

An Introduction to the Market Basket Analysis

Team 7 - Guillem Amat, Zhiwei (Calvin) Dong, Jose Luis Moscoso, Varun Prasad, Anshupriya Srivastava

IDS 705: Principles of Machine Learning

April 18, 2020

Abstract

Association rule learning is a form of unsupervised machine learning used to identify sets of variables that most commonly appear together in a database. One of the most popular association rule learning methods is the Market Basket Analysis (MBA), used often by retailers to identify products that customers frequently purchase together. The results of this analysis have many real-world consequences, such as better enabling retailers to make marketing decisions, including discounted deals, promotional pricing, and product placement. In this report, we further explore the applications of the MBA, provide an overview of the theory and mathematics behind its apriori algorithm, and demonstrate its performance on an online retail dataset.

Introduction

Association rule mining is a branch of unsupervised machine learning that identifies commonly occurring instances in a set or database. These rules indicate which items appear together with a probability much higher than from random sampling.¹ The simplicity and interpretability of association rules make them highly valuable in data science. One of association rule learning's most common methods is the Market Basket Analysis (MBA), whose name is derived from analyzing products that customers frequently purchase together in a supermarket.¹ Retailers can use the results of the MBA to adjust their inventory and improve their marketing strategies to focus on product promotions or deals that will appeal greatly to consumers. The principles behind the MBA also have applications in areas such as medicine and bioinformatics, described further in Background.²

The MBA showed its benefits in 2004, when a series of hurricanes crossed Florida. In response to the increased shopping demand, Walmart used association rule mining on their massive retail database to identify the most frequently purchased items during hurricanes. Interestingly, they found that strawberry pop tarts, and not essentials such as bottled water and flashlights, sold at seven times their normal rate. As a result, Walmart increased their stocks of pop tarts, improving their sales and providing for their customers.³ Such a response would not have been easily attainable without effective association rule mining.

In this report, we will further describe the context of association rule learning and some more applications of the MBA. We will also describe the MBA's apriori algorithm and demonstrate its performance on an online retail dataset.

Background

Association rules in the context of the Market Basket Analysis can be seen as a simple type of *Collaborative Filtering* recommender system, or an alternative to them, that tries to find interesting associations between items.⁴ They are also frequently regarded as a data mining method.

There are many reasons why association rules are used rather than more sophisticated or recent models, including the following:

1. Advanced ML methods that are customer-centric require large amounts of sparse data to train models. Association rules are item-based, meaning that transactions across many users are used more efficiently.⁵
2. Sophisticated methods require extensive use of computational power and infrastructure. Due to association rules' simple mathematical base, they are a good fit for batch operations on commodity data infrastructures such as Hadoop.⁶
3. Association rules work well with categorical data in their original state in contrast to other methods that require such data to be encoded.⁷
4. Association rules work well with unlabeled data, which is cheap and easy to collect and store.

Association rules have been widely used in industry and research. We summarize notable applications below:

1. Molecular Biology → Better identify the co-occurrences of different amino acids in proteins through association rules, improving scientists' understanding of protein structure, protein interaction, and artificial protein synthesis.⁸
2. Medicine → Improved identification of symptoms associated with diseases such as cancer, improving patient diagnosis and treatment.⁹
3. Telecommunications Industry → Investigate customer churn for the activation of mobile services based on the analysis of call center transaction data.¹⁰
4. Economics → Discover and study agglomerations of companies or individuals.¹¹
5. Banking Industry → Identify the preferences of different groups of customers to tailor recommendations.¹²

Methods

An overview of the apriori algorithm used in MBA is described in this section. The methodology for determining association rules and establishing their significance will be explained using the simple transaction list shown in Figure 1.

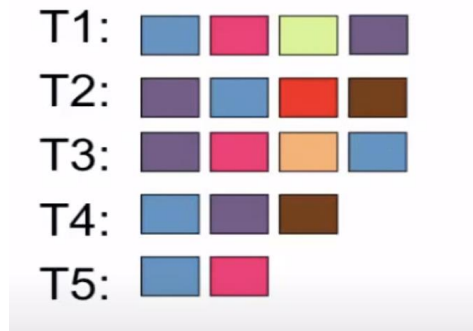


Figure 1: List of simple transactions to explain MBA metrics and apriori algorithm.¹⁵ There are five transactions: $T = \{T1, T2, T3, T4, T5\}$. Each transaction consists of a list of elements from $X = \{\text{Blue, Pink, Yellow, Purple, Brown, Peach, Red}\}$. The goal is to find the values of X that appear most frequently together in T , which is accomplished using association rules.

An association rule can be represented as an if-then statement with the following form: *If event A has occurred, then event B is likely to occur*. The first part of this statement is known as the antecedent, and the second part is known as the consequent. These rules can also be written as the following: $A \rightarrow B$. Here, A, the antecedent, and B, the consequent, are each a disjointed list of items. While creating these associations it is imperative to remember that the association implies co-occurrence, not causality.¹⁶

Prior to creating the association rules, we generate frequent itemsets, which represent the collection of items that appear above the frequency threshold, defined as the minimum number of occurrences required to categorize an itemset as frequent. For example, if a threshold is set to 2, only itemsets that appear more than twice are considered. Figure 2 illustrates this concept.

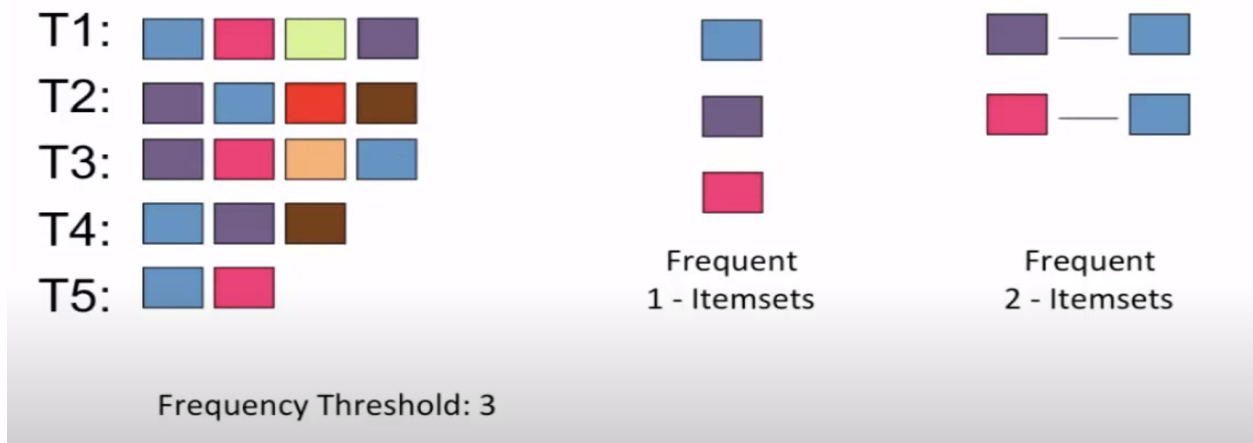


Figure 2: Collection of frequent itemsets based on a frequency threshold of 3.¹⁵ 1-Itemsets {Blue}, {Purple}, and {Pink} occur in 5, 4, and 3 transactions, respectively, so they are the most frequent 1-itemsets. {Purple, Blue} and {Pink, Blue} occur in 4 and 3 transactions, respectively, so they are the most frequent 2-itemsets.

Based on the frequent 2-itemsets generated in Figure 2, our association rules are the following:

$$\{\text{Purple}\} \rightarrow \{\text{Blue}\}; \{\text{Pink}\} \rightarrow \{\text{Blue}\}$$

There are several key metrics that establish the strength of the association between the antecedent and consequent. These are the support, confidence, and lift.

Support is a measure of the frequency of an itemset in the transactions T, represented as follows:

$$\text{Support}(\{X\} \rightarrow \{Y\}) = \frac{\text{Transactions containing both } X \text{ and } Y}{\text{Total number of transactions}}$$

$$\text{Support}(\{\text{Purple}\} \rightarrow \{\text{Blue}\}) = \frac{\text{Transactions containing both Purple and Blue}}{\text{Total number of transactions}} = \frac{4}{5}$$

Confidence is the conditional probability of occurrence of the consequent if the antecedent is present in the transaction. For the association rule $\{\text{Purple}\} \rightarrow \{\text{Blue}\}$ it represents the fraction of transactions with Purple that also have Blue. The formula is written as follows:

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

$$\text{Confidence}(\{\text{Purple}\} \rightarrow \{\text{Blue}\}) = \frac{\text{Support}(\{\text{Purple} \rightarrow \text{Blue}\})}{\text{Support}(\{\text{Purple}\})} = 1$$

Unlike support, confidence is not a commutative measure. This is demonstrated in Figure 3 and in the equation below when we calculate the confidence for the association rule $\{\text{Blue}\} \rightarrow \{\text{Purple}\}$:

$$\text{Confidence}(\{\text{Blue}\} \rightarrow \{\text{Purple}\}) = \frac{\text{Support}(\{\text{Blue} \rightarrow \text{Purple}\})}{\text{Support}(\{\text{Blue}\})} = \frac{4}{5}$$

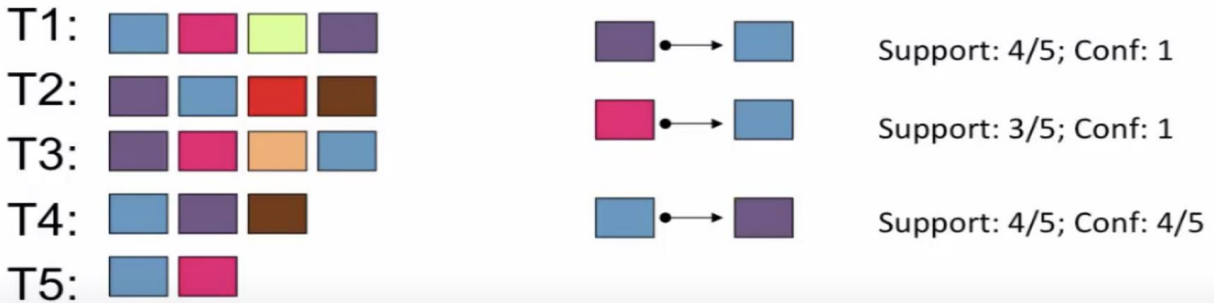


Figure 3: Sample calculations of support and confidence for different itemsets.¹⁵ It can be clearly seen that while the support for $\{\text{Blue}\} \rightarrow \{\text{Purple}\}$ is the same as that for $\{\text{Purple}\} \rightarrow \{\text{Blue}\}$, the confidence values are different due to condition probability. The confidence for the rule $\{\text{Purple}\} \rightarrow \{\text{Blue}\}$ indicates that if transactions with Purple are selected, Blue will appear in all those transactions. However, for the rule $\{\text{Blue}\} \rightarrow \{\text{Purple}\}$, transactions with Blue will only contain Purple 80% of the time.

Lift is used to determine if an association rule is truly important by considering the support for both items. A rule with high confidence may only occur because both items are common and not because they are associated with each other. Mathematically, the lift is the ratio between the confidence of the rule and the fraction of transactions containing the rule:

$$Lift(\{X\} \rightarrow \{Y\}) = \frac{Confidence(\{X\} \rightarrow \{Y\})}{Support(\{Y\})}$$

Lift is not commutative either. It represents the increase in probability of element Y being a part of the transaction if X is already a part of the transaction. Hence, a lift value greater than 1 indicates that having X in the transaction increases the probability of having Y in the same transaction. A lift equal to 1 means the items are independent. Examples are shown below:

$$Lift(\{Purple\} \rightarrow \{Blue\}) = \frac{Confidence(\{Purple\} \rightarrow \{Blue\})}{Support(\{Blue\})} = \frac{1}{1} = 1$$

$$Lift(\{Purple\} \rightarrow \{Brown\}) = \frac{Confidence(\{Purple\} \rightarrow \{Brown\})}{Support(\{Brown\})} = \frac{0.5}{0.4} = 1.25$$

These metrics are utilized by the apriori principle to efficiently develop association rules based on the frequent itemsets. The principle states the following:

“All subsets of a frequent itemset must be frequent”¹⁷

In other words, if an n-itemset is not a frequent itemset, then any (n+1)-itemset containing all elements in the n-itemset cannot be frequent. This allows us to prune off itemsets and arrive at the final frequent itemsets list quickly.

For example, in Figure 1, we set the frequent threshold to 3. As a result, the minimum support threshold is $\frac{3}{5} = 0.6$. Based on this assumption, we derived 3 1-itemsets:

$$\begin{aligned} Support(\{Blue\}) &= 1 \\ Support(\{Purple\}) &= 0.8 \\ Support(\{Pink\}) &= 0.6 \end{aligned}$$

Based on these results, we can establish that any itemset that does not contain these elements cannot be a frequent itemset. This leaves us with 3 possible 2-itemsets:

$$\begin{aligned} Support(\{Blue, Pink\}) &= 0.6 \\ Support(\{Purple, Pink\}) &= 0.4 \\ Support(\{Blue, Purple\}) &= 0.8 \end{aligned}$$

Since our minimum support threshold is 0.6, we have only two frequent 2-itemsets and can eliminate {Purple, Pink} from all future calculations. Based on the two frequent 2-itemsets we get only one 3-itemset:

$$\text{Support}(\{Blue, Pink, Purple\}) = 0.4$$

Since our minimum support threshold is 0.6, the 3-itemset is not frequent and our association rules will be created based on the 2 frequent 2-itemsets.

Once these frequent itemsets are generated, we need to identify the association rules whose confidence is greater than the minimum confidence threshold. We prune these association rules in a way similar to pruning itemsets.

Consider the itemset {Blue, Purple} with support 0.8. For the 2 rules for this itemset, we get:

$$\begin{aligned}\text{Confidence}(\{Blue\} \rightarrow \{Purple\}) &= \frac{4}{5} = 0.8 \\ \text{Confidence}(\{Purple\} \rightarrow \{Blue\}) &= \frac{4}{4} = 1\end{aligned}$$

This implies that both these association rules are strong, especially $\{Purple\} \rightarrow \{Blue\}$ since whenever Purple is present in a transaction, Blue is always present in the same transaction.

Similarly, for the itemset {Blue, Pink} with support 0.6:

$$\begin{aligned}\text{Confidence}(\{Blue\} \rightarrow \{Pink\}) &= \frac{3}{5} = 0.6 \\ \text{Confidence}(\{Pink\} \rightarrow \{Blue\}) &= \frac{3}{3} = 1\end{aligned}$$

When we calculate lift for the association rules $\{Pink\} \rightarrow \{Blue\}$ and vice versa, the value is 1, so the rules do not hold any significance. In general, larger datasets are needed to establish realistic relationships.

To summarize, the steps involved in the apriori algorithm are the following:

1. Identifying frequent itemsets
2. Setting a frequency threshold
3. Calculating support value for the itemsets
4. Selecting itemsets with support value greater than frequency threshold
5. Creating Association Rules based on selected itemsets
6. Calculating confidence to identify strong and weak associations
7. Calculating lift to explain the occurrence of an item Y when an item X is present in a transaction
8. Return rules that follow confidence and lift conditions

Using these principles, we will identify important rules in an online retail dataset.

Examples

Data Background

We used a dataset from a UK online retailer to demonstrate association analysis (<http://archive.ics.uci.edu/ml/datasets/Online+Retail>). The transnational dataset contains all the transactions occurring between 01/12/2010 and 09/12/2011. The dataset contains 8 variables shown in Table 1. There are 539,209 purchased items and 23,502 unique transactions in the dataset.

Table 1: Variables in the online retail dataset

<i>Variable Name</i>	<i>Data Type</i>	<i>Description</i>
InvoiceNo	String	Invoice number; a 6-digit integral number uniquely assigned to each transaction
StockCode	String	Product (item) code; a 5-digit integral number (plus a character) uniquely assigned to each distinct product
Description	String	Product (item) name; Ex: HERB MAKER BASIL
Quantity	Number	The quantities of each product (item) per transaction
UnitPrice	Number	Product price per unit in sterling; £45.23
InvoiceDate	Date	The day and time when each transaction was generated; 31/05/2011 15:59
CustomerID	String	Customer ID; a 5-digit integral number uniquely assigned to each customer
Country	String	Delivery address country; England

Data Preprocessing

We removed all records with missing values or abnormal values (wrong descriptions in *Description* column). Invoice numbers starting with “C” indicate cancelled transactions and were also removed. The final dataset used for analysis has 529,509 records, 19,935 unique transactions, and 4093 unique items. The distribution of these items is shown in Figure 4.

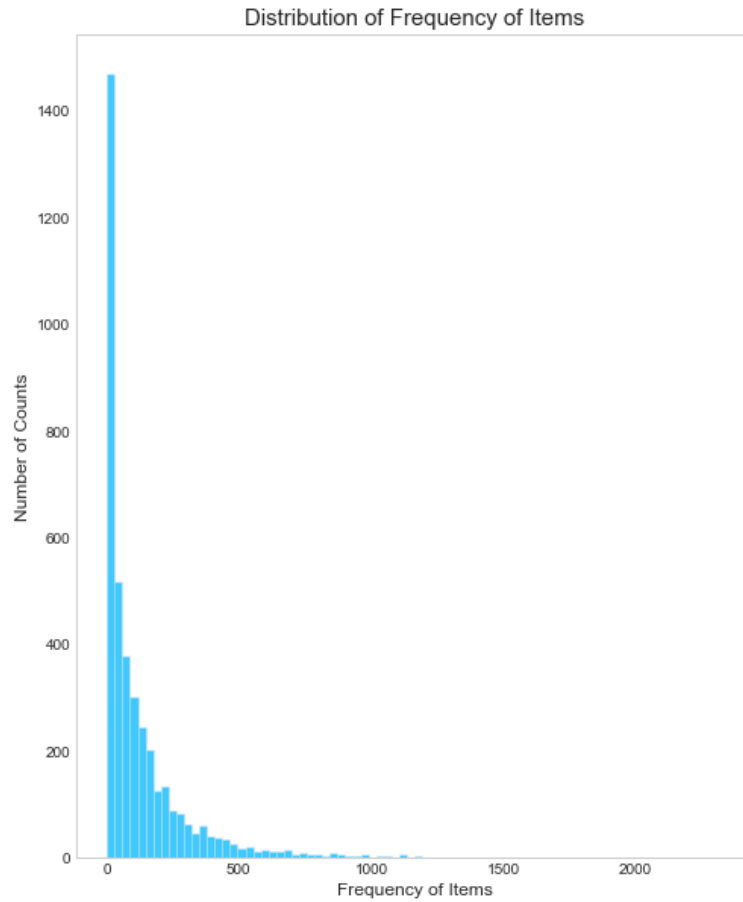


Figure 4: Distribution of Frequency of Items. Approximately 61.5% of items were purchased fewer than 100 times.

Table 2 presents the top five most frequently purchased items in the dataset.

Table 2: Top 5 Most Frequently Purchased Items

Item	Number of Unique Transactions with Item	Percent of Unique Transactions with Item
White Hanging Heart T-Light Holder	2327	11.67%
Jumbo Bag Red Retrospot	2115	10.61%
Regency Cakestand 3 Tier	2019	10.13%
Party Bunting	1707	8.56%
Lunch Bag Red Retrospot	1594	7.99%

Association Rule Model

To create the model, we used the apriori function from the apyori package in Python 3. The only necessary variables were *InvoiceNo* (transactions) and *Description* (purchased items). To avoid extreme examples, we set a minimum support as 0.0045, which means only items or item combinations that appear more than 90 times will be counted in final analysis. The minimum confidence and minimum lift were set as 0.2 and 3 in order to obtain a thorough distribution of potentially important rules. In some cases, multiple rules can be generated from the same combination of items, so only rules with the highest confidence and lift scores are kept. This concept is illustrated in the Appendix 1.1.

Model Performance Evaluation

We used confidence and lift to measure model performance. The distribution of confidence and lift of the results are presented below in Figures 5 and 6.

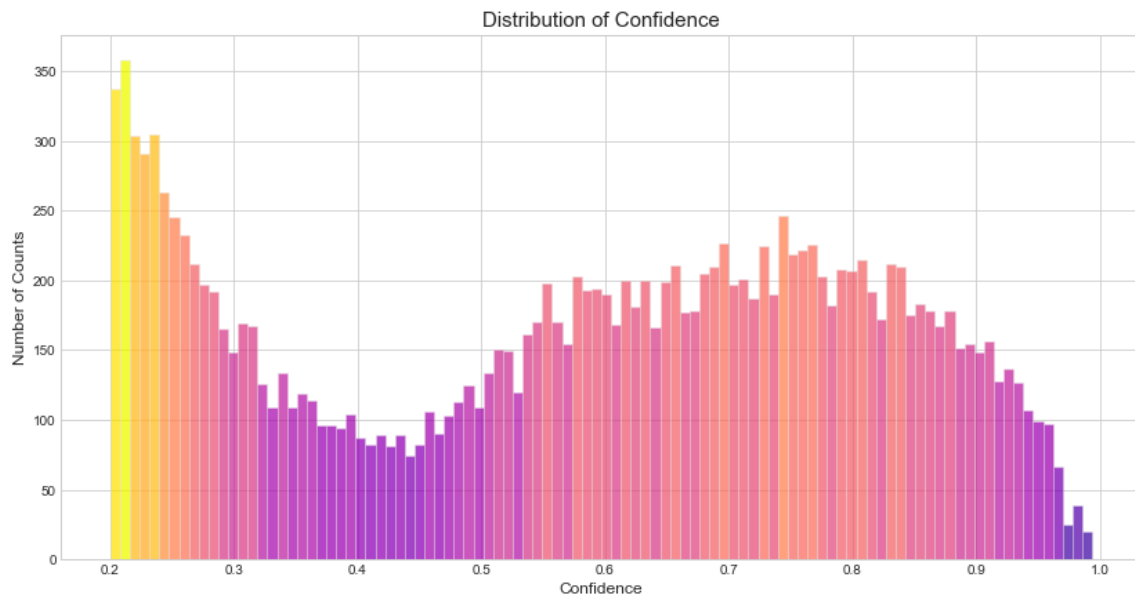


Figure 5: Distribution of Confidence (only rules with confidence larger than 0.2 are shown)

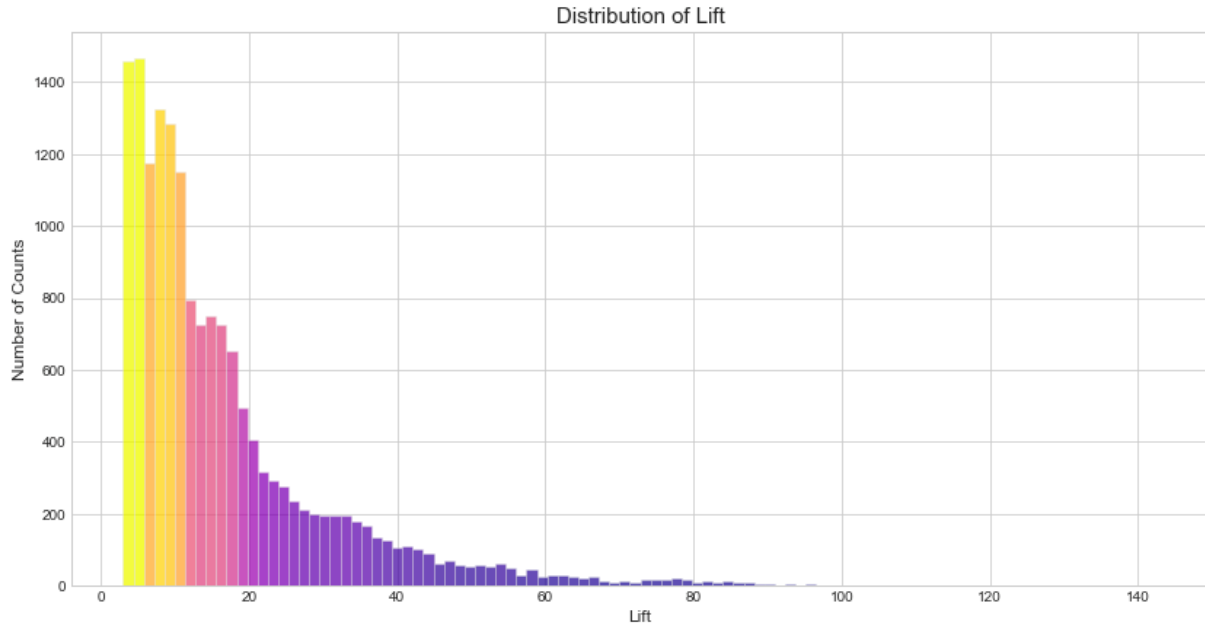


Figure 6: Distribution of Lift (only rules with lift larger than 3 are shown)

The results from Figures 5 and 6 have the following summary statistics:

1. 52% of rules have confidence greater than 0.6 and 7% have a confidence greater than 0.9
2. 59% of rules have lifts greater than 10 and 8% have lifts greater than 40.

Table 3 shows 10 example rules from the model along with their support, confidence, and lift values.

Table 3: Sample Rules from MBA of Online Retail Dataset

	Rule*	Confidence	Lift	Support
1	HERB MARKER BASIL → HERB MARKER THYME, HERB MARKER MINT, HERB MARKER ROSEMARY, HERB MARKER CHIVES, HERB MARKER PARSLEY	99.37%	89	0.79%
2	REGENCY MILK JUG PINK → REGENCY TEA PLATE PINK, REGENCY SUGAR BOWL GREEN, REGENCY TEA PLATE GREEN, REGENCY TEA PLATE ROSES, REGENCY TEAPOT ROSES	98.91%	118	0.46%
3	CARD SUKI BIRTHDAY → SKULL SHOULDER BAG, SUKI SHOULDER BAG, FLORAL FOLK STATIONERY SET	97.89%	91	0.47%
4	HERB MARKER BASIL → HERB MARKER MINT, HERB MARKER THYME	91.26%	76	0.94%
5	DOLLY GIRL CHILDRENS BOWL → SPACEBOY CHILDRENS CUP, SPACEBOY CHILDRENS BOWL, DOLLY GIRL CHILDRENS CUP	98.10%	143	0.52%
6	BISCUITS SMALL BOWL LIGHT BLUE → SMALL DOLLY MIX DESIGN ORANGE BOWL, SMALL CHOCOLATES PINK BOWL, SMALL MARSHMALLOWS PINK BOWL	98.02%	70	0.50%
7	ALARM CLOCK BAKELIKE GREEN, CHARLOTTE BAG PINK POLKADOT → STRAWBERRY CHARLOTTE BAG, CHARLOTTE BAG SUKI DESIGN, RED RETROSPOT CHARLOTTE BAG, PACK OF 72 RETROSPOT CAKE CASES	98.90%	75	0.45%
8	CHRISTMAS TREE DECORATION WITH BELL → CHRISTMAS TREE STAR DECORATION, CHRISTMAS TREE HEART DECORATION	95.19%	113	0.50%
9	RED RETROSPOT SHOPPER BAG → JUMBO BAG PINK POLKADOT, JUMBO BAG RED RETROSPOT	86.39%	8	0.73%
10	CHRISTMAS RETROSPOT ANGEL WOOD → CHRISTMAS RETROSPOT STAR WOOD	51.81%	52	0.50%

From the Table 3's results, rules with high confidence and lift generally fit into several types:

1. Products with different colors → seen in rule 2, where different colors are purchased together.
2. Products with different designs or styles → seen in rule 1, where different spices are purchased together.
3. Products that are part of a set → seen in rules 5 and 8, where a set of items are purchased together.

Rules 1 and 4 show that people prefer to buy all colors or designs together rather than several of them. Rules 9 and 10 are rules with high confidence but low lift and high lift but low confidence, respectively. This happens when the support of consequent is too high or too low. Therefore, it is important to check both confidence and lift when searching for rules.

Model Review

The results of the MBA on this dataset highlight the method's strengths and weaknesses. We did not require numeric data and we could also remove many unnecessary variables, such as *Quantity*, *Unit Price* and *Transaction Date*, while still identifying important relationships. In addition, confidence and lift are evaluation metrics that are convincing and easy to interpret. Many of the rules themselves also show associations that intuitively make sense given consumer purchase habits.

However, this model also highlighted several weaknesses of using MBA. The apriori algorithm uses breadth-wide search and bottom-up exploration to find all possible combinations of items in the dataset, making it very computationally expensive for large datasets.¹⁸ In this example, when we set a support threshold of 0.0045, it took two hours to find all subsets. Other algorithms such as FP-Growth and Eclat have demonstrated good performance in some special datasets but still generally struggled with large datasets.¹⁹ In addition, MBA can result in potentially false rules if two items happen to be popular but are not actually related to each other. Even though we account for this using lift, it is possible that some identified rules may still not truly be accurate. Finally, the results of the MBA are primarily used for inference and cannot be used to make predictions of future purchases in contrast to other ML algorithms such as Decision Trees.

Summary

The Market Basket Analysis is a simple yet powerful algorithm that is especially useful in the process of data mining. It operates on transactions or itemsets, and it aims to create simple if-then rules between itemsets of the form $A \rightarrow B$, where A is the *Antecedent* and B is the *Consequent*.

The algorithm behind the MBA has an intuitive mathematical foundation based on three main concepts:

1. *Support* is the proportion of times an itemset appears in all transactions
2. *Confidence* is the *Support* of the itemset divided by the *Support* of the *Antecedent* i.e. conditional probability
3. *Lift* is the *Confidence* of the itemset divided by the *Support* of the *Consequent*.

Apriori is the algorithm that creates association rules in a computationally feasible way. It generates itemsets, filters them based on a minimum support threshold, generates rules on the resulting itemsets, filters them based on a minimum confidence level, and identifies maximum frequent itemsets.

The advantages of MBAs include their ease of use, their interpretability, and the fact they do not require labeled data. They also have many industrial applications and can be used to create cross-selling strategies, compute churn, or identify patterns in the data. As disadvantages, they might produce false associations with a small dataset, they have a high computational cost with large datasets, and they require categorical variables to work. MBAs also have some limitations in their use and application. Their output is a table with best rules and cannot be used to make predictions. The process of creating rules also needs human input to define appropriate thresholds.

In the future, we would like to further explore the topic of association rule learning by looking at other algorithms such as *FP-growth* or *Eclat*, as well as further comparing results with more traditional supervised and unsupervised learning algorithms.

Roles

1. Anshupriya: Worked on the theoretical aspects of the algorithm. Researched the implementation and wrote the Methods section of the report.
2. Calvin: Worked on examples in the video showcase. Worked on modeling based on example dataset in python. Took charge of the Example section. Researched different algorithms of Association Rules and tested in code. Researched recommendation systems.
3. Guillem: Worked on the video presentation, specifically in developing and improving the visual content of the scenes. In charge of the Summary and Background sections. Reviewed the report, synthesized and made improvements. General desk research and investigation.
4. Jose: Provided cases examples and baseline for video presentation in slide format. In the report, developed content for advantages and constraints of the method, as well as current use cases. General editing support.
5. Varun: Wrote script, provided content, and recorded audio for video showcase. Researched topic and primarily wrote Abstract and Introduction. Helped review and revise report.

References

1. Trevor Hastie, Robert Tibshirani, & Jerome Friedman. (2016). *The Elements of Statistical Learning: Data Mining, Inference, and Prediction* (2nd ed.). Springer.
2. Rajak, A., & Gupta, M. K. (2008, February). *Association rule mining: applications in various areas*. In Proceedings of International Conference on Data Management, Ghaziabad, India (pp. 3-7).
3. Revathi Gopalakrishnan, & Avinash Venkateswarlu. (2018). *Machine Learning for Mobile: Practical guide to building intelligent mobile applications powered by machine learning*. Packt Publishing.
4. Kordík, P. (2019, December 15). Machine Learning for Recommender systems—Part 1 (algorithms, evaluation and cold start). Medium. <https://medium.com/recombee->

blog/machine-learning-for-recommender-systems-part-1-algorithms-evaluation-and-cold-start-6f696683d0ed

5. Recommendation System using Association Rule Mining for Implicit Data. (n.d.). Retrieved April 18, 2020, from <https://medium.com/datadriveninvestor/recommendation-system-using-association-rule-mining-for-implicit-data-6fba0f6c5012>
6. Verma, N., & Singh, J. (2017). A comprehensive review from sequential association computing to Hadoop-MapReduce parallel computing in a retail scenario. *Journal of Management Analytics*, 4(4), 359–392. <https://doi.org/10.1080/23270012.2017.1373261>
7. Zhang, C., Chen, Y., Yang, J., & Yin, Z. (2019). An association rule based approach to reducing visual clutter in parallel sets. *Visual Informatics*, 3(1), 48–57. <https://doi.org/10.1016/j.visinf.2019.03.006>
8. Gupta, N., Mangal, N., Tiwari, K., & Mitra, P. (2006). *Mining Quantitative Association Rules in Protein Sequences*. In G. J. Williams & S. J. Simoff (Eds.), *Data Mining* (Vol. 3755, pp. 273–281). Springer Berlin Heidelberg. https://doi.org/10.1007/11677437_21
9. Serban, G., Czibula, I. G., & Campan, A. L. I. N. A. (2006). *A programming interface for medical diagnosis prediction*. *Studia Universitatis "Babes-Bolyai", Informatica*, LI (1), pp (21-30).
10. Petrovic, N. (2018, November 29). Adopting Data Mining Techniques in Telecommunications Industry: Call Center Case Study.
11. Cecaj, A., Mamei, M. Investigating economic activity concentration patterns of co-agglomerations through association rule mining. *J Ambient Intell Human Comput* 10, 463–476 (2019). <https://doi.org/10.1007/s12652-017-0665-3>
12. Suchacka, G., Chodak, G. Using association rules to assess purchase probability in online stores. *Inf Syst E-Bus Manage* 15, 751–780 (2017). <https://doi.org/10.1007/s10257-016-0329-4>
13. *Associative Rule Mining*. (n.d.). Retrieved April 18, 2020, from https://www.youtube.com/watch?v=_Z5tZuVskaQ&t=332s
14. *Complete guide to Association Rules (1/2)—Towards Data Science*. (n.d.). Retrieved April 18, 2020, from <https://towardsdatascience.com/association-rules-2-aa9a77241654>
15. *Complete guide to Association Rules (2/2)—Towards Data Science*. (n.d.). Retrieved April 18, 2020, from <https://towardsdatascience.com/complete-guide-to-association-rules-2-2-c92072b56c84>
16. Beamer, S., Asanovic, K., & Patterson, D. (n.d.). *Direction-Optimizing Breadth-First Search*. 10.
17. Heaton, J. (2017). Comparing Dataset Characteristics that Favor the Apriori, Eclat or FP-Growth Frequent Itemset Mining Algorithms. *ArXiv:1701.09042 [Cs]*. <http://arxiv.org/abs/1701.09042>

Appendix

Appendix 1.1: Comparing Multiple Association Rules from Same Combination of Items

Concept

Let us consider the 3 itemset {Blue, Pink, Purple} from Figure 1 with support 0.4. This itemset will not generate any useful rules based on our minimum support threshold but is useful to look at to determine calculation of the confidence of association rules. When confidence is calculated, the denominator is the support of the antecedent. The numerator remains the same for an itemset since it is the support of the itemset in the transactions. As a result, when we remove an element from the antecedent and move it to the consequent the denominator value cannot decrease. Hence,

$$\text{Confidence}(\{a, b\} \rightarrow \{c\}) \geq \text{Confidence}(\{a\} \rightarrow \{b, c\})$$

This rule can be used to improve the efficiency of calculating strong association rules since we can ignore the association rules $\{a, b\} \rightarrow \{c\}$ if $\{a\} \rightarrow \{b, c\}$ is less than the minimum confidence threshold. This will ensure that we reduce the number of association rules considered exponentially.

Example: Online Retail Dataset

Consider the following items from the online retail dataset that are all highly associated with each other (each item is denoted as a letter for this example):

1. PINK REGENCY TEACUP AND SAUCER (A)
2. ROSES REGENCY TEACUP AND SAUCER (B)
3. GREEN REGENCY TEACUP AND SAUCER (C)

There are three possible rules:

1. $\{A, B\} \rightarrow \{C\}$
2. $\{B, C\} \rightarrow \{A\}$
3. $\{A, C\} \rightarrow \{B\}$

However, the rule with highest confidence and lift is the following is Rule 1. Therefore, only this rule is considered in the final analysis to explain the association between these items.