

Association Rule Mining using Apriori Algorithm on Grocery Shopping Data

KHADKA PILOT^{1*} AND POUDEL KSHITIZ^{2*}

¹Institute of Engineering, Thapathali Campus, Thapathali, Nepal (e-mail: pilot.076bct025@tcioe.edu.np)

²Institute of Engineering, Thapathali Campus, Thapathali, Nepal (e-mail: kshitizpoudel18@gmail.com)

Corresponding author: Khadka Pilot (e-mail: pilot.076bct025@tcioe.edu.np).

* Authors contributed equally

ABSTRACT In the modern era of data-driven decision-making, business enterprises amass vast volumes of data from their daily operations. One prime example is the extensive troves of customer purchase data collected at grocery store checkout counters on a daily basis. This data deluge offers opportunities to extract meaningful insights. One such transformative technique is Association Rule Mining, a cornerstone of data mining, which uncovers intricate relationships and patterns hidden within the data. Association rule mining thrives in deciphering the underlying connections among items within large datasets. It allows understanding the co-occurrences and dependencies between products, behaviors, or events. From market basket analysis to recommendation engines, this technique fuels applications across various sectors, allowing efficient inventory management, personalized recommendations, and optimized business strategies.

INDEX TERMS Association rule mining, Apriori

I. INTRODUCTION

IN the realm of data mining and machine learning, association rule mining stands as a fundamental technique that unveils hidden relationships within large datasets. The essence of this method lies in its ability to uncover associations, correlations, and patterns in data that might otherwise remain obscured. These extracted associations, often presented in the form of association rules, shed light on the co-occurrence and dependencies among various items in a dataset. Association rule mining finds its significance in a multitude of real-world applications, ranging from market basket analysis in retail to recommendation systems, healthcare analytics, and beyond. By identifying meaningful connections between items, it enables businesses and researchers to make informed decisions, develop personalized recommendations, optimize supply chain management, and enhance customer experiences.

At its core, the process involves sifting through transactional or relational datasets to discover associations among items that tend to occur together. The most prominent algorithm for association rule mining is the Apriori algorithm, which efficiently generates frequent itemsets and derives association rules. The output of association rule mining often comprises a list of rules, each consisting of an antecedent (the condition) and a consequent (the outcome). These rules are usually ranked based on their support and confidence values,

allowing analysts to focus on the most relevant and trustworthy associations. As datasets continue to grow in size and complexity, association rule mining techniques have evolved to handle challenges like scalability and noise in the data. Advanced methods, including FP-Growth and Eclat, offer efficient alternatives to Apriori for discovering associations. Moreover, the integration of association rule mining with machine learning and deep learning techniques has opened up new avenues for extracting intricate patterns from diverse data sources.

II. METHODOLOGY

A. THEORY

In the context of market basket analysis, we define $I = \{i_1, i_2, \dots, i_d\}$ to represent the set encompassing all available items. Correspondingly, $T = \{t_1, t_2, \dots, t_N\}$ denotes the collection of transactions. Each individual transaction t_i consists of a subset of items selected from the item set I .

An 'itemset' refers to a collection of items, which can be composed of any number of elements, even none. Specifically, an itemset containing k distinct items is termed a 'k-itemset'. A transaction t_j is considered to encompass an itemset X if X is a subset of t_j .

Support count is the frequency of occurrence of an itemset. Mathematically, the 'support count' $\sigma(X)$ for an itemset X is

defined as:

$$\sigma(X) = |\{t_i \mid X \subseteq t_i, t_i \in T\}| \quad (1)$$

An association rule relates the two itemset and takes the form of $X \rightarrow Y$, where X and Y are sets of items that have no common elements ($X \cap Y = \emptyset$). The association rule is measured by its support and confidence. Support signifies how often a rule is applicable within a dataset, while confidence reveals how often items in set Y appear in transactions containing set X . This is so because it helps find out whether a rule occurs by random chance. A rule with extremely low support might arise coincidentally. Additionally, from a business standpoint, a rule with low support could lack practical value. For instance, promoting items rarely purchased together might not yield profitable outcomes.

The formula for support and confidence are:

$$s(X \rightarrow Y) = \frac{\sigma(X \cup Y)}{N} \quad (2)$$

$$c(X \rightarrow Y) = \frac{\sigma(X \cup Y)}{\sigma(X)} \quad (3)$$

Association rule mining involves two steps:

- 1) **Frequent Itemset Generation:** I involves identifying all itemsets that meet the minimum support requirement. Those itemset are called 'frequent itemset'.
- 2) **Rule Generation:** Extract high confidence rules from the generated itemset. These extracted rules are called strong rules.

B. FREQUENT ITEMSET GENERATION

In the brute-force method, the process involves iterating through the itemset list to calculate the support and confidence values for each conceivable rule. Nevertheless, this approach becomes significantly computationally intensive. For a dataset with d distinct items:

- Total count of itemsets 2^d
- Total count of possible association rules: $R = 3^d \times 2^{d+1} + 1$

To address the challenges posed by the brute-force approach, the Apriori method was developed. This method operates on a fundamental principle: if an itemset is frequent, then all its subsets must also be frequent.

In the context of the Apriori method,

let C_k denote the collection of candidate k -itemsets, and F_k denote the collection of frequent k -itemsets.

The algorithm calculates the support of each individual item once. After this step, the set of frequent 1-itemsets, F_1 , is generated. The algorithm then enters an iterative process wherein it generates new candidate k -itemsets by utilizing the frequent $(k - 1)$ -itemsets discovered in the previous iteration. Next, the algorithm counts the support of candidates to determine candidate itemset in C_k are contained in each transaction. It then removes candidate itemset with support less than the minimum support. This is followed iteratively until there are no frequent itemset being generated.

C. RULE GENERATION

An association rule can be extracted by partitioning the itemset Y into two non-empty subsets, X and $Y - X$, such that $X \rightarrow Y - X$ satisfies the confidence threshold.

Initially, all the high-confidence rules that have only one item in the rule consequent are extracted. These rules are then used to generate new candidate rules.

D. FACTORS INFLUENCING THE COMPLEXITY OF THE APRIORI ALGORITHM

When working with the Apriori algorithm, several factors contribute to its overall complexity and performance:

- 1) **Minimum Support Threshold:** Adjusting the support threshold impacts the algorithm's outcomes. Lowering the threshold leads to discovering more frequent itemsets, yet it might increase the count of candidate itemsets and extend the length of frequent itemsets.
- 2) **Dimensionality of the Data Set:** The dimensionality, or the total number of distinct items in the dataset, plays a role in complexity. With more items, additional storage is needed to maintain support counts. Moreover, if frequent item counts rise, computational and I/O costs could escalate.
- 3) **Database Size:** As the Apriori algorithm makes multiple passes over the data, a larger number of transactions in the database may elongate its runtime.
- 4) **Average Transaction Width:** Datasets with denser patterns exhibit increased average transaction width. Consequently, this broader width could lead to longer frequent itemsets and greater hash tree traversals, due to the rise in subsets within each transaction.

E. ALGORITHM

Algorithm 1 Frequent itemset generation using Apriori algorithm

```

1:  $k \leftarrow 1$ 
2:  $F_k \leftarrow \{i \mid i \in I \text{ and } \sigma(\{i\}) \geq N \times \text{minsup}\}$ 
3: repeat
4:    $k \leftarrow k + 1$ 
5:    $C_k \leftarrow \text{apriori-gen}(F_{k-1})$ 
6:   for all  $t \in T$  do
7:      $C_t \leftarrow \text{subset}(C_k, t)$ 
8:     for all  $c \in C_t$  do
9:        $\sigma(c) \leftarrow \sigma(c) + 1$ 
10:    end for
11:  end for
12:   $F_k \leftarrow \{c \mid c \in C_k \text{ and } \sigma(c) \geq N \times \text{minsup}\}$ 
13: until  $F_k = \emptyset$ 
14: Result  $\leftarrow \bigcup F_k$ 

```

F. SYSTEM BLOCK DIAGRAM

Appendix

Algorithm 2 Rule generation

```

1: for each frequent  $k$ -itemset  $F_k, k \geq 2$  do
2:    $H_1 \leftarrow \{i \mid i \in F_k\}$   $\triangleright$  1-item consequents of the rule
3:   ap-genrules( $F_k, H_1$ )  $\triangleright$  Call ap-genrules function
4: end for

```

Algorithm 3 Procedure ap-genrules(F_k, H_m)

```

1:  $k = |F_k|$   $\triangleright$  Size of frequent itemset.
2:  $m = |H_m|$   $\triangleright$  Size of rule consequent.
3: if  $k > m + 1$  then
4:    $H_{m+1} = \text{apriori-gen}(H_m)$ 
5:   for each  $h_{m+1} \in H_{m+1}$  do
6:      $\text{conf} = \frac{\sigma(F_k)}{\sigma(F_k - h_{m+1})}$ 
7:     if  $\text{conf} \geq \text{minconf}$  then
8:       output the rule  $(F_k - h_{m+1}) \rightarrow h_{m+1}$ 
9:     else
10:      delete  $h_{m+1}$  from  $H_{m+1}$ 
11:    end if
12:  end for
13:  ap-genrules( $F_k, H_{m+1}$ )
14: end if

```

G. INSTRUMENTATION

The following libraries were instrumental in completion of this project :

- **pd.read_csv**: This function is employed to read the CSV groceries dataset.
- **plt.violinplot**: Utilized for generating a violin plot that showcases the distribution of transaction sizes across various itemset counts.
- **plt.boxplot**: This function aids in identifying outliers within the grocery dataset by visualizing their distribution.
- **plt.hist**: It serves the purpose of plotting a histogram to depict the distribution of transaction sizes.
- **plt.barh**: Employed to create horizontal bar plots, useful for visualizing the ten most and least frequent itemsets.
- **mlxtend.preprocessing.TransactionEncoder**: This utility is used to encode the transactional data from a Python list of lists into a NumPy array.
- **mlxtend.frequent_patterns.apriori**: Applied to execute the Apriori algorithm on the transaction database, extracting frequent itemsets.
- **mlxtend.frequent_patterns.association_rules**: Utilized to generate association rules from the frequent itemsets obtained.

III. WORKING PRINCIPLE**A. DATA LOADING AND EXPLORATION**

The grocery dataset was loaded using the `pd.read_csv` function. This dataset comprises 9,835 customer transactions involving grocery shopping. It encompasses a total of 169 distinct items, reflecting products commonly found in a grocery store. Among the items included are butter, milk, yogurt,

cream cheese, spread cheese, rolls/buns, bottled water, soda, newspapers, and more.

B. ITEMSET ANALYSIS

Itemset analysis is a pivotal step in data mining, especially when applying algorithms like Apriori to transactional datasets such as grocery shopping records. This analysis focuses on understanding how often items appear together in transactions, revealing patterns in item combinations.

Visualization plays a critical role in itemset analysis as it helps to intuitively grasp the distribution and patterns of these item combinations. Common visualization techniques include:

- **Violin Plots**: These plots provide insights into transaction size distribution for varying itemset counts. They help identify variations in transaction sizes and their relationship with the number of items purchased.
- **Box Plots**: By illustrating transaction size distribution, box plots highlight potential outliers. These outliers could indicate unusually large or small transactions.
- **Histograms**: Histograms depict transaction size distribution, shedding light on the frequency of specific transaction sizes. This offers insights into transaction behaviors.

In our dataset, we observed that up the top six-item itemsets are more frequent, as revealed by the violin plot. Similarly, by visualizing transaction distribution, we could identify the most frequently occurring items. We found our 10 most frequent itemset to be: We identified the ten most frequent itemsets, which are as follows:

1) Most frequent itemset

- Whole milk
- Other vegetables
- Rolls/buns
- Soda
- Yogurt

These items are the ones most commonly purchased together in transactions. This finding aligns with common grocery shopping behavior where staples like milk, vegetables, rolls/buns, soda, and yogurt are often bought together. Customers might frequently purchase these items in a single shopping trip due to their everyday consumption or complementary nature.

2) Least Frequent Itemsets

The least frequent items in the dataset include:

- Sound storage medium
- Baby food
- Preservation products

These items have relatively low occurrence frequencies in transactions. Possible reasons for their infrequent appearance could be niche or specialized products. Customers might not need these items as frequently as others, leading to fewer occurrences in transactions.

Top One-Item Products

The top five one-item products, which are the most frequently purchased individual items, are:

- Canned beer
- Soda
- Whole milk
- Bottled beet
- Rolls/buns

These items suggest that beverages like beer, soda, and milk, along with staple food items like rolls/buns, have a high individual purchase rate. These items are commonly consumed and likely to be part of customers' regular grocery shopping lists.

The results suggest that the dataset's shopping patterns are characterized by the frequent purchase of staple items like milk, vegetables, rolls/buns, and beverages such as soda. Niche or specialized items like sound storage media, baby food, and preservation products appear less frequently, likely due to their specific use cases. These insights can help retailers optimize their product placements, marketing strategies, and inventory management to cater to customer preferences and boost sales.

TRANSACTION ENCODING

Transaction encoding is a fundamental preprocessing step in data mining, particularly when working with algorithms like the Apriori algorithm. It involves transforming transactional data into a suitable format that algorithms can process effectively. The primary objective is to convert the transactional information into a binary representation that reflects the presence or absence of specific items in each transaction.

In transactional datasets, each transaction represents a customer's purchase or interaction, and it contains a list of items bought together. For instance, in a grocery store dataset, each transaction could list the items a customer purchased in a single shopping instance.

Need for Encoding

Algorithms like Apriori work with structured data formats. However, transactional data is inherently unstructured, with varying items in each transaction. Transaction encoding addresses this disparity by transforming the dataset into a consistent matrix-like structure.

Binary Representation

Transaction encoding involves converting the transactional data into a binary matrix, where rows correspond to transactions, and columns correspond to unique items. Each cell in the matrix contains a binary value, indicating whether a specific item is present (1) or absent (0) in the corresponding transaction.

SUPPORT CALCULATION STEP

The support of an itemset in a dataset is a measure of how frequently that itemset appears in the transactions of the

dataset. It is a fundamental metric in association rule mining, as it helps identify which itemsets are frequent and potentially interesting for generating association rules.

Here's how the support calculation step works:

- 1) **Transaction Count:** Count the total number of transactions in the dataset. This value is denoted as N .
- 2) **Itemset Count:** Count the number of transactions that contain the specific itemset. This count is denoted as n , where n represents the number of transactions that include the given itemset.
- 3) **Support Calculation:** The support of the itemset is calculated as the ratio of the number of transactions containing the itemset to the total number of transactions in the dataset:

$$\text{Support (s)} = \frac{n}{N}$$

The support value lies between 0 and 1. A higher support value indicates that the itemset is more frequently occurring in the transactions.

- 4) **Threshold Comparison:** Compare the calculated support value with a predefined minimum support threshold. This threshold determines the minimum level of support that an itemset must have to be considered frequent. If the calculated support is greater than or equal to the threshold, the itemset is considered frequent; otherwise, it's considered infrequent.

ASSOCIATION RULE GENERATION STEP

The association rule generation step is a crucial part of the Apriori algorithm, which aims to extract meaningful rules from frequent itemsets identified in the dataset. An association rule is an implication of the form "if X , then Y ," where X and Y are itemsets, often referred to as the antecedent and consequent, respectively.

Here's how the association rule generation step works:

- 1) **Frequent Itemsets:** First, a set of frequent itemsets using the minimum support threshold are identified. These frequent itemsets act as the basis for generating association rules.
- 2) **Rule Evaluation Metrics:** An evaluation metrics is needed to determine the strength and significance of potential association rules. Common metrics include support, confidence, and lift. Confidence is particularly important. It measures the proportion of transactions containing the antecedent that also contain the consequent. Mathematically, it's defined as:

$$\text{Confidence (c)} = \frac{\text{Support (X} \cup \text{Y)}}{\text{Support (X)}}$$

The higher the confidence, the more likely the consequent will occur when the antecedent is present.

- 3) **Rule Generation:** For each frequent itemset with at least two items, association rules are generated by considering all possible combinations of the itemset as antecedent and the remaining items as consequent. If

the confidence of a rule exceeds a predefined minimum confidence threshold, it is considered a valid association rule.

- 4) **Pruning:** Not all generated rules are necessarily interesting or useful. To refine the rule set, further filtering based on additional metrics, such as lift or conviction can be applied to ensure that the rules discovered are meaningful.

RULE ANALYSIS FOR GROCERY DATASET

After applying the Apriori algorithm to a grocery store dataset and generating association rules, let's consider a few rules for analysis:

RULE 1

"If {Milk, Bread}, then {Eggs}"

- Support: 0.1 (10%)
- Confidence: 0.6 (60%)
- Lift: 1.5

Interpretation: This rule suggests that when customers buy both Milk and Bread, there's a 60% chance they will also buy Eggs. The lift of 1.5 indicates a positive association, meaning the occurrence of Milk, Bread, and Eggs together is 1.5 times more likely than if they were independent. The positive lift suggests a potential cross-promotion strategy for Milk, Bread, and Eggs.

IV. RESULT ANALYSIS

The preceding two plots offer insightful observations:

Transaction Distribution: The presented plots in Figure(3,4) highlight that a majority of transactions consist of a single item. This trend suggests that customers often prefer purchasing individual items rather than buying in bulk.

Most Bought Itemsets: The histogram in figure(5) provides a visual representation of the most frequently bought itemsets, along with their corresponding support counts. This visualization enables a quick identification of the pivotal items in the grocery store.

Least Frequent Itemsets: The plotted data of the least frequent itemsets in figure(6) implies that these particular sets may not hold significant importance. Consequently, these infrequent itemsets could potentially be considered for removal from the grocery store offerings.

Dominant Single Item: Despite whole milk being the most frequently bought item overall as seen in figure(5), an intriguing insight emerges when considering single-item purchases as seen in figure(7). Specifically, the item "canned beer" emerges with the highest support count for such cases.

Variation with Minimum Support: The graph in figure(8) illustrates how the number of generated rules changes with varying minimum support levels (10%, 8%, 6%, 4%, 1%, and 0.5%). This visual aids in estimating the appropriate minimum support count for generating frequent itemsets, as it depends on the nature and size of the database. Notably, as the minimum support threshold increases, the number of

generated rules experiences rapid growth. A higher support threshold ensures that only rules involving frequent items are generated. Conversely, setting a lower minimum support threshold results in numerous rules with high confidence. While such a scenario captures a multitude of associations, some of these rules may have lower frequency and lesser relevance for decision-making.

The insights drawn from these analyses contribute to a better understanding of customer behaviors, preferences, and shopping patterns, enabling informed business strategies and enhancing the overall shopping experience.

Support	Itemsets	Length
0.256	(whole milk)	1
0.193	(other vegetables)	1
0.184	(rolls/buns)	1
0.174	(soda)	1

TABLE 1. Frequent Itemsets generated with minimum support = 0.15

Support	Itemsets	Length
0.256	(whole milk)	1
0.193	(other vegetables)	1
0.184	(rolls/buns)	1
0.174	(soda)	1
0.140	(yogurt)	1
0.111	(bottled water)	1
0.109	(root vegetables)	1
0.105	(tropical fruit)	1
0.099	(shopping bags)	1
0.094	(sausage)	1
0.089	(pastry)	1
0.083	(citrus fruit)	1
0.081	(bottled beer)	1
0.080	(newspapers)	1
0.078	(canned beer)	1
0.076	(pip fruit)	1
0.075	(whole milk, other vegetables)	2
0.072	(fruit/vegetable juice)	1
0.072	(whipped/sour cream)	1
0.065	(brown bread)	1
0.063	(domestic eggs)	1
0.059	(frankfurter)	1
0.059	(margarine)	1
0.058	(coffee)	1
0.058	(pork)	1
0.057	(whole milk, rolls/buns)	2
0.056	(yogurt, whole milk)	2
0.055	(butter)	1
0.053	(curd)	1
0.052	(beef)	1
0.052	(napkins)	1

TABLE 2. Frequent Itemsets generated with minimum support = 0.05

Two distinct figures (9,10) showcase the generated association rules resulting from different frequent itemsets, which emerge due to varying choices of the minimum support threshold. Once these rules are generated, they can undergo further filtration based on a confidence threshold.

Key Concepts in Association Rule Analysis:

- 1) **Antecedents:** Refers to items or itemsets on the left-hand side (LHS) of the rule.

- 2) Consequents: Refers to items or itemsets on the right-hand side (RHS) of the rule.
- 3) Antecedent Support: Indicates how often the antecedents appear individually in the transactional database.
- 4) Consequent Support: Indicates how often the consequents appear individually in the transactional database.
- 5) Support: Denotes the frequency of appearance of both antecedents and consequents together.
- 6) Confidence: Measures how frequently the rule has proven true. Calculated as the support of the entire rule divided by the support of the antecedent.
- 7) Lift: Indicates how much more likely antecedent and consequent items are to appear together compared to if they were independent. $Lift > 1$ implies a positive association.
- 8) Leverage: Represents the difference between observed and expected co-occurrence of antecedents and consequents.
- 9) Conviction: Measures the strength of the rule's prediction reliance on dependency.
- 10) Zhang's Metric: A composite measure encompassing support, confidence, and the rule's structure.

Understanding these concepts aids in comprehending the significance of generated association rules. Further analysis and selection of rules based on these metrics empower businesses to uncover meaningful patterns and make informed decisions to optimize their strategies.

Visualization of Association Rules

The visualizations offer valuable insights into the generated association rules:

- 1) Lift-based Bar Chart: The horizontal bar chart in figure(11), portrays the top N association rules, arranged by their lift values. Higher lift values indicate more robust associations. Rules with elevated lift values imply that customers who purchase the antecedent items are more likely to also purchase the consequent items.
- 2) Network Graph: Figure(12) presents a network graph illustrating the association rules. Each node within the graph represents an item, and the edges between nodes signify associations between those items, as derived from the generated association rules. The thickness of the edges (represented by lift values) serves as an indicator of the strength of the association between the items.

V. DISCUSSION

A. BRUTE FORCE APPROACH VS APRIORI ALGORITHM

While both the brute force approach and the Apriori algorithm yield the same frequent itemsets, the Apriori algorithm boasts greater speed. The brute force approach involves considering all possible itemsets and computing their support and confidence values. On the other hand, the Apriori algorithm

capitalizes on the Apriori property to decrease the count of itemsets under consideration, rendering it faster. Let's assume the following specifications:

- N : Total number of transactions in the dataset.
- M : Total number of unique items in the dataset.
- min_sup : Minimum support threshold.
- W : Maximum size of itemsets.

For the brute force approach, the complexity involves generating 2^n possible itemsets and scanning the database for each itemset, resulting in a complexity of $O(N \cdot M \cdot W)$. For the Apriori approach, candidate itemsets are generated by self-joining frequent itemsets from the previous level. The time and space complexity is lower due to this efficient generation process.

C) RULE GENERATION

Once frequent itemsets are derived, the generation of association rules becomes essential. For a dataset containing d items, using the brute force rule generation approach, the total number of generated rules (R) is given by:

$$R = 3^d - 2^{d+1} + 1.$$

The goal of frequent itemset generation is to reduce this number significantly by considering only those itemsets meeting the minimum support count criterion. For each frequent itemset with 'k' items, only $2^k - 2$ rules need to be generated, avoiding empty antecedents or consequents. This approach is faster due to the significant reduction in the number of itemsets. Generating rules involves calculating the confidence value for each rule. This task is also streamlined by not considering redundant rules. Additionally, generating rules from frequent itemsets eliminates the need for extra scans on the database, as the support count for both antecedents and consequents has already been calculated during frequent itemset generation. This efficiency arises from the fact that antecedents and consequents are subsets of frequent itemsets, enabling the calculation of their support counts in advance.

VI. CONCLUSION

Association analysis using the Apriori algorithm provides valuable insights into the purchasing patterns and relationships among items in the grocery dataset. Through the Apriori algorithm, we were able to identify frequent itemsets, which are combinations of items that occur together in a significant number of transactions. These frequent itemsets unveil hidden patterns and associations that can be leveraged for decision-making and business strategies.

The Apriori algorithm's efficiency in pruning the search space through the Apriori property played a crucial role in speeding up the process of generating frequent itemsets and association rules. By considering only the relevant itemsets and rules, we were able to gain meaningful insights without exhaustive calculations.

Visualizations, such as the lift-based bar chart and network graph, added an extra layer of understanding, enabling us to

grasp the strength of associations and relationships among different items. These visualizations make it easier to interpret the results and identify items that are often purchased together.

In business contexts, the generated association rules can be used to enhance marketing strategies, optimize product placements, and even personalize recommendations for customers. By understanding the relationships between items, businesses can drive sales, improve customer satisfaction, and optimize inventory management. Association analysis using the Apriori algorithm is a powerful technique for uncovering valuable insights from transactional data, contributing to more informed and data-driven decision-making in various domains.

REFERENCES

- [1] Goodfellow, I., Bengio, Y., & Courville, A. (2016). *"Deep Learning."* MIT Press.
- [2] Daniel Jurafsky & James H. Martin. *Speech and Language Processing*.
- [3] M. Narasimha Murty, V. Susheela Devi. *Pattern Recognition: An Algorithmic Approach*.
- [4] Rakesh Agrawal, Tomasz Imieliński, Arun Swami, *Mining Association Rules between Sets of Items in Large Databases*, *ACM SIGMOD Record*, Vol. 22, Issue 2, June 1993, Publisher: Association for Computing Machinery.
- [5] Rakesh Agrawal, Ramakrishnan Srikant, *Fast Algorithms for Mining Association Rules in Large Databases*, In *Proceedings of the 20th International Conference on Very Large Data Bases*, ISBN: 1558601538, Publisher: Morgan Kaufmann Publishers Inc., Address: San Francisco, CA, USA, Year: 1994.



PILOT KHADKA is a student at Institute of Engineering, Thapathali Campus. He is expected to graduate in Bachelor of Computer Engineering in 2024. During his time at Thapathali Campus, Pilot has actively engaged in various academic and extracurricular activities. He has participated in coding competitions, collaborated on software development projects, and demonstrated a keen interest in exploring new technologies.



KSHITIZ POUDEL is a student at Institute of Engineering, Thapathali Campus. He is expected to graduate in Bachelor of Computer Engineering in 2024. His relentless pursuit of knowledge and dedication to uplifting and empowering others make him an exceptional contributor and an invaluable asset to the academic community.

VII. APPENDIX

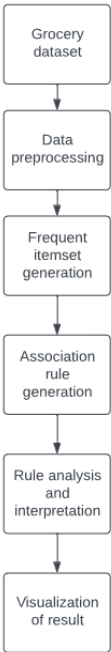


FIGURE 1. System Block Diagram

Time(s)	Time 1	Time 2	Time 3	Time 4	Time 5	Time 6	Time 7	Time 8	Time 9	Time 10	Time 11	Time 12	Time 13	Time 14	Time 15	Time 16	Time 17	Time 18	Time 19	Time 20	Time 21	Time 22	Time 23	Time 24	Time 25	Time 26	Time 27	Time 28	Time 29	Time 30	Time 31	Time 32
0	4	citrus fruit	beans flavored protein	margarine	ready soups	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN
1	3	tropical fruit	yogurt	coffee	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN
2	1	whole milk	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN
3	4	app fruit	yogurt	cream cheese	meat spreads	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN
4	4	other vegetables	whole milk	condensed milk	very fine bakery product	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN

FIGURE 2. Snapshot of dataset

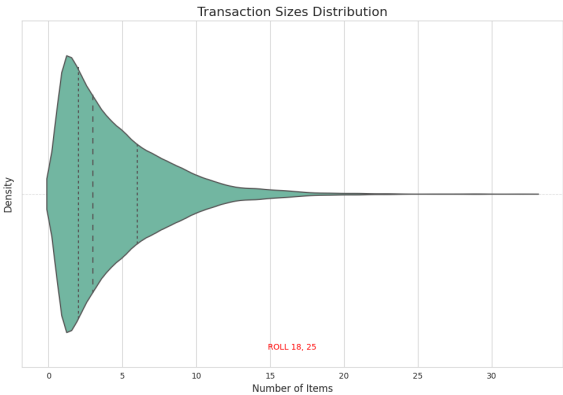


FIGURE 3. violin plot of transaction size distribution

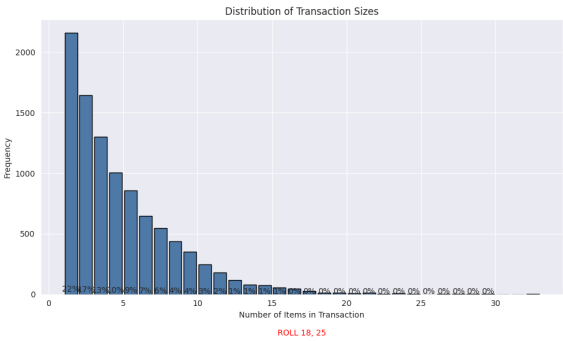


FIGURE 4. Frequency distribution of transaction size

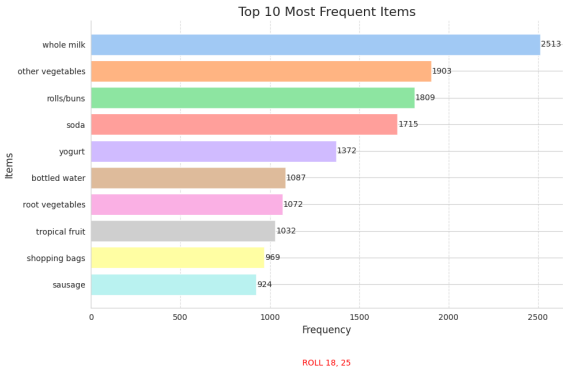


FIGURE 5. Top 10 most frequent items

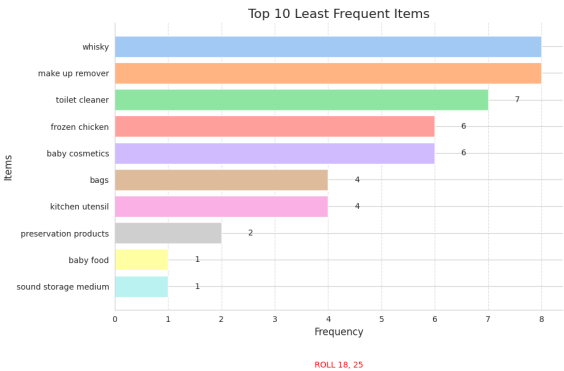


FIGURE 6. Top 10 least frequent items

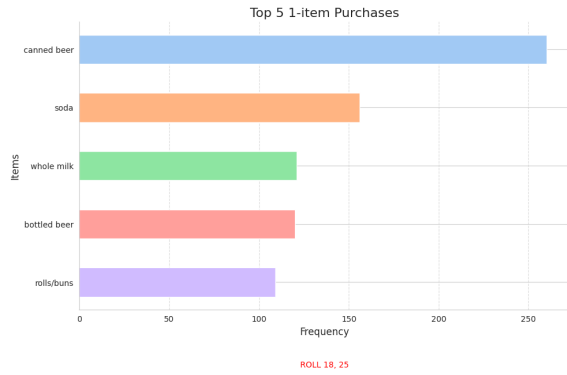


FIGURE 7. Top 5 1 item purchases

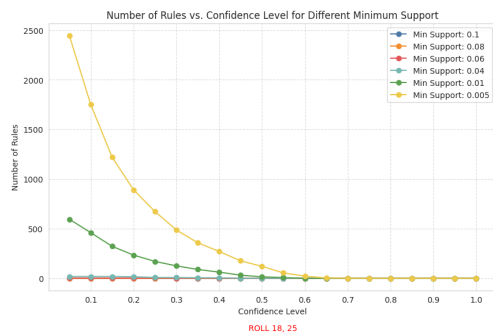


FIGURE 8. Number of rules vs confidence level

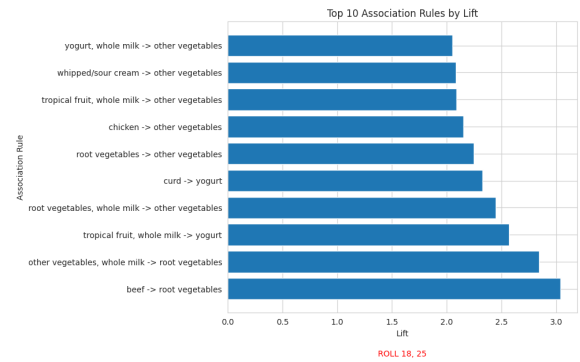


FIGURE 11. Top 10 association rules by lift

	antecedents	consequents	antecedent support	consequent support	support	confidence	lift	leverage	conviction	zhangs_metric
0	(yogurt)	(whole milk)	0.139502	0.255516	0.056024	0.401603	1.571735	0.020379	1.244132	0.422732
1	(other vegetables)	(whole milk)	0.193493	0.255516	0.074835	0.386758	1.513634	0.025394	1.214013	0.420750
2	(rolls/buns)	(whole milk)	0.183935	0.255516	0.056634	0.307905	1.205032	0.009636	1.075696	0.208496

FIGURE 9. Association rules for minimum support=0.05

	antecedents	consequents	antecedent support	consequent support	support	confidence	lift	leverage	conviction	zhangs_metric
0	(beef)	(root vegetables)	0.052466	0.108998	0.017387	0.331395	3.040367	0.011668	1.332828	0.708251
1	(other vegetables, whole milk)	(root vegetables)	0.074835	0.108998	0.023183	0.309783	2.842062	0.015026	1.290900	0.700572
2	(tropical fruit, whole milk)	(yogurt)	0.042288	0.139502	0.015150	0.358173	2.567516	0.009249	1.340701	0.637483
3	(root vegetables, whole milk)	(other vegetables)	0.048907	0.193493	0.023183	0.474012	2.449770	0.013719	1.533320	0.622230
4	(curd)	(yogurt)	0.053279	0.139502	0.017285	0.324427	2.325615	0.008853	1.273732	0.602085
5	(root vegetables)	(other vegetables)	0.108998	0.193493	0.047382	0.434701	2.246605	0.026291	1.426893	0.622764
6	(chicken)	(other vegetables)	0.042908	0.193493	0.017895	0.417062	2.155439	0.008593	1.383521	0.560090
7	(tropical fruit, whole milk)	(other vegetables)	0.042288	0.193493	0.017082	0.403846	2.087140	0.008898	1.352851	0.543880
8	(whipped/sour cream)	(other vegetables)	0.071683	0.193493	0.028876	0.402837	2.081924	0.015006	1.350565	0.558803
9	(yogurt, whole milk)	(other vegetables)	0.056024	0.193493	0.022267	0.387459	2.054131	0.011427	1.338511	0.543633
10	(tropical fruit, yogurt)	(whole milk)	0.028283	0.255516	0.015150	0.517361	2.024770	0.007698	1.542528	0.521384
11	(other vegetables, yogurt)	(whole milk)	0.043416	0.255516	0.022267	0.512881	2.007235	0.011174	1.528340	0.524577
12	(butter)	(whole milk)	0.055414	0.255516	0.027555	0.497248	1.946053	0.013395	1.480617	0.514659
13	(beef)	(other vegetables)	0.052466	0.193493	0.019725	0.375969	1.943066	0.009574	1.292416	0.512224
14	(pork)	(other vegetables)	0.057651	0.193493	0.021657	0.375661	1.941476	0.010502	1.291779	0.514595

FIGURE 10. Association rules for minimum support =0.02

Network Graph of Association Rules

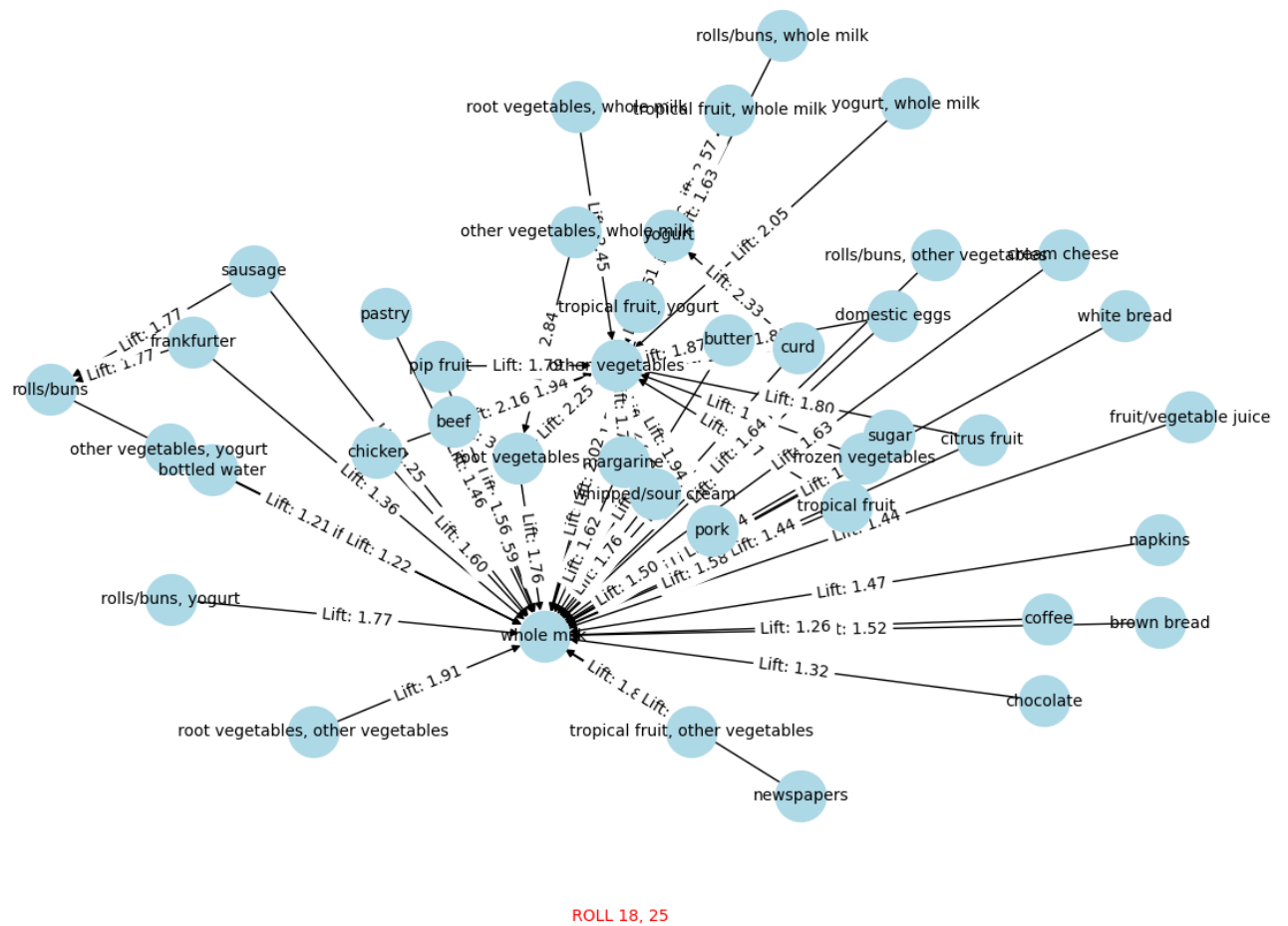


FIGURE 12. Network graph of association rules

VIII. CODE

```

1  # -*- coding: utf-8 -*-
2  """apriori.ipynb
3
4  Automatically generated by Colaboratory.
5
6  Original file is located at
7      https://colab.research.google.com/drive/1
8      a8msdNSGQpLGM3lYE-6H2hct8_psDmcV
9  """
10 import pandas as pd
11 import numpy as np
12 import math
13
14 from google.colab import drive
15 drive.mount('/content/drive')
16
17 path = '/content/drive/MyDrive/Colab_Notebooks/
18     groceries2.csv'
19
20 data_df = pd.read_csv(path)
21 data_df
22
23 new_df = data_df.head(50)
24 # Display the new DataFrame
25 new_df
26
27 new_df.drop('Item(s)', axis=1, inplace=True)
28 new_df.head()
29
30 """#Data visualization"""
31
32 import matplotlib.pyplot as plt
33 import seaborn as sns
34 num_items_per_transaction = data_df["Item(s)"]
35
36 # Set the style using Seaborn
37 sns.set_style("whitegrid")
38 plt.figure(figsize=(10, 7))
39
40 # Create a violin plot with custom colors
41 colors = ['#66c2a5', '#fc8d62', '#8da0cb', '#
42     e78ac3']
43 sns.violinplot(x=num_items_per_transaction,
44     palette=colors, inner="quartile")
45
46 # Customize the plot for better presentation
47 plt.title("Transaction Sizes Distribution",
48     fontsize=16)
49 plt.xlabel("Number of Items", fontsize=12)
50 plt.ylabel("Density", fontsize=12)
51 plt.xticks(fontsize=10)
52 plt.yticks(fontsize=10)
53 plt.grid(axis="y", linestyle="--", alpha=0.7)
54
55 # Add extra legend
56 extra_legend = "ROLL 18, 25"
57 plt.text(0.5, 0.05, extra_legend, ha='center',
58     color='red', transform=plt.gca().transAxes)
59
60 # Show the plot
61 plt.tight_layout()
62 plt.show()
63
64 import matplotlib.pyplot as plt
65 plt.figure(figsize=(8, 6))
66
67 # Use patch_artist=True to fill the box plot with
68     colors
69 plt.boxplot(num_items_per_transaction, vert=False,
70     patch_artist=True)

```

```

65
66 # Customize the plot for better presentation
67 plt.title("Distribution of Transaction Sizes")
68 plt.xlabel("Number of Items")
69 plt.grid(axis="x") # Show grid lines on the x-
70     axis for better readability
71 plt.yticks([]) # Remove the labels for the y-axis
72
73 # Add extra legend
74 extra_legend = "ROLL 18, 25"
75 plt.text(0.5, -0.15, extra_legend, ha='center',
76     color='red', transform=plt.gca().transAxes)
77
78 plt.tight_layout()
79
80 # Show the plot
81 plt.show()
82
83 import matplotlib.pyplot as plt
84 import seaborn as sns
85 num_items_per_transaction = data_df["Item(s)"]
86 min_size = min(num_items_per_transaction)
87 max_size = max(num_items_per_transaction)
88
89 # Set the style using Seaborn
90 sns.set_style("darkgrid")
91 plt.figure(figsize=(10, 6))
92
93 # Create bins with a size of 1
94 bins = list(range(min_size, max_size + 2))
95
96 # Create histograms for transaction sizes with a
97     custom color
98 colors = ['#4e79a7']
99 plt.hist(num_items_per_transaction, bins=bins,
100     edgecolor="black", color=colors, rwidth=0.85)
101 plt.title("Distribution of Transaction Sizes")
102 plt.xlabel("Number of Items in Transaction")
103 plt.ylabel("Frequency")
104
105 # Calculate the percentage of each number of items
106     in the "Item(s)" column
107 item_count = data_df["Item(s)"].value_counts()
108 total_transactions = len(data_df["Item(s)"])
109 percentage_item_purchases = (item_count /
110     total_transactions) * 100
111
112 # Add percentage labels above bars
113 for i, percentage in enumerate(
114     percentage_item_purchases):
115     plt.text(bins[i] + 0.5, percentage + 2, f"{
116         round(percentage)}%", ha='center', fontsize
117         =10)
118
119 # Add extra legend
120 extra_legend = "ROLL 18, 25"
121 plt.text(0.5, -0.15, extra_legend, ha='center',
122     color='red', transform=plt.gca().transAxes)
123
124 plt.tight_layout()
125
126 # Show the plot
127 plt.show()
128
129 import matplotlib.pyplot as plt
130 import seaborn as sns
131
132 # Set the style using Seaborn
133 sns.set_style("whitegrid")
134 plt.figure(figsize=(10, 6))
135
136 # Calculate the frequency of each item in the
137     DataFrame
138 item_counts = data_df.iloc[:, 1:].stack().

```

```

value_counts()
128
129 # Select the top 10 most frequent items
130 top_10_items = item_counts.head(10)
131
132 # Create a horizontal bar plot with a custom color
133 colors = sns.color_palette("pastel")
134 plt.barh(top_10_items.index, top_10_items.values,
135          color=colors)
136
137 # Customize the plot for better presentation
138 plt.gca().invert_yaxis() # Invert the y-axis
139 plt.title("Top 10 Most Frequent Items", fontsize
140          =16)
141 plt.xlabel("Frequency", fontsize=12)
142 plt.ylabel("Items", fontsize=12)
143 plt.xticks(fontsize=10)
144 plt.yticks(fontsize=10)
145 plt.grid(axis="x", linestyle="--", alpha=0.7)
146 plt.gca().spines["right"].set_visible(False) #
147     Remove right border
148 plt.gca().spines["top"].set_visible(False) #
149     Remove top border
150 plt.tight_layout()
151
152 # Add values inside the bars
153 for index, value in enumerate(top_10_items.values)
154 :
155     plt.text(value + 3, index, str(value), va="
156     center", fontsize=10)
157
158 # Add extra legend
159 extra_legend = "ROLL 18, 25"
160 plt.text(0.5, -0.2, extra_legend, ha='center',
161         color='red', transform=plt.gca().transAxes)
162
163 # Show the plot
164 plt.show()
165
166 # Set the style using Seaborn
167 sns.set_style("whitegrid")
168 plt.figure(figsize=(10, 6))
169
170 # Calculate the frequency of each item in the
171 DataFrame
172 item_counts = data_df.iloc[:, 1:].stack().
173     value_counts()
174
175 # Select the bottom 10 least frequent items
176 bottom_10_items = item_counts.tail(10)
177
178 # Create a horizontal bar plot with a custom color
179 colors = sns.color_palette("pastel")
180 bars = plt.barh(bottom_10_items.index,
181                bottom_10_items.values, color=colors)
182
183 # Customize the plot for better presentation
184 plt.gca().invert_yaxis() # Invert the y-axis
185 plt.title("Top 10 Least Frequent Items", fontsize
186          =16)
187 plt.xlabel("Frequency", fontsize=12)
188 plt.ylabel("Items", fontsize=12)
189 plt.xticks(fontsize=10)
190 plt.yticks(fontsize=10)
191 plt.grid(axis="x", linestyle="--", alpha=0.7)
192 plt.gca().spines["right"].set_visible(False) #
193     Remove right border
194 plt.gca().spines["top"].set_visible(False) #
195     Remove top border
196 plt.tight_layout()
197
198 # Add values inside the bars at the edge
199 for bar, value in zip(bars, bottom_10_items.values
200 ):

```

```

201     plt.text(bar.get_width() + 0.5, bar.get_y() +
202             bar.get_height() / 2, str(value), va="center",
203             fontsize=10)
204
205 # Add extra legend
206 extra_legend = "ROLL 18, 25"
207 plt.text(0.5, -0.2, extra_legend, ha='center',
208         color='red', transform=plt.gca().transAxes)
209
210 # Show the plot
211 plt.show()
212
213 import matplotlib.pyplot as plt
214 import seaborn as sns
215
216 # Set the style using Seaborn
217 sns.set_style("whitegrid")
218 plt.figure(figsize=(10, 6))
219
220 # Extract the item columns
221 item_columns = data_df.columns[1:33]
222
223 # Create a new DataFrame to store the standalone
224 purchases (transactions with only one item)
225 standalone_purchases = data_df[data_df["Item(s)"]
226 == 1][item_columns]
227
228 # Item Frequency Analysis for Standalone Purchases
229 standalone_item_counts = standalone_purchases.
230     stack().value_counts()
231 top_standalone_items = standalone_item_counts.head
232 (5) # Get the top 5 most frequent standalone
233 items
234
235 # Create a horizontal bar plot with a custom color
236 colors = sns.color_palette("pastel")
237 plt.barh(top_standalone_items.index,
238         top_standalone_items.values, color=colors,
239         height=0.5)
240
241 # Customize the plot for better presentation
242 plt.gca().invert_yaxis() # Invert the y-axis
243 plt.title("Top 5 1-item Purchases", fontsize=16)
244 plt.xlabel("Frequency", fontsize=12)
245 plt.ylabel("Items", fontsize=12)
246 plt.xticks(fontsize=10)
247 plt.yticks(fontsize=10)
248 plt.grid(axis="x", linestyle="--", alpha=0.7)
249 plt.gca().spines["right"].set_visible(False) #
250     Remove right border
251 plt.gca().spines["top"].set_visible(False) #
252     Remove top border
253 plt.tight_layout()
254
255 # Add extra legend
256 extra_legend = "ROLL 18, 25"
257 plt.text(0.5, -0.2, extra_legend, ha='center',
258         color='red', transform=plt.gca().transAxes)
259
260 # Show the plot
261 plt.show()
262
263 # Initialize a list to store transactions
264 transactions = []
265
266 # Iterate through each row in the DataFrame
267 for index, row in new_df.iterrows():
268     # Convert row to a list and remove 'NaN' items
269     transaction = [item for item in row if pd.
270         notna(item)]
271     transactions.append(transaction)
272
273 # Display the first two transactions
274 transactions[:2]

```

```

247
248 # Calculate the number of non-missing values in
    each transaction
249 num_non_missing_per_transaction = data_df.iloc[:,
    1:33].count(axis=1)
250
251 num_items_per_transaction = data_df["Item(s)"]
252
253 is_complete_transaction = (
    num_non_missing_per_transaction ==
    num_items_per_transaction)
254
255 # Count the number of rows with missing values or
    mismatches in item counts
256 num_incomplete_transactions = (~
    is_complete_transaction).sum()
257
258 # Print the number of incomplete transactions
259 print("Number of Incomplete Transactions:",
    num_incomplete_transactions)
260
261 transactions_df = pd.DataFrame(transactions)
262 transactions_df
263
264 """##Frequent Itemset Generation using Apriori
    Algorithm"""
265
266 from itertools import combinations
267 from collections import defaultdict
268
269 # Minimum support count
270 min_support = 2
271
272 # Generate candidate 1-itemsets
273 item_counts = defaultdict(int)
274 for transaction in transactions:
275     for item in transaction:
276         item_counts[item] += 1
277
278 # Prune infrequent items
279 frequent_1_itemsets = {item for item, count in
    item_counts.items() if count >= min_support}
280
281 # Generate frequent itemsets of size k
282 def generate_frequent_itemsets(itemsets, k):
283     candidates = set()
284     for itemset in itemsets:
285         for item in frequent_1_itemsets:
286             if item not in itemset:
287                 candidates.add(itemset + (item,))
288
289     frequent_itemsets = set()
290     for candidate in candidates:
291         count = sum(1 for transaction in
    transactions if all(item in transaction for
    item in candidate))
292         if count >= min_support:
293             frequent_itemsets.add(candidate)
294
295     return frequent_itemsets
296
297 # Generate frequent itemsets of increasing size
298 k = 2
299 frequent_itemsets = {(item,) for item in
    frequent_1_itemsets}
300 while frequent_itemsets:
301     print(f"Frequent {k-1}-itemsets: {
    frequent_itemsets}")
302     k += 1
303     candidate_itemsets =
    generate_frequent_itemsets(frequent_itemsets,
    k)
304     frequent_itemsets = candidate_itemsets
305
306 frequent_itemsets
307
308 """##Rule Generation"""
309
310 # Generate association rules from frequent
    itemsets
311 def generate_association_rules(frequent_itemset):
312     rules = []
313     for itemset in frequent_itemset:
314         for i in range(1, len(itemset)):
315             for subset in combinations(itemset, i)
    :
316                 remaining = tuple(item for item in
    itemset if item not in subset)
317                 support_itemset = sum(1 for
    transaction in transactions if all(item in
    transaction for item in itemset))
318                 support_subset = sum(1 for
    transaction in transactions if all(item in
    transaction for item in subset))
319                 confidence = support_itemset /
    support_subset
320                 if confidence >= 0.1: # Minimum
    confidence threshold
321                     rules.append((subset,
    remaining, confidence))
322     return rules
323
324 # Generate association rules from frequent
    itemsets
325 association_rules = []
326 for frequent_itemset in frequent_itemsets:
327     rules = generate_association_rules(
    frequent_itemset)
328     association_rules.extend(rules)
329
330 # Print association rules
331 for antecedent, consequent, confidence in
    association_rules:
332     print(f"Rule: {antecedent} => {consequent},
    Confidence: {confidence:.2f}")
333
334 """## Library Implementation"""
335
336 # Extract the item columns
337 item_columns = data_df.columns[1:33]
338
339 # Convert the data into a list of transactions
340 transactions = data_df[item_columns].apply(lambda
    row: row.dropna().tolist(), axis=1).tolist()
341
342 # Create a one-hot encoded DataFrame for the
    transactions
343 onehot_transactions = pd.DataFrame(transactions)
344
345 # Apply one-hot encoding
346 onehot_encoded = pd.get_dummies(
    onehot_transactions.unstack()).groupby(level
    =1).max()
347
348 from mlxtend.preprocessing import
    TransactionEncoder
349 from mlxtend.frequent_patterns import apriori,
    association_rules
350
351 # List of minimum support values
352 min_support_values = [0.1, 0.08, 0.06, 0.04, 0.01,
    0.005] # Adding more values
353
354 # Confidence levels to evaluate
355 confidence_levels = list(np.arange(0.05, 1.05,
    0.05))
356
357 # Empty lists to store results

```



```

358 num_rules_lists = []
359
360 # Calculate and store the number of rules for each
    combination of minimum support and confidence
    level
361 for min_support in min_support_values:
362     frequent_itemsets = apriori(onehot_encoded,
        min_support=min_support, use_colnames=True)
363     rules_list = []
364     for confidence_level in confidence_levels:
365         rules = association_rules(
            frequent_itemsets, metric="confidence",
            min_threshold=confidence_level)
366         num_rules = len(rules)
367         rules_list.append(num_rules)
368     num_rules_lists.append(rules_list)
369
370 # Plot the results
371 plt.figure(figsize=(10, 6))
372
373 colors = ["#4e79a7", "#f28e2b", "#e15759", "#76
    b7b2", "#59a14f", "#edc948"] # Changed line
    colors
374
375 for i, min_support in enumerate(min_support_values
    ):
376     plt.plot(confidence_levels, num_rules_lists[i
        ], marker="o", color=colors[i], alpha=0.9,
        label=f"Min Support: {min_support}")
377
378 plt.xlabel("Confidence Level")
379 plt.ylabel("Number of Rules")
380 plt.title("Number of Rules vs. Confidence Level
    for Different Minimum Support")
381
382 # Set the desired x-axis labels
383 plt.xticks([0.10, 0.20, 0.30, 0.40, 0.50, 0.60,
    0.70, 0.80, 0.90, 1])
384
385 # Add grid lines for better readability
386 plt.grid(True, linestyle="--", alpha=0.3)
387
388 plt.legend()
389
390 # Add extra legend
391 extra_legend = "ROLL 18, 25"
392 plt.text(0.5, -0.15, extra_legend, ha='center',
    color='red', transform=plt.gca().transAxes)
393
394 plt.show()
395
396 # Run Apriori algorithm with a minimum support
    threshold of 0.05
397
398 from mlxtend.frequent_patterns import apriori,
    association_rules
399 frequent_itemsets = apriori(onehot_encoded,
    min_support=0.015, use_colnames=True)
400
401 # Sort the frequent itemsets DataFrame by 'support
    ' in descending order
402 sorted_frequent_itemsets = frequent_itemsets.
    sort_values(by="support", ascending=False)
403
404 # Calculate the length of each itemset and add it
    as a new column 'length'
405 sorted_frequent_itemsets["length"] =
    sorted_frequent_itemsets["itemsets"].apply(len
    )
406
407 # Display the sorted frequent itemsets DataFrame
    with the desired formatting
408 with pd.option_context("display.max_rows", None, "
    display.max_columns", None, "display.

```

```

    float_format", '{:.3f}'.format):
    print(sorted_frequent_itemsets)
409
410 # Generate association rules
411 association_rules_df = association_rules(
    frequent_itemsets, metric="confidence",
    min_threshold=0.3)
412
413 # Sort the association rules by 'lift' metric in
    descending order and reset the index
414 sorted_association_rules = association_rules_df.
    sort_values(by="lift", ascending=False).
    reset_index(drop=True)
415
416 # Display the sorted association rules DataFrame
417 print("\nAssociation Rules:")
418 sorted_association_rules
419
420
421
422 association_rules_df = association_rules_df.
    sort_values(by='lift', ascending=False)
423
424 # Choose the top N rules to display
425 top_n = 10
426 top_df = association_rules_df.head(top_n)
427
428 # Convert frozenset objects to strings for
    concatenation
429 top_df['rule'] = top_df['antecedents'].apply(
    lambda x: ', '.join(x)) + ' -> ' + top_df['
    consequents'].apply(lambda x: ', '.join(x))
430
431 plt.figure(figsize=(10, 6))
432 plt.barh(top_df['rule'], top_df['lift'])
433 plt.xlabel('Lift')
434 plt.ylabel('Association Rule')
435 plt.title(f'Top {top_n} Association Rules by Lift'
    )
436
437 plt.tight_layout()
438 extra_legend = "ROLL 18, 25"
439 plt.text(0.5, -0.15, extra_legend, ha='center',
    color='red', transform=plt.gca().transAxes)
440 plt.show()
441
442 import networkx as nx
443 G = nx.DiGraph()
444 # Add nodes and edges
445 for index, row in association_rules_df.iterrows():
446     antecedent = ', '.join(row['antecedents'])
447     consequent = ', '.join(row['consequents'])
448     G.add_node(antecedent)
449     G.add_node(consequent)
450     G.add_edge(antecedent, consequent, weight=row[
        'lift'])
451
452 # Position nodes using a spring layout algorithm
453 pos = nx.spring_layout(G)
454 plt.figure(figsize=(12, 8))
455
456 # Draw nodes and edges
457 nx.draw(G, pos, with_labels=True, font_size=10,
    node_size=1000, node_color='lightblue')
458 nx.draw_networkx_edge_labels(G, pos, edge_labels
    ={(u, v): f'Lift: {d["weight"]:.2f}' for u, v,
        d in G.edges(data=True)})
459 plt.title('Network Graph of Association Rules')
460 extra_legend = "ROLL 18, 25"
461 plt.text(0.5, -0.05, extra_legend, ha='center',
    color='red', transform=plt.gca().transAxes)
462
463 plt.show()

```

...