

Graph Theory and its Applications

Assignment 4

Movie Recommender System

Abhishek Patil	PES1UG20CS010	'A' Section
Amogh S Vithapnavar	PES1UG20CS035	'A' Section

About the Dataset :

For building the recommender system, we use three main datasets

- 'movies.csv'
- 'ratings.csv'
- 'genome scores.csv'

provided by [GroupLens](#).

We further break and perform some operations on these 3 datasets to get a total of 7 datasets

- 3 of them are for creating nodes

a. Users

b. Movies

c. Genres

- 4 are for relationships

a. WATCHED

b. SIMILAR

c. GENRE

d. FAVORITE

Step 1. Import the Datasets into Neo4j Graph

A. Creating Nodes

Users :

```
1 LOAD CSV WITH HEADERS FROM "file:///users.csv" AS row
2 FIELDTERMINATOR '|'
3 CREATE (:Users {userId: row.userId});
```

neo4j\$ MATCH (n:Users) RETURN n LIMIT 25

Overview

Node labels

(25) Users (25)

Displaying 25 nodes, 0 relationships.

Movies:

```
1 LOAD CSV WITH HEADERS FROM "file:///movies.csv" AS row
2 FIELDTERMINATOR '|'
3 CREATE (:Movies {movieId: row.movieId, title: row.title, rating_mean:
  row.rating_mean});
```

neo4j\$ MATCH (n:Movies) RETURN n LIMIT 25

Overview

Node labels

(25) Movies (25)

Displaying 25 nodes, 0 relationships.

Genres:

```
1 LOAD CSV WITH HEADERS FROM "file:///genres.csv" AS row
2 FIELDTERMINATOR '|'
3 CREATE (:Genres {genres: row.genres});
```



```
neo4j$ MATCH (n:Genres) RETURN n LIMIT 25
```

Overview

Node labels

(19) Genres (19)

Displaying 19 nodes, 0 relationships.

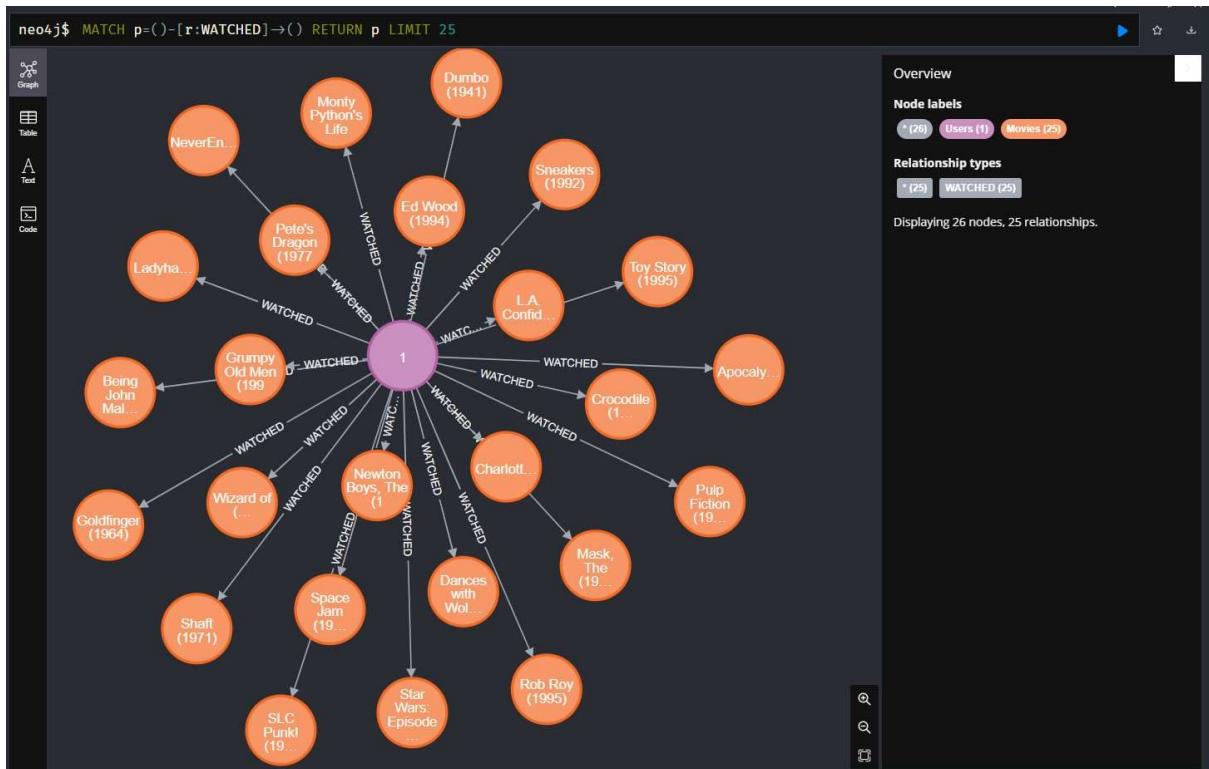
Create INDEX for nodes :

```
1 CREATE INDEX FOR (Users:userId) ON (Users.userId);
2 CREATE INDEX FOR (Movies:movieId) ON (Movies.movieId);
```

B. Relations

user-movie : WATCHED

```
1 CREATE INDEX FOR (Movies:movieId) ON (Movies.movieId);
2 LOAD CSV WITH HEADERS FROM "file:///users_movies.csv" AS row
3 FIELDTERMINATOR '|'
4 MATCH (user:Users {userId: row.userId})
5 MATCH (movie:Movies {movieId: row.movieId})
6 MERGE (user)-[:WATCHED {rating: row.rating}]->(movie);
7 CREATE INDEX FOR (Movies:movieId) ON (Movies.movieId);
```

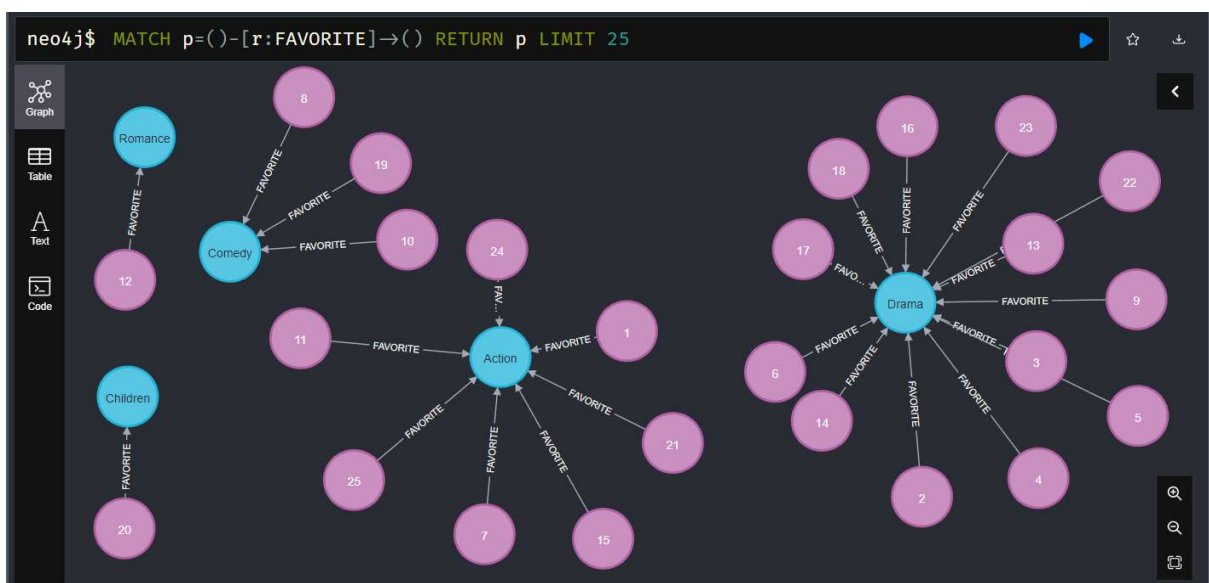


user-genre : FAVORITE

```

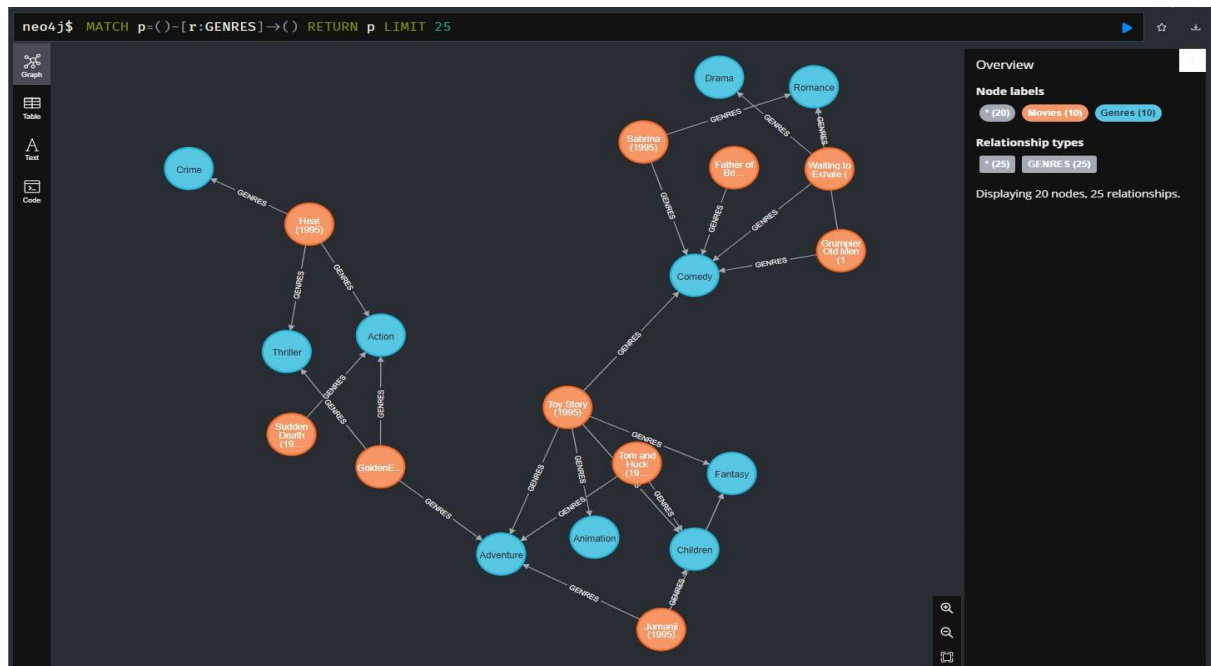
1 LOAD CSV WITH HEADERS FROM "file:///users_genres.csv" AS row
2 FIELDTERMINATOR '|'
3 MATCH (user:Users {userId: row.userId})
4 MATCH (genres:Genres {genres: row.genre})
5 MERGE (user)-[:FAVORITE]->(genres);

```



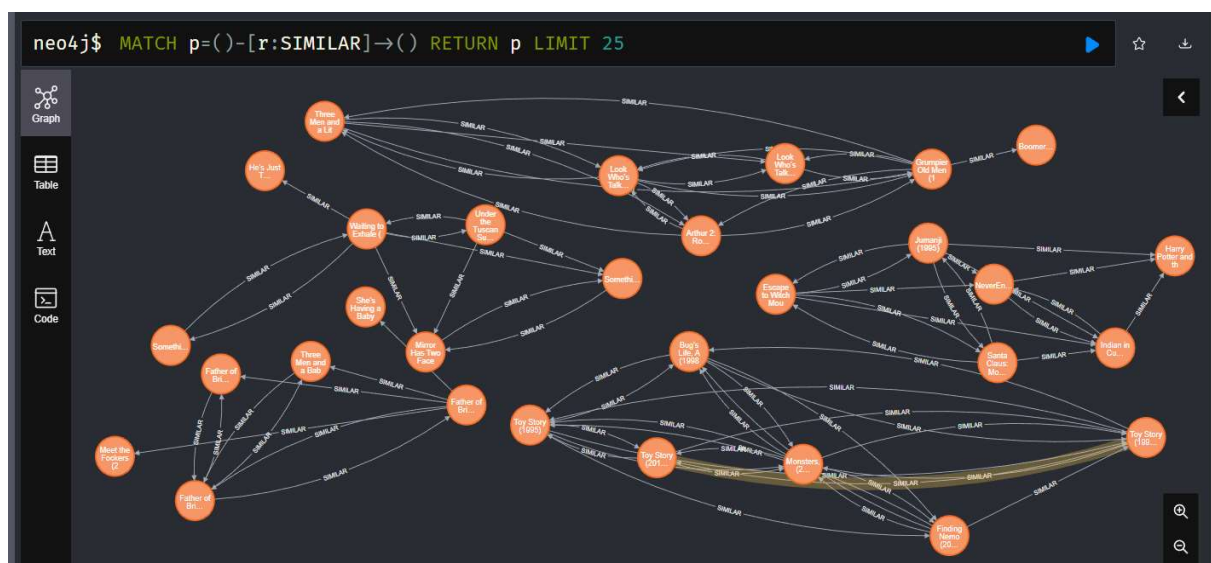
movie - genre : **GENRE**

```
1 CREATE INDEX FOR (Movies:movieId) ON (Movies.movieId);
2 LOAD CSV WITH HEADERS FROM "file:///movies_genres.csv" AS row
3 FIELDTERMINATOR '|'
4 MATCH (movie:Movies {movieId: row.movieId})
5 MATCH (genres:Genres {genres: row.genres})
6 MERGE (movie)-[:GENRES]→(genres);
```



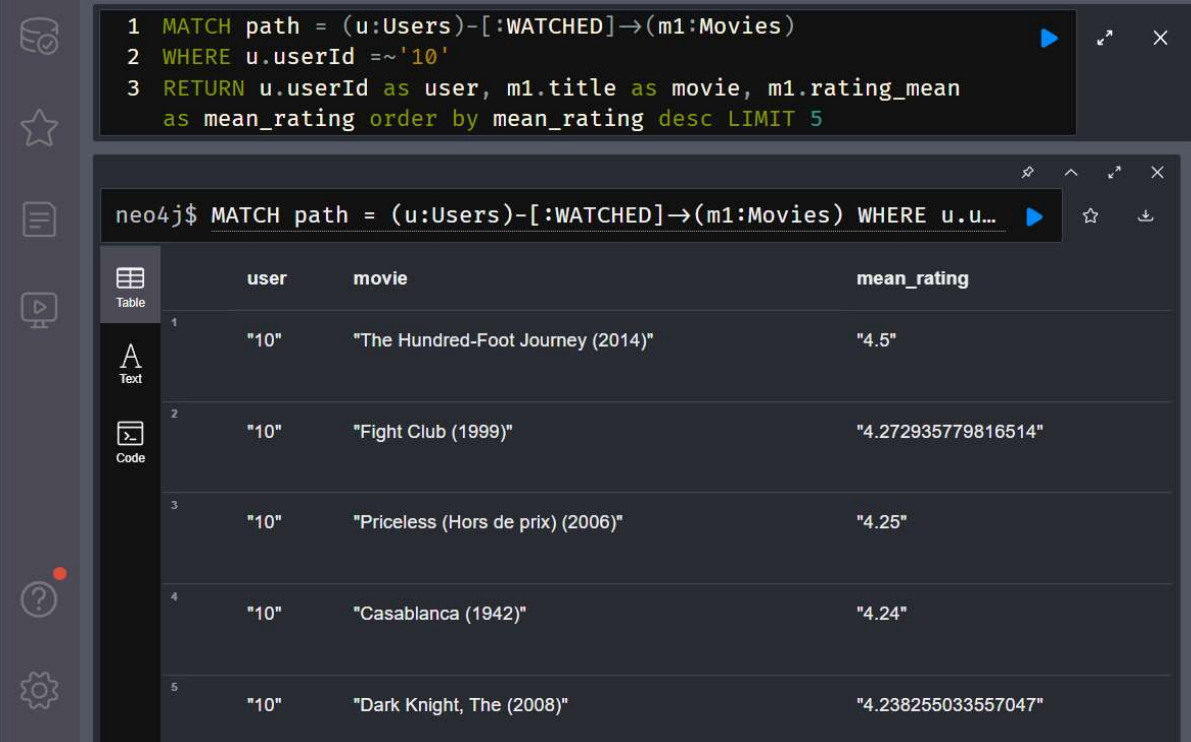
movie-movie : SIMILAR

```
1 CREATE INDEX FOR (Movies:movieId) ON (Movies.movieId);
2 LOAD CSV WITH HEADERS FROM "file:///movies_similarity.csv" AS row
3 FIELDTERMINATOR '|'
4 MATCH (movie1:Movies {movieId: row.movieId})
5 MATCH [movie2:Movies {movieId: row.sim_movieId}]
6 MERGE (movie1)-[:SIMILAR {relevance: row.relevance}]->(movie2);
```



Step 2. Basic Queries

Top 5 rated- movies watched by user 10 :



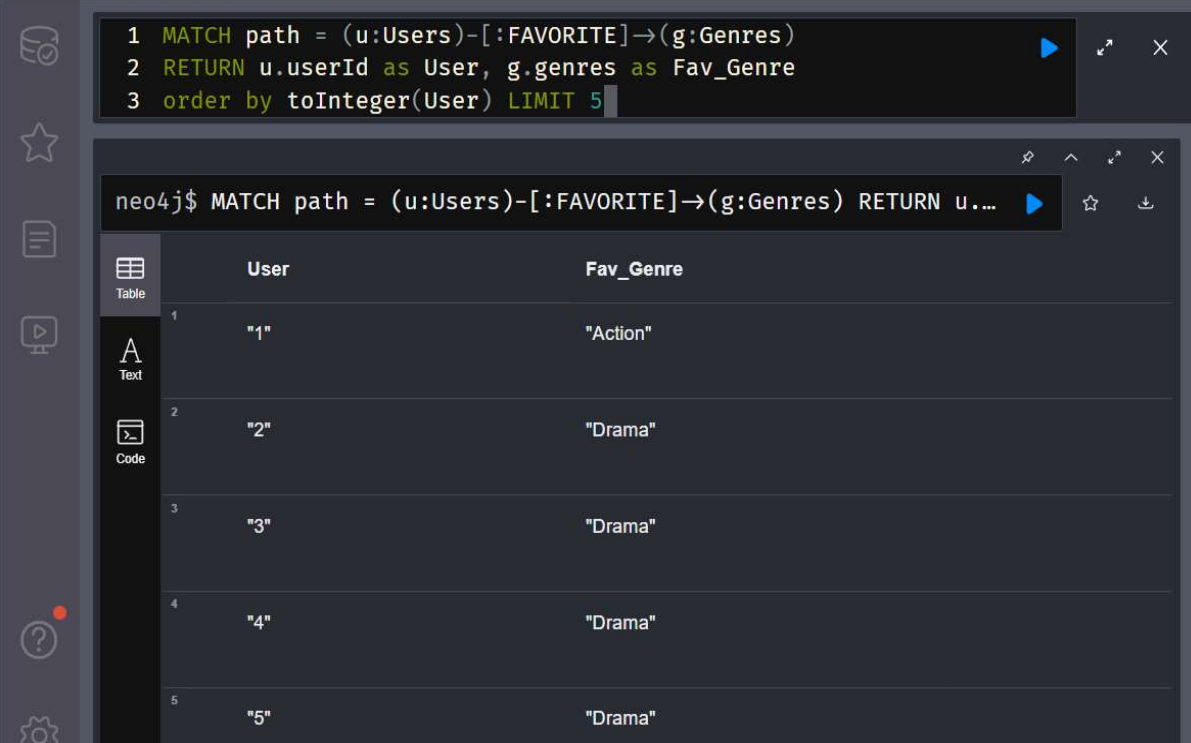
The image shows the Neo4j query interface. The top panel contains a Cypher query:

```
1 MATCH path = (u:Users)-[:WATCHED]→(m1:Movies)
2 WHERE u.userId =~ '10'
3 RETURN u.userId as user, m1.title as movie, m1.rating_mean
   as mean_rating order by mean_rating desc LIMIT 5
```

The bottom panel shows the query results in a table view:

	user	movie	mean_rating
1	"10"	"The Hundred-Foot Journey (2014)"	"4.5"
2	"10"	"Fight Club (1999)"	"4.272935779816514"
3	"10"	"Priceless (Hors de prix) (2006)"	"4.25"
4	"10"	"Casablanca (1942)"	"4.24"
5	"10"	"Dark Knight, The (2008)"	"4.238255033557047"

Favorite Genre of first 5 users:



The image shows the Neo4j query interface. The top panel contains a Cypher query:

```
1 MATCH path = (u:Users)-[:FAVORITE]→(g:Genres)
2 RETURN u.userId as User, g.genres as Fav_Genre
3 order by toInteger(User) LIMIT 5
```

The bottom panel shows the query results in a table view:

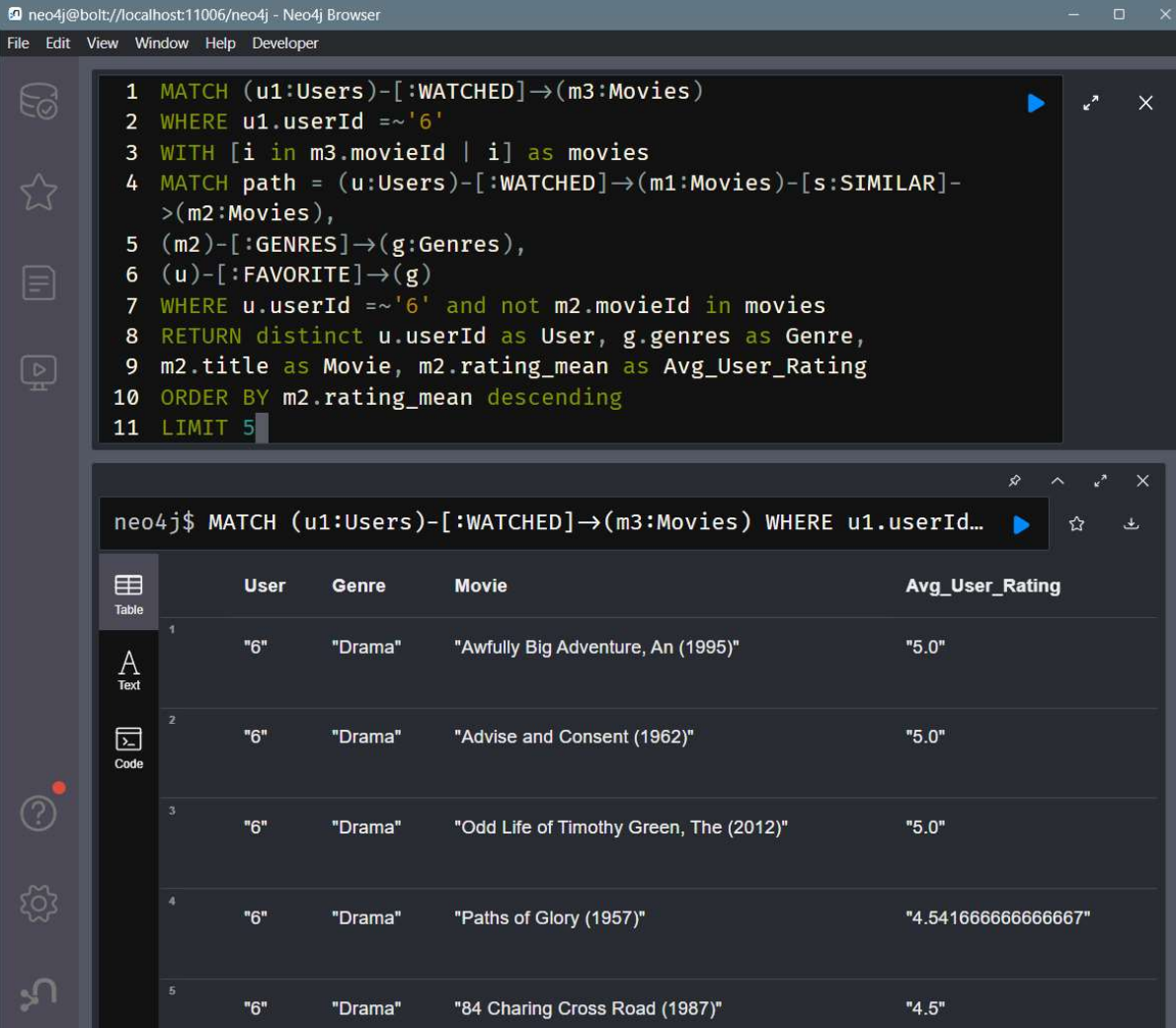
	User	Fav_Genre
1	"1"	"Action"
2	"2"	"Drama"
3	"3"	"Drama"
4	"4"	"Drama"
5	"5"	"Drama"

Step 3. Content Based Filtering:

Content-based filtering uses the category of the movies to recommend other movies which fall under the categories that the user likes, based on the ratings of their previous watched movies through movie ratings.

The following are the steps and the demonstration of Content based filtering to recommend movies to user **6**:

- a. Get the list of movies **watched** by User 6 - **List A**
- b. For List A, get the list of **similar** movies - **List B**
- c. From List B, select the movies which fall under the **favorite** genre of the user
- d. Order them by their **Avg ratings** given by the users and recommend top 5 among them.



The screenshot shows the Neo4j Browser interface. The top panel displays a Cypher query for content-based filtering. The query finds movies watched by user 6, identifies similar movies, filters by the user's favorite genre (Drama), and returns the top 5 by average rating.

```
1 MATCH (u1:Users)-[:WATCHED]→(m3:Movies)
2 WHERE u1.userId =~ '6'
3 WITH [i in m3.movieId | i] as movies
4 MATCH path = (u:Users)-[:WATCHED]→(m1:Movies)-[s:SIMILAR]->(m2:Movies),
5 (m2)-[:GENRES]→(g:Genres),
6 (u)-[:FAVORITE]→(g)
7 WHERE u.userId =~ '6' and not m2.movieId in movies
8 RETURN distinct u.userId as User, g.genres as Genre,
9 m2.title as Movie, m2.rating_mean as Avg_User_Rating
10 ORDER BY m2.rating_mean descending
11 LIMIT 5
```

The bottom panel shows the results of the query in a table format:

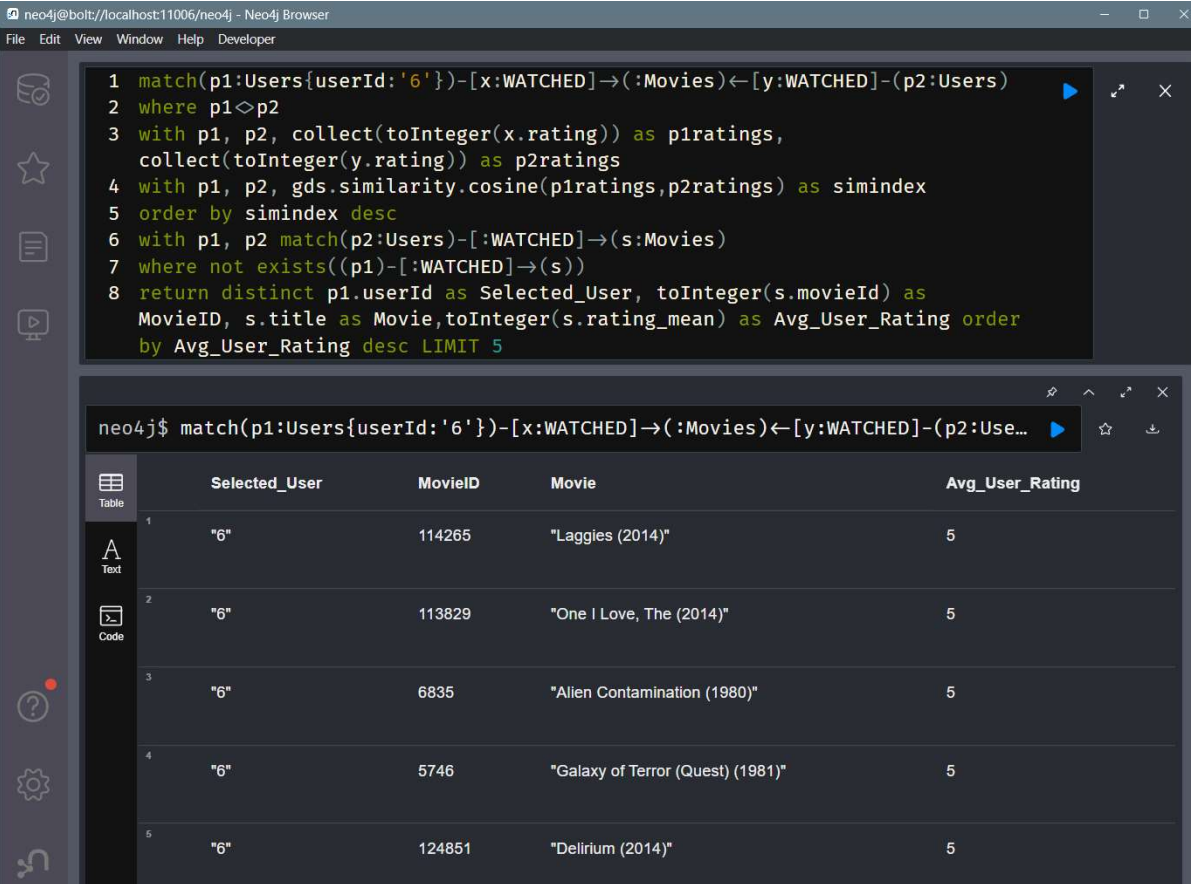
	User	Genre	Movie	Avg_User_Rating
1	"6"	"Drama"	"Awfully Big Adventure, An (1995)"	"5.0"
2	"6"	"Drama"	"Advise and Consent (1962)"	"5.0"
3	"6"	"Drama"	"Odd Life of Timothy Green, The (2012)"	"5.0"
4	"6"	"Drama"	"Paths of Glory (1957)"	"4.541666666666667"
5	"6"	"Drama"	"84 Charing Cross Road (1987)"	"4.5"

Step 4. Collaborative Filtering:

Collaborative Filtering uses **similarities** between users and movies simultaneously, to provide recommendations. This allows for serendipitous recommendations: i.e. collaborative filtering models can recommend a movie to user A based on the interests of a similar user B.

The following are the steps and the demonstration of Collaborative filtering to recommend movies to user **6**:

- a. Get the list of users who watched the **same** movies as selected user
- b. Find the **cosine similarity** between the rating of the selected user and the ratings of the other users for the movies they've watched and order them according to the **cosine similarity index**
- c. Get the list of the movies these users have watched but **aren't watched** by the selected user and recommend the **top 5** according to the ratings.



The screenshot shows the Neo4j Browser interface. The top panel displays a Cypher query for collaborative filtering. The bottom panel shows the results of the query in a table format.

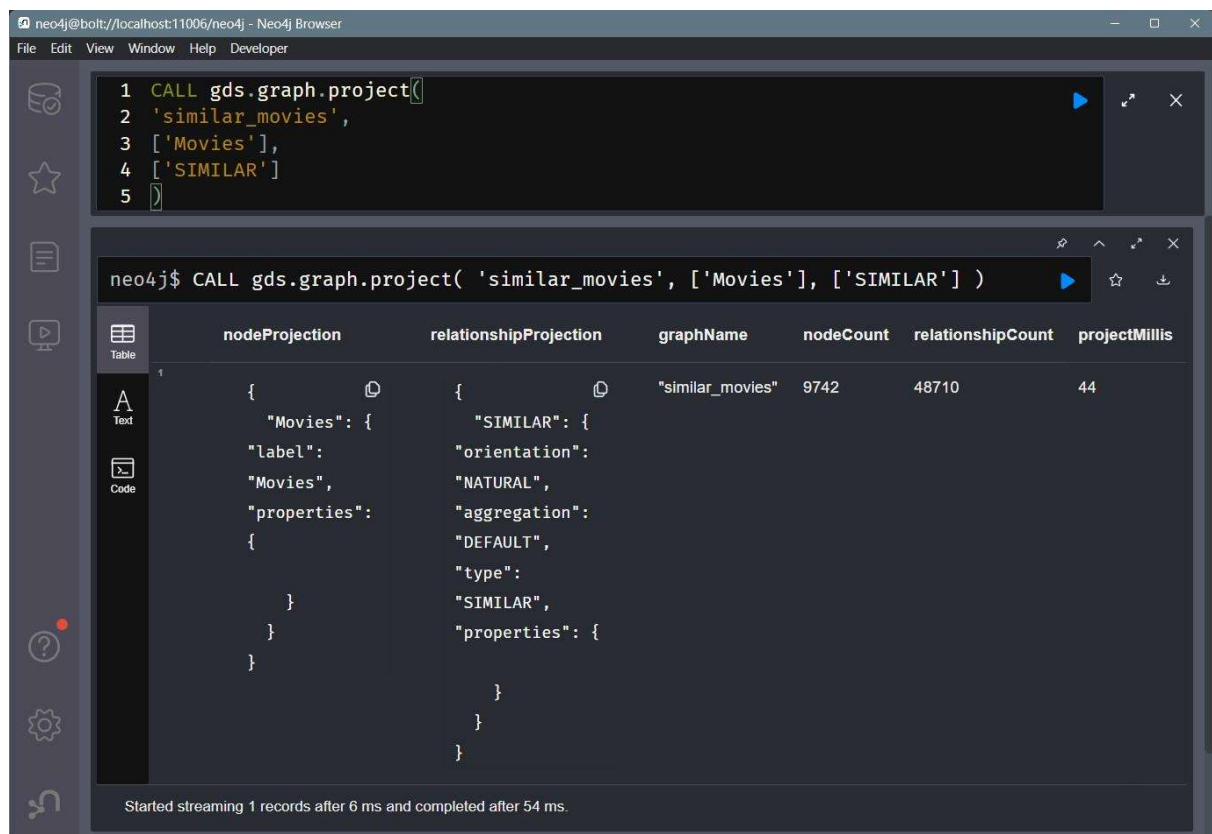
```
1 match(p1:Users{userId:'6'})-[x:WATCHED]→(:Movies)←[y:WATCHED]-(p2:Users)
2 where p1<≠p2
3 with p1, p2, collect(toInteger(x.rating)) as p1ratings,
  collect(toInteger(y.rating)) as p2ratings
4 with p1, p2, gds.similarity.cosine(p1ratings,p2ratings) as simindex
5 order by simindex desc
6 with p1, p2 match(p2:Users)-[:WATCHED]→(s:Movies)
7 where not exists((p1)-[:WATCHED]→(s))
8 return distinct p1.userId as Selected_User, toInteger(s.movieId) as
  MovieID, s.title as Movie,toInteger(s.rating_mean) as Avg_User_Rating order
  by Avg_User_Rating desc LIMIT 5
```

	Selected_User	MovieID	Movie	Avg_User_Rating
1	"6"	114265	"Laggies (2014)"	5
2	"6"	113829	"One I Love, The (2014)"	5
3	"6"	6835	"Alien Contamination (1980)"	5
4	"6"	5746	"Galaxy of Terror (Quest) (1981)"	5
5	"6"	124851	"Delirium (2014)"	5

Step 5. Betweenness Centrality:

Betweenness Centrality is a measure of centrality in a graph based on shortest paths. For each node, it is the **number of shortest paths** (between pairs of nodes) that pass through the node.

First create a **graph** of **similar** movies:



As seen in the above image, the nodeCount for the graph is **9742**.

We calculated the betweenness centrality measure for each node, then normalized the measures using the **normalization factor = $(n-1)(n-2)$** , here $n = 9742$ and fetched the **top 5** nodes with the highest betweenness centrality:

