

The Entertainment Knowledge Graph

An IR project based on Knowledge Graph Structures for movies and actors/actresses analysis.

Group Members

1. Abhinav Talari (S20180010003)
2. Utkarsh Ajay Aditya (S20180010182)

Project Details

Project Title: The Entertainment Knowledge Graph

Domain of Project: Entertainment (Movies, television programs and celebrity content)

Project Description

This project focuses on end-to-end utilization of knowledge graphs and semantics in the field of entertainment. Our entities include movies, tv shows, actors/actresses and all their attributes and relationships. Our project is based on entity ranking, entity linking, entity-based retrieval models, entity recommendation, document filtering and knowledge graph population. The knowledge graph will represent movies and tv shows as well as actors, genres and the complex interrelationships among them, which will help us retrieve more informational value compared to free text.

Algorithms/Approaches Used

Individual Breakdown for each member.

1. Abhinav Talari(S20180010003) Tasks

A. Web Scraping(Actors/Actresses data)

Associated Files: Scraper_For_Actors.py

This python code file is used for scraping actor data from Wikipedia and storing it. Scraped Data from Wikipedia by Even going through the underlying links which increased the corpus exhaustively

The Functions Used :

1. **wiki_scrape(page)**

This Looks for a page on Wikipedia extracts it,parses it and the underlying link and keeps on fetching data from those links.

For Example for a query as `wiki_scrape("Robert Downey Jr.")` Made it crawl 800+ links

2. `wiki_page(page)`

This just fetches the data from the very first hit it gets on Wikipedia
the query for this function needs to be very precise

For Example `wiki_page("The Avengers")` will not fetch the movie synopsis of the MCU's Avengers but just links and names of other movies/entities with the same name or context.

B. Knowledge Graph for Actors

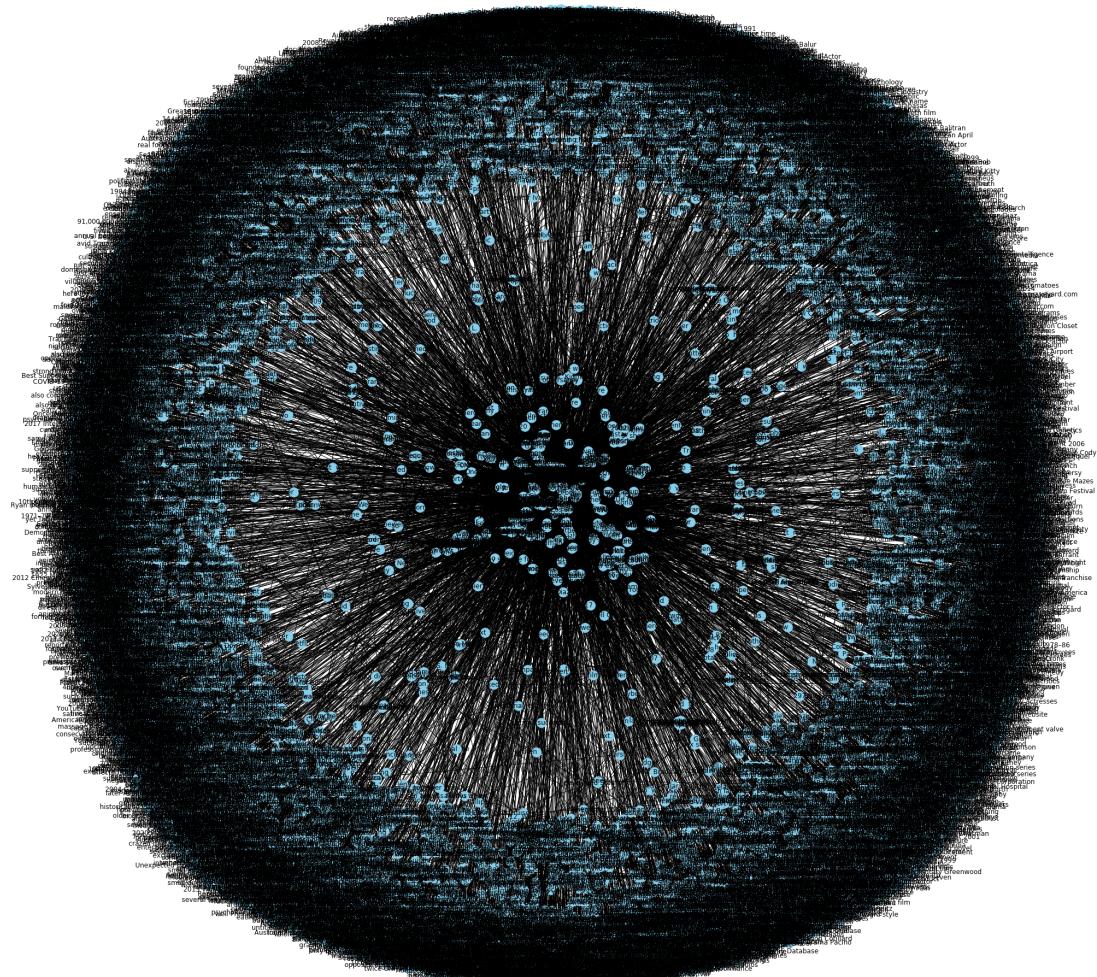
Associated Files: `Knowledge_Graph_for_Actors.py`

Based on the data crawled and collected for Actors and Movies , I created a knowledge graph to perform "Entity-Entity pair based on Relation" Query.

The sentences were tokenized and the Entity-Relation-Entity were identified and put into the Knowledge Graph

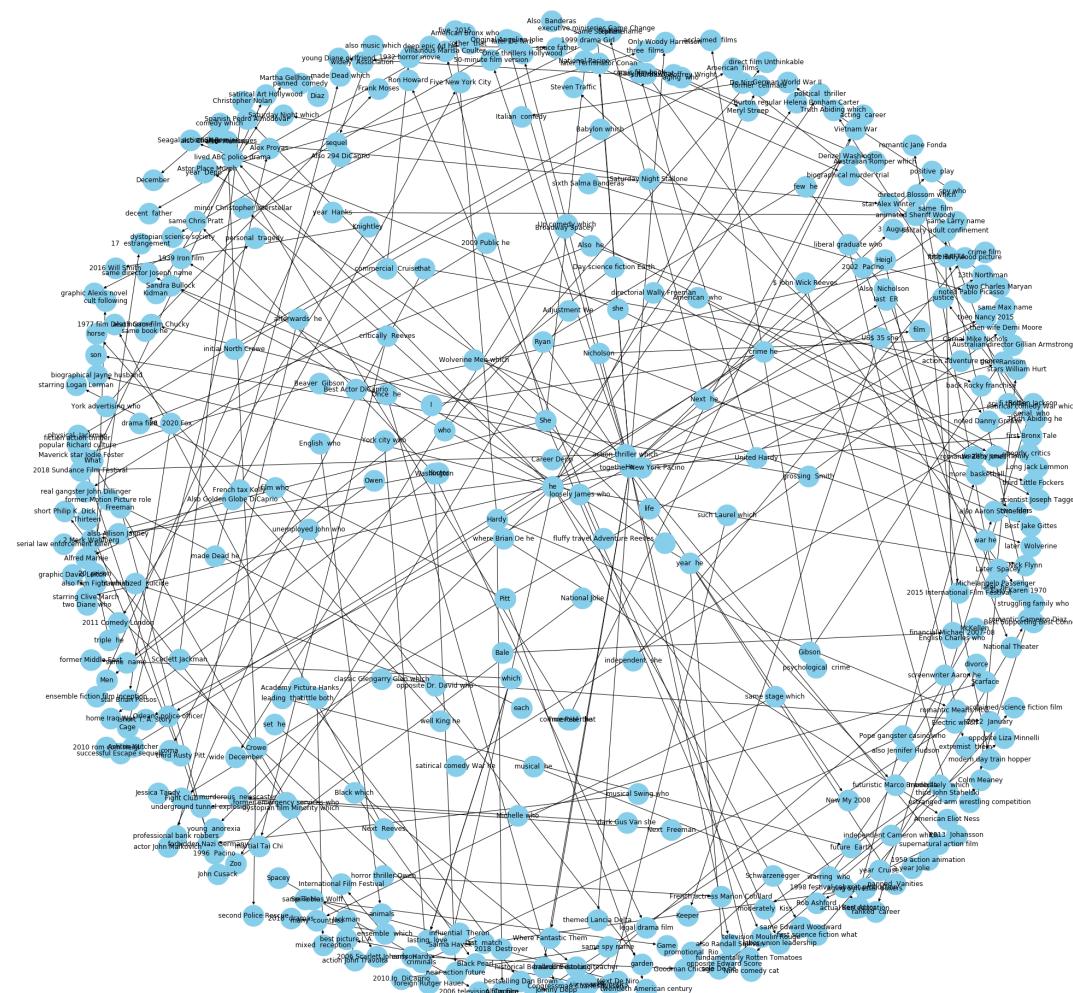
The knowledge graph was visualized using Networkx

Output: Below is the knowledge graph representing top 250 actors along with their extracted entities (subject-object pairs) and the predicates (relation between entities).



The networkx library was used to create a network from this data frame. The nodes will represent the entities and the edges or connections between the nodes will represent the relations between the nodes.

It is going to be a directed graph. In other words, the relation between any connected node pair is not two-way, it is only from one node to another. For example, "John eats pasta"



The above figure represents a sub-knowledge graph based on the relation "starred in" for actors.

C. Ranked Retrieval model(Actors)

Associated Files: RankedRetrieval_actors.py

Based on the Data Scrapped I created a Ranked Retrieval System using tf-idf ranking.

Since the Boolean Model only fetches complete matches, it doesn't address the problem of the documents being partially matched. The Vector Space Model solves this problem by introducing vectors of index items each assigned with weights. The weights are ranged from positive (if matched completely or to some extent) to negative (if unmatched or completely oppositely matched) if documents are present. Term Frequency - Inverse Document Frequency (tf-idf) is one of the most popular techniques where weights are terms (e.g. words, keywords, phrases etc.) and dimensions is number of words inside corpus.

Therefore I created a ranked retrieval based on tf-idf ranking

Output

1. Below is the output showing for the query "Angelina Jolie"

```

File Edit Selection View Go Run Terminal Help rr.py - IR Project - Visual Studio Code
documentizer.py crawl.py rr.py
rr.py sim_temp = sim_temp + Documents_dict["doc"+str(i)]|query_list[j]|Query_doc_freq[query_list[j]]
TERMINAL SQL CONSOLE DEBUG CONSOLE PROBLEMS OUTPUT 1: cmd ...
for details
C:\Users\khinav\Desktop\IR Project>python rr.py
Enter your query:
Angelina
-> 50 > 0.871929824561403
Angelina -> 7 -> 6.12280701754385964
|| q || = 0.885747758922341
Similarity scores:
{"doc1": 0.06873489640290444, "doc2": 0.05028328078667279, "doc3": 0.0927672210876984, "doc4": 0.0776882361762366, "doc5": 0.0657958025993699, "doc6": 0.0594292379015691, "doc7": 0.04137474625548831, "doc8": 0.0490008845610802, "doc9": 0.05661776365452542, "doc10": 0.05427371972209, "doc11": 0.080702046468119, "doc12": 0.08897591851623611, "doc13": 0.0740031737165263, "doc14": 0.08540278918165266, "doc15": 0.065573, "doc16": 0.050874183072047, "doc17": 0.052000002422817083, "doc18": 0.0545471060872047, "doc19": 0.0590563230602, "doc20": 0.08806709408791, "doc21": 0.08802699181679, "doc22": 0.053200002422817083, "doc23": 0.0513467389516314, "doc24": 0.088191002422817083, "doc25": 0.05750408181088776, "doc26": 0.0434594162161896, "doc27": 0.0525115073976668, "doc28": 0.082759195724677024326, "doc29": 0.0465530126473806085, "doc30": 0.05706423902667664, "doc31": 0.0765761563854699, "doc32": 0.05270481456459452, "doc33": 0.05296899166811271, "doc34": 0.052714677024326, "doc35": 0.05706423902667664, "doc36": 0.0658761052976665, "doc37": 0.043259665911101496, "doc38": 0.05849984686348584, "doc39": 0.04782140711975736, "doc40": 0.06521347696622687, "doc41": 0.06707169926408979, "doc42": 0.0568897816718145, "doc43": 0.053165237104899, "doc44": 0.235328801356158582, "doc45": 0.436906113671280865, "doc46": 0.0747956326888911, "doc47": 0.05363422051578979, "doc48": 0.0761, "doc49": 0.0995338954375564, "doc50": 0.0671838523262545}
Ranking based on similarity scores:

Simon Baker.json is the DocumentRanked at 1 with Similarity Score :0.232380001359582
Jack Nicholson.json is the DocumentRanked at 2 with Similarity Score :0.11992656332305982
Vin Diesel.json is the DocumentRanked at 3 with Similarity Score :0.09953309054375564
Guy Pearce.json is the DocumentRanked at 4 with Similarity Score :0.09378029421852666
Antonio Banderas.json is the DocumentRanked at 5 with Similarity Score :0.0927672210876984
Denzel Washington.json is the DocumentRanked at 6 with Similarity Score :0.08897551851623611
George Clooney.json is the DocumentRanked at 7 with Similarity Score :0.08806709408791
John Travolta.json is the DocumentRanked at 8 with Similarity Score :0.08877446700405781
Clive Owen.json is the DocumentRanked at 9 with Similarity Score :0.08072046468119
Arnold Schwarzenegger.json is the DocumentRanked at 10 with Similarity Score :0.0776882361762366
Mel Gibson.json is the DocumentRanked at 11 with Similarity Score :0.0765761563854699
Tom Hardy.json is the DocumentRanked at 12 with Similarity Score :0.0761071652865095
Tom Cruise.json is the DocumentRanked at 13 with Similarity Score :0.0747956326888911
Meg Ryan.json is the DocumentRanked at 14 with Similarity Score :0.072759152490541
Johnny Depp.json is the DocumentRanked at 15 with Similarity Score :0.06870026695181679
Al Pacino.json is the DocumentRanked at 16 with Similarity Score :0.06873489640290444
Will Smith.json is the DocumentRanked at 17 with Similarity Score :0.0671838623320545,
Sam Rockwell.json is the DocumentRanked at 18 with Similarity Score :0.06670026695180879
Richard Gere.json is the DocumentRanked at 19 with Similarity Score :0.0665876156386695
Brad Pitt.json is the DocumentRanked at 20 with Similarity Score :0.0657958025993699
Gerard Butler.json is the DocumentRanked at 21 with Similarity Score :0.0657366886874186
Samuel L. Jackson.json is the DocumentRanked at 22 with Similarity Score :0.06521347696622687
Kevin Spacey.json is the DocumentRanked at 23 with Similarity Score :0.0639004222817803
Keisha Reeves.json is the DocumentRanked at 24 with Similarity Score :0.0636700405781
Bruce Willis.json is the DocumentRanked at 25 with Similarity Score :0.063202231180018
Robert Downey Jr.json is the DocumentRanked at 26 with Similarity Score :0.05849984686348584
Keira Knightley.json is the DocumentRanked at 27 with Similarity Score :0.057594081871088776
Orlando Bloom.json is the DocumentRanked at 28 with Similarity Score :0.0570764239026767664
Scarlett Johansson.json is the DocumentRanked at 29 with Similarity Score :0.0568897816718145
Charlize Theron.json is the DocumentRanked at 30 with Similarity Score :0.05657366886874186
Katherine Heigl.json is the DocumentRanked at 31 with Similarity Score :0.05511403653153414
Christian Bale.json is the DocumentRanked at 32 with Similarity Score :0.05427371972209741
Tom Hanks.json is the DocumentRanked at 33 with Similarity Score :0.05363422051575079
Sean Connery.json is the DocumentRanked at 34 with Similarity Score :0.0531652371042889
Nicolas Cage.json is the DocumentRanked at 35 with Similarity Score :0.05296899166811271
Nico Kidman.json is the DocumentRanked at 36 with Similarity Score :0.05270481456459452
Kerry Washington.json is the DocumentRanked at 37 with Similarity Score :0.05270481456459452
Matt Damon.json is the DocumentRanked at 38 with Similarity Score :0.05270481456459452
Ian McEwan.json is the DocumentRanked at 39 with Similarity Score :0.051934710993585566
Angelina Jolie.json is the DocumentRanked at 40 with Similarity Score :0.05028328078667279
Catherine Zeta-Jones.json is the DocumentRanked at 41 with Similarity Score :0.0490008845610802
Edward Norton.json is the DocumentRanked at 42 with Similarity Score :0.04782140711975736

```

2.Utkarsh Ajay Aditya(S20180010182) Tasks

A. Web Scraping(Movies)

Associated Files: imdb_moviedata_scrape.py

The above code is used for scraping movie details of over 83,000 movies and storing it in final_dataset.csv

The collected fields for a given movie were:

Title, Date, Run Time, Genre, Rating Score, Description, Director, Stars, VotesGross

A	B	C	D	E	F	G	H	I	J	K	L	M	N	O	P
1	0 title	original_title	year	date_published	genre	duration	country	language	director	writer	production_companies	actors	description	avg_vote	votes
2	1 Miss Jerry	Miss Jerry	1894	1894-10-09	Romance	45	USA	None	Alexander B. Alexander B. Alexander B. Blanche Baylies, Will	The adventures of a fema	5.9	154			
3	2 The Story of the Kelly Gang	The Story of the Ke	1906	26-12-1906	Biography, Crime, I	70	Australia	None	Charles Tait Charles Tait J. and N. Ta Elizabeth Tait, John	True story of notorious Ai	6.1	589			
4	3 Den sorte drøm	Den sorte drøm	1911	19-08-1911	Drama	53	Germany, Denmark	Urban Gad Urban Gad, Fotorama	Asta Nielsen, Valder	Two men of high rank are	5.8	188			
5	4 Cleopatra	Cleopatra	1912	13-11-1912	Drama, History	100	USA	English	Charles L. G. Victorien Sa Helen Gardi Helen Gardner, Pear	The fabled queen of Egypt	5.2	446			
6	5 L'Inferno	L'Inferno	1911	06-03-1911	Adventure, Drama,	68	Italy	Italian	Francesco B. Dantone Alighi Milano Film Salvatore Papa, Artu	Loosely adapted from Dai	7	2237			

B. Knowledge Graph for Movies

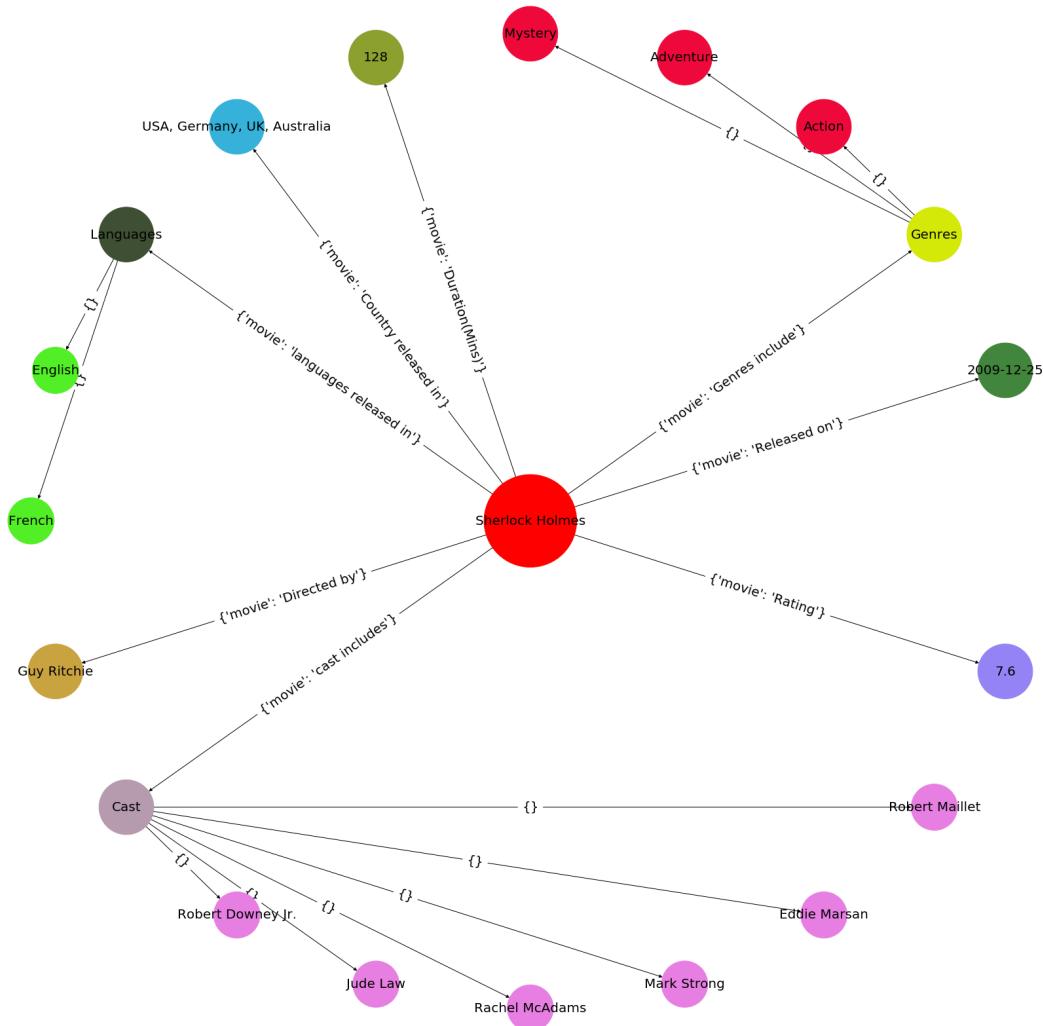
Associated Files: KnowledgeGraph_movies_recommender_system.py

It works on the basis of final_dataset_imdb.csv.

NOTE : All the graphs would be saved as .pdf file in the code directory. One has to view the pdf file where he can zoom and every node will be visible in HD resolution.

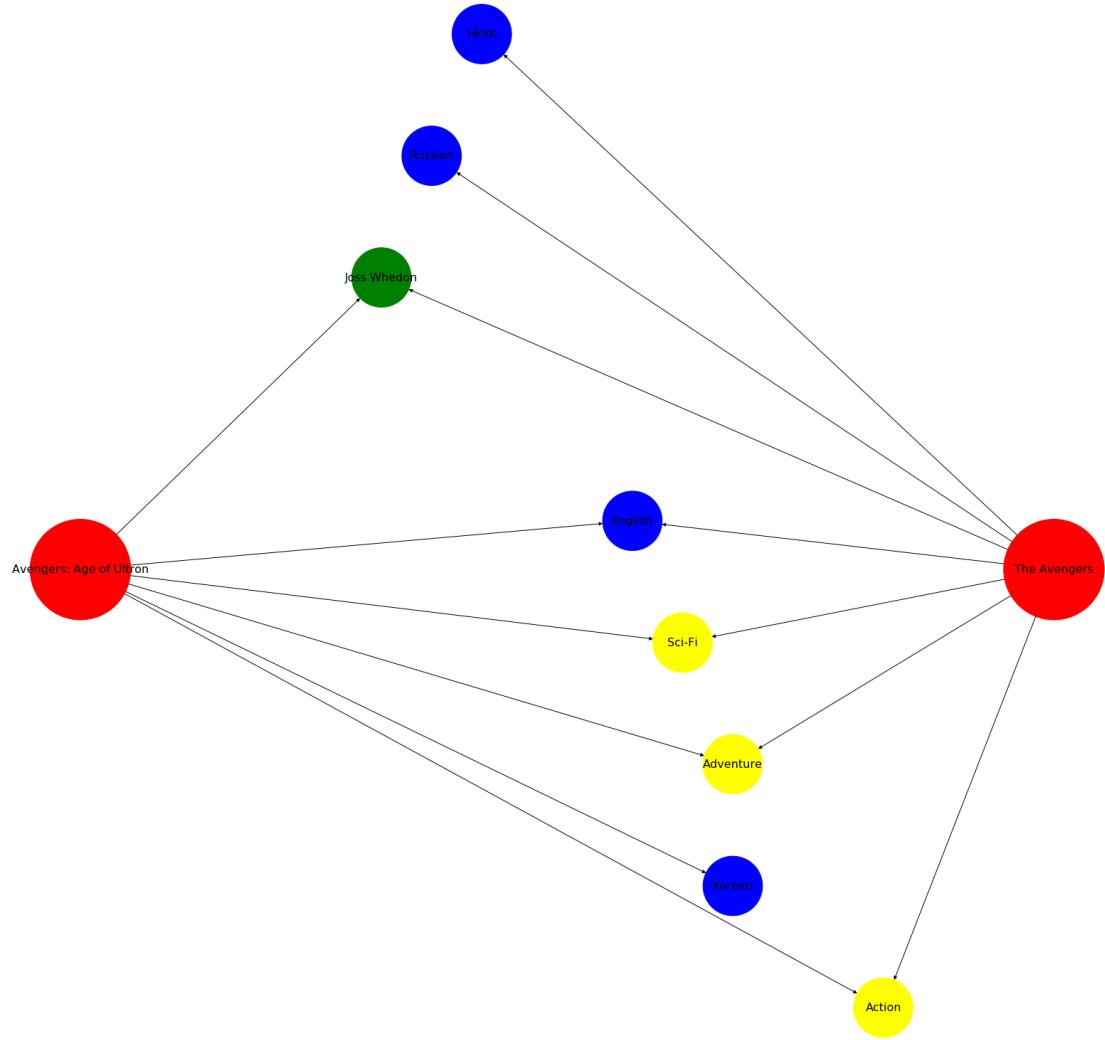
The above code file defines 3 different types of knowledge graph structures related to movies:

1. Movie detail graph for a single movie



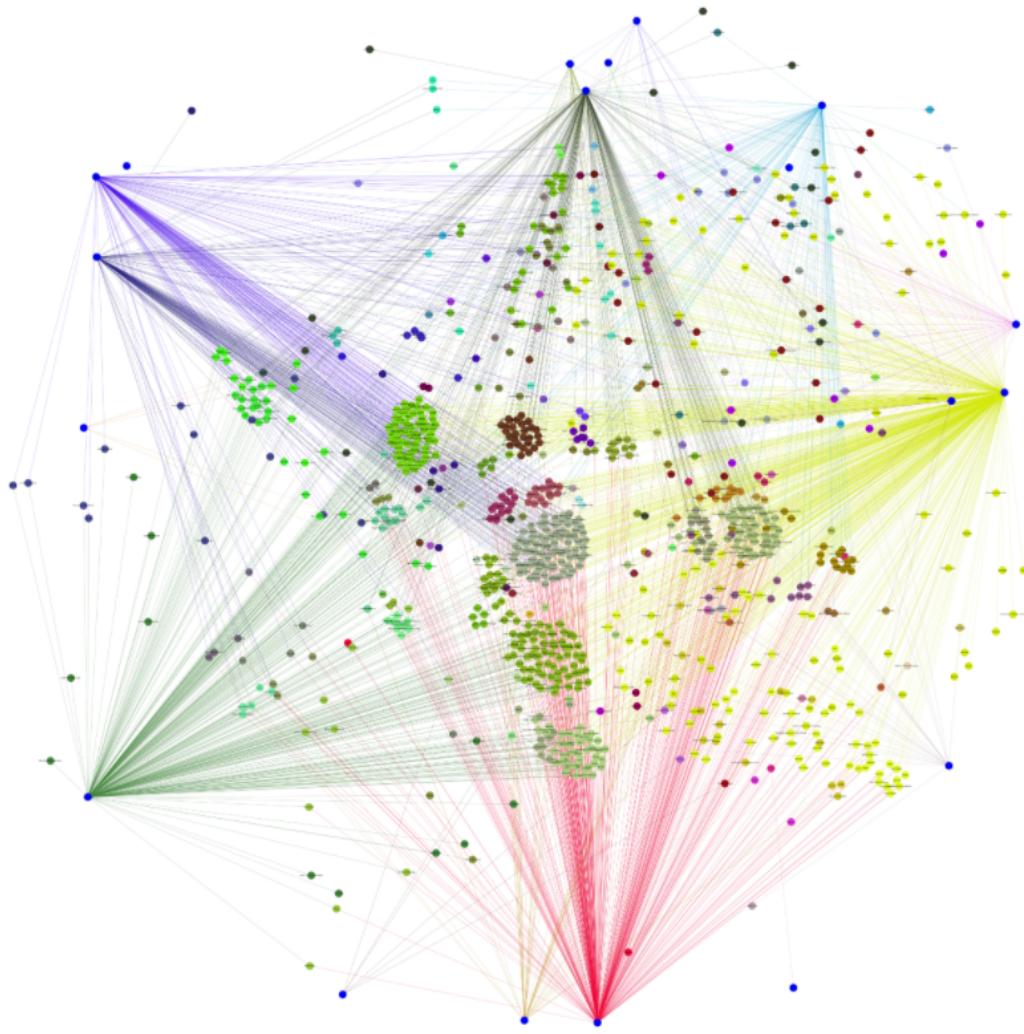
The above figure shows the detailed view of a given movie name showing all its details like director, cast, release date, etc.

2. Movie Comparison graph between two movies.



The above figure shows the comparison between 2 input movie names on basis of genre, languages and director.

3.Knowledge graph of 1000 movies.



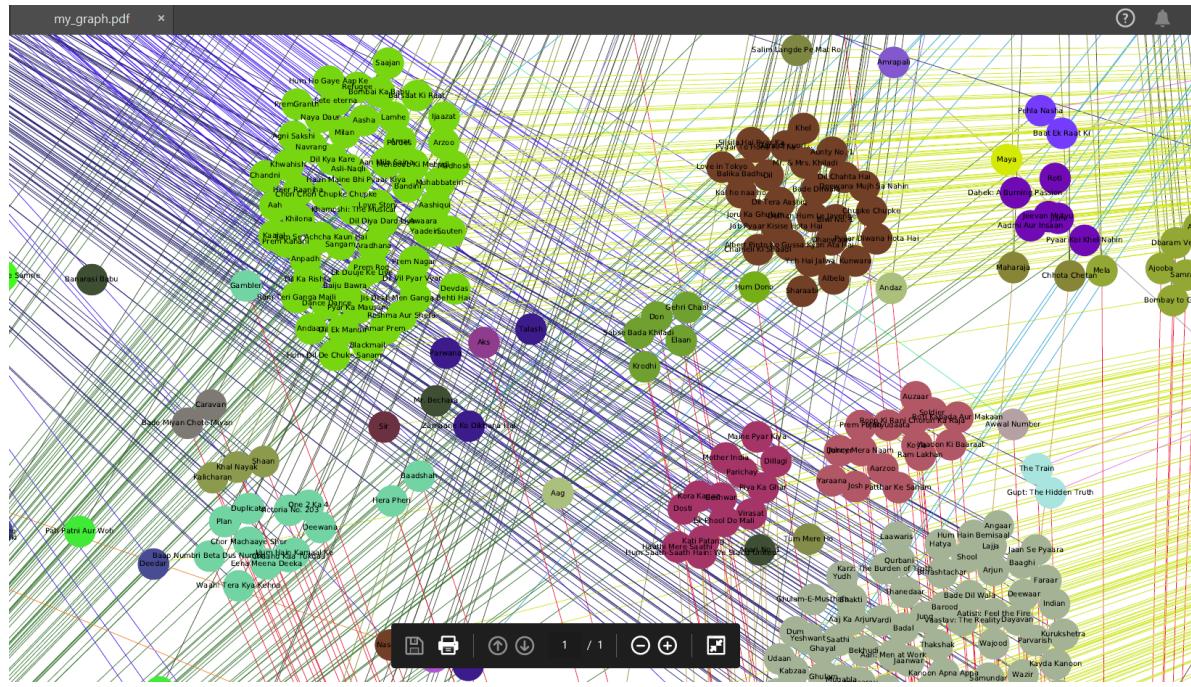
The blue nodes near the edges of the graph represent genres, edges of same color represent same genre movies.

A movie node's color is derived by the combination of its genres color. Therefore movies with similar characteristics would be of same color and in the same color.

Properties of graph :

- a. There is homogeneity between the movies belonging to the same color. (Represented by same color.)
- b. There is heterogeneity among movies from different clusters. (Represented by different colors.)

Zoomed pic of a cluster.



Here you can see the movies which belong to different colors. Each color represents a color.

Such a graph provides us insight about what movie an user can like based on few of his judgments.

C. Movie Recommendation System

Associated Files: KnowledgeGraph_movies_recommender_system.py

Given a movie name it will return a list of top 10 similar movies along with their matching scores.

It is based on an empirically derived ranking system which takes into account all the movie details like genres, crew,

director, writer, year, other user reviews, etc.

Example : For movie - The Avengers

```
(base) PS C:\Users\utkar\Desktop\Untitled Folder> python '.\KnowledgeGraph (1).py'
1. Enter 1 for movie detail graph
2. Enter 2 for movie similarity graph between 2 movies.
3. Enter 3 for knowledge graph of 1000 movies.
4. Enter 4 for Movie Recommendations(Top 10 along with score)
5. Enter 0 to exit
4
Enter Movie Name(Case Sensitive): The Avengers
Rank 1 ==> Avengers: Age of Ultron
Score: 65.0000000722684

Rank 2 ==> Avengers: Infinity War
Score: 54.0000000796486

Rank 3 ==> Iron Man 2
Score: 48.0000000706284

Rank 4 ==> Captain America: Civil War
Score: 48.0000000644241

Rank 5 ==> Avengers: Endgame
Score: 41.0000000754786

Rank 6 ==> Iron Man
Score: 39.0000000920706

Rank 7 ==> Captain America: The Winter Soldier
Score: 39.0000000719085

Rank 8 ==> Iron Man Three
Score: 38.0000000739816

Rank 9 ==> Captain America: The First Avenger
Score: 31.0000000723473

Rank 10 ==> Spider-Man: Homecoming
Score: 31.000000049406

Best movie match is : Avengers: Age of Ultron
```

Example : For movie mission impossible

```
Enter Movie Name(Case Sensitive): Mission: Impossible
Rank 1 ==> Mission: Impossible III
Score: 34.0000000032431

Rank 2 ==> Mission: Impossible - Fallout
Score: 34.00000000280981

Rank 3 ==> Mission: Impossible - Rogue Nation
Score: 33.0000000033453

Rank 4 ==> Mission: Impossible II
Score: 33.00000000308586

Rank 5 ==> Mission: Impossible - Ghost Protocol
Score: 30.00000000449833

Rank 6 ==> The A-Team
Score: 27.00000000243142

Rank 7 ==> Largo Winch
Score: 27.00000000013296

Rank 8 ==> The Hunt for Red October
Score: 26.00000000175033

Rank 9 ==> Drop Zone
Score: 26.00000000018953

Rank 10 ==> The Rock
Score: 25.00000000307112

Best movie match is : Mission: Impossible III
```

D. Spell Correction for the entertainment domain

[Associated Files](#): KnowledgeGraph_movies_recommender_system.py

Entertainment_domain_spell_correction.py is an autocorrect program file to correct incorrect

- a) movie name or TV series name
- b) actors/actresses name
- c) director/write name

Example For query "Robept Downez" and another query "emmug watsoz"

```
PS C:\Users\utkar\Desktop\Untitled Folder> python .\Entertainment_domain_spell_correction.py
Enter Query: Robept Downez

Candidate set of robept
{'robert', 'robpet', 'robetp'}

Most Suitable Candidate is robert

Candidate set of downez
{'downe', 'downer', 'downze', 'downey', 'downes', 'dornez'}

Most Suitable Candidate is downey
PS C:\Users\utkar\Desktop\Untitled Folder> python .\Entertainment_domain_spell_correction.py
Enter Query: emmg watsoz

Candidate set of emmg
{'emms', 'emmy', 'emme', 'emmi', 'emm', 'emma', 'emig'}

Most Suitable Candidate is emma

Candidate set of watsoz
{'watson', 'watszo', 'watosz'}

Most Suitable Candidate is watson
PS C:\Users\utkar\Desktop\Untitled Folder>
```

Implementation of spell checker

- 1) Store Dictionary in Trie
- 2) Query is processed word by word.
- 3) If word is in dictionary, move ahead.
- 4) If word is not in dictionary, display the candidate set.
- 5) Also display the post probable candidate by calculating the posterior probability

References

1. Wikipedia
2. IMDB
3. networkx
4. **Research Paper** : Knowledge Graphs: An Information Retrieval Perspective