



RAPPORT DE PROJET ANALYSE DES RÉSEAUX SOCIAUX



Élaborée par :
Wiem BOUKHARI

Proposé par :
DR. Fériel BEN FRAJ

I. INTRODUCTION :

Ce projet explore le réseau de collaboration académique dans le domaine de la Relativité Générale et de la Cosmologie Quantique (GR-QC), basé sur les articles soumis à la catégorie correspondante sur arXiv.

En examinant les collaborations entre auteurs sur une période de janvier 1993 à avril 2003, nous cherchons à comprendre la structure et l'évolution de cette communauté scientifique.

Notre analyse comprend l'exploration des données, la détection de communautés et la prédiction des liens potentiels entre chercheurs.
Les résultats fourniront des informations précieuses sur la dynamique des collaborations dans ce domaine de recherche fondamental.

II. COLLECTE DES DONNÉES :

1. Identification de la source de données en ligne :

J'ai sélectionné le jeu de données Arxiv GR-QC (General Relativity and Quantum Cosmology) collaboration network.

Le réseau de collaboration Arxiv GR-QC provient des e-preprints de arXiv et couvre les collaborations scientifiques entre les auteurs de documents soumis à la catégorie de la Relativité Générale et de la Cosmologie Quantique.

Si un auteur i a co-écrit un document avec l'auteur j , le graphe contient une arête non dirigée de i à j .

Si le document est co-écrit par k auteurs, cela génère un (sous)-graphe complètement connecté sur k nœuds.

Les données couvrent les documents de la période de janvier 1993 à avril 2003 (124 mois). Elles commencent quelques mois seulement après la création de arXiv, et représentent ainsi essentiellement l'histoire complète de sa section GR-QC.

2. Méthodes de collecte des données :

Les données ont été obtenues directement à partir du site web de Stanford SNAP (Stanford Network Analysis Project).

J'ai téléchargé le fichier de données au format .txt à partir de la page dédiée au jeu de données Arxiv GR-QC.

3. Identification des entités (nœuds) et des relations entre elles (liens) :

Dans notre ensemble de données, les entités sont représentées par les auteurs de documents scientifiques, et les relations entre ces entités sont établies par les collaborations entre ces auteurs.

```
print("Nombre total d'entités (nœuds) :", nombre_entites)
print("Nombre total de relations (liens) :", nombre_relations)
```

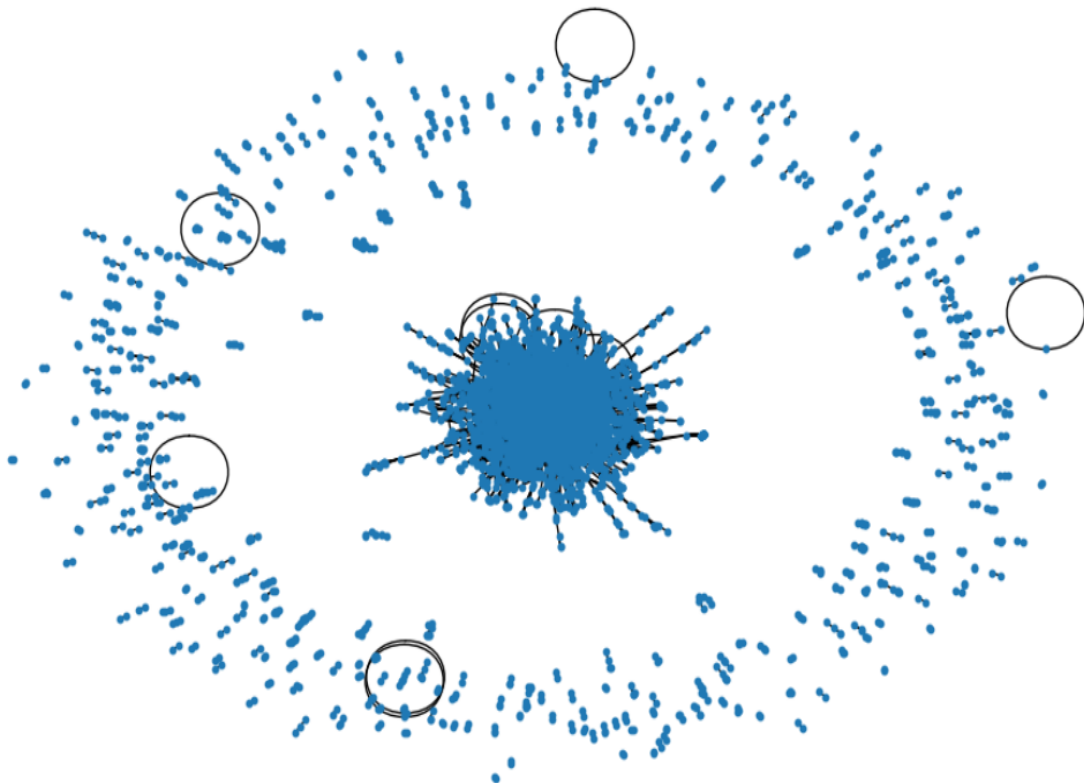
```
Nombre total d'entités (nœuds) : 5242
Nombre total de relations (liens) : 14496
```

4. Construction du réseau :

Une fois les données téléchargées, j'ai utilisé des bibliothèques Python telles que NetworkX pour construire un réseau à partir des informations de collaboration entre les auteurs. Chaque nœud dans le réseau représente un auteur, et chaque lien représente une collaboration entre deux auteurs.

En construisant ce réseau, j'ai créé une représentation graphique des interactions entre les auteurs dans le domaine de la relativité générale et de la cosmologie quantique.

Graphe de General Relativity and Quantum Cosmology collaboration network



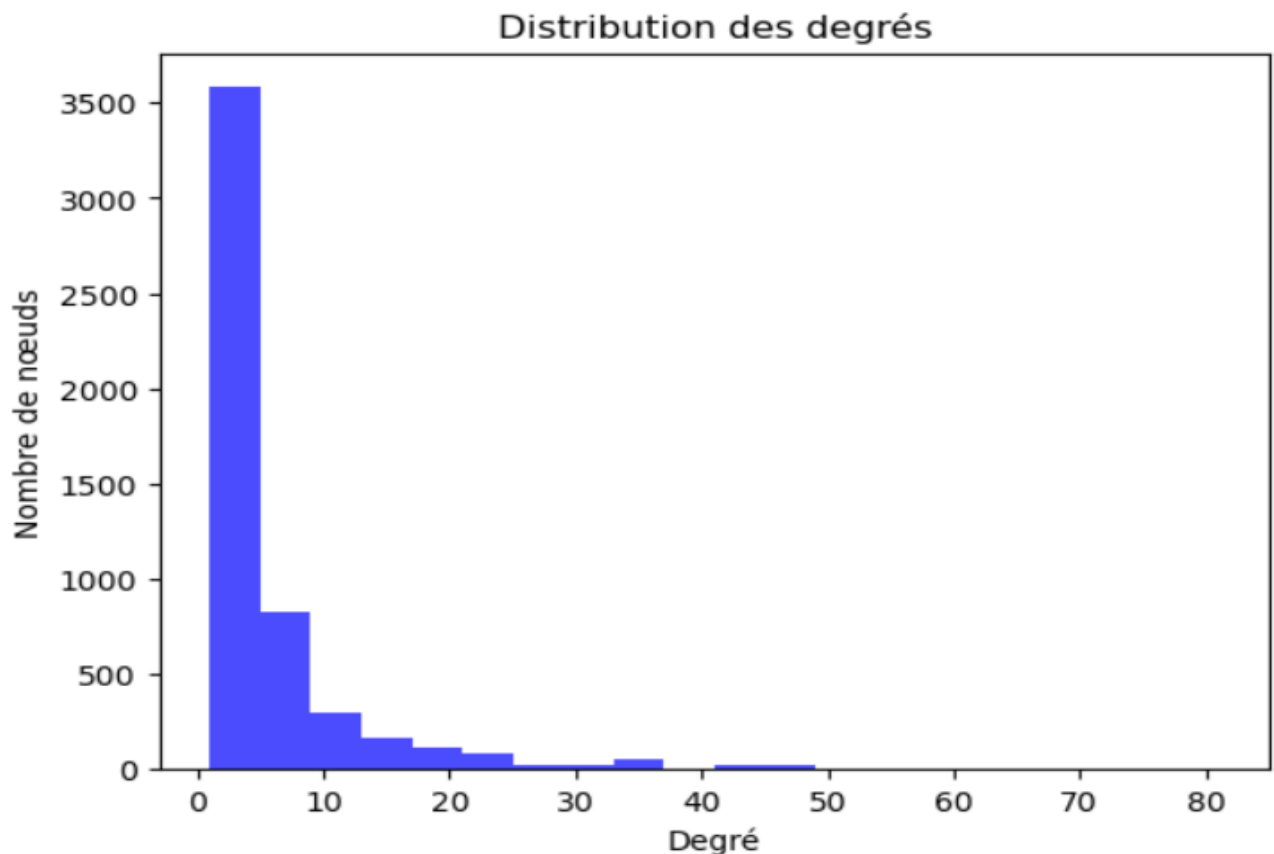
III. ANALYSE DU RÉSEAU :

Dans cette deuxième partie, je me suis concentré sur l'analyse exploratoire des données pour mieux comprendre la structure et les caractéristiques du réseau de collaboration entre les auteurs dans le domaine de la relativité générale et de la cosmologie quantique.

➤ L'analyse de la distribution des degrés :

J'ai calculé la distribution des degrés des nœuds dans le réseau, ce qui m'a permis de comprendre la répartition des collaborations entre les différents auteurs.

Une distribution des degrés étudiée peut révéler des informations sur la structure globale du réseau, telles que la présence de nœuds fortement connectés ou de nœuds isolés.



➤ L'analyse des composants connectés :

J'ai identifié les composants connectés dans le réseau, c'est-à-dire les sous-graphes dans lesquels chaque nœud est relié à tous les autres nœuds par un chemin.

Cette analyse permet de comprendre s'il existe des sous-groupes d'auteurs qui sont fortement connectés entre eux mais moins connectés au reste du réseau.

Nombre de composants connectés :	355
Taille du composant connecté 1 :	4158
Taille du composant connecté 2 :	2
Taille du composant connecté 3 :	3
Taille du composant connecté 4 :	2
Taille du composant connecté 5 :	2
Taille du composant connecté 6 :	2
Taille du composant connecté 7 :	2

➤ L'analyse des chemins :

J'ai étudié les chemins les plus courts entre les paires de nœuds dans le réseau.

Longueur moyenne des chemins les plus courts :	6.049380016182999
------------------------------------------------	-------------------

Cette analyse permet de comprendre la distance moyenne entre les auteurs et la facilité avec laquelle l'information peut se propager à travers le réseau.

Distance moyenne entre les nœuds :	3.8061249382218865
------------------------------------	--------------------

➤ Le coefficient de clustering et l'analyse de la densité :

J'ai calculé le coefficient de clustering moyen pour évaluer à quel point les nœuds dans le réseau forment des clusters ou des regroupements denses.

La densité du réseau a également été évaluée pour comprendre à quel point le réseau est dense en termes de connexions entre les nœuds.

Coefficient de clustering moyen :	0.5294338141753044
Densité du graphe :	0.0010545461581625744

➤ L'analyse de la centralité :

J'ai calculé plusieurs mesures de centralité, telles que la centralité d'intermédiation (betweenness centrality) et la centralité de proximité (closeness centrality), pour identifier les auteurs les plus centraux dans le réseau.

Ces analyses permettent de comprendre qui sont les auteurs les plus influents ou les mieux positionnés dans le réseau de collaboration.

Moyenne de la centralité d'intermédiation :	0.0006062707858585091
---------------------------------------------	-----------------------

Auteur: 3466
Betweenness Centrality: 0.001715248119663197
Closeness Centrality: 0.15250009440442658
Eigenvector Centrality: 3.489474540927225e-05

Auteur: 937
Betweenness Centrality: 0.0003493412893676944
Closeness Centrality: 0.14789649865964416
Eigenvector Centrality: 1.459657837625449e-05

Auteur: 5233
Betweenness Centrality: 0.0
Closeness Centrality: 0.14373793718636851
Eigenvector Centrality: 6.681320201491573e-06

Auteur: 8579
Betweenness Centrality: 0.00035838805830684087
Closeness Centrality: 0.14223737289668725
Eigenvector Centrality: 5.131748375726487e-06

Auteur: 10310
Betweenness Centrality: 0.00875853820067516
Closeness Centrality: 0.17089274080637024
Eigenvector Centrality: 0.0002698843812326264

IV. IDENTIFICATION DES COMMUNAUTÉS:

Pour identifier les communautés dans le réseau de collaboration, j'ai utilisé plusieurs algorithmes de détection des communautés, notamment :

Louvain : Cet algorithme est largement utilisé pour détecter les communautés dans les réseaux. Il maximise la modularité du réseau en recherchant des partitions qui optimisent le regroupement des nœuds en communautés.

Propagation des labels : Cette méthode attribue initialement une étiquette à chaque nœud, puis propage ces étiquettes à travers le réseau en fonction des nœuds voisins. Les nœuds finissent par appartenir à la même communauté s'ils partagent des étiquettes similaires.

K-clique : Cet algorithme identifie les sous-graphes complets de taille k (cliques) dans le réseau. Chaque clique est considérée comme une communauté potentielle.

Infomap : Cette méthode utilise la théorie de l'information pour détecter les structures modulaires dans le réseau. Elle cherche à minimiser le code de longueur moyenne nécessaire pour décrire les mouvements d'information sur le réseau.

J'ai implémenté ces algorithmes à l'aide de la bibliothèque NetworkX et évalué leurs performances en termes de qualité des partitions identifiées.

En utilisant les attributs supplémentaires disponibles dans le réseau, tels que les noms des auteurs ou d'autres métadonnées, j'ai cherché à donner un sens aux communautés identifiées.

L'analyse des communautés peut être étendue en explorant d'autres approches disponibles dans la bibliothèque CDlib ou en développant des méthodes personnalisées en fonction des besoins spécifiques du réseau étudié.

mesures de centralité, etc. Chaque algorithme a été ajusté et évalué en fonction de ses performances prédictives sur l'ensemble de test.

✧ Résultats de l'algorithme de K-Means :

Exactitude de prédiction: 0.5139751552795031

✧ Résultats de l'algorithme de Forêt aléatoire :

```
Exactitude de prédiction: 1.0
```

✧ Résultats de l'algorithme de Réseaux de neurones :

```
Epoch 10/10
580/580 627us/step - accuracy: 1.0000 - loss: 5.3806e-09 - val_accuracy: 1.0000 - val_loss: 9.3524e-06
182/182 444us/step - accuracy: 1.0000 - loss: 4.0916e-22
Exactitude de prédiction sur les données de test: 1.0
```

✧ Résultats de l'algorithme de K-Nearest Neighbors :

```
# Évaluer la performance du modèle
accuracy = np.mean(predictions == np.zeros(len(X_test)))
print("Exactitude de prédiction:", accuracy)
```

Exactitude de prédiction: 1.0

✧ Résultats de l'algorithme de Multi-Layer Perceptron :

```
print("Exactitude de prédiction:", accuracy)
```

Exactitude de prédiction: 1.0

✧ Résultats de l'algorithme de l'arbre de décision :

```
Exactitude de prédiction: 1.0
```

Il semble que tous les algorithmes utilisés ont une exactitude de prédiction de 1.0 sur ce ensemble de données, cela signifie que ces modèles ont réussi à prédire correctement tous les liens dans l'ensemble de test.

Une exactitude de prédiction de 1.0 peut être inhabituelle et peut indiquer plusieurs choses :

Sur-apprentissage (Overfitting) : Les modèles peuvent avoir sur-appris les données d'entraînement et mémorisé les exemples au lieu d'apprendre des motifs généraux. Cela peut se produire si les modèles sont trop complexes par rapport à la quantité de données d'entraînement disponible.

Données d'entraînement et de test similaires : Il est possible que les données d'entraînement et de test soient très similaires, ce qui permet aux modèles de généraliser efficacement à de nouvelles données.

Données de test trop simples : Les données de test peuvent être trop simples ou peu représentatives, ce qui facilite la prédiction pour les modèles