

# Investigation in Customer Value Segmentation Quality under Different Preprocessing Types of RFM Attributes

<https://doi.org/10.3991/ijes.v4i4.6532>

Nesma mahmoud Taher<sup>1</sup>, Doaa Elzanfaly<sup>1,2</sup>, Shaimaa Salama<sup>1</sup>

<sup>1</sup>Helwan University, Cairo, Egypt

<sup>2</sup>British University in Egypt, Cairo, Egypt

**Abstract**—Customer value segmentation helps retailers to understand different types of customers, develops long term relationship with them, and hence increases their value and loyalty. This study aims to evaluate the quality of customer value segmentation based on two methods of preprocessing the RFM attributes. K-means clustering algorithm is used for the customer value segmentation based on the scored RFM and the actual value of RFM. The quality of the clustering results is tested using the Sum of Squared Error (SSE). Results obtained show that using the actual value of RFM in customer segmentation reduces the clustering error (SSE) and enhances the accuracy of segmentation than using the scored RFM.

**Index Terms**—Customer Relationship Management (CRM), K-means algorithm, RFM Model (Recency, Frequency, Monetary), Sum of Squared Error (SSE).

## I. INTRODUCTION

Customers are one of the most important asset in an enterprise, business cannot success and get profit without satisfied and loyal customers, and without developing their relationship with the organization [1].

Enterprises adopt the customer relationship management (CRM) strategy that is concerned with managing the relationship with customer to learn more about customers' needs, behaviors and values to develop stronger relationship with them [2].

It's easy for small business to notice their customers' preferences and learn from past interactions to serve them better next time. However, with the huge expand in business and the rapid advance in data collection and storage technology, large enterprises found it difficult to understand customer needs and behavior specially that their employees may not interact personally with customers. That's way data mining tools are used to analyze the enterprise data of transactions and customers, and extracting useful patterns that improve CRM through improving the enterprise's understanding of its customers [3].

One significant area in today's competitive business is to analyze the value of customers through customer value segmentation, so that the enterprise can develop personalized relationship with its customers and maximize the current and the future value of them.

The RFM model, which stands for Recency, Frequency, and Monetary, is the most commonly used model for customer value segmentation. There are two methods of preprocessing RFM attributes in the literature, this paper

discusses these two methods trying to find out the better to process the RFM input attributes for more customer value segmentation quality.

In this paper, the K-means algorithm is used for customer segmentation based on two different methods of preprocessing the RFM attributes. First, segmentation will be executed using the actual value of RFM attributes. Second, the segmentation will be done using the scoring quantitative scale of RFM attributes. The segmentation quality when using these two types of RFM preprocessing will be tested using the Sum of Squared Error measure (SSE). Weka data mining software has been employed to conduct the experiments.

The rest of this paper is organized as follows. Section II presents overview of related topics. Section III introduces the dataset of this study. The research model is discussed in section IV, Finally, Section V concludes the paper and suggests some directions for future work.

## II. OVERVIEW OF RELATED TOPICS

This section provides an overview of related topics concerning CRM, Kmeans, SSE, customer value analysis and RFM model.

### A. Customer Relationship Management (CRM)

Reference [5] defined the customer relationship management CRM as an “enterprise approach to understand and influence customers' behavior through meaningful communications in order to improve customer acquisition, customer retention, customer loyalty, and customer profitability.

CRM aims to develop long-term relationship with customers through keeping their satisfaction and increasing the amount of business done with them [6.7].

CRM helps enterprises to understand customer needs, preference and values and effectively developing marketing strategies that lead to long term customer relationship and value [4].

### B. K-means

One of the descriptive data mining task is clustering analysis. Clustering is the process of grouping set of data objects in to groups or clusters, where objects within a cluster are similar to each other, and dissimilar to objects in other clusters [8.10].

There are many clustering methods represented in partitioning methods, hierarchical methods, density-based

methods, grid-based methods, and model based methods [8].

K-means is a well-known partitioning clustering algorithm that works as follows:

Step 1: Partitioning the data set (m objects) into K initial clusters.

Step 2: Assigning an item to the cluster whose centroid is nearest. Distance is computed by using Euclidean distance, and re-calculate the centroid for the cluster receiving the new item or for the cluster losing the item.

Step 3: Repeating Step 2 until no more reassigning.

C. Sum of Squared Error (SSE)

The good clustering result is composed of dense concentrations of data points around their cluster center. The large dispersion of data points from their cluster center indicate non homogeneous clustering [8]. There are a number of measures in the literature that can be used to calculate the level of internal cohesion of the resulted clusters, the Sum of Square Error (SSE) is one of them.

The Sum of Square Error (SSE) evaluates the accuracy of the data segmentation to ensure compact clusters with little deviation from the cluster centroids. As provided in [8], the SSE can be calculated as following equation:

$$SSE = \sum_{i=1}^K \sum_{x \in c_i} dist^2(m_i, x) \tag{1}$$

Where K is the number of clusters,  $m_i$  is the centroid of cluster  $c_i$ , and  $x$  is a data point of cluster  $c_i$ . The clustering result with small SSE is preferred.

D. CLV and RFM Model

Customer value is an important metric to measure customer asset, a concept which is at the heart of CRM [6]. Customer value is found in literature in different names such as Customer lifetime value “CLV”, “customer equity” or “customer profitability”.

RFM (Recency, Frequency, and Monetary) model, which was proposed by [11], is the most widely used model for CLV analysis. This model is based on the most common marketing axiom, the Pareto principle, which stating that “80% of your business comes from 20% of your customers”. RFM attributes can be defined as follows:

- R (Recency of the last purchase): refers to the interval between the latest purchasing date and the end date of a particular period.
- F (Frequency of purchasing): refers to the number of purchases in a particular period.
- M (Monetary value of purchases): refers to the total amount of money spent during a particular period.

The Basic method used in RFM model for customer behavior analysis based on [11] is determined as follows:

The RFM attributes have equal weight (i.e. 1:1:1).

(2) Define the scaling of RFM attributes, which are 5, 4, 3, 2 and 1 that refer to the customer contributions to revenue for enterprises. The ‘5’ refers to the highest customer value and ‘1’ refers to the least value.

(3) Sorting the customers’ data of three RFM attributes by descendant order. The top 20% of customers are assigned the value of 5, the value of 4 is given to the next 20% of customers and so on.

(4) Finally there at the most, 125 different scores (5x5x5) can be assigned, range from 555 (the highest), to 111 (the lowest)

Many researchers have employed the RFM model for customer value segmentation and analysis. References [12, 13, 14, 15, and 16] have used the k-means clustering algorithm and the Scored value of RFM attributes based on [11] for segmenting customers in to different customer value groups.

While the previous studies have used the quantitative numeric value of RFM attributes, reference [17] used the actual value of RFM attribute to segment customers of online retail business into various meaningful groups using the K-means clustering algorithm.

Some authors used the Weighted RFM (WRFM) instead of RFM in customer value analysis. WRFM supposed by [18] where different weights are assigned to R, F, and M depends on characteristics of the industry.

References [19, 20, and 21] have used the actual value of RFM as input to the K-means algorithm for customer segmentation, and then use the weights assigned to RFM attributes to calculate the value of each customer segment.

Some researches try to enhance the RFM model by adding some other parameters. Reference [22, 23] add “L” to RFM model, “L” refers to the interval between the first and last exchange with a customer. LRFM model is used based on the assumption that the longer a customer relationship, the higher the loyalty. Both of the two previous studies have used the actual value of LRFM model in customer segmentation.

Reference [24] uses RFM, WRFM, and LRFM to understand which method is suitable for the case study of Small and medium sized enterprises SMEs, the study uses the score value of those attributes in their analysis. Table I summarizes the previous mentioned methods based on their data preprocessing technique.

TABLE I.  
TYPE OF DATA PREPROCESSING ACCORDING TO RFM METHOD

Technique	Type of Data preprocessing	Reference	Case study
RFM	Scored RFM	[12] [13] [14] [15] [16]	Electronic company online sports store - grocery store commercial store
	Actual value of RFM	[17]	online retail business
WRFM	Actual value of RFM	[19] [20] [21]	hardware retailing company health and Beauty Company electronic flower retailing
LRFM	Actual value of LRFM	[22] [23]	textile business B2B
	Scored RFM	[24]	SME Company

III. DATA UNDERSTANDING

The data set used in this study is customer transaction data of an enterprise that sells office supplements products such as book cases, tables, office furnishing, chairs, copiers and fax, telephones, papers, pens, etc. the dataset is characterized by the following 13 attribute (Order ID, Customer ID, Customer Name, Order Date, Order Quantity, Sales amount, Discount, Unit Price, Shipping Cost, Product Category, Product Name, Ship Date). Table II shows part of the used dataset.

Our data set contains four years of transactions from 1/1/2009 to 31/12/2012 stored in excel sheet. Three time frames of transaction data are used to evaluate the quality of customer segmentation and verify results. The first time frame is 12 month of 2010 contains 600 customers with total 1747 transactions, second is 9 months of 2009 contains 550 customers and 1640 transactions, and 6 months of 2012 contains 1027 transaction generated from 435 customers.

TABLE II.  
PART OF THE USED DATA SET

Order ID	Customer ID	Customer Name	Order Date	Order Quantity	Sales amount	Discount	Unit Price	Shipping Cost	Product Category	Product Name	Ship Date
37218	21	Frank Merwin	12/19/2010	23	153.02	0.08	6.48	5.82	Office Supplies	Xerox 1998	12/20/2010
50307	40	Bill Overfelt	10/15/2010	31	542.01	0.06	17.48	1.99	Technology	Maxell Pro 80 CD-R, 10/Pack	10/15/2010
57314	51	Patrick O'Brill	10/12/2010	26	2238.53	0	85.99	2.79	Technology	6340	10/19/2010
27232	60	Ralph Arnett	8/18/2010	19	671.59	0.05	34.99	7.73	Office Supplies	Hunt Boston® Vacuum Mount	8/20/2010
22820	77	Craig Yedwab	8/17/2010	18	144.84	0	7.89	2.82	Office Supplies	Coated Paper Clips, 800/Box	8/20/2010

IV. RESEARCH MODEL

This section explains the research model and the proposed procedures for evaluating the quality of customer value segmentation using the actual and the score RFM.

First, the transaction data is prepared to get the RFM attributes, then RFM attributes are transformed based on two methods found in literature into actual and score value.

K-means clustering algorithm is used for the customer value segmentation based on the actual and the score RFM. The quality of resulted clusters is tested using the sum of squared error (SSE) to examine which method of preprocessing the RFM attributes lead to more accurate clustering results.

The research procedures as shown Figure 1 contain the following phases:

- A. Data Preprocessing
- B. Get the Actual and Scored RFM
- C. Customer Segmentation Using K-means
- D. Results Evaluation

A. Data Preprocessing

Data preprocessing is the most important and time consuming phase in any data mining project, it affects the accuracy and quality of the subsequent data mining algorithm.

The following four steps in the data preprocessing phase are applied on the data of the three selected time frames of this study:

1. Delete records contain missing values.
2. Extract transaction records corresponding to the tree times frames mentioned above from the basic dataset into other separated excel sheets.
3. Select attributes that are related to RFM model and will be used by the clustering algorithm. These data is the date of transactions and the sales amount.

4. Transform the data set with the selected attributes in to three attributes which are Recency R: most recent transaction date of each customer, frequency F: the count of purchases done by each customer within the specified time frame, and monetary M: the total amount of sales by a customer over the whole time frame of the study.

After preparing the data, the next step is to get the score and the actual value of RFM attributes.

B. Get the Actual and Scored RFM

In this phase, two methods of transforming the RFM attributes are used. The preprocessed data from the previous

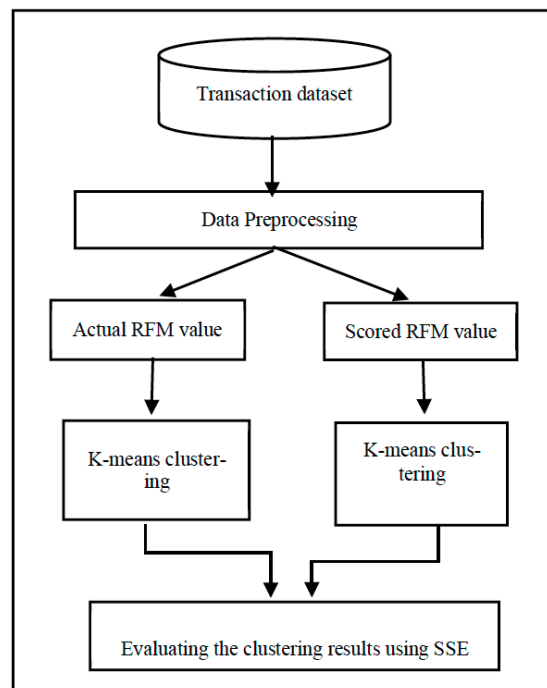


Figure 1. Actual RFM and Scored RFM Model

phase is transformed in to two forms of data, the first is the scored RFM and the second is the actual RFM as follows:

### 1. The Scored RFM

The first method of RFM transformation, is to transform the RFM attributes to quantitative scores. Four steps are followed to get the quantitative scored value of RFM attributes to be used by the clustering algorithm as follows:

1. Sort customers data based on R attribute in descending order from recent to oldest.
2. Split customers into five equal quartiles, and assign score 5 to the first 20% of customers, score 4 for the next 20%, and so on.
3. Repeat step 1 and 2 for F and M attributes by sorting customers based on F and M in descending order and assign scores.
4. Ordering F is in each Recency category, and ordering M is in each combination of Recency and Frequency categories
5. The transformed scored RFM value is stored in excel sheet and then the file is converted into format accepted by Weka [25].

Table III shows an example of Scored RFM values of customers.

### 2. Actual Value of RFM

Unlike the score RFM, this method aims to get the actual value of RFM attributes to be included in to the clustering algorithm as follow:

1. The actual value of Frequency F: is the count of purchases done by each customer within the specified time frame, and the actual value of Monetary M: is the total amount of sales by a customer over the time frame of the study.
2. The R attribute is in date format so it is need to be transformed into number format to be used with the clustering algorithm. So, the R attribute is transformed to number of days since last purchasing by subtracting the last date of the time frame from the last date of customer's purchase.
3. The customer data is sorted in ascending order based on the R at first, then F is sorted in descending order in each Recency category, and finally, M is sorted in descending order in each combination of Recency and Frequency categories.
4. The transformed actual RFM value is stored in excel sheet and then the file is converted into format accepted by Weka [25].

An example of the actual RFM value is provided in Table IV.

### C. Customer Segmentation Using K-means

After data preprocessing and transformation, the actual and scored value of RFM attributes is used for the customer segmentation. K-means clustering algorithm in Weka data mining software is employed in this study.

The clustering algorithm will run two times, first using the actual RFM and second with Scored RFM. The data set of the time frame 12 month of 2010 from 1/1/2010 to 31/12/2010 is firstly used, and then the same process is applied using the other two time frames considered in this study in order to verify results.

TABLE III.  
SCORED RFM VALUES OF CUSTOMERS

Customer ID	Customer Name	R	F	M
80	Evan Minnotte	9	5	13553
81	Alejandro Ballentine	216	5	13221
82	Hilary Holden	233	6	13106
83	Ruben Dartt	155	2	13056
84	Dave Poirier	61	8	13005
85	Joy Smith	29	6	12590
86	Bobby Elias	82	4	12549

TABLE IV.  
ACTUAL RFM VALUES OF CUSTOMERS

Customer ID	Customer Name	R	F	M
80	Evan Minnotte	5	5	5
81	Alejandro Ballentine	1	5	5
82	Hilary Holden	1	5	5
83	Ruben Dartt	2	2	5
84	Dave Poirier	4	5	5
85	Joy Smith	5	5	5
86	Bobby Elias	3	4	5

Customers are segmented into eight number of clusters (K=8) because eight ( $2 \times 2 \times 2$ ) possible combinations of inputs RFM can be obtained by assigning  $\uparrow$  or  $\downarrow$ . If the average R (F, M) value of a cluster exceeds the total average R (F, M), then an upward arrow  $\uparrow$  is shown, in the opposite case, a downward arrow  $\downarrow$  is shown.

Figure 2 presents the clustering result of the K-means clustering algorithm by Weka software using the score value of RFM for the time frame of 12 month from 1/1/2010 to 31/12/2010. The figure shows eight clusters each with the corresponding number of customers, the mean score value of RFM of each cluster and the overall average of the scored RFM attributes. Also the figure shows the SSE of clustering result.

Figure 3 shows eight clusters resulted from using the actual value of RFM in clustering for the time frame of 12 month from 1/1/2010 to 31/12/2010. The figure displays the corresponding number of customers in each resulted cluster, the mean actual value of RFM of each cluster, the overall average of actual RFM attributes, and the SSE of clustering result.

### D. Results Evaluation

The accuracy of the resulted clusters based on the actual and scored RFM is examined based on the intra cluster distance using the Sum of Squared Error measure (SSE). It can be noticed from Figure 2 and Figure3 that the Sum of Squared Error of the resulted clusters at the first time frame from 1/1/2010 to 31/12/2010 is 45.18 using the scored RFM and 14.56 using actual RFM, which means that using the actual value of RFM in customer segmentation minimize the sum of intra cluster distance between elements and their cluster center, and result in more clustering accuracy than using the score RFM.

To validate this result the same procedures were applied on two different datasets with different number of customers, different number of transactions, and with two different timeframe (9 month in 2009, and 6 months in 2012).

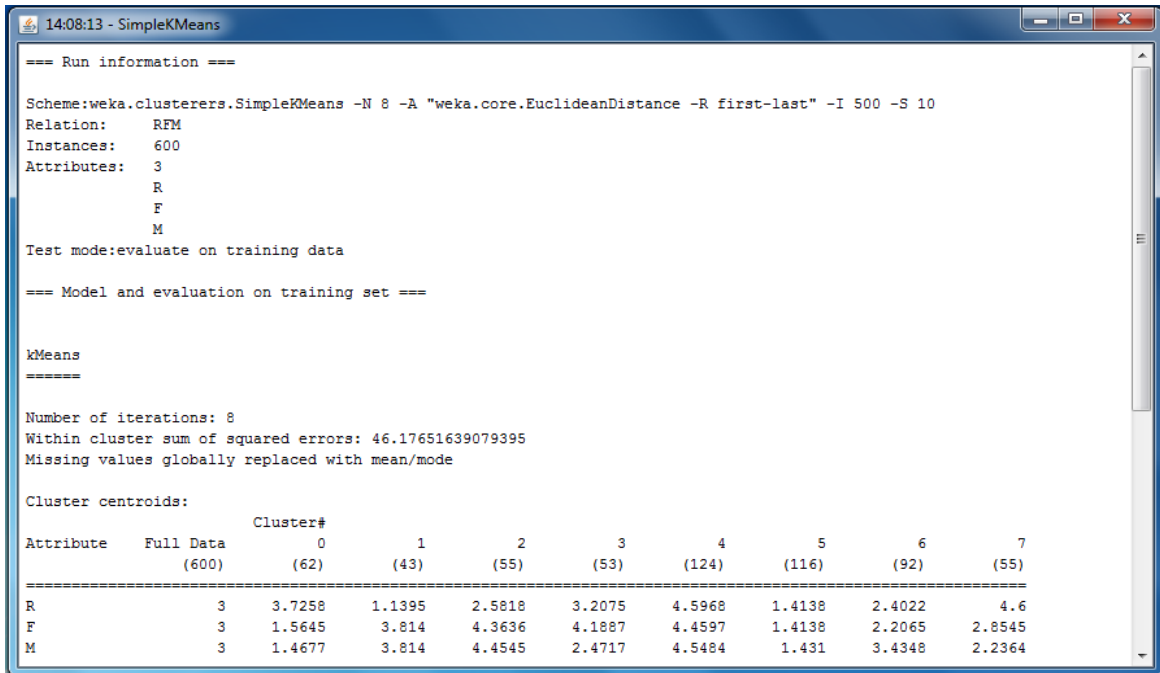


Figure 2. Scored RFM Clustering Result sand SSE for 12 month time frame from 1/1/2010 to 31/12/2010

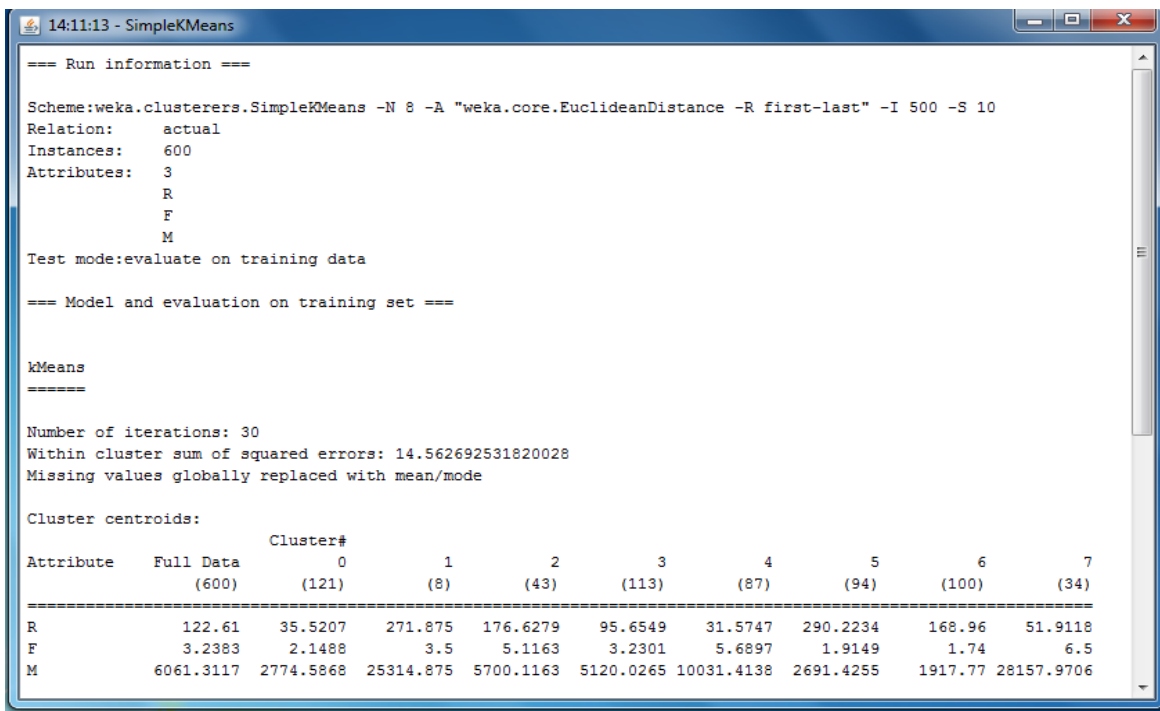


Figure 3. Actual value RFM Clustering Results and SSE for 12 month time frame from 1/1/2010 to 31/12/2010.

The resulted Sum of Squared Error of clusters at the second time frame from 1/1/2009 to 30/9/2009 is 42.51 using the score RFM and 10.87 using the actual RFM. Regarding the third time frame from 1/7/2012 till 31/12/2012, the within cluster sum of squared error using the scored RFM is 32.5 and 6.74 using the actual RFM.

Results of SSE of customer segmentation using scored RFM and actual RFM at different time frames insure our findings that using the actual value of RFM in customer segmentation result in more compact clusters with little deviation from the cluster centroids than using the scored value of RFM.

Table V provides a summary of previous results and shows a comparisons of customer segmentation quality based on the two methods of RFM preprocessing (actual and scored value) in term of the SSE.

As shown in Table V, using the actual value of RFM attributes in customer value segmentation at three time frames considered in this study reduces the sum of intra cluster distance between observations and their cluster center than using the scored value of RFM attributes, which means that the actual value of RFM reduce the clustering error and result in more compact and accurate clusters than using the scored RFM.

TABLE V.  
WITHIN CLUSTER SSE COMPARISON

Dataset		Resulted Sum of squared error (SSE)		
Time frame	Number of Customers	Number of Transactions	Score RFM	Actual RFM
12 month	600	1747	45.18	14.56
9 months	550	1640	42.51	10.87
6 months	435	1027	32.5	6.75

## V. CONCLUSION AND FUTURE WORKS

This study has evaluated the quality of customer value segmentation using the actual value of RFM and the score RFM. K-means algorithm is used for Customer value segmentation, and SSE is used as a measure for the resulted clusters accuracy. The study is applied on transaction data with three time frame to validate results. Results shows that using the actual value of RFM in customer value segmentation yield more clustering compactness and accuracy than using the score RFM. For future research, other types of clustering algorithms and clustering quality measures can be considered to assess our findings.

## REFERENCES

- [1] K. Tsipstsis, and A. Chorianopoulos, "Data Mining Techniques in CRM: Inside Customer segmentation", 1<sup>st</sup> ed, Wiley Publishing Inc, 2009, p.17.
- [2] M. Kim, J. Eun Park, A.J. Dubinsky, and S. Chaui, "Frequency of CRM implementation activities: a customer-centric view", Journal of Services Marketing, 2012m p. 84-85. <https://doi.org/10.1108/08876041211215248>
- [3] M. Berry, and G. Linoff, "Data Mining Techniques for Marketing, Sales, and Customer Relationship Management", 2<sup>nd</sup> ed, Wiley Publishing, Inc, 2004, p.3.
- [4] E. Ngai, L. Xiu, and D. Chau, "Application of data mining techniques in customer relationship management", elsevier, 2009, p.1-2.
- [5] R. Swift, "Accelerating customer relationships: Using CRM and relationship Technologies", upper saddle river. N.J.: Prentice Hall, 2001, p.12.
- [6] V. Kumar, and W. Reinartz, "Strategic Customer Relationship Management Today", Springer, 2012, p. 5-6. <https://doi.org/10.1007/978-3-642-20110-3>
- [7] A. Payne, "The Value Creation Process in Customer Relationship Management", Insight Interactive, 2002, p.2-3.
- [8] J. Han, M. Kamber, and J. Pei, "Data Mining: Concepts and Techniques", Elsevier, 2011, p.362-408.
- [9] P. Tan, M. Steinbach, and V. kumar, "Introduction to data mining", Addison-Wesley, 2005, p.317.
- [10] C. Astudillo, M. Bardeen, and N. Cerpa, "Data Mining in Electronic Commerce – Support vs Confidence", Journal of Theoretical and Applied Electronic Commerce Research, 2014, p.3.
- [11] A. Hughes, "Strategic database marketing" Probus Publishing Company, 1994.
- [12] C. Cheng, and Y. Chen, "Classifying the segmentation of customer value via RFM model and RS theory", Elsevier, 2008, p.4179-4180.
- [13] D. Birant, "Data Mining Using RFM Analysis", INTECH publisher, 2011, p.92-93.
- [14] M. Namvar, S. Khakabimamaghani, and M. Gholamian, "An approach to optimize customer segmentation and profiling using RFM, LTV, and demographic features", Journal of Electronic Customer Relationship Management, 2011, p. 227.
- [15] R. Qiasi, M. baqeri-Dehnavi, B. Minaei-Bidgoli, and G. Amooee, "Developing a model for measuring customer's loyalty and value with RFM technique and clustering algorithms", Journal of Mathematics and Computer Science, 2012, p.176-179.
- [16] P. Bunnak, S. Thammaboosadee, and S. Kiattisin, "Applying Data Mining Techniques and Extended RFM Model in Customer Loyalty Measurement", Journal of Advances in Information Technology, 2015, p.240.
- [17] D. Chen, S. Sain, and K. Guo., "Data mining for the online retail industry: A case study of RFM model-based customer segmentation using data mining", Journal of Database Marketing & Customer Strategy Management, 2012, p.199-201. <https://doi.org/10.1057/dbm.2012.17>
- [18] B. Stone, and R., Jacobs, "Successful direct marketing methods", NTC Business Books, Lincolnwood, 1995, p.37-57.
- [19] D. Liu, and Y. Shih, "Hybrid approaches to product recommendation based on customer lifetime value and purchase preferences", Elsevier, 2005, p.185-186.
- [20] M. Khajvand, K. Zolfaghar, S. Ashoori, and S. Alizadeh, "Estimating customer lifetime value based on RFM analysis of customer purchase behavior", Elsevier, 2011, p. 60-62.
- [21] Z. Tabaei, and M. Fathian, "Developing W-RFM Model for Customer Value: An Electronic Retailing Case Study", IEEE, 2011, p.305-306.
- [22] D. Li, W. Dai, and W. Tseng, "A two-stage clustering method to analyze customer characteristics to build discriminative customer management: A case of textile manufacturing business", Elsevier, 2011, p. 7187-7189.
- [23] A. Parvaneh, H. Abbasimehr, and M. Tarokhc, "Integrating AHP and Data Mining for Effective Retailer Segmentation Based on Retailer Lifetime Value", Journal of Institutional Economics, 2012, p. 28-30.
- [24] A. Mesforoush, and M.J. Tarokh, "Customer Profitability Segmentation for SMEs Case Study: Network Equipment Company", International Journal of Research in Engineering and Technology, 2013, p.34-37.
- [25] S. Aksenova, "Machine Learning with WEKA", California State University, 2004, p.4-6.

## AUTHORS

**Nesma mahmoud Taher** is with Helwan University in Egypt, Cairo, Egypt (it\_nesmamahmoud@live.com).

**Doaa Elzanfaly** is with Helwan University in Egypt, Cairo, Egypt and British University in Egypt, Cairo, Egypt (Doaa.elzanfaly@bue.edu.eg).

**Shaimaa Salama** is with Helwan University in Egypt, Cairo, Egypt (shaimaa.salama@fci.helwan.edu.eg).

Submitted 31 October 2016. Published as resubmitted by the authors 03 December 2016.