

# **Movie Recommendation System Report**

# 1. Brief Overview of the Dataset and Problem Statement

**Dataset Overview:** The project utilizes movie-related datasets which consists of movies, ratings, and tags. The movie dataset includes attributes like movie ID, title, genres, and tags, while the ratings dataset contains user ID, movie ID, rating, and timestamp. These datasets were merged and cleaned which ensured consistency and usability for building recommendation models.

**Problem Statement:** The proposed movie recommendation system combines Content-Based Filtering (CBF) and Collaborative Filtering (CF) to offer personalized movie suggestions. CBF recommends movies based on similarities in tags, genres, and metadata, while CF leverages historical ratings from similar users or items. By hybridizing both approaches, the system improves accuracy and diversity in recommendations. This method enhances user experience by reducing search time, discovering niche content, and providing more contextual suggestions based on user preferences and collaborative patterns.

## 2. Challenges Faced & Solutions

- 1. **Data Preparation:**
  - **Merging and Cleaning:** Combining multiple datasets (movies, ratings, tags) posed challenges due to inconsistencies and missing values.
  - **Solution:** Data cleaning scripts were developed to handle missing values and merge datasets seamlessly.
- 2. **High Dimensionality:**
  - **Genres and Tags:** The datasets contained numerous features from genres and tags, leading to high dimensionality and complexity.
  - **Solution:** One-hot encoding was used to transform genres and tags into binary features, making the data suitable for model training.
- 3. **Cold-Start Problem:**
  - **New Users/Items:** Content-Based Filtering struggled with recommending items for new users with no prior data.
  - **Solution:** Enhanced Collaborative Filtering with embedding layers to capture latent features and address the cold-start problem.
- 4. **Computational Complexity:**
  - **Model Training:** Training models on large datasets required significant computational resources and time.
  - **Solution:** Implemented efficient algorithms and optimized training parameters to reduce computational load.

## 3. Hypothesis Test and Their Conclusion

User-Tag Consistency Test	Rating Patterns Test
<div>Hypothesis:<ul style="list-style-type: none"><li>• H0: users are random in their tag assignments.</li><li>• H1: Users show consistency in their tag usage patterns.</li></ul></div>	<div>Hypothesis:<ul style="list-style-type: none"><li>• H0: Ratings are uniformly distributed.</li><li>• H1: Ratings show significant patterns or biases.</li></ul></div>

<p>Methods:</p> <ul style="list-style-type: none"> <li>Group tags by user.</li> <li>Calculate the tag repetition rate for each user.</li> <li>Perform a one-sample t-test against a random repetition rate of 0.0.</li> </ul>	<p>Methods:</p> <ul style="list-style-type: none"> <li>Perform a normality test on the rating distribution.</li> <li>Analyze hourly rating patterns.</li> <li>Conduct ANOVA on ratings grouped by hour of the day.</li> </ul>
<p>Results:</p> <ul style="list-style-type: none"> <li>Statistic: 86.034</li> <li>p-value: 0.0</li> <li>Mean Repetition Rate: 0.106</li> </ul>	<p>Results:</p> <ul style="list-style-type: none"> <li>Normality Test p-value: 0.0</li> <li>Time Pattern ANOVA F-statistic: 161.476</li> <li>Time Pattern ANOVA p-value: 0.0</li> <li>Mean Rating: 3.543</li> <li>Rating Standard Deviation: 1.064</li> <li>Number of Ratings: 33,776,103</li> <li>Rating Range: [0.5, 5.0]</li> </ul>
<p>Conclusion:</p> <p>The low p-value (<math>&lt; 0.05</math>) indicates that users show significant consistency in their tag usage patterns, rejecting the null hypothesis.</p>	<p>Conclusion:</p> <p>The significant results from the normality test and ANOVA indicate that movie ratings are not uniformly distributed and show significant patterns, rejecting the null hypothesis.</p>

## 4. Performance Comparison of Trained Models

### Training and Evaluation on Original Dataset

Content-Based Filtering (Model 1)	Collaborative Filtering (Model 2)
<p><b>Training Details:</b></p> <ul style="list-style-type: none"> <li>Data split: 80% training, 20% testing.</li> <li>Batch size: 64.</li> <li>Optimizer: Adam, learning rate 0.001.</li> <li>Loss function: MSE.</li> <li>Training epochs: 35.</li> </ul>	<p><b>Training Details:</b></p> <ul style="list-style-type: none"> <li>Data split: 80% training, 20% testing.</li> <li>Batch size: 64.</li> <li>Optimizer: Adam, learning rate 0.001.</li> <li>Loss function: MSE.</li> <li>Training epochs: 10.</li> </ul>
<p><b>Running Time:</b></p> <ul style="list-style-type: none"> <li>Per epoch: ~0.6 minutes.</li> <li>Total: ~21 minutes.</li> </ul>	<p><b>Running Time:</b></p> <ul style="list-style-type: none"> <li>Per epoch: ~1 minute.</li> <li>Total: ~10 minutes.</li> </ul>

<p><b>Performance:</b></p> <ul style="list-style-type: none"> <li>• <b>Training Accuracy:</b> Loss: 0.0972, Accuracy: 53.93%.</li> <li>• <b>Testing Accuracy:</b> RMSE: 1.0912, Accuracy (<math>\pm 0.5</math> tolerance): 36.99%.</li> </ul>	<p><b>Performance:</b></p> <ul style="list-style-type: none"> <li>• <b>Training Accuracy:</b> Loss: 0.0743, Accuracy: 93.20%.</li> <li>• <b>Testing Accuracy:</b> RMSE: 1.0536, Accuracy (<math>\pm 0.5</math> tolerance): 39.74%.</li> </ul>
---	---

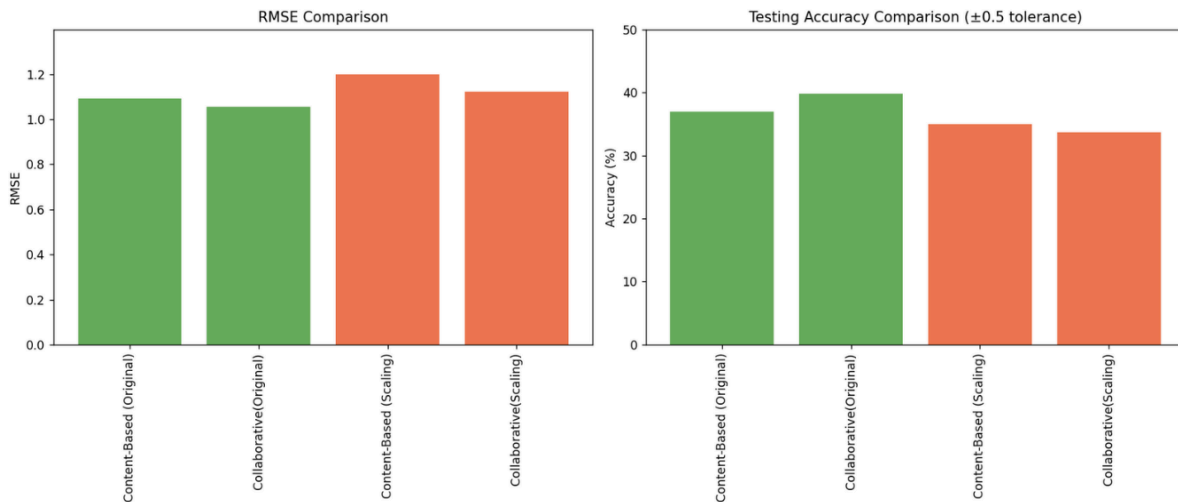
### Training and Evaluation on Scaled Dataset (Random Sampling)

Content-Based Filtering (Model 1)	Collaborative Filtering (Model 2)
<p><b>Training Details:</b></p> <ul style="list-style-type: none"> <li>• Data split: 80% training, 20% testing.</li> <li>• Batch size: 64.</li> <li>• Optimizer: Adam, learning rate 0.001.</li> <li>• Loss function: MSE.</li> <li>• Training epochs: 35</li> </ul>	<p><b>Training Details:</b></p> <ul style="list-style-type: none"> <li>• Data split: 80% training, 20% testing.</li> <li>• Batch size: 64.</li> <li>• Optimizer: Adam, learning rate 0.001.</li> <li>• Loss function: MSE.</li> <li>• Training epochs: 10.</li> </ul>
<p><b>Running Time:</b></p> <ul style="list-style-type: none"> <li>• Per epoch: ~0.5 minutes.</li> <li>• Total: ~17.5 minutes.</li> </ul>	<p><b>Running Time:</b></p> <ul style="list-style-type: none"> <li>• Per epoch: ~1 minute.</li> <li>• Total: ~10 minutes.</li> </ul>
<p><b>Performance:</b></p> <ul style="list-style-type: none"> <li>• <b>Training Accuracy:</b> Final Epoch (35/35): Loss: 0.1032, Accuracy: 55.33%.</li> <li>• <b>Testing Accuracy:</b> RMSE: 1.1986, Accuracy (<math>\pm 0.5</math> tolerance): 35.00%.</li> </ul>	<p><b>Performance:</b></p> <ul style="list-style-type: none"> <li>• <b>Training Accuracy:</b> Final Epoch (10/10): Loss: 0.0283, Accuracy: 98.79%.</li> <li>• <b>Testing Accuracy:</b> RMSE: 1.1231, Accuracy (<math>\pm 0.5</math> tolerance): 33.71%.</li> </ul>

### Discussion

- Training Performance:**
  - **Content-Based Filtering:** Shows moderate training accuracy, with slight degradation when using scaled datasets, and a higher testing RMSE (1.1986).
  - **Collaborative Filtering:** Performs significantly better in terms of training accuracy (up to 98.79%), but shows a slight decline in testing performance with scaled datasets.
- Testing Accuracy:**
  - **Content-Based Filtering:** Tends to have a higher RMSE in both original and scaled datasets, indicating less predictive power for unseen data.
  - **Collaborative Filtering:** Shows relatively consistent performance, with minor degradation under randomized sampling.
- Advantages of Scaling (Random Sampling):**
  - **Speeding Up Experimentation:** Allows faster feedback and quicker testing cycles.
  - **Ensuring Representativeness:** Reduces the risk of bias.

- **Preventing Overfitting:** Helps models generalize better by exposing them to varied data subsets.
- **Reducing Computational Cost and Time:** Enables efficient parameter tuning and model evaluation.



## 5. Conclusion & Future of the Project

**Conclusion:** The project demonstrated that Collaborative Filtering is a more effective approach for movie recommendation systems due to its ability to learn user preferences and uncover latent patterns. Collaborative Filtering's higher training accuracy and lower RMSE compared to Content-Based Filtering indicate its superior performance in providing personalized recommendations. However, Content-Based Filtering remains useful in scenarios with limited user data, such as when new movies are added to the system or when there is insufficient user interaction data. Both methods have their strengths and can be employed depending on the specific requirements and context of the recommendation system.

### Future Work:

1. **Integrate Additional Data:**
  - Include user demographics and temporal behavior to enhance recommendation quality. By incorporating data such as age, gender, and viewing history over time, the models can be further refined to offer more tailored recommendations.
2. **Incorporate More Evaluation Metrics:**
  - Use metrics like Mean Absolute Error (MAE) and precision for a comprehensive assessment. These additional metrics will provide a more nuanced understanding of model performance and help in identifying areas for improvement.
3. **Implement Dynamic Learning:**
  - Recommendations must be adapted from real-time user interactions in order to improve relevance and user satisfaction. This means the models would be continuously monitored and updated to reflect the latest preferences and behaviors of the users, ensuring the recommendation system is effective and user-centric. Further iterations may include hybrid models that combine the best aspects of both Content-Based and Collaborative Filtering models to bring the benefits of both to make the recommendations even more accurate and engaging.

**Source of Dataset:**

**<https://grouplens.org/datasets/movielens/latest/>**