

# **Boosting Retail Revenue with Low-Cost High Utility Itemsets: An LCIM Algorithm Application**

# TABLE OF CONTENTS

1. Introduction .....	6
2. Related work .....	8
3. Background.....	10
3.1 Shopping cart analysis.....	10
3.2 Frequent itemset mining .....	10
3.3 High utility itemset mining.....	10
3.4 LCIM Algorithm: .....	10
3.5 The Search procedure .....	12
3.6 The Construct process .....	12
3.7 Trade-off value .....	13
4. Data Exploration.....	14
5. Proposed model .....	18
6. Experimental Result and Analysis.....	20
6.1 Pre-processing .....	20
6.2 LCIM algorithm and evaluator.....	22
6.3 Pattern analysis .....	24
7. Conclusion .....	26
REFERENCE .....	28

# TABLE OF FIGURES

Figure 1: LCIM Algorithm.....	11
Figure 2: Sales Module.....	14
Figure 3: Product Module.....	15
Figure 4: Product Sales Insight.....	16
Figure 5: The proposed methodology framework.....	18

# TABLE OF TABLES

Table 1: Table of Abbreviation glossary .....	5
Table 2: A part of original dataset .....	21
Table 3: Parameters value of LCIM algorithm.....	22
Table 4: The result of using LCIM algorithm .....	22
Table 5: The trade-off of each pattern .....	24

# ABBREVIATION GLOSSARY

*Table 1: Table of Abbreviation glossary*

FIM	Frequent Itemset Mining
HUIM	High Utility Itemset Mining
LCIM	Low Cost High Utility Itemsets Mining

## 1. Introduction

The exponential growth of data has given rise to an ever-increasing need for a comprehensive grasp of dataset patterns, ultimately bolstering the efficacy of decision-making processes. This prompted the inception of data mining, a strategic approach aimed at unearthing intricate patterns and valuable insights from datasets through a diverse array of mining techniques and algorithms. Within the realm of data mining lies frequent pattern mining, a specialized subset that focuses on pinpointing recurring sets of data elements within the dataset. However, this technique is confined to merely assessing the prevalence of itemsets, often overlooking other critical criteria that significantly impact a company's revenue and profitability.

In response to this limitation, the concept of High Utility Itemset Mining (HUIM) was conceived to transcend the confines of conventional frequency-based methodologies. HUIM goes beyond the conventional approach by not only considering the frequency of itemsets but also weighing their potential contribution to a company's profit margins. Its applicability spans a multitude of business sectors, with the realm of retail being particularly fertile ground for its application due to the readily identifiable item sets.

Notwithstanding its utility, HUIM is not without its drawbacks. By predominantly focusing on frequency and utility, it runs the risk of neglecting the associated costs and efforts required by a company to attain a given level of utility. Consequently, the outcomes of this technique may highlight itemsets with elevated frequency and utility, yet these achievements might come at a substantial cost, leading to the inadvertent dismissal of itemsets that could yield relatively high benefits with minimal expenditure. This incongruity in the information generated could potentially misguide decision-makers in devising business strategies.

Addressing this concern, the Low Cost High Utility Itemsets Mining (LCIM) algorithm emerged as a potential solution. LCIM was conceived with the aim of uncovering itemsets within a transactional dataset that possess attributes encompassing both cost and benefit. This algorithm, by incorporating cost considerations into the analysis, endeavors to unearth

itemsets that not only exhibit high utility but also do so at a manageable cost, offering decision-makers a more balanced and accurate perspective on which to base their strategic choices.

This study will apply LCIM algorithm in transactional database to find low cost high utility itemsets and analyze the result to generate insightful information about the itemsets that is both beneficial and non-costly for business. The study also use Java to handle the data, transform it into suitable form, then implement the LCIM algorithm on the input file to produce desired output.

This report will include 7 sections: Section 1 introduces the overview of the research, with purpose of this research. The related work in the following section summarizes all the papers that are involved in the subject of the report. Section 3 describes the background of the report with the important concepts, then the data exploration part will be presented to give a better understanding about the dataset. With section 5, the model is fully illustrated with 5 stages. Next, experiment results using the model will be performed in section 6 with discussion. Section 7 will conclude the report and the future work.

## 2. Related work

As the volume of data grows, extracting concealed patterns from frequent itemsets becomes progressively laborious. Thus, a more efficient algorithm is imperative to swiftly uncover hidden patterns within frequent itemsets and minimize memory requirements, especially as data volume continues to expand over time in (Chee et al., 2019). (Agrawal et al., 1994) represents the Apriori algorithm, which is a fundamental algorithm for discovering the frequent itemsets based on the association rules. (Agrawal et al., 1994) lays the foundation for the FIM data mining techniques and helps to discover the relationship between items in the large transactional database. After Apriori, (Han et al., 2000) proposed FP-Growth, is a technique that extracts frequent itemsets without the need for an expensive generation of candidate sets. Then, EClat algorithm appeared in (Zaki et al., 2000) achieves efficient extraction of frequent itemsets through utilization of the vertical data structure. Through the utilization of the FPM algorithm, individuals can enhance decision-making speed, optimize performance and functioning, and boost their organizations' profits (Chee et al., 2019).

To develop the FIM algorithm, (Chan et al., 2003) integrated the utility factor to the frequent itemsets mining to solve the drawbacks of FIM. FIM only generates the output with the frequent itemsets without providing information about profit of those itemsets. The objective of HUIM is to identify highly significant patterns by taking into account both the quantity of purchases and the benefits associated with products (Nellutla et al., 2022). (Fournier-Viger et al., 2019). High utility itemset mining extends the concept of frequent itemset mining by considering the quantities and values associated with items. An often-applied scenario for high utility itemset mining involves identifying collections of items that customers buy together, resulting in substantial profits. This algorithm is based on the Apriori algorithm but can create the outcome of itemset with high utility and high frequency. HUIM adds value to the FIM algorithm in terms of utility. Numerous algorithms for mining high utility itemsets have been suggested, including UMining in (Yao et al., 2006) that enables users to assign numerical values to express their preferences



regarding the usefulness of itemsets. The UMining algorithm is capable of effectively identifying all valuable itemsets within a database, a task that cannot be accomplished by techniques such as frequent itemset mining, convertible constraint-based mining, or share-based mining. Two Phase algorithm proposed in (Liu et al., 2005) is also a HUIM algorithm. It effectively reduces the candidate count while accurately deriving the full collection of high utility itemsets. More and more high utility itemset mining algorithms are discovered to find out itemset that represents the statistical association among items, and its portrayal of the semantic importance attributed to the items.

However, these research investigations only concentrate on the advantages offered by patterns, neglecting details pertaining to their expenses (such as time, labor, financial resources, or other consumable assets). (Nawaz et al., 2022) will address this problem by using the LCIM algorithm that incorporates frequency, utility, and cost to produce the most frequent itemsets that have both high benefits and low expense. It is imperative to employ not only an upper limit on utility but also a lower threshold on cost to achieve efficient pattern mining. It also introduces an Average Cost Bound (ACB), which sets a minimum threshold for the average cost, aimed at diminishing the exploration scope, along with a cost-list data organization.

### **3. Background**

#### *3.1 Shopping cart analysis*

Shopping cart analysis refers to a method in data mining employed to reveal connections, correlations, and links between items that customers commonly buy in tandem. This process entails scrutinizing the contents of customer shopping carts or transaction histories to pinpoint items that are commonly acquired together. Such analysis yields valuable understandings into customer actions, inclinations, and buying routines, offering businesses opportunities for diverse applications such as marketing, optimizing inventory, and informed strategic choices.

#### *3.2 Frequent itemset mining*

Frequent itemset mining is an approach within data mining where the aim is to detect sets of items that commonly occur together within a provided dataset. This technique emphasizes the discovery of item groups that display notable instances of appearing alongside each other in a collection of transactions or records. The central objective of frequent itemset mining is to unveil patterns of item combinations that exhibit a higher occurrence rate than what random chance would account for.

#### *3.3 High utility itemset mining*

High utility itemset mining is a sophisticated and valuable data mining methodology that delves into the realm of uncovering specific sets of items residing within a dataset. These sets are not just ordinary assemblages of elements; rather, they possess a remarkable attribute of jointly bestowing upon the dataset a considerable and noteworthy utility value. This utility value encapsulates the overarching importance, significance, or desirability that these itemsets bring to the table.

#### *3.4 LCIM Algorithm:*

The LCIM algorithm, which stands for Low Cost High Utility Itemset Mining, is a specialized computational technique designed to unearth and unravel intricate patterns that exhibit a unique blend of characteristics: a notably high average utility and concurrently, a remarkably low average cost.

The input of the algorithm is transactional database with information the purchased items with their utility, cost per items and the minsup, minutil and maxcost thresholds. The result of LCIM is a collection of value groupings that frequently co-occur (referred to as itemsets), characterized by being both economically efficient and substantially beneficial. Identifying these patterns can subsequently contribute to comprehending the underlying data.

Algorithm 1: The LCIM algorithm

input :  $D$ : a transaction database,  
minsup, minutil, maxcost: the user-specified thresholds  
output : all the set of low cost itemsets  
1 Scan  $D$  to calculate the support  $s(\{i\})$  of each item  $i$ ;  
2  $I^* \leftarrow$  each item  $i$  such that  $s(\{i\}) \geq \text{minsup}$ ;  
3 Let  $\succ$  be the total order of support ascending values on  $I^*$ ;  
4 Scan  $D$  to build the cost-list of each item  $i \in I^*$ ;  
5  $I^{**} \leftarrow$  each item  $i \in I^*$  such that  $\text{acb}(\{i\}) \leq \text{maxcost}$  according to  $L(\{i\})$ ;  
6 Search ( $I^{**}$ , minsup, minutil, maxcost);

*Figure 1: LCIM Algorithm*

The initial step of LCIM involves scanning the database to compute the support level for each individual item. This process aids in identifying the collection  $I^*$ , which comprises all the frequent items surpassing the defined minimum support threshold (minsup). Following this, LCIM utilizes this gathered data to establish a sequencing criterion  $\succ$  for items. This sequence corresponds to arranging items in ascending order based on their support values. Subsequently, LCIM conducts another pass through the database to generate costlists for all items within the set  $I^*$ . These costlists are essential for the computation of the ACB lower bound for each item  $i$  belonging to  $I^*$ . Any item with a lower bound value not exceeding the specified maximum cost (maxcost) is included in the set  $I^{**}$ . The subsequent

step entails invoking the recursive process Search using the set  $I^{**}$  to explore itemsets with low costs.

### 3.5 The Search procedure

The Search process takes a set of itemsets  $P$  and their cost-lists, along with minsup, minutil, and maxcost thresholds as input. It outputs low cost itemsets within  $P$  or extensions of  $P$ . For each itemset  $X$  in  $P$ , its average utility and cost from cost-list  $L(X)$  are calculated. If  $X$ 's average utility is at least minutil and its average cost is at most maxcost,  $X$  is output as a low cost itemset. ExtensionsOf $X$  is initialized to store potential extensions of  $P$  with one additional item compared to  $X$ .

For each itemset  $Y$  in  $P$  that can be combined with  $X$ , an extension  $Z = X \cup Y$  is formed, along with its cost-list  $L(Z)$  using the Construct process. The cost-list  $L(Z)$  provides support  $s(Z)$  and ACB lower bound  $acb(Z)$  directly, without database scanning. If  $s(Z)$  is at least minsup and  $acb(Z)$  is at most maxcost,  $Z$  is added to ExtensionsOf $X$ , as  $Z$  and its recursive extensions could be low cost itemsets. The search process is then recursively applied to ExtensionsOf $X$  to explore extensions of  $X$ . Once all itemsets in  $P$  have been looped through, the output comprises low cost itemsets within  $P$  and extensions of  $P$ .

### 3.6 The Construct process

The algorithm designated as Construct, as outlined in Algorithm 3, facilitates the formation of an extension  $Z = X \cup Y$  by combining the cost-lists  $L(X)$  and  $L(Y)$  of two distinct itemsets,  $X$  and  $Y$ . The resulting output of this procedure is the cost-list  $L(Z)$ . Initially, the cost-list  $L(Z)$  is initialized, such that  $L(Z).utility$  is set to 0,  $L(Z).tids$  is assigned an empty state, and  $L(Z).costs$  is also initialized as empty.

A loop is then executed to examine each transaction  $Tw$  present in  $L(X).tids$ , checking whether it is also present in  $L(Y).tids$ . For each transaction  $Tw$  that fulfills this criterion, it is added to  $L(Z).tids$ , and the corresponding utility  $u(Tw)$  is incorporated into  $L(Z).utility$ .

Furthermore,  $L(Z).costs$  is updated by summing the costs associated with  $Z$  in transaction  $Tw$ . This summation is performed using a procedure termed Merge, which involves adding the costs of  $X$  and  $Y$  within  $Tw$  (this information is sourced from  $L(X)$  and  $L(Y)$ ).

Upon the completion of this loop, the total cost  $L(Z).cost$  is calculated as the cumulative sum of all values within  $L(Z).costs$ . Subsequently, the final result, denoted as  $L(Z)$ , is returned by the procedure. This algorithmic process essentially constructs the cost-list  $L(Z)$  for the extension  $Z$  by aggregating and manipulating information from the input cost-lists  $L(X)$  and  $L(Y)$ .

### *3.7 Trade-off value*

The efficiency of a pattern can be determined by calculating its trade-off value, which is the ratio between the average cost and the average utility of the pattern. This trade-off value serves as a measure of the pattern's effectiveness. A pattern with a lower trade-off value is considered cost-effective, as it delivers utility while requiring minimal resources to achieve a desirable outcome (Fournier-Viger et al., 2020).

$$tf(p) = \sum_{p \subseteq S_s \wedge S_s \in SEL} (c(p, S_s) / su(S_s))$$

## 4. Data Exploration

The AdventureWorks dataset, initially crafted by Microsoft, serves as a model database. It mimics a fictional bicycle enterprise named Adventure Works Cycles and encompasses information spanning sales, clientele, items, and personnel.

The AdventureWorks dataset encompasses two primary sections: Sales and Product. These segments offer pertinent data regarding sales activities and product details. Together, they furnish a fitting foundational dataset for applying the LCIM algorithm.

Within the transactional records, there are 334 distinct products, each associated with unit quantities and corresponding costs per transaction.

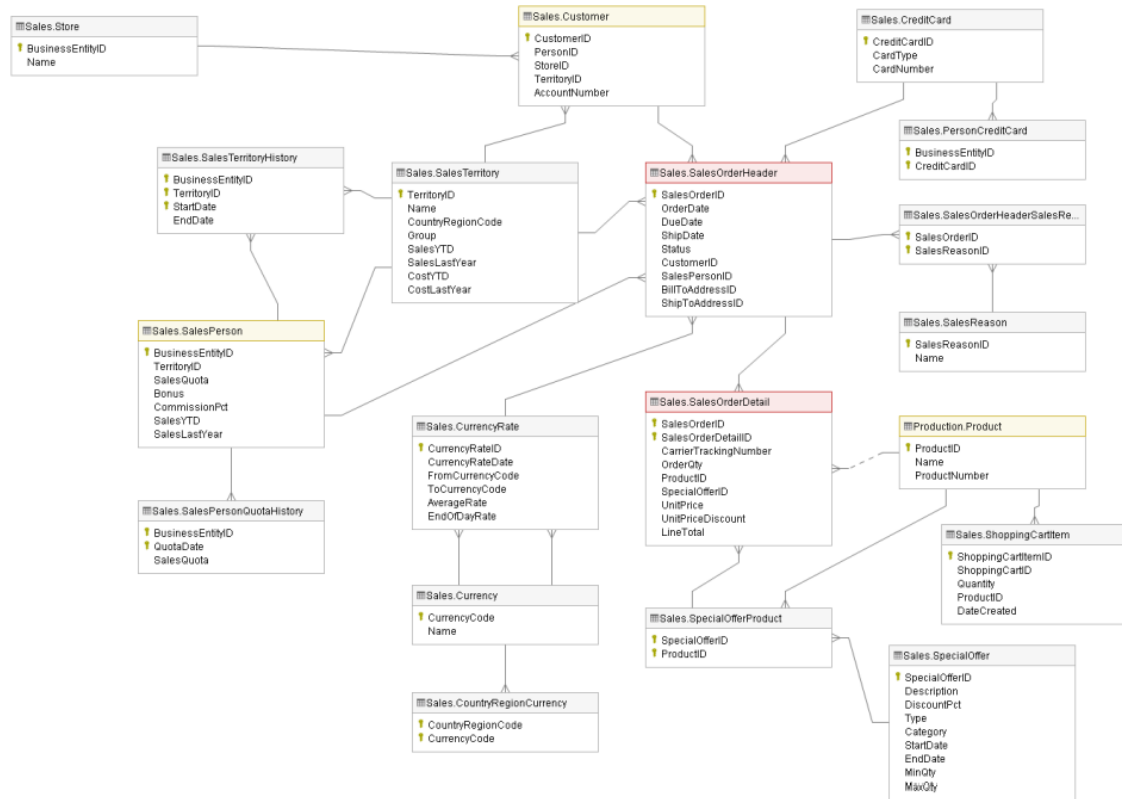
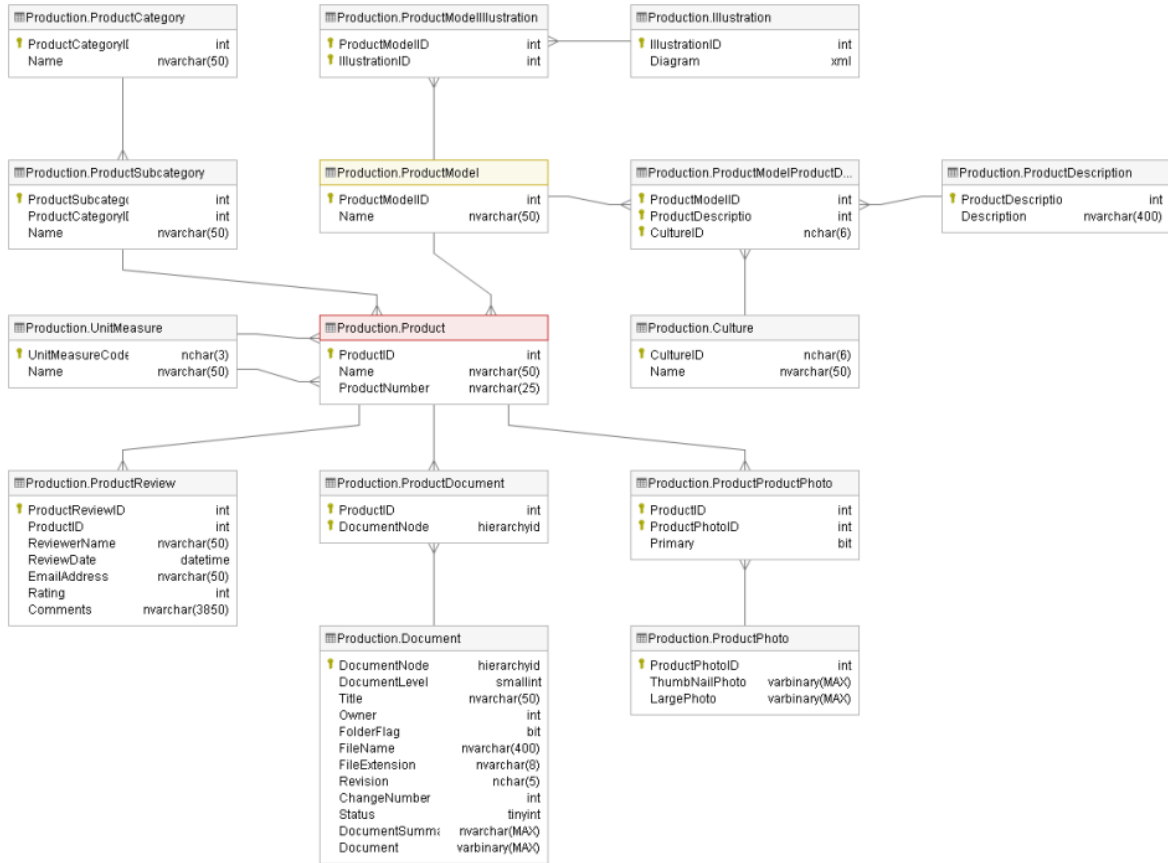


Figure 2: Sales Module

The sales module contained within the AdventureWorks dataset serves as an all-encompassing and constantly evolving portrayal of sales-centric endeavors occurring within an imaginative establishment operating within either the retail or manufacturing sector. Crafted with the intention of emulating genuine commercial activities, this

particular module furnishes an extensive collection of data markers, collectively presenting a treasure trove of invaluable understandings spanning a multitude of dimensions within the intricate realm of sales procedures. By granting the capability for meticulous scrutiny, refinement, and well-informed deliberations, this module profoundly empowers the undertaking of dissecting, enhancing, and ultimately making judicious choices in the convoluted arena of sales processes.



*Figure 3: Product Module*

Situated within the expansive realm of the AdventureWorks dataset, the product module stands as a meticulously crafted and multifaceted representation, offering an elaborate and intricate depiction of a wide-ranging assortment of products that hold pivotal roles within the intricate workings of a simulated retail or manufacturing establishment. Every facet of this module has been thoughtfully architected to emulate the intricate subtleties inherent in

the domain of real-world product management and representation, thereby serving as a dynamic and immersive reservoir of insights.

As users venture into the virtual landscape of this product module, they are greeted with a meticulously curated compilation of data elements, each contributing to the rich tapestry of product-related dynamics. These data elements converge to present a comprehensive and multifaceted view encompassing an assortment of attributes, encompassing details such as product classifications, categories, pricing configurations, and historical evolution. This collection of information transcends mere numbers; it encapsulates the very essence of products as they traverse their lifecycles, capturing the ebb and flow of market trends, consumer preferences, and technological innovations.

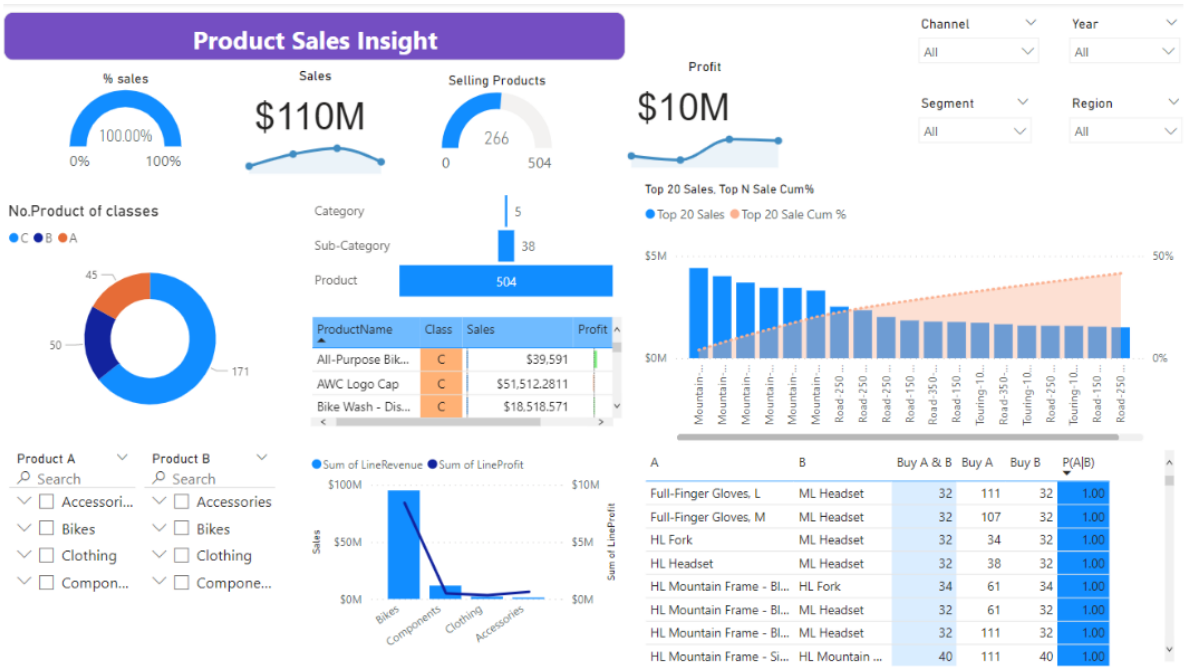


Figure 4: Product Sales Insight

The "Product Sales Insight" Figure offers a comprehensive view of the enterprise's product-related business performance. It presents the total sales and profits, along with a breakdown of revenue by product, highlighting the top 20 most frequently purchased items. The report employs the ABC classification approach to illustrate how products contribute to revenue:



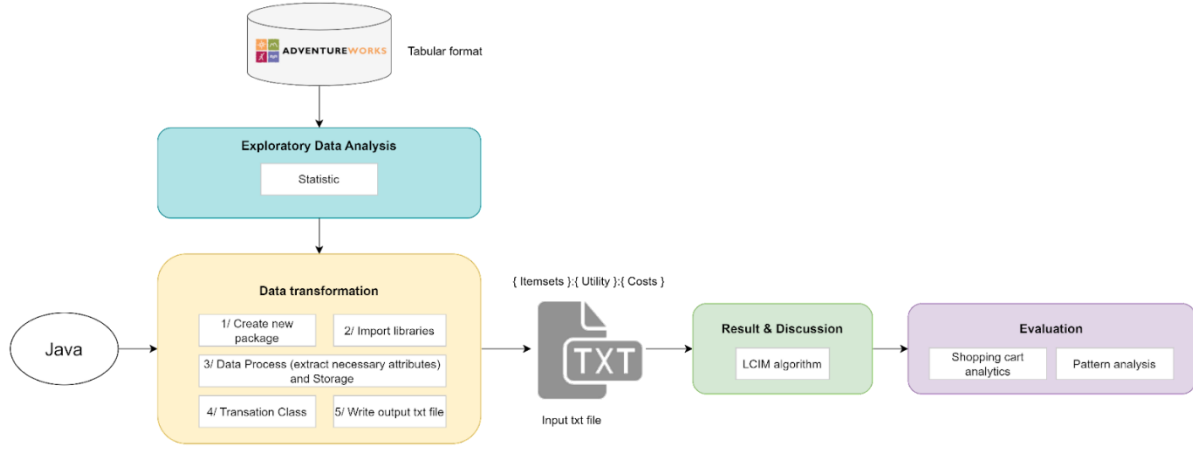
Class A: Represents high-revenue products, accounting for 70% of the total revenue.

Class B: Encompasses medium-revenue products, contributing 20% of the total revenue.

Class C: Includes low-revenue products, making up 10% of the total revenue.

Additionally, the report addresses the question of product co-purchases by utilizing correlation analysis to determine the likelihood of customers buying specific pairs of products together. To further enhance this probability, the report employs LCIM algorithms to identify itemsets that offer both high utility and low cost.

## 5. Proposed model



*Figure 5: The proposed methodology framework*

As depicted in Figure 5, this section delineates the comprehensive methodology employed to accomplish the research objectives. The methodology consists of four meticulous stages, each thoughtfully crafted and executed to ensure the integrity and precision of the results. These stages encompass data selection, exploratory data analysis (EDA), data transformation, algorithmic application, and culminate in conclusive shopping cart and pattern analyses.

*Stage 1: Data Selection and Acquisition.* The initial stage of this study involved the meticulous selection of a suitable dataset to serve as the foundation for analysis. After thoroughly evaluating available options, the AdventureWork2019 dataset was identified as the most pertinent due to its comprehensive nature and alignment with the research goals. This dataset, which encapsulates a wealth of business-related information, was deemed ideal for investigating the patterns and relationships of interest.

*Stage 2: Exploratory Data Analysis (EDA).* Following the acquisition of the AdventureWork2019 dataset, an in-depth exploratory data analysis (EDA) was conducted. EDA is a critical phase in understanding the characteristics and structure of the dataset. This involved a systematic examination of data distributions, central tendencies, variabilities, and preliminary insights into potential trends or anomalies. By uncovering

these foundational insights, the subsequent stages of the methodology were informed and shaped.

*Stage 3: Data Transformation.* In this stage, the tabular format of the dataset was transformed into a structured text (txt) format. This transformation was implemented with meticulous care to ensure the integrity and fidelity of the data. This txt format was chosen to enhance the compatibility of the dataset with the subsequent application of LCIM (Low-Cost High Utility Itemsets) algorithms. The transformation process was executed while adhering to best practices for data formatting and preservation.

*Stage 4: Application of LCIM Algorithms and Analytical Evaluation.* The fourth and final stage of the methodology encompassed two parallel processes: the application of (LCIM) algorithms and subsequent analytical evaluation. LCIM algorithms were deployed to uncover Low-Cost High Utility Itemsets and associations within the transformed dataset. This involved the identification of relationships that might not be immediately apparent through traditional analytical techniques.

Subsequently, the evaluation phase focused on two primary techniques: shopping cart analysis and pattern analysis. Shopping cart analysis was employed to unveil correlations between various items and products, shedding light on consumer behavior and preferences. Pattern analysis, conversely, was directed towards identifying the optimal trade-off where the itemset achieves the most favorable outcome in terms of Low-Cost High Utility Cost (LCIM).

## 6. Experimental Result and Analysis

### 6.1. Pre-processing

The data transformation process comprises the following steps:

#### *Step 0: Creation of a Package and Java File within the Package*

The initial step involves creating a designated package and a Java file within it to facilitate organized code structuring.

#### *Step 1: Importing Relevant Libraries*

Appropriate libraries are imported to enable the necessary functionalities. These include `java.io.FileReader`, `java.io.FileWriter`, `java.io.IOException`, `java.io.Reader`, `java.io.BufferedReader`, `java.util.HashMap`, `java.util.Map`, and `java.util.TreeMap`.

#### *Step 2: Parsing CSV Data*

The program reads the CSV file named "RetailSales.csv". The Apache Commons CSV library is utilized to parse the CSV file's syntax using `CSVParser`. The default format and headers are employed for this parsing.

#### *Step 3: Data Processing and Result Storage*

Each record (row) in the CSV file is processed individually. Essential fields (such as `SalesOrderNumber`, `ProductID`, `ExtendedAmount`, `TotalProductCost`) are extracted from the `CSVRecord`. The program executes necessary transformations and calculations, saving the results in a `Map<Integer, Transaction>` named `transactionsMap`.

#### *Step 4: Transaction Class Creation*

A `Transaction` class is established to manage and process data for each transaction, identified by `SalesOrderID`. This class encompasses methods to aggregate total lines and costs for different products, compute transaction utility, retrieve `ProductID` lists, and obtain costs of individual products within a transaction.

#### *Step 5: Writing Results to Output File*

Upon processing all records in the CSV file, the program writes the outcomes into an output file named "input\_file.txt" using FileWriter. The results are formatted specifically as "ProductIDs:TransactionUtility:TotalCostOfEachItem".

By meticulously following these steps, the data transformation phase ensures the conversion of the initial tabular data into a structured text format that is optimally suited for subsequent analysis using LCIM (Low-Cost High Utility Cost) algorithms. This transformation process ensures data integrity and appropriateness for the overarching analytical objectives of the study.

*Table 2: A part of original dataset*

<b>OrderNumber</b>	<b>ProductID</b>	<b>Quantity</b>	<b>UnitPrice</b>	<b>Cost</b>
SO65283	359	2	1377	2504
SO65283	474	8	42	209
SO65283	475	3	42	79
SO65283	476	1	42	26

As indicated in Table X, the initial dataset is presented in a tabular format. However, following the completion of the transformation stage, the dataset undergoes a conversion process resulting in a text-based format, such as the example:

"359 474 475 476:3258:2504 209 79 26".

The significance of the data line is outlined below:

- "359 474 475 476" constitutes a set of items representing nominal values, specifically denoting the ProductID associated with each item within the transaction.
- "3258" represents a utility value expressed as an integer. This value is attributed to the transaction, signifying the total revenue generated from the transaction.

- "2504 209 79 26" indicates cost values expressed as integers corresponding to each individual item within the transaction. In precise terms, these values denote the cost associated with each respective ProductID within the transaction.

Through this comprehensive explanation, the intricate components of the transformed dataset become clear. The adoption of this structured text format serves to facilitate subsequent analysis employing the Low-Cost High Utility Cost (LCIM) algorithms, allowing for the identification of latent patterns and associations within the dataset.

## 6.2. LCIM algorithm and evaluator

After we got the converted input txt file with standard format, we continue to set the appropriate parameters.

*Table 3: Parameters value of LCIM algorithm*

Parameter	Value
minutility	40000
maxcost	10000
minsup	0.1

As depicted in Table 3, this entails that the resulting outcome needs to meet the following criteria:

- The average utility must not fall below the minimum threshold of 40,000 monetary units.
- The average cost must not exceed the maximum threshold of 10,000 monetary units.
- The support (frequency of occurrence) must be greater than or equal to the minimum threshold of 10% within the dataset.

*Table 4: The result of using LCIM algorithm*

No.	Patterns	Avg. Util	Avg. Cost	Sup
1	456	40742.81	152.28	418
2	456, 233	42500.67	322.07	384
3	456, 458	41703.15	312.08	392
4	456, 224	42594.28	187.12	385
5	233, 458	40403.68	313.84	408
6	233, 458, 224	40859.32	503.68	398
7	458, 224	40499.08	180.82	410

As shown in table 4, the outcomes encompass all the category groups representing several effective product clusters within sales. These groups consist of product combinations that lead to high profitability and do not entail significant expenditure within a designated timeframe.

For a more thorough clarification of the outcomes, consider the instance presented in the sixth row:

Data Line: "233 458 224 #AUTIL: 40859.324120603014 #ACOST: 503.68090452261305 #SUP: 398"

- An "itemset" constitutes a collection of "productID" entities, as demonstrated by "233 458 224".
- The "Utility" of an "itemset" represents the cumulative revenue generated by the given "itemset". In this instance, the value of "AUTIL" stands at "40,859", computed by aggregating the utilities of all "itemsets" and subsequently dividing it by its own "support".
- The "Cost" associated with this "itemset" pertains to the expenses incurred for each individual product. "ACOST", quantified at "503", represents the average cost of

this specific "itemset". This is determined by dividing the total cost by its own "support".

- Consequently, the "support" metric of this "itemset" serves as the count of transactions that include this particular "itemset". In this context, the value is "398", indicating the presence of this "itemset" within a total of 398 transactions.

This illustrative example offers a clear breakdown of how the various components interrelate within the context of an "itemset", emphasizing the role of "utility", "cost", and "support" in comprehending its significance within the dataset.

### 6.3. Pattern analysis

In the pursuit of enhanced decision-making and operational efficacy, the process of pattern analysis unveils invaluable insights by exploring trade-offs between costs and benefits within datasets. This report delves into the fundamental concept of trade-off within the context of pattern analysis, aiming to illuminate how it influences strategic decision-making.

Comparing the effectiveness of various datasets necessitates a comprehensive evaluation of the trade-off between costs and benefits. For a given dataset (p), the trade-off is calculated as the ratio between its average cost and its average benefit. In essence, it quantifies the equilibrium between the expenses incurred and the benefits reaped. A dataset with a low trade-off value is regarded favorably, as it signifies a remarkable balance between generating substantial revenue and minimizing expenditures.

*Table 5: The trade-off of each pattern*

No.	Patterns	Avg. Util	Avg. Cost	Sup	Trade-off
1	456	40742.81	152.28	418	<b>0.0037</b>
2	456, 233	42500.67	322.07	384	0.0076
3	456, 458	41703.15	312.08	392	0.0075



4	456, 224	42594.28	187.12	385	0.0044
5	233, 458	40403.68	313.84	408	0.0078
6	233, 458, 224	40859.32	503.68	398	0.0123
7	458, 224	40499.08	180.82	410	0.0045

As evident from the data presented in Table 5, it is notable that the first data line possesses the lowest trade-off value, measuring at an impressively minimal "0.0037". This value epitomizes excellence and stands as an exemplar of optimal equilibrium between cost and benefit. The significance of this observation lies in its ability to yield the most favorable outcome, aligning seamlessly with the principles of Low-Cost High Revenue strategy. This unique attribute positions the identified data line as a prime candidate for strategic consideration by businesses aiming to elevate their revenue streams. The data's inherent potential to deliver high revenue while maintaining prudent expenditure places it in a strategic advantage, urging enterprises to carefully deliberate upon this particular itemset with the prospect of fostering significant revenue growth.

## 7. Conclusion

In conclusion, this study has effectively demonstrated the significance of three pivotal components in advancing the realms of data analysis and business strategy within the retail domain.

Firstly, the utilization of Java programming has been pivotal in achieving seamless data transformation from raw datasets into a structured text format that aligns with the specific algorithmic requirements. This technical prowess is vital as it ensures that the input data is optimally configured for subsequent algorithmic analysis. By employing Java, we have successfully bridged the gap between raw data and algorithmic compatibility, thereby enhancing the accuracy and effectiveness of subsequent analyses.

Secondly, the grasp and application of shopping cart analytics, coupled with the intricate LCIM algorithm, have greatly enriched the understanding of consumer behavior and pattern identification within the retail landscape. By delving into the depths of shopping cart analytics, this study has unraveled latent trends and correlations that are otherwise concealed within complex retail datasets. The adoption of the LCIM algorithm, specifically tailored for the retail domain, has empowered the identification of patterns that lead to a Low-Cost High Revenue strategy. This holistic approach has the potential to revolutionize retail operations by providing actionable insights that can steer revenue growth while simultaneously managing costs judiciously.

Thirdly, the concept of trade-off emerges as a central pillar in the realm of pattern analysis. It underscores the interplay between costs and benefits, illuminating datasets that exhibit the optimal equilibrium. As businesses navigate the complexities of modern markets, the identification and cultivation of datasets with low trade-off values can serve as a pivotal strategy for achieving the twin objectives of revenue maximization and cost minimization. This report lays the foundation for strategic decision-making, emphasizing the potency of low trade-off datasets as catalysts for informed and impactful business maneuvers.

In summary, the integration of Java-driven data transformation and the sophisticated LCIM algorithm-driven shopping cart analytics serves as a potent toolset for retail enterprises

seeking to extract maximum value from their datasets. The amalgamation of technical acumen and strategic understanding propels businesses towards a data-informed future, fostering efficiency, innovation, and a sustainable competitive edge in the dynamic retail landscape.

While this study leveraged Java-driven data transformation and LCIM algorithm-powered shopping cart analytics to uncover valuable retail insights, potential future enhancements could involve incorporating sequence and time elements for a more nuanced understanding of customer behavior. Additionally, broadening the analysis to include high-frequency items, even if they don't strictly meet the utility and cost criteria, could offer a more holistic view of product popularity. Ensuring scalability for real-time implementation, integrating external factors like market trends, and considering the impact of temporal dynamics are promising directions to enhance the applicability and robustness of the methodologies in a dynamic retail landscape.

## REFERENCE

- Chee, C. H., Jaafar, J., Aziz, I. A., Hasan, M. H., & Yeoh, W. (2019, December 1). Algorithms for frequent itemset mining: a literature review. *Artificial Intelligence Review*. Springer Netherlands. <https://doi.org/10.1007/s10462-018-9629-z>
- Agrawal, R., Srikant, R.: Fast algorithms for mining association rules. In: *Proceedings of VLDB*. vol. 1215, pp. 487–499 (1994)
- Han, J., Pei, J., & Yin, Y. (2000). Mining frequent patterns without candidate generation. *SIGMOD Record (ACM Special Interest Group on Management of Data)*, 29(2), 1–12. <https://doi.org/10.1145/335191.335372>
- Zaki, M. J. (2000). Scalable algorithms for association mining. *IEEE Transactions on Knowledge and Data Engineering*, 12(3), 372–390. <https://doi.org/10.1109/69.846291>
- Chan, R., Yang, Q., Shen, Y.: Mining high utility itemsets. In: *Proceedings of ICDM*. pp. 19–26 (2003)
- Nellutla, A., & Srinivasan, N. (2022). A Survey on Analysis of Data Mining Algorithms for High Utility Itemsets. *El-Cezeri Journal of Science and Engineering*, 9(3), 1085–1100. <https://doi.org/10.31202/ecjse.1075528>
- Fournier-Viger, P., Chun-Wei Lin, J., Truong-Chi, T., & Nkambou, R. (2019). A Survey of High Utility Itemset Mining. In *Studies in Big Data (Vol. 51, pp. 1–45)*. Springer Science and Business Media Deutschland GmbH. [https://doi.org/10.1007/978-3-030-04921-8\\_1](https://doi.org/10.1007/978-3-030-04921-8_1)
- Yao, H., & Hamilton, H. J. (2006). Mining itemset utilities from transaction databases. *Data and Knowledge Engineering*, 59(3 SPEC. ISS.), 603–626. <https://doi.org/10.1016/j.datak.2005.10.004>
- Liu, Y., Liao, W. K., & Choudhary, A. (2005). A two-phase algorithm for fast discovery of high utility itemsets. In *Lecture Notes in Computer Science (including subseries Lecture Notes in Artificial Intelligence and Lecture Notes in Bioinformatics) (Vol. 3518 LNAI, pp. 689–695)*. Springer Verlag. [https://doi.org/10.1007/11430919\\_79](https://doi.org/10.1007/11430919_79)

Nawaz, M. S., Fournier-Viger, P., Alhusaini, N., He, Y., Wu, Y., & Bhattacharya, D. (2022). LCIM: Mining Low Cost High Utility Itemsets. In Lecture Notes in Computer Science (including subseries Lecture Notes in Artificial Intelligence and Lecture Notes in Bioinformatics) (Vol. 13651 LNAI, pp. 73–85). Springer Science and Business Media Deutschland GmbH.

[https://doi.org/10.1007/978-3-031-20992-5\\_7](https://doi.org/10.1007/978-3-031-20992-5_7)

Bramer, M. (2020). Principles of data mining - 4E. Drug Safety (Vol. 30, pp. 621–622).

Fournier-Viger, P., Li, J., Lin, J. C. W., Chi, T. T., & Uday Kiran, R. (2020). Mining cost-effective patterns in event logs. Knowledge-Based Systems, 191. <https://doi.org/10.1016/j.knosys.2019.105241>