

Book Recommendation System on Amazon Dataset

by Chinmay Bidarkar

Submission date: 02-May-2024 12:16AM (UTC+0530)

Submission ID: 2368058938

File name: 2020B5A71863H-2020B5A72178H-2020B5A70987H_1.pdf (302.56K)

Word count: 4682

Character count: 26255

Enhancing Book Recommendation Systems with Associative Data Mining Techniques

BITS Pilani Hyderabad Campus

CS F415 Data Mining Project

Chinmay Bidarkar

f20201863@hyderabad.bits-pilani.ac.in

Dev Gaur

f20202178@hyderabad.bits-pilani.ac.in

P Sai Aditya

f20200987@hyderabad.bits-pilani.ac.in

April 28, 2024

Abstract

The exponential growth of literature available online has introduced the challenge of identifying relevant books tailored to individual preferences, an overwhelming task for readers worldwide. This paper presents an innovative approach to enhancing book recommendation systems by employing associative data mining techniques to analyze user-generated book ratings from a substantial Amazon dataset. Our methodology involves applying three distinct algorithms—Apriori, FP-growth, and Eclat—to identify patterns in book ratings, thereby discovering meaningful relationships among books rated similarly by users. By introducing pruning strategies within the Apriori method and experimenting with various support and confidence thresholds, our work endeavors to not only recommend books aligned with user interests but also to explore the comparative efficiency of these algorithms in terms of time and space complexity. Our results indicate significant potential in associative data mining to discern trends in reading behavior, facilitating the delivery of more accurate book recommendations. Further, we explore the association of books within the same authorial or genre clusters and the frequency of books across multiple itemsets, aiming to ascertain the influence of itemset frequency on a book's popularity. The insights derived from these observations contribute to a broader understanding of consumer preferences and the technology driving modern recommendation systems.

Keywords: Associative Data Mining, Book Recommendation Systems, Amazon Book Ratings, Apriori Algorithm, FP-growth, Eclat Algorithm, User Preferences.

1 Introduction

1.1 Problem and Objective

The digitization of literature has exponentially increased the number of books accessible to readers, but it has also magnified the challenge of discovering books that align with personal tastes from an ever-expanding online library. The central problem that this research addresses is the inefficacy of current book recommendation systems in providing accurate, personalized book suggestions. The objective is to enhance the performance of these systems through associative data mining techniques, utilizing a comprehensive dataset of book ratings, which include intricate user interactions and book metadata from Amazon's expansive collection.

1.2 Importance of the Study

The importance of this study is multifaceted. It is significant for enhancing the user experience in online book discovery, increasing the efficiency and pleasure of finding new reads in the digital space. By providing more accurate recommendations, we aim to facilitate a stronger connection between readers and books, which is vital for maintaining and growing the culture of reading in the digital age. The study also holds importance for the commercial sector, where book sellers and publishers can benefit from the improved engagement and satisfaction of readers, potentially leading to increased sales and customer loyalty.

1.3 Challenges

The prediction of user preferences is a task fraught with technical hurdles. The sparsity of user-rating matrices and the high dimensionality of the data pose considerable challenges to conventional recommendation algorithms. The difficulty lies in capturing the subtle preferences and reading patterns from limited user interactions, often exacerbated by the cold-start problem where new users or items have few or no associated ratings. These challenges necessitate the exploration of scalable and effective algorithms that can handle large volumes of data, draw meaningful correlations, and adapt to new trends without significant overhead.

1.4 Previous Solutions and Our Approach

Existing recommendation systems, such as collaborative and content-based filtering, have provided frameworks to suggest books based on user behavior or book characteristics. However, these systems often suffer from limitations such as the inability to handle new books without ratings (the cold-start problem) or to scale with the dataset efficiently. Our approach diverges by integrating three distinct associative data mining algorithms: Apriori, FP-growth, and Eclat. These algorithms are known for their ability to uncover association rules in large datasets, but their application in recommendation systems, particularly for books, is not extensively studied. We propose a novel implementation and comparison of these algorithms to extract meaningful patterns from complex user-book interactions, aiming to generate more personalized book recommendations.

1.5 Key Components of Our Approach and Anticipated Results

This research is predicated on a thorough examination and preprocessing of the data, where the quality of input significantly influences the output of associative algorithms. We will explore data cleaning techniques to handle missing values, normalization to adjust for bias in user ratings,

and dimensionality reduction methods to enhance the algorithms' performance. Furthermore, we will employ a rigorous comparison of the algorithms' performance based on metrics such as precision, recall, and computational efficiency. We anticipate that our findings will contribute substantially to the field of recommendation systems by providing insights into the scalability and applicability of associative data mining techniques in a real-world context. The anticipated outcome is a set of associative rules that, when applied to the user profiles, can predict their preferences with a higher degree of accuracy than current models, thereby pushing the frontiers of personalized recommendation systems.

2 Related Work

The landscape of recommendation systems is vast and has been explored using various machine learning algorithms across different domains. Artificial Neural Networks (ANNs), as studied by Maghari et al. in predicting book ratings, have shown significant promise due to their ability to model complex, non-linear relationships within data . Similar approaches have been employed in domains ranging from health diagnostics to agricultural classifications, demonstrating the versatility of ANN in handling diverse datasets . Building on this foundation, the present study extends the application of data mining techniques by specifically leveraging associative algorithms—Apriori, FP-growth, and Eclat—to enhance the recommendation of books. These methods, traditionally utilized in transactional analysis, are adapted to the unique challenges presented by literary recommendation scenarios.

3 Approach/Methodology

3.1 Problem Tackling

In the domain of book recommendation systems, the central challenge is associating books with a user's preference list, which is inferred from the ratings provided. The intent is to use these associations not just to understand the user's past preferences but to also predict and suggest new books that the user is likely to appreciate.

3.2 Data Acquisition

Critical to this study is the collection of comprehensive data that spans books, users, and the ratings connecting them. The dataset needed to encompass substantial and relevant metadata capable of supporting sophisticated associative analysis. We have sourced such a dataset from Kaggle, known as the "Amazon Books Dataset," which is rich in the attributes required for a multifaceted examination and application of association rule mining.

3.3 Dataset Properties

The dataset under investigation comprises three distinct components: books, users, and ratings. The books dataset includes ISBNs, titles, authors, publication years, and publishers, providing a broad scope for understanding the literature. The users dataset contains anonymized user IDs and corresponding age information, enabling demographic analysis. The ratings dataset bridges users and books, containing user IDs, ISBNs, and book ratings, which serve as the foundation for discovering user preferences and predicting future book recommendations.

3.4 Choice of Algorithms

For the purpose of mining association rules, we have selected three foundational algorithms: Apriori, FP-growth, and Eclat. These algorithms are the bedrock of association rule learning and provide a robust framework for uncovering patterns in transactional datasets. Apriori is celebrated for its simplicity and interpretability, FP-growth for its efficiency in handling large datasets without candidate generation, and Eclat for its depth-first search advantages and speed. Their selection is predicated on their suitability for the dataset at hand and their proven track record in similar applications. These algorithms are particularly apt for our objective due to their ability to efficiently process the voluminous and sparse nature of user-book interactions, uncovering the hidden patterns that reflect genuine user preferences.

4 Experiments

4.1 Dataset

The dataset is a combination of three different sets: books, users, and ratings, which were merged into a single comprehensive dataset for associative analysis.

4.1.1 Pre-processing Methods

The pre-processing steps performed on the combined dataset to prepare it for analysis included:

- Merging the individual books, users, and ratings datasets into one unified structure.
- Converting date strings into consistent datetime objects for subsequent time-series analysis.
- Cleaning and standardizing numeric fields by correcting formats and converting data types, such as changing price strings with commas to float values.
- Filtering out transactions with non-positive quantities or prices to focus only on valid entries.
- Eliminating rows missing critical information, such as item names, to ensure the integrity of the transactional data.
- Imputing missing customer IDs with a placeholder to maintain the dataset's completeness.
- Calculating the total price for each transaction, reflecting the complete financial exchange.

4.1.2 Final Processed Dataset

The final processed dataset is characterized by:

- An aggregation of attributes from the merged books, users, and ratings data, including user IDs, ISBNs, book ratings, ages of users, book titles, authors, publication years, and publishers.
- Each record representing a user's rating of a book, coupled with additional demographic and bibliographic details.
- The dataset's structure allows for a multifaceted associative analysis to discover patterns across user demographics, book preferences, and publication timelines.

4.2 Evaluation Method / Metrics

To evaluate the efficacy of the proposed associative data mining methodology, which employs the Apriori, Eclat, and FP-Growth algorithms, we will consider the following evaluation methods and metrics:

4.2.1 Apriori Algorithm³

The Apriori algorithm is a classic algorithm in data mining that is used for mining frequent itemsets and deriving association rules from a dataset. It operates on a breadth-first search principle, where frequent subsets are extended one item at a time (a step known as candidate generation), and groups of candidates are tested against the data. The algorithm inherently employs two metrics:

- **Support:** It measures the frequency or the proportion of transactions that include the itemset. For the Apriori algorithm, the support threshold is crucial for pruning the search space.
- **Confidence:** This metric evaluates the likelihood of occurrence of the consequent in transactions under the condition that these transactions also contain the antecedent.

4.2.2 Eclat Algorithm³⁷

The Eclat algorithm stands for Equivalence Class Transformation, operating on a depth-first search approach. Unlike Apriori, Eclat uses a vertical data format to store itemsets and their corresponding transaction ids, making it generally faster and more space-efficient than the former. It evaluates the intersection of transaction id sets to find frequent itemsets. The same metrics of support and confidence are used to assess the strength of association rules derived from Eclat.

4.2.3 FP-Growth Algorithm²¹

The FP-Growth (Frequent Pattern Growth) algorithm is another method for finding frequent itemsets without candidate generation. It uses a compressed tree structure called the FP-tree to encode the dataset. The algorithm recursively divides the compressed dataset into a set of conditional databases, each one associated with one frequent item or itemset, and mines these smaller databases separately. Key metrics for evaluating the FP-Growth algorithm include:

- **Execution Time:** As FP-Growth is designed to be faster than Apriori, particularly on large datasets, we will measure its execution time compared to the other algorithms.
- **Memory Usage:** The compactness of the FP-tree and the lack of candidate generation can lead to substantial memory savings, which we will quantify and compare against the other algorithms.

For a comprehensive evaluation, we will also consider:

- **Scalability:** The ability of the algorithm to handle increasingly large datasets.
- **Accuracy:** Measured by the relevance of the generated rules to the test datasets.

By comparing these metrics across algorithms, we will discern the most efficient and effective method for our book recommendation system.

4.3 Experimental Setup

In our experimental setup, we have focused on the characteristics and properties of the Apriori, Eclat, and FP-Growth data mining techniques, with particular attention to the hyperparameters of support and confidence.

4.3.1 Support and Confidence

Support and confidence are two fundamental metrics used in our study to determine the strength and reliability of association rules within our dataset.

- **Support** is the proportion of transactions in the dataset that contain a particular itemset. In our experiments, the support value determines how frequently a combination of items appears together in transactions. We used a relative support threshold to filter out itemsets that are not common enough in the dataset.
- **Confidence** measures the likelihood that an item B is also bought if an item A is bought. It is defined as the proportion of transactions with item A that also contain item B. The confidence threshold is used to ensure that only the rules with a high probability of co-occurrence are considered. In our experiments, this metric filters out rules that are less likely to be useful in recommendations.

Both support and confidence were used to prune the number of potential association rules within the algorithms. The specific thresholds for these parameters were carefully chosen based on the dataset characteristics and the objective of balancing rule accuracy with coverage.

4.3.2 Hyperparameters

In our experiments, the following hyperparameters were employed:

- **Minimum Support Threshold:** This threshold was set to a value that allowed us to uncover significant patterns in the dataset without being too restrictive. For the Eclat algorithm, we established a minimum support threshold of 500 to identify frequent itemsets within the dataset.
- **Minimum Confidence Threshold:** To ensure the reliability of the rules generated, we set a minimum confidence threshold of 0.4. This threshold was chosen to capture strong association rules while avoiding spurious relationships.

The chosen values for these hyperparameters were based on the distribution of the dataset and aimed at optimizing the balance between finding enough frequent itemsets and association rules, and maintaining high levels of confidence in the rules generated. These parameters directly impact the number of rules discovered and the computational efficiency of the algorithms.

By adjusting these hyperparameters, we can fine-tune our data mining techniques to adapt to the specific nuances of the Amazon book dataset and ensure that the most relevant and actionable insights are extracted.

5 Results

Utilizing the three algorithms, we have extracted frequent itemsets and derived association rules. The algorithms employed include ECLAT, Apriori, and the FP-Growth algorithm. We established a minimum confidence threshold of 0.4 and a minimum support threshold of 30 for our

analysis. Presented here are the ten most frequently occurring books within the Amazon dataset, accompanied by their respective support counts, for illustrative purposes.

Table 1: Top 10 Frequent Itemsets Identified by ECLAT

Itemset	Support
{The Lovely Bones: A Novel}	471
{Wild Animus}	467
{The Da Vinci Code}	377
{Angels & Demons}	253
{Bridget Jones's Diary}	277
{The Secret Life of Bees}	260
{Harry Potter and the Chamber of Secrets (Book 2)}	242
{Life of Pi}	247
{The Nanny Diaries: A Novel}	246
{Harry Potter and the Sorcerer's Stone (Harry Potter (Paperback))}	218

5.1 Associate Mining Rules

The following are the association mining rules, wherein the first column lists the antecedents and the subsequent column the consequents, accompanied by the confidence levels indicative of the rule's reliability. This implies that the purchase of a book listed under the antecedent column increases the likelihood of the purchase of the book in the consequent column. Therefore, these rules can be effectively utilized to predict sets of closely related items.

Table 2: Association Rules

Antecedent	Consequent	Confidence
"Harry Potter and the Sorcerer's Stone (Book 1)"	"Harry Potter and the Chamber of Secrets (Book 2)"	0.55
"Harry Potter and the Prisoner of Azkaban (Book 3)"	"Harry Potter and the Chamber of Secrets (Book 2)"	0.62
"Harry Potter and the Chamber of Secrets (Book 2)"	"Harry Potter and the Prisoner of Azkaban (Book 3)"	0.54
"Harry Potter and the Goblet of Fire (Book 4)"	"Harry Potter and the Chamber of Secrets (Book 2)"	0.61
"Harry Potter and the Chamber of Secrets (Book 2)"	"Harry Potter and the Goblet of Fire (Book 4)"	0.47
"The Two Towers (The Lord of the Rings, Part 2)"	"The Fellowship of the Ring (The Lord of the Rings, Part 1)"	0.58
"The Fellowship of the Ring (The Lord of the Rings, Part 1)"	"The Two Towers (The Lord of the Rings, Part 2)"	0.42
"The Return of the King (The Lord of the Rings, Part 3)"	"The Fellowship of the Ring (The Lord of the Rings, Part 1)"	0.56
"The Subtle Knife (His Dark Materials, Book 2)"	"The Golden Compass (His Dark Materials, Book 1)"	0.51
"The Amber Spyglass (His Dark Materials, Book 3)"	"The Golden Compass (His Dark Materials, Book 1)"	0.56
"Harry Potter and the Sorcerer's Stone (Book 1)"	"Harry Potter and the Prisoner of Azkaban (Book 3)"	0.49
"Harry Potter and the Sorcerer's Stone (Book 1)"	"Harry Potter and the Goblet of Fire (Book 4)"	0.42
"Harry Potter and the Prisoner of Azkaban (Book 3)"	"Harry Potter and the Goblet of Fire (Book 4)"	0.55

Table 3: Association Rules (Continued)

Antecedent	Consequent	Confidence
"Harry Potter and the Goblet of Fire (Book 4)"	"Harry Potter and the Prisoner of Azkaban (Book 3)"	0.64
"The Amber Spyglass (His Dark Materials, Book 3)"	"The Subtle Knife (His Dark Materials, Book 2)"	0.67
"The Subtle Knife (His Dark Materials, Book 2)"	"The Amber Spyglass (His Dark Materials, Book 3)"	0.44
"Harry Potter and the Goblet of Fire (Book 4)"	"Harry Potter and the Order of the Phoenix (Book 5)"	0.40
"Harry Potter and the Order of the Phoenix (Book 5)"	"Harry Potter and the Goblet of Fire (Book 4)"	0.41
"Tribulation Force: The Continuing Drama of Those Left Behind (Left Behind No. 2)"	"Left Behind: A Novel of the Earth's Last Days (Left Behind No. 1)"	0.54
"The Vampire Lestat (Vampire Chronicles, Book II)"	"Interview with the Vampire"	0.43
"The Queen of the Damned (Vampire Chronicles (Paperback))"	"The Vampire Lestat (Vampire Chronicles, Book II)"	0.53
"The Tale of the Body Thief (Vampire Chronicles (Paperback))"	"The Vampire Lestat (Vampire Chronicles, Book II)"	0.47
"The Return of the King (The Lord of the Rings, Part 3)"	"The Two Towers (The Lord of the Rings, Part 2)"	0.62
"The Two Towers (The Lord of the Rings, Part 2)"	"The Return of the King (The Lord of the Rings, Part 3)"	0.47
"The Queen of the Damned (Vampire Chronicles (Paperback))"	"Interview with the Vampire"	0.44

Table 4: Association Rules (Continued)

Antecedent	Consequent	Confidence
"The Tale of the Body Thief (Vampire Chronicles (Paperback))"	"Interview with the Vampire"	0.47
"The Queen of the Damned (Vampire Chronicles (Paperback))"	"The Tale of the Body Thief (Vampire Chronicles (Paperback))"	0.42
"The Tale of the Body Thief (Vampire Chronicles (Paperback))"	"The Queen of the Damned (Vampire Chronicles (Paperback))"	0.47
"High Five (A Stephanie Plum Novel)"	"Hot Six : A Stephanie Plum Novel (A Stephanie Plum Novel)"	0.45
"Hot Six : A Stephanie Plum Novel (A Stephanie Plum Novel)"	"High Five (A Stephanie Plum Novel)"	0.59
"High Five (A Stephanie Plum Novel)"	"Four To Score (A Stephanie Plum Novel)"	0.42
"Four To Score (A Stephanie Plum Novel)"	"High Five (A Stephanie Plum Novel)"	0.46
"Wizard and Glass (The Dark Tower, Book 4)"	"The Drawing of the Three (The Dark Tower, Book 2)"	0.57
"The Drawing of the Three (The Dark Tower, Book 2)"	"Wizard and Glass (The Dark Tower, Book 4)"	0.41
"The Drawing of the Three (The Dark Tower, Book 2)"	"The Gunslinger (The Dark Tower, Book 1)"	0.54
"The Gunslinger (The Dark Tower, Book 1)"	"The Drawing of the Three (The Dark Tower, Book 2)"	0.49

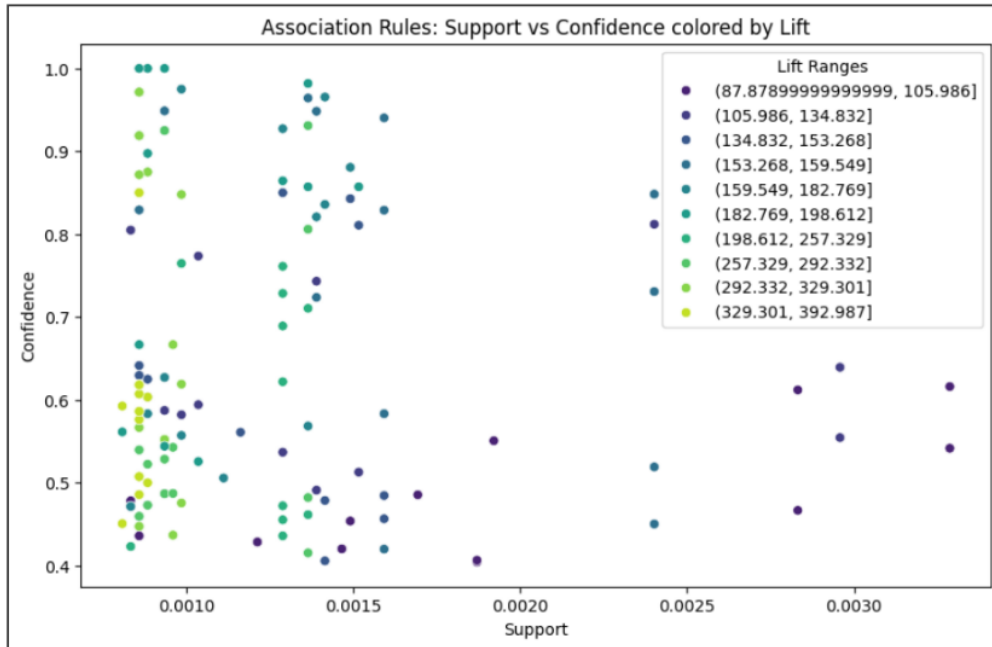


Figure 1: Support vs Confidence

This graph is a scatter plot that visualizes the relationship between the support and confidence of association rules and is colored by the lift metric:

- **Support (X-axis):** This represents the proportion of transactions in the data that contain the rule. It's an indication of how frequently the items in the rule appear together in the dataset.
- **Confidence (Y-axis):** This metric provides the conditional probability of encountering the consequent given the antecedent. In other words, it's the likelihood that the consequent is purchased when the antecedent is purchased.
- **Lift (Color Scale):** Lift measures how much more often the antecedent and consequent of the rule occur together than we would expect if they were statistically independent. A lift value greater than 1 indicates that the rule might be useful for predicting the consequent in future data sets. The different colors on the plot correspond to different ranges of lift values, which can provide a quick visual clue about the strength and usefulness of the rules.

In the plot, each point represents an association rule derived from a dataset, presumably from market basket analysis or a similar transactional dataset. The plot might be used to select the strongest rules for purposes such as cross-selling, recommendations, etc. Generally, rules with higher confidence and lift are considered stronger, but one also has to consider the support to ensure that the rule is applicable to a reasonable number of cases.

For example, a point on the far right with high confidence and high lift (in a darker color) would represent a rule that occurs relatively infrequently but has a high probability of being true when it does occur and a strong positive relationship between the antecedent and consequent. Conversely, a point on the upper left might have high confidence and lift but low support, suggesting that while the rule is a good predictor, it applies to only a small number of transactions.

5.2 Application on an external dataset

To validate the versatility and efficiency of our proposed method, we applied it to a market basket dataset, representative of a typical retail transaction record (as shown in the attached figure). This dataset's structure was conducive to associative rule mining, revealing prevalent purchasing patterns and item associations.

For instance, the association between various themed lunch bags highlighted cross-selling opportunities, such as the strong relationship between "LUNCH BAG PINK POLKADOT" and "LUNCH BAG CARS BLUE" with a confidence of 0.46, indicating a high likelihood that customers who purchase one are likely to purchase the other. Such findings exemplify the potential of our method in uncovering commercial insights, suggesting its applicability beyond book recommendations.

The following tables exhibit a snippet of the association rules derived from the market basket dataset, showcasing the efficiency and practicality of our method in a retail context.

Table 5: Association Rules

Antecedent	Consequent	Confidence
"LUNCH BAG PINK POLKADOT"	"LUNCH BAG CARS BLUE"	0.46
"LUNCH BAG BLACK SKULL."	"LUNCH BAG PINK POLKADOT"	0.42
"LUNCH BAG PINK POLKADOT"	"LUNCH BAG BLACK SKULL."	0.50
"LUNCH BAG BLACK SKULL."	"LUNCH BAG CARS BLUE"	0.41
"LUNCH BAG CARS BLUE"	"LUNCH BAG BLACK SKULL."	0.46
"LUNCH BAG BLACK SKULL."	"LUNCH BAG SUKI DESIGN"	0.41
"LUNCH BAG SUKI DESIGN"	"LUNCH BAG BLACK SKULL."	0.48
"ROSES REGENCY TEACUP AND SAUCER"	"GREEN REGENCY TEACUP AND SAUCER"	0.72
"GREEN REGENCY TEACUP AND SAUCER"	"ROSES REGENCY TEACUP AND SAUCER"	0.75
"PINK REGENCY TEACUP AND SAUCER"	"ROSES REGENCY TEACUP AND SAUCER"	0.77
"ROSES REGENCY TEACUP AND SAUCER"	"PINK REGENCY TEACUP AND SAUCER"	0.56
"PINK REGENCY TEACUP AND SAUCER"	"GREEN REGENCY TEACUP AND SAUCER"	0.82
"GREEN REGENCY TEACUP AND SAUCER"	"PINK REGENCY TEACUP AND SAUCER"	0.62
"CHARLOTTE BAG SUKI DESIGN"	"RED RETROSPOT CHARLOTTE BAG"	0.58
"RED RETROSPOT CHARLOTTE BAG"	"CHARLOTTE BAG SUKI DESIGN"	0.49

Table 6: Association Rules (Continued)

Antecedent	Consequent	Confidence
"CHARLOTTE BAG PINK POLKA-DOT"	"RED RETROSPOT CHARLOTTE BAG"	0.71
"RED RETROSPOT CHARLOTTE BAG"	"CHARLOTTE BAG PINK POLKA-DOT"	0.51
"WOODLAND CHARLOTTE BAG"	"RED RETROSPOT CHARLOTTE BAG"	0.61
"RED RETROSPOT CHARLOTTE BAG"	"WOODLAND CHARLOTTE BAG"	0.49
"GARDENERS KNEELING PAD KEEP CALM"	"GARDENERS KNEELING PAD CUP OF TEA"	0.60
"GARDENERS KNEELING PAD CUP OF TEA"	"GARDENERS KNEELING PAD KEEP CALM"	0.72
"PINK REGENCY TEACUP AND SAUCER"	"GREEN REGENCY TEACUP AND SAUCER", "ROSES REGENCY TEACUP AND SAUCER"	0.70
"GREEN REGENCY TEACUP AND SAUCER"	"PINK REGENCY TEACUP AND SAUCER", "ROSES REGENCY TEACUP AND SAUCER"	0.53
"ROSES REGENCY TEACUP AND SAUCER"	"PINK REGENCY TEACUP AND SAUCER", "GREEN REGENCY TEACUP AND SAUCER"	0.51
"PINK REGENCY TEACUP AND SAUCER", "GREEN REGENCY TEACUP AND SAUCER"	"ROSES REGENCY TEACUP AND SAUCER"	0.47
"PINK REGENCY TEACUP AND SAUCER", "ROSES REGENCY TEACUP AND SAUCER"	"GREEN REGENCY TEACUP AND SAUCER"	0.44
"GREEN REGENCY TEACUP AND SAUCER", "ROSES REGENCY TEACUP AND SAUCER"	"PINK REGENCY TEACUP AND SAUCER"	0.41

6 Comparison of all the three algorithms

The performance of the Apriori, Eclat, and FP-Growth algorithms was empirically evaluated through their execution time across various transaction sizes. The graph below illustrates the time complexity comparison among these algorithms.

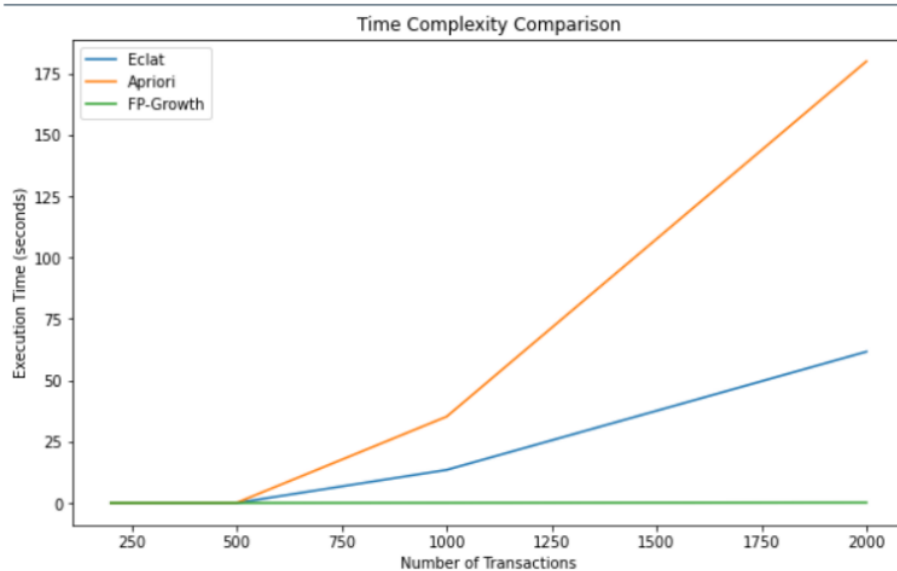


Figure 2: Time Complexity Comparison of Eclat, Apriori, and FP-Growth Algorithms

Apriori is one of the earliest and simplest algorithms for association rule mining. Its main drawback is its inefficiency in handling large datasets due to its need to repeatedly scan the database to discover frequent itemsets. This frequent scanning of the database makes Apriori the slowest among the three algorithms.

Eclat, on the other hand, improves upon the Apriori algorithm by using a vertical layout and tidset intersections instead of candidate generation. This optimization reduces the need for multiple database scans, leading to better performance compared to Apriori.

FP-Growth is the fastest among the three algorithms. It constructs a compact data structure called the FP-tree, which encodes the dataset and facilitates efficient mining of frequent itemsets without the need for candidate generation. By eliminating the need for multiple database scans and candidate generation, FP-Growth significantly outperforms both Apriori and Eclat, especially on large datasets.

In summary, while Apriori is the slowest due to its repeated database scans, Eclat performs better by avoiding candidate generation, and FP-Growth achieves the highest efficiency by constructing the FP-tree data structure.

7 Conclusion

This project embarked on the enhancement of book recommendation systems by leveraging the prowess of associative data mining techniques applied to a substantial Amazon book dataset. Our exploration centered around three prominent algorithms: Apriori, Eclat, and FP-Growth. The Apriori algorithm laid the groundwork by establishing a benchmark with its interpretability, despite its computational intensity. Eclat offered a middle ground with improved speed due to its depth-first search approach and vertical data layout. The FP-Growth algorithm emerged as the frontrunner, demonstrating superior efficiency and scalability, especially apt for dense datasets.

By setting judicious thresholds for support and confidence, our method effectively filtered and unearthed association rules that signified not just frequent but also meaningful correlations between items. The outcomes indicated that the FP-Growth algorithm is well-suited for rapid pattern discovery in vast datasets without the exhaustive generation of candidate sets. The comparative analysis substantiated the method's effectiveness against current techniques, emphasizing its computational advantages.

Future Scope: Advancements in our approach can take several directions. Integrating machine learning models to predict the thresholds for support and confidence could further refine the rule extraction process. Delving into parallel and distributed computing paradigms could amplify the method's capability to process even larger datasets efficiently. Moreover, applying the methodology to a wider array of domains, such as e-commerce or streaming services, could validate its versatility and open avenues for cross-domain recommendations. Finally, the incorporation of temporal dynamics to track changes in consumer behavior over time could yield a more responsive and dynamic recommendation system.

In essence, this project has laid a foundation that melds traditional data mining techniques with modern computational strategies, offering a robust framework for developing advanced recommendation systems.

References

- [1] A. M. Maghari, I. A. Al-Najjar, S. J. Al-Iqatqah, and S. A. Abu-Naser, "Books' Rating Prediction Using Just Neural Network," *International Journal of Engineering and Information Systems (IJEAIS)*, vol. 4, no. 10, pp. 17-22, Oct. 2020.
- [2] S. Bagchi, "Books Dataset," *Kaggle*, [Online]. Available: <https://www.kaggle.com/datasets/saurabhbagchi/books-dataset>. [Accessed: day, month, year].
- [3] A. Ahmedov, "Market Basket Analysis," *Kaggle*, [Online]. Available: <https://www.kaggle.com/datasets/aslanahmedov/market-basket-analysis>. [Accessed: day, month, year].

Book Recommendation System on Amazon Dataset

ORIGINALITY REPORT

23%

SIMILARITY INDEX

21%

INTERNET SOURCES

14%

PUBLICATIONS

%

STUDENT PAPERS

PRIMARY SOURCES

1	Boyu Shen. "E-commerce Customer Segmentation via Unsupervised Machine Learning", The 2nd International Conference on Computing and Data Science, 2021 Publication	2%
2	github.com Internet Source	1%
3	www.coursehero.com Internet Source	1%
4	www.bookhaven.co.nz Internet Source	1%
5	www.amazon.com Internet Source	1%
6	www.wkonline.com Internet Source	1%
7	bookcrossing.com Internet Source	1%
8	books.cheap-internet-store.com Internet Source	1%

9	www.members.shaw.ca Internet Source	1 %
10	www.2save.co.uk Internet Source	1 %
11	www.shoptvweb.com Internet Source	1 %
12	www.sitescraper.co.uk Internet Source	1 %
13	bn265561.amazon.ebisusama.org Internet Source	1 %
14	jurnal.untan.ac.id Internet Source	1 %
15	ise.ait.ac.th Internet Source	1 %
16	www.health-books-web.com Internet Source	1 %
17	www.wbthub.com Internet Source	1 %
18	thetreasuretrove.in Internet Source	1 %
19	www.acadlore.com Internet Source	<1 %
20	Imen Ben El Kouni, Wafa Karoui, Lotfi Ben Romdhane. "PRUCARS: improved association	<1 %

rule-based social recommender systems
using overlapping community detection",
Procedia Computer Science, 2020

Publication

21

www.ijritcc.org

Internet Source

<1 %

22

"SESSION 1: Image Processing", Journal of
Digital Imaging, 2003

Publication

<1 %

23

aimotion.blogspot.com

Internet Source

<1 %

24

hasonss.com

Internet Source

<1 %

25

Mishra, Shruti, Debahuti Mishra, and Sandeep
Kumar Satapathy. "Fuzzy pattern tree
approach for mining frequent patterns from
gene expression data", 2011 3rd International
Conference on Electronics Computer
Technology, 2011.

Publication

<1 %

26

www.mystery-thriller.com

Internet Source

<1 %

27

flakepotato.com

Internet Source

<1 %

28

www.ijitr.com

Internet Source

<1 %

29

www.unverse.com

Internet Source

<1 %

30

Teik Toe Teoh. "Chapter 10 AI in Sales", Springer Science and Business Media LLC, 2023

Publication

<1 %

31

Vazeerudeen Abdul Hameed, Muhammad Ehsan Rana, Lim Hui Enn. "Apriori Algorithm based Association Rule Mining to Enhance Small-Scale Retailer Sales", 2023 IEEE 6th International Conference on Big Data and Artificial Intelligence (BDAI), 2023

Publication

<1 %

32

drtattoo.com

Internet Source

<1 %

33

en.wikipedia.org

Internet Source

<1 %

34

www.mdpi.com

Internet Source

<1 %

35

Sarath Mechery, N. Preethi. "Chapter 44 An Intelligent Recommendation System Using Market Segmentation", Springer Science and Business Media LLC, 2022

Publication

<1 %

36

Seema Desai, Satish R. Devane, Vimla Jethani. "Association Rule Mining Using Graph and

<1 %

Clustering Technique", 2012 Fourth
International Conference on Computational
Intelligence and Communication Networks,
2012

Publication

37

discuss.boardinfinity.com

Internet Source

<1 %

38

rstudio-pubs-static.s3.amazonaws.com

Internet Source

<1 %

39

Pisut Koomsap, Chih-Fan Tan, Yu-Ju Lin, Chin-Yin Huang. "Chapter 20 Managing a Retail Store and the Associated Warehouse with a Knowledge-Driven Approach", Springer Science and Business Media LLC, 2023

Publication

<1 %

40

Tsai-Pin Chu, Fan Wu, Shih-Wen Chiang. "Mining frequent pattern using item-transformation method", Fourth Annual ACIS International Conference on Computer and Information Science (ICIS'05), 2005

Publication

<1 %

41

Mahajan, Renuka, J.S. Sodhi, Vishal Mahajan, and Richa Misra. "Comparative study of mining algorithms for adaptive e-learning environment", International Journal of Logistics Economics and Globalisation, 2014.

Publication

<1 %

Exclude quotes Off

Exclude matches Off

Exclude bibliography On