AND KRZAKALA, F. AND MOORE, C. AND ZDEBOROVÁ, L., *Inference and phase transitions in the detection of modules in sparse networks*, Phys. Rev. Lett., 107 (2011), p. 065701.
- [33] DERÉNYI, I. AND PALLA, G. AND VICSEK, T., *Clique percolation in random networks*, Phys. Rev. Lett., 94 (2005), p. 160202.
- [34] DUNBAR, R., *Neocortex Size as a Constraint on Group Size in Primates*, Journal of Human Evolution, 22 (1992), p. 469.
- [35] ERDÖS, P. AND RÉNYI, A., *On random graphs, I*, Publicationes Mathematicae, 6 (1959), p. 290.
- [36] EVANS, T.S. AND LAMBIOTTE, R., *Line graphs, link partitions, and overlapping communities*, Phys. Rev. E, 80 (2009), p. 016105.
- [37] FORRESTER, P.J., *Log-gases and Random Matrices (LMS-34)*, Princeton University Press, 2010.
- [38] FORTUNATO, S., *Community detection in graphs*, Physics Reports, 486 (2010), pp. 75–174.
- [39] FORTUNATO, S. AND BARTHÉLEMY, M., *Resolution limit in community detection*, Proc. Natl. Acad. Sci. U.S.A., 104 (2007), p. 36.
- [40] GIRVAN, M. AND NEWMAN, M.E.J., *Community structure in social and biological networks.*, Proc. Natl. Acad. Sci. U.S.A., 99 (2002), p. 7821.
- [41] GONÇALVES, B. AND PERRA, N. AND VESPIGNANI, A., *Modeling Users’ Activity on Twitter Networks : Validation of Dunbar’s Number*, PLoS ONE, 6 (2011), p. e22656.
- [42] GUIMERÀ, R. AND DANON, L. AND DÍAZ-GUILERA, A. AND GIRALT, F. AND ARENAS, A., *Self-similar community structure in a network of human interactions*, Phys. Rev. E, 68 (2003), p. 065103.
- [43] HE, D. AND LIU, J. AND LIU, D. AND JIN, D. AND JIA, Z., *Ant colony optimization for community detection in large-scale complex networks*, in 2011 Seventh International Conference on Natural Computation, 2011, p. 1151.
- [44] HÉBERT-DUFRESNE, L. AND ALLARD, A. AND MARCEAU, V. AND NOËL, P.-A. AND DUBÉ, L.J., *Structural preferential attachment : network organization beyond the link*, Phys. Rev. Lett., 107 (2011), p. 1158702.

- [45] ———, *Structural preferential attachment : Stochastic process for the growth of scale-free, modular, and self-similar systems*, Phys. Rev. E, 85 (2012), p. 026108.
- [46] HÉBERT-DUFRESNE, L. AND ALLARD, A. AND YOUNG, J.-G. AND DUBÉ, L.J., *Global efficiency of local immunization on complex networks*, Scientific reports, 3 (2013), p. 2171.
- [47] ———, *Percolation on random networks with arbitrary  $k$ -core structure*, Phys. Rev. E, 88 (2013), p. 062820.
- [48] ———, *Universal growth constraints of human systems*, arXiv :1310.0112, (2013).
- [49] HÉBERT-DUFRESNE, L. AND NOËL, P.-A. AND MARCEAU, V. AND ALLARD, A. AND DUBÉ, L.J., *Propagation dynamics on networks featuring complex topologies*, Phys. Rev. E, 82 (2010), p. 036115.
- [50] HOLLAND, P.W. AND LEINHARDT, S., *Transitivity in structural models of small groups.*, Comparative Group Studies, (1971).
- [51] HRIC, D. AND DARST, R.K. AND FORTUNATO, S., *Community detection in networks : structural clusters versus ground truth*, arXiv :1406.0146, (2014).
- [52] KARRER, B. AND NEWMAN, M.E.J., *Random graphs containing arbitrary distributions of subgraphs*, Phys. Rev. E, 82 (2010), p. 066118.
- [53] KAWAMOTO, T. AND ROSVALL, M., *The map equation and the resolution limit in community detection*, arXiv :1402.4385, (2014).
- [54] KERNIGHAN, B.W. AND LIN, S., *An efficient heuristic procedure for partitioning graphs*, Bell System Technical Journal, 49 (1970), p. 291.
- [55] KLIMT, B. AND YANG, Y., *Introducing the Enron Corpus.*, in CEAS, 2004.
- [56] KUHN, T.S., *The Structure of Scientific Revolutions*, University of Chicago Press, 1962.
- [57] LAARHOVEN, T.V. AND MARCHIORI, E., *Graph clustering with local search optimization : The resolution bias of the objective function matters most*, Phys. Rev. E, 87 (2013), p. 012812.
- [58] LANCICHINETTI, A. AND FORTUNATO, S. AND RADICCHI, F., *Benchmark graphs for testing community detection algorithms*, Phys. Rev. E, 78 (2008), p. 046110.
- [59] LANCICHINETTI, A. AND RADICCHI, F. AND RAMASCO, J.J. AND FORTUNATO, S., *Finding statistically significant communities in networks*, PLoS One, 6 (2011), p. e18961.
- [60] LEE, C. AND REID, F. AND MCDAID, A. AND HURLEY, N., *Detecting Highly Overlapping Community Structure by Greedy Clique Expansion*, arXiv :1002.1827v2, (2010).
- [61] LESKOVEC, J. AND LANG, K.J. AND DASGUPTA, A. AND MAHONEY, M.W., *Community structure in large networks : Natural cluster sizes and the absence of large well-defined clusters*, Internet Mathematics, 6 (2009), p. 29.

- 
- [62] MACQUEEN, J., *Some methods for classification and analysis of multivariate observations*, in Proceedings of the fifth Berkeley symposium on mathematical statistics and probability, 1967, p. 281.
- [63] MACULA, A.J., *Covers of a finite set*, Mathematics Magazine, 67 (1994), p. 141.
- [64] ———, *Lewis Carroll and the enumeration of minimal covers*, Mathematics Magazine, 68 (1995), p. 269.
- [65] MOLLOY, M. AND REED, B., *A critical point for random graphs with a given degree sequence*, Random structures & algorithms, 6 (1995), p. 161.
- [66] MOTTER, A.E. AND ALBERT, R., *Networks in Motion*, Physics Today, 4 (2012), p. 43.
- [67] NADAKUDITI, R.R. AND NEWMAN, M.E.J., *Graph Spectra and the Detectability of Community Structure in Networks*, Phys. Rev. Lett., 108 (2012), p. 188701.
- [68] ———, *Spectra of random graphs with arbitrary expected degrees*, Phys. Rev. E, 87 (2013), p. 012803.
- [69] NEWMAN, M.E.J., *Clustering and preferential attachment in growing networks*, Phys. Rev. E, 64 (2001), p. 025102.
- [70] ———, *Spread of epidemic disease on networks*, Phys. Rev. E, 66 (2002), p. 016128.
- [71] ———, *Analysis of weighted networks*, Phys. Rev. E, 70 (2004), p. 056131.
- [72] ———, *Coauthorship networks and patterns of scientific collaboration*, Proc. Natl. Acad. Sci. U.S.A., 101 (2004), p. 5200.
- [73] ———, *Finding community structure in networks using the eigenvectors of matrices*, Phys. Rev. E, 74 (2006), p. 036104.
- [74] ———, *Modularity and community structure in networks*, Proc. Natl. Acad. Sci. U.S.A., 103 (2006), p. 8577.
- [75] ———, *Networks : An Introduction*, Oxford University Press, 2010.
- [76] ———, *Communities, modules and large-scale structure in networks*, Nature Physics, 8 (2012), p. 25.
- [77] NEWMAN, M.E.J. AND GIRVAN, M., *Finding and evaluating community structure in networks*, Phys. Rev. E, 69 (2004), p. 026113.
- [78] NEWMAN, M.E.J. AND PARK, J., *Why social networks are different from other types of networks*, Phys. Rev. E, 68 (2003), p. 036122.
- [79] NEWMAN, M.E.J. AND STROGATZ, S.H. AND WATTS, D.J., *Random graphs with arbitrary degree distributions and their applications*, Phys. Rev. E, 64 (2001), p. 026118.

- [80] NOËL, P.-A. AND ALLARD, A. AND HÉBERT-DUFRESNE, L. AND MARCEAU, V. AND DUBÉ, L.J., *Spreading dynamics on complex networks : a general stochastic approach*, J. Math. Bio., (2013), p. 1.
- [81] PALLA, G. AND DERÉNYI, I. AND FARKAS, I.J. AND VICSEK, T., *Uncovering the overlapping community structure of complex networks in nature and society*, Nature, 435 (2005), p. 814.
- [82] PALLA, G. AND FARKAS, I.J. AND POLLNER, P. AND DERÉNYI, I. AND VICSEK, T., *Directed network modules*, New Journal of Physics, 9 (2007), p. 186.
- [83] ———, *Fundamental statistical features and self-similar properties of tagged networks*, New Journal of Physics, 10 (2008), p. 123026.
- [84] PEIXOTO, T.P., *Eigenvalue spectra of modular networks*, Phys. Rev. Lett., 111 (2013), p. 098701.
- [85] PETERSEN, K.B., *The matrix Cookbook*, 2008.
- [86] PRESS, W.H. AND FLANNERY, B.P. AND TEUKOLSKY, S.A. AND VETTERLING, W.T., *Numerical Recipes : The Art of Scientific Computing*, Cambridge University Press, 1986.
- [87] RADICCHI, F. AND CASTELLANO, C. AND CECCONI, F. AND LORETO, V. AND PARISI, D., *Defining and identifying communities in networks*, Proc. Natl. Acad. Sci. U.S.A., 101 (2004), p. 2658.
- [88] RAVASZ, E. AND BARABÁSI, A.-L., *Hierarchical organization in complex networks*, Phys. Rev. E, 67 (2003), p. 026112.
- [89] RAVASZ, E. AND SOMERA, A.L. AND MONGRU, D. A. AND OLTVAI, Z.N AND BARABÁSI, A.-L., *Hierarchical Organization of Modularity in Metabolic Networks*, Science, 297 (2002), p. 1551.
- [90] RIOLO, M.A. AND NEWMAN, M.E.J., *First-principles multiway spectral partitioning of graphs*, Journal of Complex Networks, 2 (2014), pp. 124–140.
- [91] RIPEANU, M. AND FOSTER, I., *Mapping the Gnutella Network : Macroscopic Properties of Large-Scale Peer-to-Peer Systems*, in Peer-to-Peer Systems, Springer, 2002, pp. 85–93.
- [92] RONHOVDE, P. AND NUSSINOV, Z., *Local resolution-limit-free Potts model for community detection*, Phys. Rev. E, 81 (2010), p. 046114.
- [93] ROSVALL, M. AND BERGSTROM, C.T., *An information-theoretic framework for resolving community structure in complex networks*, Proc. Natl. Acad. Sci. U.S.A., 104 (2007), p. 7327.
- [94] ———, *Maps of random walks on complex networks reveal community structure*, Proc. Natl. Acad. Sci. U.S.A., 105 (2008), p. 1118.

- 
- [95] ROZENBLATT-ROSEN, O. AND DEO, R.C. AND PADI, M. AND ADELMANT, G. AND CALDERWOOD, M.A. AND ROLLAND, T. AND GRACE, M. AND DRICOT, A. AND ASKENAZI, M. AND TAVARES, M., *Interpreting cancer genomes using systematic host network perturbations by tumour virus proteins*, Nature, 487 (2012), p. 491.
- [96] SESHADHRI, C. AND KOLDA, T. G. AND PINAR, A., *Community structure and scale-free collections of Erdős-Rényi graphs*, Phys. Rev. E, 85 (2012), p. 056109.
- [97] SHEN, H.-W. AND CHENG, X.-Q., *Spectral methods for the detection of network community structure : a comparative analysis*, Journal of Statistical Mechanics : Theory and Experiment, 2010 (2010), p. P10020.
- [98] SIMON, H.A., *A behavioral model of rational choice*, The quarterly journal of economics, (1955), p. 99.
- [99] ———, *The architecture of complexity*, Proceedings of the American Philosophical Society, 106 (1962), p. 467.
- [100] STANLEY, R.P., *Enumerative Combinatorics*, vol. 1, Cambridge University Press, 1997.
- [101] ———, *Enumerative Combinatorics*, vol. 2, Cambridge University Press, 1999.
- [102] STROGATZ, STEVEN H., *Exploring complex networks*, Nature, 410 (2001), p. 268.
- [103] SZU, H. AND HARTLEY, R., *Fast simulated annealing*, Phys. Lett. A, 122 (1987), p. 157.
- [104] TALBI, E.-G., *Metaheuristics : From Design to Implementation*, John Wiley & Sons, 2009.
- [105] TIBÉLY, GERGELY, *Criteria for locally dense subgraphs*, Physica A, 391 (2012), p. 1831.
- [106] TRAAG, V.A. AND DOOREN, P.V. AND NESTEROV, Y., *Narrow scope for resolution-limit-free community detection*, Phys. Rev. E, 84 (2011), p. 016114.
- [107] WANG, X.F. AND CHEN, G., *Complex networks : small-world, scale-free and beyond*, Circuits and Systems Magazine, IEEE, 3 (2003), p. 6.
- [108] WASSERMAN, S., *Social Network Analysis : Methods and Applications*, Cambridge University Press, 1994.
- [109] WATTS, D.J., *Small Worlds : The Dynamics of Networks Between Order and Randomness*, Princeton University Press, 1999.
- [110] WATTS, D.J. AND STROGATZ, S.H., *Collective dynamics of 'small-world' networks*, Nature, 393 (1998), p. 440.
- [111] WEIDLICH, W. AND HAAG, G., *Concepts and Models of a Quantitative Sociology : The Dynamics of Interacting Populations*, Springer, 1982.
- [112] WHITNEY, H., *Congruent Graphs and the Connectivity of Graphs*, American Journal of Mathematics, 54 (1932), p. 150.

- [113] W.W. ZACHARY, *An information flow model for conflict and fission in small groups*, Journal of Anthropological Research, 33 (1977), p. 452.
- [114] XIE, J. AND KELLEY, S. AND SZYMANSKI, B.K., *Overlapping community detection in networks : The state-of-the-art and comparative study*, ACM Computing Surveys, 45 (2013), p. 43.
- [115] YANG, J. AND LESKOVEC, J., *Overlapping community detection at scale*, in Proceedings of the sixth ACM international conference on Web search and data mining - WSDM '13, 2013, p. 587.
- [116] YOUNG, J.-G. AND ALLARD, A. AND HÉBERT-DUFRESNE, L. AND DUBÉ, L.J., *Unveiling Hidden Communities Through Cascading Detection on Network Structures*, arXiv :1211.1364, (2012).
- [117] ZHANG, X. AND NADAKUDITI, R.R. AND NEWMAN, M.E.J., *Spectra of random graphs with community structure and arbitrary degrees*, Phys. Rev. E, 89 (2013), p. 042816.
- [118] ZHANG, Z. AND JORDAN, M.I., *Multiway spectral clustering : A margin-based perspective*, Statistical Science, 23 (2008), p. 383.
- [119] ZVONKIN, A., *Matrix integrals and map enumeration : an accessible introduction*, Mathematical and Computer Modelling, 26 (1997), p. 281.