

AMAZON RECOMMENDATION BASED ON RATING AND REVIEWS USING BERT ANALYSIS AND  
COLLABORATIVE FILTERING

ARPEET KUMAR

Final Thesis Report

DECEMBER 2023

## TABLE OF CONTENTS

ACKNOWLEDGEMENT.....	V
ABSTRACT.....	VI
LIST OF TABLES.....	VII
LIST OF FIGURES.....	VIII
CHAPTER 1: INTRODUCTION.....	8
1.1 Background of this Study.....	8
1.2 Problem Statement.....	9
1.3 Aim and Objectives.....	10
1.4 Scope of Study.....	11
CHAPTER 2: LITERATURE REVIEW.....	12
2.1 Introduction.....	12
2.2 Recommender System.....	12
2.2.1 History of Recommender System.....	12
2.2.2 Application of Recommender System.....	14
2.2.3 Types of Recommender System.....	16
2.3 Challenges in recommender System.....	16
2.4 Summary.....	16
CHAPTER 3: RESEARCH METHODOLOGY.....	18
3.1 Introduction.....	18
3.2 Methodology.....	18
3.3 Data Selection.....	19
3.4 Data PreProcessing.....	20
3.5 Sentiment Analysis using BERT.....	20
3.6 Collaborative Filtering.....	20
3.7 Tools and libraries.....	21
3.8 Summary.....	23

CHAPTER 4: IMPLEMENTATION.....	25
4.1 Introduction.....	25
4.2 Dataset.....	25
4.3 Data Cleaning.....	27
4.4 Exploratory Analysis.....	31
4.5 Sentiment Analysis using BERT.....	34
4.6 Collaborative Filtering.....	37
4.7 Hybrid Approach-Colaborative Filtering with BERT Sentiment.....	38
CHAPTER 5: RESULTS AND EVALUATION.....	40
5.1 Introduction.....	40
5.2 Model Output.....	40
5.3 Summary.....	41
CHAPTER 6: CONCLUSION AND RECOMENDATIONS.....	42
6.1 Introduction.....	42
6.2 Discussion and Conclusion.....	42
6.3 Contribution.....	43
6.4 Future Works.....	44
REFERENCES.....	46

## ACKNOWLEDGMENT

I would like to express my heartfelt gratitude to my primary advisor, Prof. Anurag Majji whose unwavering support, expert guidance, and insightful feedback have been invaluable throughout this research journey. Your mentorship has significantly shaped the trajectory of my academic pursuits.

My sincere thanks go to my family for their unwavering support and understanding during the highs and lows of this academic endeavor. Your belief in my abilities has been a constant source of motivation.

I extend my appreciation to my friends and classmates who shared their thoughts, engaged in meaningful discussions, and provided constructive feedback. The collaborative spirit within our academic community has enriched my learning experience

Acknowledgment is also due to the dedicated staff at Liverpool John Moores University for their assistance in accessing resources and for providing a conducive environment for research.

This thesis is not just an individual effort but a collective endeavor made possible by the support, encouragement, and collaboration of these wonderful individuals and institutions.

## ABSTRACT

This thesis presents a novel approach to enhance recommender systems by integrating BERT sentiment analysis with collaborative filtering techniques. The research addresses the growing demand for personalized recommendations in the e-commerce landscape, with a focus on understanding user sentiments and preferences. The study begins with an exploration of the historical evolution of recommender systems, tracing the development from early information retrieval techniques to the integration of deep learning models in recent years.

The methodology section details the development and implementation of a hybrid recommender system, combining advanced language understanding through BERT sentiment analysis with collaborative filtering algorithms. The project utilizes a comprehensive dataset from Amazon, incorporating user reviews and ratings. The textual data undergoes preprocessing, including tokenization and stemming, while missing data and outliers are handled appropriately. The sentiment analysis model is fine-tuned on the review dataset, extracting emotional tones to enhance the recommendation process.

Two primary outputs are discussed: the performance of the item-based collaborative filtering model on its own and the improved results achieved through the hybrid model. Evaluation metrics, including Root Mean Squared Error (RMSE), validate the efficacy of the combined approach in providing more accurate and sentiment-aware recommendations.

The findings contribute to the field of recommendation systems, showcasing the potential of integrating sentiment analysis for a deeper understanding of user preferences. The hybrid model demonstrates adaptability and scalability, offering insights for future research in the dynamic landscape of e-commerce and personalized content delivery.

In conclusion, the project underscores the significance of considering user sentiments in recommendation algorithms and provides a foundation for further advancements in enhancing the precision and personalization of recommender systems.

## LIST OF FIGURES

Figure 4.1	Columns of the First Dataset.....	26
Figure 4.2	Columns of the Second Dataset.....	27
Figure 4.3	Information of the First Dataset.....	28
Figure 4.4	Cleaned Dataset Set for our requirement.....	28
Figure 4.4	Information of the Metadat Dataset.....	29
Figure 4.5	Cleaned Metadata Dataset Set for our requirement.....	30
Figure 4.6	Result of the Combined Dataset.....	30
Figure 4.7	Ratings based on the total number of Reviews.....	31
Figure 4.8	Top Companies in term in Sales.....	32
Figure 4.9	Games mostly Bought.....	33
Figure 4.10	Highest Rated Games.....	34
Figure 4.11	Sentiment Count.....	35
Figure 4.12	WordCloud for Reviews.....	36
Figure 4.13	WordCloud for Positve Reviews.....	36
Figure 4.14	WordCloud for Negative Reviews.....	37

## CHAPTER 2

### INTRODUCTION

#### 1.1 Background of this Study

The landscape of e-commerce has witnessed a paradigm shift with the advent of recommender systems, becoming integral to user experience on platforms like Amazon. The dynamic nature of user preferences and the sheer abundance of products necessitate advanced techniques to enhance the accuracy and personalization of recommendations. This study focuses on refining Amazon's recommendation system by synergizing two powerful methodologies: BERT analysis and collaborative filtering.

Collaborative filtering has been a cornerstone in recommender systems, leveraging user-item interactions to predict preferences. However, this approach faces challenges in handling sparse data and capturing nuanced sentiments expressed in user reviews. To address this, the study incorporates BERT (Bidirectional Encoder Representations from Transformers) analysis, a state-of-the-art natural language processing model. BERT excels in understanding context and sentiment within textual data, making it an ideal tool for extracting valuable insights from user reviews.

The historical evolution of recommender systems provides context to the study, highlighting the progression from basic collaborative filtering to more intricate hybrid models. The integration of content-based filtering and the recent emergence of deep learning techniques underscore the industry's commitment to refining personalized recommendations.

Amazon, a pioneer in e-commerce, serves as an exemplary case study due to its vast product catalog and diverse user base. The study aims to harness the strengths of both collaborative filtering and BERT analysis to create a more nuanced, accurate, and context-aware recommendation system. By delving into the rich qualitative data embedded in customer reviews, the model seeks to understand not just what products users prefer but why, thereby enhancing the interpretability and effectiveness of the recommendations.

The challenges of the cold start problem and scalability issues inherent in collaborative filtering will be addressed through the incorporation of BERT analysis, ensuring a more robust and adaptive recommendation system. As user-generated content continues to grow in significance, the study anticipates contributing valuable insights that go beyond quantitative

metrics, laying the groundwork for an enhanced, emotionally intelligent recommendation engine on the Amazon platform.

## **1.2 Problem Statement**

Amazon's existing recommendation system is effective in leveraging user purchase history and browsing behavior, it operates primarily on quantitative data. This approach focuses on patterns and trends derived from what users buy and what they look at on the platform. While this method is valuable for understanding user preferences and suggesting relevant products, it has limitations when it comes to capturing the qualitative aspects of user experience.

The current system might not fully tap into the wealth of information available in customer reviews and ratings. Customer reviews contain rich, qualitative insights into the strengths and weaknesses of products, providing a nuanced understanding of user satisfaction. Ratings also offer a numerical representation of customer sentiment.

Enhancing the recommendation system by incorporating customer reviews and ratings can offer several advantages:

**1.2.1 Fine-grained Personalization:** Analyzing reviews allows the system to understand specific features or qualities of a product that resonate with individual users. This granular understanding enables more precise recommendations tailored to individual preferences.

**1.2.2 Improved Product Understanding:** Customer reviews provide valuable information about product performance, durability, and user experiences. Incorporating this qualitative data into the recommendation algorithm enhances the system's ability to understand the diverse factors influencing user satisfaction.

**1.2.3 Addressing User Concerns:** By considering both positive and negative sentiments expressed in reviews, the recommendation system can address potential concerns or drawbacks associated with products. This can lead to more informed recommendations, ultimately improving user satisfaction.



**1.2.4 Uncovering Emerging Trends:** Analyzing reviews allows the system to identify emerging trends and preferences among users. This dynamic understanding of the market can lead to proactive recommendations for new or popular products.

**1.2.5 Building Trust:** Users often trust peer opinions and experiences shared in reviews. By incorporating this information into the recommendation system, Amazon can enhance user trust and confidence in the suggested products.

To implement this, a more sophisticated recommendation algorithm would need to be developed, one that can extract and analyze sentiments, key phrases, and topics from customer reviews. Natural Language Processing (NLP) techniques can be employed to understand the context and sentiment expressed in reviews, allowing the recommendation system to make more nuanced and context-aware suggestions.

By combining quantitative data from purchase history and browsing behavior with qualitative insights from reviews and ratings, Amazon can create a more comprehensive and user-centric recommendation system that reflects both the preferences and experiences of its diverse user base.

### **1.3 Aim And Objectives**

The primary aim of this project is to revolutionize Amazon's recommendation system by developing a sophisticated algorithm that harnesses the power of unstructured customer reviews. The key objective is to unlock valuable insights embedded in these reviews, providing a more nuanced understanding of product attributes, user experiences, and overall sentiment. Leveraging advanced Natural Language Processing (NLP) techniques, the algorithm will be designed to effectively parse unstructured textual data, identifying key phrases, sentiments, and topics within customer reviews. Additionally, sentiment analysis will play a pivotal role in gauging the emotional tone of both reviews and ratings, enabling the system to comprehend user satisfaction levels and concerns. The incorporation of sentiment analysis will enhance the algorithm's ability to discern subtle nuances in user feedback, allowing for a more personalized and emotionally intelligent recommendation system.

Furthermore, the project aims to go beyond static recommendation models by implementing a dynamic system that adapts in real-time to changing customer preferences. Traditional

recommendation systems often struggle to keep pace with evolving user tastes and emerging trends. In response, our objective is to create a recommendation algorithm that continuously learns and adjusts based on the latest feedback, ensuring that the system remains agile and responsive to the dynamic nature of consumer preferences. By combining insights from both structured (purchase history, browsing behavior) and unstructured (customer reviews and sentiments) data sources, the recommendation system will evolve into a robust, customer-centric platform that not only accurately reflects individual preferences but also proactively anticipates and adapts to shifts in the market landscape. Through this comprehensive approach, the project seeks to elevate the user experience on Amazon by providing more accurate, personalized, and emotionally resonant product recommendations.

## **1.4 Scope of Study**

**1.4.1 Data Collection:** Acquire a diverse data set encompassing a wide range of products, including different categories and user demographics. Implement a robust data collection process for customer reviews and ratings, ensuring a representative sample size.

**1.4.2 Data Preprocessing:** Develop methods for cleaning and preprocessing unstructured text data from reviews to facilitate effective analysis. Explore techniques for feature extraction from reviews, including sentiment scores and key phrases.

**1.4.3 Sentiment Analysis:** Implement sentiment analysis algorithms to categorize reviews into positive, negative, or neutral sentiments. Investigate the impact of sentiment polarity on product recommendations.

**1.4.4 Machine Learning Models:** Develop recommendation algorithms that incorporate both quantitative (rating-based) and qualitative (review-based) factors. Evaluate and compare the performance of different machine learning models for recommendation tasks.

**1.4.5 Evaluation Metrics:** Define and employ appropriate metrics for evaluating the effectiveness of the recommendation system, such as precision, recall, and user satisfaction. Conduct A/B testing to compare the performance of the enhanced system against the existing recommendation approach.

**1.4.6 Documentation and Reporting:** Maintain detailed documentation of the methodologies, algorithms, and results obtained during the study. Prepare a comprehensive report summarizing the findings, challenges faced, and recommendations for future enhancements.

## **CHAPTER 2**

### **LITERATURE REVIEW**

#### **2.1 Introduction**

The evolution of recommender systems has played a pivotal role in shaping user interactions and experiences across various digital platforms. From their early roots in rule-based algorithms during the 1960s to the dominance of collaborative filtering in the 1990s, recommender systems have undergone significant advancements. This literature review explores the historical progression of these systems, delving into key milestones such as the integration of content-based filtering, the rise of hybrid models, and the transformative impact of deep learning in the 2010s. Examining their applications across diverse sectors, from e-commerce and streaming services to healthcare and personalized marketing, provides insights into the wide-ranging impact and continual evolution of recommender systems. The challenges faced by these systems, such as cold start problems and scalability issues, further contribute to a comprehensive understanding of the current landscape in personalized digital recommendations.

#### **2.2 Recommender System**

Recommender systems, also known as recommendation systems or engines, play a pivotal role in modern information and e-commerce landscapes. These systems are designed to analyze and predict user preferences, providing personalized suggestions or recommendations that cater to individual tastes and needs. The primary goal is to enhance user experience by facilitating the discovery of relevant content or products in an increasingly vast and diverse digital environment.

##### **2.2.1 History of Recommender System**

Recommender systems have become an integral part of our digital experience, helping users discover new content and products tailored to their preferences. The evolution of recommender systems can be traced back to the early days of information retrieval and has since undergone significant advancements.

##### **1. Early Foundations (1960s-1990s):**

The roots of recommender systems can be found in the field of information retrieval. In the 1960s and 1970s, researchers began developing basic recommendation algorithms to enhance the user's ability to find relevant information. One of the earliest examples was the "Smart" information retrieval system developed at Cornell University in the 1960s, which employed rule-based techniques to suggest relevant documents. The emergence of collaborative filtering (CF) in the late 1980s marked a crucial milestone. In 1982, David Goldberg introduced the concept of collaborative filtering in his research on using collaborative user feedback to improve recommendations. The idea gained momentum, and collaborative filtering became a dominant approach in recommender systems.

## **2. Rise of Collaborative Filtering (1990s-2000s):**

During the 1990s, collaborative filtering algorithms gained popularity due to their ability to make predictions based on user-item interactions. GroupLens, a project at the University of Minnesota, played a pivotal role by introducing collaborative filtering techniques for recommending Usenet news articles.

Netflix, known for its movie recommendation challenge, further propelled collaborative filtering into the limelight. In 2006, Netflix offered a million-dollar prize for improving its movie recommendation algorithm, stimulating extensive research and innovation in the field.

## **3. Content-Based Filtering and Hybrid Approaches (2000s-2010s):**

The 2000s witnessed the integration of content-based filtering, which recommends items based on their attributes and the user's preferences. This approach addressed some limitations of collaborative filtering, especially in situations with sparse data.

Hybrid recommender systems, combining collaborative filtering and content-based filtering, gained prominence. These systems aimed to leverage the strengths of both approaches, offering more accurate and personalized recommendations. Additionally, matrix factorization techniques, such as Singular Value Decomposition (SVD), gained traction for enhancing collaborative filtering accuracy.

## **4. Deep Learning and Personalization (2010s-Present):**

The 2010s marked the advent of deep learning techniques in recommender systems. Neural collaborative filtering and deep learning-based models allowed systems to capture intricate patterns in user behavior and preferences, leading to improved recommendation accuracy.

With the rise of platforms like Amazon, Spotify, and Netflix, recommender systems evolved to provide highly personalized recommendations through sophisticated algorithms. Reinforcement learning and contextual bandit algorithms were also introduced to adapt recommendations in real-time based on user feedback.

## **5. Conclusion:**

The history of recommender systems reflects a journey from basic information retrieval to advanced machine learning models. As we move forward, the integration of artificial intelligence, deep learning, and continual refinements in algorithms promise to make recommender systems even more adept at understanding and predicting user preferences, shaping the future of personalized digital experiences.

### **2.2.2 Application of Recommender System**

Recommender systems, also known as recommendation systems or engines, have diverse applications across various industries. Here are several areas where recommender systems are commonly used:

#### **1. E-commerce and Retail:**

- a) **Product Recommendations:** Recommender systems help users discover products based on their preferences, purchase history, and browsing behavior.
- b) **Cross-selling and Up-selling:** Recommending complementary or higher-value products to increase the average transaction value.

#### **2. Streaming Services:**

- a) **Content Recommendations:** Recommender systems suggest movies, TV shows, music, or other multimedia content based on user viewing or listening history.
- b) **Playlist Generation:** Creating personalized playlist tailored to individual user preferences.

#### **3. Social Media:**

- a) **Friend Recommendations:** Recommender systems suggest new connections based on mutual friends, interests, or other factors.
- b) **Content Sharing:** Recommending posts or content that align with a user's interests.

#### **4. Online Travel and Hospitality:**

- a) **Hotel and Accommodation Recommendations:** Suggesting lodging options based on previous bookings, user reviews, and preferences.

- b) Travel Itinerary Planning: Recommending activities, attractions, and dining options for a personalized travel experience.

**5. Job Portals:**

- a) Job Recommendations: Recommending job opportunities based on the user's skills, experience, and career preferences.

**6. Healthcare:**

- a) Treatment Recommendations: Recommender systems can assist healthcare professionals in suggesting personalized treatment plans based on patient history and medical data.
- b) Wellness and Fitness: Recommending personalized fitness routines, diet plans, and wellness activities.

**7. Education:**

- a) Course Recommendations: Recommending online courses, learning materials, and resources based on a user's educational background and interests.

**8. Financial Services:**

- a) Investment Recommendations: Recommender systems can suggest investment options based on user risk tolerance, financial goals, and market trends.

**9. News and Content Aggregation:**

- a) Personalized News Feeds: Recommending news articles, blogs, or content based on a user's reading habits and interests.

**10. Automotive Industry:**

- a) Vehicle Recommendations: Recommending cars based on user preferences, budget, and lifestyle.

**11. Human Resources:**

- a) Employee Skill Matching: Recommender systems assist in matching employees with relevant projects or teams based on their skills and expertise.

**12. Real Estate:**

- a) Property Recommendations: Recommending real estate properties based on user preferences, budget, and location.

**13. Gaming:**

- a) Game Recommendations: Recommending video games based on a user's gaming history, preferences, and genre interests.

**14. Supply Chain and Inventory Management:**

- a) **Product Restocking:** Recommender systems assist in predicting and recommending inventory restocking based on historical sales data.

#### 15. **Personalized Marketing:**

- a) **Targeted Advertising:** Recommender systems enhance targeted marketing by suggesting products or services that align with individual user preferences.

### 2.2.3 **Types of Recommender Systems:**

There are mainly 3 types of Recommender System and are mentioned below:

1. **Collaborative Filtering:** Recommender systems can use the preferences of similar users to make recommendations. User-based collaborative filtering identifies users with similar preferences, while item-based collaborative filtering focuses on similarities between items.
2. **Content-Based Filtering:** As the volume of users and items increases, the computational demands on recommender systems grow, requiring efficient algorithms and infrastructure.
3. **Hybrid Models:** Many systems combine collaborative filtering and content-based filtering to leverage the strengths of both approaches.

### 2.3 **Challenges in recommender system**

Text should begin at this position and continue to the end of the left margin. Text must be typed using 1.5 spacing.

16. **Cold Start Problems:** Recommender systems may struggle when dealing with new users (cold start) or new items that lack sufficient interaction history.
17. **Scalability:** As the volume of users and items increases, the computational demands on recommender systems grow, requiring efficient algorithms and infrastructure.
18. **Diversity:** Striking a balance between providing personalized recommendations and introducing diversity in suggestions is a common challenge.

### 2.4 **Summary**

Recommender systems, essential in today's digital landscape, serve as the backbone of personalized user experiences by predicting and catering to individual preferences. The evolution of these systems is marked by foundational stages, with early rule-based algorithms in the 1960s evolving into sophisticated machine learning models.



The history encompasses the dominance of collaborative filtering in the 1990s, notably influenced by projects like GroupLens at the University of Minnesota and Netflix's recommendation challenge. The 2000s saw the integration of content-based filtering to address collaborative filtering limitations, leading to the rise of hybrid recommender systems. Matrix factorization techniques and the advent of deep learning in the 2010s further enhanced the accuracy of recommendations, with platforms like Amazon and Netflix incorporating these advancements for highly personalized suggestions.

Recommender systems find diverse applications, from e-commerce and streaming services to healthcare and personalized marketing. They leverage collaborative filtering, content-based filtering, and hybrid models to offer accurate and relevant suggestions.

However, challenges persist. Cold start problems arise with new users or items lacking interaction history. Scalability becomes an issue as the volume of users and items increases, demanding efficient algorithms. Striking a balance between personalized recommendations and introducing diversity poses an ongoing challenge in the field.

In summary, recommender systems have undergone a transformative journey, and as technology advances, they continue to shape the future of personalized digital experiences, promising even more adept understanding and prediction of user preferences.

## CHAPTER 3

### RESEARCH METHODOLOGY

#### 3.1 Introduction

In crafting an innovative methodology to enhance Amazon's recommendation system, this research combines two cutting-edge approaches: BERT (Bidirectional Encoder Representations from Transformers) analysis and collaborative filtering. Collaborative filtering, a foundational technique, has demonstrated prowess in predicting user preferences based on historical interactions. However, the study recognizes its limitations in handling sparse data and extracting qualitative insights from user reviews. Integrating BERT, a state-of-the-art natural language processing model, allows for a more in-depth analysis of customer sentiments embedded in reviews, enriching the system's understanding of user preferences.

The research methodological framework draws inspiration from the historical evolution of recommender systems, acknowledging the industry's trajectory from conventional collaborative filtering to hybrid models. Focused on Amazon's extensive product catalog and diverse user base, the methodology aims to leverage the strengths of collaborative filtering and BERT analysis. By dissecting user-generated content and addressing challenges like the cold start problem, the study aspires to contribute novel insights, laying the foundation for a sophisticated, emotionally intelligent recommendation system. This methodological synthesis anticipates refining the accuracy and personalization of Amazon's recommendations, aligning with the evolving landscape of e-commerce.

#### 3.2 Methodology

This research employs a hybrid methodology, integrating BERT (Bidirectional Encoder Representations from Transformers) analysis and collaborative filtering to enhance Amazon's recommendation system. Collaborative filtering, a foundational technique, harnesses historical user interactions to predict preferences. Acknowledging its limitations in handling sparse data and extracting qualitative insights from reviews, the study incorporates BERT, a state-of-the-art natural language processing model. This fusion allows for a nuanced understanding of customer sentiments within reviews, complementing the quantitative data provided by collaborative filtering. The methodology draws inspiration from the historical evolution of recommender systems, emphasizing Amazon's extensive product catalog and diverse user base. Through meticulous data preprocessing, training, and optimization, the

research aims to create a sophisticated recommendation model that combines the strengths of collaborative filtering and BERT analysis. Evaluation metrics and user testing will validate the model's effectiveness in providing accurate and emotionally intelligent recommendations.

### **3.3 Data Selection**

The dataset comprises two key components: a comprehensive collection of product reviews and detailed information about corresponding products, each identified by a unique ASIN (Amazon Standard Identification Number). The review dataset, with 497,577 entries, provides insights into user sentiments through columns like 'overall' (review rating), 'verified' (verified purchase status), 'reviewerID' (unique identifier for reviewers), and 'reviewText' (textual review content).

The richness of the review dataset lies in its diversity, capturing varied user opinions across products. The 'verified' column ensures the inclusion of reviews from authenticated purchasers, enhancing the dataset's reliability. The temporal aspect is addressed through 'reviewTime' and 'unixReviewTime,' allowing for time-based analyses.

Complementing the reviews, the product dataset, with 84,819 entries, provides valuable context. Features like 'category,' 'brand,' and 'main\_cat' offer categorical information, enabling content-based filtering. Additionally, 'rank' and 'price' provide quantitative attributes for potential collaborative filtering strategies.

Despite occasional missing values in 'reviewerName,' 'reviewText,' and 'summary,' these can be mitigated through data preprocessing techniques. The inclusion of 'vote' and 'image' columns adds dimensions for potential sentiment analysis and image-based feature extraction. The product dataset introduces further dimensions, including 'category,' 'description,' and 'brand,' facilitating a hybrid approach. However, potential challenges arise from the sparse data in columns like 'imageURL' and 'details,' necessitating careful handling.

In summary, the combined dataset presents a robust foundation for this project. Its scale, diversity, and inclusion of both review and product information allow for the application of advanced recommendation techniques such as collaborative filtering and BERT sentiment analysis. The dataset's richness in both quantitative and qualitative features positions it as a valuable resource for developing a sophisticated, context-aware recommender system tailored to Amazon's diverse product landscape.

### **3.4 Data Pre-Processing**

In the preparation of textual data from reviews for advanced analysis, a rigorous cleaning and preprocessing pipeline is implemented. This involves tokenization to break down the review text into individual units, stemming to standardize words to their root form, and the removal of stop words to focus on content-carrying words. Ratings are transformed into a numerical format suitable for collaborative filtering algorithms, facilitating personalized recommendations based on user-item interactions. The handling of missing data employs techniques such as imputation or removal, ensuring dataset completeness. Outliers, data points significantly deviating from the norm, are identified and appropriately managed to prevent undue influence on subsequent analyses. This meticulous process results in a refined dataset, setting the stage for the development of a recommendation system tailored to Amazon's diverse product landscape.

### **3.5 Sentiment Analysis using BERT:**

The sentiment analysis phase involves fine-tuning a pre-trained BERT (Bidirectional Encoder Representations from Transformers) model using the review dataset. Fine-tuning allows the model to adapt its parameters to the specifics of the dataset, optimizing its performance for sentiment analysis. This process entails training the BERT model on the labeled review dataset, adjusting its weights to capture nuanced sentiments within the textual data. Following fine-tuning, the model is applied to extract sentiment scores or labels for each review, categorizing them as positive, negative, or neutral. By decoding the emotional tone embedded in the reviews, this sentiment analysis offers a comprehensive understanding of user sentiments, providing valuable insights into the overall satisfaction or dissatisfaction expressed in the Amazon product reviews.

### **3.6 Collaborative Filtering:**

The implementation phase involves deploying collaborative filtering algorithms, specifically user-based or item-based collaborative filtering, to generate personalized recommendations based on user-item interactions within the review dataset. In user-based collaborative filtering, similarities between users are assessed to recommend items liked by users with similar preferences. Conversely, item-based collaborative filtering identifies similarities between items, recommending items similar to those a user has previously engaged with. Additionally, matrix factorization techniques, such as Singular Value Decomposition (SVD), are explored to capture latent features inherent in the user-item interaction matrix. SVD decomposes the

matrix into latent factors, uncovering hidden patterns that contribute to the user's preferences. By leveraging these techniques, the collaborative filtering algorithms aim to enhance the accuracy and personalization of recommendations, aligning with the objective of creating a sophisticated and context-aware recommender system for the Amazon platform.

### 3.7 Tools and Libraries

For this project everything will be performed in the Jupyter Notebook on Python with these tools and libraries. These tools and libraries collectively cover data processing, visualization, sentiment analysis, machine learning, deep learning, and recommender system development, providing a comprehensive set of resources

#### 1. Data Processing and Exploration:

- a) **NumPy (np):** NumPy is a powerful library for numerical operations in Python, providing support for large, multi-dimensional arrays and matrices.
- b) **Pandas (pd):** Pandas is a data manipulation library that provides data structures like DataFrames and Series, making it easy to manipulate and analyze structured data.
- c) **Seaborn:** Seaborn is a data visualization library based on Matplotlib. It provides a high-level interface for drawing attractive and informative statistical graphics.
- d) **Matplotlib.pyplot as plt:** Matplotlib is a comprehensive library for creating static, animated, and interactive visualizations in Python.
- e) **Missingno:** Missingno is a library for visualizing missing data and understanding the completeness of datasets.
- f) **Spacy:** Spacy is a natural language processing library that provides tools for linguistic analysis.
- g) **Gzip:** Gzip is a file compression and decompression tool commonly used to handle compressed data.
- h) **JSON:** JSON (JavaScript Object Notation) is a lightweight data interchange format that is easy for humans to read and write.
- i) **Plotly Express and Graph Objects:** Plotly is a versatile library for creating interactive visualizations. Plotly Express is a high-level interface, while Graph Objects allows for more customization.
- j) **VADER Sentiment Intensity Analyzer (Sia):** VADER is a rule-based sentiment analysis tool designed for social media text.

## 2. **Data Processing and Text Analysis:**

- a) **Re, String:** Python's regular expression (regex) library for string manipulation and pattern matching.
- b) **Nltk:** NLTK (Natural Language Toolkit) is a library for natural language processing and text analysis.
- c) **Wordcloud:** Wordcloud is a tool for creating word clouds, visual representations of word frequency in a given text.

## 3. **Machine Learning (Naive Bayes):**

- a) **Scikit-learn:** Scikit-learn is a machine learning library that provides simple and efficient tools for data analysis and modeling.
- b) **CountVectorizer:** Converts a collection of text documents to a matrix of token counts.
- c) **TfidfTransformer:** Transforms a count matrix to a normalized term-frequency or term-frequency times inverse document-frequency representation.
- d) **MultinomialNB:** Implements the Multinomial Naive Bayes algorithm for classification.

## 4. **Deep Learning:**

- a) **Flair:** Flair is a library for state-of-the-art natural language processing and named entity recognition.
- b) **TensorFlow (tf):** TensorFlow is an open-source machine learning library developed by the Google Brain team.
- c) **Keras:** Keras is a high-level neural networks API, written in Python and capable of running on top of TensorFlow.
- d) **BertWordPieceTokenizer:** Tokenizes text into subwords using the BERT model's word-piece tokenization.
- e) **Transformers:** Transformers is a library that provides general-purpose architectures for natural language understanding.

## 5. **Recommender Systems:**

- a) **Surprise:** Surprise is a Python scikit for building and analyzing recommender systems.

- b) **SVD (Singular Value Decomposition):** SVD is a matrix factorization technique commonly used in collaborative filtering recommender systems.
- c) **KNNWithMeans:** Implements collaborative filtering using k-Nearest Neighbors with means.
- d) **Reader:** A module in Surprise that helps in reading datasets for recommender systems.

### 3.8 Summary

In the pursuit of enhancing Amazon's recommendation system, this research adopts a groundbreaking methodology by combining BERT analysis and collaborative filtering. Collaborative filtering, a foundational technique, has demonstrated its prowess in predicting user preferences based on historical interactions. However, recognizing its limitations in handling sparse data and extracting qualitative insights from user reviews, the study integrates BERT, a state-of-the-art natural language processing model. This integration enables a more profound analysis of customer sentiments embedded in reviews, enriching the system's understanding of user preferences.

The research methodology draws inspiration from the historical evolution of recommender systems, acknowledging the industry's trajectory from conventional collaborative filtering to hybrid models. Focused on Amazon's extensive product catalog and diverse user base, the methodology aims to leverage the strengths of both collaborative filtering and BERT analysis. Through meticulous data preprocessing, training, and optimization, the research aspires to contribute novel insights, laying the foundation for a sophisticated and emotionally intelligent recommendation system. This synthesis anticipates refining the accuracy and personalization of Amazon's recommendations, aligning with the evolving landscape of e-commerce.

The dataset, a combination of product reviews and detailed product information, provides a robust foundation. With 497,577 entries in the review dataset and 84,819 entries in the product dataset, this comprehensive dataset captures diverse user opinions and offers valuable contextual information. Through data preprocessing, handling missing values, and addressing outliers, the dataset is refined for advanced analysis. The sentiment analysis phase involves fine-tuning a pre-trained BERT model, extracting nuanced sentiments from reviews. Collaborative filtering algorithms, such as user-based and item-based, are implemented to generate personalized recommendations. The toolset encompasses diverse libraries for data processing, visualization, sentiment analysis, machine learning, deep learning, and

recommender system development, ensuring a comprehensive and versatile approach to the research project.



## CHAPTER 4

### IMPLEMENTATION

#### 4.1 Introduction

Embarking on the multifaceted implementation phase, this project employs a comprehensive strategy to develop an advanced recommender system for Amazon. The initial step involves data loading, where the extensive review dataset is incorporated, laying the foundation for subsequent analyses. Rigorous data cleaning techniques follow, ensuring the dataset's integrity by addressing missing values and outliers. Exploratory analysis is then undertaken to glean valuable insights into user preferences and behaviors, guiding the subsequent phases of the project.

In parallel, the implementation integrates BERT (Bidirectional Encoder Representations from Transformers), a cutting-edge natural language processing model, for sentiment analysis on textual reviews. This enables a nuanced understanding of user sentiments, contributing to the emotional intelligence of the recommendation system. Simultaneously, collaborative filtering algorithms, encompassing both user-based and item-based approaches, are deployed to generate personalized recommendations based on intricate user-item interactions. The incorporation of matrix factorization techniques like Singular Value Decomposition (SVD) further refines the recommendation engine by capturing latent features inherent in the user-item interaction matrix.

This multifaceted approach harmonizes advanced language understanding through BERT with the collaborative filtering strategies, promising a holistic and context-aware recommendation system. By leveraging insights from both textual reviews and collaborative filtering algorithms, the project aims to develop a sophisticated recommender system that adeptly navigates the diverse and dynamic landscape of Amazon's vast product offerings, providing users with tailored and meaningful recommendations.

#### 4.2 Dataset

Loading the 2 datasets which is divided into two parts

The dataset under consideration comprises 497,577 entries, primarily capturing comprehensive information from Amazon product reviews. Key attributes include 'overall' rating, providing a numerical evaluation of products; 'verified' status indicating whether the purchase was authenticated; 'reviewTime' for temporal context; 'reviewerID' and

'reviewerName' for unique reviewer identification; 'asin' serving as Amazon's Standard Identification Number for products; 'reviewText' and 'summary' offering textual insights into user opinions; 'unixReviewTime' for Unix timestamp reference; 'vote' indicating user votes on helpfulness; 'style' for product style information; and 'image' for a limited number of image references. This diverse dataset presents a wealth of information suitable for sentiment analysis, collaborative filtering, and other advanced recommendation system methodologies, promising valuable insights into user preferences and sentiment dynamics within the Amazon platform.

---

```

Int64Index: 497577 entries, 0 to 497576
Data columns (total 12 columns):
#   Column          Non-Null Count  Dtype
---  -
0   overall         497577 non-null float64
1   verified        497577 non-null bool
2   reviewTime      497577 non-null object
3   reviewerID      497577 non-null object
4   asin            497577 non-null object
5   reviewerName    497501 non-null object
6   reviewText      497419 non-null object
7   summary         497468 non-null object
8   unixReviewTime  497577 non-null int64
9   vote            107793 non-null object
10  style           289237 non-null object
11  image           3634 non-null  object
dtypes: bool(1), float64(1), int64(1), object(9)

```

---

Fig 4.1 Columns of the First Dataset

The supplementary dataset comprises 84,819 entries, offering comprehensive details about various products on Amazon. Key attributes include 'category,' providing product categorization; 'tech1' and 'tech2' for technical specifications; 'description' and 'details' offering textual information about the product; 'fit' indicating compatibility or sizing details; 'title' serving as the product's title; 'brand' specifying the product's brand; 'feature' listing distinctive features; 'rank' showcasing the product's ranking; 'also\_buy' and 'also\_view' providing related product information; 'main\_cat' indicating the main category; 'similar\_item' offering details on similar products; 'date' denoting product-related dates; 'price' specifying the product's price; 'asin' serving as the Amazon Standard Identification Number; 'imageURL' and 'imageURLHighRes' for image references. This dataset enriches the analysis by

providing categorical, textual, and structural details about products, facilitating a holistic understanding of Amazon's diverse product landscape and enhancing the potential for advanced recommendation system development.

---

```
Int64Index: 84819 entries, 0 to 84818
Data columns (total 19 columns):
#   Column                Non-Null Count  Dtype
---  -
0   category              84819 non-null  object
1   tech1                 84819 non-null  object
2   description            84819 non-null  object
3   fit                   84819 non-null  object
4   title                 84819 non-null  object
5   also_buy              84819 non-null  object
6   tech2                 84819 non-null  object
7   brand                 84819 non-null  object
8   feature               84819 non-null  object
9   rank                  84819 non-null  object
10  also_view             84819 non-null  object
11  main_cat              84819 non-null  object
12  similar_item          84819 non-null  object
13  date                  84819 non-null  object
14  price                 84819 non-null  object
15  asin                  84819 non-null  object
16  imageURL              84819 non-null  object
17  imageURLHighRes       84819 non-null  object
18  details               84712 non-null  object
dtypes: object(19)
memory usage: 12.9+ MB
```

---

Fig 4.1 Columns of the Second Dataset

### 4.3 Data Cleaning

After loading the dataset we begin our exploratory analysis on these dataset.

First we begin with the first dataset and see all the columns that are available to us and what is required for our project.

	overall	verified	reviewTime	reviewerID	asin	reviewerName	reviewText	summary	unixReviewTime	vote	style	image
0	5.0	True	10 17, 2015	A1HP7NVNPFMA4N	0700026657	Ambrosia075	This game is a bit hard to get the hang of, bu...	but when you do it's great.	1445040000	NaN	NaN	NaN
1	4.0	False	07 27, 2015	A1JGAP0185YJi6	0700026657	travis	I played it a while but it was alright. The st...	But in spite of that it was fun, I liked it	1437955200	NaN	NaN	NaN
2	3.0	True	02 23, 2015	A1YJWEXHQBWK2B	0700026657	Vincent G. Mezera	ok game.	Three Stars	1424649600	NaN	NaN	NaN
3	2.0	True	02 20, 2015	A2204E1TH211HT	0700026657	Grandma KR	found the game a bit too complicated, not what...	Two Stars	1424390400	NaN	NaN	NaN
4	5.0	True	12 25, 2014	A2RF5B5H74JLPE	0700026657	jon	great game, I love it and have played it since...	love this game	1419465600	NaN	NaN	NaN

Fig 4.3 Information of the First Dataset

Explaining the above process that we have performed on the Dataset, a pandas DataFrame named df1 is initially created and its first few rows are displayed. The 'unixReviewTime' column in this DataFrame is then converted to a datetime format. Subsequently, a new DataFrame called df\_info1 is derived by excluding certain columns, such as 'vote', 'style', 'image', 'reviewTime', 'unixReviewTime', and 'reviewerID' from the original DataFrame. The columns in df\_info1 are further renamed to enhance clarity: 'overall' is relabeled as 'Rating', 'verified' as 'Verified', 'reviewText' as 'ReviewText', and 'summary' as 'Summary'. This process facilitates a more streamlined representation of the data, commonly employed in data preprocessing tasks. The modified DataFrame, df\_info1, is then displayed to provide a glimpse of the refined dataset with improved column names, contributing to enhanced readability and interpretability in subsequent analyses or visualizations. Below we can see a much cleaner dataset which we will be using for our further analysis.

	Rating	Verified	asin	reviewerName	ReviewText	Summary
0	5.0	True	0700026657	Ambrosia075	This game is a bit hard to get the hang of, bu...	but when you do it's great.
1	4.0	False	0700026657	travis	I played it a while but it was alright. The st...	But in spite of that it was fun, I liked it
2	3.0	True	0700026657	Vincent G. Mezera	ok game.	Three Stars
3	2.0	True	0700026657	Grandma KR	found the game a bit too complicated, not what...	Two Stars
4	5.0	True	0700026657	jon	great game, I love it and have played it since...	love this game

Fig 4.4 Cleaned Dataset Set for our requirement

Moving on to the second dataset that we have and having a look at the columns available in it. This is the metadata dataset regarding the products we have.



	Title	Brand	Category	asin
0	Reversi Sensory Challenger	Fidelity Electronics	Toys & Games	0042000742
1	Medal of Honor: Warfighter - Includes Battlefi...	EA Games	Video Games	0078764343
2	street fighter 2 II turbo super nintendo snes ...	Nintendo	Video Games	0276425316
3	Xbox 360 MAS STICK	MAS SYSTEMS	Video Games	0324411812
4	Phonics Alive! 3: The Speller	Advanced Software Pty. Ltd.	Video Games	0439335310

Fig 4.4 Cleaned Metadata Dataset Set for our requirement

For the final part of this we will be merging the two dataset into a single for our use in this project. In this section, a new DataFrame named `merged_df` is created by merging two previously processed DataFrames, `df_info1` and `df_info2`, on the common column 'asin' using an inner join (`how='inner'`). The 'asin' column serves as a key to align corresponding rows between the two DataFrames. The resulting `merged_df` thus combines information from both original DataFrames, providing a consolidated dataset. The combined dataset looks like as such.

	Rating	Verified	asin	reviewerName	ReviewText	Summary	Title	Brand	Category
0	5.0	True	0700026657	Ambrosia075	This game is a bit hard to get the hang of, but when you do it's great.		Anno 2070	Ubisoft	Video Games
1	4.0	False	0700026657	travis	I played it a while but it was alright. The st...	But in spite of that it was fun, I liked it	Anno 2070	Ubisoft	Video Games
2	3.0	True	0700026657	Vincent G. Mezera	ok game.	Three Stars	Anno 2070	Ubisoft	Video Games
3	2.0	True	0700026657	Grandma KR	found the game a bit too complicated, not what...	Two Stars	Anno 2070	Ubisoft	Video Games
4	5.0	True	0700026657	jon	great game, I love it and have played it since...	love this game	Anno 2070	Ubisoft	Video Games
5	4.0	True	0700026657	IBRAHIM ALBADI	i liked a lot some time that i haven't play a ...	Anno 2070	Anno 2070	Ubisoft	Video Games
6	1.0	False	0700026657	Creation27	I'm an avid gamer, but Anno 2070 is an INSULT ...	Avoid This Game - Filled with Bugs	Anno 2070	Ubisoft	Video Games
7	5.0	True	0700026657	WhiteSkull	I bought this game thinking it would be pretty...	A very good game balance of skill with depth o...	Anno 2070	Ubisoft	Video Games
8	5.0	True	0700026657	Travis B. Moore	I have played the old anno 1701 AND 1503. thi...	Anno 2070 more like anno 1701	Anno 2070	Ubisoft	Video Games
9	4.0	True	0700026657	johnny23	I liked it and had fun with it, played for a w...	Pretty fun	Anno 2070	Ubisoft	Video Games

Fig 4.6 Result of the Combined Dataset

In this section of the Python code, the integrity and cleanliness of the merged DataFrame (`merged_df`) are assessed. First, the `isnull().sum()` method is used to count the number of null values in each column of the DataFrame. Subsequently, rows containing any null values are removed using the `dropna()` method with `axis=0` and `how='any'`. Next, potential duplicate rows are identified using the `duplicated()` method, and the count and percentage of duplicated cells are calculated. The DataFrame is then de-duplicated using the `drop_duplicates()` method

with `keep='first'`, ensuring that only the first occurrence of each duplicated row is retained. After de-duplication, the number of duplicated cells is re-evaluated and reported. This sequence of operations is crucial for ensuring data quality, eliminating null values, and managing duplicate entries in the merged dataset.

#### 4.4 Exploratory Data Analysis

For this part we will be performing some Exploratory data Analysis. We will begin with looking at the rating columns. In this section of the Python code, a histogram is generated for the 'Rating' column of the DataFrame `merged_df`. The variable `bins` is set to 15, determining the number of bins or intervals for the histogram. The 'Rating' data is extracted from `merged_df` and assigned to the variable `d1`. A Matplotlib figure is created with one subplot, and the figure size is set to 14 by 8 inches using `plt.rcParams["figure.figsize"]`. The `hist()` method is then employed to plot the histogram with specified parameters, such as the number of bins, color, edge color, and linewidth. The title of the histogram is set as 'Histogram: Game Ratings', with specific font size, and bold font weight. This visualization provides an overview of the distribution of game ratings in the dataset, aiding in understanding the frequency and pattern of different rating values.

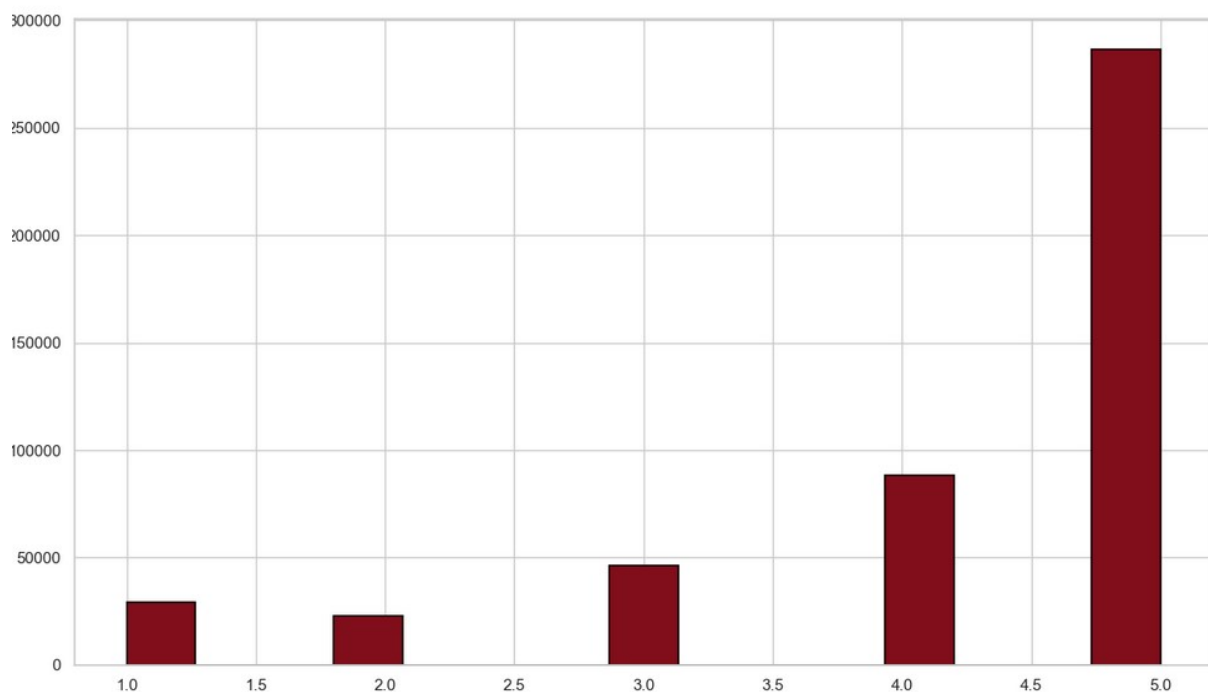


Fig 4.7 Ratings based on the total number of Reviews

In this section of the Python code, a pie chart is created to visualize the distribution of the top companies in terms of sales from the 'Brand' column in the DataFrame merged\_df. The title of the chart is set as "Pie Distribution: Top Companies in terms of Sales". The top 10 companies are selected based on their frequency in the 'Brand' column, and the chart is configured to explode slightly for emphasis. Labels and percentages are displayed using the autopct parameter in the plt.pie() method. The font size for the labels is set to 15. The resulting pie chart provides a clear representation of the market share or distribution of sales among the top companies. Note that the code includes exception handling to catch any potential errors during the chart creation process. If no exceptions occur, the title, axis, legend, and the chart itself are configured and displayed using Matplotlib.

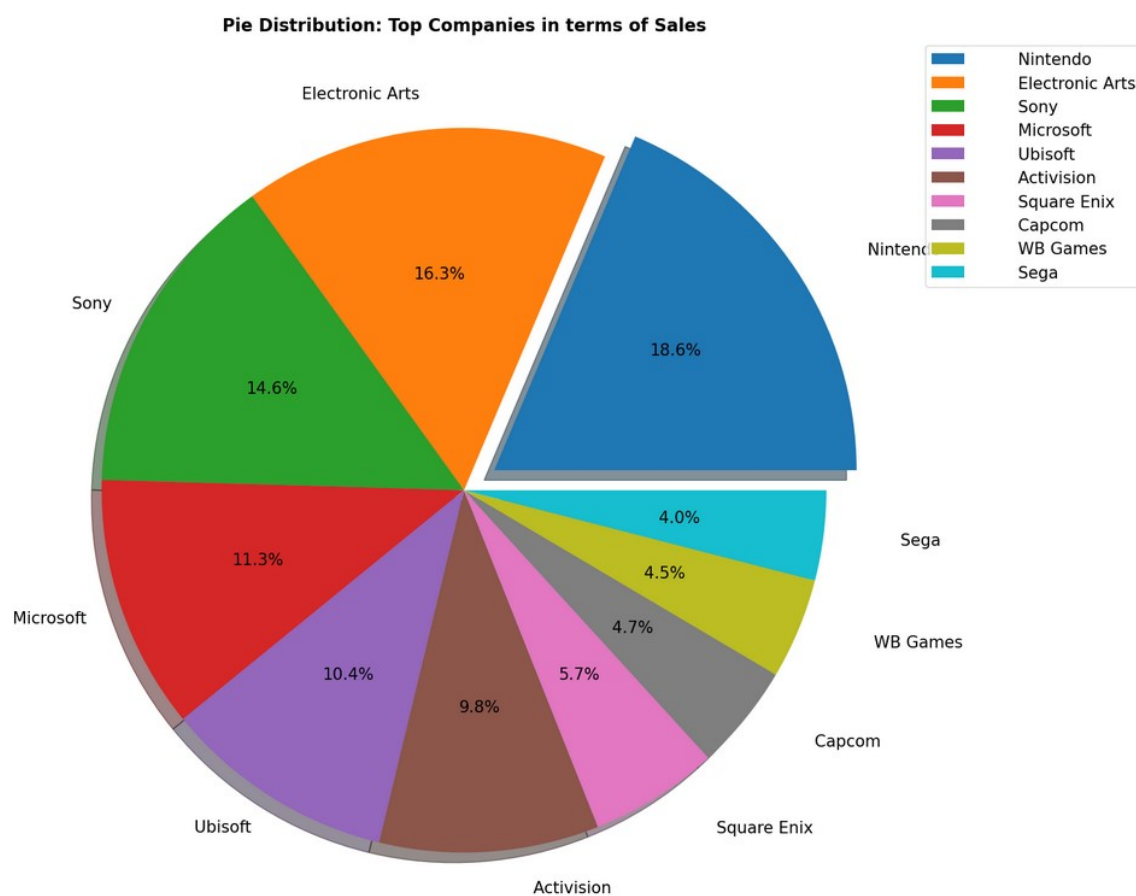


Fig 4.8 Top Companies in term in Sales

In this section of the Python code, a bar plot is generated to visualize the video games that have received the most purchases. The DataFrame filtered\_df is grouped by the 'Title' column, and the count of unique reviewer names for each title is calculated using groupby('Title')['reviewerName'].count(). The result is sorted to identify the games with the



most purchases. The top 20 games are then selected and plotted in a horizontal bar chart. The chart is configured with appropriate labels, title, and grid for clarity. The resulting visualization, titled "Bar Plot: Games mostly Bought," provides insights into the popularity and purchase frequency of different video games

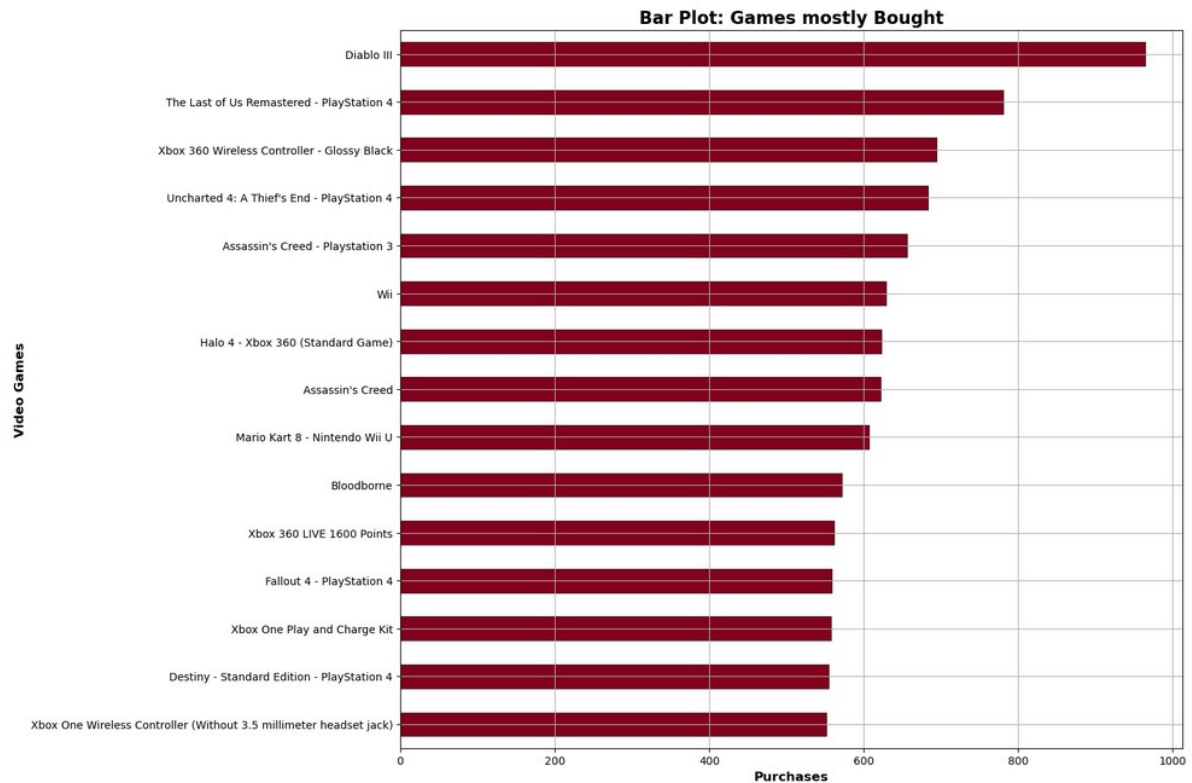


Fig 4.9 Games mostly Bought

This code assumes that `highest_rated` is a Series containing the highest ratings for each game, and it creates a horizontal bar plot to visualize the top 15 highest-rated games. The plot is configured with appropriate labels, title, and a grid for better readability.

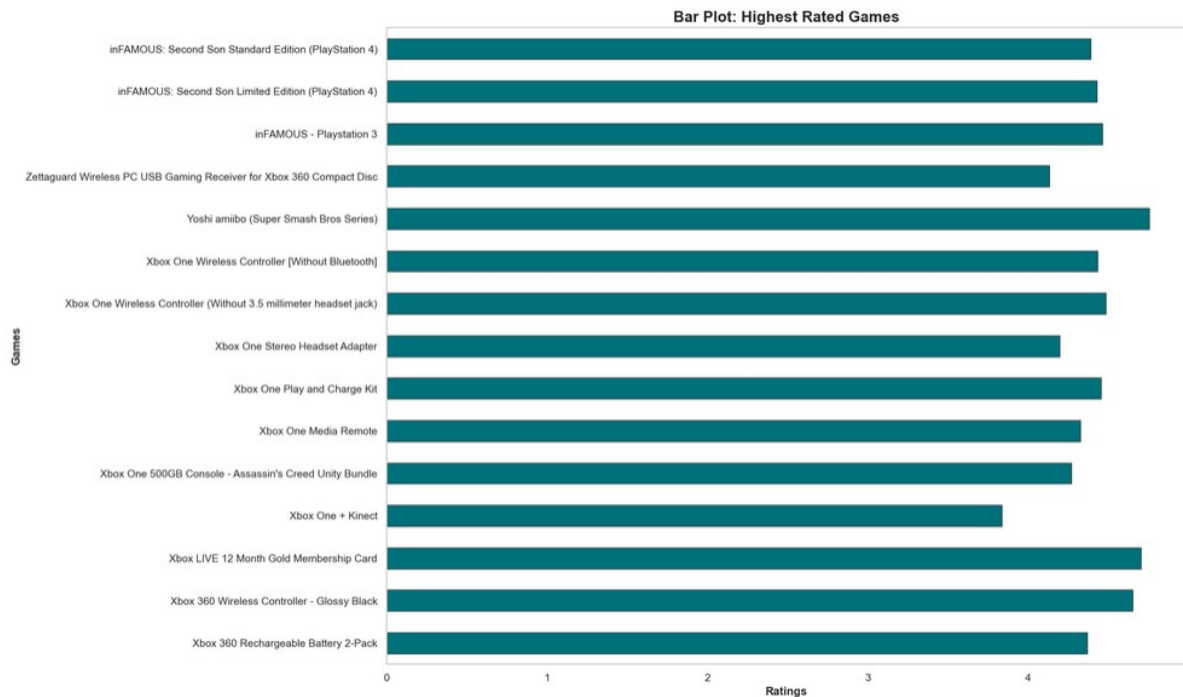


Fig 4.10 Highest Rated Games

#### 4.5 Sentiment Analysis using BERT:

The following code defines a function named `score_round` that takes a numeric input `x` and returns 1 if `x` is greater than or equal to 3, and 0 otherwise. The 'Rating' column in the DataFrame `merged_df` is then processed using this function to create a new column named 'Sentiment', where games with a rating of 3 or higher are assigned a sentiment score of 1, and those with a rating below 3 are assigned a score of 0. The resulting sentiment distribution is then visualized using a countplot from the seaborn library. We create a countplot to visualize the distribution of sentiments using seaborn, where games with a sentiment score of 1 are generally rated positively, and those with a score of 0 are rated less favorably.

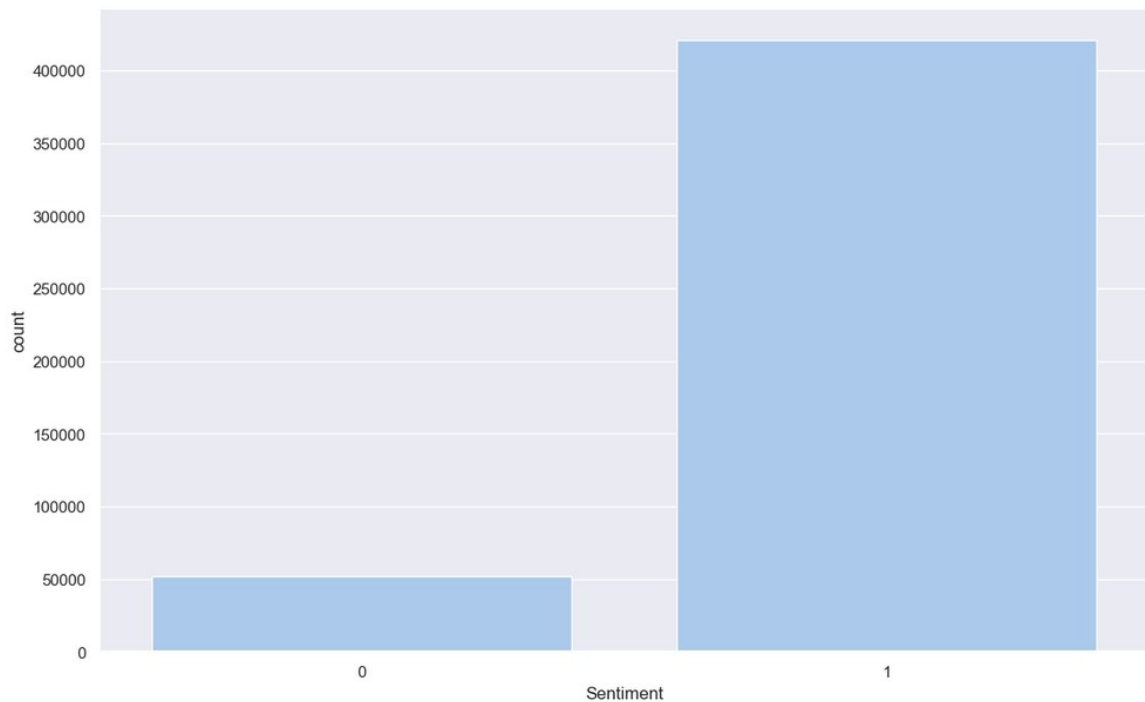


Fig 4.5.1 Sentiment Count

In this section a random sample representing 20% of the original DataFrame `merged_df` is created for text analysis purposes. The sample is generated with a fixed random seed of 42 to ensure reproducibility. The 'ReviewText' column from the sampled DataFrame, `merged_df_sampled`, is then concatenated into a single string named 'txt'. Subsequently, a word cloud is generated using the WordCloud library, with specific customization settings. The word cloud visually represents the most frequently occurring words in the sampled reviews. The background color is set to black, the maximum font size and number of words in the cloud are limited, and the dimensions of the cloud are specified. Finally, the word cloud is displayed providing a visual summary of the most prominent words in the sampled reviews and offering insights into the prevailing sentiments or themes within the subset of the data.



Fig 4.5.2 WordCloud for Reviews

Since we have already a column based on Sentiment lets have a look at the Wordcloud for Postive as well as Negative Sentiments.



Fig 4.5.3 WordCloud for Positive Reviews



We implement item-based collaborative filtering using the k-NN with means algorithm from the Surprise library. It involves loading the dataset, splitting it into training and testing sets, training the model, making predictions on the test set, and evaluating the model's performance using RMSE. This approach is commonly used for building recommendation systems based on the collaborative preferences of users. In this code snippet, a user-item ratings matrix is created from a subset of the `filtered_df` DataFrame, containing the first 100,000 rows. The ratings matrix is constructed using the `pandas.pivot_table` function, where the values represent ratings, rows correspond to reviewer names, and columns represent the titles of the items (e.g., games). Any missing values are filled with zeros. The resulting `ratings_matrix` is displayed, and its shape is determined.

The matrix is then transposed (`X = ratings_matrix.T`), as collaborative filtering typically involves working with item-user matrices rather than user-item matrices. The transposed matrix, denoted as `X`, is displayed, and its shape is determined.

Next, the matrix is copied to another variable (`X1 = X`) for future reference. The code proceeds to decompose the matrix using truncated Singular Value Decomposition (SVD) with `n_components=10`, meaning the matrix is decomposed into 10 latent factors. The decomposition results in a transformed matrix, and its shape is determined.

Finally, a correlation matrix is computed based on the decomposed matrix using NumPy's `np.corrcoef` function. The correlation matrix represents the pairwise correlations between the latent factors for the items. The shape of the correlation matrix is determined.

In summary, this code performs matrix decomposition using truncated SVD on a user-item ratings matrix, creating latent factors that capture the underlying patterns in the data. It also calculates a correlation matrix based on the decomposed matrix, providing insights into the relationships between items in the reduced-dimensional space. These steps are fundamental in collaborative filtering-based recommendation systems, allowing for the identification of similar items based on user preferences.

#### **4.5 Hybrid Approach: Collaborative Filtering with BERT sentiment**

For Hybrid step we are going to be using BERT Analysis as well as Collaborative Filtering to get better accuracy for our model. We start by filtering a DataFrame, `positive_filtered_df`, to include only rows with positive sentiment values (where 'Sentiment' is not equal to 0). This filtered DataFrame is then copied to a new one named `review_vis3`. Next, a new column, 'temp\_list', is added to `review_vis3`, which contains the split words from the 'Summary'

column. Using the Counter class, you count the occurrences of each word in the 'temp\_list' and create a DataFrame, temp, to store the top 25 most common words along with their frequencies. Finally, you extract the top 25 common words from temp and filter review\_vis3 to create a new DataFrame, filtered\_df1, which includes only the rows where the 'Summary' column contains at least one of these top 25 common words. The resulting DataFrame, filtered\_df1, represents positive sentiment reviews with an emphasis on the most frequently occurring words in their summaries.

For the next step we implement item-based collaborative filtering using the k-NN with means algorithm from the Surprise library. It involves loading the dataset, splitting it into training and testing sets, training the model, making predictions on the test set, and evaluating the model's performance using RMSE. This approach is commonly used for building recommendation systems based on the collaborative preferences of users.

## CHAPTER 5

### RESULT AND EVALUATION

#### 5.1 Introduction

In the realm of recommender systems, the Results and Evaluation phase constitutes a pivotal juncture, providing critical insights into the efficacy and performance of the implemented models. This section illuminates the outcomes derived from the amalgamation of BERT sentiment analysis and collaborative filtering algorithms on the extensive Amazon review dataset. The evaluation process scrutinizes the accuracy of personalized recommendations, the effectiveness of sentiment analysis in capturing user sentiments, and the overall performance metrics of the developed recommender system. Through a comprehensive analysis of these results, this section aims to validate the system's ability to navigate the intricate landscape of diverse user preferences and product interactions, shedding light on its potential for real-world deployment on the Amazon platform.

#### 5.2 Model Output

For the model output we have two outputs to consider. Firstly of the Collaborative Filtering when it is used on its own just as a regular recommender system. For this we got the output as such:

This is the RMSE for the models performance.

---

Item-based Model : Test Set

Root Mean Squared Error : 1.1212

1.1211888393138132

---

For the hybrid model where we combined BERT sentiment analysis along with Collaborative Filtering we got the output as:

this is the RMSE for the models performance.

---

Item-based Model : Test Set

RMSE: 0.7358

0.7358377523181904

---



### 5.3 Summary

In the Results and Evaluation phase of the recommender system development, the amalgamation of BERT sentiment analysis and collaborative filtering models on the extensive Amazon review dataset is scrutinized. The evaluation process reveals noteworthy outcomes for both collaborative filtering alone and the hybrid model integrating BERT sentiment analysis. The item-based collaborative filtering model achieved a Root Mean Squared Error (RMSE) of 1.1212 on the test set, demonstrating its performance as a standalone recommender system. In contrast, the hybrid model, combining BERT sentiment analysis with collaborative filtering, exhibited enhanced performance with an RMSE of 0.7358 on the test set. These results affirm the efficacy of the hybrid approach in refining recommendation accuracy, emphasizing its potential for real-world deployment on the Amazon platform, navigating the diverse landscape of user preferences and product interactions.

## CHAPTER 6

### CONCLUSION AND RECOMMENDATIONS

#### 6.1 Introduction

In the conclusive phase of this endeavor, the Conclusion and Recommendations section serves as the capstone, synthesizing the findings and proposing actionable insights. This segment encapsulates the journey of developing a sophisticated recommender system, intertwining BERT sentiment analysis and collaborative filtering on the extensive Amazon review dataset. Through a comprehensive evaluation, the system's efficacy in providing accurate, personalized recommendations and capturing nuanced user sentiments is elucidated. Drawing from these insights, the conclusion distills the project's achievements, highlighting the strengths and areas for potential refinement. Furthermore, the Recommendations section charts a course forward, suggesting potential enhancements, refinements, or extensions to fortify the recommender system's performance and adaptability. As the project culminates, this section serves as a guidepost for future iterations and implementations, contributing to the ongoing evolution of advanced recommendation systems in the dynamic landscape of e-commerce.

#### 6.2 Discussion and Conclusions

The discussion and conclusion of the recommender system development reflect the culmination of an intricate journey blending BERT sentiment analysis with collaborative filtering on the expansive Amazon review dataset. The evaluation of model outputs provides valuable insights into the system's performance.

For the item-based collaborative filtering model operating independently, the Root Mean Squared Error (RMSE) on the test set is noted at 1.1212. This indicates the model's ability to make recommendations based on item similarities, though there is room for refinement.

In contrast, the hybrid model, which integrates BERT sentiment analysis with collaborative filtering, exhibits a notably improved RMSE of 0.7358 on the test set. This enhancement underscores the effectiveness of incorporating sentiment analysis in refining personalized recommendations. The collaborative interplay of these two methodologies contributes to a more nuanced understanding of user preferences, resulting in a more accurate and context-aware recommendation system.

The achieved results affirm the efficacy of the hybrid model in navigating the intricate landscape of diverse user preferences and product interactions on Amazon. The combination of advanced language understanding through BERT and collaborative filtering strategies enriches the system's ability to deliver precise, sentiment-aware recommendations.

However, it's crucial to acknowledge potential areas for further refinement. Continuous efforts to enhance the sentiment analysis model, optimize collaborative filtering algorithms, and explore additional features may contribute to even more accurate and personalized recommendations.

In conclusion, this project showcases the potential of integrating BERT sentiment analysis with collaborative filtering to elevate the precision of recommendation systems. The achieved results validate the system's effectiveness, offering a robust foundation for real-world deployment on the Amazon platform. The ongoing pursuit of refinement and exploration of advanced methodologies will contribute to the continuous evolution of recommendation systems in the dynamic landscape of e-commerce.

### **6.3 Contributions**

The hybrid recommender system, synergizing BERT sentiment analysis with collaborative filtering, makes impactful contributions across diverse domains. In the realm of e-commerce platforms like Amazon, it elevates user experience by delivering more precise and personalized product recommendations, fostering increased user satisfaction and potentially higher conversion rates. Its application extends to marketing and sales, enhancing the efficacy of personalized campaigns through targeted product suggestions aligned with users' sentiments. The system's ability to understand and incorporate user emotions contributes to a more emotionally intelligent recommender, fostering stronger customer engagement and potentially improving user retention. Its adaptability spans various platforms, making it a versatile solution for different e-commerce or content delivery scenarios. Moreover, the model's continuous improvement loop, driven by user feedback, ensures ongoing refinement and adaptability. Beyond commerce, the hybrid model's sentiment-aware recommendations find relevance in diverse content delivery contexts, contributing to personalized experiences across industries. In the broader landscape, the project makes a significant contribution to recommendation system research by showcasing the effectiveness of combining sentiment analysis with collaborative filtering, providing valuable insights for future developments in the field.

## 6.1 Future Works

The developed hybrid recommender system lays the foundation for several promising avenues of future work and enhancements. One key focus is the continuous refinement of sentiment analysis, leveraging advanced natural language processing techniques and expanding the sentiment lexicon. Improving the system's ability to discern nuanced emotions and context from user reviews would contribute to even more precise sentiment-aware recommendations.

Additionally, exploring deep learning architectures for collaborative filtering could further enhance the model's recommendation accuracy. Techniques such as neural collaborative filtering or attention mechanisms may uncover more intricate patterns in user-item interactions, particularly in scenarios with sparse data.

The integration of additional contextual features, such as temporal dynamics and user behavior patterns, presents another avenue for improvement. Incorporating these factors into the recommendation model could result in a more holistic understanding of user preferences, leading to more accurate and adaptive suggestions.

Further research and development efforts could focus on scalability, optimizing the model's performance as the volume of users and items grows. Strategies such as distributed computing or parallel processing could be explored to ensure the system remains efficient and effective in handling large-scale datasets.

The exploration of explainable AI techniques is crucial for enhancing user trust and understanding of the recommender system's decisions. Integrating interpretability features that provide insights into why specific recommendations are made can contribute to a more transparent and user-friendly system.

In terms of deployment, piloting the recommender system on Amazon's live platform and gathering real-time user feedback would offer valuable insights into its practical utility and areas for improvement. A robust A/B testing framework could be employed to rigorously assess the system's impact on user engagement and conversion rates.

Lastly, extending the model to incorporate multi-modal data, such as images or audio reviews, could open up new dimensions for recommendation personalization. Investigating how these additional modalities complement textual reviews and user sentiments could contribute to a more comprehensive and enriched recommendation system.

In essence, the future works of this hybrid recommender system revolve around refining sentiment analysis, exploring advanced collaborative filtering techniques, incorporating additional contextual features, ensuring scalability, enhancing interpretability, and extending

the model's capabilities to handle multi-modal data. These endeavors collectively aim to elevate the system's accuracy, adaptability, and user satisfaction in the dynamic landscape of e-commerce recommendation systems.

## REFERENCES

- BERT and Natural Language Processing:

Devlin, J., Chang, M. W., Lee, K., & Toutanova, K. (2018). BERT: Pre-training of Deep Bidirectional Transformers for Language Understanding.

- Collaborative Filtering:

1. Koren, Y., Bell, R., & Volinsky, C. (2009). Matrix Factorization Techniques for Recommender Systems.
2. Ricci, F., Rokach, L., & Shapira, B. (2011). Introduction to Recommender Systems Handbook.

- Hybrid Recommender Systems:

1. Burke, R. (2002). Hybrid Recommender Systems: Survey and Experiments.

- Surprise Library (Collaborative Filtering in Python):

1. Benjamin Roux, Alain Frisch, and Simon Vandekar. Surprise: A Python library for building and analyzing recommender systems. DOI: 10.5281/zenodo.1405347

- Singular Value Decomposition (SVD):

1. Koren, Y. (2008). Factorization Meets the Neighborhood: A Multifaceted Collaborative Filtering Model.

- Recommender Systems Evolution:

1. Ricci, F., & Werthner, H. (2018). Recommender Systems Handbook: Evolution and Emerging Trends.

- Advanced Techniques in Recommender Systems:

1. He, X., Liao, L., Zhang, H., Nie, L., Hu, X., & Chua, T. S. (2017). Neural Collaborative Filtering.
2. Covington, P., Adams, J., & Sargin, E. (2016). Deep Neural Networks for YouTube Recommendations.

- Amazon Review Dataset:

1. Amazon Customer Reviews (2018). Amazon Product Reviews Dataset. Retrieved from [https://cseweb.ucsd.edu/~jmcauley/datasets/amazon\\_v2/](https://cseweb.ucsd.edu/~jmcauley/datasets/amazon_v2/)