#### Dana Baranova 106654

## Maria Baranovskaya 130974

## **Report on Collaborative Filtering in Recommendation Systems**

#### Introduction

Recommendation systems play an essential role in online services, helping users find products, movies, music, and other resources of interest. One of the most popular and effective techniques used to create such systems is collaborative filtering. This method recommends items to users based on the preferences of other users with similar tastes. This report discusses the basic principles of collaborative filtering, its types, methods of measuring similarity between items and users, as well as examples of its application in real-world services.

# 1. Basic Principles of Collaborative Filtering

Collaborative filtering relies on the idea that if two users have shown similar interests in the past, their preferences in the future are likely to be similar as well. The key steps in collaborative filtering are as follows:

- **Identifying Similar Users**: The system analyzes the preferences of various users and identifies those with similar tastes. In this case, one user's preferences can be used to predict what another user with similar interests might like.
- **Predicting Ratings**: Based on the analysis of similar preferences, the system predicts what rating a user might give to an item that they have not yet rated. For example, if user A rates movie X highly, and users B and C, who are similar to A, also rate this movie highly, the system predicts that user A would rate a similar movie Y highly as well.

# 2. Types of Recommendation Systems

There are several types of recommendation systems, which can be used depending on the available data and the nature of the tasks:

- Content-Based Filtering: This approach focuses on the characteristics of the items that interest the user, such as genre, author, actors, and other attributes. For instance, if a user frequently reads books on a specific topic, the system will recommend other books related to that topic.
- Collaborative Filtering: Unlike content-based filtering, this method works based on useritem interactions (ratings, purchases, etc.). For example, if a user frequently watches comedy movies, the system will recommend other movies in the same genre based on the preferences of other users with similar interests.
- **Hybrid Methods**: These methods combine content-based and collaborative filtering, enhancing the accuracy of recommendations and overcoming the weaknesses of each method. For instance, the system can consider both the attributes of items and the similarity between users.

## 3. Application of Collaborative Filtering in Real-World Services

Collaborative filtering is widely used in many popular online services. Here are some examples:

- **Netflix**: Netflix uses collaborative filtering to recommend movies and TV shows. The system analyzes the behavior of similar users and suggests content that may interest the current user. For example, if a user frequently watches comedy movies, the system will suggest other films in the same genre.
- Amazon: On Amazon, collaborative filtering helps recommend products. For example, if a user buys a book, the system may recommend items that other users who bought this book also purchased, or products that are often bought together with it.
- **Spotify**: Spotify uses collaborative filtering to create personalized playlists. The system analyzes the user's music preferences and suggests new songs, albums, or artists based on the interests of other users with similar tastes.
- **Instagram**: Instagram applies collaborative filtering to recommend content. The system analyzes the users with whom a person interacts, which posts they like, and suggests new accounts and posts that may be of interest based on the preferences of other users.

## 4. Methods of Collaborative Filtering

There are two main approaches to collaborative filtering:

• **Memory-Based**: This method relies on storing data on interactions between users and items, analyzing them to predict interests. The advantage of this approach is its simplicity and intuitiveness; however, it may be inefficient with large data volumes.

**Model-Based**: In this approach, mathematical models are created to predict preferences. These methods use more complex algorithms, such as matrix factorization and latent factor models, to handle sparse interaction matrices and predict preferences.

### 5. Advantages and Disadvantages of Collaborative Filtering Methods

### Advantages:

- **Personalized Recommendations**: Collaborative filtering provides personalized recommendations based on user behavior.
- **Discovery of Hidden Patterns**: The system can uncover users' interests and preferences that are not immediately obvious from their ratings alone.
- **Flexibility**: This method does not require information about the items themselves (e.g., movie genre), making it universally applicable to different types of data.

### **Disadvantages**:

- **Cold Start Problem**: For new users or items with insufficient data, the system cannot accurately predict interests.
- **Data Sparsity**: In real-world systems, users rarely interact with every available item, leading to sparse data.
- **Scalability Issues**: With a large number of users and items, memory-based methods may face computational difficulties.

### 6. Memory-Based Methods

To make collaborative filtering effective, it is necessary to properly measure the similarity between users and items. Several common metrics include:

• Cosine Similarity: This method measures the angle between the rating vectors of users. The cosine of the angle between two vectors A and B is calculated as:

$$\text{similarity} = \cos(\theta) = \frac{A \cdot B}{\|A\| \|B\|} = \frac{\sum\limits_{i=1}^{n} A_i \times B_i}{\sqrt{\sum\limits_{i=1}^{n} (A_i)^2} \times \sqrt{\sum\limits_{i=1}^{n} (B_i)^2}}$$

where A and B are the rating vectors of the users. The closer the value is to 1, the more similar the users' interests are.

• **Pearson Correlation**: This metric measures the linear relationship between users, ignoring their individual biases in ratings. It is calculated as:

$$r_{xy} = \frac{\sum (x_i - M_x) (y_i - M_y)}{\sqrt{\sum (x_i - M_x)^2 (y_i - M_y)^2}}$$

• **Manhattan and Euclidean Distance**: These methods measure the distance between users or items based on their ratings:

Manhattan Distance (L1 norm):

$$\operatorname{distance}(A,B) = \sum_i |R[A,i] - R[B,i]|$$

Euclidean Distance (L2 norm):

$$\operatorname{distance}(A,B) = \sqrt{\sum_i (R[A,i] - R[B,i])^2}$$

The smaller the distance value, the more similar the items or users are.

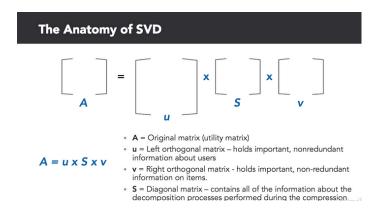
### 7. Model-Based Collaborative Filtering

For preference modeling, various algorithms can be used. The main ones include:

• Matrix Factorization: The goal of matrix factorization is to decompose the original interaction matrix R into two smaller matrices: the user matrix U and the item matrix V, such that their product approximates the original matrix:

# Examples of matrix factorization algorithms:

o **SVD (Singular Value Decomposition)** — a classical matrix decomposition method that extracts key components of user-item interactions.



- o **ALS (Alternating Least Squares)** a method that optimizes the matrices U and V alternately to minimize the prediction error.
- Latent Factor Models: These models analyze the interactions between users and items to find hidden factors that can explain user preferences. For example, latent factors might include movie genres or user preferences regarding specific attributes.

### Advantages:

- **Generalization**: These models can reveal hidden patterns in data, making them suitable for sparse matrices.
- **Scalability**: Such methods are effective for processing large datasets and can be applied to real-world services with millions of users and items.
- **Handling Latent Factors**: Complex models like neural networks can account for nonlinear dependencies between users and items.

#### **Disadvantages:**

- **High Computational Cost**: Methods like deep neural networks require substantial computational resources for training.
- **Dependence on Data Quality**: If the data is incomplete or noisy, it can negatively impact model performance.
- Cold Start Problem: New users or items with no data pose a challenge for the model, as it cannot predict their preferences accurately.

## 8. Strengths and Weaknesses of Collaborative Filtering

**Strengths**: Collaborative filtering provides personalized recommendations based on user behavior without requiring domain knowledge or item attributes. It discovers hidden patterns, adapts over time with more data, and can be applied across various domains, even leveraging data from multiple areas for cross-domain recommendations.

**Weaknesses**: It struggles with the cold start problem for new users or items, faces data sparsity issues, and has scalability challenges in large systems. It also tends to favor popular items, lacks context awareness, risks over-specialization, raises privacy concerns, and is sensitive to noise in the data.

### Conclusion

Collaborative filtering is one of the most effective and widely used methods in recommendation systems. It allows the creation of personalized recommendations based on user interaction data. Despite challenges like the cold start problem and data sparsity, the method continues to evolve, with more advanced algorithms and models becoming available each year. It is important to note that combining different approaches (e.g., hybrid methods) can significantly improve recommendation quality and address existing issues.

#### **Sources**

Bennett, J., & Lanning, S. (2007). "The Netflix Prize.": https://www.cs.uic.edu/ $\sim$ liub/KDD-cup-2007/NetflixPrize-description.pdf

Linden, G., Smith, B., & York, J. (2003). "Amazon.com recommendations: Item-to-item collaborative filtering.": https://www.cs.umd.edu/~samir/498/Amazon-Recommendations.pdf

Johnson, C. (2014). "Logistic Matrix Factorization for Implicit Feedback Data.": https://web.stanford.edu/~rezab/nips2014workshop/submits/logmat.pdf

Real Python. "Building a Recommendation Engine with Collaborative Filtering.": https://realpython.com/build-recommendation-engine-collaborative-filtering/

CodezUp. "Building Recommendation System Using Collaborative Filtering.": https://codezup.com/building-recommendation-system-collaborative-filtering/

Koren, Y., Bell, R., & Volinsky, C. (2009). "Matrix Factorization Techniques for Recommender Systems." IEEE Computer, 42(8), 30-37.