Improving Personalisation in Recommendation Engines with User-Item Interactions using Neural Collaborative Filtering, Implicit Feedback and Attention Mechanism

A thesis submitted to the University of Birmingham for the degree of MSc Data Science

Ashwin Malekar | 2587143

School of Computer Science

College of Engineering and Physical Sciences

University of Birmingham

# Abstract

This work investigates the improvement of recommender systems through the use of Neural Collaborative Filtering (NCF) models in conjunction with attention and implicit feedback conversion methods. The main goal is to employ implicit feedback obtained from user behaviour data to increase the accuracy of forecasting user behaviours, such as clicks, views, and purchases. This work highlights the significance of converting explicit user actions into implicit feedback, which better captures the complex preferences and behaviours of users than typical recommender systems that mostly rely on explicit input. As part of the methodology, raw interaction data is cleaned and processed to better reflect user intent during a thorough data pre-processing step. The development of neural network models, which use methods like negative sampling and log loss optimization, comes next. To further enhance the model's emphasis on essential user-item interactions and guarantee that the most significant aspects are given greater weight throughout the prediction process, an attention mechanism is also included. These models' efficacy is assessed in comparison to more established collaborative filtering methods by using applicable benchmarks and metrics like Hit Ratio@10 (HR). The expected results include proving that NCF models beat conventional techniques in user-item interaction prediction, especially when enhanced with implicit feedback and attention mechanisms. The study also intends to demonstrate the benefits of deeper neural network designs and show that the addition of an attention mechanism can greatly improve suggestion accuracy and provide a more customised user experience.

Keywords: Recommender Systems; Neural Collaborative Filtering (NCF); Neural Network; Attention Mechanism; Implicit Feedback; Hit Ratio@10.

# Acknowledgments

# Abbreviations

RS……………………………………...Recommender Systems

NCF…………………………………….Neural Collaborative Filtering

CF…………………………………...Collaborative Filtering

MF…………………………………..Matrix Factorization

UCF………………………………….User-based Collaborative Filtering

HR@10..............................................Hit Ratio@10

EDA………………………………….Exploratory Data Analysis

AM…………………………………...Attention Mechanism

ReLU…………………………………...Rectified Linear Unit

# Contents

# List of Figures

# List of Tables

# Introduction

Recommender systems (RS) play a critical role in improving user experiences on a variety of platforms, including streaming services, by accurately identifying and recommending goods that users are likely to interact with. However, typical collaborative filtering algorithms, which rely on explicit input like user ratings, often face severe challenges due to the inherent limitations of sparse and limited data. Insufficient user-item interactions caused by data sparsity make it challenging for systems to anticipate user preferences with any degree of accuracy. When handling the cold-start issue, new users or items have little to no prior data, this becomes extremely challenging and makes it more difficult to create recommendations that are both relevant and tailored. Implicit feedback adds complexity to modelling and interpretation, but it also offers a deeper source of user interaction data—views, clicks, or any kind of interaction with the item, for example. This study uses the Neural Collaborative Filtering (NCF), *He et al. (2017),* framework to forecast user actions in recommender systems more accurately, in an attempt to overcome these obstacles. NCF models can identify non-linear patterns in user-item interactions that standard approaches might miss by utilising deep learning techniques. The MovieLens 20M Dataset is the main source of explicit feedback data used in this project, which is transformed into implicit feedback in the form of clicks, views. Through an examination of the relationships between "*userId*" and "*movieId,*" the study seeks to determine how well NCF models perform in terms of improving user satisfaction and recommendation accuracy when contrasted with traditional collaborative filtering methods. Additionally, this research introduces an attention mechanism into the NCF framework to further enhance the model's capability to focus on the most important user-item interactions. The attention mechanism allows the model to prioritise certain interactions that are more indicative of user preferences, potentially leading to more personalised and accurate recommendations. The overarching goal of this research is to improve the accuracy and personalization of recommendation engines through the effective use of implicit feedback data and attention mechanisms.

## 1.1   Research Motivation

Recommender systems must advance beyond conventional techniques due to the growing complexity and volume of user interactions on digital platforms. Although collaborative filtering has proven fundamental, it faces several challenges, such as data sparsity and a strong reliance on explicit feedback that frequently falls short of capturing the entire range of user preferences. The promise of Neural Collaborative Filtering (NCF) models to overcome these drawbacks by utilising deep learning to reveal more complex patterns in user behaviour is what spurs this study.

To fully exploit the potential of NCF models, however, merely implementing them is insufficient. This study incorporates a multi-head attention mechanism that enables the model to rank the most pertinent interactions in order to improve their efficacy. This strategy is especially relevant when user interactions become more varied and intricate, requiring more advanced technologies in order to provide reliable recommendations.

Considering these factors, the main research question is: *How to improve personalisation and handle complex user-item interactions considering data sparsity and cold start problem in recommendation systems?* The answer to this topic is essential for improving user experiences on a variety of digital platforms as well as for the advancement of recommender system technology.

## 1.2    General Limitations of this Research

This work has limitations even though it uses implicit feedback, attention mechanisms, and NCF models to advance the field of recommender systems. A notable limitation is the dependence on the MovieLens 20M Dataset, which, although its extensive scope, might not accurately capture the variety of user interactions prevalent in different domains or datasets. Furthermore, there is an inherent interpretive component to the conversion of explicit feedback into implicit feedback, which could lead to biases or inaccurate representations of actual user preferences. Even though the used deep learning models are strong, they also demand a lot of computational capacity and could be difficult to scale in practical applications. Lastly, although informative, the offline evaluation criteria used in this study may not accurately reflect user satisfaction in a live system.

## 1.3    Aim and Objectives

The aim of this work is to design and assess NCF models that efficiently leverage implicit feedback and attention mechanisms, with the goal of improving recommender system accuracy and personalization. In comparison to conventional collaborative filtering techniques, this study aims to show how incorporating these advanced methods can greatly enhance the prediction of user activities, such as clicks, views, resulting in more precise and targeted suggestions.

- To investigate and translate explicit user feedback from the MovieLens 20M 'rating' Dataset into implicit signs of feedback, like views, clicks, and movie ratings.
- To build and implement NCF framework that captures the complex relationship between users and items utilising deep learning approach.
- To integrate an attention mechanism within the NCF framework to prioritize and enhance important user-item interactions.

- To evaluate the performance of the proposed models with unseen (test) data using metrics such as Hit Ratio @10 (HR@10) and compare with traditional collaborative filtering models.
- To gauge the effect of attention mechanism and implicit feedback conversion on the recommender system's overall accuracy and user satisfaction.
- To outline some obstacles and restrictions when using NCF models with attention mechanisms in practical settings and make recommendations for further study.

## 1.4 Thesis Overview

The methodology as shown in Figure 1. 1, outlines the structured approach taken in this thesis to enhance personalization in recommendation systems by addressing problems like data sparsity and the cold start problem. The process begins with data pre-processing, which includes steps like data cleaning, exploratory data analysis (EDA), and the conversion of explicit feedback to implicit feedback. The dataset used is the MovieLens 20M dataset, which is split into training and testing sets.

Following data preparation, the model building & training phase is initiated, where a Neural Collaborative Filtering (NCF) model with an integrated attention mechanism is developed. This involves embedding user and item data, applying the attention layer to focus on significant features, and utilizing fully connected (ReLU) layers to capture complex interactions. The trained model is then tested and validated.

Finally, the Evaluation phase involves comparing the proposed model's performance against traditional collaborative filtering methods using the Hit Ratio @10 metric. The structured approach ensures a comprehensive analysis and validation of the proposed model's effectiveness in improving recommendation accuracy.

*Figure 1. 1 Thesis Methodology*

## 1.5  Legal, Social, Ethical, and Professional Issues

The development of recommender systems involves a number of legal, social, ethical, and professional considerations. Legally, this study has to abide by data privacy requirements such as GDPR, guaranteeing that user data is managed sensibly and compliantly with applicable legislation. Socially, efforts are undertaken to ensure fairness and diversity in recommendations, acknowledging the potential for recommender systems to reinforce prejudices or create echo chambers. The study considers the ethical dangers of taking advantage of users or influencing their choices, stressing the significance of openness and regard for user autonomy, particularly when employing implicit feedback. The study complies with strict guidelines for academic integrity, guaranteeing that the results advance the discipline in a morally and constructive way. These factors guarantee that the research is socially and legally accountable, as well as technically robust.

## 1.6  Thesis Structure

Chapter 1 (Introduction) provides a detailed overview of the research problem, including background information, research motivation, and the significance of the study, along with the aim, objectives, and a summary of the thesis structure. Chapter 2 (Background) delves into the history and challenges of the topic,

discussing evaluation metrics and techniques used in the study. Chapter 3 (Literature Review) examines existing research on recommender systems, implicit feedback, and the application of deep learning and attention mechanisms. Chapter 4 (Data) details data sourcing, collection, exploratory data analysis (EDA), and the pre-processing and feature engineering performed for model compatibility. Chapter 5 (Methodology) outlines the research methodology, including model development, data pre-processing, and experimental design. Chapter 6 (Results and Discussion) presents the research findings and interprets them in light of the study's objectives and existing literature. Chapter 7 (Conclusion and Future Work) summarizes the key contributions, discusses limitations, and challenges, and suggests potential future research directions. The thesis concludes with References, listing all cited materials, and an Appendix that includes supporting materials such as data, code, and additional results.

To sum up, the importance of improving recommender systems using cutting-edge methods like neural collaborative filtering, implicit feedback, and attention processes has been demonstrated in the Introduction chapter. Setting the foundation for the in-depth investigation that comes next, it has described the goals and objectives of the research as well as the more general considerations.

# Background

## 2.1 Background of the topic

### 2.1.1 Recommender Systems

Recommender systems have become an essential element of contemporary digital experiences, offering consumers personalised content and product recommendations across a variety of platforms, like social media, streaming, and e-commerce. The primary goal of these systems is to anticipate user preferences and recommend items that are likely to be interacted with, thereby increasing user satisfaction and engagement.

### 2.1.2 Collaborative Filtering – Matrix Factorization and User-based methods

Collaborative filtering (CF) has been the conventional method of developing recommender systems. This approach is predicated on the similarities between users and items in order to produce recommendations.

CF algorithm is used to find the item likeliness, as shown in Figure 2. 1 in a detailed view. The algorithm first predicts $P_{aj}$ which is the predicted score of an item 'j' for user 'a', then, we simply recommend a list of *Top-N* items that this active user 'a' will like the most. *(Bokde, Girase, and Mukhopadhyay, 2015)*
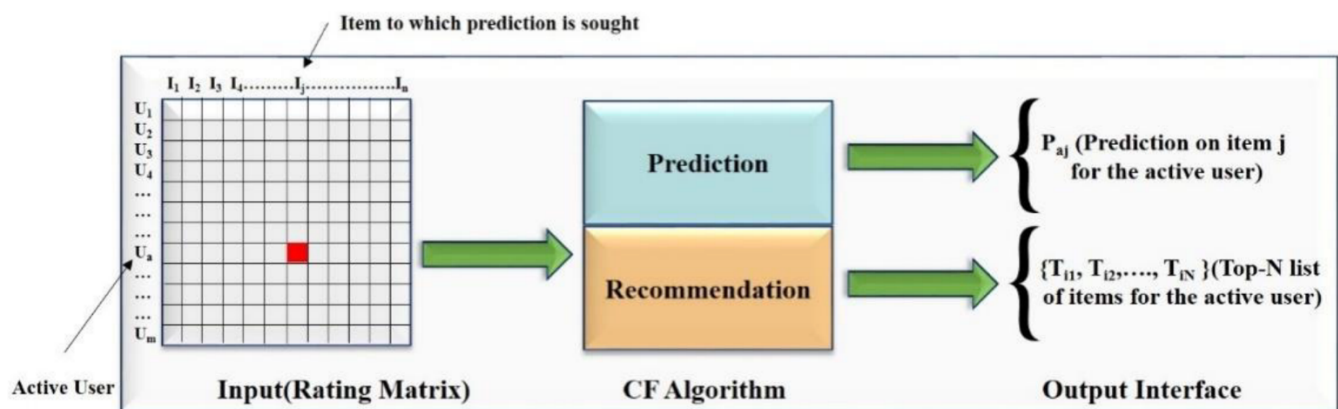


*Figure 2. 1 Collaborative Filtering, (Bokde, Girase, and Mukhopadhyay, 2015)*

User-based and item-based are the two primary categories into which CF techniques are typically classified. User-based CF predicts a user's interest in an item based on the preferences of other users with similar patterns:

$$\hat{R}_{ui} = \bar{R}_u + \frac{\sum_{v \in Neighbours(u)} sim(u, v) \times (R_{vi} - \bar{R}_v)}{\sum_{v \in Neighbours(u)} |sim(u, v)|}$$

where,

$\bar{R}_u$ is the average rating of user u,

$sim(u, v)$ is the similarity between user u and v,

$R_{vi}$ is the rating of user v for item i

To address some limitations of traditional CF, Matrix Factorization (MF) techniques, such as Singular Value Decomposition (SVD), have been widely adopted. MF decomposes the user-item interaction matrix into latent factors, enabling the system to understand hidden patterns and relationships between users and items. However, while MF has been effective in improving recommendation accuracy, it struggles to model complex, non-linear interactions between users and items. The linear nature of MF means it can miss more subtle patterns in the data, which limits its ability to fully capture the richness of user preferences.

### 2.1.3 Data Sparsity and Cold Start Problem in CF methods

Despite their extensive use, traditional CF approaches, including MF, suffer considerable hurdles, particularly data sparsity and the cold start issue. Data sparsity arises when there is insufficient user-item interaction data, resulting in a sparse user-item matrix where most entries are empty. This sparsity significantly hampers the performance of CF algorithms, as they rely on the availability of substantial interaction data to make accurate predictions. Consequently, the recommendations may be inaccurate or overly generalised.

The cold start problem further exacerbates these challenges, especially for new users or items that have little to no historical interaction data. In such cases, CF techniques, including MF, struggle to provide meaningful recommendations due to the lack of sufficient data points to base their predictions on. These limitations highlight the need for more advanced models that can better handle data sparsity and capture non-linear relationships in user-item interactions, leading to more personalised and accurate recommendations.

### 2.1.4 Explicit and Implicit Feedback

An essential aspect of recommender systems is the type of feedback they utilize to model user preferences. Explicit feedback refers to direct user input, in this case, the ratings, where users explicitly express their preferences. In contrast, implicit feedback is derived from user behaviour, such as clicks, views, or purchase history, which indirectly indicates user preferences. While explicit feedback is often more straightforward to interpret, it is typically less abundant than implicit feedback, which is naturally generated as users interact

with digital platforms. This makes implicit feedback particularly valuable, as it can provide a more comprehensive picture of user behaviour.

In this research, the MovieLens dataset, which primarily contains explicit ratings given by users, will be leveraged. However, to fully exploit the potential of the dataset and address the limitations of sparse and limited explicit data, these ratings will be converted into implicit feedback signals, such as views, clicks, and inferred interest levels. This conversion is crucial because implicit feedback is often more abundant and can help mitigate the data sparsity and cold start problems that frequently hinder collaborative filtering approaches. By transforming explicit ratings into implicit signals, the recommender system can better capture the subtle, non-linear interactions between users and items, leading to more accurate and personalized recommendations.

## 2.2 Utilised Algorithms and Metrics

### 2.2.1 Neural Collaborative Filtering

To address these constraints, the area has increasingly relied on advanced machine learning approaches, specifically deep learning, which can model complex and nonlinear user-item interactions. Neural Collaborative Filtering (NCF) as shown in Figure 2. 2 from *He et al. (2017)*, is a technology that uses neural networks to increase suggestion accuracy.



*Figure 2. 2 Neural Collaborative Filtering Framework (He et al., 2017)*

NCF models combine the characteristics of matrix factorization (a prominent latent factor model in CF) and deep neural networks, allowing for the detection of subtle patterns in user-item interactions that older

approaches may miss. Recent advances in the field have centred on hybrid models that incorporate both explicit and implicit feedback to improve the resilience and accuracy of suggestions. Furthermore, the use of knowledge graphs and advanced embedding methods has enhanced recommender system performance by giving more contextual details and better handling of sparse data.

The current research and development in deep learning-based recommender systems aims to address the remaining issues and push the limits of these systems' capabilities. The integration of complex models and different data sources has the potential to give highly tailored and accurate suggestions, considerably improving user experience across several digital or streaming platforms.

### 2.2.2 Attention Mechanism (AM)

To further enhance the performance of the recommender system, an attention mechanism (AM), as described in Figure 2. 3 *Vaswani et al. (2017)*, will be integrated into the model. AM is a powerful tool that allows the model to selectively focus on the most relevant aspects of the user-item interactions.



*Figure 2. 3 Multi-head Attention Mechanism, (Vaswani et al., 2017)*

In the context of recommender systems, the attention mechanism can identify which past interactions are most indicative of a user's current preferences and weigh them more heavily when generating recommendations. This approach not only improves the ability of the model to capture complex, non-linear relationships but also enhances the interpretability of the recommendations by highlighting the key factors driving each suggestion. By incorporating an AM, this research aims to push the boundaries of recommendation accuracy and personalization, offering users more tailored and contextually relevant suggestions.

### 2.2.3 Hit Ratio@10 (HR@10) as a Metric

*Alsini, Huynh, and Datta (2020),* proposed Hit Ratio, which is the performance metric that measures how well the recommender system performs.

HR@10 for a single user *u,* can be defined as:

$$HR@10(u) = \begin{cases} 1 & if\ T_u \in R_u \\ 0 & if\ T_u \notin R_u \end{cases}$$

The overall HR@10 is then averaged over all users:

$$HR@10 = \frac{1}{|U|} \sum_{u \in U} HR@10(u)$$

where,

$|U|$ = total number of users,

$\sum_{u \in U} HR@10(u)$ = sum of hit ratios for all users

It basically selects 99 items at random that the user hasn't interacted with. After that, we add the user's interacting item to these 99 examples, totalling 100 items. We rank these 100 items according to the model's predicted probabilities after running the model on them. If the test item is among the top 10 items chosen based on rank, it is considered successful. The average hits are used to calculate the Hit Ratio, and this process is repeated for each user.

# Literature Review

## 3.1 Related Work

Integrating deep learning—more specifically, neural collaborative filtering, or NCF—into recommender systems has shown to be a dependable technique for predicting user behaviour, particularly when implicit input is included. This review of the literature describes the advancements and methods for using Deep Learning models, including the findings from five significant studies on the topic. For instance, *He et al. (2017)*, suggested a universal architecture that combines matrix factorization (MF) with multi-layer perceptrons (MLP) to portray complicated user-item interactions. The NCF framework was built upon this architecture. Because NCF can learn non-linear interaction functions, their study shows that it is more effective at handling implicit feedback than more traditional MF strategies. In the paper, *Zhang et al. (2016)*, added more layers and embedding techniques to the NCF model to enhance learning of user-item interactions. This article demonstrates that deeper network architectures can dramatically improve recommendation accuracy, underscoring the importance of non-linear transformations in capturing complicated user preferences. In the paper *Guo (2012)*, investigates the persistent issues in traditional collaborative filtering (CF) systems related to data sparsity and cold start. The study claims that these problems are caused by the limited quantity of user reviews, which makes it more challenging to offer accurate recommendations for new clients or goods. There are several solutions available, however these issues are still not well addressed in current CF techniques. *Shi, Larson, and Hanjalic (2014)*, provide a comprehensive review of collaborative filtering (CF) that extends beyond the traditional user-item (U-I) matrix. They also highlight the cold start and data sparsity problems. *Torkashvand, Jameii, and Reza (2023)*, investigate the value of deep learning-based collaborative filtering recommender systems in addressing the cold start problem. They show how applying neural network architectures and implicit feedback to recommender systems improves prediction accuracy and user interaction modelling. The research paper, in which *Jiang et al. (2020)*, present approaches that are essential for predicting individualized user-item interactions. It utilizes AutoEncoder (AE) and Latent Factor Analysis (LFA) models to handle high-dimensional, sparse data, highlighting the use of known ratings to improve computational efficiency and prevent data skewing. Moreover, the model's accuracy and scalability are greatly improved by integrating deep structures and parallelization techniques. These approaches are highly relevant to creating neural collaborative filtering models with implicit feedback. The research paper *Rodpysh, Mirabedini, and Banirostam (2021)*, which provides an overview of relevant approaches to address cold start and sparse data issues, makes use of a multi-level singular value decomposition (SVD) technique that combines user and item context information. The suggested method efficiently captures contextual similarities and raises recommendation accuracy by generating user context feature matrices (UCFM), item context

feature matrices (ICFM), and context similarity matrices (CSM). The momentum stochastic gradient descent included in this method improves forecast accuracy and efficiency even further. *Hernando et al. (2017)*, present methods for forecasting customized interactions between users and items. It focuses on unregistered users that face the cold-start problem and simulates forward reasoning for suggestions using graphical probabilistic models and uncertainty criteria. This method lets people view suggestion paths using item-based trees and deduce their own preferences. Personalized recommendation engines can benefit greatly from neural collaborative filtering models that incorporate implicit feedback, as these probabilistic modelling approaches prioritize interpretability and user participation. Methodologies relevant to predicting unique user-item interactions in recommendation engines are introduced in the research paper, *Natarajan et al. (2020)*. In order to solve the cold start problem and enhance matrix factorization methods for data sparsity, the paper suggests utilizing Linked Open Data (LOD). The model improves recommendation accuracy by including semantic data from LOD into the matrix factorization procedure. It combines collaborative filtering with LOD and matrix factorization. An innovative method for predicting individualized user-item interactions is presented in the research paper, *Zhou et al. (2020)*. It uses autoregressive flows and Bayesian inference to develop the Collaborative Autoregressive Flows (CAF) model, which uses invertible transformations to change basic distributions into complicated ones. This methodology addresses the bias in current Bayesian recommendation models by enabling flexible and tractable probabilistic density estimates. CAF greatly increases recommendation accuracy by leveraging implicit feedback and optimizing latent factor representation, which is in line with the goals of neural collaborative filtering models that use implicit feedback in personalized recommendation engines. *Xia, Luo, and Liu, (2021)*, integrated attention mechanisms with Neural Collaborative Filtering (NCF) enhances the model's ability to dynamically weigh the importance of various latent factors, leading to improved accuracy in predictions based on explicit feedback such as user ratings. This approach effectively addresses challenges like data sparsity and the cold start problem, resulting in a more robust recommendation system.

Even However, previous research has not adequately tackled the drawbacks of conventional collaborative filtering models, especially when it comes to managing data sparsity and the cold-start issue.

## 3.2 Knowledge Gap

This work seeks to address this *knowledge gap* by *improving the personalization and robustness of user-item interaction predictions with the help of NCF model, implicit user feedback and attention mechanism*. Previous studies demonstrate notable progress in recommendation systems, especially when deep learning methods like Neural Collaborative Filtering (NCF) are used. Nonetheless, a significant body of knowledge has to be added to adequately tackle the enduring issues of data sparsity and cold start concerns in collaborative filtering systems. Although early NCF models and conventional CF methods have demonstrated potential, they

frequently suffer from a lack of user interaction data, which makes it challenging to produce precise suggestions for new users or objects. Moreover, a lot of these methods fall short in utilizing implicit feedback's full potential, which can offer more detailed insights into user preferences. By improving the personalization and resilience of user-item interaction predictions, our study seeks to close this gap. Our study aims to narrow this gap by enhancing the personalization and resilience of user-item interaction predictions.

# Data

## 4.1 Data Source & Data Collection

For this research, the primary data source utilised is the MovieLens 20M dataset, a widely recognized dataset in the field of recommender systems, made available by the GroupLens Research Lab at the University of Minnesota. The MovieLens dataset is a rich resource that includes user ratings, movie information, and various other metadata, making it an ideal choice for developing and evaluating recommender system models.

The MovieLens 20M dataset contains approximately 20 million ratings applied to over 27,000 movies by 138,000 users. The dataset is publicly accessible and can be sourced directly from the official MovieLens website MovieLens Dataset. The ratings are collected on a scale from 0.5 to 5.0, with increments of 0.5, representing explicit feedback from users on the movies they have watched. Each rating is associated with a specific user and a movie, along with a timestamp indicating when the rating was provided.

Given the large size of the MovieLens 20M dataset, and to manage computational resources effectively on my local machine, a decision was made to work with a subset of the data. Specifically, a random sample comprising 40% of the users from the entire dataset was extracted, as shown in Figure 4. 1. This subset was chosen to ensure that the research remains manageable in terms of processing time and memory usage while still providing a sufficiently large and representative sample of the overall dataset.

```
[10]:  #Querying for collecting random 40% user data.
       rating_df = rating_df.loc[rating_df['userId'].isin(random_userIds)]

       print('Current shape: {} rows of data from {} users'.format(len(rating_df), len(random_userIds)))
       Current shape: 8022066 rows of data from 55397 users
```

*Figure 4. 1 Collect data from 40% of Random Users*

## 4.2 Exploratory Data Analysis

Exploratory Data Analysis (EDA) is a crucial step in the data science process that involves summarizing and visualizing the main characteristics of a dataset, often using graphical techniques. EDA helps in understanding the underlying structure, patterns, and relationships within the data, as well as identifying any anomalies or outliers. In the context of this research, EDA is particularly important as it provides valuable insights into the distribution of user ratings, the frequency of interactions, and the overall behaviour of users within the MovieLens dataset. By performing EDA, we can ensure that the data is well-understood and appropriately

prepared for the development and evaluation of recommender system models, ultimately leading to more accurate and effective recommendations.

## 4.2.1 Distribution of Ratings

The bar chart in Figure 4. 2, represents the distribution of user ratings across different rating values within the MovieLens 20M dataset.



*Figure 4. 2 Distribution of User Ratings*

The skewed distribution towards higher ratings suggests that the dataset is imbalanced. This imbalance could impact the performance of recommendation algorithms, potentially causing them to overestimate user satisfaction or predict higher ratings more frequently. The predominance of high ratings may lead to a lack of diversity in the recommendations if the system heavily relies on these ratings. Understanding this distribution is crucial for converting explicit ratings into implicit feedback. For instance, high ratings might be interpreted as strong positive signals, while lower ratings, despite being fewer, could still provide valuable information about user preferences that should not be overlooked.

## 4.2.2 Average Rating per Movie

The histogram, in Figure 4. 3, illustrates the distribution of average ratings for movies in the MovieLens 20M dataset. Understanding the distribution of average movie ratings can assist in normalising rating data, ensuring that the recommender system does not bias towards movies with unusually high or low average ratings.

*Figure 4. 3 Average User Rating per Movie*

The spread of ratings indicates varying user preferences, which emphasises the need for personalization in the recommendation algorithms to cater to diverse tastes. The concentration of ratings in the 3.0 to 4.0 range might reflect a selection bias where users prefer to watch and rate movies, they believe they will enjoy, thus avoiding movies with historically lower ratings.

For recommender systems, this might necessitate exploring techniques like diversity enhancement to ensure users are exposed to a broader range of movies, which could improve user engagement and satisfaction by discovering hidden gems.

### 4.2.3 Correlation between Movie Popularity and Average Rating

The scatter plot, Figure 4. 4, visualises the relationship between movie popularity (as measured by the number of ratings) and the average rating received by each movie in the MovieLens 20M dataset.

The absence of a strong correlation between popularity and average rating highlights the importance of considering both metrics separately in recommendation algorithms. Relying solely on popularity might cause the system to miss out on recommending high-quality but less popular movies.

*Figure 4. 4 Movie Popularity and Average Rating Correlation*

Recommender systems should be designed to balance the recommendation of popular content with the discovery of lesser-known, highly-rated movies to provide a more diverse and satisfying user experience.

The clustering of ratings for popular movies might introduce a bias in recommendation systems that favour these movies, leading to a "rich-get-richer" effect where popular movies become even more popular. This can reduce the diversity of recommendations.

Strategies such as incorporating diversity in recommendation model or focusing on user-specific interests rather than global popularity could help mitigate this bias and ensure a more personalised experience.

### 4.2.4 Density Plots: Number of Ratings Per User and Per Movie

The density plots, Figure 4. 5, illustrate the distribution of the number of ratings per user and per movie in the MovieLens 20M dataset. The left plot represents the number of ratings given by each user, while the right plot represents the number of ratings received by each movie.

The skewness in both plots highlights the issue of data sparsity, where most user-movie interactions are missing. This sparsity poses challenges for traditional collaborative filtering methods, which rely on having sufficient data points to make accurate predictions. For users who have rated very few movies, it becomes difficult to accurately model their preferences, leading to the cold start problem. Similarly, movies with very few ratings may not be recommended as frequently, potentially limiting their exposure to users.

17

*Figure 4. 5 Density Plots of Number of Ratings (a)Per User and (b)Per Movie (item)*

The long tail in both distributions suggests that the recommender system needs to be carefully designed to handle this imbalance. Algorithms that can generalize well from sparse data, such as those incorporating implicit feedback or employing matrix factorization with regularization, may perform better in this context. The recommender system should also account for the popularity bias, where movies with many ratings might dominate the recommendations. Introducing mechanisms to promote diversity and novel content could help mitigate this effect.

## 4.2.5 Discovering Sparsity in User-Item Matrix

The plot, Figure 4. 6, visualises the sparsity of the user-item interaction matrix in the MovieLens dataset. The matrix is overwhelmingly sparse, with a very small fraction of the possible user-item pairs having actual interactions. This is visually evident as the vast majority of the matrix is white (indicating no interaction), with only a few scattered blue dots representing ratings.



*Figure 4. 6 User-Item Sparsity*

The matrix is overwhelmingly sparse, with a very small fraction of the possible user-item pairs having actual interactions. This is visually evident as the vast majority of the matrix is white (indicating no interaction), with only a few scattered blue dots representing ratings.

The vertical and horizontal streaks of dots suggest that there are certain items that are significantly more popular or certain users who are more active in rating, contributing to higher density in specific parts of the matrix.

The high level of sparsity in the user-item interaction matrix poses a significant challenge for traditional collaborative filtering methods, which rely on shared user-item interactions to make recommendations. With so few data points available, it becomes difficult to find meaningful similarities between users or items, leading to less accurate recommendations.

The sparsity can also exacerbate the cold start problem, where new users or items with few interactions struggle to be accurately integrated into the recommendation system. This makes it harder to generate relevant recommendations for new users or items that have not yet accumulated sufficient interaction data. The use of an attention mechanism in neural collaborative filtering models can help focus on the most relevant user-item interactions, making better use of the sparse data available.

## 4.3 Data Pre-processing and Feature Engineering

In preparation for developing a robust recommender system, a series of data pre-processing and feature engineering steps were undertaken to ensure the dataset was well-suited for the task. These steps are crucial for transforming raw data into a format that can be effectively utilised by the recommendation models, particularly when dealing with large-scale datasets like MovieLens 20M.

### 4.3.1 Data Sampling

Given the size of the MovieLens 20M dataset, which contains approximately 20 million ratings, working with the full dataset on a local machine could be computationally challenging. Therefore, a random sampling technique was employed to retain 40% of the user data from the full dataset. This sampling was performed to maintain a manageable dataset size while still ensuring that the subset was representative of the entire dataset's diversity and distribution. The selected subset retained the essential characteristics necessary for training and testing the recommendation models effectively.

### 4.3.2 Splitting the Data: Leave-One-Out Method

To evaluate the performance of the recommender system, the dataset was split into training and testing sets using the leave-one-out (LOO) method. In this approach, the most recent rating (based on the timestamp) from each user was left out of the training set and used as the test set. The advantage of this method is that it closely simulates a real-world scenario, where the model is required to predict a user's next interaction based on their historical behaviour. By leaving out the most recent interaction, the model's ability to generalise to unseen data is tested more rigorously, providing a more accurate assessment of its performance.

### 4.3.3 Converting Explicit Feedback to Implicit Feedback

The MovieLens dataset primarily consists of explicit feedback in the form of user ratings, which range from 0.5 to 5.0. However, to train a recommender system to increase personalisation and handle complex non-linear relationships, which predicts user interactions rather than ratings, this explicit feedback was converted into implicit feedback. The conversion process involved binarizing the ratings so that any rating was assigned a value of '1', indicating that the user interacted with the item. This transformation shifts the focus from predicting the exact rating a user might give to a movie, to predicting whether the user is likely to interact with a movie at all. This reframing is crucial because it aligns with the practical goal of recommending items that a user is most likely to engage with, which is often more relevant in commercial and real-world applications.

### 4.3.4 Negative Sampling

After binarizing the dataset, a challenge arose: all the samples now belonged to the positive class, representing movies that users had interacted with. To effectively train the model, it was necessary to introduce *negative samples*, in Figure 4. 7, which represent movies that users have not interacted with and, by assumption, are not interested in. While this assumption may not always hold true (as users may simply be unaware of certain movies), it is a common practice in recommender system development that generally yields good results.

|   | userId | movieId | interacted |
|---|--------|---------|------------|
| 0 | 5070 | 1270 | 1 |
| 1 | 5070 | 2297 | 0 |
| 2 | 5070 | 108501 | 0 |
| 3 | 5070 | 126591 | 0 |
| 4 | 5070 | 69278 | 0 |

*Figure 4. 7 Negatively Sampled Data Frame*

To generate negative samples, a technique called *negative sampling* was employed. For each positive sample in the dataset, four negative samples were generated, maintaining a ratio of 4:1 for negative to positive samples. This method ensures that the model learns to distinguish between movies that a user is likely to interact with and those they are not, improving the overall accuracy and relevance of the recommendations.

The data pre-processing and feature engineering steps outlined above were essential in preparing the MovieLens dataset for effective modelling. By sampling the data, splitting it in a way that simulates real-world usage, converting explicit feedback into implicit feedback, and performing negative sampling, the dataset was transformed into a format that facilitates the training of a robust and realistic recommender system. These steps lay the foundation for accurate and personalised movie recommendations, which will be evaluated in the subsequent stages of the research.

# CHAPTER 5

# Methodology

This chapter delves into the architecture of each model used in the research. The models include a traditional user-based collaborative filtering model, a matrix factorization model, a neural collaborative filtering model, and a proposed NCF model augmented with an attention mechanism. Each model was carefully designed to capture different aspects of user-item interactions, allowing for a robust comparison of performance across varying methodologies.

## 5.1 Defining Model Architecture

### 5.1.1 Neural Collaborative Filtering – Model Architecture

The Neural Collaborative Filtering (NCF) model, in Figure 5. 1 *(He et al., 2017)*, is a significant evolution in recommendation systems, blending the power of neural networks with collaborative filtering techniques to capture complex, non-linear interactions between users and items. The architecture of the NCF model, as depicted in the provided diagram, is designed to leverage user and item embeddings in a neural network framework, allowing for the learning of more intricate relationships that traditional methods might overlook.
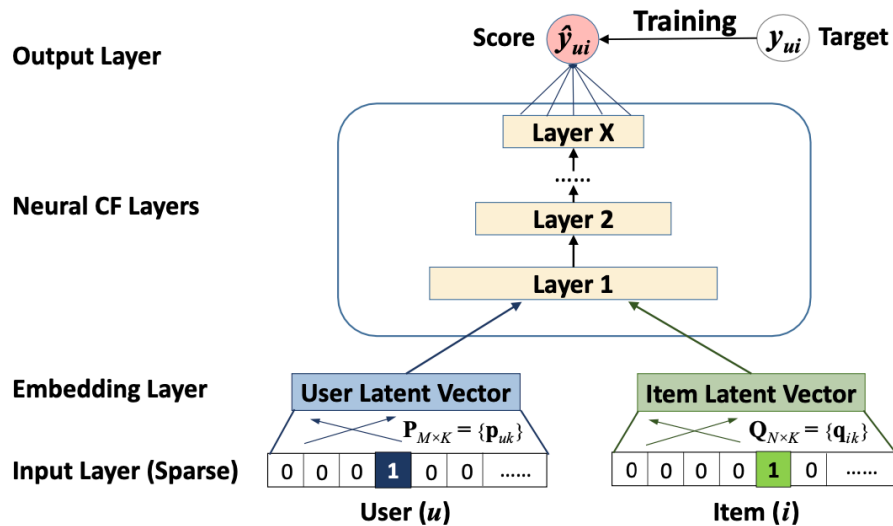


*Figure 5. 1 Neural Collaborative Filtering, (He et al., 2017)*

The NCF model is composed of several key layers that work together to transform user and item inputs into a predictive score indicating the likelihood of user interaction with an item. The model's architecture can be broken down into the following components:

**Input Layer**

- The input to the NCF model consists of user IDs and item (movie) IDs. These IDs are initially represented as sparse one-hot encoded vectors. For instance, a user ID like `userId 3` might be represented as `[0, 0, 1, ..., 0]`, where the position of the `1` indicates the specific user. Similarly, a movie ID like `movieId 1` is represented as `[1, 0, 0, ..., 0]`.

**Embedding Layer**

- The first transformation applied to these input vectors is through embedding layers. An embedding layer maps the high-dimensional sparse input vectors into lower-dimensional dense vectors, capturing latent features that represent underlying characteristics of users and items.

- For example, the user embedding layer transforms `userId 3` into an embedded user vector such as `(0.52, 0.18, 0.68, …, 0.22)`, while the item embedding layer similarly transforms `movieId 1` into an embedded item vector such as `(0.41, 0.74, 0.17, …, 0.27)`. These embeddings are the cornerstone of the NCF model, as they encapsulate the latent factors that drive user preferences and item attributes.

**Concatenation Layer**

- Once the user and item embeddings are obtained, they are concatenated into a single vector. This concatenated vector combines the information from both the user and the item, creating a unified representation that the subsequent layers of the network can use to learn interactions. The concatenated vector might look like `((0.52, 0.18, 0.68, …, 0.22), (0.41, 0.74, 0.17, …, 0.27))`, effectively merging the latent features of the user and the item.

**Fully Connected Layers**

- The concatenated vector is then passed through a series of fully connected (dense) layers. These layers apply non-linear transformations to the input vector, enabling the model to learn complex patterns and interactions that might not be apparent from the original embeddings. Each layer in the network applies a set of weights and a bias to the input, followed by an activation function such as ReLU (Rectified Linear Unit) to introduce non-linearity.

- These layers are crucial for the model's ability to generalize from the training data to unseen user-item pairs, as they allow the model to capture higher-order interactions between users and items that go beyond simple linear relationships.

**Output Layer**

- The final layer in the NCF model is the output layer, which produces a prediction vector, typically a single value, that represents the predicted interaction between the user and the item. This value can be

interpreted as the probability that the user will interact with the item, such as clicking, watching, or purchasing.

- The output is often passed through a sigmoid activation function to map it to a range between 0 and 1, making it suitable for interpretation as a probability.

The model is trained using a dataset of user-item interactions, where each instance in the dataset corresponds to a specific user-item pair and the associated interaction (e.g., rating or click). During training, the embeddings and network weights are iteratively updated using optimization techniques like stochastic gradient descent (SGD) or Adam, which adjust the parameters to reduce the prediction error.

One of the key advantages of the NCF model is its flexibility in capturing a wide range of interaction patterns. By adjusting the depth and width of the fully connected layers, the model can be tuned to balance complexity with generalization, making it suitable for various recommendation scenarios.

The NCF model's architecture is highly adaptable, allowing it to be used in different recommendation contexts. For instance, it can be employed in implicit feedback scenarios where the goal is to predict whether a user will interact with an item, rather than predicting explicit ratings. The ability to learn non-linear interactions between users and items gives the NCF model a significant edge over traditional matrix factorization techniques, especially in cases where user preferences are complex and not easily captured by simple linear models.

## 5.1.2 Proposed Model – NCF + Attention Mechanism

The proposed model architecture, in Figure 5. 2, integrates an attention mechanism into the Neural Collaborative Filtering (NCF) framework, enhancing its ability to capture complex interactions between users and items. This fusion of NCF and attention mechanisms is designed to focus the model's capacity on the most relevant aspects of user-item interactions, thus improving the overall accuracy and personalization of recommendations.

The architecture of the proposed model can be dissected into several critical layers and components, each playing a specific role in processing the input data and generating predictions. The diagram provided illustrates this layered approach, starting from sparse user and item inputs, through embedding layers, an attention layer, and finally fully connected layers leading to the output prediction vector.

*Figure 5. 2 Proposed Model Architecture (NCF + Attention Mechanism)*

**Input Layer**

- The model begins with sparse inputs for both users and items. The user ID and item (movie) ID are represented as one-hot encoded vectors, where each vector has a length equal to the number of unique users and items in the dataset, respectively. For instance, a user ID like `userId = X` might be represented as `[0, 1, 0, ..., 0]`, and an item ID like `movieId = Y` as `[0, 0, 0, ..., 1]`. These inputs serve as the foundation for the embedding process.

**Embedding Layer**

- The first major transformation occurs in the embedding layers. Each user and item ID is mapped to a dense, lower-dimensional vector through the embedding process. These embeddings capture latent features that abstractly represent user preferences and item attributes. For example, `userId X` is converted into an embedded user vector, and `movieId Y` into an embedded item vector.
- These embedding vectors are central to the model's operation, as they distil the high-dimensional input space into a more manageable form while preserving the underlying relationships between users and items.

**Attention Layer**

- The attention layer is the novel addition that differentiates this model from traditional NCF architectures. After obtaining the user and item embeddings, the attention mechanism is applied to

weigh different components of these embeddings according to their relevance in predicting user-item interactions.

- The attention layer works by assigning weights to the features within the user and item embeddings. These weights are learned during training and indicate which features are more important for making accurate predictions. For instance, if certain aspects of a movie's genre or a user's past preferences are more predictive of interaction, the attention mechanism will focus more heavily on these features.

- The output of the attention layer is a refined version of the user and item embeddings, where each feature's contribution is adjusted based on its relevance. This allows the model to concentrate on the most critical information, potentially improving the accuracy of the predictions.

**Concatenation and Fully Connected Layers**

- The next step involves concatenating the attention-adjusted user and item embeddings into a single vector. This concatenated vector represents the combined user-item interaction in the latent feature space, incorporating the insights gained from the attention mechanism.

- The concatenated vector is then passed through a series of fully connected (dense) layers. These layers apply non-linear transformations to the input, enabling the model to capture complex interactions between users and items that are not linear. Each dense layer is followed by a ReLU (Rectified Linear Unit) activation function, which introduces non-linearity and helps the model learn more intricate patterns.

- The depth and width of these layers can be adjusted depending on the complexity of the task and the size of the dataset. The fully connected layers further process the input to refine the predictions and enhance the model's ability to generalize to new, unseen data.

**Output Layer**

- The final component of the architecture is the output layer, which produces a prediction vector. This vector typically contains a single value per user-item pair, representing the probability of interaction (such as clicking, purchasing, or watching). This output is generated by applying a final activation function, such as the sigmoid function, which scales the output to a range between 0 and 1.

- This predicted probability indicates the model's confidence in the likelihood of the user engaging with the item. The predictions can be ranked to provide personalized recommendations, presenting the items with the highest predicted interaction likelihood to the user.

The training process uses backpropagation to adjust the embeddings, attention weights, and fully connected layers to minimize a loss function, binary cross-entropy. The model is trained on user-item interaction data, iteratively refining its parameters to improve the accuracy of its predictions.

The integration of the attention mechanism allows the model to dynamically adjust its focus based on the specific context of each user-item pair. This flexibility makes the model particularly well-suited for complex recommendation scenarios where user preferences and item characteristics may not be straightforward or easily captured by traditional methods.

## 5.1.3 Matrix Factorization Model

Matrix Factorization (MF) is one of the most effective and widely used techniques in collaborative filtering, especially in recommendation systems. It aims to factorize the user-item interaction matrix into two lower-dimensional matrices, capturing the latent factors that represent underlying user preferences and item characteristics. The key advantage of matrix factorization is its ability to uncover complex patterns and relationships within the data that are not immediately apparent from the raw interaction matrix.

**Latent Factor Embeddings**

- The first step in the Matrix Factorization model is to represent both users and items as vectors in a latent factor space. Each user and item are associated with a latent vector, which is of lower dimensionality than the original interaction matrix. These vectors capture the underlying features or factors that influence user preferences and item attributes.
- The `MatrixFactorization` class is initialized with two key components: `user_embedding` and `item_embedding`. These are embedding layers that learn the latent representations of users and items. The number of latent dimensions (denoted by `latent_dim`) determines the complexity of these embeddings and is a critical hyperparameter of the model.

**Dot Product of Latent Vectors**

- The core idea behind matrix factorization is that the interaction between a user and an item can be approximated by the dot product of their respective latent vectors. The dot product effectively measures the similarity between the user's preferences and the item's attributes in the latent factor space.
- In the `forward` method of the model, the user and item embeddings are multiplied element-wise (dot product) and then summed to produce a predicted interaction score (`y_hat`). This score represents the model's estimate of how much the user will like or interact with the item.

**Regularization**

- Regularization is a technique used to prevent overfitting, which can occur when the model learns too much from the training data and performs poorly on unseen data. In this implementation, regularization is applied to the embedding weights of both users and items by scaling them with factors `(1 –`

`regs[0])` and `(1 - regs[1])`, respectively. This ensures that the latent vectors do not become too large, which could lead to overfitting.

**Loss Function**

- The model is trained using a loss function, Mean Squared Error (MSE), which measures the average squared difference between the predicted interaction scores and the actual ratings provided by users. The goal is to minimize this loss, thereby improving the accuracy of the model's predictions.
- In the provided implementation, the loss function is defined using `nn.MSELoss()`, which is a standard loss function for regression tasks in PyTorch.

**Optimization**

- The optimization process involves adjusting the user and item embeddings to reduce the loss. The Adam optimizer (`optim.Adam`) is used in this implementation, which is an adaptive learning rate optimization algorithm that adjusts the learning rate for each parameter individually, leading to faster and more stable convergence.
- During each epoch of training, the model makes predictions for the user-item pairs in the training data, computes the loss, and updates the embeddings using backpropagation.

While Matrix Factorization (MF) models have been a cornerstone in the development of recommender systems, they have notable limitations that the NCF + Attention model seeks to address. MF models are inherently linear, meaning they can struggle to capture the non-linear and complex interactions between users and items. Additionally, MF models tend to rely heavily on the assumption that user and item interactions can be fully explained by latent factors, which may not always hold true, particularly in scenarios with sparse data or where user preferences are highly individualized.

### 5.1.4 User-Based Collaborative Filtering Model

The user-based collaborative filtering model operates on the principle that users with similar past behaviour are likely to have similar future preferences. This model computes the similarity between users based on their interaction histories, typically using metrics like cosine similarity.

**Input,** The input to the model is a sparse matrix where rows represent users and columns represent items, with entries corresponding to the ratings users have given to items.

**Similarity Calculation,** The core of the User-Based CF model is the calculation of similarity between users. This similarity is determined using a metric such as cosine similarity, which measures the cosine of the angle between two vectors in a multi-dimensional space. Here, each user's interaction history is treated as a vector.

The cosine similarity is computed for all pairs of users in the training set, resulting in a user similarity matrix. This matrix captures the degree of similarity between every pair of users, with values ranging from -1 (completely dissimilar) to 1 (identical)

**Prediction,** For a given user, the predicted rating for an item is computed as a weighted average of the ratings given by similar users, with the weights being the similarity scores.

The User-Based CF model has several advantages, particularly its simplicity and interpretability. By focusing on user similarity, the model can provide recommendations that are easy to understand and justify; for example, it can explain recommendations by stating that "similar users liked this item."

However, this approach also has limitations. The reliance on user similarity means that the model can struggle with data sparsity, as it requires sufficient overlap in user interactions to calculate meaningful similarity scores. Additionally, the model may not perform well in the cold start scenario, where new users or items have little to no interaction history.

## 5.2 Model Evaluation & Testing with Unseen Data

To assess the performance of the models, they were evaluated on the testing dataset using the Hit Ratio @10 (HR@10) metric. This metric measures the frequency with which the true item that a user interacted with appears in the top 10 recommendations made by the model.

**Evaluation Methodology**

- Testing Dataset: The models were evaluated on the testing dataset created using the leave-one-out method. For each user, the most recent interaction (based on timestamp) was used as the test instance, and the model was tasked with predicting whether this item would be in the top 10 recommendations.
- Hit Ratio @10 (HR@10): The HR@10 is calculated by checking if the true item is among the top 10 items recommended by the model.
- Comparison: The HR@10 scores of the user-based model, matrix factorization model, NCF model, and the proposed NCF model with attention were compared. The model with the highest HR@10 score is considered to have the best performance, as it most accurately predicts the items that users are likely to interact with.

# CHAPTER 6

# Results and Discussion

## 6.1 Results

The performance of the various models was evaluated using the Hit Ratio at 10 (HR@10) metric, which measures the frequency with which the true item that a user interacted with appears within the top 10 recommendations made by the model. This metric is particularly useful for assessing the practical effectiveness of recommender systems in ranking items in a way that aligns with user preferences.

HR@10 Evaluation Results, in Figure 6. 1 (UCF), Figure 6. 2 (MF), Figure 6. 3 (NCF),Figure 6. 4 (Proposed):

```
[18]: # Evaluate the model
      hit_ratio = hit_ratio_at_10(test_data, predicted_ratings, top_k=10)
      print(f'Hit Ratio@10 for User-based CF Model: {hit_ratio:.2f}')

      Hit Ratio@10 for User-based CF Model: 0.03
```

*Figure 6. 1 HR@10 for User-based CF model*

Evaluating: 100%|████████████████████| 20772
Final Hit Ratio@5: 0.0542

*Figure 6. 2 HR@10 for Matrix Factorization model*

100%  ███████████████████ 55:
Hit Ratio @ 10 is 0.59

*Figure 6. 3 HR@10 for NCF + Implicit Feedback model*

100%  ███████████████ 55397/55397 [03:16
Hit Ratio @ 10 is 0.80

*Figure 6. 4 HR@10 for NCF + Attention + Implicit Feedback model (Proposed Solution)*

These results as observed in Table 6. 1, clearly indicate that the proposed NCF + Attention model significantly outperforms the other models, achieving an HR@10 of 0.80. The standard NCF model also performs well, with an HR@10 of 0.59, though it lags behind the proposed model. Both the User-Based Collaborative

Filtering and Matrix Factorization models show much lower performance, with HR@10 values of 0.03 and 0.05, respectively.

*Table 6. 1 HR@10 Observation for all models*

| Model | HR@10 | Training Time |
|---|---|---|
| **User-based CF model** | 0.03 | 25 minutes |
| **Matrix Factorization** | 0.0542 | 20 minutes |
| **NCF + Implicit Feedback** | 0.59 | 1.5 – 2 hours |
| **NCF + Attention (Proposed)** | 0.80 | 23 minutes |

## 6.2 Discussion

The results of this evaluation highlight the strengths and limitations of different recommendation algorithms, particularly in how they handle the complexities of user-item interactions.

### 6.2.1 NCF + Attention Mechanism + Implicit Feedback (Proposed Model)

The proposed NCF + Attention model's HR@10 score of 0.80 demonstrates its superior ability to make accurate recommendations. This high performance can be attributed to the model's integration of an attention mechanism, which enhances the standard NCF framework by allowing the model to selectively focus on the most relevant features of user and item embeddings. The attention layer effectively highlights the critical aspects of the interaction, enabling the model to better capture complex, non-linear relationships that are crucial for making precise predictions.

By dynamically weighting the importance of different features, the model can adapt to the nuances of each specific user-item pair, resulting in more personalized and accurate recommendations.

### 6.2.2 NCF + Implicit Feedback

The NCF model, with an HR@10 of 0.59, also demonstrates strong performance, though it does not reach the level of accuracy achieved by the proposed model. The NCF model's ability to capture non-linear interactions between users and items through its neural network architecture allows it to outperform traditional models like User-Based Collaborative Filtering and Matrix Factorization.

However, the absence of an attention mechanism means that the NCF model treats all features in the user and item embeddings equally, without the ability to prioritize the most relevant features. This limitation likely accounts for the lower HR@10 compared to the NCF + Attention model.

### 6.2.3 Matrix Factorization

The Matrix Factorization model performs slightly better than the User-Based CF model, with an HR@10 of 0.05. While MF is capable of uncovering latent factors that represent underlying user preferences and item attributes, it is inherently limited by its linear nature. This limitation prevents MF from capturing the non-linear interactions that are often critical for accurate recommendations.

Moreover, the MF model does not account for the varying importance of different latent factors in the same way that an attention mechanism does, which likely contributes to its relatively low HR@10 score.

### 6.2.4 User-based CF model

The User-Based Collaborative Filtering model shows the weakest performance, with an HR@10 of 0.03. This result reflects the challenges faced by user-based approaches, particularly in dealing with data sparsity and the cold start problem. Because this model relies on finding users with similar interaction histories, it struggles when there is insufficient overlap between users, leading to poor recommendation accuracy.

Additionally, the User-Based CF model's reliance on linear similarity measures further limits its ability to capture the more complex relationships that might exist in the data, resulting in lower HR@10.

# Conclusion and Future Work

This thesis aimed to address the critical question: *How to improve personalization and handle complex user-item interactions considering data sparsity and cold start problem in recommendation systems?* Through the development, implementation, and evaluation of various recommendation models, this study provides significant insights and contributions toward answering this question.

The results from the model evaluation, particularly the comparison between the proposed NCF + Attention model and other established models like Neural Collaborative Filtering (NCF), Matrix Factorization (MF), and User-Based Collaborative Filtering, reveal that the integration of attention mechanisms within a neural network framework significantly enhances the model's ability to capture complex, non-linear user-item interactions. The proposed model achieved a (HR@10 of 0.80), outperforming the NCF model (HR@10 of 0.59), Matrix Factorization (HR@10 of 0.05), and User-Based CF (HR@10 of 0.03).

These findings demonstrate that the NCF + Attention model effectively addresses the challenges of data sparsity and the cold start problem by focusing on the most relevant features of user-item interactions. The attention mechanism allows the model to dynamically prioritize the latent factors that are most indicative of user preferences, thereby improving the accuracy of recommendations even in sparse datasets where traditional methods struggle.

Moreover, the superior performance of the NCF + Attention model underscores the importance of leveraging deep learning techniques and attention mechanisms in recommendation systems to enhance personalization. By capturing more intricate patterns of user behaviour and item attributes, the model is better equipped to make personalized recommendations that align with individual user preferences, thereby improving overall user satisfaction.

## 7.1 Thesis Contributions

This thesis makes several important contributions to the field of recommender systems, particularly in the context of improving personalization and handling complex user-item interactions in the face of data sparsity and the cold start problem:

*Development of an Enhanced NCF Model*

- The thesis introduces a novel enhancement to the standard Neural Collaborative Filtering model by integrating an attention mechanism. This addition allows the model to selectively focus on the most critical aspects of user-item interactions, leading to significantly improved recommendation accuracy, as evidenced by the HR@10 metric.

*Addressing Data Sparsity and Cold Start Challenges*

- The research demonstrates that the NCF + Attention model is better equipped to handle the challenges of data sparsity and cold start scenarios compared to traditional methods like Matrix Factorization and User-Based Collaborative Filtering. The ability to focus on relevant latent factors helps the model make more accurate predictions even when interaction data is limited.

*Comprehensive Model Evaluation*

- A thorough evaluation of various models, including NCF, MF, and User-Based CF, was conducted, providing valuable insights into their strengths and weaknesses. This comparison highlights the limitations of traditional models in dealing with complex interactions and underscores the effectiveness of the proposed model.

*Contribution to the Field of Personalized Recommendations*

- By demonstrating the effectiveness of combining deep learning techniques with attention mechanisms, this thesis contributes to the ongoing advancement of personalized recommendation systems. The findings suggest new directions for future research, particularly in exploring more sophisticated attention-based models to further enhance personalization in recommendation systems.

In conclusion, this thesis successfully achieved all of the objectives outlined in Chapter 1, providing a clear pathway for improving personalization in recommendation systems by addressing the inherent challenges of data sparsity and complex user-item interactions. The integration of implicit feedback and attention mechanisms within the NCF framework was demonstrated to be an effective approach, leading to the development of more user-centric and accurate recommendation systems. Each objective was met through the translation of explicit feedback into implicit signals, the implementation and enhancement of the NCF model, and a thorough evaluation of the proposed models, confirming their superiority over traditional methods.

## 7.2 Challenges and Limitations

While this thesis has made significant contributions to the field of recommender systems, several limitations and challenges were encountered during the research process. Acknowledging these challenges is essential for understanding the constraints of the study and identifying areas for future improvement.

*Computational Resources and Model Complexity*

- One of the primary challenges faced was the computational intensity required for training deep learning models, particularly the NCF + Attention model. The complexity of the model, combined with the large size of the MovieLens 20M dataset, required significant processing power and memory. Although efforts were made to manage these demands, including reducing the dataset size by selecting a random 40% of users, the computational limitations restricted the ability to explore even more complex models or larger datasets.

- This limitation also affected the ability to perform extensive hyperparameter tuning. While some tuning was conducted, the resource constraints meant that only a limited number of configurations could be tested, possibly leaving some performance gains unexplored.

*Data Sparsity and Cold Start Problem*

- Although the proposed NCF + Attention model was designed to address the issues of data sparsity and the cold start problem, these challenges persisted throughout the research. The nature of the MovieLens dataset, like many real-world datasets, is inherently sparse, with many users having interacted with only a small subset of the available items. While the attention mechanism improved the model's performance, achieving a more balanced solution to these issues remains a challenge.

- Additionally, the cold start problem, particularly for new users and items, was only partially mitigated by the model. New users or items with no interaction history still present a significant challenge, as the model relies on existing data to make predictions. Further research is needed to develop methods that can better handle these scenarios.

*Generalization to Other Domains*

- Another limitation is the generalization of the findings to domains beyond the scope of the MovieLens dataset. While the results were promising within the context of movie recommendations, the applicability of the NCF + Attention model to other types of recommendation systems (e.g., e-commerce, music, or news) was not tested. Different domains may present unique challenges, such as varying user behaviour patterns or different types of interactions, which could impact the model's effectiveness.

*Interpretability of the Model*

- The integration of an attention mechanism adds complexity to the model, which, while beneficial for performance, also reduces its interpretability. Understanding the exact reasons why the model makes certain recommendations can be difficult, particularly when deep learning techniques are involved.

This lack of transparency can be a drawback in situations where explain ability is crucial, such as in financial or healthcare-related recommendations.

- Efforts were made to interpret the attention weights and their influence on the predictions, but a more in-depth analysis and tools for interpreting these models are needed. This remains an ongoing challenge in the deployment of deep learning-based recommendation systems.

*Evaluation Metrics and Real-World Impact*

- While HR@10 is a widely recognized metric for evaluating recommendation systems, it is not without its limitations. This metric primarily focuses on the accuracy of the top-10 recommendations, but it does not capture other important aspects of recommendation quality, such as diversity, novelty, or user satisfaction. As a result, the evaluation may not fully reflect the model's performance in a real-world setting.

- Furthermore, the real-world impact of the recommendations, such as how they influence user behaviour or satisfaction, was not directly measured in this study. Future research should consider incorporating additional metrics and real-world user feedback to provide a more holistic evaluation of recommendation system performance.

## 7.3 Future Work

While this thesis has enhanced personalization and handled complex user-item interactions in recommendation systems, there are several avenues for future research that could further improve the effectiveness and applicability of the proposed model. The following areas highlight potential directions for future work:

*Incorporating Sequential Data to Capture Temporal Dynamics*

- One of the key limitations of the current model is its static treatment of user preferences and item characteristics. However, in reality, a user's tastes and preferences can evolve over time due to various factors such as trends, life changes, or simply exposure to new content. Incorporating sequential data into the model could allow it to capture these temporal dynamics more effectively.

- Future work could involve developing models that incorporate sequence-aware techniques, such as Recurrent Neural Networks (RNNs), Long Short-Term Memory (LSTM) networks, or Transformer models, which are capable of learning from the order of interactions.

*Integrating Explainable AI (XAI) Algorithms*

- As the complexity of recommendation models increases, so does the need for transparency and interpretability. Users and stakeholders often need to understand why certain recommendations are made, especially in sensitive domains like healthcare, finance, or legal contexts. To address this, future work could explore the integration of Explainable AI (XAI) algorithms, such as SHAP (SHapley Additive exPlanations) and LIME (Local Interpretable Model-agnostic Explanations).

*Expanding to Multi-Modal Data*

- Another promising direction for future research is the incorporation of multi-modal data into the recommendation process. Currently, the model primarily relies on user-item interaction data, but other types of data, such as textual reviews, images, audio, and video content, could provide additional insights into user preferences.
- By integrating data from multiple modalities, the model could develop a richer and more nuanced understanding of both users and items. For example, combining text embeddings from user reviews with visual features extracted from item images could lead to more comprehensive recommendations. Deep learning models like Convolutional Neural Networks (CNNs) for images and Transformers for text could be employed to process these different types of data.

*Addressing Cold Start with Hybrid Models*

- The cold start problem, particularly for new users and items, remains a significant challenge in recommendation systems. Future work could explore the development of hybrid models that combine collaborative filtering with content-based filtering and other techniques to mitigate this issue.
- For instance, using demographic information, item metadata, or contextual information (such as time, location, or social connections) could provide additional signals that help the model make accurate recommendations even when interaction data is sparse. Incorporating such auxiliary information could enhance the model's ability to make predictions in cold start scenarios.

The proposed directions for future work aim to build upon the foundation established in this thesis, addressing current limitations, and exploring new opportunities for enhancing recommendation systems. By incorporating sequential data, integrating XAI algorithms, expanding to multi-modal inputs, and addressing cold start challenges, future research can further improve the personalization, accuracy, and transparency of recommendations. These advancements will contribute to the development of more effective and user-centric recommendation systems across various domains.

# References

1. Xiangnan He, Lizi Liao, Hanwang Zhang, Liqiang Nie, Xia Hu, and Tat-Seng Chua. 2017. Neural Collaborative Filtering. In Proceedings of the 26th International Conference on World Wide Web (WWW '17). International World Wide Web Conferences Steering Committee, Republic and Canton of Geneva, CHE, 173–182. https://doi.org/10.1145/3038912.3052569

2. Fuzheng Zhang, Nicholas Jing Yuan, Defu Lian, Xing Xie, and Wei-Ying Ma. 2016. Collaborative Knowledge Base Embedding for Recommender Systems. In Proceedings of the 22nd ACM SIGKDD International Conference on Knowledge Discovery and Data Mining (KDD '16). Association for Computing Machinery, New York, NY, USA, 353–362. https://doi.org/10.1145/2939672.2939673

3. Guo, G. (2012). Resolving Data Sparsity and Cold Start in Recommender Systems. In: Masthoff, J., Mobasher, B., Desmarais, M.C., Nkambou, R. (eds) User Modeling, Adaptation, and Personalization. UMAP 2012. Lecture Notes in Computer Science, vol 7379. Springer, Berlin, Heidelberg. https://doi.org/10.1007/978-3-642-31454-4_36

4. Yue Shi, Martha Larson, and Alan Hanjalic. 2014. Collaborative Filtering beyond the User- Item Matrix: A Survey of the State of the Art and Future Challenges. ACM Comput. Surv. 47, 1, Article 3 (July 2014), 45 pages. https://doi.org/10.1145/2556270

5. Torkashvand, A., Jameii, S.M. & Reza, A. Deep learning-based collaborative filtering recommender systems: a comprehensive and systematic review. *Neural Comput & Applic* **35**, 24783–24827 (2023). https://doi.org/10.1007/s00521-023-08958-3

6. Jiajia Jiang, Weiling Li, Ani Dong, Quanhui Gou, Xin Luo, A Fast Deep AutoEncoder for high-dimensional and sparse matrices in recommender systems, Neurocomputing, Volume 412,2020, Pages 381-391, ISSN0925-2312, https://doi.org/10.1016/j.neucom.2020.06.109

7. Keyvan Vahidy Rodpysh, Seyed Javad Mirabedini, Touraj Banirostam, Resolving cold start and sparse data challenge in recommender systems using multi-level singular value decomposition, Computers & Electrical Engineering, Volume 94,2021,107361, ISSN 00457906, https://doi.org/10.1016/j.compeleceng.2021.107361

8. Antonio Hernando, Jesús Bobadilla, Fernando Ortega, Abraham Gutiérrez, A probabilistic model for recommending to new cold-start non-registered users, Information Sciences, Volume 376,2017, Pages 216-232, ISSN 0020-0255, https://doi.org/10.1016/j.ins.2016.10.009

9. Senthilselvan Natarajan, Subramaniyaswamy Vairavasundaram, Sivaramakrishnan Natarajan, Amir H. Gandomi, Resolving data sparsity and cold start problem in collaborative filtering recommender system using Linked Open Data, Expert Systems with Applications, Volume 149,2020,113248, ISSN0957-4174, https://doi.org/10.1016/j.eswa.2020.113248

10. Fan Zhou, Yuhua Mo, Goce Trajcevski, Kunpeng Zhang, Jin Wu, Ting Zhong,Recommendation via Collaborative Autoregressive Flows, Neural Networks, Volume 126, 2020, Pages 52-64, ISSN 0893-6080, https://doi.org/10.1016/j.neunet.2020.03.010

11. Xia, H., Luo, Y. & Liu, Y. Attention neural collaboration filtering based on GRU for recommender systems. *Complex Intell. Syst.* **7**, 1367–1379 (2021). https://doi.org/10.1007/s40747-021-00274-4

12. Vaswani, A., Shazeer, N., Parmar, N., Uszkoreit, J., Jones, L., Gomez, A.N., Kaiser, Ł., and Polosukhin, I., 2017. Attention is All You Need. In *Proceedings of the 31st International Conference on Neural Information Processing Systems (NIPS 2017)*, pp. 6000-6010. Available at: https://doi.org/10.48550/arXiv.1706.03762

13. Alsini, A., Huynh, D.Q., and Datta, A., 2020. Hit Ratio: An Evaluation Metric for Hashtag Recommendation. *arXiv preprint*. Available at: https://doi.org/10.48550/arXiv.2010.01258

14. Dheeraj Bokde, Sheetal Girase, Debajyoti Mukhopadhyay,Matrix Factorization Model in Collaborative Filtering Algorithms: A Survey, Procedia Computer Science, Volume 49, 2015, Pages 136-146, ISSN 1877-0509, https://doi.org/10.1016/j.procs.2015.04.237

# Appendix

The work provided in this thesis is supported by the extensive material and extra resources found in this appendix. The techniques for developing models and doing exploratory data analysis (EDA) are covered in the information that follows. PyTorch Lightning is used to train models in an effective and scalable manner. The GitHub & GitLab (University Server) repository that is linked below contains the code for these procedures.

*'rating_eda.ipynb'*
The MovieLens 20M dataset is thoroughly analysed in the EDA notebook. It contains preliminary findings about user-item interactions, data cleansing, and visualization.

*'UCF.ipynb'*
This notebook describes how the user-based collaborative filtering model is implemented, including how user similarities are determined and how user-item interactions are predicted.

*'mf_re.ipynb'*
This notebook covers the development of the matrix factorization model, focusing on the decomposition of the user-item interaction matrix into latent factors.

*'ncf_re.ipynb'*
This notebook documents the creation of the NCF model, emphasizing the use of deep learning techniques to capture complex user-item relationships.

*'attn_ncf.ipynb'*
This notebook showcases the implementation of the proposed NCF model with an integrated attention mechanism, designed to enhance the model's focus on the most relevant features of user-item interactions.

*GitHub & GitLab, and Dataset*

The full code and notebooks are available in the GitHub repository linked below:

GitHub: https://github.com/Ashwin-create/project-2023-24-axm20781
(in case the GitLab page does not work)

GitLab (University Server): https://git.cs.bham.ac.uk/projects-2023-24/axm2078

Dataset: MovieLens 20M (rating dataset)

To install necessary libraries, run the following commands on Google Colab or Jupyter Notebook:

- ➢ *!pip install pytorch-lightning*
- ➢ *!pip install numpy pandas matplotlib scikit-learn*