

TRAVAUX PRATIQUE

Graph Mining

Advanced Data Mining Report

MASTER DEGREE IN BIG DATA MANAGEMENT AND ANALYTICS

Prepared by:

Supervised by:

Ferreira, Jessé

RAMEL, Jean-Yves

ACADEMIC YEAR : 2019-2020



UFR DES SCIENCES ET TECHNIQUES - SITE DE BLOIS

Contents

1	TP Questions	2
---	--------------	---

List of Figures

1	Graph de type Mixte (Oriente et non oriente)	2
2	The top 10 highest degree nodes.	4

1 TP Questions

1. Question: Si on considère un graphe G pour lequel :
 1. les nœuds sont les personnages et les comics
 2. les arêtes sont les liens (personnage, comic)

Ce graphique ressemble à un graphique mixte, puisqu'un personnage peut apparaître dans plusieurs comics et plusieurs comics peuvent montrer (graph orienté) et aussi faire référence à plusieurs personnage, par contre ce personnage peut n'apparaître pas dans les comics qu'il était cité (graph non orienté).

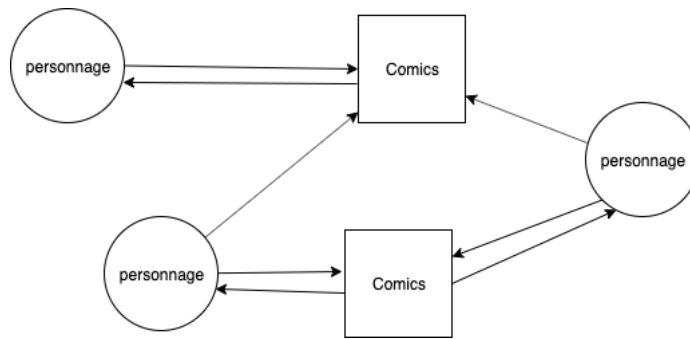


Figure 1: Graph de type Mixte (Oriente et non orienté)

Pour transformer G^* en un graphe de type 'réseau social' à partir du jeu de donnée 'Marvel Universe', on doit change les arêtes (personnage, comic) par (personnage, personnage) et attribué le comic à les attribut des personnage. Ainsi, le graph sera un graph orienté.

Donc, le poids de G^* pourrait etre associé à le nombre de fois qu'un personnage apparaît dans un comics.

2. Question: La réponse est sur le code.
3. Question: La réponse est sur le code.
4. Question: La réponse est sur le code.
5. Question: La réponse est sur le code.

6. Question: Reponse: Cela permet de mieux voir les patterns et la distribution des nœuds dans le canvas en améliorant la visualisation.
7. Question: la réponse est sur le code.
8. Question: la réponse est sur le code.
9. Question:

```
pop_heros = []

foreach pop_hero in edges:
    if pop_hero.weight >= 50:
        popularity = sum of unique weights
        of pop_hero edges
        pop_heros.push({ pop_hero , popularity })

return sortedDesc(pop_heros)
```

10. Question: la réponse est sur le code.
11. Question: Il y a encore des nœuds qui empêchent d'avoir une visualisation propre de l'autre nœud sélectionné, dont il faut supprimer les nœuds sans edges (target et source).
12. Question: Cette fonction récupérera les autorités en vérifiant les statistiques de hits dans chaque nœud. L'objectif est de rechercher les nœuds les plus pertinents. Ensuite, la liste des autorités est triée pour donner la couleur aux top nœuds sélectionnés.

Les 5 personnages les plus importants, selon le critère HITS sont:

“

HITS 0 : CAPTAIN AMERICA
HITS 1 : HAWK
HITS 2 : ANT-MAN/DR. HENRY J.
HITS 3 : THING/BENJAMIN J. GR
HITS 4 : HULK/DR. ROBERT BRUCE

“

13. Question: les 10 personnages les plus importants, selon le critère des degrés élevés sont:

“

HITS 1 : CAPTAIN AMERICA
HITS 2 : HAWK

HITS 3 : ANT-MAN/DR. HENRY J.
HITS 4 : THING/BENJAMIN J. GR
HITS 5 : HULK/DR. ROBERT BRUC
HITS 6 : QUICKSILVER/PIETRO M
HITS 7 : HUMAN TORCH/JOHNNY S
HITS 8 : WONDER MAN/SIMON WIL
HITS 9 : MR. FANTASTIC/REED R
HITS 10 : IRON MAN/TONY STARK
““

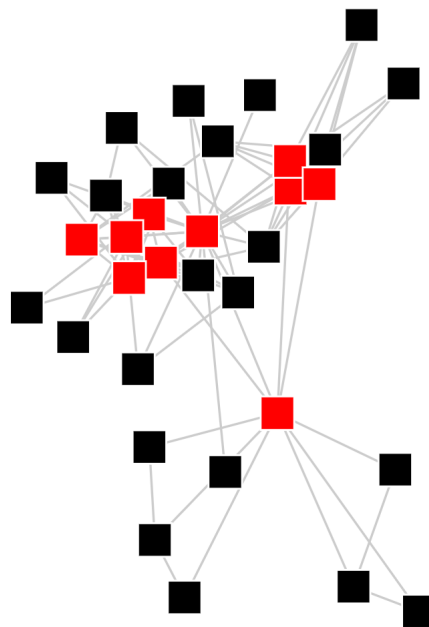


Figure 2: The top 10 highest degree nodes.

14. Question:

Part 1:

Les couleur des noudes sont attribué en basant sur les classes des communautés de chaque noudes qui est generé par l'index calculé sur le methode statistc du

louvain.

Part 2:

les repartitions sont ordonnées et mapped dans l'ordre croissant donc l'algorithme de Louvain reorganise les graphs en basant sur le clustering hiérarchique, en regroupant les nœuds voisins en communautés.

15. Question:

À un moment donné peut être exister deux graphes $G_1(N_a, E_a)$ et $G_2(N_b, E_b)$ avec éventuellement un nombre différent de nœuds et d'arêtes, et ils peuvent avoir des correspondances entre graph, c'est-à-dire, avoir la même quantité de nœuds voisins, dont les voisins aussi peuvent avoir le même nombre de prochains voisins et etc. Donc, pour mesurer la similarité on peut comparer les ensembles de voisins de chaque G (G_1 et G_2), par contre, ça peut devenir coûteux.

L'intérêt en ce type d'analyse est qu'on peut trouver des héros qui peuvent avoir de mêmes amis et aussi ennemis, et donc ils peuvent être ensemble pour sauver la vie humaine contre Thanos, comme dans le film "Avengers: Endgame".

```
voisans1 = getNeighbours(sig_GA, nodeId_GA)
voisans2 = getNeighbours(sig_GB, nodeId_GB)
```

```
for i in voisans1.nodes:
    for j in voisans2.nodes:
        if i != j:
            dissimilarity_flag = 0
        else:
            similarity_flag = 1
```

16. Question:

Là l'idée est un peu pareil à cela sur la dernière question, mais là on cherche les communautés qui se rassemblent à la fois de façon concentrée de l'analyse sur chaque individu, maintenant on cherche la similarité de sub-graph. On peut faire cela après la partition hiérarchique, et donc, parcourir les ensembles de sub-graph en demandant et comparant les labels.

L'intérêt en ce type d'analyse est de trouver de mêmes caractéristiques entre

des communautés.

```
for i in graphA.nodes:
    for j in graphB.nodes:

        if node[i].louvain != node[j].louvain
            dissimilarity_flag = 0
        else
            similarity_flag = 1
```

17. Question 17 : Après la repartitions de communautés, l'algorithme va mesurer les poids entre nœuds de une communauté et entre communautés. Où le densité intra-cluster c'est le nombre de arcs interne à une communauté par quantité de possibles nœuds, dont le global internal_density sera la moyenne entre les repartitions. Et inter-cluster c'est la maximum distance entre les communautés.

$intracluster = (Sum_{des} |E| / (|V| * |V| - 1)) / \text{nb repartitions total}$

$intercluster = |E| - |ExternalE| / (|V| * |V| - 1) - \text{quantité total de possibles nœuds}$

Nb internal edges = 80

Internal density = 0.25774198568316214

External density = 0.0016790597265531303

18. Question: La réponse est sur le code. Attention, il y a un bug dans cette fonction. Pour l'utiliser, il faut d'abord cliquer dans le respectif bouton et donc choisir deux nœuds à voir la meilleure distance. Puis, il faut cliquer à nouveau pour afficher le résultat sur console et sur le tab.