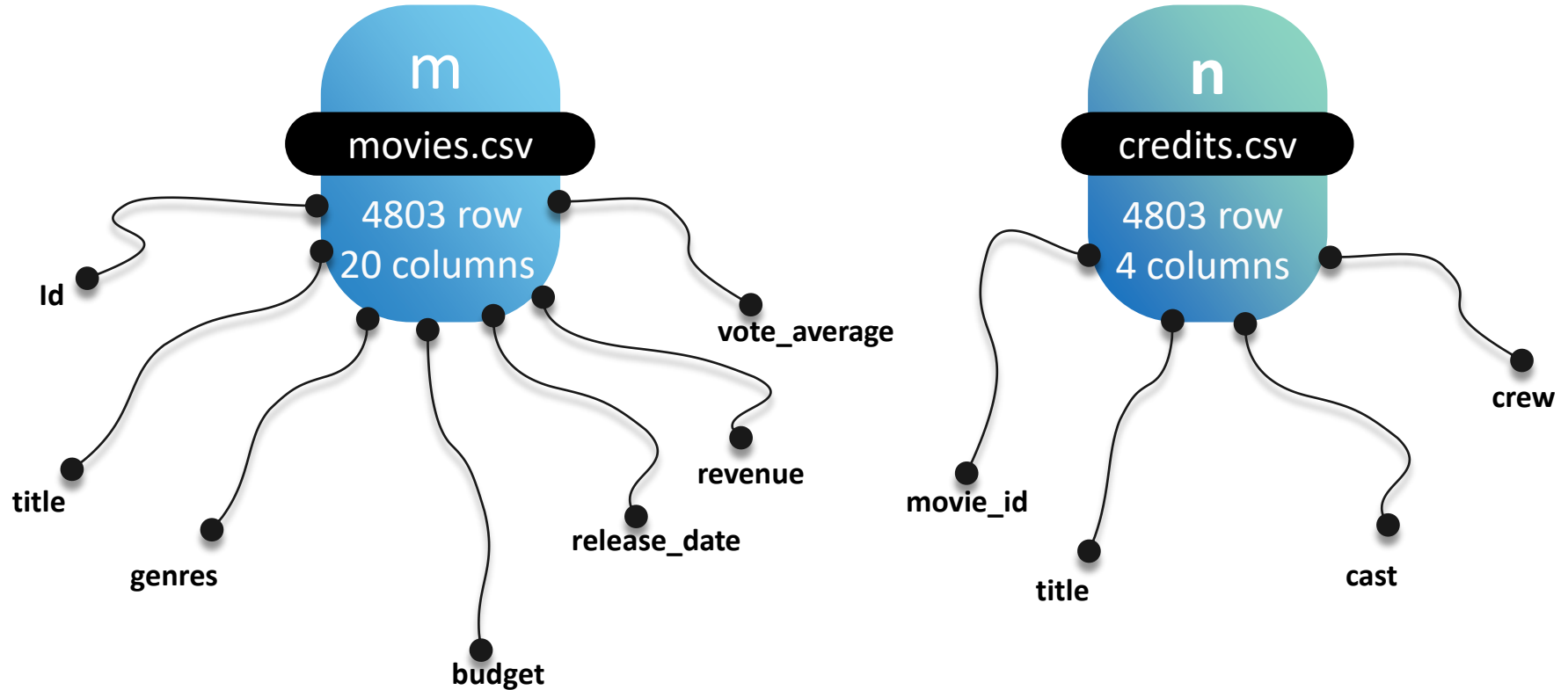# Evolution of the movie industry

## Team 04
Julien Berger - Jérémy Jayet - Hana Samet - Mathieu Shiva

# The initial dataset

## Two csv files

# Data cleaning

## Data cleaning
Remove redundancy in the crew set
➢ from 193655 to 87106 unique people

## Data expansion
From Json format to a clearer dataframe
➢ people dataframe: crew /cast

Remove movies with missing information
➢ From 4803 to 3229 movies

Removing people that worked in less then 5 movies

## Adding computed data
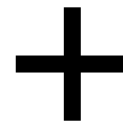➢ ROI column
➢ Success column

01

02

03

# Budget/revenue Analysis
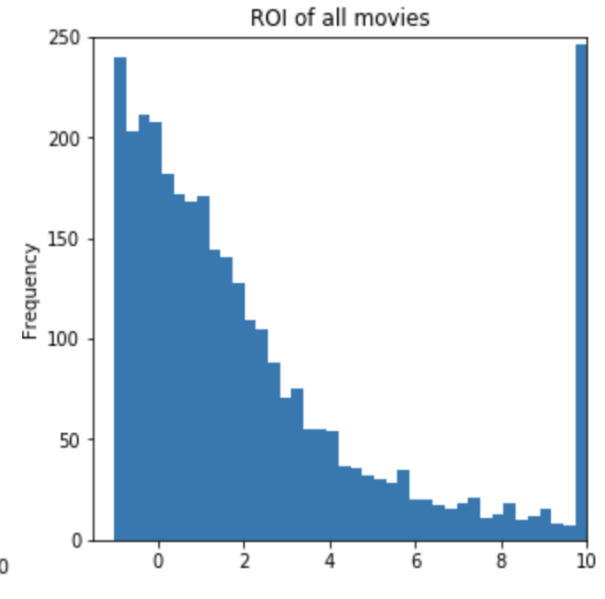
Tails-like distributions

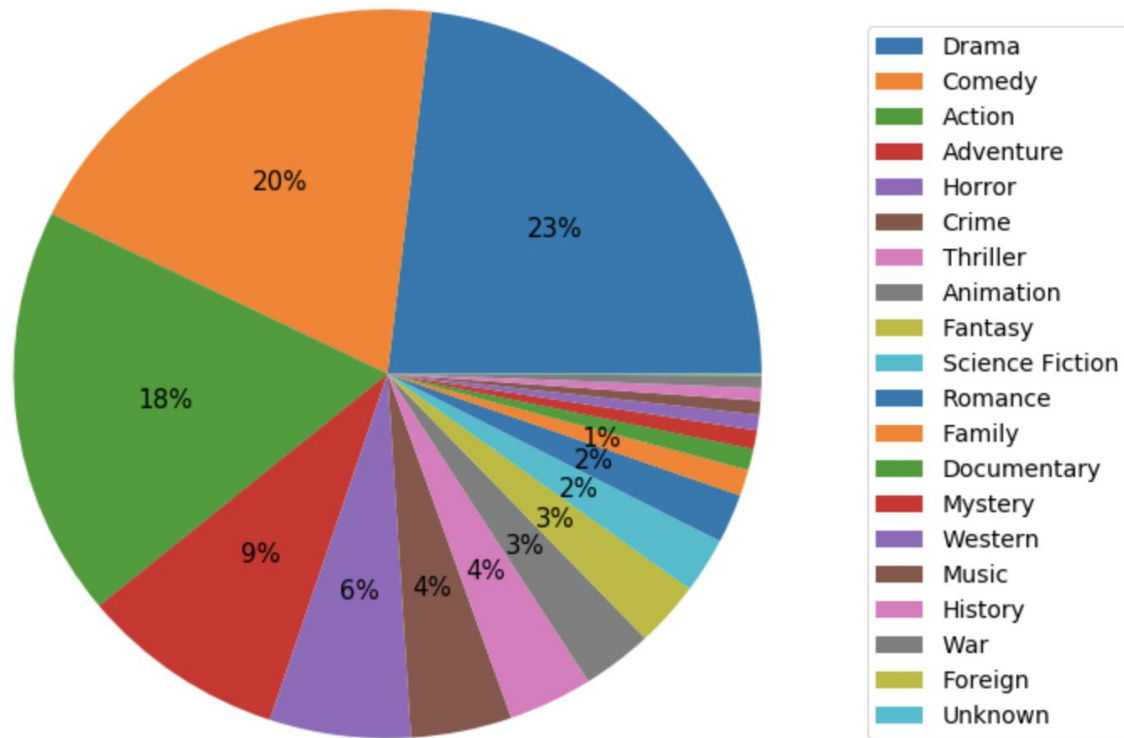A lot of movies had revenues that were smaller than their budget.

Only 23 movies had very high ROIs

# Genres



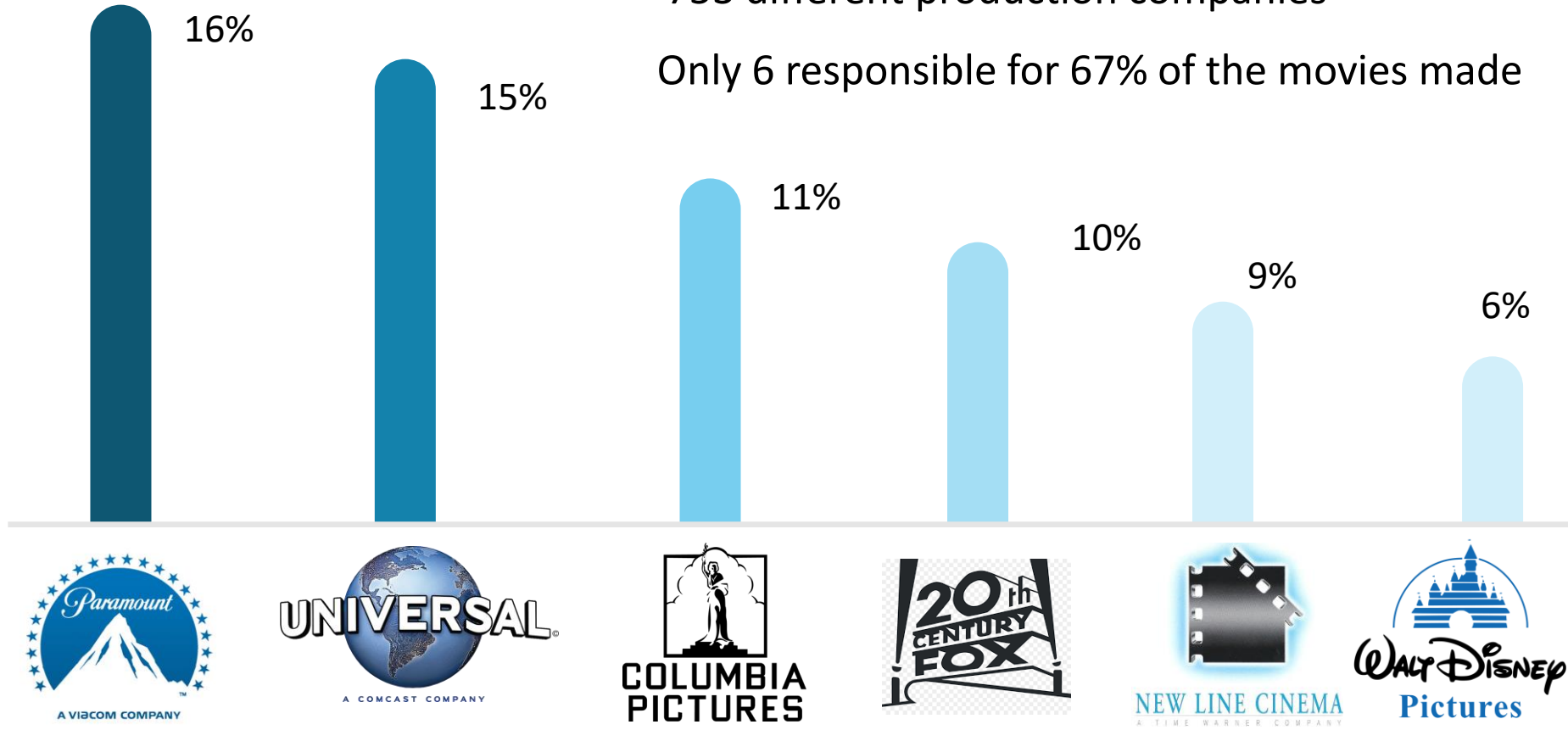Percentage of the genres in the whole dataset

Drama, Comedy and Action represent 61% of the movies

Legend:
- Drama — 23%
- Comedy — 20%
- Action — 18%
- Adventure — 9%
- Horror — 6%
- Crime — 4%
- Thriller — 4%
- Animation — 3%
- Fantasy — 3%
- Science Fiction — 2%
- Romance — 2%
- Family — 1%
- Documentary
- Mystery
- Western
- Music
- History
- War
- Foreign
- Unknown

# Production companies

753 different production companies

Only 6 responsible for 67% of the movies made

16%

15%

11%

10%

9%

6%

# Adjacency matrix - People based

## <k> = 4
### Actors only

# Adjacency matrix - People based

## <k> = 188
### Actors and crew

# Adjacency matrix - People based



Movies connected by at least 2 common actors, the colors represent the different movie genres
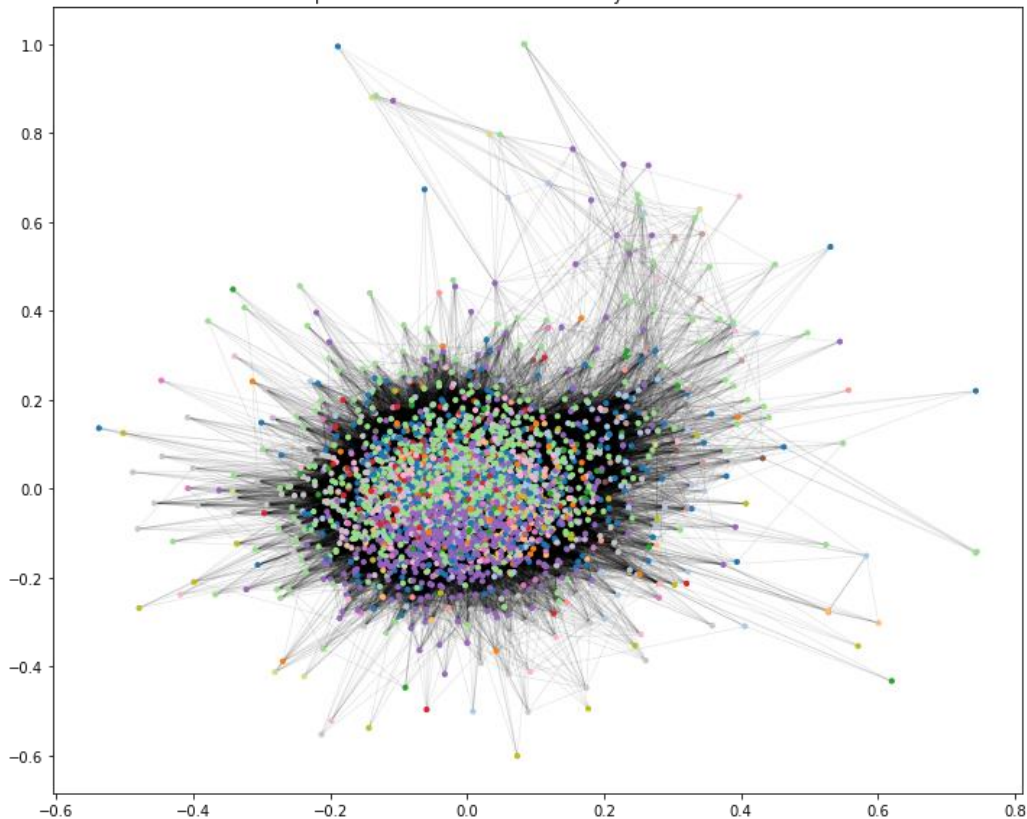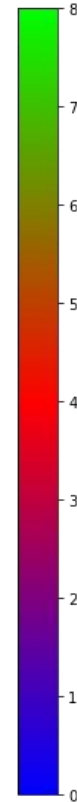
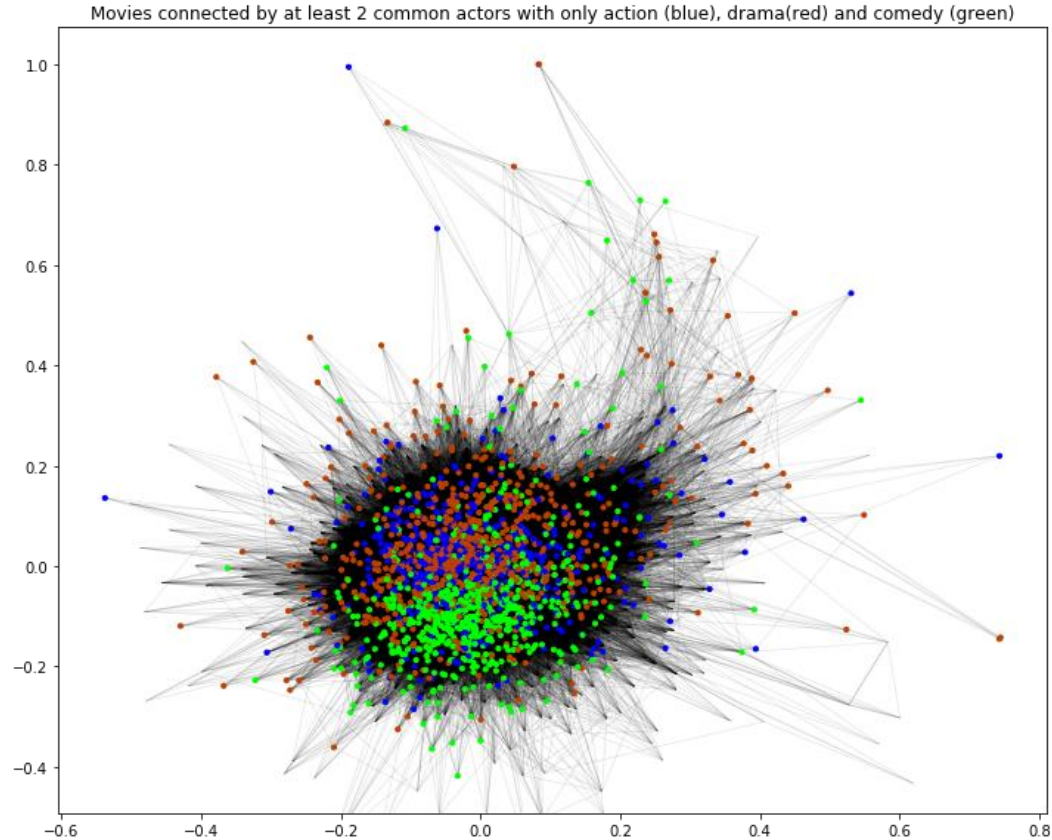2883 movies are connected

# Adjacency matrix - People based



Giant component of the movies connected by at least 2 common actors

2883 movies are connected

# Adjacency matrix - People based



Movies connected by at least 2 common actors with only action (blue), drama(red) and comedy (green)

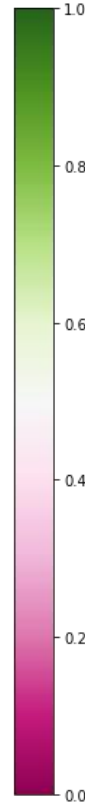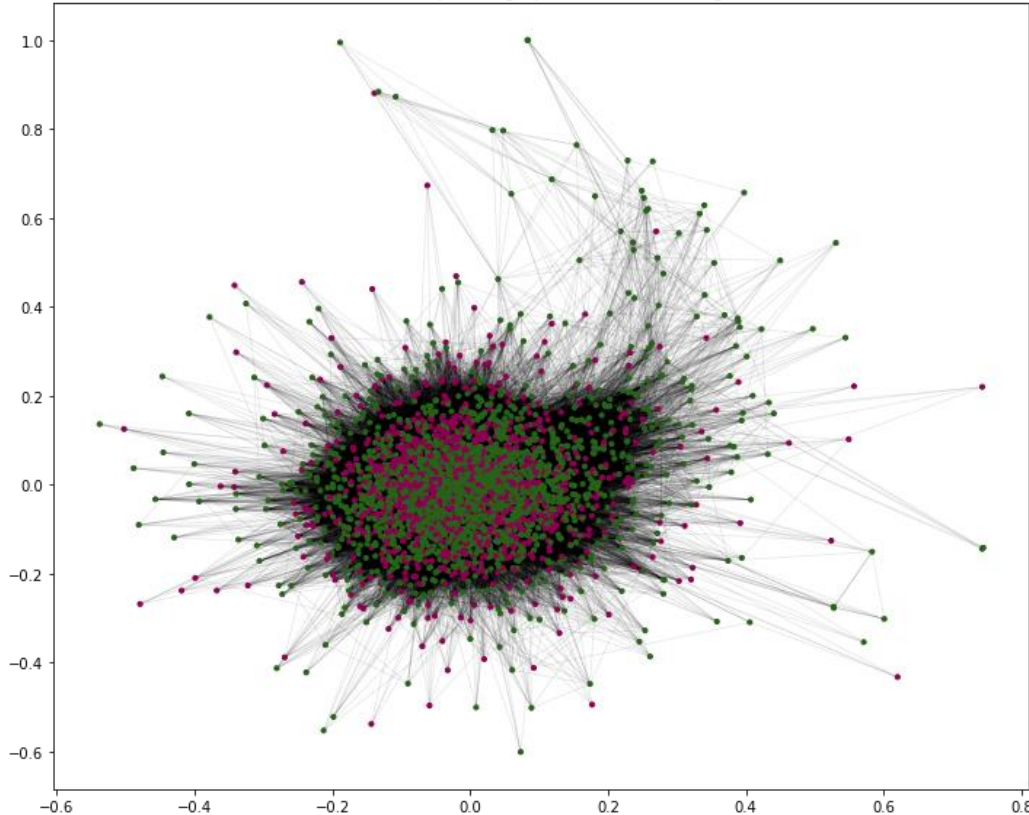2883 movies are connected

No communities can be seen

Separation of the graph in subgraphs is required

# Adjacency matrix - People based

## Success rate based on the ROI
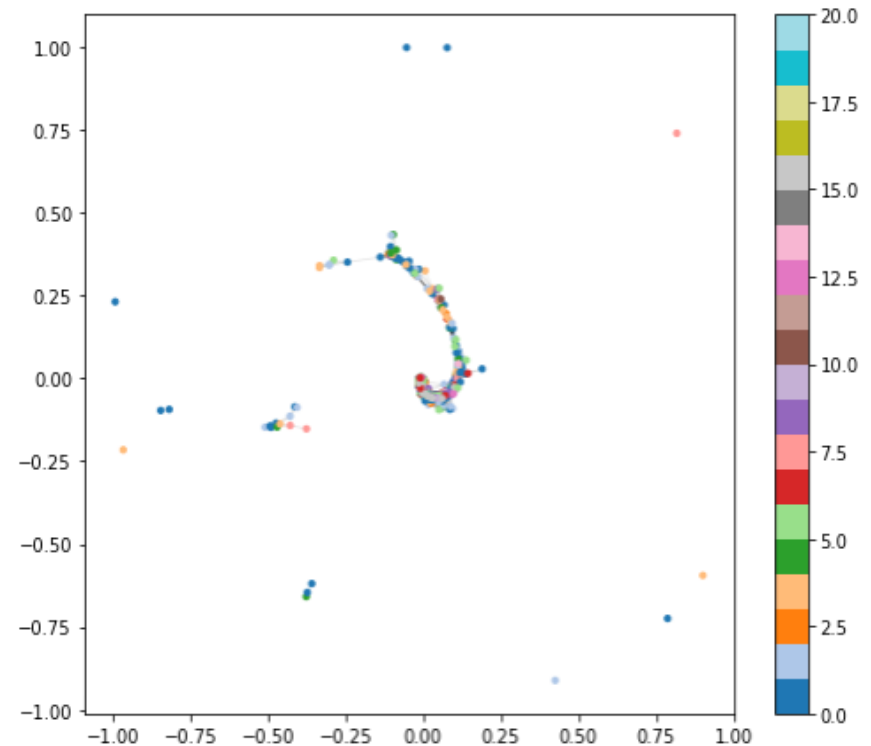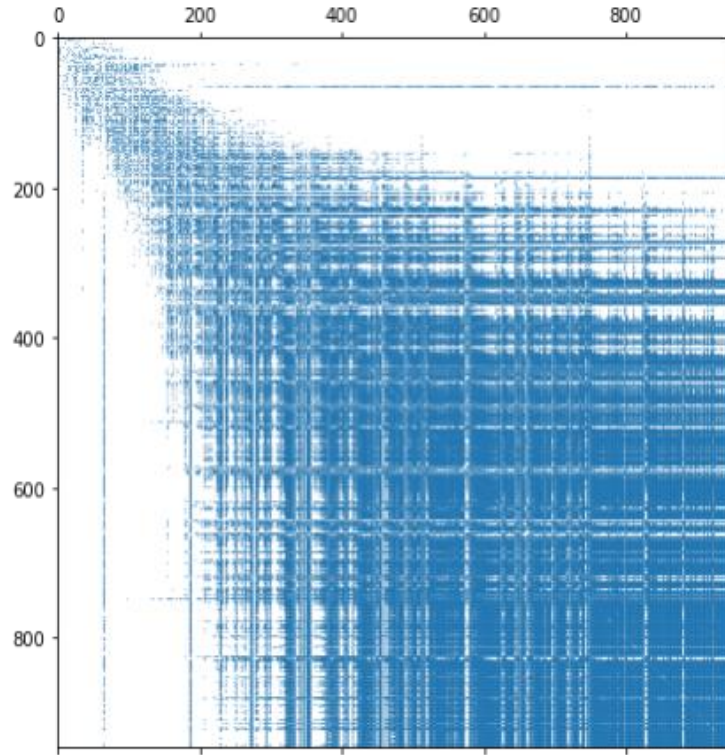


Giant component graph of the vote average

2883 movies are connected

No communities can be seen
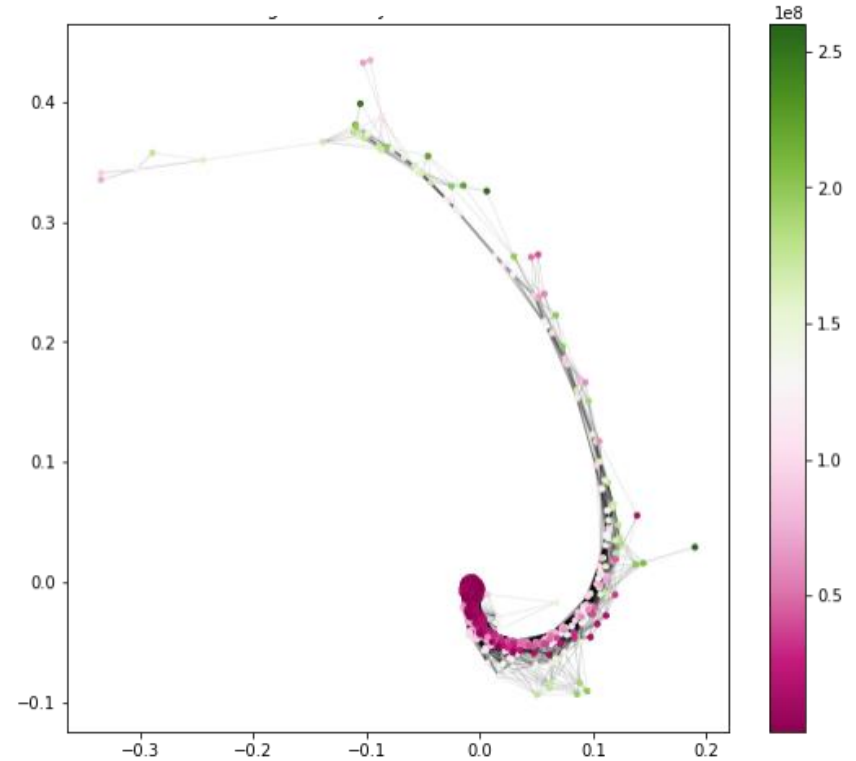
Separation of the graph in subgraphs is required
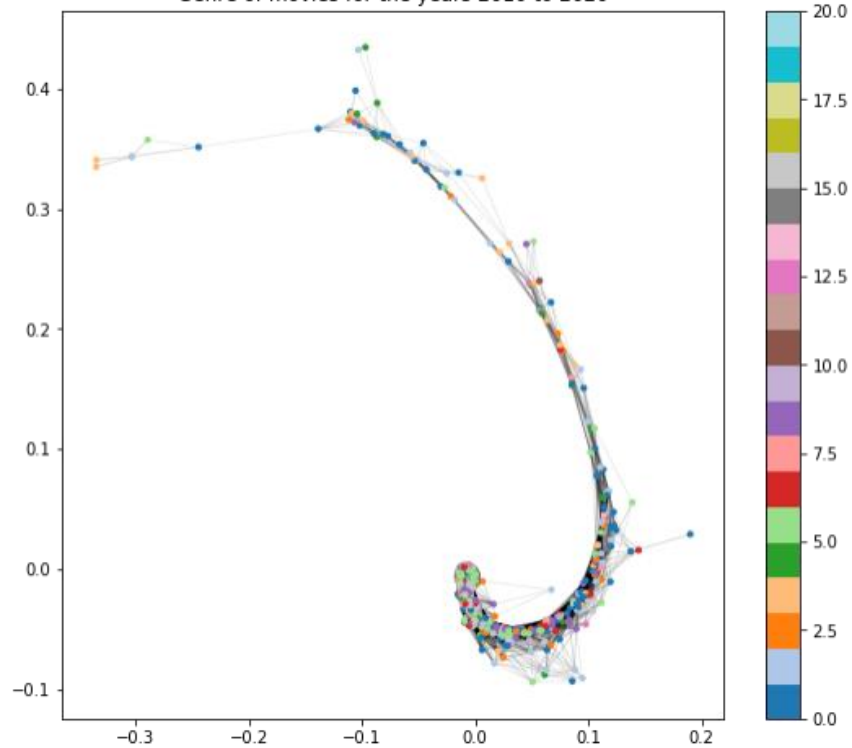
# Adjacency matrix - Budget based

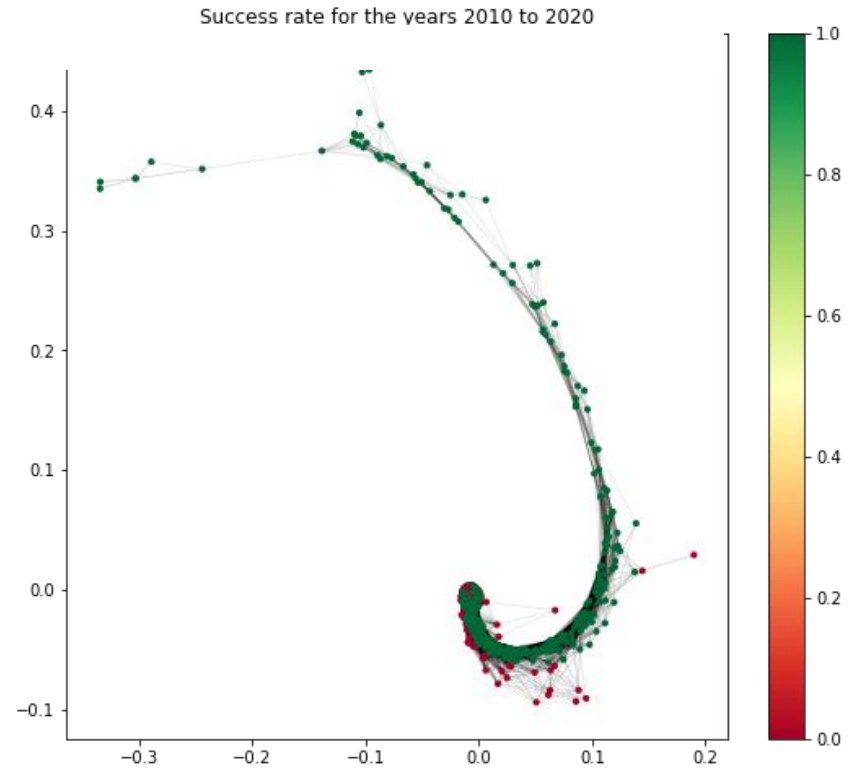## Euclidean distances between movies based on the budgets and revenues
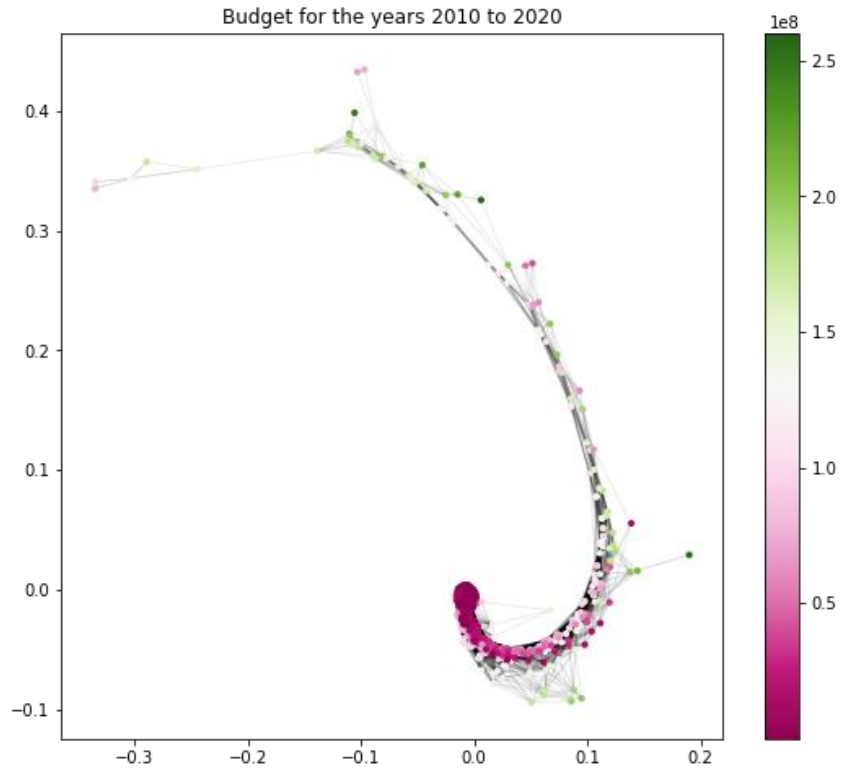
# Adjacency matrix - Budget based



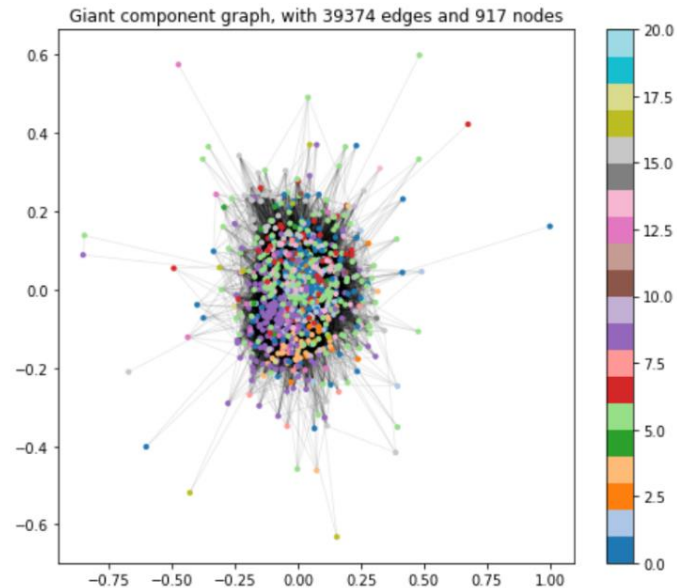Genre of movies for the years 2010 to 2020

# Adjacency matrix - Budget based

# Separation into decades

- 6 decades: 1960-2020
- Movies adjacency



Complete graph of common actors, with 30 connected components



Giant component graph, with 39374 edges and 917 nodes

# Correlation Heat Maps



The heatmap for the years 1980-1990

The heatmap for the years 2000-2010

# Correlation with the global economy



Evolution of the revenue, the budget and the GSPC

75% correlation (average revenue/GSPC) from 1967 to 2017

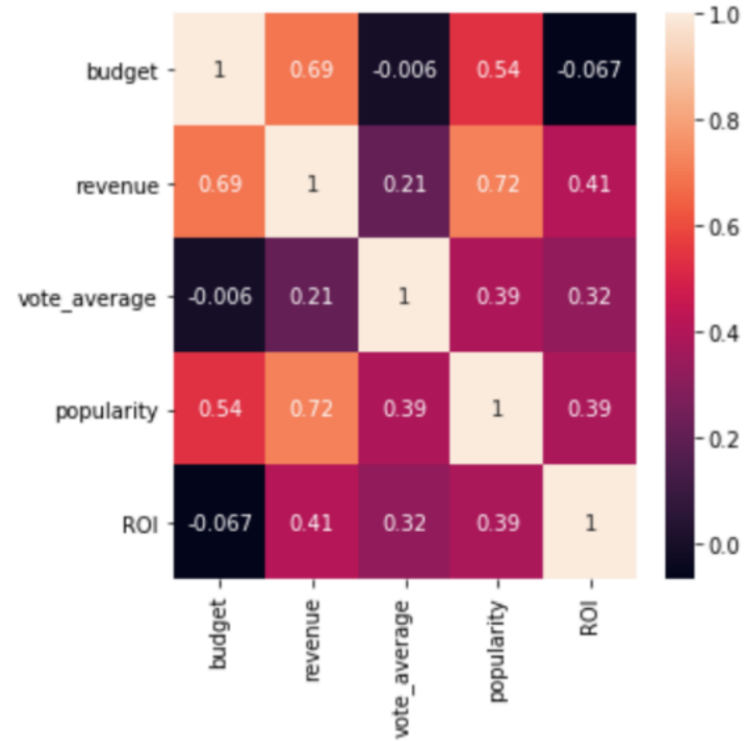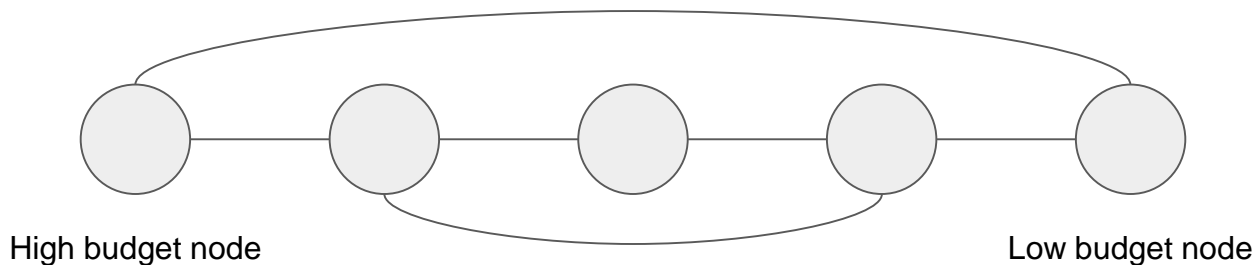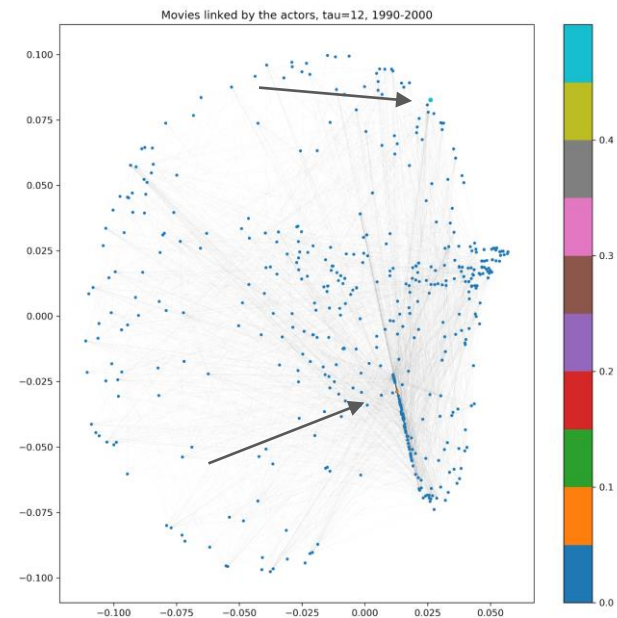# Budget-based graph vs actors-crew based graph

The actors-crew consists mostly in its giant components which is connected.

But how is it connected ?



High budget node                                    Low budget node

# Diffusion of a Dirac impulse on the graph



Movies linked by the actors, positions given by the budgets, 1990-2000

Movies linked by the actors, tau=12, 1990-2000

# Qualitative analysis



Movies linked by the actors, positions given by the budgets, 2010-2020



Movies linked by the actors, tau=12, 2010-2020

# Qualitative analysis

The signal spreads across the whole graph

Smoothness on the actors based graph does not mean smoothness on the budget based graph.

There are no communities based on the budget class of the films



Movies linked by the actors + 1 dirac, 2010-2020
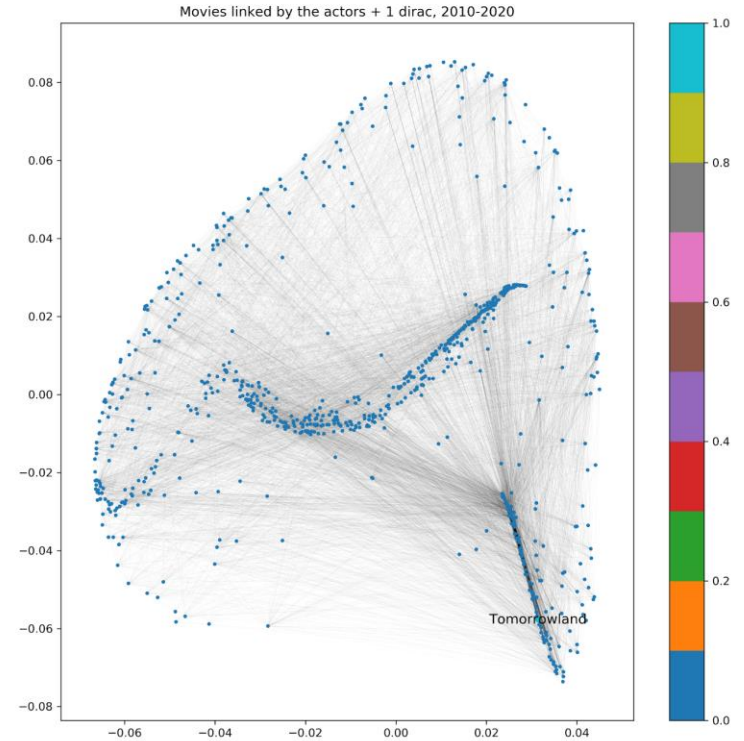
Tomorrowland

# Conclusion

- We did not find any communities in the actors-based graph.

- The budget-based graph does not show any community either and the budget appears as a very smooth signal over the graph.

- There is a high correlation between the popularity and the ROI.

- The budget based graph has a very different structure than the actors based graph. Therefore we can conclude that there is no segregation between actors playing in low budget movies and those playing in high budget movies.