

Mini-projet de TP : analyse de reseaux reels

SYSTEMES COMPLEXES AVANCES

M1 Informatique - semestre d'automne 2023-2024

UNIVERSITÉ CÔTE D'AZUR

Christophe Crespelle

`christophe.crespelle@univ-cotedazur.fr`

Le but de ce mini-projet de TP est d'analyser deux reseaux complexes provenant de donnees reelles : un de petite taille (entre 3K et 12K noeuds) et un de plus grande taille (entre 35K et 105K noeuds). Vous trouverez les deux reseaux qui vous ont ete attribues, ainsi qu'un lien pour telecharger les donnees brutes, sur la page https://www.i3s.unice.fr/~ccrespelle/enseignements/22-23-aut_M1SC/TP/TP1/affectation-reseaux.html.

Travail a effectuer.

La premiere chose que vous devez faire de filtrer les donnees brutes de ces deux reseaux et d'en calculer les statistiques de base (comme fait lors du TP1). Ensuite, vous choisirez une ou deux directions d'approfondissement que vous developperez en menant les analyses necessaires. Ces approfondissements doivent tenter de repondre a une ou des questions que vous choisirez de vous poser sur la structure de ces deux reseaux. Je vous propose ci-dessous plusieurs themes d'approfondissement possibles. Vous pouvez en choisir un parmi eux, ou bien mixer et croiser les questions entre plusieurs des themes proposes ou, encore mieux, choisir une direction que vous aurez vous meme imaginee.

Rendu.

Votre rendu sera constitue du rapport que vous rendrez et qui comprendra les resultats de vos analyses. Ce rapport sera obligatoirement au format PDF, il fera au plus 10 pages (tout compris : titre, references, etc.), avec une police de caractere d'au moins 11pt. Il pourra contenir des annexes qui seront lues a la discretion du correcteur. Mon conseil est que vous selectionnez les analyses qui donnent les resultats les plus interessants parmi celles que vous aurez faites et que vous mettiez les autres en annexe. **Vous devez formuler explicitement et tres clairement les questions auxquelles tentent de repondre vos analyses** et decrire ce que contiennent, objectivement, les courbes que vous tracez. Portez un soin particulier a la lisibilite des legendes des axes. Vous devez aussi produire une analyse critique des resultats mis en evidence par ces courbes. Enfin, pour chaque statistique que vous calculez, vous devez donner en annexe les outils, les lignes de commandes ou le code que vous avez utilises pour produire ces statistiques. Si vous avez ecrit du code, envoyez les fichiers sources en plus du rapport.

1 Prise en main et filtrage des donnees

La premiere etape de votre projet consistera a filtrer les donnees brutes de vos deux reseaux et a les mettre au format decrit dans le TP1. Vous realiserez toutes les analyses demandees dans le TP1 : calculez les quatre proprietes fondamentales de vos reseaux et comparez les a celles des modeles d'Erdo-Renyi et de configuration de memes parametres. Vous pouvez, si vous le souhaitez, vous affarncir du calcul de la distribution du coefficient de clustering local par noeud (Question 9). Si vous rencontrez des difficultes de calcul pour la distance moyenne, pensez que vous pouvez vous contenter de faire de l'echantillonnage (https://www.youtube.com/watch?v=E0-_aSMkwGw).

Vous trouverez ci-dessous plusieurs suggestions d'approfondissements possibles. Vous devez explorer au moins une direction d'approfondissement, ou deux. Mais vous n'etes pas obliges de la choisir parmi celles proposees. Vous pouvez a la place croiser les questions des differentes directions proposees ou proposer vous meme une autre direction. N'hesitez pas, l'originalite des analyses menees sera fortement valorisee dans la notation.

2 Distances

Cette partie propose d'approfondir les proprietes metriques du reseau, c'est a dire des distances et des chemins dans le reseau. Vous pourrez, par exemple, vous poser les questions suivantes.

Question 1. Calculez la distribution des distances dans le reseau.

Definition 1 (Excentricite). *L'excentricite d'un sommet u est le maximum des distances entre u et chacun des autres sommets du graphe : $exc(u) = \max_{v \in V \setminus \{u\}} \{dist(u, v)\}$.*

Question 2. Calculez la distribution de l'excentricite d'un sommet.

Question 3. Quelle est l'excentricite minimum d'un sommet, combien de sommets realisent ce minimum ?

Question 4. Combien de couples realisent le diametre ? Combien de sommets sont impliquees dans un tel couple ?

Question 5. Quel est le nombre minimum de sommets a enlever du graphe pour faire chuter le diametre de 1 ? de k ?

3 Centralites

Une possibilite est d'etudier et de comparer differentes metriques de centralite dans le reseau.

Question 6. Calculez au moins deux notions de centralite des noeuds et etudiez les correlations entre leurs scores.

Question 7. Etudiez les correlations des classements qu'elles donnent sur les noeuds.

Question 8. Sur quel(s) type(s) de noeuds sont elles le plus d'accord ? le moins d'accord ?

Question 9. Pour les mesures définies aussi sur les liens (betweenness, eigenvector, PageRank), étudiez la corrélation entre la centralité d'un lien et celles de ses deux nœuds incidents.

4 Communautés

Une autre direction possible est l'étude de la structure communautaire du réseau.

Question 10. Comparez le résultat de différentes méthodes de partition en communautés. Sont-elles ressemblantes ? Où sont-elles différentes ?

Question 11. Étudiez la stabilité de la méthode de Louvain lorsqu'on change l'ordre dans lequel sont traités les sommets. Les différentes partitions obtenues sont-elles ressemblantes ? Où sont-elles identiques ? Où sont-elles différentes ?

Question 12. Étudiez la structure du réseau entre les communautés. Combien d'arêtes y a-t-il entre les communautés ? Combien à l'intérieur des communautés ? Quelle est la distribution du nombre d'arêtes entre deux communautés ? La distribution du degré d'une communauté ?

Question 13. Étudiez les interfaces entre les communautés. Quels sont les nœuds à l'interface de plusieurs communautés ? Y en a-t-il beaucoup ? Quelles sont leurs caractéristiques ?

Question 14. Analysez la structure communautaire à l'intérieur de la plus grosse communauté.

5 Cœur/periphérie du réseau

La question poursuivie ici est de dégager la structure du cœur (k -core en anglais) et de la périphérie du réseau et de savoir en quoi celles-ci diffèrent de celle du réseau tout entier.

Definition 2 (voir [https://en.wikipedia.org/wiki/Degeneracy_\(graph_theory\)](https://en.wikipedia.org/wiki/Degeneracy_(graph_theory))).
Le k -cœur d'un réseau est l'ensemble des nœuds qui restent après avoir retiré itérativement tous les nœuds de degré au plus $k-1$, en retirant aussi ceux qui deviennent de degré au plus $k-1$ au cours de ce procédé. Le cœur du réseau est le k -cœur obtenu pour le k maximum tel que le k -cœur est non vide.

Question 15. Déterminez les k -cœurs du réseau. Comment évolue le nombre de nœuds dans le k -cœur en fonction de k ?

Question 16. Où se situent les arêtes du réseau par rapport à ses différentes couches de k -cœur ? On pourra par exemple :

- faire la distribution de l'écart dans ces couches entre les deux extrémités d'une arête,
- tracer l'évolution en fonction de k du nombre d'arêtes dans le k -cœur,
- tracer l'évolution en fonction de k du nombre d'arêtes incidentes à un nœud du k -cœur.

Le coefficient de *rich club* (voir https://en.wikipedia.org/wiki/Rich-club_coefficient) mesure à quel point les nœuds de fort degré sont connectés entre eux.

Question 17. Déterminez le coefficient de *rich club* du réseau.

Question 18. Y a-t-il un rapport entre cœur et rich club ?

Question 19. Le cœur ou le rich club ont-ils la même structure que le réseau entier ?

Question 20. Quelles sont les propriétés de la périphérie du réseau, c'est à dire la partie qui n'est pas dans le cœur ?

6 Préviation de liens

On pourra également s'intéresser à la structure des liens les plus prédictibles dans le réseau et aux rapports possibles avec d'autres propriétés structurelles du réseau comme ses communautés, son cœur, la centralité de ses nœuds, etc.

Question 21. Utilisez une méthode de prédiction de liens pour identifier les couples de nœuds les plus enclins à former une nouvelle arête. Où sont situés ces couples par rapport aux communautés et au cœur du réseau ?

Question 22. Classez les liens existant dans le réseau par niveau de prédictibilité décroissant. Pour cela, retirez un lien du réseau et calculez le score qui lui est attribué par la méthode de prévision.

Question 23. Comment évoluent les propriétés du réseau lorsqu'on ajoute ses liens un par un par prédictibilité décroissante ? Le réseau des liens les plus prédictibles a-t-il la même structure que le réseau tout entier ?

Question 24. Y a-t-il un rapport entre prédictibilité des liens du réseau et importance des nœuds auxquels ils sont rattachés ?