

# University School of Information, Communication and Technology



## MASTER OF COMPUTER APPLICATION (Software Engineering) (2024 - 2026)

Mid Evaluation: NUES (IT 666)

### Predicting Movie Sequel Success : A Data Analysis Project (ML Model Training)

**Mentor:**

Dr. Arshi Hussain

**Mentee:**

**Name:** Ananya Aggarwal

**Enrollmentno:**

03816404524

# Contents

<b>Acknowledgement</b>	<b>3</b>
<b>1 History and Background</b>	<b>4</b>
1.1 History of Movie Sequels . . . . .	4
1.2 The Role of Data Analysis in the Film Industry . . . . .	4
1.3 Machine Learning in Predicting Sequel Success . . . . .	5
1.4 Modern Applications and Future Trends . . . . .	5
<b>2 Introduction</b>	<b>5</b>
<b>3 Purpose of the Study</b>	<b>6</b>
<b>4 Objectives</b>	<b>6</b>
<b>5 Scope of the Study</b>	<b>7</b>
<b>6 Methodology</b>	<b>7</b>
6.1 Data Collection . . . . .	7
6.2 Data Preprocessing . . . . .	7
6.3 Exploratory Data Analysis (EDA) . . . . .	9
<b>7 Findings and Insights</b>	<b>9</b>
7.1 Data Collection and Preprocessing . . . . .	9
7.2 Feature Engineering . . . . .	9
7.3 Exploratory Data Analysis (EDA) . . . . .	10
<b>8 Learnings and Challenges</b>	<b>10</b>
8.1 Learnings . . . . .	10
8.2 Challenges . . . . .	11
<b>9 Outputs from EDA</b>	<b>11</b>
9.1 Visualizations . . . . .	11
9.2 Key Insights . . . . .	13

<b>10 Technologies and Tools</b>	<b>14</b>
<b>11 Expected Outcome</b>	<b>14</b>
<b>12 Conclusion</b>	<b>14</b>

# Acknowledgement

I would like to express my heartfelt gratitude to my project guide, **Dr. Arshi Hussain**, for their constant support, guidance, and encouragement throughout this project. Their expertise and feedback have been invaluable in shaping the direction of my work.

I am also thankful to my institution, **University School of Information, Communication and Technology**, for providing me with the resources and platform to work on this project. Special thanks to my peers and friends who offered their insights and suggestions during the development phase.

Lastly, I am deeply grateful to my family for their unwavering support and motivation, which kept me focused and determined to complete this project successfully.

# 1 History and Background

The evolution of sequels, data analysis, and machine learning has shaped the modern entertainment and technology landscapes. This section explores how these elements have developed and intersected over time.

## 1.1 History of Movie Sequels

The concept of sequels dates back to the early 20th century when filmmakers realized the potential of expanding successful narratives. The first known sequel, *The Fall of a Nation* (1916), followed *The Birth of a Nation* (1915), marking an early attempt at franchise storytelling. However, during Hollywood's Golden Age (1930s-1950s), sequels were not a primary focus, and standalone films dominated the industry.

By the 1960s and 1970s, franchises like *James Bond* set the stage for long-running sequels. The success of *The Godfather Part II* (1974) proved that sequels could surpass their predecessors in quality and storytelling. The late 20th century saw blockbusters like *Star Wars: The Empire Strikes Back* (1980) and *Jurassic Park: The Lost World* (1997), cementing sequels as a major industry trend.

In the 21st century, the rise of cinematic universes, led by Marvel Studios' interconnected storytelling, revolutionized sequel production. With franchises like *The Fast Furious* and *Harry Potter*, studios increasingly relied on sequels as a financial safety net. However, not all sequels succeed, leading to a growing interest in predictive analytics to assess sequel potential. [6]

## 1.2 The Role of Data Analysis in the Film Industry

Data analysis has become a crucial tool for decision-making in the entertainment industry. Studios leverage vast datasets, including box office performance, audience demographics, and online engagement, to determine whether a sequel is viable. Streaming platforms like Netflix and Disney+ analyze viewing patterns to decide which content should receive a continuation.

Techniques such as sentiment analysis on social media, trend forecasting, and revenue prediction models have transformed how sequels are planned. For instance, analyzing

past box office trends helps studios estimate a sequel's potential revenue, while viewer feedback guides script modifications to enhance appeal. [7]

### 1.3 Machine Learning in Predicting Sequel Success

Machine learning (ML) has introduced a data-driven approach to predicting the success of sequels. By analyzing variables such as budget, cast, director, genre, and audience reception, ML models identify patterns that indicate potential success.

Popular ML techniques in the film industry include:

- **Regression Models:** Predict box office revenue based on historical data.
- **Classification Algorithms:** Determine whether a sequel is likely to be a hit or a flop.
- **Natural Language Processing (NLP):** Analyze audience reviews and social media discussions to gauge reception.
- **Recommendation Systems:** Streaming services use ML to suggest sequels and related content to viewers based on their preferences.

### 1.4 Modern Applications and Future Trends

With advancements in artificial intelligence and big data, sequel prediction is becoming increasingly accurate. Streaming platforms are now using reinforcement learning to optimize content recommendations. Studios are also experimenting with generative AI to develop storylines based on audience preferences.

As this project explores the application of machine learning in predicting sequel success, it builds upon these historical and technological foundations. Understanding the evolution of sequels, data analysis, and ML provides valuable insights into how the entertainment industry continues to innovate and adapt to audience demands.

## 2 Introduction

The entertainment industry is a dynamic and ever-evolving field, with movie sequels playing a significant role in driving revenue. While sequels are often seen as a safe bet for

studios, their success is far from guaranteed. Some sequels outperform their predecessors, while others fail to meet audience expectations, resulting in financial losses. [12]

This project aims to address the unpredictability of movie sequel success by developing a **machine learning model** that can predict the commercial viability of a sequel based on various factors such as the original movie's performance, cast, director, budget, and genre. By leveraging data from the **TMDB 5000 Movie Dataset** [10], this study seeks to uncover patterns and trends that contribute to a sequel's success or failure.

The project combines **data science**, **machine learning**, and **web development** to create a practical and interactive solution for predicting movie sequel success. The ultimate goal is to provide film studios and investors with a **data-driven tool** that can help them make informed decisions and minimize financial risks.

### 3 Purpose of the Study

The primary purpose of this study is to **analyze the factors** that influence the success of movie sequels and to develop a **predictive model** that can forecast the commercial performance of a sequel. By understanding the key drivers of sequel success, this project aims to provide actionable insights for film studios, producers, and investors.

Additionally, the study seeks to bridge the gap between **creativity** and **analytics** in the entertainment industry. While creativity is essential for producing engaging content, data-driven insights can help studios make informed decisions about which projects to greenlight. This project aims to combine these two aspects by using machine learning to predict sequel success based on historical data.

### 4 Objectives

The primary objectives of this project are:

- To analyze factors that influence the success of movie sequels.
- To preprocess and clean the TMDB 5000 Movie Dataset.
- To apply machine learning techniques to predict sequel success.
- To visualize trends and insights using data analysis tools.

- To develop a web-based interface for users to check predictions.

## 5 Scope of the Study

This project focuses on:

- Movie data from the TMDB 5000 dataset. [10]
- Features such as budget, revenue, genre, actors, directors, and previous movie performance.
- Machine learning models for prediction.
- Data visualization and deployment using web technologies.

The study is limited to the data available in the TMDB 5000 dataset, and the predictions are based on historical trends and patterns. While the model aims to provide accurate predictions, it is important to note that external factors such as market trends, audience preferences, and competition may also influence a sequel's success.

## 6 Methodology

The project is being carried out in the following steps:

### 6.1 Data Collection

I started by downloading the **TMDB 5000 Movie Dataset** [10] from Kaggle. This dataset contains information about 5000 movies, including details such as budget, revenue, genres, actors, directors, and more. The dataset serves as the foundation for the analysis and machine learning model training.

### 6.2 Data Preprocessing

Before analyzing the data, I had to preprocess and clean it to ensure consistency and accuracy. The preprocessing steps included [2]:



- **Removing NA Values:** I noticed that the dataset had missing values, especially in the **budget** and **revenue** columns. Since these features are critical for predicting sequel success, I decided to remove rows with missing values in these columns. This ensured that the dataset used for analysis was complete and reliable.
- **Encoding Categorical Variables:** Since machine learning models require numerical input, I encoded categorical variables such as **genres**, **cast**, and **directors** using **one-hot encoding**. This allowed me to include these features in the model.
- **Normalizing Numerical Data:** I normalized numerical features such as **budget** and **revenue** to ensure that they are on the same scale, which improves model performance. [8]
- **Feature Engineering:** I created new features based on existing data to provide additional insights. For example:

- **Director Impact:** I calculated the average revenue for each director using the formula:

$$\text{director\_avg\_revenue} = \frac{\sum \text{revenue of movies by the director}}{\text{number of movies by the director}}$$

This feature captures the historical performance of directors and helps predict the success of sequels directed by them.

- **Cast Impact:** Similarly, I calculated the average revenue for each cast member using the formula:

$$\text{cast\_avg\_revenue} = \frac{\sum \text{revenue of movies featuring the cast member}}{\text{number of movies featuring the cast member}}$$

This feature captures the historical performance of actors and helps predict the success of sequels featuring them.

- **Saving the Processed Dataset:** After preprocessing and feature engineering, I saved the final processed dataset as **final\_processed\_data.csv** for use in model training.

## 6.3 Exploratory Data Analysis (EDA)

EDA was a critical step in understanding the dataset and identifying patterns and trends. The following steps were involved in EDA [3]:

- **Identifying Key Trends:** I analyzed the distribution of key variables such as **budget**, **revenue**, and **genre** to understand their impact on sequel success.
- **Correlation Analysis:** I examined the relationships between different variables to identify potential predictors of sequel success. For example, I looked at the correlation between **revenue** and **budget**, as well as the impact of **cast** and **directors** on movie performance.
- **Data Visualization:** I used visualization tools such as **Matplotlib** and **Seaborn** to create charts and graphs that helped me understand the data. For example, I created **bar charts** to show the distribution of movie genres and **scatter plots** to reveal the relationship between budget and revenue. [9]

## 7 Findings and Insights

### 7.1 Data Collection and Preprocessing

The **TMDB 5000 Movie Dataset** [10] provided a comprehensive set of features, but it had several missing values, especially in the **revenue** and **budget** columns. Since these features are critical for predicting sequel success, I decided to remove rows with missing values in these columns. This ensured that the dataset used for analysis was complete and reliable.

I also encoded categorical variables such as **genres**, **cast**, and **directors** using **one-hot encoding** to make them suitable for machine learning models.

### 7.2 Feature Engineering

I performed feature engineering to create new features that capture the historical performance of directors and cast members [4]:

- **Director Impact:** I calculated the average revenue for each director and added it as a new feature (**director\_avg\_revenue**). This feature helps quantify the impact of a director's track record on the success of a sequel.
- **Cast Impact:** Similarly, I calculated the average revenue for each cast member and added it as a new feature (**cast\_avg\_revenue**). This feature helps quantify the impact of an actor's track record on the success of a sequel.

These new features were merged back into the dataset, resulting in a final processed dataset saved as **final\_processed\_data.csv**.

## 7.3 Exploratory Data Analysis (EDA)

- **Budget and Revenue Correlation:** I found a positive correlation between a movie's budget and its revenue. Higher-budget movies tend to generate higher revenue, but this relationship is not linear, as some low-budget movies also performed exceptionally well.
- **Genre Analysis:** Certain genres, such as **Action**, **Adventure**, and **Sci-Fi**, were more likely to have sequels compared to others like **Drama** or **Romance**. These genres also tended to have higher budgets and revenues.
- **Director and Cast Impact:** Movies directed by well-known directors or featuring popular actors had a higher likelihood of success. For example, I observed that movies directed by **Christopher Nolan** or featuring actors like **Leonardo DiCaprio** consistently performed well at the box office.
- **Original Movie Performance:** The success of the original movie (measured by revenue/budget ratio) was a strong indicator of sequel success. Sequels of highly successful original movies were more likely to perform well.

# 8 Learnings and Challenges

## 8.1 Learnings

- **Data Preprocessing is Critical:** I learned that data cleaning and preprocessing are essential steps in any data science project. Handling missing values, encoding

categorical variables, and normalizing data are crucial for building accurate machine learning models.

- **Exploratory Data Analysis (EDA) is Key to Understanding Data:** EDA helped me uncover hidden patterns and trends in the dataset. By visualizing data using tools like **Matplotlib** and **Seaborn**, I was able to identify correlations and distributions that guided feature selection and model building.
- **Understanding the Domain is Essential:** Working on this project gave me a deeper understanding of the **entertainment industry**, particularly the dynamics of movie sequels. I learned that while financial factors like budget and revenue are important, other elements like **storytelling, audience reception, and market trends** also play a significant role in a sequel's success. [1]

## 8.2 Challenges

- **Handling Missing Data:** One of the biggest challenges I faced was dealing with missing values in the dataset. I had to carefully decide whether to fill in these missing values or remove them, as they could significantly impact the model's performance.
- **Feature Engineering:** Creating new features based on existing data was challenging but rewarding. I had to think creatively to come up with features that could provide meaningful insights, such as the **success ratio** of the original movie.

# 9 Outputs from EDA

## 9.1 Visualizations

I created several visualizations to understand the relationships and trends in the dataset [5]:

- **Budget vs. Revenue Scatter Plot:**

The scatter plot (Figure 1) shows a positive correlation between a movie's budget and its revenue. Higher-budget movies tend to generate higher revenue, but this

relationship is not linear, as some low-budget movies also performed exceptionally well.

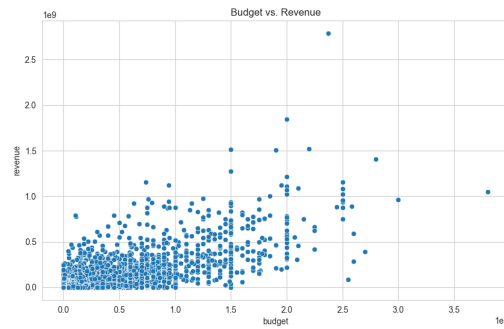


Figure 1: Budget vs. Revenue Scatter Plot

### Top Genres by Average Revenue:

The bar chart (Figure 2) shows the top 10 genres by average revenue. Genres like **Fantasy**, **Science Fiction**, and **Action** dominate the list, indicating that these genres tend to generate higher revenues compared to others.

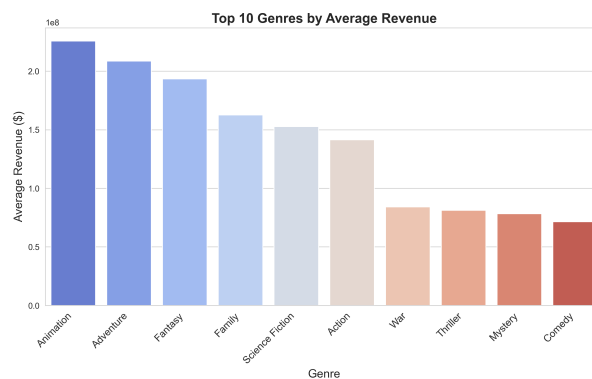


Figure 2: Top 10 Genres by Average Revenue

### Top Actors by Average Revenue:

The bar chart (Figure 3) shows the top 10 actors by average revenue. Actors like **Idina Menzel**, **Jonathan Groff**, and **Neel Sethi** are among the top performers, indicating that movies featuring these actors tend to generate higher revenues.

**Top Directors by Revenue:** The bar chart (Figure 4) shows the top 10 directors by revenue. Directors like **Christopher Nolan**, **James Cameron**, and **Steven Spielberg** are among the top performers, indicating that movies directed by them tend to generate higher revenues.

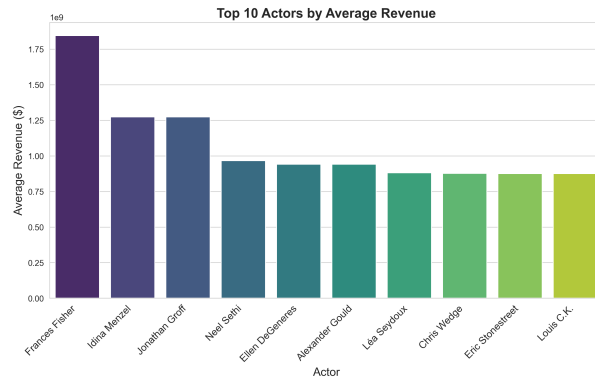


Figure 3: Top 10 Actors by Average Revenue

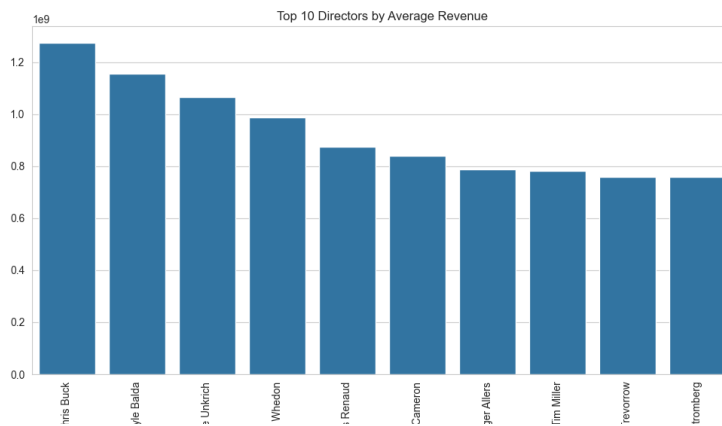


Figure 4: Top 10 Directors by Average Revenue

## 9.2 Key Insights

- **Budget and Revenue Correlation:** There is a positive correlation between a movie's budget and its revenue. However, some low-budget movies also perform exceptionally well, indicating that factors like storytelling and audience reception play a significant role.
- **Genre Impact:** Genres like **Fantasy**, **Science Fiction**, and **Action** tend to generate higher revenues, making them popular choices for sequels.
- **Actor and Director Impact:** The presence of popular actors and directors significantly influences a movie's success. For example, movies featuring actors like **Idina Menzel** or directed by **Christopher Nolan** consistently perform well at the box office.
- **Sequels are not guaranteed successes:** Even with high budgets, popular ac-

tors, and successful original movies, some sequels fail to meet expectations. This highlights the importance of other factors, such as **story quality, audience reception, and market competition**. [12]

## 10 Technologies and Tools

The following technologies and tools were used in this project:

- **Programming Language:** Python
- **Libraries:** Pandas, NumPy, Matplotlib, Seaborn [11]
- **Database:** SQL (MySQL/PostgreSQL)
- **Machine Learning Framework:** Scikit-learn
- **Web Framework:** Flask or Streamlit

## 11 Expected Outcome

By the end of this project, the following outcomes are expected:

- A trained machine learning model that predicts the success of a movie sequel.
- A web-based tool to input movie details and check predictions.
- Insights into factors affecting sequel success.

## 12 Conclusion

This project aims to provide valuable insights into the factors that influence the success of movie sequels. By combining **data science, machine learning, and web development**, the project offers a practical solution for predicting sequel performance. The mid-term evaluation focuses on **data collection, preprocessing, and exploratory data analysis**, which are critical steps in building an accurate and reliable predictive model.

## References

- [1] Cinema of the United States - Wikipedia — en.wikipedia.org. [https://en.wikipedia.org/wiki/Cinema\\_of\\_the\\_United\\_States](https://en.wikipedia.org/wiki/Cinema_of_the_United_States). [Accessed 24-03-2025].
- [2] Data Preprocessing - an overview — ScienceDirect Topics — sciencedirect.com. <https://www.sciencedirect.com/topics/engineering/data-preprocessing#:~:text=Data%20preprocessing%20is%20the%20concept,executing%20it%20to%20the%20algorithm>. [Accessed 24-03-2025].
- [3] Exploratory data analysis - Wikipedia — en.wikipedia.org. [https://en.wikipedia.org/wiki/Exploratory\\_data\\_analysis](https://en.wikipedia.org/wiki/Exploratory_data_analysis). [Accessed 24-03-2025].
- [4] Feature Engineering Explained — Built In — builtin.com. <https://builtin.com/articles/feature-engineering#:~:text=Feature%20Engineering%20Definition-,Feature%20engineering%20is%20the%20process%20of%20selecting%20C%20manipulating%20and%20transforming,exploratory%20data%20analysis%20and%20benchmarking>. [Accessed 24-03-2025].
- [5] Introduction to Matplotlib - GeeksforGeeks — geeksforgeeks.org. <https://www.geeksforgeeks.org/python-introduction-matplotlib/>. [Accessed 24-03-2025].
- [6] Sequel - Wikipedia — en.wikipedia.org. <https://en.wikipedia.org/wiki/Sequel>. [Accessed 24-03-2025].
- [7] The Role of Data Analytics in Cinema — cloudthat.com. <https://www.cloudthat.com/resources/blog/the-role-of-data-analytics-in-cinema#:~:text=Data%20analytics%20helps%20scriptwriters%20compare,blockbusters%20C%20identifying%20winning%20plot%20points.&text=Predictive%20analytics%20guide%20producers%20in,and%20choosing%20optimal%20filming%20locations>. [Accessed 24-03-2025].
- [8] What is Data Preprocessing? Key Steps and Techniques — techtarget.com. <https://www.techtarget.com/searchdatamanagement/definition/data-preprocessing>. [Accessed 24-03-2025].



- [9] What is Exploratory Data Analysis? — IBM — ibm.com. <https://www.ibm.com/think/topics/exploratory-data-analysis>. [Accessed 24-03-2025].
- [10] Tmdb 5000 movie dataset, 2024. Available at: <https://www.kaggle.com/datasets/tmdb/tmdb-movie-dataset>.
- [11] eBooks.com. Python for Data Analysis (3rd ed.) — ebooks.com. [https://www.ebooks.com/en-in/book/210644288/python-for-data-analysis/wes-mckinney/?affId=WES398681F&\\_c=1](https://www.ebooks.com/en-in/book/210644288/python-for-data-analysis/wes-mckinney/?affId=WES398681F&_c=1). [Accessed 24-03-2025].
- [12] Simone Murray. *The Sequel: A History of the Hollywood Remake*. Palgrave Macmillan, 2015.