# An Analysis of the Transformations in the Movie Industry

University ID: 2222991

January 10, 2024

## Abstract

This paper gives a data analysis of the movie industry from 1980 to 2020 in the area of revenue trends and cultural changes. By using datasets from IMDb and Netflix, this project studies the changes of movie industry and confirms several hypothesises. The increasing trend in budgets and gross revenues over the years is shown, attributed to globalisation and the expansion of global movie market. The average length of popular movies shows an upward trend even under the era of human's decreasing attention spans. There is also a trend of popular movies becoming more concentrated in genres. Classification models have been applied to identify key factors contributing to a movie's success, such as runtime, budget, and genre. These findings offer insights about the evolution of movie industry from 1980 to 2020, and can serve as a foundation for potential future research areas, such as the changes in movie industry after COVID-19 pandemic and in the age of streaming platforms.

## 1 Introduction

SINCE the late 19th century, movies have evolved from a form of entertainment to a major cultural and artistic medium in human life, by their unique visual storytelling ability. As progressed into 20th and 21st centuries, people's lives have been rapidly changed due to the technology revolution and globalisation. These also profoundly impacted the movie industry, which not only altered significantly on the way how movies are produced and consumed, but also transformed the role of movies in society and culture, which leaded different opinions. Some argue that these shifts represent the decline of movie industry. However, others suggested that these changes are just merely adaptive evolution. They argue that the continuous evolution of films, reflected the resilience of movie industry, just the same case like the past when filmmakers embraced the transition of film production from analog to digital media.

In this context, this project aims to explore the changes in the movie industry, by analysing datasets of movies from past decades. It seeks to understand how changes of environmental factors have influenced the movie industry through the data of popular films from 1980 to 2020.

## 2 Background

Before analysing the changes in the cinematic world over the years, it is crucial to understand following background factors.

### 2.1 Globalisation

Starting from post-cold war era, the increased interactions between nations were marked as a new age of globalisation. This had a significant impact and influenced the movie industry, as the markets outside Western sphere, such as Asian countries like China and India, began to expand constantly. While Hollywood studios have remained to be the major film producers globally over the years, their revenue sources have been shifted from U.S. to more international, where the global box office now contributes a larger share of the total revenue, which surpassing U.S. revenue. Figure 1 shows the international box office revenues accounted for over 60% of the global total from 2010 to 2014, and Figure 2 shows the international box office revenues accounted for over 70% of the global total from 2015 to 2019. This shift towards the international markets had transformed the movie industry, enabling those successful films now to achieve higher profitability on a worldwide scale, which keeps on enlarging the scale of Hollywood movie productions.
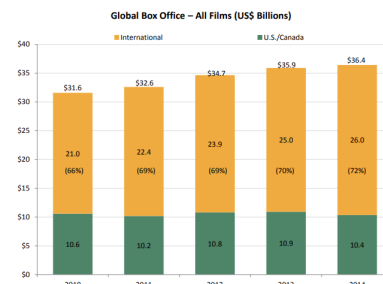


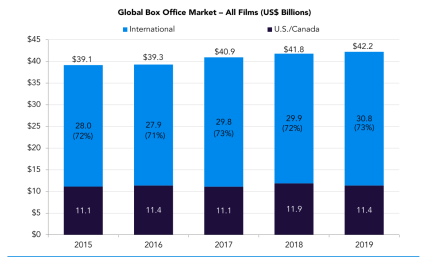Figure 1: Global Box Office – All Films (US$ Billions) 2010 - 2014 from 2014 Theatrical Statistics Summary [1]

Figure 2: Global Box Office Market – All Films (US$ Billions) 2015 - 2019 from "THEME REPORT", 2019 [2]

## 2.2 Internet and Streaming Platform

As the most important invention of the 20th century, internet has dramatically changed the world, including the movie industry. The first major problem presented was the issue of piracy, while consumers began downloading films online illegally. Fortunately, as public awareness of intellectual property rights grew, this issue was gradually mitigated. Nowadays, viewers are more likely to pay for different online streaming services to watch movies or TV contents online. This makes another significant shift in the movie industry's revenue models. Unlike traditional box office income, which is generated and calculated from ticket sales, streaming platforms rely on subscription-based models. Filmmakers have to adapt to this shift in financial models where box office returns are no longer the only factor in revenue generation, which increased the difficulty in measuring the success of films, or even the changes of strategies on content production, in order to keep the viewers who watch on streaming platform to be satisfied.

## 2.3 Internet and Streaming Platform

As the most important invention of the 20th century, internet has dramatically changed the world, including the movie industry. The first major problem presented was the issue of piracy, while consumers began downloading films online illegally. Fortunately, as public awareness of intellectual property rights grew, this issue was gradually mitigated. Nowadays, viewers are most likely to pay for different online streaming services to watch movies or TV contents on various devices. This makes another significant shift in the movie industry's revenue models. Unlike traditional box office income, which is generated and calculated from ticket sales, streaming platforms rely on subscription-based models. Filmmakers have to adapt to this shift in financial models where box office returns are no longer the only factor in revenue generation, which increased the difficulty in measuring the success of films, or even the changes of strategies on content production, in order to keep the viewers who watch on streaming platform rather than in cinema to be satisfied.

## 2.4 Short Attention Spans

A significant decrease in human attention spans was led by the use of technology and social media in nowadays. According to the study [3], it is suggested that people in nowadays generally lose their concentration in just eight seconds, primarily due to the heavy use of mobile devices and especially social media. Although this finding of the reduction of attention is task-dependent, it still shows the deep influence of technology on our behaviour. When it comes to movies, the rise of short video format media platforms, such as TikTok, is also transforming audience expectations on watching movies. Filmmakers now have to face the challenge of adapting the movies to this new era of audiences with shortened attention spans. Innovative storytelling techniques, more attractive content, or even shorten the length of movies have to be considered to capture the viewer's attention.

## 2.5 Change in Film Genres

Another major transformation is the shift in the popularity of film genres over the past years. From 1930s to 1950s, the mid-20th century was marked as the golden age of musical films in Hollywood, and in contrast there are only a few musical films on today's market. Modern viewers now prefer films like blockbusters with action element or fantasy narratives. This reflects changing cultural trends, while the movie industry on the one hand is driving the evolution of popular genres, but on the other hand filmmakers have to responding to it, by adapting the change in order to produce content which align the different expectations of global audience. This adaptation is also a key to be successful in the global market of nowadays movies.

# 3 The Data

## 3.1 Movies from 1980 to 2020

To analyse the transformation of the movie industry over past decades, a comprehensive dataset of past movies which focusing on revenues is being used.

This dataset [4] was publicised on Kaggle and GitHub, which sourced from Internet Movie Database(IMDb), a reputable online movie database that not only provides different information of films and television programmes, but also serves as a user rating platform in which makes the scores a valuable indicator of a movie's popularity and quality.

The dataset contains 7,668 popular films ranging from 1986 to 2020, with about 220 films featured per year. It has 15 attributes across three categories: movie information (e.g., name, rating, genre, director), revenue statistics (e.g., budget, gross), and audience reception (e.g., votes, score). This rich combination of data offers an comprehensive view of the movie industry, allow for a

detailed data analysis on its changes in terms of financial success and audience preferences along the years.

## 3.2 Netflix Movies and TV Shows

To explore the relationship between the movie industry and streaming platforms, another additional dataset that covers Netflix Movies and TV Shows was used. This dataset [5] was also publicised on Kaggle and sourced from Netflix, a global leading streaming platform known for providing on-demand content, which also includes its original productions.

This dataset contains 8,807 entries of which 6,131 are movies and the remaining are TV shows. This dataset has 12 attributes, most of them are textual, such as title, director, cast, country, and description. This dataset acts as a secondary focus of the project, which emphasising the study of the streaming platform and its relationship of the movie industry.

# 4 Hypothesis

In the context of the evolving movie industry as discussed in the background section, the following hypothesises have been drawn to be explored in this project.

## 4.1 Increasing Budgets and Grossing Over Time

As discussed earlier, globalisation has facilitated Hollywood movies to expand beyond Western nations into the global markets, such as those in Asia. This expansion has not only escalated the revenues, but such financial gains also implied an increase in both film production budgets and gross earnings. Evidence of such trend can be seen from those blockbusters in recent years: The use of state-of-the-art production technology, the involvement of famous actors and actresses, and the huge advertising campaigns. These factors act as indicators of the growing financial capacity of the movie industry. Therefore, this posits a hypothesis that there has been a increasing trend in both budgets for film production and gross earnings along the years.

## 4.2 Increasing Average Film Length Over Time

Despite filmmakers are having challenges in creating movies under the trend of dominating short media formats, which caused by collapsing attention spans, some studies show a contrary result that the film industry appears to be in the opposite direction. The news article by John Renda shows that the average duration of the top 25 popular films from 1931 to 2013 is actually on a upward trend [6]. Renda suggested in the article that the reason

behind this is because films now are providing an immersive experiences rather than just products for viewers. The longer duration helps to enhance the satisfaction of audience in their immersive experiences. Therefore, this hypothesis assumes actually there has been an increasing trend in the average length of popular films over time, to prove the shift of the approach Renda suggested of the movie industry which to meet audience expectations for more engaging cinematic experiences.

## 4.3 Increasing Concentration in High-grossing Film Genres Over Time

While movie genres of popular films have evolved over the years, there is another potential trend which is the increasing concentration of high-grossing films, turns into fewer genres. A wide range of genres such as drama, musicals, comedy and romance have been welcomed by viewers in history, and in contrast nowadays audiences seems to enjoy films in genres like action and anime, showing a narrower of focus in genres. This hypothesis assumes the increasing of concentration in genres for popular films over the years and the project aims to verify this hypothesis by analysing the data of movie genres.

## 4.4 Identifying Characteristics of Successful Films

By analysing the scores provided on IMDb, which serve as audience rating and a measure of the film's reputation, it is possible to identify different characteristics of successful films. This project will explore different conditions and patterns appear in those high-rating movies, aiming to find those hidden factors contributing to the success. This hypothesis assumes there are certain attributes, which significantly influencing the likelihood for a film to have a high reputation. Those attributes may not be obvious, but this project tends to discover those key elements by different classification methods.

# 5 Statistical software

The data analysis process in this project used a variety of software and tools, and each of them served specific purposes.

1. Sublime Text

   Sublime Text is an open-source text editor which is famous among developers for its speed and powerful text-editing features, such as multiple selections and text replacement. Due to the efficiency of Sublime Text in handling large text files, it was used for rapid text search and global text replacement in the dataset.

2. Excel

Microsoft Excel is a well-known spreadsheet editing software with different computational functions. Though it was not used for in-depth data analysis during the project, it helped on quick data manipulation and visualisation.

3. MySQL

MySQL is an open-source relational database management system (RDBMS) which utilising Structured Query Language (SQL) for data management. It helped for quick querying of the results after importing the dataset into the database.

4. R/RStudio

R is a famous programming language specifies for statistical computing. R served as a major tool in this project, with it's extensive libraries for statistical computing and data visualisation. On the other hand, RStudio provided an easy-to-use interface for these functions.

5. Weka

As a Java-based software for data analysis and machine learning, Weka was used on tasks such as classification analysis and visualisation during this project. Its user-friendly interface was beneficial during the classification process.

# 6 Data Cleaning and Imputation

To ensure the suitability of the dataset for study, data cleaning and preprocessing are crucial before the data analysis phase.

Initially, there are 189 records identified with null value for gross revenue. As this is a key variable in analysis, those records were removed, which is about 2.5% of the total 7,668 entries of the dataset.

However, 2,171 records were also discovered that there are missing value of budget information, which are approximately 29% of the remaining 7,479 entries. Given the large proportion of these records, simply deleting those entries was not suitable as it could cause potential data imbalance and bias in the results. To address this issue, regression imputation was used.

Regression imputation is recognised as one of the best approaches to effectively estimate the missing values in a dataset [7], especially suitable in this case after a positive correlation was established between budget and gross revenue. A linear regression model was developed by using complete records with budget and gross revenue, then the same model served to estimate and fill in the missing budget values of the 2,171 entries. This helps to keep the integrity of the dataset and minimise the biases in the analysis.
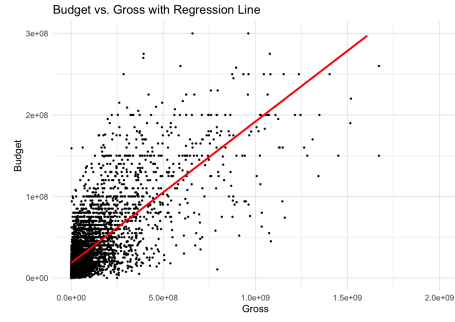


Figure 3: Budget vs. gross revenue with regression line

Additional and minor data cleaning tasks were performed to ensure the compatibility of the dataset with those different software and tools. Such tasks included correcting different typo of textual values, as such errors of quotation marks could affect the data processing giving the CSV format of the dataset. This ensure of textual and format accuracy of the dataset was crucial, as those minor mistakes could cause misinterpretation by software and tools during analysis, and affect the integrity and reliability of the result.



Figure 4: Typo of quotation marks in the original dataset

# 7 Data Analysis

## 7.1 Overall Analysis

After the completion of data cleaning and preprocessing, the dataset was ready to be analysed. First of all, the dataset was explored in overall for preliminary findings. By using R and ggplot2 package, a heatmap was created to display the correlations among all numerical attributes, including runtime, gross, budget, votes, score, and year. This visualisation helps to identify relationships for in-depth analysis afterwards.
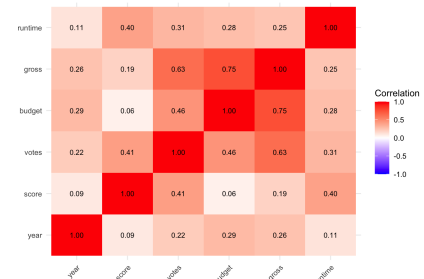


Figure 5: Heat-map of correlations among all numerical attributes of the dataset

As observed, attributes 'gross' and 'budget' have the highest positive correlation (0.75) in the dataset, indicating a strong relationship between films' budget and gross revenue. 'votes' and 'gross' also show a positive correlation (0.63) which suggests high grossing films tend to have more viewers to vote on IMDb. On the other hand, variable 'year' shows positive correlation with 'budget' (0.29), 'gross' (0.26), 'votes' (0.22) and 'runtime' (0.11).

This heatmap showed a comprehensive overview of how variables interact with others. These preliminary findings aligned with those hypotheses, and will be further explored in subsequent sections. Overall, this step was an crucial foundation to ensure the direction of the study.

## 7.2 Analysis of Movies on Streaming Platform

To analyse movies available on streaming platforms, two datasets were integrated together. This step involved the using of MySQL, where both datasets which in CSV format, were imported as tables into the database. Following SQL query has been executed for the merging:

```sql
CREATE TABLE result AS
SELECT *
FROM movies
INNER JOIN netflix_content
ON movies.name = netflix_content.title AND
    movies.year = netflix_content.release_year
WHERE netflix_content.type = 'Movie'
```

JOIN operation has been used to merge data based on the common attributes. To extract the exact movies in both datasets, not only the movie name was being used as the identifier, same release year was used as another condition for the joining. Moreover, only movie type of records are extracted in the table of Netflix's content. The result then was stored in a result table for follow-up study.

852 records were created in the result table, and in those records, most of the movies are from 2000 to 2020. This indicated that the popular films in the first dataset, which are also available for streaming on Netflix, are mostly released after 2000. While this may due to copyright constraints, as only Netflix was used as an example, it still highlighted the trend where streaming platform focus on modern popular films rather than classical films.
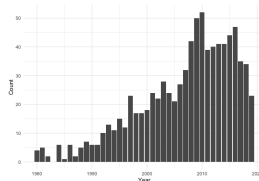


Figure 6: Number of films in the first dataset which are on Netflix according to the release year

## 7.3 Correlation Analysis

### 7.3.1 Year vs Budget and Year vs Gross Revenue

A low but positive correlation between year and both budgets and gross revenues was observed, supporting the hypothesis. This result suggests an upward trend in both budgets and gross revenues over the years. Refer to figure 7 and 8, which show the average budgets and gross revenues from each year plotted, along with the linear regression model created. These illustrated how the industry has seen a rise in both financial investments and returns from 1980 to 2020.
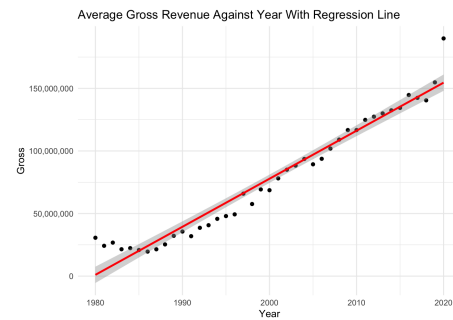


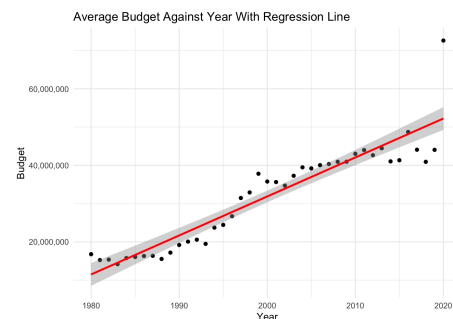Figure 7: Average Gross Revenue Against Year With Regression Line



Figure 8: Average Budget Against Year With Regression Line

### 7.3.2 Year vs Runtime

Though a relatively low correlation coefficient (0.11) between attribute 'year' and 'runtime' shown above, the trend becomes apparent when the plot of average runtime of those popular films against year was visualised. This gradual increase of average runtime over the years helped to prove the hypothesis that the shift of practices to make popular films becoming a long cinematic experiences for audiences. Moreover, to address this finding, actually there are other studies [8] show the change is not on the length of the movies, but on the movie's pace: the increase of short-duration scenes and the number of

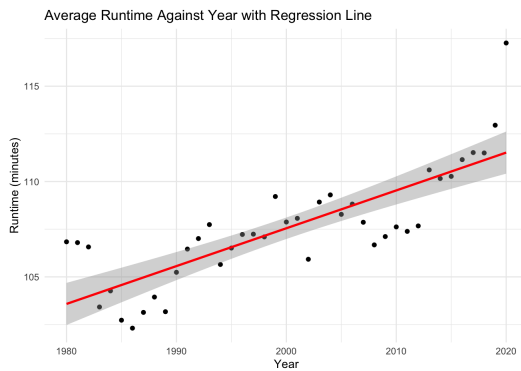cuts. These are the true factors which engaging viewers in nowadays movies.



Figure 9: Average Runtime Against Year With Regression Line

### 7.3.3 Genre Diversity Over Time

To address the hypothesis of the concentration of film genres over time, it was necessary to covert the textual genre data in the dataset into numerical data. First of all, the number of occurrences of each genre by year was grouped and counted by using R. This helped to transform genre data into a quantifiable format, and enabling the follow-up statistical analysis.



Figure 10: Genre information of films in the dataset after the transformation

The next step is to normalise the data, as the total number of movies per year are different. To compare the data across different years, the genre counts were converted to proportions, which means the percentage representing the share of films in each genre per year. This is essential for making the comparisons.

To explore genre dispersion, standard deviation of genre proportions is calculated. This statistical measure helped to assess and compare the diversity of genres across the years. In such case, a low value of standard deviation shows high concentration of films in genres, suggesting the lack of diversity. On the other hand, a high value of standard deviation value suggested diverse genre distribution. The following figure illustrated the genre concentration trends over time with the regression model.



Figure 11: Standard Deviation of Genres Against Year With Regression Line

Despite the relatively large confidence interval around the smooth line of the regression model indicating some uncertainty, the decreasing trend in standard deviation from 1980 to 2020 suggests an increasing concentration in genres of nowadays movies, helps to support the hypothesis.

To further study the transformation of the movie genres, following figure reveals some some key genres during the period of time. The rise in action movies, the decline in comedy movies and the steady increase of animated movies were shown. The consistent proportion of drama genre along the years also shows its importance throughout the cinematic history.



Figure 12: Proportion of Particular Genres Over Years

## 7.4 Classification

To study the relationship between various movie attributes and the success of movies, different classification models were used in Weka for analysis. The dataset was first refined to divide films into 'Good' and 'Not Good' categories based on their IMDb scores. Given that the mean score of all movies in the dataset is 6.399 and the third quartile (Q3) is 7.100, movies were categorised as 'Good' with scores above 7.100. This distinguish of top 25% of the records to be 'Good' movies is a reasonable approach for analysis.



Figure 13: Summary of the score attribute in the dataset

Then, attributes had to be strategically chosen or excluded to build the classification model. First of all, the attribute 'score' is closely related to the attribute 'votes', which has to be omitted as otherwise the results would be rely heavily on this attribute. Secondly, text attributes such as movie name, director, writer, star, country and company were excluded. Only meaningful variables were focused such as rating, genre, year, budget, gross revenue and runtime.

By using Weka with the setting of 66% training split, various classification models were built and tested. The Naive Bayes classification gave 76.3965% accuracy rate, while Support Vector Machine (SVM) classification slightly improved the accuracy rate to 77.144%. The SVM classification also gave us information that 'gross', 'runtime', and 'budget' are the three most important factors (Figure 14). While the Reduced Error Pruning Tree (RepTree) model gave 77.38% accuracy rate, if the model was simplified to a 3-level depth the accuracy rate will be reduced to 76.7506, but the result visualised which was more readable. It highlighted 'runtime' and 'genre' as key determinants in classifying a movie to be 'Good' (Figure 15). These models offered valuable insights which attributes such as runtime, budget and genres are the most influential in success of movies in order to achieve high score, which means high reputation.
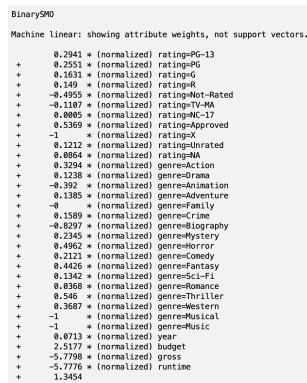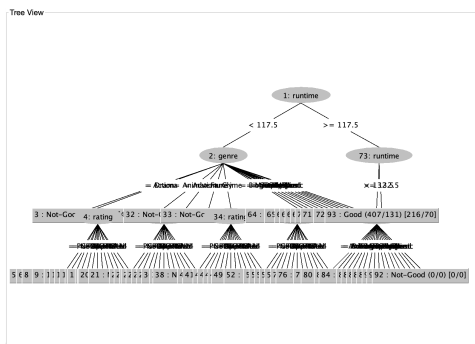


Figure 14: Result of SVM classification



Figure 15: Result of RepTree classification with 3-level depth

# 8 Conclusion

This project provides a comprehensive, data-driven analysis of the movie industry's trends and culture transformation from 1980 to 2020, validating different hypotheses with statistical evidence.

It shows a positive correlation between the years and both the budget and gross revenues in the ear of globalisation. Opposite to the trend of shrinking attention spans nowadays, the study shows there is a slight increase in the average length of popular movies throughout the years. It also shows the increase of concentration of movies genres compared with earlier time. Finally, different classification models were used to find the critical factors which make a movie's success. These findings offer different insights and can be used to guide film productions to align the evolving movie trends.

## 8.1 Limitation

This project's major limitation comes from the incompleteness of the datasets. For example, the data of 2020 includes only 15 films, which is much fewer compared to other years. This discrepancy makes the films from 2020 to be outliers, especially when recent data plays a crucial role in trend analysis.

Additionally, the dataset has not been updated since 2020. This limits the study's scope, and most importantly it excludes the COVID-19 pandemic period and the impact on the movie industry when cinema have been shutdown during the period. These limit the findings in understanding the future trend of movie industry based on recent data.

## 8.2 Recommendations and Future Study Area

Despite the findings of the study, the impact of COVID-19 pandemic on the movie industry is huge, along with the accelerating growth of streaming platform. According to some studies [9], it suggests that the change of public viewing habits during the pandemic period is not temporally and will influence the future of movie industry in long-term. Under the context of the validated hypotheses, to study how the movies industry has be reformed after the pandemic is valuable for future research, and the exploration will definitely offer other valuable insights of the evolving movie industry.

# References

[1] C. the Motion Picture Association, "2014 theatrical statistics summary." [Online]. Available: https://www.motionpictures.org/wp-content/uploads/2015/03/MPAA-Theatrical-Market-Statistics-2014.pdf

[2] ——, "Theme report," 2019. [Online]. Available: https://www.motionpictures.org/wp-content/uploads/2020/03/MPA-THEME-2019.pdf

[3] K. McSpadden, "You now have a shorter attention span than a goldfish," Time, 05 2015. [Online]. Available: https://time.com/3858309/attention-spans-goldfish/

[4] D. GRIJALVA, "Movie industry," www.kaggle.com, 07 2021. [Online]. Available: https://www.kaggle.com/datasets/danielgrijalvas/movies

[5] S. BANSAL, "Netflix movies and tv shows," www.kaggle.com. [Online]. Available: https://www.kaggle.com/datasets/shivamb/netflix-shows

[6] J. Renda, "Trends: Despite shorter attention spans, the most popular movies are getting longer," The Dartmouth, 10 2023. [Online]. Available: https://www.thedartmouth.com/article/2023/10/trends-movies-are-longer

[7] A. B. Gelman and J. Hill, *Data analysis using regression and multilevel/hierarchical models.* Cambridge University Press, 2009.

[8] J. E. Cutting, "The evolution of pace in popular movies," *Cognitive Research: Principles and Implications*, vol. 1, 12 2016.

[9] A. Lai, "Post-pandemic media consumption: Online streaming accelerates a new content experience," Forrester, 06 2021. [Online]. Available: https://www.forrester.com/blogs/post-pandemic-media-consumption-online-streaming-accelerates-a-new-content-experience/?utm_source=forbes&utm_medium=pr&utm_campaign=b2cm