

**PREFERENCE LEARNING BASED APPROACH TO ADDRESS COLD START  
ISSUE IN RECOMMENDATION SYSTEM**

A Main project thesis submitted in partial fulfillment of requirements for the award of  
degree for VIII semester

**BACHELOR OF TECHNOLOGY  
IN  
COMPUTER SCIENCE AND ENGINEERING**

by

G. SAI NIKHIL (20131A0570)

G. HIMANSHU (20131A0573)

G. SANJAY (20131A0578)

K. NAGA RAJU (21135A0510)

Under the esteemed guidance of

**Dr. P. ARAVIND,**

Associate Professor,

Department of Computer Science and Engineering



**GAYATRI VIDYA PARISHAD COLLEGE OF ENGINEERING  
(AUTONOMOUS)**

Approved by AICTE & Affiliated to Andhra University, Visakhapatnam from 2022-23

(Affiliated to JNTU-K, Kakinada upto 2021-22)

**VISAKHAPATNAM**

**2023-2024**



## **GAYATRI VIDYA PARISHAD COLLEGE OF ENGINEERING**

**(Autonomous)**  
Approved by AICTE & Affiliated to Andhra University, Visakhapatnam from 2022-23  
(Affiliated to JNTUK, Kakinada upto 2021-22)  
Accredited twice by NAAC with 'A' Grade with a CGPA of 3.47/4.00  
Madhurawada, Visakhapatnam - 530048



COLLEGE OF ENGINEERING  
(AUTONOMOUS)

### **CERTIFICATE**

This is to certify that the main project entitled “Preference learning based approach to address cold start issue in Recommendation system” being submitted by

G. SAI NIKHIL (20131A0570)

G. HIMANSHU (20131A0573)

G. SANJAY (20131A0578)

K. NAGA RAJU (21135A0510)

in partial fulfilment for the award of the degree “Bachelor of Technology” in Computer Science and Engineering to the Jawaharlal Nehru Technological University, Kakinada is a record of bonafide work done under my guidance and supervision during VIII semester of the academic year 2023-2024.

The results embodied in this record have not been submitted to any other university or institution for the award of any Degree or Diploma.

#### **Guide**

Dr. P. Aravind  
Associate Professor  
Department of CSE  
GVPCE (A)

#### **Head of the Department**

Dr. D. Uma Devi  
Associate Professor &  
Associate Head of CSE (AI & ML) &  
I/c Head of CSE  
Department of CSE  
GVPCE (A)

#### **External Signature**

## DECLARATION

We hereby declare that this project entitled “**PREFERENCE LEARNING BASED APPROACH TO ADDRESS COLD START ISSUE IN RECOMMENDATION SYSTEM**” is a bonafide work done by us and submitted to “**Department of Computer Science and Engineering, G.V.P College of Engineering (Autonomous) Visakhapatnam**”, in partial fulfilment for the award of the degree of B. Tech and it is not submitted to any other university or has not been published anytime before.

PLACE: VISAKHAPATNAM

G.SAI NIKHIL (20131A0570)

DATE:

G. HIMANSHU (20131A0573)

G. SANJAY (20131A0578)

K. NAGA RAJU (21135A0510)

## ACKNOWLEDGEMENT

We would like to express our deep sense of gratitude to our esteemed institute **Gayatri Vidya Parishad College of Engineering (Autonomous)**, which has provided us an opportunity to fulfill our cherished desire.

We express our sincere thanks to our principal **Dr. A. B. Koteswara Rao, Gayatri Vidya Parishad College of Engineering (Autonomous)** for his encouragement to us during this project, giving us a chance to explore and learn new technologies in the form of projects.

We express our deep sense of gratitude to **Dr. D. Uma Devi, Associate Professor, Associate Head of B. Tech CSE with AI & ML and Incharge Head of the Department of Computer Science and Engineering, Gayatri Vidya Parishad College of Engineering (Autonomous)** for giving us an opportunity to do the project in college.

We express our profound gratitude and our deep indebtedness to our guide **Dr. P. Aravind**, whose valuable suggestions, guidance and comprehensive assessments helped us a lot in realizing our project.

We also thank our coordinator, **Dr. Ch. Sita Kumari**, Associate Professor, and our Assistant Coordinators, **Ms. P. Pooja Ratnam** and **Ms. K. Swathi**, Assistant Professors, Department of Computer Science and Engineering, for the kind suggestions and guidance for the successful completion of our project work.

G. SAI NIKHIL (20131A0570)

G. HIMANSHU (20131A0573)

G. SANJAY (20131A0578)

K. NAGA RAJU (21135A0510)

## ABSTRACT

In the realm of recommendation systems, the "cold start issue" presents a formidable challenge, particularly in the face of diverse user preferences and sparse data. Our research endeavors to tackle this challenge head-on with a preference learning-based approach, aimed at enhancing the effectiveness of recommendation systems. Focusing on neighborhood-based collaborative filtering methods, we employ sophisticated techniques such as Bhattacharyya similarity to glean insights from limited user-item interactions, thus mitigating the impact of sparse data on recommendation accuracy.

Through meticulous experimentation and analysis, we demonstrate the efficiency of our methodology in overcoming the cold start issue, thereby elevating the quality of recommendations provided to users. By amalgamating traditional wisdom with cutting-edge techniques, our approach sets a new standard for personalized recommendation systems. This research represents a significant step forward in the quest to make recommendation systems more robust, adaptive, and ultimately more valuable to users in various domains.

**Key Words:** Bhattacharyya Similarity, Collaborative Filtering, User-item Matrix, Data Preprocessing, Data Analysis, User Preferences, Data Visualization.

# INDEX

<b>CHAPTER 1. INTRODUCTION.....</b>	<b>8</b>
1.1 Objective.....	8
1.2 About the Algorithm.....	9
1.3 Purpose.....	11
1.4 Scope of the project.....	11
<b>CHAPTER 2. SRS DOCUMENT.....</b>	<b>13</b>
2.1 Functional Requirements.....	13
2.2 Non-functional Requirements.....	13
2.3 Minimum Hardware Requirements.....	14
2.4 Minimum Software Requirements.....	14
<b>CHAPTER 3. ALGORITHM ANALYSIS.....</b>	<b>15</b>
3.1 Existing Algorithm.....	15
3.2 Proposed Algorithm.....	16
3.3 Feasibility Study.....	17
3.4 Cost Benefit Analysis.....	19
<b>CHAPTER 4. SOFTWARE DESCRIPTION .....</b>	<b>21</b>
4.1 Google Colab.....	21
4.2 Python .....	21
4.3 Scikit-learn .....	21
4.4 NumPy .....	22
4.5 Matplotlib.....	22
4.6 Streamlit .....	22

<b>CHAPTER 5. PROJECT DESCRIPTION .....</b>	<b>23</b>
5.1 Problem Definition.....	23
5.2 Project Overview.....	23
5.3 Module Description.....	24
5.3.1 Python and Streamlit.....	24
5.3.2 Model .....	24
<b>CHAPTER 6. SYSTEM DESIGN.....</b>	<b>25</b>
6.1 Introduction to UML.....	25
6.2 Building Blocks of the ML.....	27
6.3 UML Diagrams.....	32
<b>CHAPTER 7. IMPLEMENTATION.....</b>	<b>35</b>
7.1 Datasets used.....	35
7.2 Sample Code.....	36
7.3 Results.....	44
<b>CHAPTER 8. TESTING.....</b>	<b>46</b>
8.1 Introduction of Testing.....	46
<b>CHAPTER 9. CONCLUSION .....</b>	<b>50</b>
FUTURE SCOPE.....	51
REFERENCE LINKS.....	53

# 1. INTRODUCTION

The explosion of data in recommendation systems demands ever-more accurate methods for predicting user preferences. Collaborative filtering (CF) has been a cornerstone of recommendation technology, but it faces challenges with sparse datasets where users have rated few items in common. To address this, neighborhood-based CF focuses on finding users with similar taste profiles based on their past interactions.

This project explores the use of Bhattacharya similarity, a measure that captures the shared characteristics of rating patterns between users, even for items they haven't both rated. Unlike traditional similarity metrics that only consider co-rated items, Bhattacharya similarity leverages all user ratings, leading to more robust recommendations in sparse data scenarios. By implementing neighborhood-based CF with Bhattacharya similarity, we aim to create a more accurate and effective recommendation system that can personalize suggestions even for new users or items with limited ratings.

## 1.1. OBJECTIVE

While traditional matrix factorization plays a significant role in recommendation systems, it often struggles with capturing intricate user preferences and hidden relationships within sparse datasets. This project proposes an alternative approach using neighborhood-based collaborative filtering (CF) with Bhattacharya similarity.

Bhattacharya similarity offers a distinct advantage by considering all user ratings, even for unshared items. This allows for a more nuanced understanding of user taste profiles, leading to more accurate recommendations, especially when dealing with limited data. Through this project, we aim to develop a robust and efficient recommendation system that leverages the strengths of neighborhood-based CF and Bhattacharya similarity to deliver personalized suggestions even in sparse data scenarios. Our ultimate goal is to contribute to the advancement of CF techniques with real-world applications, particularly in domains like e-commerce and personalized content recommendation.



## 1.2. ABOUT THE ALGORITHM

### 1.2.1 Neighborhood Based Collaborative Filtering (NCF)<sup>[2]</sup>

Neighborhood-based collaborative filtering (CF) is a popular technique in recommendation systems that leverages the power of similar users. Here's how it works:

1. **Data Collection:** The system gathers user-item interaction data, typically in the form of ratings or purchases.
2. **User Similarity Calculation:** The core concept lies in identifying users with similar taste profiles. Various similarity metrics can be used, like Pearson correlation or cosine similarity. These metrics consider co-rated items (items both users have interacted with) and quantify how closely their ratings align.
3. **Neighborhood Formation:** Based on the chosen similarity metric, a neighborhood is formed around a target user. This neighborhood consists of the most similar users (k-nearest neighbors) identified from the overall user base. The value of k (number of neighbors) is a hyperparameter that can be tuned for optimal performance.
4. **Recommendation Generation:** For a target user, the system wants to predict their preference for an item they haven't interacted with yet. Here's where the neighborhood comes in:
  - The system retrieves the ratings of the neighbors for the target item.
  - It then weights these neighbor ratings based on their similarity to the target user (more similar neighbors contribute more weight).
  - Finally, it aggregates the weighted neighbor ratings to predict the target user's potential rating for the item. This predicted rating serves as the recommendation score, with higher scores indicating a greater likelihood of the user enjoying the item.

#### Advantages of Neighborhood-based CF:

- **Interpretability:** The underlying logic is easy to understand, making it a good starting point for recommender systems.
- **Efficiency:** For smaller datasets, neighborhood-based CF can be computationally efficient, especially compared to more complex matrix factorization techniques.

- **Cold Start Problem Mitigation:** It can address the "cold start" problem to some extent, where new users have limited interaction data. By finding similar users with established preferences, the system can still generate recommendations for them.

#### **Disadvantages of Neighborhood-based CF:**

- **Scalability Challenges:** As the user base and item catalog grow, finding relevant neighbors and generating recommendations can become computationally expensive.

#### **1.2.2. Bhattacharya coefficient**

The Bhattacharya coefficient is a metric used to quantify the similarity between two probability distributions. Here's how it's calculated:

##### **For Discrete Distributions:**

Let P and Q represent the two probability distributions you want to compare.

Let  $p(i)$  and  $q(i)$  represent the probabilities of the  $i$ -th event or sample within their respective distributions.

The Bhattacharya Coefficient ( $BC(P, Q)$ ) is then calculated as the sum, across all possible events ( $i$ ), of the square root of the product of the corresponding probabilities in each distribution ( $p(i)$  and  $q(i)$ ).

Here's the mathematical formula:

$$BC(P, Q) = \sum(\sqrt{p(i) * q(i)})$$

##### **For Continuous Distributions:**

The concept remains similar, but instead of summing across events, you'd integrate across the entire range of the variable:

$$BC(P, Q) = \int(\sqrt{p(x) * q(x)}) dx$$

##### **Interpretation:**

The Bhattacharya coefficient ranges from 0 to 1.

A value of 1 indicates perfect similarity between the two distributions (they completely overlap).

A value of 0 indicates no overlap between the distributions (they are completely different). So, the higher the coefficient, the more similar the two distributions are.

### **1.3. PURPOSE**

#### **Understanding User Preferences Beyond Common Ground:**

Neighborhood-based collaborative filtering (CF) relies on user similarity to generate recommendations. Traditionally, similarity metrics like Pearson correlation focus only on co-rated items (items both users have interacted with). This can be limiting, especially with sparse datasets where users haven't rated many items in common.

#### **Bhattacharya Similarity: Embracing All the Data:**

Bhattacharya similarity steps in to address this limitation. It goes beyond co-rated items and considers all user ratings, even for unshared items. By incorporating the overall rating patterns, it paints a richer picture of user preferences.

#### **Benefits for Sparse Data Scenarios:**

In sparse datasets, traditional metrics struggle to find meaningful similarities due to limited shared data points. Bhattacharya similarity, by leveraging all ratings, provides a more nuanced understanding of user taste, even with limited co-rated items. This leads to more accurate recommendations, particularly when dealing with new users or items with few ratings.

### **1.4. SCOPE**

The scope of the project on Bhattacharyya similarity-based recommendation systems involves designing and implementing an innovative recommendation approach to address the cold start issue prevalent in traditional recommendation systems. This entails developing a preference learning-based model that leverages Bhattacharyya similarity to generate accurate recommendations in scenarios with sparse data and diverse user preferences. The project aims to contribute to advancements in recommendation systems by providing a more effective solution to the cold start problem.

1. **Data Collection:** Gathering relevant datasets containing user-item interactions, ratings, and item characteristics necessary for recommendation tasks. The dataset will be diverse and representative, ensuring it captures a wide range of user preferences and item attributes.
2. **Feature Engineering:** Preprocessing and encoding the dataset to extract meaningful features, including user IDs, item IDs, and rating information. Additionally, creating interaction features to capture relationships between users and items, enhancing the model's ability to make accurate recommendations.
3. **Model Development:** Designing and implementing a recommendation model based on Bhattacharyya similarity. This model will calculate the similarity between items using Bhattacharyya coefficient, effectively capturing the intrinsic relationships between items based on user preferences.
4. **Model Evaluation and Validation:** Assessing the performance of the developed recommendation model using appropriate metrics such as accuracy, precision, and recall. Validation will involve testing the model on independent datasets to ensure its robustness and effectiveness in making accurate recommendations.
5. **Deployment and Integration:** Deploying the trained recommendation model into a production-ready environment, such as a web server or cloud platform, to enable real-time recommendation services. Integration with existing recommendation systems or platforms will be conducted to seamlessly incorporate the Bhattacharyya similarity-based approach into existing infrastructures.

## 2. SRS DOCUMENT

A software requirements specification (SRS) is a document that describes what the software will do and how it will be expected to perform.

### 2.1 FUNCTIONAL REQUIREMENTS

A Functional Requirement (FR) is a description of the service that the model must offer. It describes a system or its component. A function is nothing but inputs to the system, its behavior, and outputs. It can be a calculation, data manipulation, business process, user interaction, or any other specific functionality which defines what function a system is likely to perform. Functional Requirements are also called Functional Specification.

- **Bhattacharyya Similarity Calculation:** The system must accurately compute Bhattacharyya similarity between items based on their user interaction data. This involves analyzing the distribution of user preferences for each item and determining the degree of similarity between items' preference distributions.
- **Personalized Recommendation Generation:** System should utilize Bhattacharyya similarity scores to generate personalized recommendations for users. By considering each user's unique preferences and behaviors, the system can tailor recommendations to individual users, thereby enhancing user satisfaction.
- **Scalability and Efficiency:** The system must be scalable and efficient to handle large-scale datasets and deliver timely recommendations, even as the user base grows. It should leverage optimized algorithms and distributed computing techniques to process data swiftly and maintain responsiveness, ensuring a seamless user experience.

### 2.2 NON-FUNCTIONAL REQUIREMENTS

Non-functional requirements specify the quality attribute of a software system. They judge the system based on Responsiveness, Usability, Security, Portability. Non-functional requirements are called qualities of a system, there are as follows:

- **Performance:** Ensure efficient processing of large-scale user-item data for timely recommendations.

- **Robustness:** Handle data noise and outliers to maintain system stability and reliability.
- **Accuracy:** Deliver highly accurate recommendations reflecting user preferences consistently.
- **Reliability:** Provide dependable performance under varying conditions, ensuring consistent service availability.

## **2.3 MINIMUM HARDWARE REQUIREMENTS**

- CPU: Multi-core processors (e.g., Intel Xeon, AMD EPYC)
- RAM: 16GB or more
- Storage: SSDs or HDDs with multiple terabytes

## **2.4 MINIMUM SOFTWARE REQUIREMENTS**

Python based Machine Learning and Deep Learning libraries will be needed for the development and experimentation of the project.

- Operating System: Linux or Windows
- Framework: NumPy, Pandas, Streamlit, Matplotlib
- Programming Language: Python
- Platforms: Jupyter Notebook or Google Colab

## 3. ANALYSIS

### 3.1. EXISTING ALGORITHM

Collaborative Filtering<sup>[3]</sup> is a recommendation technique that relies on user-item interactions to make predictions. It works by finding similarities between users or items based on their past interactions and uses these similarities to recommend items to users.

Content-based filtering is a recommendation technique that relies on the attributes or features of items and the user's preferences to make recommendations. Instead of relying on past user interactions like collaborative filtering, content-based filtering recommends items that are similar to those the user has liked in the past, based on the content or characteristics of the items themselves.

There are several similarity methods that are in use in the case of Neighborhood based collaborative approach to compute similarity in between users or in between items in which some are:

1. Cosine Similarity: Measures the cosine of the angle between two vectors, indicating their similarity based on the orientation, regardless of their magnitude.
2. Pearson Correlation Coefficient: Measures the linear correlation between two variables, indicating the strength and direction of their relationship.
3. Jaccard Similarity: Measures the similarity between two sets by comparing the intersection over the union of their elements, commonly used for binary data like user-item interactions.
4. Euclidean Distance: Measures the straight-line distance between two points in a multi-dimensional space, indicating their similarity based on proximity.

#### 3.1.1. DRAWBACKS OF EXISTING ALGORITHM

The existing collaborative filtering (CF)<sup>[4]</sup> system has several disadvantages:

- **Poor performance in sparse datasets:** Traditional similarity measures used in CF systems perform poorly in sparse rating datasets, leading to inaccurate recommendations.

- **Inability to handle the cold-start problem:** The system struggles to make accurate recommendations for new users or items with limited ratings data, known as the cold-start problem.
- **Limited utilization of non-co-rated items:** Some approaches only take into account ratings on a pair of non-co-rated items if the similarity between them is maximum, leading to a limited utilization of available data.
- **Domain-specific measures:** Some proposed measures for CF systems are highly domain-specific, limiting their applicability to broader contexts.
- **Inability to capture context information:** Traditional measures cannot capture context information in collaborative team environments, limiting their effectiveness in such scenarios.

These disadvantages highlight the need for improved similarity measures and collaborative filtering approaches to address these limitations.

## 3.2. PROPOSED ALGORITHM

We came up with the idea of using Bhattacharya Coefficient<sup>[5]</sup> which computes similarity between two provided probability distributions and gives the similarity to the user to know the similarity between the items he is computing for here Our proposed approach injects a new lens into neighborhood-based CF by incorporating Bhattacharya similarity. This metric goes beyond co-rated items and analyzes users' entire rating patterns. By capturing these nuanced details, Bhattacharya similarity helps us identify similar users with a deeper understanding of their tastes, even when they haven't rated many of the same items.

In datasets with limited shared ratings, traditional CF methods struggle to find meaningful similarities. Bhattacharya similarity, by considering all user interactions, provides a more robust measure of user taste, leading to more accurate recommendations, particularly for new users or items with few ratings. This allows us to personalize user experiences even with limited data.

### 3.2.1. ADVANTAGES OF PROPOSED MODEL

- Improved Accuracy in Sparse Data



- Enhanced Cold Start Handling
- Deeper User Understanding
- Interpretability
- Computational Efficiency

### 3.3 FEASIBILITY STUDY

A feasibility study is an analysis that takes all a project's relevant factors into account including economic, technical, legal, and scheduling considerations to ascertain the likelihood of completing the project successfully. A feasibility study is important and essential to evaluate any proposed project is feasible or not. A feasibility study is simply an assessment of the practicality of a proposed plan or project.

**The main objectives of feasibility are mentioned below:**

To determine if the product is technically and financially feasible to develop, is the main aim of the feasibility study activity. A feasibility study should provide management with enough information to decide:

- Whether the project can be done.
- To determine how successful your proposed action will be.
- Whether the final product will benefit its intended users.
- To describe the nature and complexity of the project.
- What are the alternatives among which a solution will be chosen (During subsequent phases)
- To analyze if the software meets organizational requirements. There are various types of feasibility that can be determined. They are:

**Operational** - Define the urgency of the problem and the acceptability of any solution, includes people-oriented and social issues: internal issues, such as manpower problems, labor objections, manager resistance, organizational conflicts, and policies; also, external issues, including social acceptability, legal aspects, and government regulations.

**Technical** - Is the feasibility within the limits of current technology? Does the technology exist at all? Is it available within a given resource?

**Economic** - Is the project possible, given resource constraints? Are the benefits that will accrue from the new system worth the costs? What are the savings that will result from the system, including tangible and intangible ones? What are the development and operational costs?

**Schedule** - Constraints on the project schedule and whether they could be reasonably met.

### **3.3.1. ECONOMIC FEASIBILITY:**

Economic analysis could also be referred to as cost/benefit analysis. It is the most frequently used method for evaluating the effectiveness of a new system. In economic analysis the procedure is to determine the benefits and savings that are expected from a candidate system and compare them with costs. Economic feasibility study related to price, and all kinds of expenditure related to the scheme before the project starts. This study also improves project reliability. It is also helpful for the decision-makers to decide the planned scheme processed latter or now, depending on the financial condition of the organization. This evaluation process also studies the price benefits of the proposed scheme. Economic Feasibility also performs the following tasks.

- Cost of packaged software.
- Cost of doing full system study.
- Is the system cost Effective?

### **3.3.2. TECHNICAL FEASIBILITY:**

A large part of determining resources has to do with assessing technical feasibility. It considers the technical requirements of the proposed project. The technical requirements are then compared to the technical capability of the organization. The systems project is considered technically feasible if the internal technical capability is sufficient to support the project requirements. The analyst must find out whether current technical resources can be where the expertise of system analysts is beneficial, since using their own experience and their contact with vendors they will be able to answer the question of technical feasibility. Technical Feasibility also performs the following tasks.

- Is the technology available within the given resource constraints?
- Is the technology have the capacity to handle the solution

- Determines whether the relevant technology is stable and established.
- Is the technology chosen for software development has a large number of users so that they can be consulted when problems arise, or improvements are required?

### **3.3.3. OPERATIONAL FEASIBILITY:**

Operational feasibility is a measure of how well a proposed system solves the problems and takes advantage of the opportunities identified during scope definition and how it satisfies the requirements identified in the requirements analysis phase of system development. The operational feasibility refers to the availability of the operational resources needed to extend research results beyond on which they were developed and for which all the operational requirements are minimal and easily accommodated. In addition, the operational feasibility would include any rational compromises farmers make in adjusting the technology to the limited operational resources available to them. The operational Feasibility also perform the tasks like

- Does the current mode of operation provide adequate response time?
- Does the current of operation make maximum use of resources?
- Determines whether the solution suggested by the software development team is acceptable.
- Does the operation offer an effective way to control the data?
- Our project operates with a processor and packages installed are supported by the system.

### **3.4. COST BENEFIT ANALYSIS**

The financial and the economic questions during the preliminary investigation are verified to estimate the following:

- The cost of the hardware and software for the class of application being considered.
- The benefits in the form of reduced cost.
- The proposed system will give the minute information, as a result.
- Performance is improved which in turn may be expected to provide increased profits.

- This feasibility checks whether the system can be developed with the available funds.
- This can be done economically if planned judiciously, so it is economically feasible.
- The cost of the project depends upon the number of man-hours required.

## **4. SOFTWARE DESCRIPTION**

### **4.1. Google Colab**

Google Colab, short for Colaboratory, is a cloud-based platform provided by Google that enables users to write, execute, and share Python code collaboratively. It integrates with Google Drive and offers free access to computational resources, including CPU, GPU, and TPU, making it ideal for running machine learning and deep learning experiments without the need for expensive hardware. Colab provides a Jupyter Notebook-like environment, allowing users to write code, add explanatory text, and visualize results in a single document. It supports popular Python libraries such as TensorFlow, PyTorch, and scikit-learn, along with built-in support for data visualization tools like Matplotlib and Seaborn. Colab also offers seamless integration with Google Drive, enabling users to import and export data and notebooks easily.

### **4.2. Python**

Python is an interpreted, object-oriented, high-level programming language with dynamic semantics. Its high-level built-in data structures, combined with dynamic typing and dynamic binding, make it very attractive for Rapid Application Development, as well as for use as a scripting or glue language to connect existing components together. Python's simple, easy to learn syntax emphasizes readability and therefore reduces the cost of program maintenance. Python supports modules and packages, which encourages program modularity and code reuse. The Python interpreter and the extensive standard library are available in source or binary form without charge for all major platforms, and can be freely distributed.

### **4.3. Scikit-learn**

Scikit-learn[5] is a free software machine learning library for the Python programming language. It features various classification, regression and clustering algorithms including support vector machines, random forests, gradient boosting, k-means and DBSCAN, and is designed to interoperate with the Python numerical and scientific libraries NumPy and SciPy.

#### **4.4. NumPy**

NumPy is a library for the Python programming language, adding support for large, multidimensional arrays and matrices, along with a large collection of high-level mathematical functions to operate on these arrays.

#### **4.5. Matplotlib**

Matplotlib is a popular open-source plotting library for Python. It provides a wide range of visualization options for data analysis, statistics, and machine learning applications. It provides an object-oriented API for embedding plots into applications using general purpose GUI toolkits.

#### **4.6. Streamlit**

Streamlit empowers data scientists and machine learning engineers to create user-friendly web apps directly from their Python code. Forget juggling HTML, CSS, and Javascript - Streamlit lets you write everything in Python, streamlining the development process. Perfect for rapid prototyping, you can build interactive data apps in minutes, ideal for showcasing your projects or validating ideas. Streamlit integrates seamlessly with popular data visualization libraries like matplotlib and seaborn, allowing you to present your data insights with clarity. The deployment options are flexible as well - share your apps locally for internal review or leverage Streamlit's free Community Cloud for wider accessibility. By handling the web development aspects, Streamlit empowers you to focus on the core functionality of your data app, bridging the gap between data science and user interaction.

## 5. PROBLEM DESCRIPTION

### 5.1. PROBLEM DEFINITION

Our project proposes a solution to these limitations by incorporating Bhattacharya similarity into the neighborhood-based CF approach. This metric offers a distinct advantage:

- **Considering All the Data:** Unlike traditional methods, Bhattacharya similarity goes beyond co-rated items. It analyzes users' entire rating patterns, including items they haven't interacted with in common. This allows us to capture a richer picture of user taste and identify subtle similarities that might be missed otherwise.
- **Enhanced Performance in Sparse Data:** By leveraging all user ratings, Bhattacharya similarity provides a more robust measure of user similarity, even with limited shared data points. This leads to more accurate recommendations, especially when dealing with new users or items with few ratings.

By incorporating Bhattacharya similarity, we aim to develop a recommendation system that offers a deeper understanding of user preferences and delivers more accurate and personalized suggestions, even in the face of data sparsity.

### 5.2. PROJECT OVERVIEW

Our project, titled "Enhancing Recommendation Systems with Bhattacharyya Similarity<sup>[5]</sup>," addresses the challenges posed by the cold start issue in recommendation systems. We propose a novel approach based on Bhattacharyya similarity to effectively tackle sparse data and diverse user preferences. By leveraging this measure, which quantifies the similarity between item profiles, we aim to improve recommendation accuracy, particularly in scenarios with limited user-item interactions.

The steps involved in the project are: -

1. Data Collection and Preprocessing.
2. Similarity Computation.
3. Recommendation Generation.
4. Evaluation and Optimization:

5. Deployment and Integration.
6. Monitoring and Maintenance.

Refer to the image 5.1 for the description of algorithm. Overall, we aim to demonstrate the efficacy of Bhattacharyya similarity as a powerful tool for improving recommendation accuracy and addressing the cold start issue in recommendation systems.

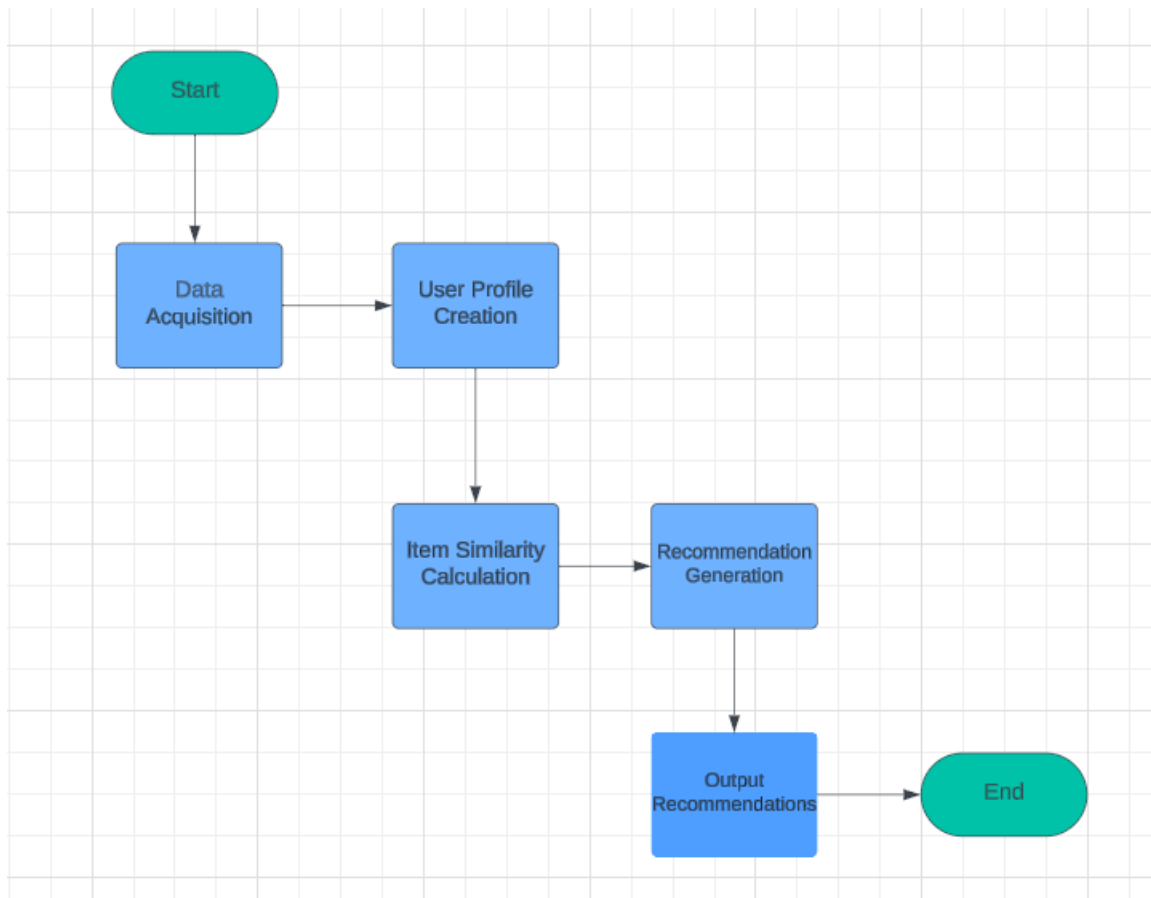
## **5.3. MODULE DESCRIPTION**

### **5.3.1. PYTHON AND STREAMLIT**

Streamlit is a Python library for creating interactive web applications with ease. It simplifies the process of building data-focused web apps by providing a simple syntax and a wide range of interactive widgets. With Streamlit, users can rapidly prototype and iterate on their ideas, integrating visualizations, machine learning models, and data analysis seamlessly. Additionally, Streamlit handles deployment, allowing users to share their applications effortlessly. Overall, Streamlit empowers users to communicate insights and deploy models more effectively through interactive web interfaces.

Python is an interpreted, object-oriented, high-level programming language with dynamic semantics. Its high-level built-in data structures, combined with dynamic typing and dynamic binding, make it very attractive for Rapid Application Development, as well as for use as a scripting or glue language to connect existing components together. Python's simple, easy to learn syntax emphasizes readability and therefore reduces the cost of program maintenance. Python supports modules and packages, which encourages program modularity and code reuse. The Python interpreter and the extensive standard library are available in source or binary form without charge for all major platforms, and can be freely distributed.





**Fig 5.1 - Flowchart describing the model**

## **6. SYSTEM DESIGN**

System design is the phase that bridges the gap between problem domain and the existing system in a manageable way. It is the process of defining elements of a system like modules, architecture, components and their interfaces and data for a system based on the specified requirements. This involves defining, developing and designing systems which satisfies the specific needs and requirements.

In the phase of system design, the SRS document is converted into a format that can be implemented and decides how the system will operate. A designer uses the modelling languages to express the information and knowledge in a structure of system that is defined by a consistent set of rules and definitions. The designs can be defined in graphical or textual modelling languages.

### **6.1 INTRODUCTION TO UML**

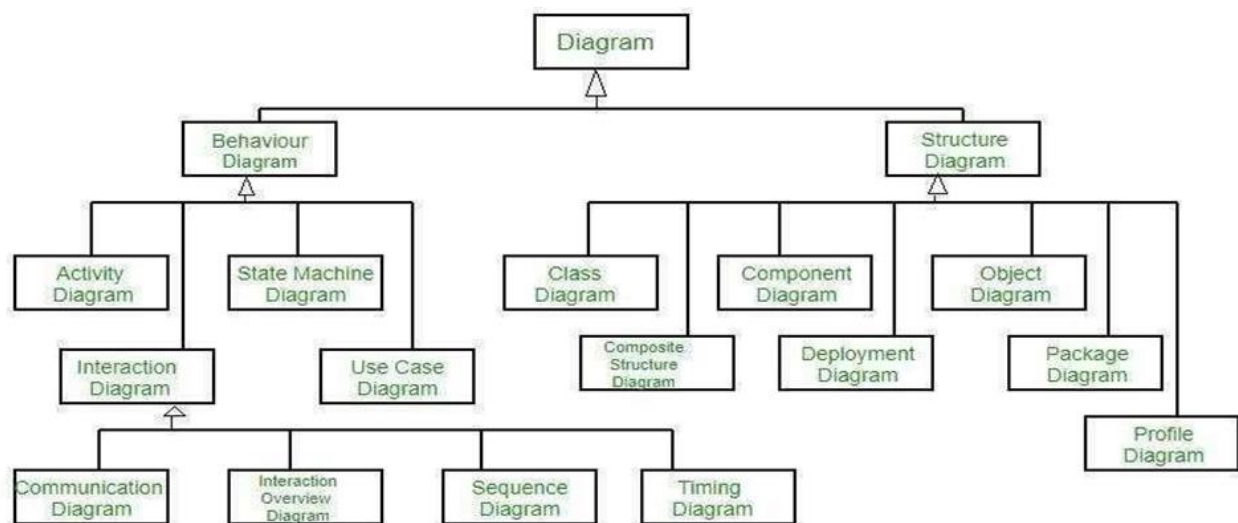
Unified Modeling Language (UML)<sup>[6]</sup> is a general-purpose modeling language. The main aim of UML is to define a standard way to visualize the way a system has been designed. It is quite like blueprints used in other fields of engineering. UML is not a programming language; it is rather a visual language. We use UML diagrams to portray the behavior and structure of a system. UML helps software engineers, businessmen and system architects with modeling, design and analysis. The Object Management Group (OMG) adopted Unified Modelling Language as a standard in 1997. It's been managed by OMG ever since. International Organization for Standardization (ISO) published UML as an approved standard in 2005. UML has been revised over the years and is reviewed periodically.

#### **Why we need UML**

1. Complex applications need collaboration and planning from multiple teams and hence require a clear and concise way to communicate amongst them.
2. Businessmen do not understand code. So, UML becomes essential to communicate with nonprogrammers' essential requirements, functionalities and processes.
3. A lot of time is saved down the line when teams can visualize processes, user interactions and static structure of the system. UML is linked with object-oriented

design and analysis. UML makes the use of elements and forms associations between them to form diagrams. Diagrams in UML can be broadly classified as:

- **Structural Diagrams** – Capture static aspects or structure of a system. Structural Diagrams include Component Diagrams, Object Diagrams, Class Diagrams and Deployment Diagrams.
- **Behavior Diagrams** – Capture dynamic aspects or behavior of the system. Behavior diagrams include Use Case Diagrams, State Diagrams, Activity Diagrams and Interaction Diagrams. Building Blocks of the UML Building Blocks of the UML



**Fig 6.1 Building Block in UML**

## 6.2 Building Block of the UML

The vocabulary of the UML encompasses three kinds of building blocks:

- Things
- Relationships
- Diagrams

Things are the abstractions that are first-class citizens in a model; diagrams group interesting collections of things. Refer to the image 5.1 for the sample block diagram.

## Things in the UML

There are four kinds of things in the UML:

- Structural things
- Behavioral things
- Grouping things
- Annotational things

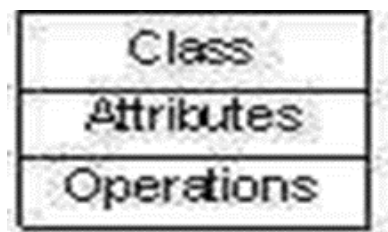
These things are the basic object-oriented building blocks of the UML. You use them to write well-formed models.

### Structural Things

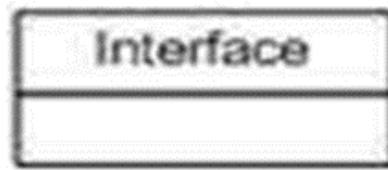
Structural Things are the nouns of UML models. These are the mostly static parts of a model, representing elements that are either conceptual or physical. Collectively, the structural things are called classifiers.

A class is a description of a set of objects that share the same attributes, operations, relationships, and semantics. A class implements one or more interfaces. Graphically, a class is rendered as a rectangle, usually including its name, attributes, and operations

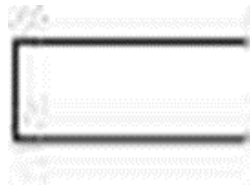
**Class** - A Class is a set of identical things that outlines the functionality and properties of an object. It also represents the abstract class whose functionalities are not defined. Its notation is as follows



**Interface** - A collection of functions that specify a service of a class or component, i.e., Externally visible behavior of that class.



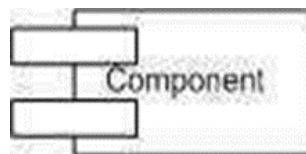
**Collaboration** - A larger pattern of behaviors and actions. Example: All classes and behaviors that create the modeling of a moving tank in a simulation.



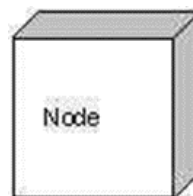
**Use Case** - A sequence of actions that a system performs that yields an observable result. Used to structure behavior in a model. Is realized by collaboration.



**Component** - A physical and replaceable part of a system that implements a number of interfaces. Example: a set of classes, interfaces, and collaborations.



**Node** - A physical element existing at run time and represents are source.



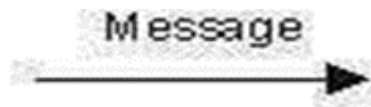
## Behavioral Things

Behavioral things are the dynamic parts of UML models. These are the verbs of a model, representing behavior over time and space. In all, there are three primary kinds of behavioral things

- Interaction
- State machine

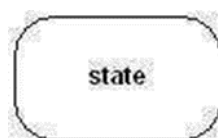
### Interaction

It is a behavior that comprises a set of messages exchanged among a set of objects or roles within a particular context to accomplish a specific purpose. The behavior of a society of objects or of an individual operation may be specified with an interaction. An interaction involves a number of other elements, including messages, actions, and connectors (the connection between objects). Graphically, a message is rendered as a directed line, almost always including the name of its operation.



### State machine

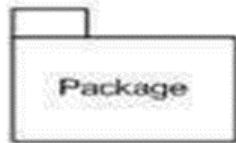
State machine is a behaviour that specifies the sequences of states an object or an interaction goes through during its lifetime in response to events, together with its responses to those events. The behaviour of an individual class or a collaboration of classes may be specified with a state machine. A state machine involves a number of other elements, including states, transitions (the flow from state to state), events (things that trigger a transition), and activities (the response to a transition). Graphically, a state is rendered as a rounded rectangle, usually including its name and its substates.



## Grouping Things

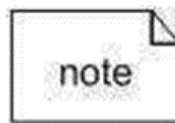
Grouping things can be defined as a mechanism to group elements of a UML model together. There is only one grouping thing available.

**Package** – Package is the only one grouping thing available for gathering structural and behavioral things.



## Annotational Things

Annotational things are the explanatory parts of UML models. These are the comments you may apply to describe, illuminate, and remark about any element in a model. There is one primary kind of annotational thing, called a note. A note is simply a symbol for rendering constraints and comments attached to an element or a collection of elements.



## Relationships in the UML

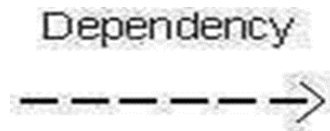
Relationship is another most important building block of UML. It shows how the elements are associated with each other and this association describes the functionality of an application.

There are four kinds of relationships in the UML:

- Dependency
- Association
- Generalization
- Realization

## Dependency

It is an element (the independent one) that may affect the semantics of the other element (the dependent one). Graphically, a dependency is rendered as a dashed line, possibly directed, and occasionally including a label.



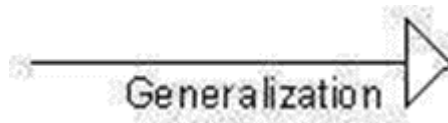
## Association

Association is basically a set of links that connects the elements of a UML model. It also describes how many objects are taking part in that relationship.



## Generalization

It is a specialization/generalization relationship in which the specialized element (the child) builds on the specification of the generalized element (the parent). The child shares the structure and the behavior of the parent. Graphically, a generalization relationship is rendered as a solid line with a hollow arrowhead pointing to the parent.



## Realization

Realization can be defined as a relationship in which two elements are connected. One element describes some responsibility, which is not implemented and the other one implements them. This relationship exists in case of interfaces.

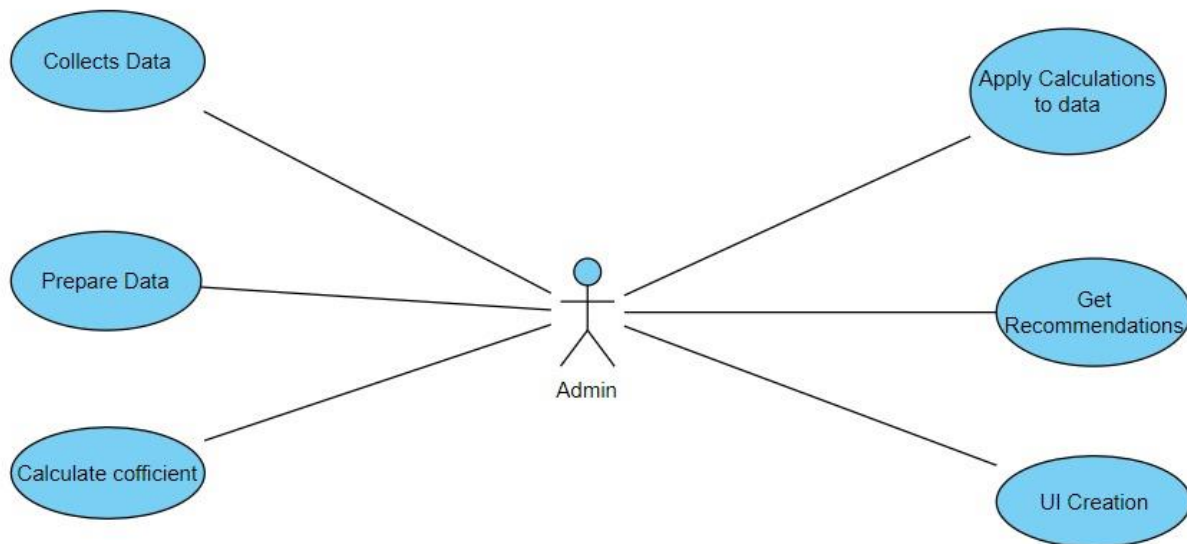


## 6.3 UML DIAGRAMS

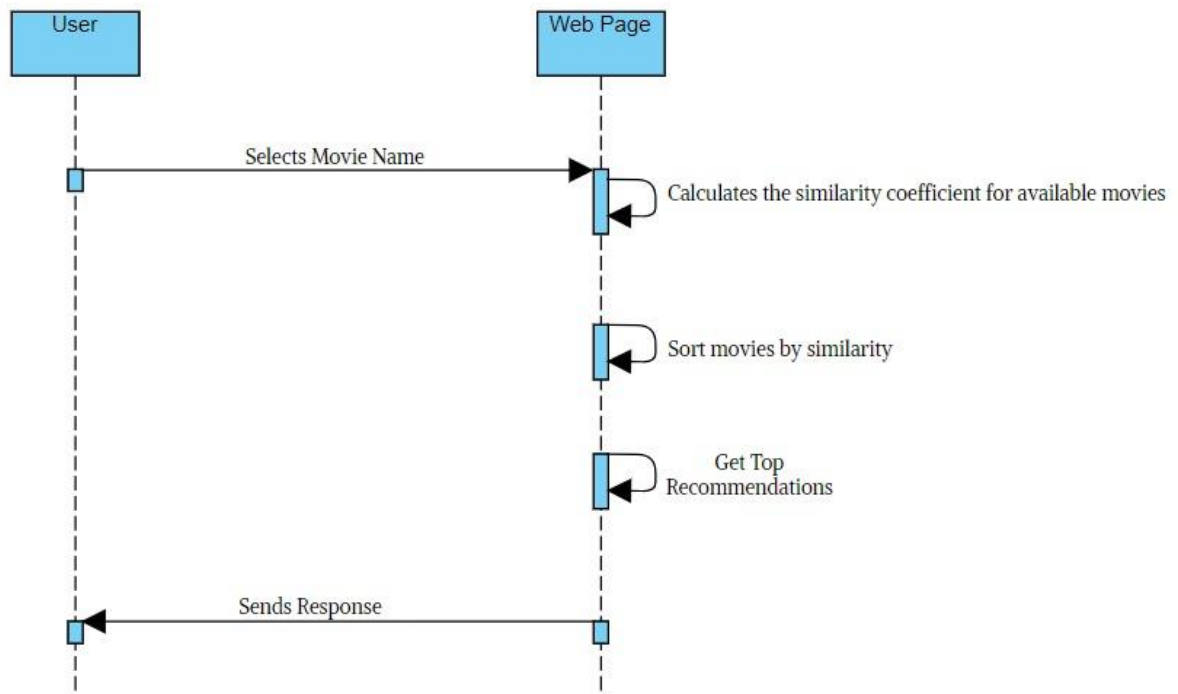
UML is a modern approach to modeling and documenting software. It is based on diagrammatic representations of software components. It is the final output, and the diagram represents the system<sup>[6]</sup>.

UML includes the following

- Class diagram
- Object diagram
- Component diagram
- Composite structure diagram
- Use case diagram (Refer to image 6.2)
- Sequence diagram (Refer to image 6.3)
- Communication diagram
- State diagram
- Activity diagram



**Fig 6.2 Use-Case Diagram**



**Fig 6.3 Sequence Diagram**

## 7. IMPLEMENTATION

### 7.1 Data

- The dataset utilized for our project is sourced from Kaggle's Movie dataset, a rich repository containing a vast collection of movie ratings and user preferences.
- This dataset encompasses ratings given by users to a wide array of movies across different genres and time periods, providing a comprehensive view of viewer preferences and behaviors.
- Covering a diverse range of movies spanning various genres such as action, comedy, drama, and more, the dataset offers ample opportunities for building a robust collaborative filtering recommendation system.
- Structured into distinct user-item interaction matrices, the dataset enables seamless implementation of collaborative filtering techniques, facilitating the prediction of user preferences based on similar user behavior.
- With thousands of users and movies represented in the dataset, our project benefits from a sizable volume of data, allowing for the development of accurate and personalized recommendation models.
- Each interaction within the dataset, comprising user ratings for specific movies, serves as valuable input for training collaborative filtering algorithms, empowering the system to suggest relevant and engaging movie recommendations tailored to individual user tastes.

#### **Data set Link:**

<https://www.kaggle.com/datasets/rounakbanik/the-movies-dataset>

## 7.2 Sample Code

### 7.2.1 Data Collection and Preprocessing:

```
# Importing Libraries
import numpy as np
import pandas as pd
import matplotlib.pyplot as plt
import seaborn as sns
from ipywidgets import *
sns.set_style('white')
%matplotlib inline

# Load Data
column_names = ['user_id', 'item_id', 'rating', 'timestamp']
df = pd.read_csv('dataset.csv', sep='\t', names=column_names)
movie_titles = pd.read_csv('movieIdTitles.csv')
df.head()

movie_titles = pd.read_csv('movieIdTitles.csv')
movie_titles.head()

# Merge Data
df = pd.merge(df, movie_titles, on='item_id')
df.head()

df.groupby('title')['rating'].mean().sort_values(ascending = False).head()

df.groupby('title')['rating'].count().sort_values(ascending = False).head()

ratings = pd.DataFrame(df.groupby('title')['rating'].mean())
ratings.head()

ratings['numOfRatings'] =
pd.DataFrame(df.groupby('title')['rating'].count())
ratings.head()
# Create Pivot Table
```

```

moviemat = df.pivot_table(index='user_id', columns='title',
values='rating')
df.head()

ratings = pd.DataFrame(index=moviemat.columns,
columns=['FirstMovieRecommendation', 'SecondMovieRecommendation',
'ThirdMovieRecommendation', 'FourthMovieRecommendation'])
df.head()

```

### 7.2.2 Similarity Computation:

```

def bhattacharyya_coefficient(p, q):
    # Filter out ratings greater than 0
    p_nonzero = p[p > 0]
    q_nonzero = q[q > 0]

    # Compute relative frequencies
    max_rating = int(np.max(np.concatenate((p_nonzero, q_nonzero)))) #
Determine the maximum rating
    p_freq = np.array([np.sum(p_nonzero == i) for i in range(1, max_rating
+ 1)]) / len(p_nonzero)
    q_freq = np.array([np.sum(q_nonzero == i) for i in range(1, max_rating
+ 1)]) / len(q_nonzero)

    # Compute Bhattacharyya coefficient
    bc = np.sum(np.sqrt(p_freq * q_freq))
    return bc

def get_movie_recommendations_for_single_movie(movie_name, moviemat):
    if movie_name not in moviemat.columns:
        return "Movie not found in the dataset"

    # Get ratings for the specified movie
    movie_ratings = moviemat[movie_name]

    # Calculate s(U,V) with all other movies
    similar_movies = {}
    for other_movie in moviemat.columns:
        if movie_name != other_movie:
            other_ratings = moviemat[other_movie]
            similarity = bhattacharyya_coefficient(movie_ratings,
other_ratings)

```

```

        similar_movies[other_movie] = similarity

    # Sort the similar movies by similarity
    similar_movies_sorted = sorted(similar_movies.items(), key=lambda x:
x[1], reverse=True)

    # Get top recommendations
    top_recommendations = [movie for movie, similarity in
similar_movies_sorted[:4]]

    # Return recommendations
    return top_recommendations

```

### 7.2.3 Recommendation Generation:

```

# Dynamically enter the movie name
movie_name_input = input("Enter the movie name: ")

# Get recommendations for the entered movie
recommendations =
get_movie_recommendations_for_single_movie(movie_name_input, moviemat)

# Print recommendations in the desired format
print("Recommendations for", movie_name_input, ":")
for recommendation in recommendations:
    print(recommendation)

```

### Output:

```

Enter the movie name: Á köldum klaka (Cold Fever) (1994)
Recommendations for Á köldum klaka (Cold Fever) (1994):
All Over Me (1997)
All Things Fair (1996)
Angela (1995)
B. Monkey (1998)

```

### 7.2.4 Model Evaluation and Optimization:

While the provided code implements a movie recommendation system, directly measuring its accuracy isn't straightforward<sup>[7]</sup>. Here's why:

- **Collaborative Filtering Approach:** The code uses collaborative filtering, which recommends movies based on user ratings for similar movies. It doesn't predict exact ratings, but rather suggests movies users might enjoy based on past preferences.
- **Missing Ground Truth:** There's no "ground truth" data available to compare the recommendations with. We don't know for sure if users would actually like the recommended movies because they haven't rated them yet.

However, we can evaluate the system's effectiveness indirectly using these approaches:

#### 1. Filtering Metrics:

- **Recall:** This measures how well the system recommends movies a user has already rated highly. You could calculate recall by comparing the number of recommended movies a user has rated highly (above a certain threshold) to the total number of such movies the user has watched. Higher recall suggests the system is picking up on user preferences effectively.
- **Precision:** This measures how relevant the recommendations are. You could calculate precision by dividing the number of recommended movies a user actually likes (based on ratings) by the total number of recommendations made. Higher precision indicates the system is suggesting movies with a high chance of being enjoyed.

#### 2. User Engagement:

If this system is implemented in a real-world application, you can track user behavior:

- How often users accept recommendations?
- Do they watch the recommended movies?
- Do they rate the recommendations highly?
- Higher engagement suggests the recommendations are relevant and useful to users.

### 3. A/B Testing:

You could compare the recommendation system with a baseline approach (e.g., random recommendations) or a different recommendation algorithm.

Track which system leads to higher user engagement metrics like clicks, views, or ratings. Incorporating these methods can provide valuable insights into the effectiveness of your recommendation system.

#### 7.2.5 UI Creation using Streamlit:

```
import streamlit as st
import numpy as np
import pandas as pd
import requests
import re

# Load Data
column_names = ['user_id', 'item_id', 'rating', 'timestamp']
df = pd.read_csv('dataset.csv', sep='\t', names=column_names)
movie_titles = pd.read_csv('movieIdTitles.csv')
df = pd.merge(df, movie_titles, on='item_id')
moviemat = df.pivot_table(index='user_id', columns='title',
values='rating')

# TMDb API Details
API_KEY =
'eyJhbGciOiJIUzI1NiJ9.eyJhdWQiOiJmMjEwZWZzOGJjNDhkZTQyZDQ1YzMzMzUxMGZkODI1M
yIsInN1YiI6IjY1ZmQzMTEwMzc4MDYyMDE3ZTg3MWM2ZCIsInNjb3BlcyI6WyJhcGlfcmlhZCIj
LCJ2ZXJzaW9uIjoxfQ.ifeHqFnSL7MLwhmIbn7mMN0fwRdD6DmWTwmvkQpZqQQ'
BASE_URL = 'https://api.themoviedb.org/3'
IMAGE_BASE_URL = 'https://image.tmdb.org/t/p/w500'

def bhattacharyya_coefficient(p, q):
    p_nonzero = p[p > 0]
    q_nonzero = q[q > 0]
    max_rating = int(np.max(np.concatenate((p_nonzero, q_nonzero))))
    p_freq = np.array([np.sum(p_nonzero == i) for i in range(1, max_rating
+ 1)]) / len(p_nonzero)
```





```

params = {
    "query": title,
    "year": year, # Add year as a parameter
    "include_adult": False,
    "language": "en-US",
    "page": 1
}
print("API Request Params:", params)
response = requests.get(url, headers=headers, params=params)
print("API Response:", response.json())
data = response.json()
if 'results' in data and data['results']:
    poster_path = data['results'][0].get('poster_path')
    if poster_path:
        base_url = "https://image.tmdb.org/t/p/w500" # Adjust the size
as needed
        poster_url = f"{base_url}{poster_path}"
        return poster_url
return None

```

```

def get_movie_recommendations_for_single_movie(movie_name, moviemat):
    if movie_name not in moviemat.columns:
        return "Movie not found in the dataset"

    movie_ratings = moviemat[movie_name]
    similar_movies = {}
    for other_movie in moviemat.columns:
        if movie_name != other_movie:
            other_ratings = moviemat[other_movie]
            similarity = bhattacharyya_coefficient(movie_ratings,
other_ratings)
            similar_movies[other_movie] = similarity
    similar_movies_sorted = sorted(similar_movies.items(), key=lambda x:
x[1], reverse=True)
    top_recommendations = [movie for movie, similarity in
similar_movies_sorted[:4]]
    return top_recommendations

add_bg_from_url()

```

```

# Streamlit App
st.title("Movie Recommendation System")

```

```

# Movie Selection
selected_movie = st.selectbox("Select a movie:", df['title'].unique())

# Display Recommendations with Posters
if st.button("Get Recommendations"):
    recommendations =
get_movie_recommendations_for_single_movie(selected_movie, moviemat)
    st.write("Top Recommendations for", selected_movie)

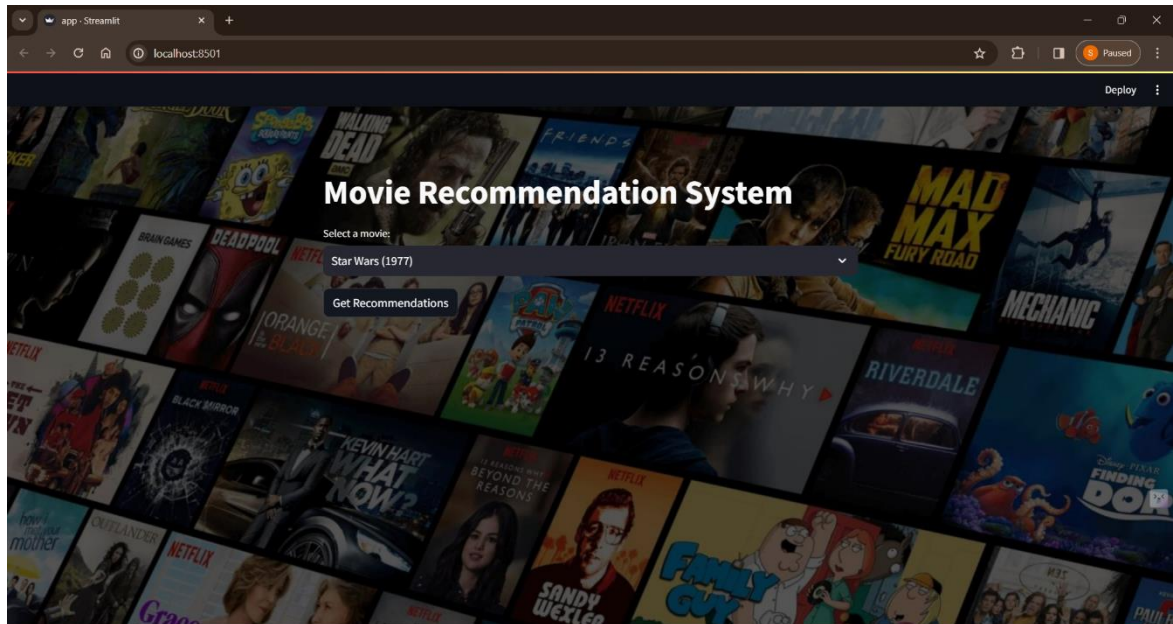
# Display recommendations in a 1x4 grid
col1, col2, col3, col4 = st.columns(4) # Create four columns
for i, title in enumerate(recommendations):
    try:
        poster_url = fetch_poster_from_tmdb_by_title(title)
    except KeyError:
        poster_url = None

    if poster_url is not None:
        # Display each poster in a smaller size
        with globals()[f"col{i+1}"]:
            st.image(poster_url, caption=title, use_column_width=True)
    else:
        # If no poster is available, display a placeholder
        with globals()[f"col{i+1}"]:
            st.write(f"No poster available for {title}")

```

## 7.3 Results

## 1. Navigating to Home page

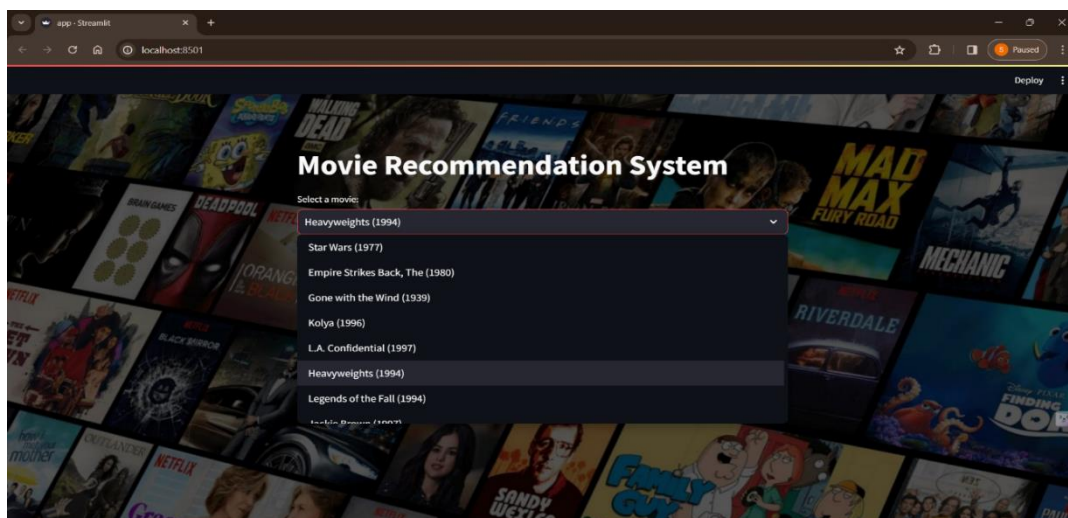


## 7.1 Initial Web Page

To start you can click on the drop-down menu or else you can search for the movie you have seen before. (Refer to image 7.1)

## 2. Services

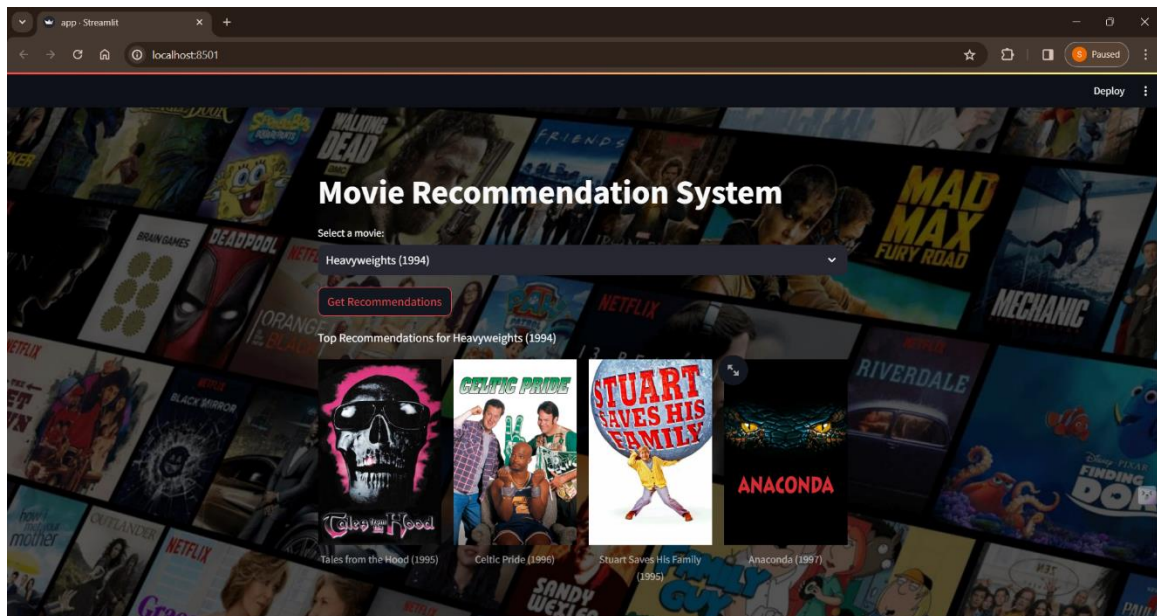
### i. Select or search for the movie



## 7.2 Selecting a movie

- Select the movie from the drop-down menu.
- Click get recommendations. (Refer to image 7.2)

## ii. Get recommendations



## 7.3 Output of recommended movies

## **8. TESTING**

### **8.1 INTRODUCTION TO TESTING**

Software Testing is defined as an activity to check whether the actual results match the expected results and to ensure that the software system is Defect free. It involves the execution of a software component or system component to evaluate one or more properties of interest. It is required for evaluating the system. This phase is the critical phase of software quality assurance and presents the ultimate view of coding.

#### **Importance of Testing**

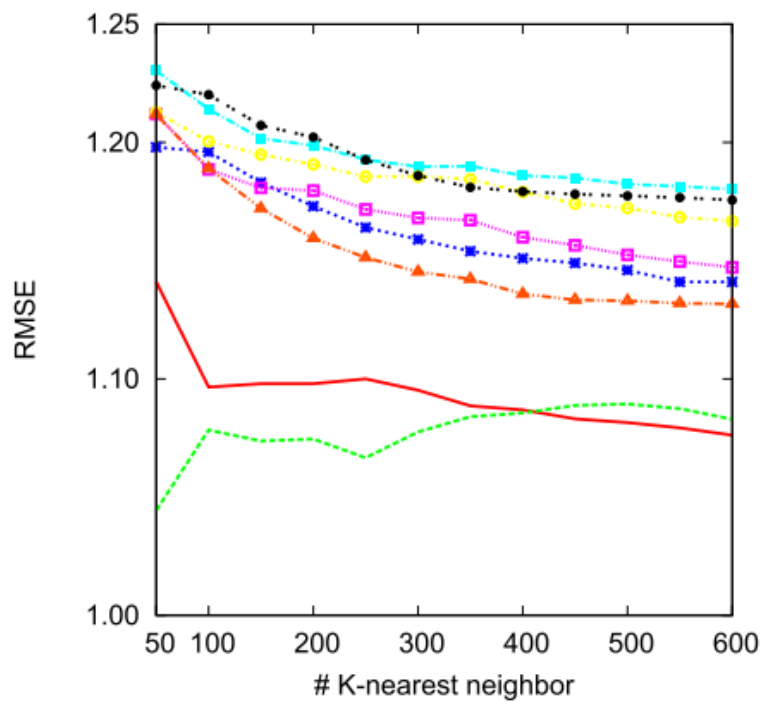
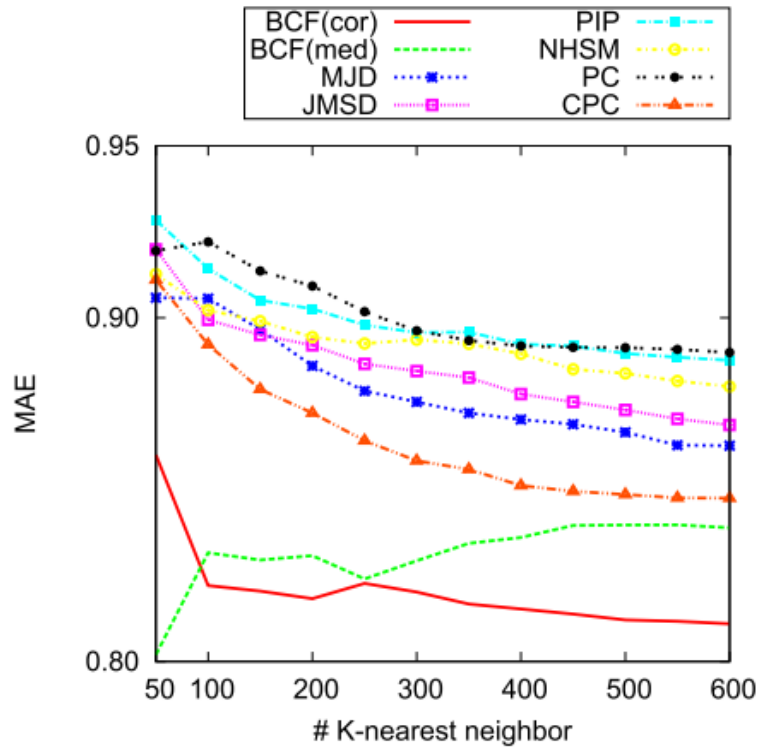
The importance of software testing is imperative. A lot of times this process is skipped, therefore, the product and business might suffer. To understand the importance of testing, here are some key points to explain

- Software Testing saves money
- Provides Security
- Improves Product Quality
- Customer satisfaction

Testing is of different ways The main idea behind the testing is to reduce the errors and do it with a minimum time and effort.

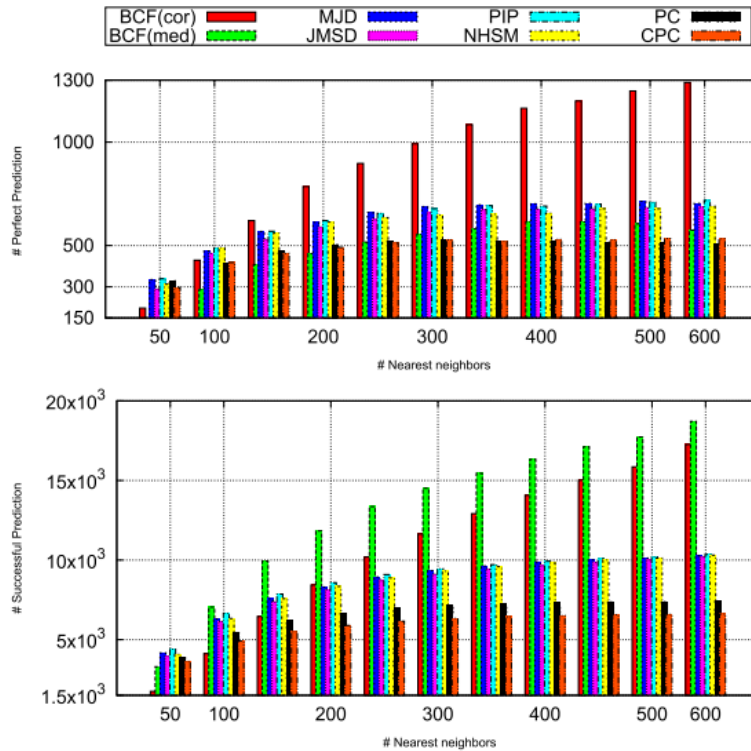
#### **Benefits of Testing**

- **Cost-Effective:** It is one of the important advantages of software testing. Testing any IT project on time helps you to save your money for the long term. In case if the bugs caught in the earlier stage of software testing, it costs less to fix.
- **Security:** It is the most vulnerable and sensitive benefit of software testing. People are looking for trusted products. It helps in removing risks and problems earlier.
- **Product quality:** It is an essential requirement of any software product. Testing ensures a quality product is delivered to customers.
- **Customer Satisfaction:** The main aim of any product is to give satisfaction to their customers. UI/UX Testing ensures the best user experience.



8.1 Graphs for MAE and RMSE results

The MAE<sup>[1]</sup> and RMSE obtained from Net2 subset are plotted and shown in above figure (Refer to figure 8.1) Traditional PC based CF produces highest MAE. However, another traditional CPC based CF shows better performance than CFMJD; CFJMSD; CFPIP; CFNHSM; CFPC in both metrics (MAE, RMSE). It may be noted that CPC based CF makes least number of valid predictions (below figure). The CPC makes 6613 times valid predictions whereas other non-traditional measures (PIP, MJD) make close to 11,000 times valid predictions. The CFBCF (med) can perform close to 20,000 successful predictions out 24,408 requested predictions. Proposed BCF based collaborative filtering CFBCF (corr); CFBC (med) outperform all CFs including CPC based CF in both accuracy metrics. The CFBCF (corr); is the best and it makes lowest MAE (0.811) and RMSE (1.076) with K = 600. We report number of successful predictions and number of perfect predictions performed by each CF in the below figure (Refer to figure 8.2) and similar trends are found.



## 8.2 Prediction Trends



The F1 measure<sup>[1]</sup> provide qualitative importance of recommendations suggested by a CF. F1 measures of all CFs are obtained after getting executed them on YM subset. Here, we found that our proposed BCF measure-based CFs provide better reliable recommendations compared to the other CFs (Fig. 15). The proposed CFBCF (med) provides remarkably good recommendations ( $F1 = 0.68$ ) with average user ratings less than 4 (Table 4). The CFBC (med) retrieves 10% more accurate good items (relevant) compared to MJD based CF. It can be noted that CPC based CF produces second lowest F1 measure. Other BCF measure based CFBCF (corr) is also found to be outperforming existing measure-based CFs and it attains its max value of 0.63 at  $K = 400$ .

## 9. CONCLUSION

The quest for ever-more personalized recommendations in recommender systems necessitates innovative approaches for capturing user preferences, particularly in the face of sparse data. Our project addresses this challenge by leveraging neighbourhood-based collaborative filtering (CF) with Bhattacharya similarity.

This approach offers a distinct advantage over traditional CF methods. By incorporating all user ratings, even for unshared items, Bhattacharya similarity paints a richer picture of user taste profiles. This refined understanding translates to more accurate and relevant recommendations, especially in sparse data scenarios.

Beyond enhanced accuracy, our project fosters improved user experiences. By delving deeper into user preferences, the system generates recommendations that resonate more effectively with individual needs and preferences. This, in turn, can lead to increased user satisfaction, engagement, and loyalty within the recommendation platform.

Furthermore, neighbourhood-based CF offers a level of interpretability. We can analyse the ratings of similar users to understand the reasoning behind recommendations, which can be valuable for debugging and improving the system. This fosters trust and confidence in the system by providing users with insights into the recommendation process.

In conclusion, our project on neighbourhood-based CF with Bhattacharya similarity presents a robust and innovative approach for recommendation systems. By addressing data sparsity and promoting user understanding, this project paves the way for a more personalized and engaging user experience across various recommendation domains.

## 10. FUTURE SCOPE

Our project on neighborhood-based CF with Bhattacharya similarity offers a promising foundation for personalized recommendations. Here are some exciting avenues for future exploration:

### 1. Hybrid Recommendation Approaches:

**Incorporating Side Information:** Currently, the project focuses on user ratings. We can explore integrating user/item attributes (e.g., demographics, genre) to enrich user profiles and potentially improve recommendation accuracy.

**Combining CF with Content-Based Filtering:** Consider a hybrid approach that leverages both user-item interactions (CF) and inherent item properties (content-based filtering) within the neighborhood-based framework. This can provide a more comprehensive understanding of user preferences and item characteristics.

### 2. Dynamic Recommendation Systems:

**Time-Aware Recommendations:** User preferences can evolve over time. We can explore incorporating time-decay weights into the similarity calculation to prioritize recent user ratings and capture evolving preferences.

**Contextual Adaptation:** Consider extending the system to incorporate contextual information (e.g., time of day, location) to tailor recommendations dynamically based on the user's current context.

### 3. Scalability and Interpretability:

**Efficient Algorithm Design:** As datasets grow, explore parallelization techniques or efficient data structures to maintain scalability of the neighborhood-based CF approach with Bhattacharya similarity.

**Explainability of Recommendations:** While neighborhood-based CF offers some interpretability, delve deeper into techniques that explain why specific neighbors are chosen for a user, providing more transparency into the recommendation process.

By exploring these future directions, we can further enhance the effectiveness and personalization capabilities of our neighborhood-based CF system with Bhattacharya similarity, ensuring its adaptability and scalability in real-world recommender system applications.

## 11. REFERENCE LINKS

- [1] Patra, B. K., Launonen, R., Ollikainen, V., & Nandi, S. (2015). A new similarity measure using Bhattacharyya coefficient for collaborative filtering in sparse data. Knowledge-Based Systems. <https://doi.org/10.1016/j.knosys.2015.03.001>
- [2] Bobadilla, J., Ortega, F., Hernando, A., & Bernal, J. (2012). A collaborative filtering approach to mitigate the new user cold start problem. Knowledge-Based Systems, 26, 225–238. <https://doi.org/10.1016/j.knosys.2011.07.021>
- [3] Lendave, V. (2021, September 24). Cold-Start Problem in Recommender Systems and its Mitigation Techniques. Analytics India Magazine. <https://analyticsindiamag.com/cold-start-problem-in-recommender-systems-and-its-mitigation-techniques/>
- [4] Turing. (2022, July 29). How does collaborative filtering work in recommender systems? <https://www.turing.com/kb/collaborative-filtering-in-recommender-system>
- [5] Dixit, V. S., & Jain, P. C. (2019). Proposed similarity measure using Bhattacharyya coefficient for context aware recommender system. Journal of Intelligent and Fuzzy Systems, 36(4), 3105–3117. <https://doi.org/10.3233/jifs-18341>
- [6] Nishadha. (2022, September 28). UML Diagram Types Guide: Learn About All Types of UML Diagrams with Examples. <https://creately.com/blog/diagrams/uml-diagram-types-examples/>
- [7] Deutschman, Z. (2023, August 7). Recommender Systems: Machine learning metrics and business metrics. neptune.ai. <https://neptune.ai/blog/recommender-systems-metrics>