Finding Key Influential Actor in Movie Industry

BADS 7201, 7202

Potchara V Ativit C

Vinitwattanakoon Chaninchodeuk

6210412003 6220412019

Outline

- Introduction
- Dataset
- Analysis
- Conclusion

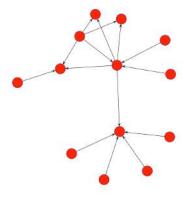
Introduction

Issues for film makers

Too many actors to choose from

Some Co-Lead actors has no chemistry

Difficultly choosing the right actors for the movie



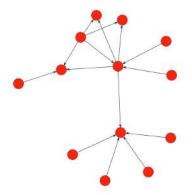
- Node = Actors
- Edge = Co-starring

Objective

Too many actors to choose from Ranking actors

Some Co-Lead actors has no chemistry ———— Choose actors within same community

Difficultly choosing the right actors for the movie



- Node = Actors
- Edge = Co-starring

Process Flow

Data preprocessing

Data exploratory

Filter data, transform data

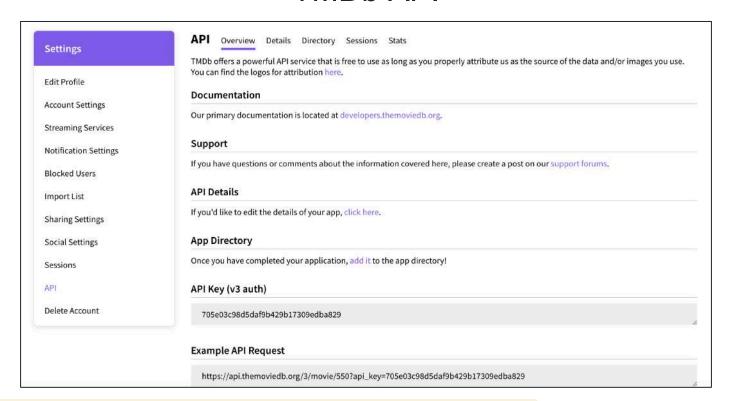
- Create network graph
- Analysis

Centrality

Community Detection

Dataset

TMDb API



TMDb API

- tmdbsimple Version 2.6.6
- API Key

Register and verify an account to get API key

API key = <u>'705e03c98d5daf9b429b17309edba829'</u>

Finding

Movie Dataset: tmdb id, title, revenue, budget, product country, imdb rating, imdb votecount

Crew Dataset : crew id, name, department, job, character, popularity

Create Movie.csv

```
1 import tmdbsimple as tmdb
 2 import requests
 3 import locale
 4 import pandas as pd
 5 imdb id list = []
 6 tmdb id list = []
 7 revenue list = []
 8 budget list = []
 9 title list = []
10 product country list = []
11 vote_average_list = []
12 vote count list = []
13 Keyerror = []
15 def parsing(obj):
       mystr = '
       j = 0
18
       for i in obj:
19
           if i != 0:
20
               mystr += ','
           mystr += i['iso 3166 1']
           j += 1
       return mystr
```

```
27 for tmdb id lists in imdb id :
28
     trv :
29
30
        url = 'https://api.themoviedb.org/3/movie/'+ tmdb id lists
31
        +'?api key=705e03c98d5daf9b429b17309edba829&language=en-US'
32
33
34
        print('fetching: ',tmdb id lists)
35
36
                                                                                      Jason Format
        response = requests.get(url)
        response json = response.json()
        imdb_id = response_json("imdb_id")
                                                                          "production_companies": [
        tmdb id = response json("id")
                                                                             "id" - 508
        revenue = response_json["revenue"]
                                                                             "logo_path": "/7PzJdsLGLR7oW4J0J5Xcd0pHGRq.png".
        budget = response_json["budget"]
                                                                             "name": "Regency Enterprises",
4.3
        title = response json["title"]
                                                                             "origin_country": "US"
       product_country = response_json['production_countries']
        vote average = response json['vote average']
                                                                             "id": 711.
        vote count = response json['vote count']
                                                                             "logo_path": null.
47
                                                                             "name": "Fox 2000 Pictures",
48
                                                                             "origin_country": ""
        product country str = parsing(product country)
        imdb_id_list.append(imdb_id)
        tmdb id list.append(tmdb id)
                                                                             "logo_path": null.
        revenue list.append(revenue)
                                                                             "name": "Taurus Film",
        budget list.append(budget)
                                                                             "origin_country": ""
       title list.append(title)
        product country list.append(product country str)
                                                                             "id": 54050
        vote average list.append(vote average)
                                                                             "logo_path": null.
        vote_count_list.append(vote_count)
                                                                             "name": "Linson Films".
58
                                                                             "origin_country": ""
59
     except KeyError :
        print('KeyError-----', tmdb id lists)
                                                                             "id": 54051
61
       C = Keyerror.append(tmdb id lists)
                                                                             "logo_path": null.
                                                                             "name": "Atman Entertainment",
                                                                             "origin_country": ""
```

Create Movie.csv

	tmdb_id	title	product_country	startYear	imdb_votecount
0	400531.0	Frivolinas	ES	2014	15.0
1	11232.0	Kate & Leopold	US	2001	78160.0
2	602986.0	The Tango of the Widower and Its Distorting Mi	CL	2020	76.0
3	299782.0	The Other Side of the Wind	FR,IR,US	2018	5995.0
4	339338.0	Toula, or the Genie of the Water	NE,DE	2017	25.0
***			(64.6)	***	***
159527	655187.0	Akelarre	AR,ES,FR	2020	97.0
159528	622812.0	The Secret of China	CN	2019	11.0
159529	NaN	NaN	NaN	2019	5.0
159530	NaN	NaN	NaN	2013	NaN
159531	NaN	NaN	NaN	2017	NaN
59532 rc	ws × 5 colu	mns			

Create Crew.csv

```
1 import tmdbsimple as tmdb
 2 import requests
 3 import locale
 4 import pandas as pd
 5 id list = []
 6 imdb id list = []
 7 name list = []
 8 popularity list = []
 9 Keverror = []
10 Valueerror = []
11 Typeerror = []
12 count = 0
14 for crew id in my list:
15 try:
16
17
      url = 'https://api.themoviedb.org/3/person/'+ crew id
18
      +'?api key=705e03c98d5daf9b429b17309edba829&language=en-US'
19
20
       count += 1
21
      print('fetching: ', crew id ,'counting' , count)
22
```

```
26
      response = requests.get(url)
27
      response json = response.json()
28
      id = response json["id"]
      imdb id = response json["imdb id"]
30
      name = response json["name"]
31
      popularity = response json["popularity"]
      gender = response json["gender"]
32
33
34
      imdb id list.append(imdb id)
35
      id list.append(id)
      name list.append(name)
36
37
       popularity list.append(popularity)
38
      gender list.append(gender)
39
40
    except KeyError :
41
      print('KeyError----> ', crew id)
42
      C = Keyerror.append(crew id)
43
       pass
```

Create Crew.csv

	movie_id	crew_id	name	department	job	character	popularit
0	11232	5344	Meg Ryan	Acting	Actor	Kate McKay	8.36
1	11232	6968	Hugh Jackman	Acting	Actor	Leopold	17.61
2	11232	23626	Liev Schreiber	Acting	Actor	Stuart Besser	7.49
3	11232	33654	Breckin Meyer	Acting	Actor	Charlie McKay	7.79
4	11232	10871	Natasha Lyonne	Acting	Actor	Darci	5.35
	***	252			***	we	
125134	615177	11677	Terry Chen	Acting	Actor	Officer Steve Choi	3.24
125135	615177	62709	Garfield Wilson	Acting	Actor	Steve (Angry Man)	1.09
125136	615177	1406026	Arielle Tuliao	Acting	Actor	Sharon	1.47
125137	615177	228893	Jag Bal	Acting	Actor	Snarky IT Guy	1.40
125138	615177	1904663	Sonia Sunger	Acting	Actor	Présentatrice du journal télévisé	0.60

Dataset

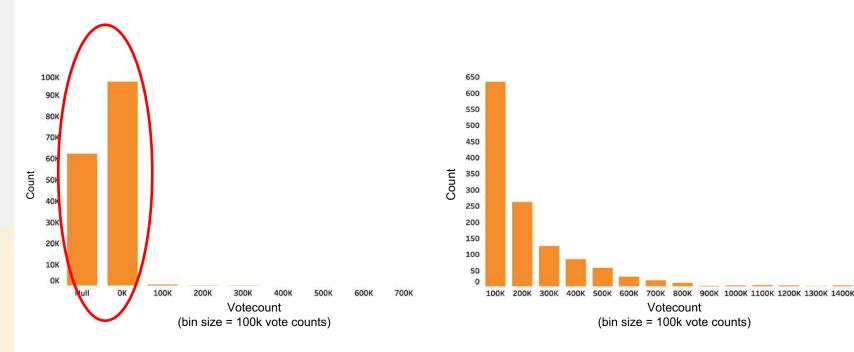
movies dataset Movie.csv

	tmdb_id	title	product_country	startYear	imdb_votecount
0	400531.0	Frivolinas	ES	2014	15.0
1	11232.0	Kate & Leopold	us	2001	78160.0
2	602986.0	The Tango of the Widower and Its Distorting Mi	CL	2020	76.0
3	299782.0	The Other Side of the Wind	FR,IR,US	2018	5995.0
4	339338.0	Toula, or the Genie of the Water	NE,DE	2017	25.0

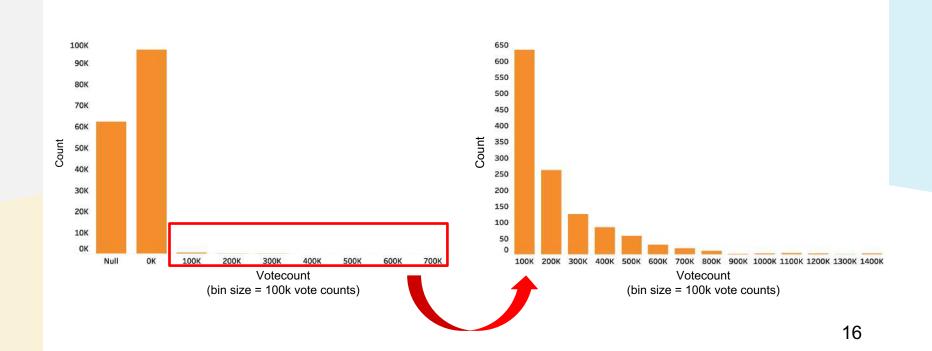
actors dataset Crew.csv

		movie_id	crew_id	name	department	job	character	popularity
	1311	503919	10965	Chris Columbus	Production	Executive Producer	NaN	2.871
	7114	503919	5293	Willem Dafoe	Acting	Actor	Thomas Wake	15.189
t	13884	503919	40142	Robert Fernandez	Sound	Sound Re-Recording Mixer	NaN	0.600
	31106	503919	376	Arnon Milchan	Production	Executive Producer	NaN	1.702
	31791	503919	1197157	Paul Rutledge	Crew	Stunts	NaN	1.380
	32006	503919	11288	Robert Pattinson	Acting	Actor	Thomas Howard	11.369
	32351	503919	11183	Eleanor Columbus	Production	Executive Producer	NaN	2.407

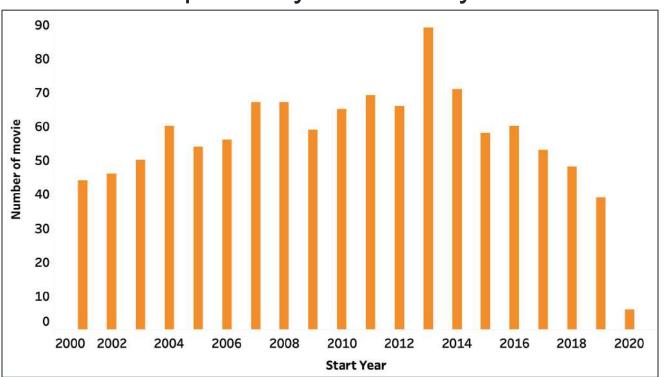
Exploratory Data Analysis







Exploratory Data Analysis



What movies are filtered?

Vote Count

Select a vote count of movie that has more than 100,000 of total user votes.

Product Country

Only select produced in US.

Too many actors

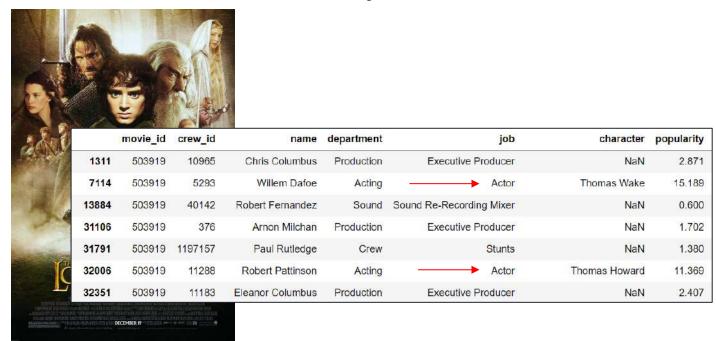


The Lord Of The Rings
The Fellowship Of The Ring.

popularity	name	
13.101	Elijah Wood	139
13.451	Ian McKellen	140
7.959	Liv Tyler	141
10.152	Viggo Mortensen	142
8.999	Sean Astin	143
12.751	Cate Blanchett	144
7.129	John Rhys-Davies	145
2.978	Billy Boyd	146
7.423	Dominic Monaghan	147
6.675	Orlando Bloom	148
10.467	Christopher Lee	149
6.959	Hugo Weaving	150
22.236	Sean Bean	151
4.658	Ian Holm	152
8.840	Andy Serkis	153

244	Tim McLachlan	1.013
245	Liz Merton	0.600
246	Arnold Montey	0.994
247	Dean Morganty	0.694
248	Greg Morrison	0.600
249	Blair Morton	0.600
250	Andrew Munro	0.728
251	David J. Muzzerali	0.600
252	Shane Rangi	2.588
253	Chris Reid	0.600
254	Steve Reinsfield	0.600
255	Larry Rew	0.600
256	Grant Roa	1.220
257	Thomas Robins	0.694
258	Vincent Roxburgh	0.600
259	Chris Ryan	0.600
260	Paul Shapcott	0.600
261	Samuel E. Shore	0.840

Too many actors



The Lord Of The Rings
The Fellowship Of The Ring.

What actors are filtered?

Popularity of actors

Select actor with a popularity value of more than or equal to 7

Job

Only select job = actor on Crew.csv

Dataset

Original data base on periods 20 years.

159,532 movies

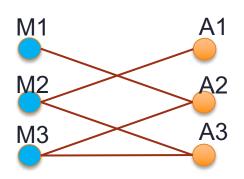
125,139 actors

Filtered (Product Country, Vote count, Popularity)

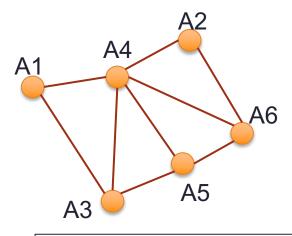
868 movies

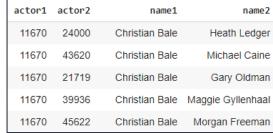
1,091 actors

Data preprocessing, data tranformation



movie_id	crew_id	name	character
49026	3894	Christian Bale	Bruce Wayne / Batman
49026	3895	Michael Caine	Alfred Pennyworth
49026	64	Gary Oldman	Commissioner James Gordon
49026	1813	Anne Hathaway	Selina Kyle / Catwoman
49026	2524	Tom Hardy	Bane
49026	8293	Marion Cotillard	Miranda Tate / Talia al Ghul
49026	24045	Joseph Gordon-Levitt	John Blake







Convert Bipartite to Unipartite Graph

```
1 #create [movie-actor] object
2 mov actor = df.groupby(['movie'])['actor 1'].apply(','.join)
3 mov actor
```

```
movie
                                      31209,34569,65181
                         62656,39289,15026,44008,11159
                         26925, 26237, 49865, 39159, 30826
                                29526,34012,35906,18593
                                 31209,34569,6929,46072
        11379,3048,61695,36637,12194,34084,26129,32054
3532
3533
                                      10867,11220,23840
3541
                                         152,8879,36313
3545
                                                  54887
3551
                                52866,18722,43620,12747
Name: actor 1, Length: 1091, dtype: object
```



```
1 #find combination between 2 pair of actors from [mov-actor] object
2 mylist = []
3 for i in mov actor:
     actor list = i.split(',')
     comb = combinations(actor list, 2)
     for i in list(comb):
         #convert combination inside into integer so it can be sorted
         edge pair = list(map(int, i))
         mylist.append(sorted(edge pair))
```

	actor1	actor2
0	31209	34569
1	31209	65181
2	34569	65181
3	39289	62656
4	15026	62656

Create weight from sum of repeated edge pairs

```
#create ['repeat'] col that will be used as weight of edges
temp = df_merged2.groupby('concat').agg({'concat':'size'}).rename(columns={'name':'repeat'})
temp.columns = ['repeat']
print(temp.shape)
temp.sort values('repeat', ascending=False)
(10801, 1)
                                repeat
                        concat
     (Chris Evans, Stan Lee)
                                    13
  (Robert Downey Jr., Stan Lee)
  (Samuel L. Jackson, Stan Lee)
 (Scarlett Johansson, Stan Lee)
 (Alan Rickman, Robbie Coltrane)
 (David Duchovny, Vince Vaughn)
 (David Duchovny, Winona Ryder)
```

Create weight from sum of repeated edge pairs

Weight \

	acton1	acton2	nama1	nama?	concat	noncat	ohunch
	actori	actor2	name1	name2	Concat	repeat	ebunch
1777	11379	57874	Chris Evans	Stan Lee	(Chris Evans, Stan Lee)	13	[Chris Evans, Stan Lee, 13]
1713	52680	57874	Robert Downey Jr.	Stan Lee	(Robert Downey Jr., Stan Lee)	10	[Robert Downey Jr., Stan Lee, 10]
1733	54954	57874	Samuel L. Jackson	Stan Lee	(Samuel L. Jackson, Stan Lee)	10	[Samuel L. Jackson, Stan Lee, 10]
3637	55565	57874	Scarlett Johansson	Stan Lee	(Scarlett Johansson, Stan Lee)	9	[Scarlett Johansson, Stan Lee, 9]
958	52549	54022	Robbie Coltrane	Rupert Grint	(Robbie Coltrane, Rupert Grint)	8	[Robbie Coltrane, Rupert Grint, 8]
4033	24198	29526	Helena Bonham Carter	Jim Carrey	(Helena Bonham Carter, Jim Carrey)	1	[Helena Bonham Carter, Jim Carrey, 1]
4034	17795	29526	Dustin Hoffman	Jim Carrey	(Dustin Hoffman, Jim Carrey)	1	[Dustin Hoffman, Jim Carrey, 1]
4035	27132	29526	Jane Lynch	Jim Carrey	(Jane Lynch, Jim Carrey)	1	[Jane Lynch, Jim Carrey, 1]
4036	7007	43388	Billy Connolly	Meryl Streep	(Billy Connolly, Meryl Streep)	1	[Billy Connolly, Meryl Streep, 1]
10800	12747	43620	Clémence Poésy	Michael Caine	(Clémence Poésy, Michael Caine)	1	[Clémence Poésy, Michael Caine, 1]
10801 ro	ws × 7 co	lumns					

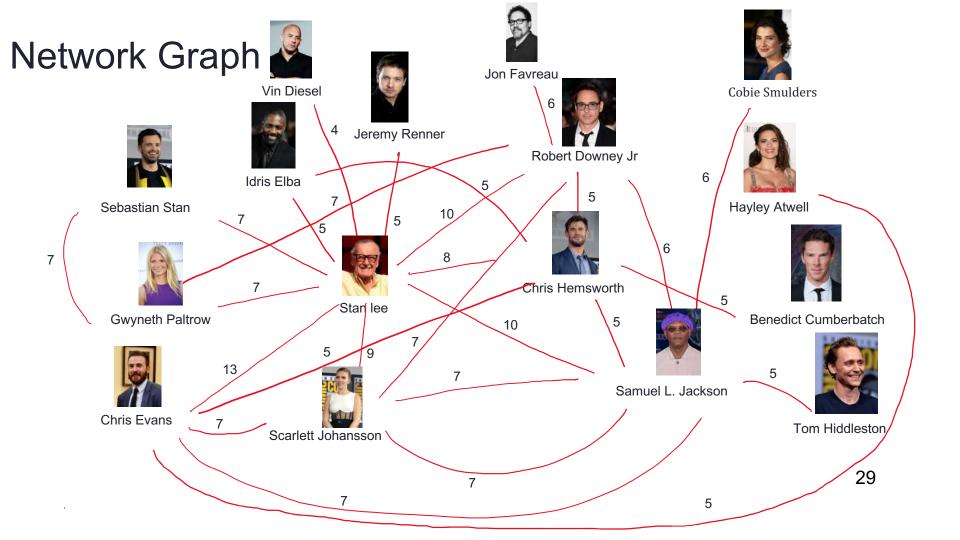
Create networkx graph

```
[47] #create function to pass the edge pair and weight into nx.add edges from
     def edge generator(reader):
         reader
         for row in reader:
              row[2] = float(row[2]) # convert numeric weight to float
             yield (row[0],
                     row[1],
                     dict(weight=row[2]))
                                                                  actor
                                                                            name
                                                                      Aaron Eckhart
                                                                         Aaron Paul
                                                                   152
[48] #create undirected graph
                                                                   214 Abbie Cornish
     G = nx.Graph()
                                                                   282 Abigail Breslin
     #add_nodes
                                                                   548 Adam Sandler
     G.add nodes from(actor table.name.tolist())
     #add edge pairs with weight from edge generator function
     G.add edges from(edge generator(df edge.ebunch.tolist()))
```

Graph Information

```
print(nx.info(G))

Is G graph connected: True
Name:
Type: Graph
Number of nodes: 864
Number of edges: 10800
Average degree: 25.0000
Actor
```



Network Graph Analysis

Centrality

Centrality Concept

Centrality	Main idea
Degree Centrality	Ratio number of degree link of a node
Closeness Centrality	Ratio of path length of the shortest path between the node and all other nodes (Euclidean Distance)
Betweenness Centrality	Ratio number of times a node acts as a bridge along the shortest path between two other nodes

Centrality Concept

Centrality	Main idea	
Degree Centrality	Ratio number of degree link of a node	
Closeness Centrality	Ratio of path length of the shortest path between the node and all other nodes (Euclidean Distance)	
Betweenness Centrality	Ratio number of times a node acts as a bridge along the shortest path between two other nodes	

Betweenness Centrality

$$C_B(v) = \sum_{s,t \in V} \frac{\sigma(s,t|v)}{\sigma(s,t)}$$

V is a set of vertex

 $\sigma(s,t)$ = total number of shortest part between node s and t

 $\sigma(s,t|v)$ = total number of shortest part between node s and t that pass through v

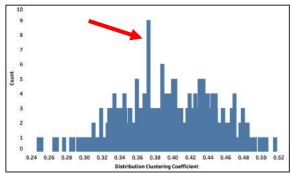
Betweenness Centrality Score

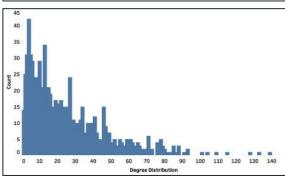
#compute betweenness centrality
btw = nx.betweenness_centrality(G,k=100)
dict_btw = dict(btw)

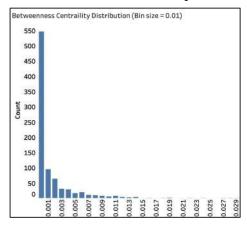
Name	Betweenness Centrality Score
Samuel L. Jackson	0.0317
Hugh Jackman	0.0265
Natalie Portman	0.0248
Liam Neeson	0.0248
Dwayne Johnson	0.0206
Stan Lee	0.0189
Ryan Reynolds	0.0182
Jason Statham	0.0165
Tilda Swinton	0.0165
Scarlett Johansson	0.0157
Vince Vaughn	0.0152
Bruce Willis	0.0143
Matt Damon	0.0139
Ben Affleck	0.0138
Channing Tatum	0.0132
Ben Stiller	0.0130
Jason Bateman	0.0126
Jastin Theroux	0.0120
Josh Brolin	0.0117
Mark Wahlberg	0.0115

Node Distribution:

Degree, Clustering Coefficient, Betweenness Centrality







	degree	clust_coeff	btw_cent
max	139	0.52	0.0298
min	1	0.25	0.0000
mean	25	0.39	0.0019
std	22	0.05	0.0035

Community Detection

Louvain Algorithm

$$\Delta Q(i \to C) = \left[\frac{\sum_{in} + k_i, in}{2m} - \left(\frac{\sum_{tot} + k_i}{2m} \right)^2 \right] - \left[\frac{\sum_{tot} in}{2m} - \left(\frac{\sum_{tot} tot}{2m} \right)^2 - \left(\frac{k_i}{2m} \right)^2 \right]$$

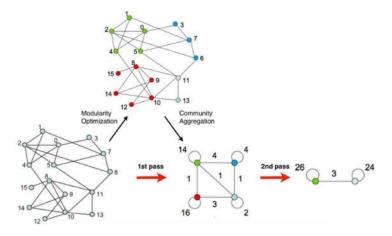
 $\sum in$ = sum of link weights between nodes in C

 $\sum tot$ = sum of all link weights of nodes in C

 k_i , in = sum of link weights between node i and C

 k_i = sum of all link weights of node i

$$\Delta Q = \Delta Q(i \rightarrow C) + \Delta Q(D \rightarrow i)$$



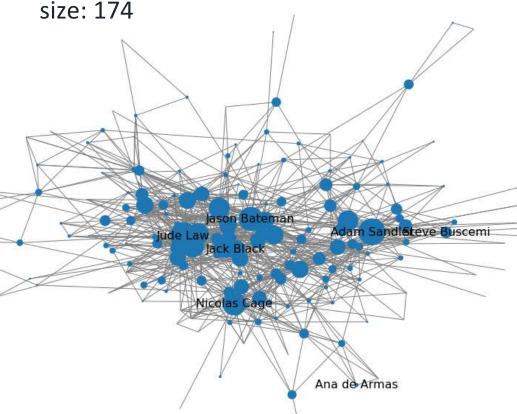
Louvain Algorithm

"All Stars" Community group 0 size: 71 Angelina Jolie Christopher Lloyd Clint Eastwood Tom Hardy Leonardo DiCaprio Liam Neeson Charging Tatum Levitt Michael Caine Anne Hathaway 40 Manuel Ferrara

"Support & Comedy Stars"

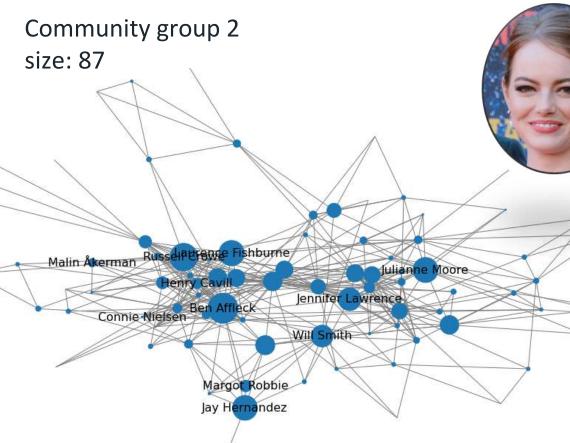
Community group 1

size: 174





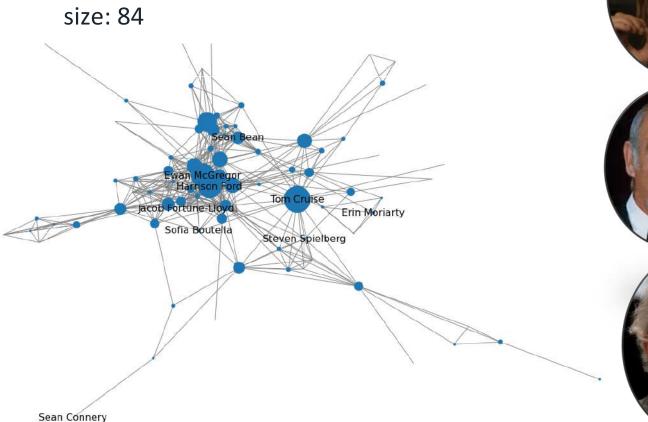
"Lead Stars 1"

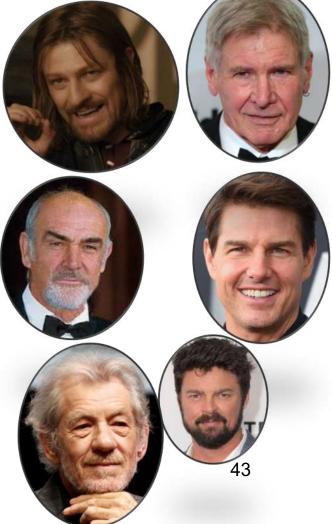


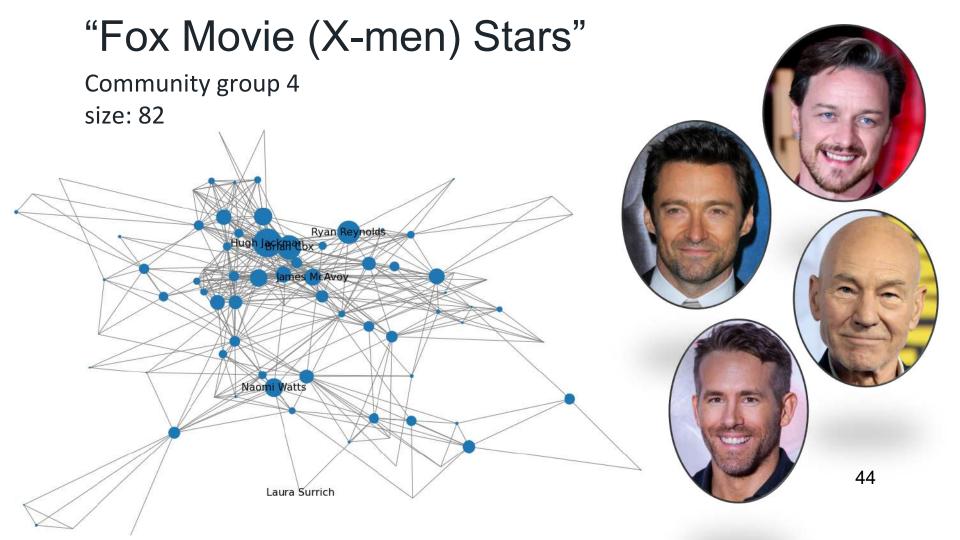


"Old Inactive Stars"

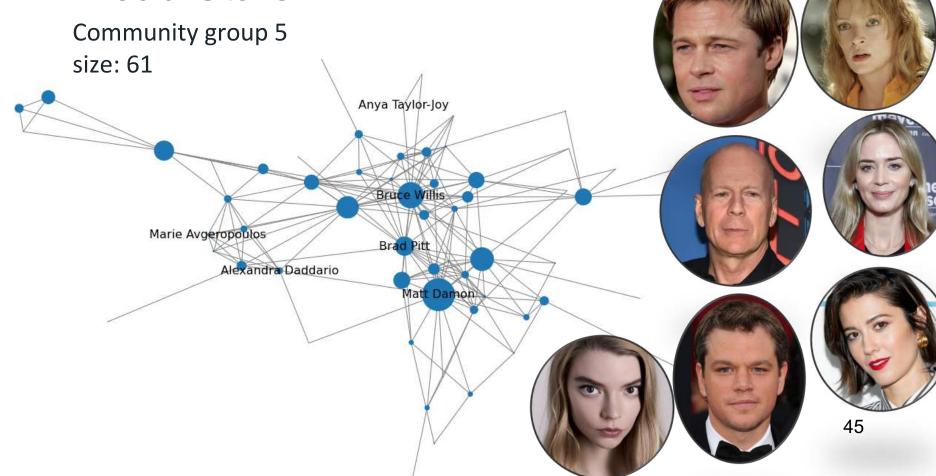
Community group 3

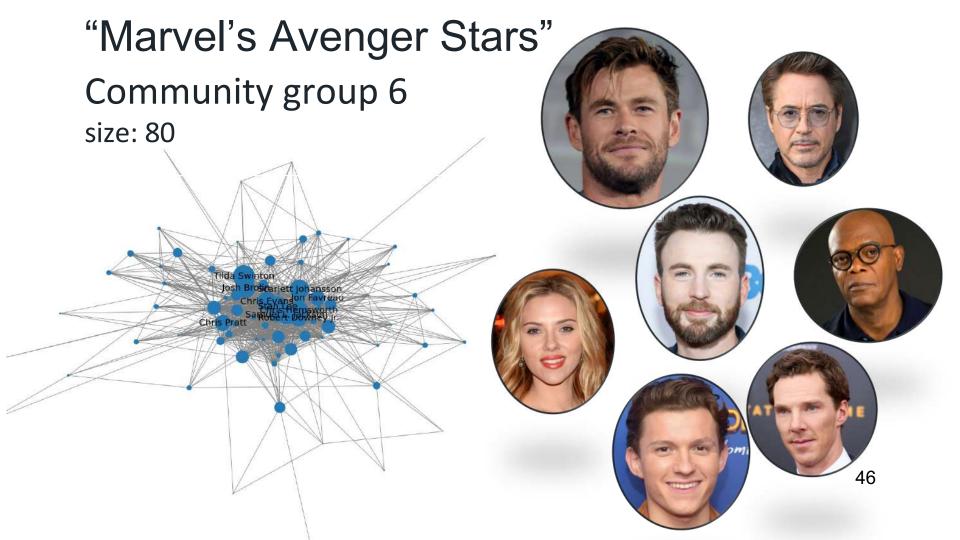






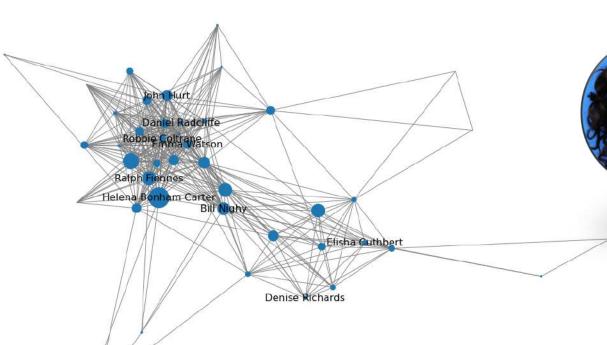
"Lead Stars 2"





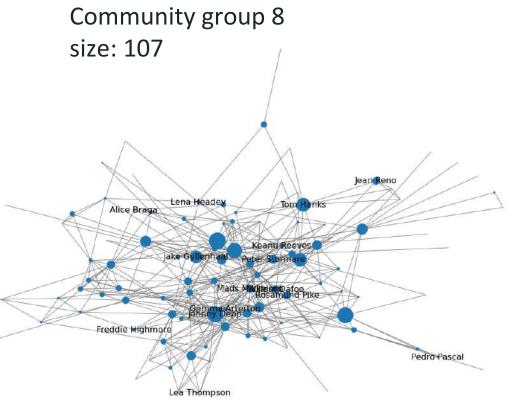
"English Stars"

Community group 7 size: 44





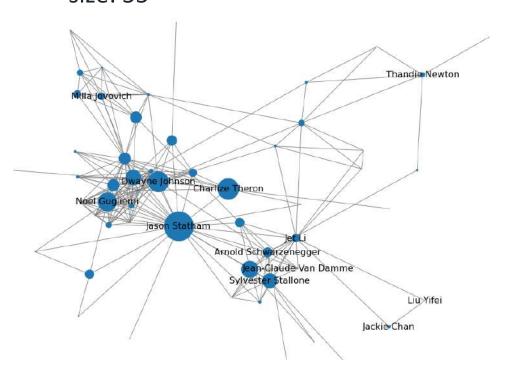
"Well-known Niche Stars"





"Action Stars"

Community group 9 size: 55



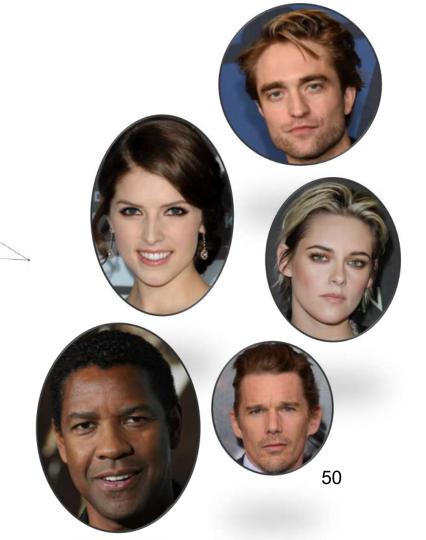


"Niche Stars"

Community group 10

Denzel Washington

size: 19



Summary Between Community

Mean

Community Group	Gr. Name	Amount Movie Played	Popularity	degree	btw_cent
0	All Stars	12	21	52	0.0062
1	Support Actors + Comedy	6	15	23	0.0031
2	Lead Actor/Actress 1	10	19	39	0.0048
3	Old Inactive Actors	7	19	34	0.0043
4	Fox Moive (X-men)	10	17	41	0.0061
5	Lead Actor/Actress 2	8	23	27	0.0022
6	Avenger	17	20	85	0.0095
7	English Actors	10	17	42	0.0018
8	Well-known Niche Actors	10	23	34	0.0029
9	Action Stars	8	24	27	0.0055
10	Niche Actors	8	11	23	0.0010

Summary Between Community

Mean

Community Group	Gr. Name	Amount Movie Played	Popularity	degree	btw_cent
0	All Stars	12	21	52	0.0062
1	Support Actors + Comedy	6	15	23	0.0031
2	Lead Actor/Actress 1	10	19	39	0.0048
3	Old Inactive Actors	7	19	34	0.0043
4	Fox Moive (X-men)	10	17	41	0.0061
5	Lead Actor/Actress 2	8	23	27	0.0022
6	Avenger	17	20	85	0.0095
7	English Actors	10	17	42	0.0018
8	Well-known Niche Actors	10	23	34	0.0029
9	Action Stars	8	24	27	0.0055
10	Niche Actors	8	11	23	0.0010

Conclusion

Centrality (Betweenness)

- Tell the versatility of actors and how likely they will accept role having to work with different groups of actors

Community

- Help understand the behavior of accepting roles in movie of each actors

Limitation

- Could not cluster in a group of genre based
- The model does not take into account of the presence (screen times) of actors within movie

THANKS!