

# r/marvelstudios: Community and Content Analysis after the release of Avengers: Endgame

Roberto Ferrari 852220, Davide Prati 845926, Marco Sallustio 906149

**Abstract**

This report analyzes interaction within the r/marvelstudios subreddit, focusing on the discussion of Avengers: Endgame, one of the most popular films in the Marvel Cinematic Universe (MCU). With the objective of identifying relationships, influential users, sentiments, and communities, the research employs Social Network Analysis, Sentiment Analysis and Topic Modeling. The report's structure comprehends data collection and exploration, social network analysis, sentiment analysis, and topic modeling. The approach is centered on understanding dynamics within an online community engaged in passionate discussions on a popular topic like Avengers: Endgame.

**Keywords**

Graph Theory - Community Detection - Sentiment Analysis - Topic Modeling - Reddit - Social Media Analytics

**Contents**

<b>Introduction</b>	<b>1</b>
Objective . . . . .	1
Structure of the report . . . . .	2
<b>1 Data Collection</b>	<b>2</b>
1.1 Data Preprocessing . . . . .	2
<b>2 Data Exploration</b>	<b>3</b>
<b>3 Social Network Analysis</b>	<b>3</b>
3.1 Graph . . . . .	3
3.2 Metrics . . . . .	4
3.3 Community Detection . . . . .	5
<b>4 Social Content Analysis</b>	<b>6</b>
4.1 Sentiment Analysis . . . . .	6
4.2 Topic Modeling . . . . .	8
Text Representation • The LDA model • Results	
<b>Conclusions and future developments</b>	<b>10</b>
<b>References</b>	<b>10</b>

In particular, we decided to analyze the subreddit r/marvelstudios, that, at the time of writing this report and conducting the project, with 3.5 million members, it ranks 168th among the world’s most popular communities [1]. This subreddit is dedicated to discussing Marvel Studios’ films and series and anything else related to the Marvel Cinematic Universe (MCU). We chose this subreddit because this is a rather popular topic among young people and therefore very popular in social networks like Reddit; furthermore, the selection of this subreddit was driven by the fact that Marvel movies often elicit diverse and controversial opinions among users, leading to intriguing discussions that would be particularly interesting to analyze.

For this reason we decided to analyze the interactions generated for one of the most popular Marvel films: Avengers Endgame.

This film in fact marks the end of Marvel’s epic Infinity Saga and achieved great success with critics and audiences, setting numerous box office records, also becoming the highest grossing film in the history of cinema.

**Objective**

Based on the considerations made previously we have identified some of the main objectives of our analysis:

- Identify the relationships generated within the subreddit and more specifically within our discussion of interest;
- Identify the most influential users through accurate Social Network Analysis;
- Understand the feedback that the users of the subreddit have had through the sentiment generated by each one;
- Visualize the main topics discussed in the community

**Introduction**

Reddit is a social media platform based on content sharing, discussion, and community voting.

It is organized into 'subreddits', which are thematic communities created by users. Each subreddit focuses on a specific topic, such as news, hobbies, or entertainment.

Users, commonly referred to as 'redditors', can subscribe to one or more subreddits based on their interests, allowing them to view and participate in discussions within those subreddits. They can engage by creating posts or commenting on existing posts.

and highlight groups of words that tend to relate to a specific topic rather than another;

- Identify any communities present and their characteristics.

## Structure of the report

The report is organized as follows:

1. **Data Collection:** in this stage, we accessed the Reddit API through PRAW to obtain the dataset we will be working on. This dataset will undergo a preprocessing phase, during which we will apply various steps to make it suitable for the purposes and models we intend to apply it.
2. **Data Exploration:** in this section, we present some graphs and visualizations that provide an overview of the construction and characteristics of the preprocessed dataframe.
3. **Social Content Analysis**, divided in
  - **Sentiment Analysis:** section in which we perform sentiment analysis on the subreddit, classifying comments as positive, negative, or neutral using specific scores to assign each comment to the correct category. The section also includes analysis and visualizations of the dataset, this time divided into the three types of sentiment, compared against each other.
  - **Topic Modeling:** after a preparation phase, we used the Latent Dirichlet Allocation (LDA) algorithm to compute topic modeling. To improve interpretation, we represented the results visually through various figures.
4. **Social Network Analysis:** We conducted a social network analysis, creating a graphical representation of the network, calculating and studying its associated metrics. Additionally, we computed the modularities of the communities, based on which we performed community detection on the six most significant ones.

## 1. Data Collection

We chose to use PRAW to extract our data from Reddit. PRAW (Python Reddit API Wrapper) [2] is a Python library that simplifies access to the Reddit API.

With PRAW, you have the capability to develop Python scripts and applications for interacting with Reddit data. This includes tasks such as reading posts, posting comments, retrieving user information, and a variety of other functionalities.

Furthermore, we decided to create our dataset by extracting the comments for each post in order to obtain as much data as possible and therefore to more easily extract discussions regarding

our topic of interest. To obtain useful comments for our analysis, we performed a search for posts in the “marvelstudios” subreddit that contain the keywords ‘avengers’ or ‘endgame’.

We thus obtained a dataset containing 112 567 columns. To avoid having results that were too sparse and not very relevant to our analysis, we filtered our dataset so that it only contained comments created between April 26, 2019 and April 26, 2020; therefore starting from the release date of the film in the USA for one year, obtaining a final dataset of 67 764 rows.

We decided to extract for each comment the following features:

- `author_submission`: author of the post
- `title_post`: title of the post
- `score_post`: score number of votes received by the post, indicates the “popularity” of the post
- `url_post`: URL of the post
- `created_utc_post`: date and hour of the creation of the post
- `author_comment`: author of the comment
- `body`: text of the comment
- `created_utc_comment`: date and hour of the creation of the comment
- `score_comment`: score number of votes received by the comment, indicates the “popularity” of the comment

### 1.1 Data Preprocessing

Before proceeding further, we applied a significant preprocessing phase to the data obtained and stored in our dataset, in order to make the extracted text as usable as possible for the algorithms and models we intend to apply.

For this purpose, a series of operations were carried out, mainly conducted through the use of regex rules.

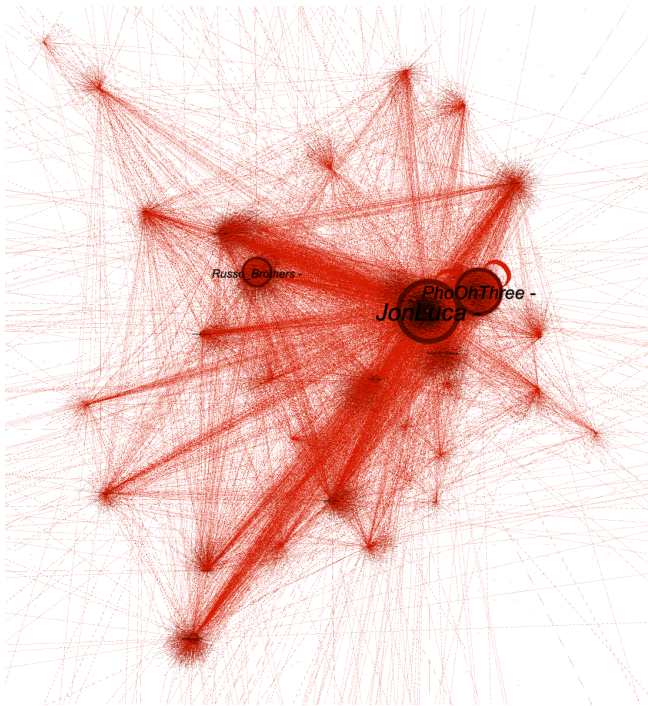
The operations are as follows:

- Lemmatization;
- Removal of links;
- Removal of e-mail addresses;
- Removal of numbers;
- Removal of extra-white spaces;
- Removal of punctuation;
- Removal of emoji;
- Transforming text to lowercase;



optimization process. The default value is 0.8: this means that, during the optimization process, every arc with a length greater than 80% of the length of the longest arc present is cut. This allows clusters to separate, as long edges with high weight may have too much influence on distant clusters.

The figure below shows the graph in which we decided to adjust the size of the nodes and labels (name of the account associated with the node) based on the weighted degree of each node, which is the sum of the weights of the incoming and outgoing edges. In practical terms, weighted degree is useful for understanding the overall amount of 'influence' or 'strength' of a node in the network, taking into account the strength of its connections.



**Figure 3.** Graph regulated by the weighted degree

### 3.2 Metrics

At this point, we proceeded by calculating the values of some standard metrics regarding the graphs, applying them specifically to our graph. This provides us with important information about the nature of our network. In the table below, we present the values obtained for these metrics.

Metric	Value
Density	0
Avg. Degree	1,29
Avg. Weighted Degree	88,25
Avg. Path Length	1,764

**Table 1.** Graph metrics

Analyzing these results, regarding first of all the density of the graph, it is equal to 0. This value indicates that there are very few connections in the network compared to the total number of possible connections; so the network can be clearly categorized as sparse.

Continuing the analysis, the average degree represents the average of the number of incoming and outgoing connections for each node in the network. In this case, it appears that on average each node has about 1.29 connections. This relatively low value could indicate a not very dense network, and this result is consistent with the value we found for the density.

The weighted average degree takes into account the weights of the edges in the network. This value suggests that, on average, each node has a much higher weighted sum of connections (88.25) than the unweighted degree (1.29). This could indicate that some connections are much stronger than others in the network.

Finally, the average path length is defined as the average of the shortest distances between all pairs of nodes in the network. In this case, the average path length is 1,764 edges, which indicates that the paths between nodes in the network are, on average, relatively short.

In order to understand who the main influencers of the network are, we computed two centrality metrics: Closeness Centrality and Betweenness Centrality.

These measures allow us to find the nodes that take on an important structural role within the network: the first provides information on which actors are the most central, therefore with the shortest distance from the rest of the nodes; the second instead shows which users, if any, are useful for connecting different parts of the graph and potentially putting different communities in contact.

We have reported in Table 2 the 10 nodes with the highest value of Betweenness Centrality:

User	Betweenness Centrality
<i>The_Asian_Hamster</i>	14530.1
<i>KostisPat257</i>	13641.5
<i>Flamma_Man</i>	9227.5
<i>Sisiwakanamaru</i>	4520.5
<i>PhoOhThree</i>	3600.6
<i>kahlkorver</i>	1983
<i>MangoJam18</i>	1294
<i>LordHyperBreath</i>	670.9
<i>dannys717</i>	570
<i>ENusatron</i>	392.08

**Table 2.** Betweenness centrality



As regards the closeness centrality values, we obtained many nodes with a value of 1; more precisely we have that 64% of the nodes have a value equal to 1 which represents the maximum value that a node takes on in our network.

It would not make sense to report a table with the 10 nodes with the highest value because they would all have a value equal to 1 and for this reason we decided to observe and report in the Table 3 the Closeness Centrality values for the same nodes present in Table 2 in order to make a comparison between the two values:

User	Closeness Centrality
<i>The_Asian_Hamster</i>	0.58
<i>KostisPat257</i>	0.85
<i>Flamma_Man</i>	1
<i>Sisiwakanamaru</i>	0.42
<i>PhoOhThree</i>	1
<i>kahlkorver</i>	0.5
<i>MangoJam18</i>	0.61
<i>LordHyperBreath</i>	1
<i>dannys717</i>	1
<i>ENusatron</i>	0.66

**Table 3.** *Closeness centrality*

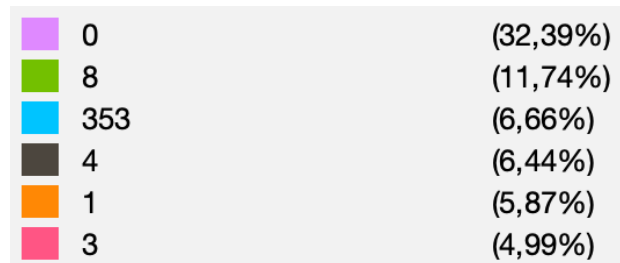
Therefore, comparing the values within the two tables we can make several considerations:

- Users like *The\_Asian\_Hamster*, *KostisPat257*, and *Flamma\_Man* are central in terms of both betweenness and closeness. These users play key roles in network connectivity and are quickly accessible to others.
- Some users show significant differences between betweenness and closeness values. For example, *Sisiwakanamaru* has a relatively low betweenness value but a lower closeness value, suggesting that it may be less involved in key pathways but is well connected with neighbors.
- Users like *ENusatron* and *MangoJam18* could be considered more peripheral in the network, as they have relatively low values of both betweenness and closeness.
- *Sisiwakanamaru* and *MangoJam18*, with relatively low closeness, may be more isolated than other users.
- Users with high betweenness, such as *The\_Asian\_Hamster* and *KostisPat257*, could play mediation and bridging roles, while those with high closeness, such as *Flamma\_Man* and *PhoOhThree*, are quickly accessible and could have a more direct impact on information dissemination.

### 3.3 Community Detection

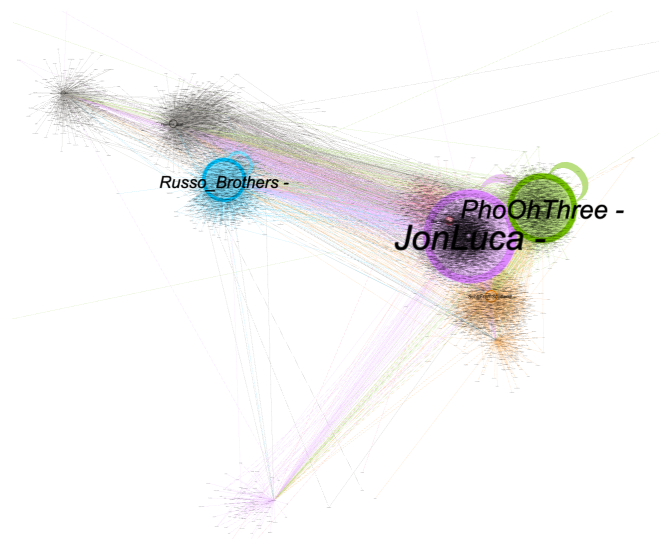
With the aim of analyzing the communities of interest, the modularity value of our network was calculated. Modularity is an important measure of network or graph theory. It was designed to measure the strength of the division of a network into modules (also called groups, clusters or communities). Networks with high modularity have dense connections between nodes within clusters but sparse connections between nodes in different clusters. Gephi to calculate the modularity value uses the Louvain algorithm [5], one of the fastest algorithms available, considered cutting-edge for its high performance. We found 485 communities with a modularity value of 0.744. This high value of modularity tells us that within the network there are dense connections between nodes belonging to the same cluster, but very sparse connections between nodes that are part of different clusters.

We show below the Community Detection graph in which only the 6 most significant community classes have been selected and displayed, i.e. those that contain the greatest number of nodes (the most populous). The 6 communities are thus shown below with the community number and associated color:



**Figure 4.** *The 6 most significant communities*

Also in this case the size of the nodes was adjusted based on the weighted degree and furthermore each node is colored differently based on the color associated with the community to which it belongs:



**Figure 5.** *Community detection*

Thus, we observe that the communities with a higher population correspond to nodes with larger dimensions, as indicated by the weighted degree value. This suggests that there may be a relationship between node size (based on weighted degree) and its membership in a community. One reason is certainly that more central nodes with a higher degree might have a greater probability of belonging to larger communities, node centrality can be reflected in the weighted degree, and central nodes are often important to the overall structure of the graph. Another rationale is that communities could grow around central nodes with a higher weighted degree, as these nodes attract more connections over time and indeed some nodes with a higher weighted degree could play a key role in maintaining the integrity and cohesion of communities, leading to the formation of larger communities around them.

Finally, to better analyze the communities found, we report the 6 communities and the associated sentiment values in the figure below to understand whether sentiment influenced the division of the communities:

Community	Populousness	Positive	Negative	Neutral
0	8131	49,61%	33,20%	17,10%
8	2948	46,80%	32,50%	20,65%
353	1673	61,92%	19,78%	18,29%
4	1817	48,67%	32,40%	19,35%
1	1474	47,15%	33,24%	19,60%
3	1253	47,72%	33,28%	18,99%

**Figure 6.** *Communities' sentiment*

As you can see, we have very similar sentiment values between these 6 communities; this certainly indicates that people within each community interact similarly, with a balance between positive, negative and neutral comments. This could be a sign of cohesion and coherence within communities. The only different values are for community 353 which has a higher positive sentiment value (almost 62%) and more unbalanced than the other communities. In fact, observing the graph in Figure 3 we know that this community includes the user 'Russo\_Brothers' whose node is in fact positioned in a more distant position than the others.

This may be because this community may cover topics or content that is significantly different from other communities. This diversity in content could influence sentiment more markedly, leading to a more distant arrangement in the layout we chose to use.

## 4. Social Content Analysis

### 4.1 Sentiment Analysis

Sentiment Analysis was performed to detect the polarity (positive, negative, neutral) of comments' texts.

The model used is VADER (Valence Aware Dictionary for Sentiment Reasoning) [6], that is sensitive to both polarity (positive/negative) and intensity (strength) of emotion. More

precisely, it relies on a dictionary that maps lexical features to emotion intensities, known as sentiment scores. The sentiment score of a text can be obtained by summing up the intensity of each word in the text, which ranges between -1 and 1. A score close to -1 indicates a very negative sentiment, a score close to 1 indicates a very positive sentiment, and a score close to 0 indicates a neutral sentiment.

Obviously this sentiment analysis phase was carried out after the pre-processing phase explained in the section 1.1.

To have even more detailed values, still using the VADER library, we have created three different functions for each polarity (positive, negative, neutral) which return the sentiment score rather than a compound score.

So, for example, for the first row of the dataset we obtained this scores which we interpret in the following way:

- **CompoundValue:** -0.9745

The compound score is very negative, indicating that the first sentence has a strong negative sentiment. Since it is close to -1, this is a very bad rating.

- **NegativeValue:** 0.328

The negative sentiment score is 32.8%, indicating that a significant portion of the sentence has a negative tone.

- **PositiveValue:** 0.065

The positive sentiment score is 6.5%, indicating that only a small part of the sentence has a positive tone.

- **NeutralValue:** 0.607

The neutral sentiment score is 60.7%, indicating that the majority of the sentence has a neutral tone.

After doing this, we decided to create a new column in our dataset that directly contained the sentiment value associated with the comment (Positive, Negative or Neutral). To do this we used the CompoundValue obtained previously and we established a specific threshold in order to convert the numerical value into the associated sentiment. We then analyzed the distribution of the sentiment generated for our dataset and obtained that:

- The percentage of negative comments is 31.62%
- The percentage of positive comments is 49.61%
- The percentage of neutral comments is 18.77%

This result is also displayed graphically in the following histogram:







observed in that case, which are 0.355 and for  $C_v$  and -4.106 for  $U_{mass}$ .

Of course, it is better to underline that this choice was solely based on our judgment and a balance we struck between the goodness of the two parameters. Other choices would have been possible with different approaches.

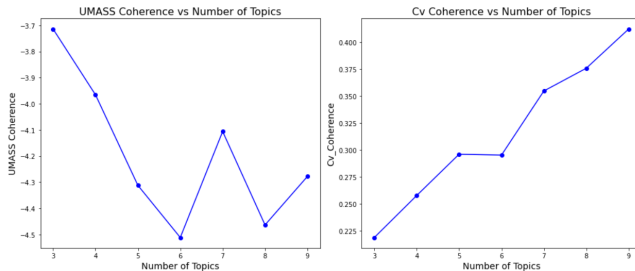


Figure 12.  $C_v$  and  $U_{mass}$  values for the LDA models

Thus, at this point, we proceeded with the model configured with the number of topics set to 7.

#### 4.2.3 Results

To visualize the results, first of all we generated an interactive graph called an Intertopic Distance Map, which illustrates the sizes and relationships between various clusters/topics. Below is an image showcasing the initial visualization, but interactively, by navigating through the notebook, one can explore the top 30 most prevalent words in each topic relative to their overall frequency in the text.



Figure 13. Intertopic distance map

But a more significant visualization for our purpose is undoubtedly that of the WordClouds, that we generated based on

the word weights for each topic. The weights were computed using the  $\gamma$  parameter, set to 0.6, which is used to balance two factors: the probability of a term within a topic against its overall probability in the corpus and the overall probability of the term in the corpus in the pyLDAvis library [9].

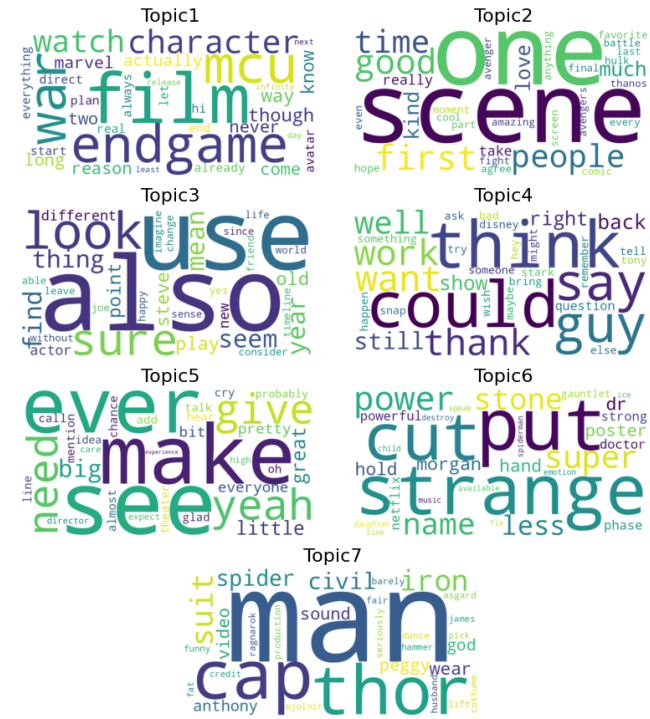


Figure 14. The WordClouds obtained for the 7 topics

At this point, having the representations of the most frequent words in each topic, we can try to understand what distinguishes one from another and give them a brief explanatory title, always based on our judgment. We proceeded in the following manner:

1. **Topic 1:** general opinions on the movie, expectations,
2. **Topic 2:** focus on particular memorable scenes, and favorite moments,
3. **Topic 3:** actors and characters' development, differences between this film and the previous ones,
4. **Topic 4:** critical reflections on the plot, actions, and thoughts of the characters,
5. **Topic 5:** astonished and sensationalistic comments on impactful moments and unexpected plot twists,
6. **Topic 6:** specific characters and their superpowers, with reference to a cut scene, likely featuring Dr. Strange as its protagonist,
7. **Topic 7:** focus on the most iconic characters: Thor, Captain America, Iron Man.

## Conclusions and future developments

At the conclusion of this project, we can say that we managed to achieve the various objectives we had set ourselves.

In particular, thanks to social network analysis techniques we were able to identify the main players in the discussion by first analyzing the betweenness centrality and then the closeness centrality; we thus obtained very high values for the two metrics which showed us how many nodes within our network are crucial nodes for the discussion.

Furthermore, we also met one of our goals regarding the need to find how our network is divided into communities; we obtained that in the largest communities found, all the nodes with a higher weighted degree value corresponded. This phenomenon leads us to support the thesis that central nodes with a high degree of weighting could have a greater probability of belonging to larger communities and that communities could grow around central nodes with a higher degree of weighting, since these nodes attract more connections over time and indeed some nodes with a higher degree of weighting could play a key role in maintaining the integrity and cohesion of communities, leading to the formation of larger communities around them.

We can say that we also managed to satisfy the objective in which we asked ourselves to understand user feedback regarding the analyzed topic of interest; to do this we used sentiment analysis techniques which we also exploited to analyze the distribution of sentiment in the different communities found. In this regard, in the future an emotion analysis could also be performed to obtain even more detailed information on user feedback.

As for topic modeling, we observe from the word clouds and our interpretation that the topics primarily revolve around opinions on the movie's quality and plot, or considerations about the story and characters. Therefore, despite the film being among the highest-grossing in cinema history [10], none of the seven detected topics seems to reference the box office success or the film's production.

Topic modeling could be further developed, considering that, as mentioned earlier, the choice to proceed with 7 topics is justified by a good balance between the values of the two parameters, but it is not the only viable option. Therefore, it might be worthwhile to run the model with a different number of topics, for instance 4 or 5, and assess the results.

Furthermore, one could explore the use of alternative models to LDA, such as LSA (Latent Semantic Analysis) or BERT (Bidirectional Encoder Representations from Transformers). In these cases, text representation strategies other than BOW, such as TF-IDF, could be considered. Implementing these models with a range of possible topics from three to nine would allow for a comparison of  $U_{mass}$  and  $C_v$  coherence values with those obtained with our LDA model.

## References

- [1] Reddit's best communities:  
<https://www.reddit.com/best/communities/1/>
- [2] GitHub - PRAW:  
<https://github.com/praw-dev/praw>
- [3] pyclld2 library:  
<https://pypi.org/project/pyclld2/>
- [4] Gephi:  
<https://gephi.org/features/>
- [5] Louvain algorithm:  
*Towards data science - Louvain Algorithm. Luis Rita, 09/04/2020*
- [6] VADER:  
<https://pypi.org/project/vaderSentiment/>
- [7] Gensim:  
<https://pypi.org/project/gensim/>
- [8] Ismail Harrando, Pasquale Lisenà and Raphael Troncy. Apples to Apples: A Systematic Evaluation of Topic Models:  
<https://aclanthology.org/2021.ranlp-1.55.pdf>
- [9] pyLDAvis:  
<https://pypi.org/project/pyLDAvis/>
- [10] movieplayer.it, film con maggiori incassi nella storia del cinema:  
<https://movieplayer.it/film/boxoffice/internazionale/di-sempre/>