# The datafication of movies

## Investigation of Cli-Fi as a new emerging genre on the platform IMDb

*Authors*

*Lina Alkowatly, Yvonne Mrukwa, Gustaf Rossi*

*and Rodrigo de Ros*

# Innehåll

# 1. Introduction

Over the past decade, a strong consensus has emerged over climate change posing risks that could end civilization as we know it (Intergovernmental Panel on Climate Change, 2014). Climate change as a topic has emerged as a dominant theme in literature and, correspondingly in fiction, creating a new generation of narrative called climate change fiction, abbreviated Cli-Fi.

The Cli-Fi genre is defined as a genre that deals with climate disasters in a fictional way (Svoboda, 2016). Hence, the genre is derived from the same word of Sci-Fi.

The newly emerged genre has been considered particularly interesting because it can be seen as a cultural response to most scientific and policy discourses that offers a way of exploring dramatic social change through the perspectives of individual and social group experiences by way of fictional narrative (Whiteley, Chiang, Einsiedel, 2016, p.28).

Some movies can as well have a powerful effect on people's consciousness and how people form beliefs. Movies as a medium can be perceived as a way to discuss current events or portray a future where the world has been destroyed. Movies tend as well to instruct people in everyday life and can influence political issues (Igartua & Barrios, 2012; Denvir, 2003). Movies in the Cli-Fi genre could then form people influences and engagement about climate issues.

With this research project, we set out to investigate the new emerging genre of Cli-fi by understanding how movie information are datafied and as well how this genre has formed from the datafication of movies. To examine this issue we used the movie database IDMb to collect information and data about a selection of movies that cover the issue of climate change or climate disasters.

With this paper we will discuss our reflections on studying the newly emerged genre and the data that makes up these movies. The first part of this paper will discuss the current state of data and datafication with the movies and how the movie information are datafied. Then we will discuss the data we used for our project and the history of IMDb as a movie database along with a history of how movies became datafied.

Finally, this paper contains a discussion and analysis of our data that we collected about the movies and places this in relation to today's society.

# 2. Current state of Data and Datafication

Datafication as a term is used to describe the practice of taking aspects of the world that have never previously been quantified and rendering them into data. Van Dijck (2014) defined datafication as a legitimate means to access, understand and monitor people's behaviour. She also states that datafication became a leading principle amongst techno- adepts and amongst also scholars who see datafication as a revolutionary research opportunity to investigate human conduct (p. 198). In the last few years, we have been able to find data about movies from different sources. Movie data is highly popular on the internet and there are several web resources dedicated to movies and many others containing movie-related information, for instance, IMDb, Rotten Tomatoes, AllMovies, to name a few.

One of them is the Internet Movie Database (IMDb) in which we have aimed to collect the sample of movies. We have agreed that the website offers a database with an organized collection that includes all the information needed to conduct this study as the website has been praised in a good number of articles (Weible 2001; Peralta, 2007; Naun & Elhard, 2005; Schneider, 2001). Cherié L. Weible (2001) has argued that IMDb is qualified to be entered in the well-known useful Internet resources available at no cost with a click of the mouse. He has further stated that one of the best descriptions of the Internet Movie Database says that it has "tons of info" and is "cross-indexed and linked up the wazoo" (p.48). In the American Libraries article, "Love Is a Many-Splendored Gizmo," IMDb is included in a list of the ten internet resources that everyone would want with them if stranded on a desert island (Schneider, 2001, p.84).

IMDb describes itself as the world's most popular and authoritative source for movie, TV and celebrity content. It is designed to help fans explore the world of movies and shows and decide what to watch (IMDb, 2019a). Moreover, on the website (http://www.imdb.com/), IMDb affirms that it provides a searchable database with a huge collection of movie information; ''We help you jog your memory about a movie, show, or person on the tip of your tongue, find the best movie or show to watch next, and empower you to share your entertainment knowledge and opinions with the world's largest community of fans.'' (IMDb, 2019a).

The website started as a hobby project by an international group of movie fans and currently belongs to the Amazon.com Company. It can be seen that IMDb tries to catalogue every pertinent detail about a movie, from who was in it, to who made it, to trivia about it, to filming locations, and even where you can find reviews and fan sites on the web. As of May 2019,

IMDb declares that the total number of data items is 336,394,426 which has approximately 6 million titles (including episodes) and 9.9 million personalities in its database, as well as 83 million registered users (IMDb, 2019b). Moreover, the website provides 49 text files in ad-hoc format (called lists) containing different characteristics about movies (e.g. actors.list or running-times.list). These 49 lists catalogue different details about the movies. Each list has a list manager, who is responsible for updating, maintaining and publishing it. List managers mostly rely on users of the IMDb to submit the content (ibid.).

Furthermore, on the website help centre it is stated that all the data comes from various sources. IMDb gathers the information and verifies the items with studios and filmmakers. The data goes through consistency checks to ensure it is as accurate and reliable as possible. At the same time, IMDb does not neglect that mistakes are inevitable due to the fact of the big volume of data and the nature of the information they are listing. However, all the data are verified and fixed when they are spotted or reported (IMDb, 2019c).

Moreover, another source of information about movies or TV shows; "released title", in IMDb is submitted by the users. From this regard, IMDb confirms also using a specific process to ensure that this data goes through consistency checks to ensure that the data is as accurate and reliable as possible, and the process goes as follows:

- Any registered user can add data to the site.

- Users add new data that we don't already have; correct/update existing data and delete inaccurate information.

- Once data is submitted, it goes through to Data Managers for processing.

- There are no maximum or minimum limits for User's additions and corrections. IMDb accepts large and small submissions (IMDb, 2019d).

The registered users on IMDb have the possibility to add their data in dozens of data sections available on the website. These sections include a variety of information about the title, to name a few: biographical information, the box office, plot information, bio trivia, articles in the media, company credits etc. Besides, the website presents a method for the users for the contribution their data; starting from the free registration on the site, to finding the preferred title to update until submitting the data for processing by IMDb data managers. Although IMDb facilitates the procedure for data contributions for its users and offers them plenty of options, there is still certain types of data, according to the website, that are not possible for users to contribute to or update (e.g. new titles, images, user reviews, parental guides) (ibid.).

Moreover, another datafication process found in IMDb is related to quantifying the general public opinion about any title in which IMDb refers to 'IMDb rating'. The website offers the possibility for its users to cast a vote on every released title in the database. IMDb incorporates the concept of the wisdom of the crowd (i.e. the collective opinion of a large group of individuals rather than of a single expert). Accordingly, the rating of a title is done by assigning a positive integer score of 10, where 10 is considered as the highest score possible. Each such rating is regarded as a 'vote' by an individual registered user. Moreover, IMDb provides the possibility for the registered users to update their votes as often as they would like, but any new vote on the same title will overwrite the previous one, so in result it is one vote per title per user. For each title, the average rating from various individual users (say, R) is computed and displayed. IMDb states that they don't use the arithmetic mean (i.e. the sum of all votes divided by the number of votes) but they automatically tally the total number of user votes and report an average rating using an "in-house formula" that has been described by IMDb as ''a consistent, unbiased formula''. IMDb does not declare anything about this hidden algorithm and it rather explains:

> The simplest way to explain it is that although we accept and consider all votes received by users, not all votes have the same impact (or 'weight') on the final rating. Various filters are applied to the raw data in order to eliminate and reduce attempts at vote stuffing by people more interested in changing the current rating of a movie than giving their true opinion of it. In order to ensure that our rating mechanism remains effective, we do not disclose the exact method used to generate the rating. However, please rest assured that the same calculations are used to generate the rating for every title listed in the database: we don't adjust the rating for individual titles. There is no bias in how votes are weighted based on which title they have been cast for (IMDb, 2019e).

Moreover, IMDb displays the breakdown of the ratings, so users can see the distribution of votes and determine how uniform or polarized the opinion of a released title is. As the data about rating is based on the users' own criteria or taste for evaluating or rating the title, IMDb title pages also include a Metacritic Score data which is linked to professional critic reviews from newspapers, magazines and other publications. Aiming from this regard to offer a variety of opinions on a title so users can make informed viewing decisions.

We have also noted that IMDb makes use form the registered and contributed data by incorporating the special popular feature 'Top 250' charts of IMDb lists. The data on these lists is processed and calculated with declared formula this time. According to IMDb, This formula

provides a true 'Bayesian estimate', which takes into account the number of votes each title has received, minimum votes required to be on the list, and the mean vote for all titles (IMDb, 2019f). Moreover, we have also found that although the data in IMDb is user contributed, IMDb still saves some data about the titles to be available only for the users who subscribe to a paid service called 'IMDb pro'. Through this service, members have access to what IMDb has called ''Expanded title database'' along with other benefits, tools, and service (IMDb, 2019g).

User reviews generally come with opinions about products, in this case movie information. These reviews also contain emotional values of human being about movies (Padme & Kulkarni, 2018). It can be seen that IMDB represents an example of the ideology of dataism that in the first hand, encloses the belief in the objective quantification and potential tracking of all kinds of human behavior and sociality through online media technologies. On the other hand, this ideology involves trust in the (institutional) agents that collect, interpret, and share (meta)data (Van Dijck, 2014, p.198)

# 3. Methodology and Data Collection

In this chapter we will discuss how we collected the empirical data in our project and the relevance of the data we collected. In order to find answers to these questions we will investigate IMDb in accordance to the Knowledge Discovery from Data (KDD) process:

> KDD is the nontrivial process of identifying valid, novel, potentially useful, and ultimately understandable patterns in data. Here, data are a set of facts (for example, cases in a database), and pattern is an expression in some language describing a subset of the data or a model applicable to the subset (Fayyad, Piatetsky-Shapiro & Smyth, 1996, p.41).
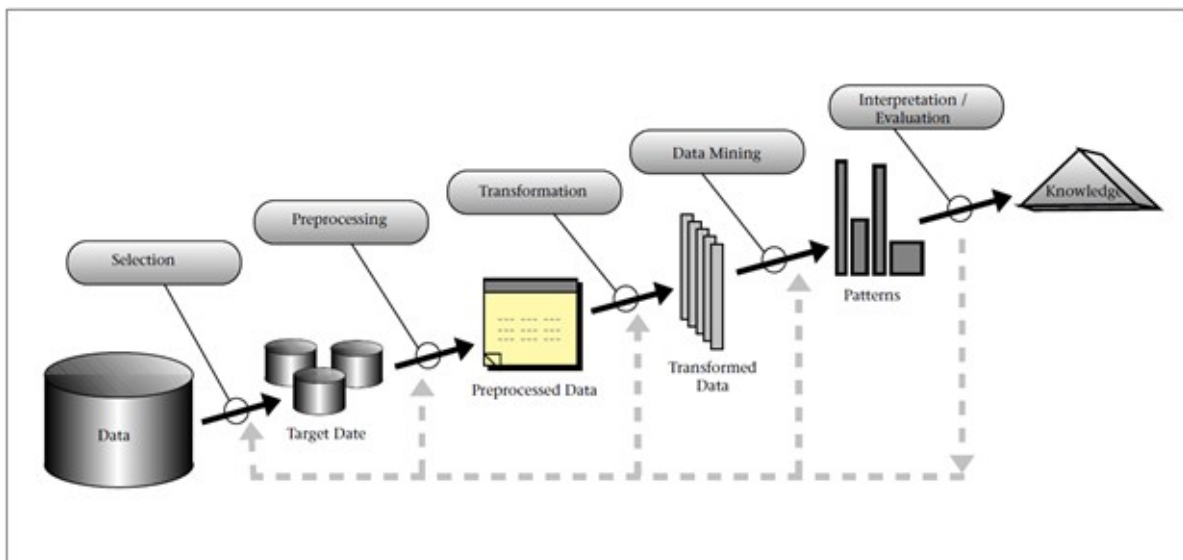


Figure 1: The five steps of the KDD process.

## **3.1** Developing an understanding of the application domain and the relevant prior knowledge and identifying the goal of the KDD:

Our starting point is the assumption that the classification of Cli-Fi movies on this movie database, can provide us not only with an interesting insight on the recognition and understanding of the Cli-Fi genre of the audience, but also in regards of what the users see as

the main driving force behind the story arc. Therefore, the goals of our data visualisations are to answer the following questions:

- How is a new emerging genre incorporated into a movie database?

- How are movies represented and rated on the platform IMDb?

- How useful is the data of IMDb for academic research?

When researching the Cli-Fi genre, many researchers of humanities and social sciences focused mainly on the movie, *The Day After Tomorrow*, as it has received far more attention than any other film, even though it is now more than 10 years old. Taking a deeper look at this particular movie, Salmose (2018) identifies the phenomena "apocalyptic sublime," as a way of actually representing the effects of climate change through the spectacle of the end of the world: "The structure of the action-adventure narrative, which frames these spectacles, negates any real impact and instead establishes a nostalgic and conservative anthropocentrism." (Salmose, 2018, p. 1417). Our first step to understand the insights that genre classifications on IMDb can provide to our research, we had a look on the process of data contribution on the platform IMDb to this particular movie, where scholars seem to have a common ground that it belongs to the cli-fi genre.

The platform itself has released very detailed information on how to contribute data, in the form of [Guides](#) and a page with collected articles called the [Contributor Zone](#). Furthermore, users can ask questions on the [IMDb Get Satisfaction community message board](#) for assistance from other users and customer service within a public forum.

In order to contribute information to a movie, users need to register with their email address. Each movie has a page where all the information is collected and displayed. Under the *Storyline* headline, users have the possibility to choose whether to add new information or to correct what is already displayed. Clicking on the "edit" button will open up forms where changes can be made (Picture 2).

Figure 2. Genre Contribution Drop Down Menu.

As a user, you can classify a movie with the help of 28 genre definitions provided by IMDb (Appendix I) - Cli-Fi, as a newly emerged genre, is not yet presented.

On their website, IMDb states:

> It should be remembered that these definitions are **guidelines** - No single definition can cover every possible eventuality. Some of these genres are objective; for the others, a little leeway is given. Either way, please try to adhere to the definitions as much as possible. Please note that the genre should relate to the main driving force behind the story arc, any sub-plots may be better represented via keywords. (IMDb, n.d.a)

So far, the displayed genres for The Day After Tomorrow are *Action*, *Adventure*, *Sci-Fi* and *Thriller*. Users can approve the genre classifications ("Correct"), delete genre classifications and add genres to a movie. For deleting ("Delete"), the submission form requires an explanation (see Picture 3).

Figure 3. Genre Contribution Submission Form.

After submitting user data, the supplied information will be processed by the IMDb team. The status of genre contributions can be checked on the Contribution History.

## 3.2 Creating a target data set focusing on data sample:

The exploration of genre classification showed that there is no possibility to add Cli-Fi as a new genre into the data input mask. Other filter options on the website is so-called keywords. IMDb writes:

> A keyword is a word (or group of connected words) attached to a title (movie / TV series / TV episode) to describe any notable object, concept, style or action that takes place during a title. The main purpose of keywords is to allow visitors to easily search and discover titles. (IMDb, n.d.b)

Looking at TDAT there are 141 total keywords. Climate related keywords are: *climate*, *climate-change*, *climate-crisis*, *climatologist*, *climatology*, *global warming*, *weather*, *natural disaster*.

When looking at the keywords, there is the possibility to rate if the keywords displayed are relevant. However, the consensus here is very little: the highest consensus among users, where relevant was rated, was n=7 with *global warming*.

As a result, we decided to go with the data sample of Svoboda's article (2016). The author used several strategies to find Cli-Fi films. On the International Movie Database (IMDb) website, keywords ("climate change" and "global warming") were searched. 16 Websites and blogs focusing on Cli-fi or climate change and popular culture were searched for the titles they listed. Other titles were found through scientific literature. And several scientists who published relevant papers were consulted via email (Svoboda, 2016). His article is a first attempt to create an overview of responses to climate change in fictional films. In his research, he was able to group 61 films produced since 1984 by the depicted climate impacts the author notes how they are unifying characteristics and distinguishing differences. Svoboda's review of the limited academic research related to Cli-Fi films reveals an emphasis of Cli-Fi moves on connecting with the lived experience of the intended audience and caution against making extreme claims or relying on appeals to fear as in his sample the elements of disaster, apocalypse and dystopia dominates.

## 3.3 Data cleaning and preprocessing:

From the sample, we decided to remove the movies which did not premiere on cinema, which made the sample end up at 28 feature films. We choose to remove these because we only wanted to focus on feature films for this research.

## 3.4 Data reduction and projection:

Different types of data can lead to different insights. We chose to aim for diverse data in order to contest the possible insights of IMDb. The following table depicts the data we aim to retrieved from the website:

| Our Data | Grouping 1 | Grouping 2 |
|---|---|---|
| Movie name | Qualitative | Nominal without order |
| Movie genre | Qualitative | Nominal without order |
| Climate change theme | Qualitative | Nominal without order |
| Release year | Quantitative | Interval |
| Rate | Quantitative | Interval from 1-10 |
| Meta score | Quantitative | Ratio from 0-100 |

Figure 4. Data variables.

## 3.5 Choose the data mining task

We chose to conduct the data by 'hand' and register in an excel sheet. In previous studies, other researchers found that:

> It is difficult to apply data mining techniques to the data in the IMDb. The data needs extensive cleaning and integration, and this consumed a large proportion of the time available for this analysis. In addition, much of the data is in textual rather than numerical format, making mining more difficult. Much of the source data could not be integrated at all, without using natural language processing techniques (Saraee, White and Eccleston, 2004, p. 351).

As a group we defined coding instructions: Go to each movie page, copy aimed data and put into an excel sheet.

## 3.6 Choose the data mining algorithm(s):

Svoboda (2016) has classified the Cli-Fi sample of movies into seven categories. Each category represents the thematic pattern which has the fictive element touched on the climate change topic which we used further for our visualisations, next to the retrieved data.

## 3.7 Data mining

During the data mining process, we experienced issues with the collection of movie genres. In a movie list, only the first 3 genres are shown with descending alphabetical order. If one clicks on the movie page, all genres are revealed. After recognising, we did an iteration and collected this data set a second time.

## 3.8 Interpreting mined patterns:

We were able to conduct a number of visualisations with our retrieved data from IMDb, here are three visualizations as an example.
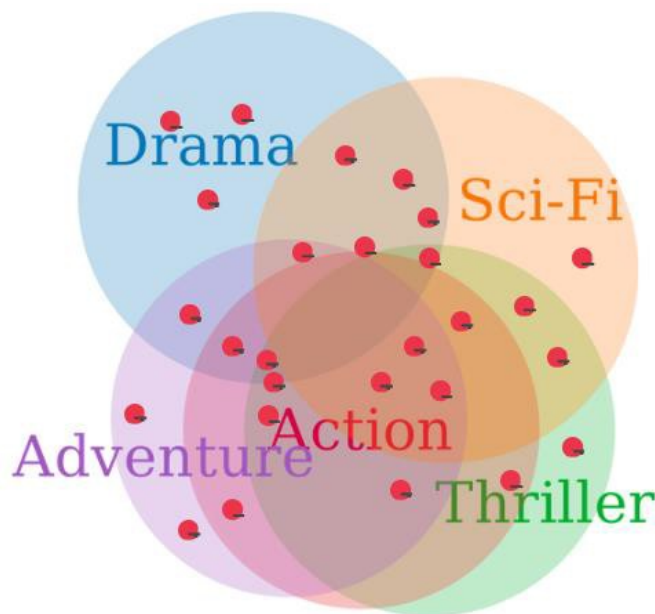


Figure 5. The combinations of genre classifications. Venn Chart.

Figure 6. The differences in movies themes over time



Figure 7. Correlation between the two variables r=0,73

## **3.9** Consolidating discovered knowledge:

Academic research in the big data area show that there is a difficulty to apply data mining techniques to the data on IMDb (see Saraee, White and Eccleston, 2004). However, extracting the data manually gave valuable insights about the user experience of the platform. Downside to this methodology is that only a small sample can be explored due to the time intensity.

Firstly, it is important to criticize the sample of Svoboda (2016). While curating the sample, the author puts a lot of emphasis on so-called experts of movies by asking only scholars or people on movie forums. This does not provide a consistent selection criteria and leaves questions of bias. The Cli-Fi movies of the sample show 1 to 6 user generated genre classifications. One can argue those movies with a high number of tags on a movie show either a discrepancy of genre understanding by the users or a very complex story arc with mixed genre elements. Furthermore, during our platform-exploration, we identified way more keywords that should be considered to investigate the climate aspect of a movie.

Most of the data we extracted is user generated data. At IMDb, users can add, review or delete genres of movies - but always in accordance to IMDb's genre definitions. Therefore, we have to revoke our assumption to get insights on the audience's recognition and understanding of the Cli-Fi genre. Firstly, Cli-Fi as a new emerging genre is not represented in the interface. So far, IMDb as a database only allows tracking movies on climate change through the keyword search, which refers to the sub-plot. As genre classification "should relate to the main driving force behind the story arc" (IMDb, n.d.a) this is excluded for climate change. This implicates that climate change can never be the main driving force of the story arc as protagonists can be, therefore it seems like the platform discriminates emerging subcultures, such as Cli-Fi, and holds on to classic genre definitions of the mainstream.

It is important to note that the mechanisms, how the IMDb team approves the contributions, is not revealed. This leads us to question the authority of the user in the data collection process. Gerlitz (2017) argues that digital media are characterized by standardization: What users can do in social media or platforms is usually pre-structured into certain forms or "grammars of action". The user data of the genre classifications do not reflect user's opinion of the driving force behind the story arc, as they were of our main interest. After deconstructing certain mechanisms of the platform it can be said that they rather reflect the user's agreement to the already existing genre definitions by IMDb, while the keyword section leaves room for individual data input for the users - thus this section can give more insights on the user's understanding of movies.

# 4. History of movie data and IMBb

As discussed in the first chapter, datafication is about quantifying what has never been quantified before or which refers to putting phenomena "in a quantified format so it can be tabulated and analyzed" (Mayer-Schonberger and Cukier, 2013, p. 77). Boyd and Crawford (2012) mention one of the first automation-based databases, where data about a specific case was collected and thus quantified. This was the US Census Bureau, which collected information and statistics about American citizens. This database was first launched as an automated data collection tool in 1890, which then quantified personal data into a large database.

The databases that was created could help with automating the data collection and usage. "Data sets that were once obscure and difficult to manage – and, thus, only of interest to social scientists – are now being aggregated and made easily accessible to anyone who is curious" (Boyd and Crawford, 2012, p. 664). Databases like IMDb would then help with organizing movie data and making data about movies more accessible to anyone who desires it.

However, with new algorithms and automation of data arising it becomes more important to understand the systems that drive these (Boyd and Crawford, 2012; Andersson, 1988). The term Big Data emerged in the 1990s (Arsenault, 2017). The first time Big Data was mentioned was in 1997 by NASA researchers Michael Cox and David Ellingsworth. They discussed the problem of big data as "data sets are generally quite large, taxing the capacities of main memory, local disk, and even remote disk. We call this the problem of big data" (Cox & Ellsworth, 1997, p.1). Big data was then first discussed when data sets became very large to be managed by any software or human input.

Different scholars agree on defining Big Data as datasets growing so large in volume that it cannot be handled or analysed by traditional analytic and data management tools (Smith, 2019; Enjolras, 2014). Enjolras (2014) also points out that Big Data emerged as a response to Web 1.0 converting into Web 2.0 where websites were more focused on interactivity. Hence, it became more common to collect data about the users of the web 2.0 platforms and datafying social situations that had not been quantified before.

Enjolras (2014) also differs between five different types of big data, these are social media data, machine-to-machine data, Big transaction data, biometric data and human-generated data. However, IMDb was created in 1990 when social media was still not a known concept. IMDb

could fit into this category of human-generated data, since it is based on the users' input on the site.

The prediction of the ability of computers to store and process data was predicted by Gordan Moore (Glass & Callahan, 2014). Known as the Moore Law, the prediction that the number of transistors on a computer chip would double every two year, allowing and ever more ability to process and store data in a cheaper way.

Within broadcasting, it has already been common to track viewing habits of how many people watch a show and for how long (Murschetz, & Prandner, 2018). By collecting data, TV broadcasting stations can collect information from the broadcasts. Today, according to Kelly (2019), Big Data is also making significant inroads within the entertainment industries. The growth of transnational services as Netflix, Hulu, Amazon Prime and others, demonstrates that Big Data is in the mind of Media Executives. Analysis of Big Data can give precise data of metrics used to lead to a greater investment in content that have higher adherence of audience (Sørensen, 2016).

Arsenault (2017) points out that with the emergence of the World Wide Web and Web 2.0, media corporations seek to balance their use of both traditional channels to using or creating new databases, where data is organized. The increase of quantification in traditional media senses, arose with the increased use of social media platforms (Arsenault, 2017). In our project we used IMDb as a database to collect data about movies. The data here is also user generated which could show IMDb as a social platform which requires some form of input of data from the user.

Boyd and Crawford (2012) discuss how big data is changing the way data is used and controlled in society. They define big data as three phenomenon. These are cultural, technological and scholarly approaches. Big data also covers the context of how data is used to control society and how it can be a powerful tool to use in society.

> On one hand, Big Data is seen as a powerful tool to address various societal ills, offering the potential of new insights into areas as diverse as cancer research, terrorism, and climate change. On the other, Big Data is seen as a troubling manifestation of Big Brother, enabling invasions of privacy, decreased civil freedoms, and increased state and corporate control. (Boyd & Crawford, 2012, p.663-664).

IMDB as a movie database shows how movies have been quantified in a similar way as described by Boyd and Crawford (2012). The movie database is similar to the big data of how society becomes quantified. Created by Col Needham (Chimielewski, 2013), the Internet Movie Database is the world's leading online source for information about movies and television shows and for celebrity news.

Oghina et. al. (2012) emphasizes that the features of a platform or database rely heavily on the underlying social patterns that form. As described by van Dijk (2014), dataism is the ideology that big data is unquestioned superior. This has been accompanied by a commercial utilization of big data. According to Richterich (2019), the ethical use of Big Data is therefore must be a key point of social interest due to power imbalances inherent to datafication.

Acording to Rohtman (2015), IMDb was launched Oct. 17, 1990. At first IMDb was organized as a private online discussion forum but was later developed to open accessed database (Chmielewski, 2013).

In 1997 Amazon.com's general counsel contacted Needham to offer a deal to integrate IMDb with Amazon. The deal was announced in April 1998. After the acquisition, IMDb was able to redesign the IMDb site and update the information daily instead of weekly. In 2002 it presented the IMDbPro subscription for entertainment industry professionals (Chmielewski, 2013).

As of May of 2019 (IMDb Press Room stats), the database has 5,980,614 titles and 9.9 million personalities. It has 83 million registered users. According to Alexa, it the 53rd site in the world according to global internet engagement. IMDb could as such be in the sector of big data with the number of users and data content it contains.

As described by William Uricchio (2017), humans have long defined, assessed, analyzed and calculated data as factors to understand how we navigate through reality. Analyzing film data is clearly distinct from analysing history and literary studies because, as stated by Olesen (2017), film scholars have tended to develop digital tools at a faster pace. Data analysis is not a new kind of practice, but the amount of data available in the recent years have "drastically increased" (van Es, Coombs & Boeshoten, 2017). This explosion of data has provided the possibility of different and more diverse social exploration practices.

# 5. The State of Data in Society

According to IMDb, user contributions are the main source of all the data available (IMDb, 2019d). Besides, the services provided on this website are based on data - starting from searching for information about a released title or casting a vote to rate the titles to the extra tools and services found in IMDb Pro. However, certain data services are only available for a monthly fee. As an example, IMDb's 'Known For' is a feature which allows IMDb Pro users to rank the top 4 relevant credits by displaying the title posters, and the default 'Known For' titles displayed at the top of an individual's IMDb page is automatically chosen through a complex weighting system. IMDb Pro members may choose these titles as they wish (IMDb, n.d.c). There are comments of users with free subscription in the IMDb community powered support questioning the credibility of this feature and its algorithm as this Pro feature could stimulate IMDB Pro signings and increase revenue. This dispute made the platform creator Col Needham himself answer that no system can be perfect (Needham, 2017).

This corresponds to what Van Dijk (2014) has argued that the metadata is generated as a result of user activity and it is believed to reflect human behaviour and the algorithms employed in this data by some sites are intrinsically selective and manipulative (p. 200). Moreover, for Burrell (2016) algorithms are "simple mathematical formulas that nobody understands'' (p.2). From this regard, scholars and practitioners across domains are increasingly concerned with algorithmic transparency and opacity, interrogating the values and assumptions embedded in automated, black-boxed systems, particularly in user-generated content platforms (Geiger, 2017, P. 1).

From this regard, we have found that the Pro features represent a particular type of "big data divide" which has been presented by Mark Andrejevic (2014). He has argued that there is a form of data divide not simply between those who generate the data and those who collect, store, and sort it, but also between the capabilities available to those two groups (p. 1674).
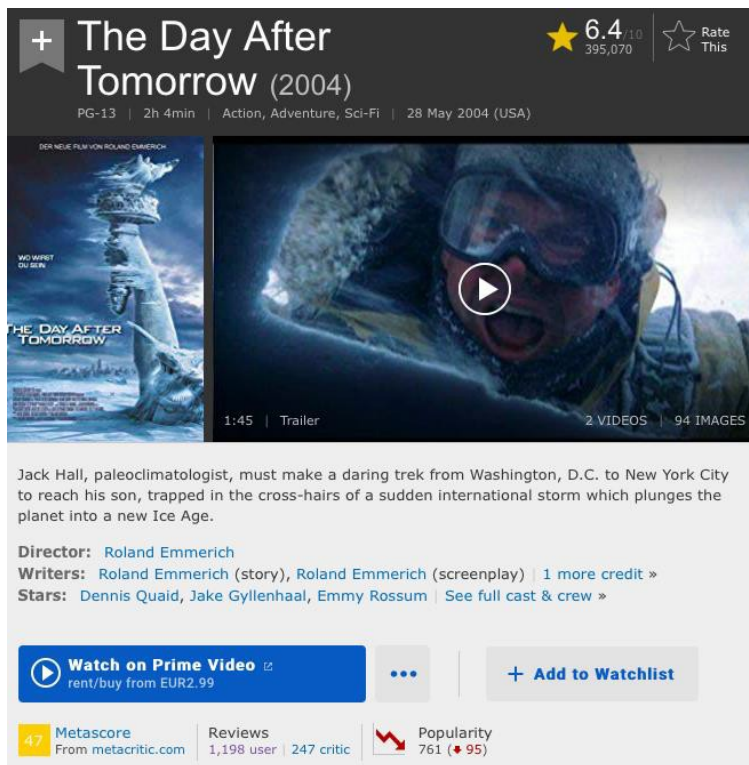
Moreover, IMDB's claim about user-generated content and the importance of users' role in IMDb is overstated as all the contributed data goes through a check-up-process by the assigned data managers of IMDb to ensure its accuracy and their role is to address the issue of the data errors (ibid.). Accordingly, this role of data managers can be seen ambiguous, since the mistakes and wrong data, according to the website itself, can still be found and evitable (IMDb, 2019c). Accordingly, Boyd & Crawford (2012) have argued that large data sets from Internet

sources are often unreliable, prone to outages and losses, and these errors and gaps are magnified when multiple data sets are used together (p. 668).

Furthermore, we have also noted that data collected on the website is not just related to the information about the released titles but, in addition, IMDb is collecting personal data from the registered users. IMDb states in its privacy notice that they collect individual's data while registration, searching, posting, participating in a contest or questionnaire or communicating with IMDb (IMDb, 2018). Moreover, the website explicitly states that all this gathered information are important to its business and to be shared with its parent corporation (Amazon.com, Inc.) (ibid).

Looking at this statement from the eye of Andrejevic (2014), we have found that users in IMDb can express the sense of ''powerlessness'' which structure the collection and use of personal information (p. 1682). Accordingly, IMDb users willingly provide all their data and information to get the access to the data they need so the sense of powerlessness here operates in two dimensions: that of ownership and control over information and communication resources, and different approaches to knowledge-based decision making (ibid.). All this again proves that the particular type of ''data divide'' which seen here is between those who are able to extract and use un-anticipatable and inexplicable findings and those who find their lives affected by the resulting decisions (ibid.).

As an Amazon company for more than 15 years, one would expect IMDb to show more obvious signs of its mother company. There are, if one looks closely but for the casual visitor, it feels very much like a separate unit without heavy Amazon branding. The connection through advertisement can be seen on the following screenshot where a direct play button for Amazon Prime, a streaming platform by Amazon, is embedded on the movie page.

Picture 8. Screenshot of Amazon-Prime-Button on IMDb

It's the subtle tip of the iceberg - the acquisition of Amazon purposely remains unveiled. IMDb disguises itself as a trustful "authoritative" information platform about movies claiming to be an authoritative source. All the emphasis on user-generated data disguises the fact that all the data is owned by Amazon, who does own movie productions with their company Amazon Studios. What users can do with their data in isolation differs strikingly from what various data collectors can do with this same data in the broader context of everyone else's data (Andrejevic 2014, p.1674).

With the help of big data practices, big data collectors can predict user sentiment and trends to ultimately incorporate them into their production and services incorporated into a marketing strategy that is successfully monetized through Amazon's famed recommendation algorithm (Andrejevic 2011 as cited in van Dijk, 2014, p.200). The deep government and industrial drive toward gathering and extracting maximal value from data, be it information that will lead to more targeted advertising, product design, traffic planning, or criminal policing (Boyd & Crawford, 2012, p.675).

Ultimately, our research shows that not only the popularization of datafication as a neutral paradigm described by van Dijck (2014, p. 206), but also the **disguise of datafication**, is carried by a belief in dataism and supported by institutional guardians of trust, gradually

yielded a view of dataveillance as a "normal" form of social monitoring is already described by van Dijck (2014, p. 206). She also recognises that metadata and data have become a regular currency for citizens to pay for their communication services and security—a trade-off that has nestled into the comfort zone of most people (van Dijk, 2014). In order to be part of the IMDb community, one has to sign up to the platform. Already by doing so the users hand over all their meta data of their movie-related activities on the platform to be capitalised upon.

## 5.1. Limitations of the work

From the sample, we decided to remove the movies which did not premiere on cinema, which made the original 61 movie sample end up at 28 feature films. Hence, we did not have a large sample to work with to create a deeper understanding of Cli-Fi and as such fully be a representative sample of the genre.

The sample we used for the paper was also constructed by another researcher, Michael Svoboda (2016). Hence the sample is very limited and not fully representative of the IMDb movie database. It is also not possible to understand if IMDb movies since it is an American website.

All movies about climate change are not included in the sample. If another researcher looks up the keyword "climate change" within IMDb, more movies come up, e.g. the catastrophe movie 2012 is one that is not included in our sample. Another limitation is that the sample of the movies does not go beyond the year 2015. This makes our sample difficult to make a just representation of climate change up until 2019. However, since we coded all the movies manually without using any external tool, 28 movies were enough to explore the possibilities of big data practices.

To further study the datafication of movie information and how genres are created, a larger sample of movies could be used, which would not only help with understanding new movie genres but also further study IMDb as a movie database.

# 6. Conclusion

We were able to unveil two major injust mechanisms that lead to what Andrejevic (2014) calls a data divide. Big Data Managers on the platform are those with access to processing power, located in an advantageous position compared to those without such access (ibid., p. 1676). Furthermore, within the user community itself, we unveiled a sort of digital hierarchy, as the Pro Users are able to buy certain data services that they can capitalize upon. This confirms the notion of "panoptic sort" which is premised on a power imbalance between those positioned to make decisions that affect the life chances of individuals and those subjected to the sorting process (Andrejevic, 2014, p. 1678).

In our data analysis, we have found that there is a correlation and similarity in scores between the movie rating and the Metascore. Besides, what we have found interesting is that usually these kinds of blockbuster movies get reviewed quite negatively by critics, but are favorably reviewed by people (See Koh, 2014), but this was not the case for the Cli-fi movies. This makes us assume that the Metascore presented on the website can be used by IMDb to prove the credibility of the results of the classifications by algorithms as they determine the visibility of movies on the platform. The strong connection to advertisement through the ownership of Amazon can lead to increased revenue. At the same time, much work is done to hide the connection to Amazon as well as the datafication processes on the platform.

Of course, we cannot generalize the assumption that the use of Metascore by the website can prove the credibility of the rating especially as the sample of movies we were using is small - however it is something to keep in mind when researching the next film to watch on IMDb.

# 7. References

Arsenault, A. H. (2017). The datafication of media: Big data and the media industries. *International Journal of Media & Cultural Politics*, *13*(1-2), pp. 7-24. DOI: 10.1386/macp.13.1-2.7_1

Anderson, C. (2008) 'The end of theory, will the data deluge makes the scientific method obsolete?', Edge, [Online] Available at: http://www.edge.org/3rd_culture/anderson08/anderson08_index.html (25 July 2011).

Boyd, D., & Crawford, K. (2012). Critical questions for big data: Provocations for a cultural, technological, and scholarly phenomenon. *Information, communication & society*, *15*(5), pp. 662-679. DOI: 10.1080/1369118X.2012.678878

Burrell, J. (2016). How the machine 'thinks': Understanding opacity in machine learning algorithms. *Big Data & Society*, 3(1), pp. 1-12. https://doi.org/10.1177/2053951715622512

Cox, M., & Ellsworth, D. (1997, October). Application-controlled demand paging for out-of-core visualization. In Proceedings. Visualization'97 (Cat. No. 97CB36155) (pp. 235-244). IEEE.

Chmielewski D. (2013). Col Needham created IMDb. Los Angeles Times, 19 January. Retrieved from: https://www.latimes.com/business/la-xpm-2013-jan-19-la-fi-himi-needham-20130120-story.html

Denvir, J. (2003). *The slotting function: how movies influence political decisions*. Vt. L. Rev., 28, 799.

Enjolras, B. (2014). Big Data og samfunnsforskning: Nye muligheter og etiske utfordringer. *Tidsskrift for samfunnsforskning*, 55(1), pp. 80-89.

Fayyad, U., Piatetsky-Shapiro, G., & Smyth, P. (1996). From data mining to knowledge discovery in databases. *AI magazine*, *17*(3), 37-37. DOI: 10.1609/aimag.v17i3.1230

Geiger, R. S. (2017). Beyond opening up the black box: Investigating the role of algorithmic systems in Wikipedian organizational culture. *Big Data & Society*, 4(2), pp. 1-14. DOI: 10.1177/2053951717730735

Gerlitz, C. (2017). Data Point Critique. In M. T. Schäfer. & K. van Es (Ed.), *The Datafied Society: Studying Culture through Data*. (pp. 241-244). Amsterdam: Amsterdam University Press.

Glass, R., & Callahan, S. (2014). *The Big Data-driven business: How to use big data to win customers, beat competitors, and boost profits*. Hoboken: John Wiley & Sons.

Igartua, J. J., & Barrios, I. (2012). Changing real-world beliefs with controversial movies: Processes and mechanisms of narrative persuasion. *Journal of Communication*, 62(3), pp. 514-531. DOI: 10.1111/j.1460-2466.2012.01640.x

IMDb (2018). *Privacy notice*. Retrieved 2019-12-21 from https://www.imdb.com/privacy?ref_=ft_pvc#auto

IMDb (2019a). *What is Imdb?*. Retrieved from: https://help.imdb.com/article/imdb/general-information/what-is-imdb/G836CY29Z4SGNMK5?ref_=helpart_nav_1#

IMDb (2019b). *IMDb Statistics*. Retrieved from: https://www.imdb.com/pressroom/stats/?ref_=helpms_ih_gi_siteindex

IMDb (2019c). *Where does the information on IMDb come from?*. Retrieved from: https://help.imdb.com/article/imdb/general-information/where-does-the-information-on-imdb-come-from/GGD7NGF5X3ECFKNN?ref_=helpart_nav_23#

IMDb (2019d). *Adding data*. Retrieved from: https://help.imdb.com/article/contribution/contribution-information/adding-data/G6BXD2JFDCCETUF4?ref_=helpart_nav_8#

IMDb (2019e). *Ratings FAQ*. Retrieved from: https://help.imdb.com/article/imdb/track-movies-tv/ratings-faq/G67Y87TFYYP6TWAV?ref_=helpart_nav_4#

IMDb (2019f). *Top Rated TV Shows*. Retrieved from: https://www.imdb.com/chart/toptv

IMDb (2019g), What is included with an IMDbPro Membership. Retrieved from: https://help.imdb.com/article/imdbpro/membership-benefits/what-s-included-with-an-imdbpro-membership/G6EFQ2AYEG3NKTM2?ref_=helpms_helpart_inline#

IMDb (n.d.a). *Genre definitions*. Retrieved from: https://help.imdb.com/article/contribution/titles/genres/GZDRMS6R742JRGAG?ref_=helpms_helpart_inline#

IMDb (n.d.b). *Keywords*. Retrieved from: https://help.imdb.com/article/contribution/titles/keywords/GXQ22G5Y72TH8MJ5?ref_=helpms_helpart_inline#

IMDb (n.d.c). *Known for the title selection*. Retrieved from:
https://help.imdb.com/article/imdb/discover-watch/known-for-title-selection/GNL2E4LJVKM9QLMD?ref_=helpsrall#

Intergovernmental Panel on Climate Change (2014). Climate change 2014: Synthesis report (Contribution of Working Groups I, II and III to the Fifth Assessment Report of the Intergovernmental Panel on Climate Change). Geneva, Switzerland: Author.

Kelly, J. P. (2019). Television by the numbers: The challenges of audience measurement in the age of big data. *Convergence*, *25*(1), pp. 113-132. DOI 10.1177/1354856517700854.

Koh, W. (2014). 'I am Iron Man': the Marvel Cinematic Universe and celeactor labour, *Celebrity Studies*, 5(4), pp. 484-500. DOI: 10.1080/19392397.2014.933675

Murschetz, P. C., & Prandner, D. (2018). 'Datafying'Broadcasting: Exploring the Role of Big Data and Its Implications for Competing in a Big Data-Driven TV Ecosystem. In *Competitiveness in Emerging Markets* (pp. 55-71). Springer, Cham.

Naun, C. C., & Elhard, K. C. (2005). Cataloguing, lies, and videotape: Comparing the IMDb and the library catalogue. *Cataloging & classification quarterly*, *41*(1), 23-43. DOI: 10.1300/J104v41n01_03

Needham, C. (2017). 'Known for' manipulation?. *IMDb community powered support for IMDb.com* [forum]. Retrieved from: https://getsatisfaction.com/imdb/topics/known-for-manipulation?topic-reply-list[settings][filter_by]=all&topic-reply-list[settings][reply_id]=19289426#reply_19289426

Oghina, A., Breuss, M., Tsagkias, M., & De Rijke, M. (2012, April). Predicting imdb movie ratings using social media. In *European Conference on Information Retrieval* (pp. 503-507). Springer, Berlin, Heidelberg.

Olesen, C. G. (2017). Towards a 'Humanistic Cinemetrics'?. In M.T. Schäfer. & K. van Es (Ed.), *The Datafied Society: Studying Culture through Data* (pp. 39-54). Amsterdam: Amsterdam University Press.

Padme, S. & Kulkarni, P. (2018). *Aspect Based Emotion Analysis on Online User-Generated Reviews*. 9th International Conference on Computing, Communication and Networking Technologies (ICCCNT), Computing, Communication and Networking Technologies (ICCCNT), 2018 9th International Conference On, 1. DOI: 10.1109/ICCCNT.2018.8494183

Peralta, V. (2007). Extraction and integration of movielens and imdb data. *Laboratoire Prisme, Université de Versailles, Versailles, France.*

Rohtman, L. (2015, October 16). How IMDb Can Be Older Than the First Web Browser, *Time*. Retrieved from: https://time.com/4068036/imdb-history-25th-anniversary/

Richterich, A. (2019). The Big Data Agenda: Data Ethics and Critical Data Studies. London: University of Westminster Press.

Salmose, N. (2018). The Apocalyptic Sublime: Anthropocene Representation and Environmental Agency in Hollywood Action-Adventure Cli-Fi Films. *The Journal of Popular Culture*, 51(6), pp. 1415-1433. DOI: 10.1111/jpcu.12742

Saraee, M., White, S., & Eccleston, J. (2004). A data mining approach to analysis and prediction of movie ratings. *WIT Transactions On Information And Communication Technologies*, *33*.

Smith, B. (2019). Big Data and Us: Human–Data Interactions. European Review, 27(3), pp. 357-377. doi:10.1017/S1062798719000048

Svoboda, M. (2016). Cli-fi on the screen(s): patterns in the representations of climate change in fictional films. WIREs Clim Change, 7:43–64. DOI: 10.1002/wcc.381

Sørensen, IE (2016) The revival of live TV: Liveness in a multiplatform context. Media, Culture & Society 38(3): 381–399. DOI: 10.1177/0163443715608260

Uricchio, W. (2017). Data, Culture and the Ambivalence of Algorithms. In M.T. Schäfer. & K. van Es (Ed.), *The Datafied Society: Studying Culture through Data* (pp. 125-137). Amsterdam: Amsterdam University Press.

van Dijck, J. (2014) Datafication, dataism and dataveillance: Big data between scientific paradigm and ideology. *Surveillance and Society,* 12(2): 197–208. DOI: 10.24908/ss.v12i2.4776

van Es, K., Coombs, N. L. & Boeschoten, T. (2017). Towards a Reflexive Digital Data Analysis. In M.T. Schäfer. & K. van Es (Ed.), *The Datafied Society: Studying Culture through Data* (pp. 171-180). Amsterdam: Amsterdam University Press.

Weible, C. L. (2001). The Internet Movie Database: A reference guide to Hollywood and beyond. *Internet reference services quarterly*, 6(2), 47-50. DOI: 10.1300/J136v06n02_05

Whiteley, A., Chiang, A., & Einsiedel, E. (2016). Climate Change Imaginaries? Examining Expectation Narratives in Cli-Fi Novels. *Bulletin of Science, Technology & Society*, 36(1), 28–37. DOI: 10.1177/0270467615622845

# Appendix I: Movie Genres by IMDb

| Genre name | Genre Definition |
| --- | --- |
| Action | Should contain numerous scenes where action is spectacular and usually destructive. Note: if a movie contains just one action scene (even if prolonged, i.e. airplane-accident) it does not qualify.Subjective. |
| Adult | Reserved for explicit works of consenting hardcore sex or sexual activity, or strong fetish material involving adults, specifically those created for the purposes of titillation or arousal. Must be used with the plot keywords of 'hardcore' and 'sex' or 'special-sexual-interest'. Documentaries about this type of material do not use the Adult genre. Subjective. |
| Adventure | Should contain numerous consecutive and inter-related scenes of characters participating in hazardous or exciting experiences for a specific goal. Not to be confused with Action, and should only sometimes be supplied with it. Subjective. |
| Animation | Over 75% of the title's running time should have scenes that are wholly, or part-animated. Any form of animation is acceptable, e.g., hand-drawn, computer-generated, stop-motion, etc. Puppetry does not count as animation, unless a form of animation such as stop-motion is also applied. Incidental animated sequences should be indicated with the keywords part-animated or animated-sequence instead. Although the overwhelming majority of video games are a form of animation it's okay to forgo this genre when adding them as this is implied by the title type. Objective. |

| | |
|---|---|
| Biography | Primary focus is on the depiction of activities and personality of a real person or persons, for some or all of their lifetime. Events in their life may be reenacted, or described in a documentary style. If re-enacted, they should generally follow reasonably close to the factual record, within the limitations of dramatic necessity. A real person in a fictional setting would not qualify a production for this genre. If the focus is primarily on events, rather than a person, use History instead. Objective. |
| Comedy | Virtually all scenes should contain characters participating in humorous or comedic experiences. The comedy can be exclusively for the viewer, at the expense of the characters in the title, or be shared with them. Please submit qualifying keywords to better describe the humor (i.e. spoof, parody, irony, slapstick, satire, black-comedy etc). If the title does not conform to the 'virtually all scenes' guideline then please do not add the comedy genre; instead, submit the same keyword variations described above to signify the comedic elements of the title. Subjective. |
| Crime | Whether the protagonists or antagonists are criminals this should contain numerous consecutive and inter-related scenes of characters participating, aiding, abetting, and/or planning criminal behavior or experiences usually for an illicit goal. Not to be confused with Film-Noir, and only sometimes should be supplied with it. Subjective. |
| Documentary | Should contain numerous consecutive scenes of real personages and not characters portrayed by actors. This does not include fake or spoof documentaries, which should instead have the fake-documentary keyword. A documentary that includes actors re-creating events should include the keyword "reenactment" so that those actors are not treated as "Himself." This genre should also |

| | |
|---|---|
| | be applied to all instances of stand-up comedy and concert performances. Objective. |
| Drama | Should contain numerous consecutive scenes of characters portrayed to effect a serious narrative throughout the title. This can be exaggerated upon to produce melodrama. Subjective. |
| Family | Should be universally accepted viewing. e.g., aimed specifically for the education and/or entertainment of children or the entire family. Note: Usually, but not always, complementary to Animation. Objective. |
| Fantasy | Should contain numerous consecutive scenes of characters portrayed to effect a magical and/or mystical narrative throughout the title. Note: not to be confused with Sci-Fi which is not usually based in magic or mysticism. Subjective. |
| Film-Noir | Typically features dark, brooding characters, corruption, detectives, and the seedy side of the big city. Almost always shot in black and white, American, and set in contemporary times (relative to shooting date). We take the view that this genre began with Underworld (1927) and ended with Touch of Evil (1958). Note: neo-noir should be submitted as a keyword instead of this genre for titles that do not fit all criteria. Objective. |
| Game-Show | Competition, other than sports, between, usually, non-professional contestants. The competition can include a physical component, but is usually primarily mental or strategic as opposed to athletic. This also includes what are known as "quiz |

| | |
|---|---|
| | shows." Talent contests staged expressly for the program are considered Game-Shows. Objective. |
| History | Primary focus is on real-life events of historical significance featuring real-life characters (allowing for some artistic license); in current terms, the sort of thing that might be expected to dominate the front page of a national newspaper for at least a week; for older times, the sort of thing likely to be included in any major history book. While some characters, incidents, and dialog may be fictional, these should be relatively minor points used primarily to bridge gaps in the record. Use of actual persons in an otherwise fictional setting, or of historic events as a backdrop for a fictional story, would not qualify. If the focus is primarily on one person's life and character, rather than events of historical scope, use Biography instead. Objective. |
| Horror | Should contain numerous consecutive scenes of characters effecting a terrifying and/or repugnant narrative throughout the title. Note: not to be confused with Thriller which is not usually based in fear or abhorrence. Subjective. |
| Musical | Should contain several scenes of characters bursting into song aimed at the viewer (this excludes songs performed for the enjoyment of other characters that may be viewing) while the rest of the time, usually but not exclusively, portraying a narrative that alludes to another Genre. Note: not to be added for titles that are simply music related or have music performances in them; e.g., pop concerts do not apply. Also, classical opera, since it is entirely musical, does not apply and should instead be treated as Music. Objective. |

| Music | Contains significant music-related elements while not actually being a Musical; this may mean a concert, or a story about a band (either fictional or documentary). Subjective. |
|---|---|
| Mystery | Should contain numerous inter-related scenes of one or more characters endeavoring to widen their knowledge of anything pertaining to themselves or others. Note: Usually, but not always associated with Crime. Subjective. |
| News | Reports and discussion of current events of public importance or interest. If the events are not current (at the time the title was initially released), use History instead. This generally includes newsreels, newsmagazines, daily news reports, and commentary/discussion programs that focus on news events. Objective. |
| Reality-TV | Often, but not always, features non-professionals in an unscripted, but generally staged or manipulated, situation. May or may not use hidden cameras; generally, but not always, in a non-studio setting. Objective. |
| Romance | Should contain numerous inter-related scenes of a character and their personal life with emphasis on emotional attachment or involvement with other characters, especially those characterized by a high level of purity and devotion. Note: Reminder, as with all genres if this does not describe the movie wholly, but only certain scenes or a subplot, then it should be submitted as a keyword instead. Subjective. |
| Sci-Fi | Numerous scenes, and/or the entire background for the setting of the narrative, should be based on speculative scientific discoveries or developments, environmental changes, space travel, or life on other planets. Subjective. |

| | |
|---|---|
| Short | Any theatrical film or made-for-video title with a running time of less than 45 minutes, i.e., 44 minutes or less, or any TV series or TV movie with a running time of less than 22 minutes, i.e. 21 minutes or less. (A "half-hour" television program should not be listed as a Short.) If known, please submit the running time if we do not have one on record. Objective. |
| Sport | Focus is on sports or a sporting event, either fictional or actual. This includes fictional stories focused on a particular sport or event, documentaries about sports, and television broadcasts of actual sporting events. In a fictional film, the sport itself can also be fictional, but it should be the primary focus of the film. Objective. |
| Talk-Show | Discussion or interviews of or with a series of guests or panelists, generally appearing as themselves in a non-fictional setting (though fictional programs that mimic the form are also included). (aka "chat show"). Objective. |
| Thriller | Should contain numerous sensational scenes or a narrative that is sensational or suspenseful. Note: not to be confused with Mystery or Horror, and should only sometimes be accompanied by one (or both). Subjective. |
| War | Should contain numerous scenes and/or a narrative that pertains to a real war (i.e., past or current). Note: for titles that portray fictional war, please submit it as a keyword only. Objective. |
| Western | Should contain numerous scenes and/or a narrative where the portrayal is similar to that of frontier life in the American West during 1600s to contemporary times. Objective. |

# Appendix II: List of figures