

Book Recommendation System

Sattwik Chowdhury
School Of Computer Engineering
Kalinga Institute Of Industrial
Technology
Odisha, India
sattwikchowdhury3@gmail.com

Abstract-The project represents a Book Recommendation System that utilize the strength of filtering techniques to give personalized book suggestion. The model contains K-Nearest Neighbors (KNN) for similar based recommendations and applies Singular Value Decomposition(SVD) and Non negative Matrix Factorization to enhance accuracy. The datasets undergoes in processing , visualization , cleaning to advance the quality of the data. Grid Search is used in hyperparameter tuning, and the model is checked by using its precision , recall , and F1- score.The model efficaciously predicts user preferences, with future betterment which aimed to be the integrating hybrid techniques and real-time updates of its enhanced recommendations

Keywords-Book Recommendation System, Collaborative Filtering, K-Nearest Neighbors (KNN), Singular Value Decomposition (SVD), Non-Negative Matrix Factorization (NMF), Matrix Factorization, Hyperparameter Tuning, Data Preprocessing, Precision, Recall, F1-Score, Personalized Recommendations, Machine Learning, User-Item Interactions.

I. INTRODUCTION

With the rapid growth of digital content, personalized recommendation systems have become essential for enhancing user experience. The **Book Recommendation System** aims to provide personalized book suggestions based on users' past interactions and preferences. Traditional methods of book recommendations, such as bestseller lists and editorial picks, do not cater to individual users' tastes. In contrast, **collaborative filtering techniques**, particularly **K-Nearest Neighbors (KNN)** and **matrix factorization methods** like **Singular Value Decomposition (SVD)** and **Non-Negative Matrix Factorization (NMF)**, leverage past user behavior to generate relevant recommendations.

This project involves preprocessing a dataset containing book ratings and user interactions, applying data cleaning techniques, and performing exploratory data analysis (EDA) to uncover insights. The recommendation model is trained using collaborative filtering, where KNN is employed to measure similarity between users and items, while SVD and NMF are used to extract latent features from user-item interactions. Model performance is evaluated using key metrics such as **precision, recall, and F1-score** to ensure the accuracy and relevance of recommendations. Additionally, **hyperparameter tuning** using **Grid Search** is implemented to optimize the model's performance.

By integrating these techniques, the system provides a robust and scalable recommendation framework, helping users discover books tailored to their reading preferences. Future enhancements include integrating **content-based filtering, deep learning approaches, and real-time updates** to further improve the recommendation quality.

II. LIBRARIES USED

The following libraries are used in the implementation of this system:

- **Pandas**: For data manipulation and preprocessing.
- **NumPy**: For numerical computations.
- **Matplotlib & Seaborn**: For data visualization.
- **Scikit-learn**: For implementing machine learning techniques, including hyperparameter tuning.
- **Surprise (Scikit-Surprise)**: For implementing collaborative filtering models such as SVD and NMF.
- **SciPy**: For advanced mathematical computations in matrix factorization.
- **Warnings**: Used to suppress warnings for cleaner outputs.
- **Random**: For generating random numbers and ensuring variability in experiments.
- **Scipy.spatial.distance (Correlation)**: Used to compute similarity measures for collaborative filtering.
- **Scipy.sparse (csr_matrix, svds)**: Helps in handling sparse matrices and performing dimensionality reduction.
- **sklearn.neighbors (NearestNeighbors)**: Used for KNN-based similarity measurements.
- **sklearn.decomposition (TruncatedSVD)**: Used for dimensionality reduction in collaborative filtering.
- **sklearn.metrics.pairwise (cosine_similarity, pairwise_distances)**: Computes similarity measures between users or items.
- **sklearn.model_selection (train_test_split, GridSearchCV, cross_validate)**: For dataset splitting, hyperparameter tuning, and cross-validation.

III. DATASET

The dataset used for this project consists of book ratings and user interactions. The key datasets include:

A. Book Dataset

- ISBN(unique for each book)
- Book-Title
- Book-Author
- Year of Publication

- Publisher

B. Ratings Dataset

- User-ID
- ISBN
- Image-URL-S
- Image-URL-M
- Image-URL-L
- Shape of Dataset - (271360, 8)
- Book-Rating
- Shape of Dataset - (1149780, 3)

C. Users Dataset

- User-ID (unique for each user)
- Location (contains city, state and country separated by commas)
- Age
- Shape of Dataset - (278858, 3)

IV. DATA PROCESSING AND CLEANING

Before training the model, the dataset undergoes preprocessing to ensure data quality:

A. Handling Missing Values:

Books or ratings with missing essential information are either filled using imputation techniques (such as mean or median values for numerical fields) or removed if the missing data is significant and cannot be reliably inferred. Missing categorical values such as book authors or publishers are sometimes replaced with a placeholder like 'Unknown' to retain as much data as possible.

B. Data Normalization:

It refers to the process of transforming rating values into a standardized scale to ensure consistency in the dataset. This is particularly important in recommendation systems where different users may have different rating behaviors. For instance:

Scaling Ratings: If different users rate books on different scales (e.g., one user tends to give ratings between 6-10 while another between 1-5), normalization ensures that all ratings are on the same scale (e.g., between 0 and 1 or mean-centered).

Reducing Bias: Some users may consistently rate higher or lower than others. Normalization helps adjust these biases to prevent them from skewing recommendations.

Improving Model Performance: Many machine learning algorithms, including matrix factorization techniques like **SVD and NMF**, perform better when input data is normalized, reducing computational complexity and improving accuracy.

Common Methods:

Min-Max Scaling: Transforms ratings to a fixed range, usually [0,1].

Z-score Normalization: Standardizes ratings by subtracting the mean and dividing by the standard deviation.

By applying data normalization, we ensure that the recommendation system treats all user ratings fairly, leading to more reliable book suggestions.

C. Removing Duplicates

It ensures that the dataset remains clean and free from repeated entries, which can negatively impact the recommendation system. Duplicates can occur when:

1. **Users rate the same book multiple times**, leading to repeated records.
2. **The same book appears more than once** due to formatting differences or data errors.
3. **Merging datasets introduces redundant entries**, causing inflated interactions.

To fix this, we:

- **Identify duplicate entries** using built-in functions like `.duplicated()`.
- **Remove unnecessary repeats** with `.drop_duplicates()`, keeping only unique records.
- **Ensure data consistency** so that recommendations remain accurate.

This step improves the system's efficiency and ensures users receive better book suggestions.

D. Filtering Sparse Data:

It helps improve the accuracy and efficiency of the recommendation system by removing users and books with very few interactions. Sparse data means that most users have rated only a small number of books, making it difficult for the model to learn meaningful patterns.

To address this:

1. **Users with very few ratings** are removed since they don't provide enough data for personalized recommendations.
2. **Books with very few ratings** are also filtered out, as they lack sufficient user feedback to make reliable suggestions.
3. **A threshold is set** (e.g., users must have rated at least 5 books, and books must have at least 10 ratings) to ensure better learning and more relevant recommendations.

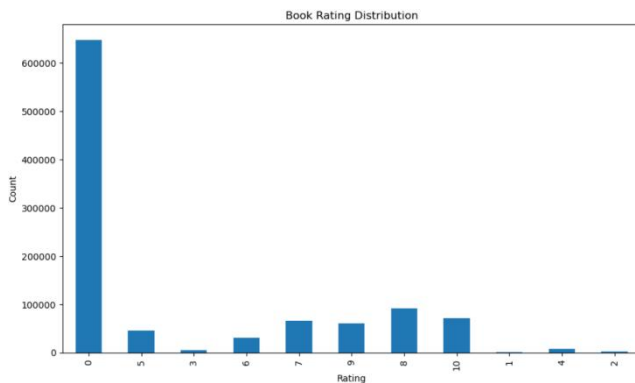
By removing sparse data, the model focuses on meaningful interactions, improving robustness and reducing computational complexity.

V. DATA VISUALIZATION

To gain deeper insights into the dataset, **Exploratory Data Analysis (EDA)** is conducted to understand user behavior, book popularity, and rating trends. This step is crucial in identifying patterns and potential anomalies in the data before building the recommendation model.

1. User Rating Distribution:

- A **histogram** is plotted to visualize how ratings are distributed across the dataset.



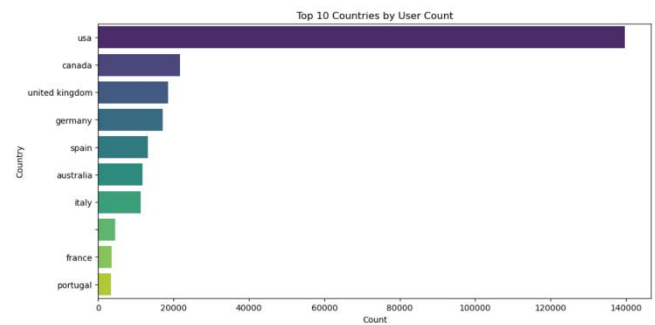
- This helps determine whether users tend to rate books positively (high ratings) or negatively (low ratings).
- It also highlights if the dataset is **skewed**, meaning whether users mostly provide extreme ratings (very high or very low) rather than balanced ratings.

2. Top Rated Books:

- The **most popular books** are identified based on the number of ratings they receive.
- Books with the highest **average ratings** and those with the most **user interactions** are analyzed separately.
- This helps in understanding **which books are consistently favored** by readers and can be useful for **popularity-based recommendations**.

3. User Activity Analysis:

- The **number of books rated per user** is visualized using bar charts or density plots.

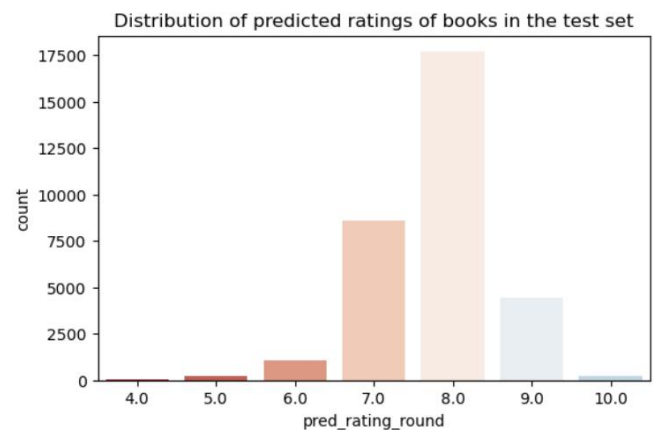
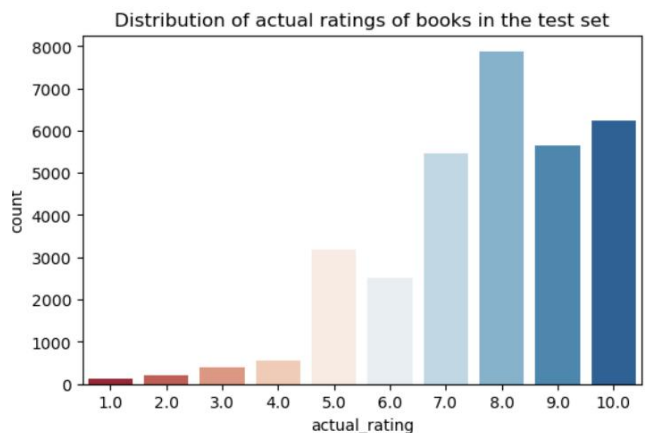


- This helps identify **active users** who contribute significantly to the dataset and **inactive users** with very few ratings.
- Understanding user engagement patterns allows for filtering out **sparse users**, ensuring that only meaningful interactions contribute to recommendations.

VI. COLLABORATIVE FILTERING (CF) MODEL

The core of the recommendation system is based on **Collaborative Filtering (CF)**, which relies on past user interactions to predict future preferences. The system uses **matrix factorization techniques** such as **SVD (Singular Value Decomposition)** and **NMF (Non-Negative Matrix Factorization)** to improve recommendation accuracy.

- A confusion matrix showing precision, recall, and F1-score for recommendation accuracy.



Why Use K-Nearest Neighbors (KNN)?

KNN is used for similarity-based recommendations. The primary reasons for using KNN in this project include:

1. **User-User Similarity:**

- KNN helps find users with **similar reading preferences** by comparing their book ratings.
- If two users have rated similar books in a similar way, KNN assumes they share the same interests and can recommend books that one user has read but the other hasn't.
- This **collaborative filtering** approach ensures that recommendations are **personalized and relevant**.

2. **Item-Item Similarity:**

- Instead of comparing users, KNN can also be used to measure the similarity between books based on how users have rated them.
- Books that are frequently rated similarly by multiple users are considered **similar** and can be recommended together.
- For example, if many users who read "**Harry Potter**" also liked "**The Hobbit**", KNN would recognize this pattern and suggest "**The Hobbit**" to users who enjoyed "**Harry Potter**".

3. **Ease of Interpretation:**

- KNN is highly **intuitive and explainable** compared to some machine learning models.
- The algorithm works by **finding a set number of neighbors (k) with the highest similarity** to a given user or book.
- This approach is easy to understand because it directly shows which users or books are most similar.

Why Use SVD and NMF?

Both **SVD (Singular Value Decomposition)** and **NMF (Non-Negative Matrix Factorization)** are used for **dimensionality reduction** and **latent feature extraction**:

- **SVD**: Decomposes the user-item interaction matrix into three matrices, capturing hidden patterns in the data. It handles sparsity well and improves recommendation quality.
- **NMF**: Ensures that the latent feature matrices contain only non-negative values, making the results more interpretable and suitable for recommendation.

VII. MODEL EVALUATION

To ensure the recommendation system provides accurate and meaningful book suggestions, several **evaluation metrics** are used. These metrics help measure how well the system predicts user preferences and how effective the recommendations are:

1. **Precision**: Measures the proportion of recommended books that were actually relevant.

- It is calculated as:

$$\text{Precision} = \frac{\text{True Positives}}{\text{True Positives} + \text{False Positives}}$$

2. **Recall**: Measures how many of the relevant books were recommended.

- It is calculated as:

$$\text{Recall} = \frac{\text{True Positives}}{\text{True Positives} + \text{False Negatives}}$$

3. **F1-Score**: Harmonic mean of precision and recall.

- It is calculated as:

$$\text{F1-Score} = 2 \times \frac{\text{Precision} \times \text{Recall}}{\text{Precision} + \text{Recall}}$$

4. **Hits@5 and Hits@10**:

- These metrics measure how often relevant books appear in the **top 5 or top 10 recommendations** for a user.
- If a relevant book appears in the top 5 (**Hits@5**) or top 10 (**Hits@10**), the system scores a **hit**.
- These metrics help assess **how well-ranked the relevant books are** within the recommendations.

Evaluation Process

To assess the effectiveness of the **Book Recommendation System**, a structured evaluation process is followed. This ensures that the model provides **reliable and personalized recommendations** based on real user preferences. The evaluation process consists of three key steps:

1. **Data Splitting**:

- The dataset is **divided into training and testing sets** to simulate real-world scenarios where the system predicts user preferences.
- Typically, **80% of the data** is used for training the model, and **20%** is reserved for testing.
- This step ensures that the model learns from past interactions while being evaluated on unseen data.

2. **Recommendation Generation**:

- Once trained, the **model generates personalized recommendations** for each user in the test set.
- The recommendations are based on the techniques implemented, including **KNN (for similarity-based recommendations)** and **SVD/NMF (for collaborative filtering)**.
- Each user's predicted preferences are compared against the actual books they interacted with.

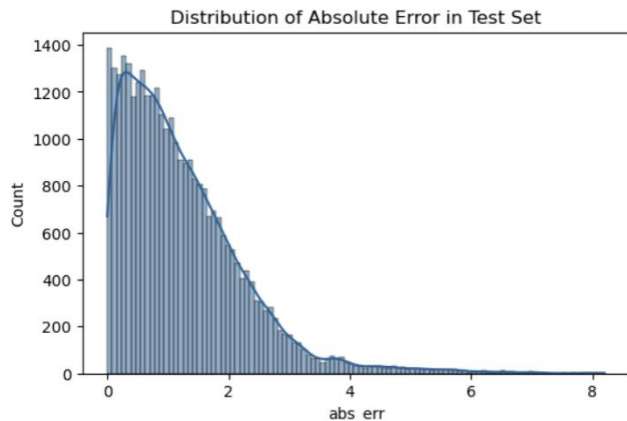
3. **Performance Measurement**:

- After generating recommendations, the system's accuracy is assessed using multiple **evaluation metrics** such as:
 - **Precision**: Measures the proportion of recommended books that were relevant.
 - **Recall**: Measures how many relevant books were successfully recommended.
 - **F1-Score**: Balances precision and recall to provide an overall performance measure.

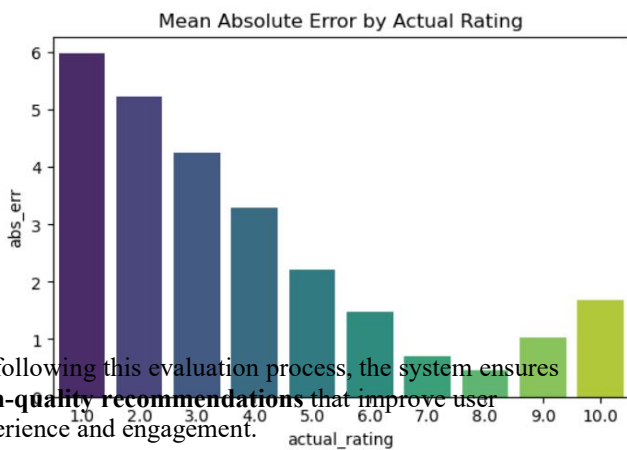
- **Hits@5 and Hits@10:** Evaluates how often relevant books appear in the **top 5 or top 10** recommendations.
- These metrics help determine how **accurate, diverse, and useful** the recommendations are.

Figure: Error Analysis of the Recommendation Model

- The graph shows the **distribution of absolute errors** in the test set, indicating how predictions deviate from actual ratings.



- The second graph visualizes the **Mean Absolute Error (MAE) by actual rating**, showing how error varies across different rating values.



By following this evaluation process, the system ensures **high-quality recommendations** that improve user experience and engagement.

VIII. HYPERPARAMETER TUNING

To enhance the performance of the **Book Recommendation System**, **Grid Search** is used for **hyperparameter tuning**. This ensures that the model operates with the best possible parameters, improving recommendation accuracy and efficiency.

- The best RMSE score achieved was 1.59 using the optimal SVD parameters: $n_factors = 80$, $n_epochs = 20$, $lr_all = 0.005$, and $reg_all = 0.2$, fine-tuned through GridSearchCV.

● KNN Parameters:

- The number of neighbors (**k**) is fine-tuned to determine how many similar users or books should be considered for recommendations.
- A higher **k-value** increases the number of recommendations but may reduce personalization, while a lower **k-value** improves specificity but may limit diversity.
- Grid Search tests multiple values of **k** to find the optimal setting for accurate recommendations.

● SVD Parameters:

- The number of **latent factors** is optimized to ensure that the **matrix factorization** effectively captures meaningful patterns in user-book interactions.
- **Regularization terms** (which prevent overfitting) are fine-tuned to balance generalization and model accuracy.
- Learning rate adjustments ensure efficient convergence of the **stochastic gradient descent (SGD) optimizer** used in SVD.

● NMF Parameters:

- The **learning rate** is adjusted to control how fast the model learns from data.
- The **number of components** (latent features) is fine-tuned to improve the quality of book recommendations.
- Since NMF ensures that all factors are **non-negative**, it provides an interpretable model, and tuning these parameters helps maintain a balance between accuracy and computational efficiency.

IX. IMPLEMENTATION DETAILS

1. Indexing the Data

- The interactions dataset is indexed for quick lookups.
- Books are merged with recommendation results for better presentation.

2. Generating Recommendations

- User preferences are analyzed using collaborative filtering.
- Top-N recommendations are provided, filtering out books the user has already interacted with.

3. User Input Handling

- The system allows users to enter their User ID to receive book recommendations.
- The recommendation engine processes the input and returns the top recommended books.

4. Model Evaluation Process

- The recommendation model is evaluated for each user in the test set.
- Precision, recall, and other metrics are computed and averaged to assess overall performance.

X. CONCLUSION

The **Book Recommendation System** is designed to provide personalized and accurate book suggestions by leveraging

advanced recommendation techniques. It effectively utilizes **collaborative filtering**, which analyzes user interactions to predict preferences. Additionally, **K-Nearest Neighbors (KNN)** is used for similarity-based recommendations, helping users discover books that align with their reading history and interests.

To further enhance the system's performance, **Singular Value Decomposition (SVD)** and **Non-Negative Matrix Factorization (NMF)** are employed for **matrix factorization**, breaking down large, sparse user-book rating matrices into meaningful patterns. These techniques help uncover hidden relationships between books and users, improving the relevance of recommendations.

The **evaluation process** includes precision, recall, F1-score, and hit rate metrics, ensuring that the model provides **accurate, diverse, and relevant book recommendations**. By filtering sparse data, removing duplicates, and optimizing hyperparameters through **GridSearchCV**, the system is fine-tuned to deliver high-quality results. This approach makes the recommendation engine both efficient and effective in assisting users in finding books that match their preferences.

XI. FUTURE IMPROVEMENTS

- **Hybrid Approach:** Combining collaborative filtering with content-based filtering for better recommendations.
- **Deep Learning Integration:** Using deep learning models to improve recommendation accuracy.
- **Real-Time Updates:** Implementing real-time updates to adapt to user preferences dynamically.

This report provides a detailed overview of the project, including the techniques used, evaluation process, and future improvements.

REFERENCES

· Dataset Source

<https://www.kaggle.com/datasets/arashnic/book-recommendation-dataset>

· Algorithm References

- · Ricci, F., Rokach, L., & Shapira, B. (2015). *Recommender Systems Handbook*. Springer. (For theoretical background on recommendation algorithms like collaborative filtering, KNN, SVD, and NMF.)

· Online Learning Resources

- · YouTube: Various tutorials on book recommendation systems and collaborative filtering.
- <https://www.youtube.com/@statquest>
- <https://www.youtube.com/@dataschool>

- <https://www.youtube.com/@freecodecamp>
- <https://www.youtube.com/@SimplilearnOfficial>
- Towards articles on recommendation systems. (<https://github.com/syedsharin/Book-Recommendation-System-Project/blob/main/BOOK%20RECOMMENDATION%20SYSTEM-Capstone%20Project%204.pdf>)

· Libraries and Tools Used

- · Pedregosa, F., et al. (2011). *Scikit-learn: Machine Learning in Python*. Journal of Machine Learning Research, 12, 2825–2830. (For *scikit-learn* library.)
- Surprise Library for collaborative filtering: <https://surpriselib.com>
- NumPy & Pandas documentation for data preprocessing.