

Partitionnement d'un réseau de sociabilité à fort coefficient de clustering

Romain Boulet, Bertrand Jouve

Institut de Mathématiques de Toulouse,
Université Toulouse II le Mirail
5 Allées Antonio Machado, 31058 Toulouse Cedex 1.
{boulet,jouve}@univ-tlse2.fr

Résumé. Afin de comparer l'organisation sociale d'une paysannerie médiévale avant et après la guerre de Cent Ans nous étudions la structure de réseaux sociaux construits à partir d'un corpus de contrats agraires. Faibles diamètres et fort clustering révèlent des graphes en petit monde. Comme beaucoup de grands réseaux d'interaction étudiés ces dernières années ces graphes sont sans échelle typique. Les distributions des degrés de leurs sommets sont bien ajustées par une loi de puissance tronquée par une coupure exponentielle. Ils possèdent en outre un club-huppé, c'est à dire un noyau dense et de faible diamètre regroupant les individus à forts degrés. La forme particulière des éléments propres du laplacien permet d'extraire des communautés qui se répartissent en étoile autour du club huppé.

1 Introduction

Les réseaux sociaux sont des systèmes complexes dont certains ont des structures maintenant bien identifiées : graphes de petits mondes et graphes sans échelle typique. Un graphe sans échelle typique est un graphe dont la distribution des degrés n'est pas groupée autour d'une valeur moyenne ; c'est le cas lorsque celle-ci suit une loi de puissance. Les études menées sur le world wide web, des réseaux de courrier électronique ou des réseaux P2P, le réseau des collaborations scientifiques, le réseau des relations sexuelles en sont des exemples (Bornholdt et Schuster, 2003). Les graphes sans échelle typique ont peu de sommets de degrés très élevés et beaucoup de faible degré, ces graphes ont la propriété de présenter des fluctuations locales des degrés d'autant plus importantes que la distribution des degrés est proche d'une loi de puissance. Si les sommets de forts degrés sont connectés entre eux on parle alors de phénomène de "club huppé"

Alors que les études ont en général été effectuées sur des réseaux sociaux contemporains nous analysons ici un réseau relatif à la paysannerie médiévale. Nous travaillons sur une base de contrats agraires signés d'une part entre 1240 et 1350 et d'autre part entre 1450 et 1520 dans une petite région du Sud-Ouest de la France. Cette base pour l'instant réduite à environ 700 actes sera amenée à plus de 8000 actes lorsque le travail de saisie et de désambiguïsation

⁰Tous les calculs ont été effectués avec le logiciel libre R.

sera terminé. Les sommets du graphes sont les paysans et ils sont liés s'ils apparaissent dans un même contrat, nous définissons ainsi deux graphes G_{av} et G_{ap} ; nous avons éclairci la base en enlevant les seigneurs de notre étude. Nous ne possédons pas de données entre 1350 et 1450, intervalle temporel correspondant à la guerre de Cent Ans.

La notion de communauté varie en fonction du réseau que l'on étudie (Palla et al., 2005; Newman, 2006). Nous supposons dans notre étude que les communautés sont constitués d'individus qui ont à la fois les mêmes liens à l'intérieur de la communauté (clique) et à l'extérieur de la communauté. Nous verrons dans la section 3 que cette définition assez contraignante permet pourtant de révéler une structuration très particulière de notre réseau.

Cet article est constitué de deux parties : nous allons tout d'abord commencer par vérifier si notre graphe partage les propriétés rencontrées dans les grands réseaux d'interaction (l'effet petit monde, la distribution des degrés et le phénomène de club huppé). Ensuite nous nous attarderons sur la détection et l'organisation des communautés grâce à des méthodes spectrales. Enfin en conclusion, nous ébaucherons une comparaison des graphes avant et après la guerre de Cent Ans.

2 Les indices des deux réseaux d'interaction

2.1 L'effet petit monde

L'effet petit monde regroupe deux propriétés : la première énonçant que la distance entre deux sommets quelconques est faible (ceci est relatif à la connectivité globale) et la deuxième que la connectivité locale est forte. Pour quantifier ces notions nous utiliserons dans le premier cas soit la moyenne des plus courts chemins $\langle l \rangle$ soit la longueur caractéristique L (médiane des moyennes des plus courts chemins de chaque sommet (Watts, 2003)) et dans le deuxième cas la moyenne C_1 des densités du graphe des voisins de chaque sommet (Watts, 2003).

Le tableau TAB. 1 résume les résultats obtenus sur les graphes des liens de sociabilités paysans avant et après la guerre de Cent Ans et les met en perspective avec d'autres exemples de réseaux dont les densités d'arêtes sont voisines. Dans le cas de nos deux réseaux signalons qu'ils ont respectivement un diamètre de 5 et de 6 et que 90% des paires de sommets sont à une distance inférieure ou égale à 3.

	n	$ E $	$\langle l \rangle$	L	densité	C_1
G_{av}	205	1928	2.38	2.37	0.092	0.85
G_{ap}	173	1044	2.52	2.47	0.070	0.85
Collaboration entre physiciens ^(a)	107	757	2.48		0.13	0.72
Réseau proie/prédateur ^(b)	134	583	2.28		0.065	0.21
C.elegans ^(c)	282	1974		2.65	0.05	0.28

TAB. 1 – Etude de l'effet petit monde sur les graphes G_{av} et G_{ap} . ^(a)Iamnitchi et al. (2004), ^(b)Montoya et Solé (2002), ^(c)Watts et Strogatz (1998).

Les coefficients de clustering de nos réseaux sont sensiblement plus forts que ceux rencontrés habituellement.

2.2 La distribution des degrés

L'objectif de cette partie est de modéliser la distribution des degrés de nos graphes.

Afin de lisser les fluctuations nous étudions la distribution cumulative des degrés $P_c(k) = \sum_{j=k}^{\infty} P(j)$ où $P(j)$ est la probabilité d'avoir un sommet de degré j . Ce choix permet aussi de repérer plus aisément un degré de coupure éventuel au delà duquel la distribution décroît plus vite (Pastor-Satorras et Vespignani, 2004). Si beaucoup de réseaux récemment étudiés montrent une distribution cumulative des degrés qui suit une loi de puissance, la présence d'un degré de coupure est le signe d'un écart à cette loi. L'ajustement de la distribution par une loi de puissance tronquée par une coupure exponentielle (TPL) peut alors permettre de mieux expliquer l'ensemble de la distribution ; citons par exemple (Amaral et al., 2000; Achard et al., 2006).

Nous testons trois lois pour ajuster la distribution : loi de puissance $P_c(k) \sim k^{-\gamma}$, loi exponentielle $P_c(k) \sim e^{-\alpha k}$ et TPL $P_c(k) \sim k^{-\gamma} e^{-\frac{k}{k_c}}$. La figure FIG. 1 présente les résultats. Dans chacun des cas nous estimons les paramètres au sens des moindres carrés (cf. TAB. 2).

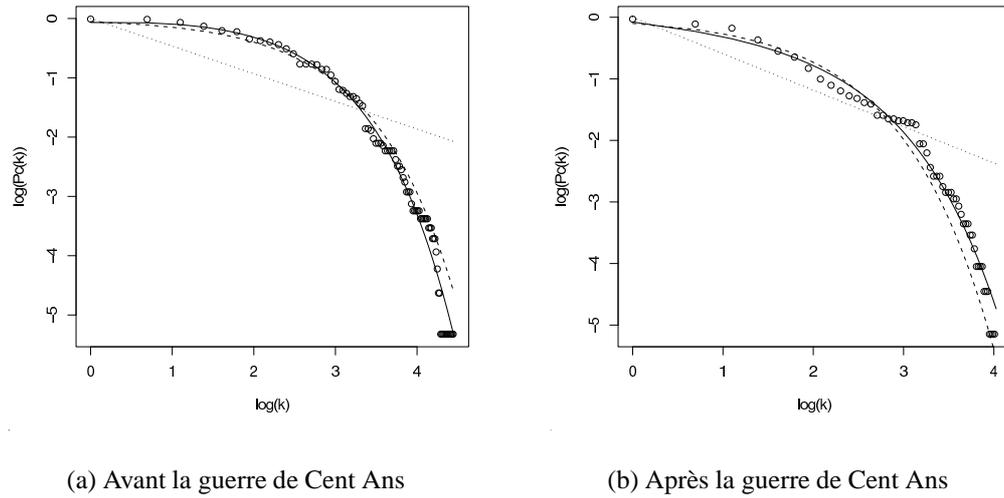


FIG. 1 – Distribution cumulative des degrés. En échelle logarithmique sont représentés en abscisse les degrés et en ordonnée la distribution cumulée P_c . Les cercles représentent les données, la ligne continue l'approximation par une TPL, la ligne discontinue par une exponentielle et la ligne en pointillés par une loi de puissance.

	puissance		exponentielle		TPL		
	erreur	γ	erreur	α	erreur	γ	k_c
G_{av}	22.73	0.466	0.81	0.054	0.30	-0.089	14.91
G_{ap}	8.96	0.592	1.38	0.099	1.11	0.109	13.04

TAB. 2 – Coefficients et erreurs quadratiques moyennes (10^{-3}) des différents modèles d'ajustement de la distribution cumulative des degrés.

Le meilleur ajustement de nos données est obtenu pour une TPL. Le graphe G_{av} échappe à une distribution des degrés en loi de puissance. Cette distribution est assez bien ajustée par une loi exponentielle même si on améliore l'erreur quadratique avec une TPL. Pour G_{ap} une loi de faible puissance donne de meilleurs résultats que pour G_{av} mais une TPL reste la meilleure modélisation.

2.3 L'effet « club huppé »

Nos deux graphes G_{av} et G_{ap} possèdent un club-huppé (Zhou et Mondragón, 2004) c'est-à-dire que les sommets de forts degrés (« les riches ») forment ensemble un sous-graphe dense. Ces individus jouant un rôle important dans l'organisation du réseau, les indices de centralité de proximité et de centralité d'intermédiarité (Degenne et Forsé, 1994) nous donnent un critère supplémentaire à celui de la densité et du degré pour le choix des individus du club-huppé. L'indice de centralité de proximité d'un sommet i est l'inverse de la somme des plus courts chemins de i aux autres sommets du graphe. L'indice de centralité d'intermédiarité d'un sommet i est le nombre des plus courts chemins du graphe passant par i . On classe les individus par ordre décroissant des degrés et on retient dans le club huppé ceux dont les indices de centralité sont élevés. Ceci nous conduit à retenir 25 individus dans G_{av} et 18 dans G_{ap} . Le diamètre des clubs huppés est de 2 et leurs densités sont respectivement 0.67 et 0.81.

3 Recherche des communautés

L'effet petit monde avec un coefficient de clustering élevé associé à une faible densité du graphe nous indique la présence de communautés ; afin de les déceler nous étudions le spectre du laplacien (non normalisé). Nous définissons nos communautés ainsi :

Définition 1 Une k -communauté d'un graphe G est une clique d'ordre k de G (clique non maximale) telle que tous les sommets de cette clique aient les mêmes voisins.

Nous utiliserons un théorème énoncé par van den Heuvel et Pejic (2000) démontrant que si le laplacien L du graphe a une valeur propre λ de multiplicité $k - 1$ dont les vecteurs propres associés possèdent exactement les mêmes k coordonnées non nulles alors ces k coordonnées correspondent aux sommets d'une k -communauté, au sens où nous l'avons définie.

Ces valeurs propres sont nécessairement entières. Considérons en effet u un tel vecteur propre nous avons $Lu = (D - A)u = \lambda u$ et il est par ailleurs facile de voir que $Au = -u$. Ainsi $Du = (\lambda - 1)u$ et λ est entier.

Nous pouvons énoncer un complément à ce théorème :

Théorème 1 (i) S'il existe une valeur propre λ de L de multiplicité $k - 1$ dont les vecteurs propres associés possèdent les mêmes $k + 1$ coordonnées non nulles et sont vecteurs propres de la matrice d'adjacence A associé à la valeur propre -1 alors il existe deux communautés d'ordre $k_1 \geq 2$ et $k_2 \geq 2$ tels que $k_1 + k_2 = k + 1$.

(ii) S'il existe une valeur propre de L de multiplicité $k - 1$ dont les vecteurs propres associés possèdent les mêmes $k + 2$ coordonnées non nulles et sont vecteurs propres de la matrice d'adjacence A associé à la valeur propre -1 alors il existe trois communautés d'ordre k_1, k_2, k_3 telles que $k_1 + k_2 + k_3 = k + 2$.

La démonstration consiste à étudier la matrice binaire $A + I$ qui est de rang 2 dans le cas (i) et de rang 3 dans le cas (ii). On procède par épuisement des cas.

En utilisant le théorème énoncé dans van den Heuvel et Pejic (2000) et le théorème 1 précédent, nous extrayons pour le graphe G_{av} 28 communautés de taille supérieure ou égale à 3 dont la plus importante est de taille 15 et pour le graphe G_{ap} 31 communautés dont la plus grande est de taille 7.

En supprimant la partie du graphe ne contenant aucune communautés (cette partie contient le club-huppé) nous obtenons un graphe à plusieurs composantes connexes. Les communautés trouvées via l'étude du spectre ne sont donc guère liées entre elles, elles sont préférentiellement liées à la partie que nous avons ôtée et notamment au club huppé. Nous avons donc une structure inter-communautaire proche de celle d'une étoile. Le centre de cette étoile contient le club-huppé dont on visualise bien à présent le rôle central qu'il joue dans l'organisation du réseau social. Ce partitionnement en club-huppé et communautés permet une bonne visualisation des graphes (FIG. 2)

4 Conclusion

Malgré le fait d'avoir enlevé les seigneurs de notre étude, nous constatons dans chacun des deux graphes G_{av} et G_{ap} la présence d'un groupe d'individus (le club huppé) possédant un rôle central. Ces deux graphes apparaissent sous forme d'une étoile de communautés. Dans G_{av} , un fort nombre de communautés de petite taille cohabitent avec un nombre significatif de communautés plus importantes, ce qui explique le bon ajustement des distributions des degrés avec une TPL. Concernant G_{ap} , la structuration est moins claire : le résidu est sensiblement plus important et la taille des communautés moins variable.

Si certaines de ces communautés correspondent à des zones géographiques comme on peut le voir dans (Hautefeuille, 2001), d'autres n'ont pour l'instant pas trouvé d'explications.

On remarquera un renouvellement quasi-complet des noms du club-huppé entre G_{av} et G_{ap} : la famille Combelcau très influente avant la guerre disparaît complètement après la guerre laissant la place à la famille Limairac, nouvelle famille qui paraît très influente.

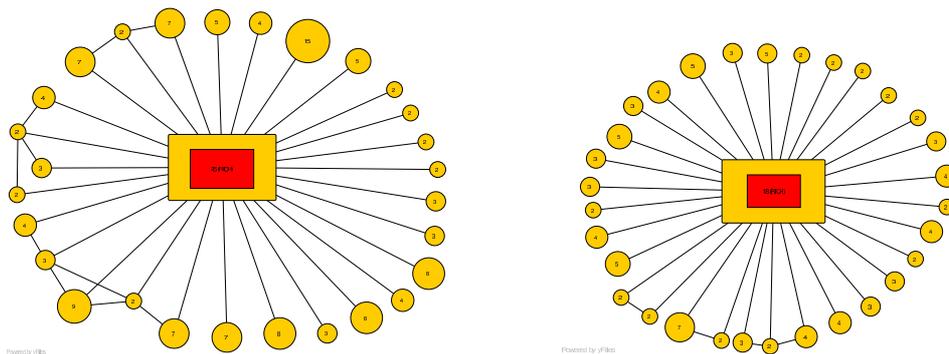


FIG. 2 – Graphe des communautés de G_{av} (à gauche) et G_{ap} (à droite). Les disques représentent les k -communautés extraites et le rectangle le reste des sommets (dont le club-huppé).

Références

- Achard, S., R. Salvador, B. Whitcher, J. Suckling, et E. Bullmore (2006). A Resilient, Low-Frequency, Small-World Human Brain Functional Network with Highly Connected Association Cortical Hubs. *The Journal of Neuroscience* 26(1), 63–72.
- Amaral, L., A. Scala, M. Barthélémy, et H. Stanley (2000). Classes of small-world network. *P.N.A.S.*
- Bornholdt, S. et H. G. Schuster (Eds.) (2003). *Handbook of Graphs and Networks - From the Genome to the Internet*. Wiley-VCH.
- Degenne, A. et M. Forsé (1994). *Les réseaux sociaux*. Armand Colin.
- Hautefeuille, F. (2001). Espace juridique, espace réel : l'exemple de la châtellenie de Castelnau-Montratier (Lot) aux XIIIe et XIVe siècles. *Habitats et territoires du sud, 126e congrès national des sociétés historiques et scientifiques.*
- Iamnitchi, A., M. Ripeanu, et I. Foster (2004). Small-world file-sharing communities. Volume 2, pp. 952–963.
- Montoya, J. et R. Solé (2002). Small world pattern in food webs. *Journal of theoretical biology* 214, 405–412.
- Newman, M. E. J. (2006). Modularity and community structure in networks. *Proc. Natl. Acad. Sci. U.S.A.* 103, 8577.
- Palla, G., I. Derenyi, I. Farkas, et T. Vicsek (2005). Uncovering the overlapping community structure of complex networks in nature and society. *Nature* 435, 814.
- Pastor-Satorras, R. et A. Vespignani (2004). *Internet, structure et évolution*. Editions Belin.
- van den Heuvel, J. et S. Pejic (2000). Using Laplacian Eigenvalues and Eigenvectors in the Analysis of Frequency Assignment Problems. *CDAM Research Report Series*.
- Watts, D. J. (2003). *Small Worlds : The Dynamics of Networks between Order and Randomness (Princeton Studies in Complexity)*. Princeton University Press.
- Watts, D. J. et S. H. Strogatz (1998). Collective dynamics of 'small-world' networks. *Nature* 393, 440–442.
- Zhou, S. et R. J. Mondragón (2004). The rich-club phenomenon in the internet topology. *IEEE Communications Letters* 8(3), 180–182.

Summary

In order to compare social organization of a medieval peasantry before and after the Hundred Years' War we study the structure of social networks built from a corpus of agrarian contracts. Low diameters and high clusterings show small-world graphs. Like many other networks studied these last years these graphs are scale-free. The distributions of the vertex degrees are fitted by a truncated power law. Moreover they have a rich-club : a dense core with a low diameter consisting of vertices with high degree. The particular shape of the laplacian spectrum allows us to extract communities that are spread along a star whose center is the rich-club.