

AAGB – TME 5

Réseaux

1 Introduction

La théorie des réseaux consiste en l'étude des propriétés statistiques qui caractérisent la structure et le comportement des systèmes définis par des réseaux. On propose d'étudier quelques mesures classiques pour les réseaux : distribution de degré, coefficient de clustering et betweenness centrality. Pour tester vos fonctions, vous pourrez par exemple utiliser les deux fichiers (reseau1, reseau2) qui listent l'ensemble des arêtes présentes dans un réseau. On peut par exemple parser ces fichiers et créer un dictionnaire avec pour clés les noeuds du réseau, et pour valeurs la liste des noeuds reliés au noeud en question.

2 Contexte biologique

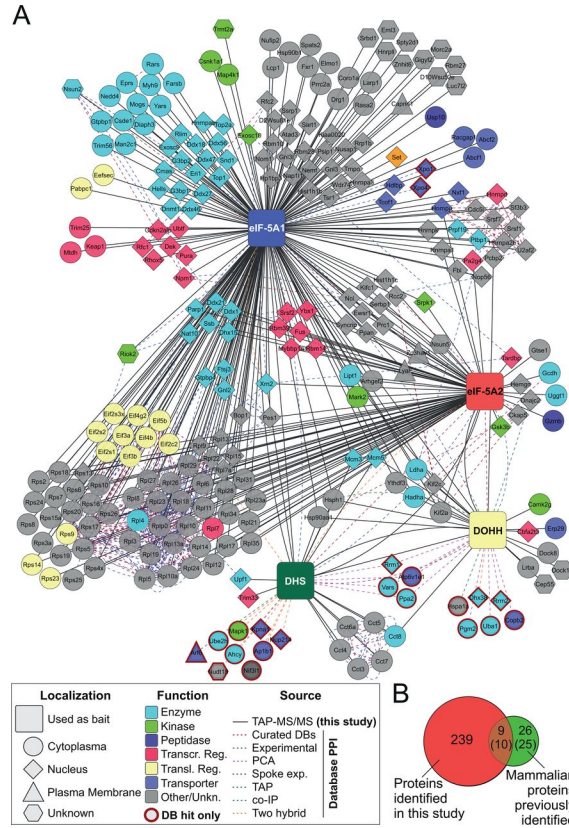
Un graphe est une structure mathématique modélisant les relations entre paires d'objets. C'est une structure très utilisée en biologie, qui permettent par exemple l'étude de :

- Réseaux d'interaction protéine-protéine (voir figure, chaque noeud correspond à une protéine, chaque arc à une interaction entre deux protéines).
- Réseaux métaboliques (ensemble des processus métaboliques et physiques qui déterminent les propriétés physiologiques et biochimiques d'une cellule).
- Réseaux d'expression génétique (chaque noeud correspond à un gène, chaque arc à une relation de co-expression).
- Réseaux phylogénétiques (!!, pourquoi?)

L'étude de ces graphes par des mesures telles que celles décrites ci-dessous est importante parce qu'elle permet d'analyser les propriétés du réseau, par exemple en prédisant la fonction d'éléments inconnus par leur relation avec des éléments bien caractérisés.

La plupart de ces réseaux ne sont pas aléatoires mais "scale-free", comme dans d'autres contextes les réseaux sociaux, internet, etc. Ces réseaux ont la particularité d'être :

- Robustes (bonne capacité à résister aux attaques aléatoires).
- Hétérogènes (les hubs, i.e. les noeuds avec le plus d'arcs les reliant aux autres noeuds, ne sont généralement pas reliés à d'autres hubs)
- Très connectés à l'intérieur du graphe, peu connectés aux extrémités.



Protein-protein-interaction Network Organization of the Hypusine Modification System (H. Sievert et al.)

3 Degree distribution

1. Définissez la distribution de degrés pour un graphe.
2. Quel est l'intérêt d'étudier la distribution des degrés dans un graphe ?
3. Tracez la distribution des degrés des réseaux donnés .
4. Qu'est ce qu'un réseau scale-free ? Parmi les deux réseaux étudiés, lequel pourrait être scale-free ? Lequel aléatoire ? Pourquoi ?

4 Clustering coefficient

Le *coefficient de clustering* donne une mesure *locale* de l'interconnectivité des voisins d'un sommet. Pour un réseau scale-free, la distribution log-log des moyennes des coefficients de clustering, $C(k)$, en fonction du degré des sommets, k , suit une loi de puissance du type $C(k) \propto k^{-\gamma}$. Cela indique que certains nœuds du réseau se regroupent en clusters (contrairement aux réseaux aléatoires).

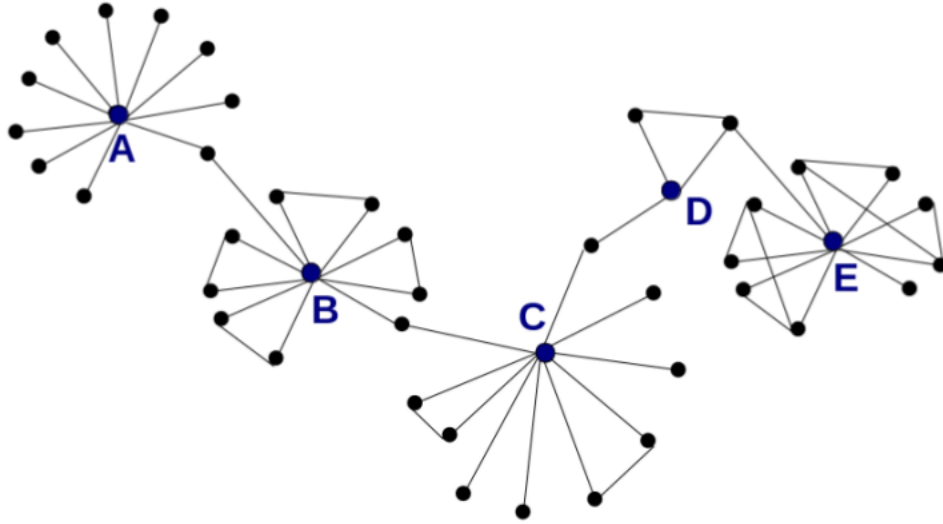
Définition 1 (Clustering Coefficient) Soit $G = (V, E)$ un graphe non orienté et $v \in V$ un sommet. Soit $N(v)$ l'ensemble des voisins de v et $E(N(v)) = \{(u_1, u_2) \in E \mid u_1, u_2 \in N(v)\}$.

Le coefficient de clustering $cc(v)$ de v est défini par

$$cc(v) = \frac{2 | E(N(v)) |}{d^\circ(v)(d^\circ(v) - 1)}$$

En d'autres termes, $cc(v)$ est le ratio entre le nombre de 'triangles' passant par le sommet v ($|E(N(v))|$) et le nombre de 'triangles' qui pourraient passer par v ($d^\circ(v)(d^\circ(v) - 1)/2$).

Considérons le réseau G suivant :



1. Calculez $cc(A)$, $cc(B)$, $cc(C)$, $cc(D)$ et $cc(E)$.
2. Ecrivez une fonction permettant de calculer la moyenne des coefficient de clustering. Que peut on en déduire pour les réseaux étudiés ?

5 Betweenness centrality

La betweenness centrality (ou centralité intermédiaire) est une propriété qui tente de capturer l'importance d'un nœud (ou d'un lien) pour la propagation d'un signal/message, autrement dit son importance en tant qu'intermédiaire dans la structure du réseau. Elle repose essentiellement sur la notion de plus court chemin (ie. le chemin comportant le moins d'arcs possible).

Par convention, on notera σ_{st} le nombre de plus courts chemins entre les nœuds s et t d'un graphe G et $\sigma_{st}(v)$ le nombre de plus courts chemins entre les nœuds s et t passant par v (on suppose ici que $s \neq t \neq v$). La betweenness centrality d'un nœud v est alors définie par la formule suivante :

$$BC(v) = \sum_{s \neq t \neq v} \frac{\sigma_{st}(v)}{\sigma_{st}}$$

Un cas particulier simple pour le calcul de $BC(v)$ se présente lorsqu'il existe exactement 1 plus court chemin pour toutes les paires de sommets (s, t) de G . Dans ce cas, $BC(v)$ vaut

exactement $BC(v) = \sum_{s \neq t \neq v} \sigma_{st}(v)$, c'est-à-dire le nombre de plus courts chemins dans G passant par v .

1. Quel type de noeuds permet de mettre en évidence la betweenness centrality ? Un réseau comportant quelques noeuds avec une forte centralité est-il robuste ? Comment utiliser cette mesure pour détecter des communautés ?
2. Ecrire une fonction permettant de calculer la betweenness centrality.