# An Algebraic Representation of Frequent Market Baskets and Association Rules

*J. Demetrovics[1], Hua Nam Son[2], Akos Guban[2]*

[1] *MTA SZTAKI, 1111 Budapest, Lágymányosi u. 11*
[2] *Budapest Business School, 1149 Budapest, Buzogány u. 11-13*

**Abstract:** *This study proposes an algebraic approach for formal representation of Market Basket (MB) model. In a more generalized model by taking into consideration the quantity of items in transactions and by using tools of lattice theory we reconsider well-known problems and show an explicit representations of frequent MBs, basic frequent MBs and association rules. As straightforward consequences, the algorithms to find them are presented.*

**Keywords:** *market basket, frequent item, association rule, lattice.*

## 1. Introduction

Great efforts have been made to discover the informations hidden in the customer transactions. The study of customer Market Baskets (MB) and mining the association rules are important in various applications, for example, in decision making and strategy determination of retail economy [1]. In those studies the market baskets (transactions) are often considered as sets of items purchased by customers. Discovering of *large itemsets* and association rules attracts the interest of researchers. One can notice that in these studies the researchers are interested in the set of items (e.g. bread, milk, ...) purchased by customers in the super market, and did not care of the quantity of each item. However, it is interesting also if we know not only that 70% of customers buy bread and milk, but we know also 50% of customers buy 1 kg bread and 2 l milk, while 1% of customers buy 10 kg bread and 1 l milk. Similar example can be found for association rules. The meaning of quantitative analysis of transactions is evident.

In this study we introduce a quantitative analysis of transactions and association rules of transactions. The quantitative analysis may reveal informations hidden in the transactions. We are interested not only in the statement "90% of customers who buy bread and milk also purchase butter", but in the statement "90% of customers who buy 1 kg bread and 2 l milk also purchase 0.5 kg butter". By dealing with the quantity of items our setting is somehow different of those in previous studies (see [1]). That is why instead of *itemsets* (see [1]) we use *market baskets* or *transactions*. The main advantage of this approach is that all transactions can be examined as elements of a lattice with natural partial order. So the lattice-theoretic methods can be applied for transactions examination.

## 2. A generalized setting for Market Basket Model

For a finite set of items $P = \{p_1, p_2, ..., p_n\}$ we consider a MB as a tube $\alpha = (\alpha[1], \alpha[2], ..., \alpha[n])$, where $\alpha[i] \in \aleph$ is the quantity of $p_i$ in the basket $\alpha$. The set of all MBs is denoted by $\Omega$.

For $\alpha, \beta \in \Omega$ where $\alpha = (\alpha[1], \alpha[2], ..., \alpha[n])$, $\beta = (\beta[1], \beta[2], ..., \beta[n])$ we write $\alpha \leq \beta$ if for all $i = 1, 2, ..., n$ we have $\alpha[i] \leq \beta[i]$. $\langle \Omega, \leq \rangle$ is a lattice with the natural partial order $\leq$. For a set $A \subseteq \Omega$ we denote

$$U(A) = \{\alpha \in \Omega \mid \forall \beta \in A : \beta \leq \alpha\},$$
$$L(A) = \{\alpha \in \Omega \mid \forall \beta \in A : \alpha \leq \beta\} .$$

We denote also

$$\sup(A) = \{\alpha \in U(A) \mid \nexists \beta \in U(A): \beta < \alpha\},$$
$$\inf(A) = \{\alpha \in L(A) \mid \nexists \beta \in L(A): \alpha < \beta\}.$$

One should remark that $\sup(A)$ and $\inf(A)$ are single elements of $\Omega$, namely $\sup(A) = u \in \Omega$, where $u[i] = \max\{\alpha[i] \mid \alpha \in A\}$ and $\inf(A) = v \in \Omega$, where $v[i] = \min\{\alpha[i] \mid \alpha \in A\}$.

For a set $A \subseteq \Omega$ and $\alpha \in \Omega$ we denote by

$$\mathrm{supp}_A(\alpha) = \frac{|\{\beta \in A \mid \alpha \leq \beta\}|}{|A|}$$

the support of $\alpha$ in $A$. In word, $\mathrm{supp}_A(\alpha)$ denotes the rate of all market baskets that exceeds the given threshold $\alpha$ (in the form of a sample market basket) to the whole $A$. The support of an market basket is a statistical index and naturally, the market baskets of more support are of more significance and attract the attention of the managers, as well as of the researchers.

One can notice that an item $p_i$ (discussed in other studies, see, for example, [1]) in our study should be identified with $U(\alpha_i)$, where $\alpha_i = (\alpha[1], \alpha[2], ..., \alpha[n])$, $\alpha[k] = 0$ if $k \neq i$ and $\alpha[i] = 1$. We should not confuse $p_i$ with $\alpha_i$.

For $\alpha, \beta \in \Omega$ where $\alpha = (\alpha[1], \alpha[2], ..., \alpha[n])$ and $\beta = (\beta[1], \beta[2], ..., \beta[n])$ we write $\gamma = \alpha \cup \beta$ if $\gamma[i] = \max\{\alpha[i], \beta[i]\}$ for all $i = 1, 2, ..., n$. We call

$\alpha \to \beta$ an *association rule* of $\beta$ to $\alpha$. By the *confidence* of $\alpha \to \beta$ in a set of MBs $A$ we understand the rate

$$\mathrm{conf}_A(\alpha \to \beta) = \frac{\mathrm{supp}_A(\alpha \cup \beta)}{\mathrm{supp}_A(\alpha)}$$

As remarked in [1] the support of MBs is a kind of statistical index, while the confidence of association rules is a measure of their "strength".

## 3. Frequent Market Baskets

For a set $A \subseteq \Omega$, $\alpha \in \Omega$ and $0 \le \varepsilon \le 1$ we say that $\alpha$ is $\varepsilon$-frequent MB, if $\mathrm{supp}_A(\alpha) \ge \varepsilon$. The set of all $\varepsilon$-frequent MBs is denoted by $\Phi_A^\varepsilon$. We have the following

**Apriori Principle.** For a set $A \subseteq \Omega$, $\alpha, \beta \in \Omega$ and $0 \le \varepsilon \le 1$, if $\alpha \le \beta$ and $\beta$ is $\varepsilon$-frequent then $\alpha$ is $\varepsilon$-frequent.

**Example 1.** Consider a set of items $P = \{a, b, c\}$ and a set of transactions $A = \{\alpha, \beta, \gamma, \delta\}$, where $\alpha = (2, 1, 0)$, $\beta = (1, 1, 1)$, $\gamma = (1, 0, 1)$, $\delta = (2, 2, 0)$. One can see that for $\sigma = (1, 1, 0)$, $\eta = (1, 2, 0)$ we have $\mathrm{supp}_A(\sigma) = \frac{3}{4}$ and $\mathrm{supp}_A(\eta) = \frac{1}{4}$. For the threshold $\varepsilon = \frac{1}{2}$ the $\varepsilon$-frequent MBs of $A$ are:

$$\Phi_A^{\frac{1}{2}} = \{(2, 1, 0), (1, 0, 1), (1, 1, 0), (2, 0, 0), (0, 0, 1), (0, 1, 0), (1, 0, 0), (0, 0, 0)\}.$$

Let us denote

$$\Phi_{A,k} = \{\alpha \in \Omega \mid \exists \alpha_1, \alpha_2, ..., \alpha_k \in A : \alpha \le \{\alpha_1, \alpha_2, ..., \alpha_k\}\}.$$

One can remark that if $k \le l$ then $\Phi_{A,k} \supseteq \Phi_{A,l}$ and $\Phi_A^\varepsilon = \Phi_{A,k}$, where $k = \lceil \varepsilon \mid A \mid \rceil$ denotes the smallest integer that is greater or equal to $\varepsilon \mid A \mid$.

We have the following

**Theorem 1.** For a set of items $P = \{p_1, p_2, ..., p_n\}$, a set of MBs $A \subseteq \Omega$ and a threshold $0 \le \varepsilon \le 1$ an MB $\alpha \in \Omega$ is $\varepsilon$-frequent iff there exist $\alpha_1, \alpha_2, ..., \alpha_k \in A$ such that $\alpha \in L(\{\alpha_1, \alpha_2, ..., \alpha_k\})$, where $k = \lceil \varepsilon \mid A \mid \rceil$.

*Proof:* If there exist $\alpha_1, \alpha_2, ..., \alpha_k \in A$, $k = \lceil \varepsilon \mid A \mid \rceil$, such that $\alpha \in L(\{\alpha_1, \alpha_2, ..., \alpha_k\})$ then $\alpha \le \alpha_i$ for all $i = 1, 2, ..., k$, i.e.,

$$\mathrm{supp}_A(\alpha) = \frac{\mid \{\beta \in A \mid \alpha \le \beta\} \mid}{\mid A \mid} \ge \frac{k}{\mid A \mid} \ge \varepsilon.$$

*Vice versa*, if $\mathrm{supp}_A \ge \varepsilon$ then $\mid \{\beta \in A \mid \alpha \le \beta\} \mid \ge \varepsilon. \mid A \mid$, i.e. there exist $\alpha_1, \alpha_2, ..., \alpha_k \in A$, $k = \lceil \varepsilon \mid A \mid \rceil$, such that $\alpha \in L(\{\alpha_1, \alpha_2, ..., \alpha_k\})$. The proof is completed.

26

By the Theorem 1 we have the following

**Algorithm 1 (Creating all $\varepsilon$-frequent MBs of a given set of transactions $A$).**

*Input.* Set of items $P$, set of MBs $A \subseteq \Omega$ and a threshold $0 \leq \varepsilon \leq 1$.

*Output.* $\Phi_A^\varepsilon$.

    *Step 1.* $\Phi_A^\varepsilon := \varnothing$.

    *Step 2.* $k = \lceil \varepsilon \, | A | \rceil$.

      For all $B \subseteq A$, $| B |= k$

      $\Phi_A^\varepsilon := \Phi_A^\varepsilon \cup L(B)$

      EndFor;

    End

Let $\quad | P |= n$, $\quad k = \lceil \varepsilon \, | A | \rceil$, $\quad m = max\{\alpha[i] \,|\, \alpha \in A, i = 1,2,...,n\}$. The algorithm requires $O\left(\binom{|A|}{k}(m+1)^n\right)$ running time.

As a consequence of the previous theorem we have the following

**Theorem 2 (Explicit representation of *large MBs*).** For a set of items $P = \{p_1, p_2, ..., p_n\}$ a set of MBs $A \subseteq \Omega$ and a threshold $0 \leq \varepsilon \leq 1$ there exist $\alpha_1, \alpha_2, ..., \alpha_s \in \Omega$, where $s = \binom{|A|}{\lceil \varepsilon |A| \rceil}$ such that

$$\Phi_A^\varepsilon = \bigcup_{i=1}^{s} L(\alpha_i).$$

*P r o o f:* Let $\alpha_1, \alpha_2, ..., \alpha_s$ be the set of all $\inf\{\beta_1, \beta_2, ..., \beta_k\}$ where $k = \lceil \varepsilon \, | A | \rceil$ and $\beta_i \in A$. By Theorem 1 we have

$$\alpha \in \Phi_A^\varepsilon \Leftrightarrow \alpha \leq \inf(\{\beta_1, \beta_2, ..., \beta_k\})$$

for some $\{\beta_1, \beta_2, ..., \beta_k\} \subseteq A$, where $k = \lceil \varepsilon \, | A | \rceil$. This implies that $\Phi_A^\varepsilon = \bigcup_{i=1}^{s} L(\alpha_i)$. The proof is completed.

We should remark that $\alpha_i \leq \alpha_j$ iff $L(\alpha_i) \subseteq L(\alpha_j)$. For a set of MBs $A$ and a given threshold $\varepsilon$ the set of MBs $\alpha_1$, $\alpha_2$, ..., $\alpha_s$ for which

(i)   $\Phi_A^\varepsilon = \bigcup_{i=1}^{s} L(\alpha_i)$,

(ii)  $\forall i, j : 0 \leq i, j \leq s$ we have $\alpha_i \nleq \alpha_j$ and $\alpha_j \nleq \alpha_i$

is called by *basic $\varepsilon$-frequent set of MBs* of $A$. It is easy to verify that for a given $A$, $\varepsilon$ the basic $\varepsilon$-frequent set of MBs of $A$ is unique, which we denote by $S_A^\varepsilon$. Since the determination of $\Phi_A^\varepsilon$ (the set of all $\varepsilon$-frequent set of MBs in $A$) is important, it is interesting to determine its basic $\varepsilon$-frequent set of MBs $S_A^\varepsilon$. We have the following

**Theorem 3.** For a set of items $P$, a threshold $0 \leq \varepsilon \leq 1$ every set of MBs $A \subseteq \Omega$ has an unique basic $\varepsilon$-frequent set of MBs $S_A^\varepsilon$.

The simple proof is omitted. The following algorithm creates the unique basic $\varepsilon$-frequent set of MBs for a given set of MBs $A \subseteq \Omega$ and a given threshold $\varepsilon$:

**Algorithm 2 (Creating the basic $\varepsilon$-frequent set of MBs $S_A^\varepsilon$).**

*Input.* Set of items $P$, Set of MBs $A \subseteq \Omega$ and a thershold $0 \leq \varepsilon \leq 1$.

*Output.* $S_A^\varepsilon$.

   *Step 1.* $S_A^\varepsilon := \varnothing$.

   *Step 2.* $k = \lceil \varepsilon \, | A | \rceil$.

     For $B \subseteq A$, $| B |= k$

       For $\alpha \in S_A^\varepsilon$

         If $\alpha \leq \inf(B)$ or $\inf(B) \leq \alpha$ then

         $S_A^\varepsilon := S_A^\varepsilon \setminus \{\min(\alpha, \inf(B))\} \cup \{\max(\alpha, \inf(B))\}$.

         else

         $S_A^\varepsilon := S_A^\varepsilon \cup \{\inf(B)\}$.

         EndIf

       EndFor

     EndFor

    End

For $| P |= n$, $k = \lceil \varepsilon \, | A | \rceil$, $m = \max\{\alpha[i] \, | \, i = 1, 2, ..., n; \alpha \in A\}$ one can see that $| S_A^\varepsilon | \leq \binom{|A|}{k}$. Therefore the algorithm requires $O\left(\binom{|A|}{k} mn\right)$ running time. One can remark also than in the case of large amount of transactions $A$ the basic $\varepsilon$-frequent set of MBs $S_A^\varepsilon$ can be generated much more quickly than the set of all $\varepsilon$-frequent set of MBs $\Phi_A^\varepsilon$.

**Example 2.** We continue the Example 1. For the set of transactions $A$ Algorithm 2 generates the basic $\dfrac{1}{2}$-frequent set of MBs $S_A^{\frac{1}{2}} = \{\rho, \theta\}$, where $\rho = (2, 1, 0)$, $\theta = (1, 0, 1)$. It means that the family of $\dfrac{1}{2}$-frequent set of MBs of $A$ is $\Phi_A^{\frac{1}{2}} = L(\rho) \cup L(\theta)$.

## 4. Association and confidence

In our generalized model of market baskets we can find all associations with given confidence. For a set of items $P$, a set of MBs $A \subseteq \Omega$ and a threshold $0 \leq \varepsilon \leq 1$ an

association $\alpha \to \beta$ is $\varepsilon$-*confident* if $\mathrm{conf}_A(\alpha \to \beta) \geq \varepsilon$. The set of all $\varepsilon$-confident associations of $A$ is denoted by $C_A^\varepsilon$. We have the following

**Theorem 4.** For a set of products $P$ a set of MBs $A \subseteq \Omega$ and $0 \leq \varepsilon \leq 1$ an association $\alpha \to \beta$ is $\varepsilon$-confident iff $\dfrac{|U(\alpha \cup \beta) \cap A|}{|U(\alpha) \cap A|} \geq \varepsilon$.

*Proof:* Remark that $\mathrm{supp}_A(\alpha \cup \beta) = \dfrac{|U(\alpha \cup \beta) \cap A|}{|A|}$ and

$\mathrm{supp}_A(\alpha) = \dfrac{|U(\alpha) \cap A|}{|A|}$. With these remarks the proof of the theorem is straightforward.

A natural question for cross marketing, store layout, ...(see, for example [1]) is to find all association rules with a given confidence. In our generalized model the following theorem shows in a sense an explicit representation of all association rules. More exactly, we show for a given MB $\alpha$ which set of MBs $\beta$ may be associated to $\alpha$ with a given threshold of confidence.

For MBs $\rho$, $\sigma$ where $\rho \leq \sigma$, let us denote
$$M(\rho, \sigma) = \{\eta \in \Omega \mid \rho \cup \eta \leq \sigma\}.$$

It should be remarked that $M(\rho, \sigma)$ can be represented explicitly. If $\rho = (\rho_1, \rho_2, ..., \rho_s)$, $\sigma = (\sigma_1, \sigma_2, ..., \sigma_s)$ then $\eta = (\eta_1, \eta_2, ..., \eta_s) \in M(\rho, \sigma)$ iff $\max(\rho_i, \eta_i)$ for all $i = 1, 2, ..., s$, i.e., $\eta_i = \sigma_i$ in the case $\rho_i \not\leq \sigma_i$ and $\eta_i \leq \sigma_i$ in the case $\rho_i = \sigma_i$.

**Theorem 5 (Explicit representation of association rules).** For a set of items $P = \{p_1, p_2, ..., p_n\}$, a set of MBs $A \subseteq \Omega$, an MB $\alpha \in \Omega$ and a threshold $0 \leq \varepsilon \leq 1$ there exist $\alpha_1, \alpha_2, ..., \alpha_k \in \Omega$ such that $\forall \beta \in \Omega : \alpha \to \beta$ is $\varepsilon$-confident association rule iff $\beta \in \bigcup_{i=1}^k M(\alpha, \alpha_i)$.

*Proof:* Put $s = \lceil \varepsilon |U(\alpha) \cap A| \rceil$ by Theorem 4 we have that $\alpha \to \beta$ is $\varepsilon$-confident association rule iff $|U(\alpha \cup \beta) \cap A| \geq s$. Let $\alpha_i$ denotes $\inf(B)$, where $B \subseteq A$, $|B| \geq s$. One can verify that $|U(\alpha \cup \beta) \cap A| \geq s$ iff $\beta \in M(\alpha, \alpha_i)$. The proof is completed.

Theorem 5 in a sense gives an explicit presentation for association rules. As a straightforward consequence, we have an algorithm to find all $\varepsilon$-confident association rules for given left side.

**Algorithm 3 (Creating all $\varepsilon$-confident association rules $\alpha \to \beta$ for given $\alpha$).**

*Input.* A set of items $P$, a set of MBs $A \subseteq \Omega$, a thershold $0 \leq \varepsilon \leq 1$ and an MB $\alpha$

*Output.* $\bigcup_{i=1}^k M(\alpha, \alpha_i)$.

*Step 1.* $B := U(\alpha) \cap A = \{\gamma \in A \mid \alpha \leq \gamma\}$.

*Step 2.* $s := \lceil \varepsilon \mid B \mid \rceil$.

  $k := \mid \{C \subseteq B \mid\mid C \mid \geq s\} \mid$

  For $C \subseteq B$, $\mid C \mid \geq s$, calculate $\alpha_i = \inf(C)$, $i = 1, 2, ..., k$.

  EndFor

*Step 3.*

  For $i = 1, 2, ..., k$ calculate $M(\alpha, \alpha_i)$

  EndFor

*Step 4.*

  Output $\bigcup_{i=1}^{k} M(\alpha, \alpha_i)$.

 End

**Example 3.** We continue the Example 1. For the set of MBs $A$ (see Example 1), the MB $\sigma = (1, 1, 0)$ and threshold $\varepsilon = \dfrac{1}{2}$ we should find all MB $\eta$ such that $\sigma \to \eta$ is $\varepsilon$-confident association rule. We can see $U(\sigma) \cap A = \{(2, 1, 0), (1, 1, 1), (2, 2, 0)\}$ and $s := \lceil \varepsilon \mid U(\alpha) \cap A \mid \rceil = 2$. By Step 2 in Algorithm 3 we have $k = 4$ and $\alpha_1 = (1, 1, 0)$, $\alpha_2 = (2, 1, 0)$. The set of all MBs $\eta$ such that $\sigma \to \eta$ is $\dfrac{1}{2}$ - confident association rule is

  $M(\sigma, \alpha_1) \cup M(\sigma, \alpha_2) = \{(1,1,0), (1,0,0), (0,1,0), (0,0,0), (2,1,0), (2,0,0)\}.$

  As a result we see that besides the trivial association rules of the form $\sigma \to \sigma'$, where $\sigma' \leq \sigma$ we got non-trivial association rules $\sigma \to (2, 1, 0)$ and $\sigma \to (2, 0, 0)$. In words, among those customers $A$ the ratio of customers who buy $a$ and $b$ also buy two $a$ and one $b$ items, as well the ratio of those who buy $a$ and $b$ also buy two $a$ items, are more than 50%.

## 5. Conclusion

In this study we have proposed an algebraic approach to consider the MB model. The well-known problems are analysed in new, more generalized setting. An explicit representation of frequent set of MBs, as well of association rules are presented. We define the set of basic frequent MBs which determines the set of frequent MBs and can be created in shorter time. We described algorithms that produces the set of frequent MBs and the set of basic frequent MBs. We described also an algorithm that produces the set of association rules for a given left side. The algebraic approach we propose here brings about a clearer representation of well-known results and appears to be a good tool for future study in market basket model.

# References

1 A g r a w a l, R., R. S r i k a n t. Fast Algorithms for Mining Association Rules. VLDB, 1994, 487-499.

2. B r ü g g e r m a n n, T., P. H e d s t r õ m, M. J o s e f s s o n,  Data Mining and Data Based Direct Marketing Activities, Book on Demand GmbH. Norderstedt, Germany, 2004.

3. H a n, J.,  M. K a m b e r. Data Mining: Concepts and Techniques. Second Edition. Morgan Kaufmann Publ., 2006.

4. M a n n i l a, H., H. T o i v o n e n. Discovering Generalized Episodes Using Minimal Occurrences. – In: Proceedings of the Second International Conference on Knowledge Discovery and Data Mining (KDD' 96). August 1996, AAAI Press, 146-151.

5. T o i v o n e n, H.  Sampling Large Databases for Association Rules. Morgan Kaufmann Publ., 1996, 134-145.

6. P a s q u i e r, N., Y. B a s t i d e, R. T a o u i l, L. L a k h a l.  Discovering Frequent Closed Itemsets for Association Rules – ICDT, 1999, 398-416.

7. H s u,  P i n g-Y u,  Y e n-L i a n g  C h e n,  C h u n-C h i n g  L i n g.   Algorithms for Mining Association Rules in Bag Databases. – Information Sciences, Vol. **166**, 2004, Issues 1-4, 31-47.