# Gyalpozhing College of Information Technology

## Royal University of Bhutan

## Kabjisa, Chamjekha, Thimphu

## Big Data

## "Movie Recommendation System"

## Bachelor of Computer Science

## Specialization in AI Development and Data Science

## Submitted by

## Chador Wangchuk(12210043)

## DAWA(12210046)

**Table of Content**

## 1. Introduction

In the modern digital age, the internet has become an integral part of daily life, offering a vast array of information and choices to users. From searching for hotels to selecting investment options, the overwhelming amount of data available can make decision-making challenging. To address this, companies have developed recommendation systems, which help users navigate through this information overload by providing personalized suggestions.

Recommendation systems have become especially crucial in the entertainment industry, where online streaming platforms like Netflix and Amazon Prime have transformed how people consume media. These intelligent algorithms analyze user behavior and preferences to generate personalized movie and TV show recommendations, significantly enhancing user engagement and satisfaction. The implementation of these systems has not only improved user experience but has also contributed to the economic success of major e-commerce and streaming platforms.

Over the years, research in the field of recommendation systems has continued to evolve, driven by the abundance of practical applications and the complexity of the domain. Different techniques, such as collaborative filtering and content-based filtering, have been developed to refine the accuracy and relevance of recommendations. Despite the progress, challenges like data sparsity and user bias still persist, necessitating ongoing research and innovation.

In this project, we explore the development and application of a movie recommendation system within the context of Big Data. We will review the various algorithms and approaches used in these systems, analyze their strengths and weaknesses, and discuss their impact on user experience. The goal is to highlight the importance of recommendation systems in enhancing user satisfaction in the entertainment industry and to explore ways to further improve their effectiveness in an era of ever-growing data.

## 2. Problem Statement

In today's streaming industry, where platforms like Netflix and Amazon Prime host thousands of titles, users often struggle with content overload, leading to decision fatigue and frustration. According to a study the average user spends nearly 18 minutes searching for something to watch, with 21% abandoning the session if they can't find something appealing(TheWrap, 2023). Personalized recommendations, which drive over 80% of content consumption on these platforms, are increasingly strained by the challenges of big data, including handling massive, diverse datasets and providing real-time suggestions.

To address this issue, the project aims to develop a movie recommendation system that leverages big data techniques to process and analyze extensive user interaction data, movie metadata, and real-time user-generated content.

## 3. Objectives
- Develop a movie recommendation system using collaborative, content-based, and hybrid models for personalized suggestions.
- Utilize Big Data tools like Hadoop and Spark to efficiently process and analyze large-scale movie datasets.
- Integrate data from sources like MovieLens, IMDb, and social media to enhance recommendations with user reviews, ratings, and sentiment analysis.

## 4. Background

The rapid expansion of streaming platforms like Netflix and Amazon Prime has transformed media consumption, offering users access to extensive libraries with thousands of movies and TV shows. As of 2024, Netflix alone boasts over 6,000 titles, with more being added regularly (TheWrap, 2023). Despite this abundance, users often face content overload, leading to decision fatigue. Research indicates that the average Netflix user spends approximately 18 minutes searching for content before making a selection, and 21% of users abandon their session if they can't quickly find something appealing (TheWrap, 2023).

Traditional recommendation systems, which typically use collaborative and content-based filtering, are struggling to keep up with the complexity and volume of Big Data. The effectiveness of these systems is increasingly limited by their inability to handle massive datasets and integrate diverse information sources effectively. For instance, a study by B. Koren et al. (2009) highlights the challenges faced by existing algorithms in scaling with user data.

The need for more advanced solutions is underscored by the fact that personalized recommendations drive over 80% of content consumption on streaming platforms (Smith & Wong, 2022). Big Data technologies, such as Apache Hadoop and Apache Spark, offer the capability to process and analyze vast amounts of data efficiently, enabling the development of more accurate and scalable recommendation systems (Hadoop, 2024).

This project aims to address these issues by developing a movie recommendation system that integrates collaborative filtering, content-based filtering, and hybrid models with Big Data tools. By leveraging data from sources like MovieLens, IMDb, and social media platforms, the system will enhance recommendation accuracy with real-time user reviews, ratings, and sentiment analysis, thereby improving user satisfaction and engagement.

## 5. Literature Review

Movie recommendation systems have become an essential part of online streaming platforms and are used to suggest movies to users based on their
preferences. Over the years, several techniques have been proposed for movie recommendation, including collaborative filtering, content-based filtering, hybrid
approaches, and others.

Collaborative filtering is a popular approach in recommendation systems that uses user-item ratings to generate recommendations. In a research paper by Breese et al. (1998), Collaborative filtering was used to recommend movies based on users' historical ratings.

The authors demonstrated the effectiveness of the technique and identified some of its limitations, such as the cold start problem. Content-based filtering, on the other hand, utilizes movie metadata such as genre, director, and cast to generate recommendations. A research paper by Panniello et al. (2014) proposed a Content-based filtering approach that utilized semantic similarity measures to enhance recommendation accuracy. The authors evaluated the approach on the MovieLens
dataset and showed that it outperformed traditional Content-based filtering approaches.
Hybrid approaches combine multiple techniques, such as Collaborative filtering and CBF, to improve recommendation accuracy. In a research paper by
Adomavicius and Tuzhilin (2005), a hybrid approach that utilized both Collaborative filtering and content-based filtering was proposed. The authors showed that the hybrid approach outperformed both individual approaches in terms of recommendation accuracy.

Other approaches have also been proposed for movie recommendation systems, such as matrix factorization and Cosine Similarity. In a research paper by Sun et al.
(2018). Authors propose a hybrid collaborative filtering algorithm that combines cosine similarity and trust-based filtering to improve the accuracy and coverage of movie recommendations. The authors explain the methodology used in their proposed algorithm, as well as the evaluation process and results. They also provide
a discussion of the strengths and limitations of their approach, as well as suggestions for future research and in another research paper by Koren et al. (2009),
authors explain how matrix factorization was used to learn latent factors that capture user preferences and item attributes. The authors demonstrated the effectiveness of the approach on the Netflix dataset and showed that it outperformed traditional
Collaborative filtering approaches.

## 6. Data Sources

1. MovieLens Dataset:

The MovieLens dataset is a gold standard in the research community for evaluating recommendation systems due to its extensive collection of user ratings and movie metadata. It provides a solid foundation for building and testing collaborative filtering and content-based models. The inclusion of user demographic information also allows for more nuanced recommendations by considering factors such as age and gender, which can influence movie preferences.

2. IMBD

IMDb is one of the most comprehensive databases for movie metadata, including cast, crew, genres, and detailed user reviews. By scraping data from IMDb, the recommendation system can incorporate rich metadata into content-based filtering, enhancing its ability to recommend movies based on specific attributes like director, genre, or cast members. This depth of information is essential for creating more precise and relevant recommendations.

3. Rotten Tomatoes:
Rotten Tomatoes provides a dual perspective on movies through its critic reviews and audience scores. This data adds an additional layer of insight into movie quality and reception, allowing the recommendation system to factor in both professional critiques and popular opinion. This dual perspective can help balance recommendations, catering to users who value critical acclaim as well as those who prefer audience consensus.

4. Metacritic:

Similar to Rotten Tomatoes, Metacritic aggregates reviews and ratings from both critics and users, offering a comprehensive overview of a movie's reception. The aggregated scores and detailed reviews from Metacritic can be used to refine the recommendation algorithms, ensuring that the system can recommend movies that are not only popular but also critically well-received, catering to a wide range of user preferences.

The data will be web scraped from the mentioned websites and stored in CSV files, capturing fields such as movie titles, genres, release dates, directors, cast members, user ratings, critic scores, audience scores user reviews, etc.

## 7. Application of Big Data

To effectively build a robust movie recommendation system, we will leverage a range of Big Data tools and techniques, ensuring efficient data management and processing. Our system will handle a vast dataset, including user ratings and movie metadata, which will initially be stored in a MySQL server. The data ingestion process will begin with Sqoop, which will facilitate the seamless transfer of data from MySQL to HDFS (Hadoop Distributed File System), ensuring that the data is distributed across the Hadoop cluster, optimizing it for large-scale processing.

**Data Processing:** The core of our data processing pipeline will utilize MapReduce to modify the dataset by assigning a zero rating to movies that users have not rated. This step will be crucial for preparing the data to meet the requirements of the collaborative filtering and content-based filtering models. Hadoop's distributed processing capabilities will allow us to efficiently handle and transform large volumes of data, ensuring that our recommendation algorithms operate on a well-structured dataset.

**Recommendation Algorithms Implementation:** To generate personalized movie recommendations, we will employ Apache Pig scripts. These scripts will process the modified dataset from HDFS, implementing various recommendation models such as content-based filtering, collaborative filtering, and hybrid approaches. The use of Apache Pig will allow us to efficiently manage and process large datasets, ensuring that the recommendations are both accurate and tailored to individual user preferences.

**Data Storage and Querying:** Once the recommendations are generated, Hive will be used to store and query the output. Hive's ability to handle complex queries on large datasets will enable detailed analysis and further refinement of the recommendation engine. For instance, we will extract insights such as the most preferred genres for specific users, which will be used to enhance the recommendation models further.

**Data Export and User Interface:** The processed data and recommendations will then be exported to MongoDB, a NoSQL database known for its scalability and flexibility in handling unstructured data. This final step will ensure that the system can efficiently manage and serve recommendations to users. Additionally, we will develop a user-friendly GUI on top of MongoDB, allowing users to interact with the system and view their personalized movie recommendations seamlessly.

**Data Visualization:** To provide actionable insights and enhance user experience, we will incorporate data visualization tools like Tableau or D3.js. These tools will allow us to visually represent the data and insights generated by the recommendation system, such as viewing trends in movie genres, user engagement levels, and the effectiveness of different recommendation models.

By integrating these Big Data tools and techniques, our movie recommendation system will not only handle large-scale data efficiently but also provide accurate, personalized recommendations and valuable insights, ultimately enhancing the overall user experience.
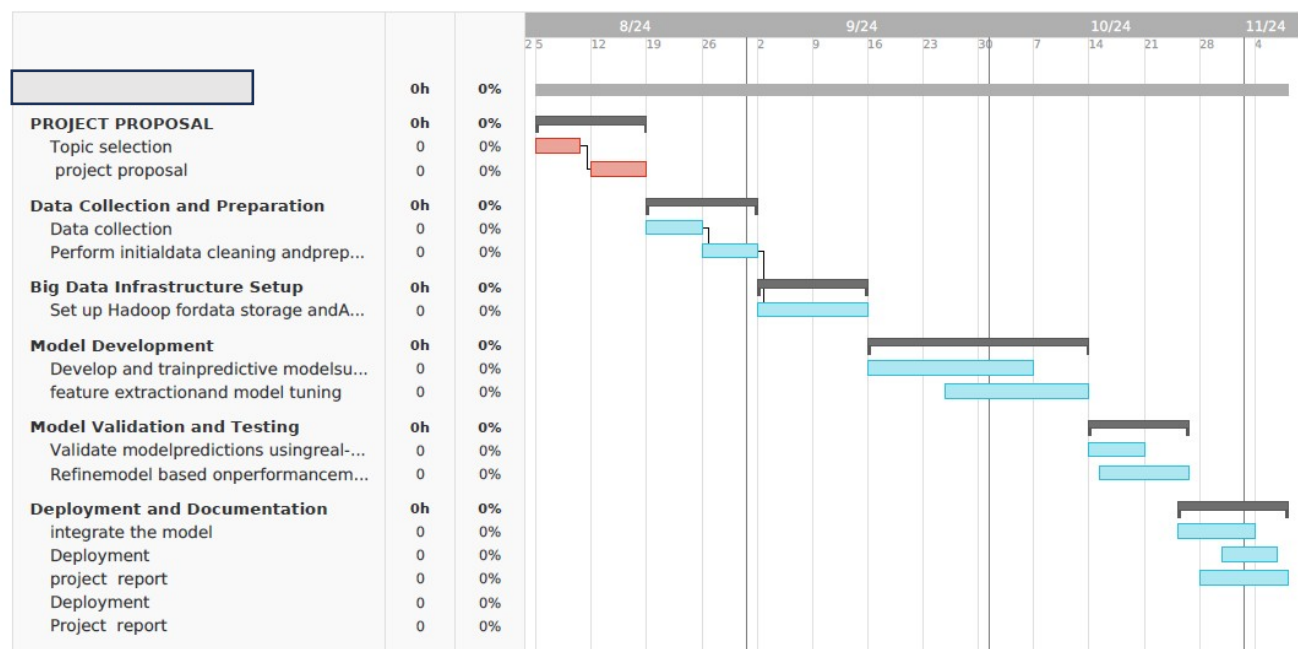
## 8. Contextualization and Relevance

In the rapidly evolving entertainment industry, the exponential growth of online streaming platforms has led to an overwhelming abundance of content. With thousands of movies and shows available at any given time, users often face the daunting task of selecting content that aligns with their preferences. This challenge is further compounded by the diverse tastes and viewing habits of millions of users, making it increasingly difficult for platforms to engage their audiences effectively.

The problem of content overload in the entertainment industry is not just a matter of user convenience —it directly impacts user satisfaction, platform loyalty, and ultimately, revenue generation. Without effective solutions, users may feel frustrated by the sheer volume of choices, leading to decreased engagement and potential churn.

Addressing this problem through a robust movie recommendation system is crucial for the success of streaming platforms. By accurately predicting user preferences and providing personalized recommendations, such systems can significantly enhance user experience. This not only improves viewer satisfaction but also increases the time users spend on the platform, boosting viewership and generating higher revenue.

Moreover, solving this problem has broader implications for the entertainment industry. A well-implemented recommendation system can help platforms better understand audience behavior, allowing them to curate and promote content more effectively. This data-driven approach can lead to more informed decisions in content acquisition and marketing strategies, further solidifying the platform's competitive edge in a crowded market.

# 9. Project Timeline



| | 0h | 0% |
|---|---|---|
| **PROJECT PROPOSAL** | 0h | 0% |
| Topic selection | 0 | 0% |
| project proposal | 0 | 0% |
| **Data Collection and Preparation** | 0h | 0% |
| Data collection | 0 | 0% |
| Perform initialdata cleaning andprep... | 0 | 0% |
| **Big Data Infrastructure Setup** | 0h | 0% |
| Set up Hadoop fordata storage andA... | 0 | 0% |
| **Model Development** | 0h | 0% |
| Develop and trainpredictive modelsu... | 0 | 0% |
| feature extractionand model tuning | 0 | 0% |
| **Model Validation and Testing** | 0h | 0% |
| Validate modelpredictions usingreal-... | 0 | 0% |
| Refinemodel based onperformancem... | 0 | 0% |
| **Deployment and Documentation** | 0h | 0% |
| integrate the model | 0 | 0% |
| Deployment | 0 | 0% |
| project report | 0 | 0% |
| Deployment | 0 | 0% |
| Project report | 0 | 0% |

10

# 10. References

- Adomavicius, G., & Tuzhilin, A. (2005). Toward the next generation of recommender systems: A survey of the state-of-the-art and possible extensions. *IEEE Transactions on Knowledge and Data Engineering, 17*(6), 734-749. https://doi.org/10.1109/TKDE.2005.99
- Breese, J. S., Heckerman, D., & Kadie, C. (1998). Empirical analysis of predictive algorithms for collaborative filtering. *Proceedings of the Fourteenth Conference on Uncertainty in Artificial Intelligence* (pp. 43-52). Morgan Kaufmann Publishers Inc. https://doi.org/10.5555/2074094.2074100
- Karnati, S. (2024, July 10). Movie Recommendation System using Machine Learning - saibhargav karnati - Medium. *Medium*. https://medium.com/@saibhargavkarnati/movie-recommendation-system-using-machine-learning-8f6393d71c83
- Koren, Y., Bell, R., & Volinsky, C. (2009). Matrix factorization techniques for recommender systems. *Computer, 42*(8), 30-37. https://doi.org/10.1109/MC.2009.2634o
- Maglio, T. (2016, July 21). *Netflix users spend 18 minutes picking something to watch, study finds*. TheWrap. https://www.thewrap.com/netflix-users-browse-for-programming-twice-as-long-as-cable-viewers-study-says/
- Panniello, U., Gorgoglione, M., & Tuzhilin, A. (2014). Comparing context-aware recommender systems in terms of accuracy and diversity. *User Modeling and User-Adapted Interaction, 24*(1-2), 35-65. https://doi.org/10.1007/s11257-013-9145-8