

ASSOCIATIVE REGRESSIVE DECISION RULE MINING FOR PREDICTING CUSTOMER SATISFACTORY PATTERNS

SHUBHA. S¹ and Dr. P. SURESH²

¹Research Scholar, Bharathiar University, Coimbatore & Asst. Prof.,
Dept. of Computer Science, GFGC,
Malleshwaram, Bangalore, Karnataka, India.
vu3gim.shubha@gmail.com

²Research Supervisor & HOD of Computer Science,
Salem Soudeshwari College, Salem, Tamilnadu, India.
sur_bhoo71@rediffmail.com

ABSTRACT

Opinion mining also known as sentiment analysis, involves customer satisfactory patterns, sentiments and attitudes toward entities, products, services and their attributes. With the rapid development in the field of Internet, potential customer's provides a satisfactory level of product/service reviews. The high volume of customer reviews were developed for product/review through taxonomy-aware processing but, it was difficult to identify the best reviews. In this paper, an Associative Regression Decision Rule Mining (ARDRM) technique is developed to predict the pattern for service provider and to improve customer satisfaction based on the review comments. Associative Regression based Decision Rule Mining performs two-steps for improving the customer satisfactory level. Initially, the Machine Learning Bayes Sentiment Classifier (MLBSC) is used to classify the class labels for each service reviews. After that, Regressive factor of the opinion words and Class labels were checked for Association between the words by using various probabilistic rules. Based on the probabilistic rules, the opinion and sentiments effect on customer reviews, are analyzed to arrive at specific set of service preferred by the customers with their review comments. The Associative Regressive Decision Rule helps the service provider to take decision on improving the customer satisfactory level. The experimental results reveal that the Associative Regression Decision Rule Mining (ARDRM) technique improved the performance in terms of true positive rate, Associative Regression factor, Regressive Decision Rule Generation time and Review Detection Accuracy of similar pattern.

KEYWORDS

Associative Regression, Decision Rule Mining, Machine Learning, Bayes Sentiment Classification, Probabilistic rules.

1. INTRODUCTION

Recently, novel method enriching semantic knowledge bases for opinion mining in big data applications has been evolved. In Opinion mining, sentiment analysis is very difficult to discover like and dislike of people. Hence by learning matrices for words, model can handle unseen word compositions.

David C. Wyld et al. (Eds) : CSITY, SIGPRO, AIFZ, NWCOM, DTMN, GRAPHHOC - 2016

pp. 121–134, 2016. © CS & IT-CSCP 2016

DOI : 10.5121/csit.2016.60411

In order to estimate their helpfulness, text mining and predictive modeling techniques toward a more complete analysis of the information captured by user-generated online reviews has been presented by many researchers. Taxonomy Aware Catalog Integration (TACI) [1] integrated products coming from multiple providers by making use of provider taxonomy information ensuring scalability that are typical on the web. However, tuning parameters does not update unless significant improvement in accuracy to avoid over fitting and it does not use any target or source taxonomy during training or application of classifier.

Tweet Analysis for Real-Time Event Detection and Earthquake (TA-RTED) [2] designed a classifier based on the keywords in order to improve the earthquake detection extracted through tweets. However, the registered location might not be current location of a tweet and it might not hold for other events such as traffic jams, accidents, and rainbows.

Many researchers have published their study of machine learning approach. Machine learning approach was developed using naïve Bayes [3] for identifying and distributing healthcare information. However, the syntactic rule-based relation extraction systems are complex based on additional tools.

Sara Hajian and Josep Domingo-Ferrer et al., [4] handles discrimination prevention in data mining and it also used for direct or indirect discrimination prevention. However, it failed to address the data distribution. An efficient algorithm was designed in [5] for detecting the top-k totally and partially unsolved sequences. This algorithm also used for reducing the running time and improving the accuracy while preserving data quality. However, it does not increase the detection accuracy of similar pattern at a required level.

Opinion mining analyzed people's opinions, sentiments, and attitudes toward entities products, services, and their attributes. Characterization of event and prediction based on temporal patterns are detected using multivariate reconstructed phase space (MRPS) [6] using fuzzy clustering unsupervised method. However, the MRPS method provides more difficult event function for different applications.

Intrinsic and extrinsic domain relevance criterion was developed in [7] aimed at improving the feasibility and effectiveness of the approach. However, it difficult to detect opinion features, including non-noun features, infrequent features, and implicit features collectively.

Probabilistic Generative Classifiers [8] used two or more classifiers resulting in the improvement of similarity measure. However, it does not address the various prior distribution investigations. The classification of trajectories on road networks was analyzed in [9] using frequent pattern-based classification which improves the accuracy. However, it does not address the pattern-based classification. The multi-class sentiment classification using that Extreme Learning Machine (ELM) methods were described in [21] for detecting their respective performance in multi-class sentiment classification of tweets. However, but it does not increases the classification accuracy effectively.

The contribution of the paper is organized as follows. Associative Regression Decision Rule Mining (ARDRM) technique is presented to predict the pattern for service owner and increasing their customer satisfaction based on their review comments. The Machine Learning Bayes Sentiment Classifier is subjected in ARDRM technique to classify the class labels for each service reviews. By applying the various probabilistic rules, the regressive factor of the opinion words and Class labels are verified between the words. This helps to increase the review detection accuracy.

The rest of the paper is organized as follows. Section 2 introduces several data mining models. Section 3 introduces our Associative Regression Decision Rule Mining technique based on the customer review comments. Section 4 presents the experimental setting and Section 5 presents the results of performance evaluation. Finally, the concluding remark is presented in Section 6.

2. RELATED WORK

In [11], a predictive model using three classification algorithms called decision tree, Bayesian classifier and back propagation neural network was presented. This model improved the diagnosis and prediction of cerebrovascular disease. Another predictive model using gradient-boosted regression tree [12] to make prediction aiming at reducing the execution flow. However the prediction accuracy did not effectively increase.

Many research works were conducted to answer top-k queries using Pareto Based Dominant Graph (DG) [10] aiming at improving the search efficiency. However, the relationship analysis remained unaddressed. Fast Distributed Mining (FDM) algorithm was designed in [13] for mining of association rules in horizontally distributed databases in a secured manner aiming at minimizing the communication rounds, communication and computational cost.

With the emergence of social media, web users have opened with a venue for expressing and sharing their thoughts and opinions related to different topics and events. Twitter feeds classification based on a hybrid approach was presented in [14] to achieve higher accuracy. However, this approach does not increase the accuracy at a required level.

In [15], a unified framework called, HOCTracker presented a novel density-based approach to identify hierarchical structure of overlapping communities. Probabilistic neural network and general regression neural network (PNN/GRNN) data mining model was planned in [16] for detect and preventing the oral cancer at earlier stage and also provides higher accuracy.

An incremental classification algorithm in [17] with feature space heterogeneity efficiently removed the outliers and extracted the relevant features at an early time period. In [18], an extensible method to mine experiential patterns from increasing game-logs was designed by utilizing the growing patterns.

An enhanced k-means clustering was applied in [19] to reduce the coefficient of variation and execution time using greedy approach for efficient discovery of patterns in health care data. In [20], random forest predictions were made using random forest algorithm to display prediction uncertainty. However, the true positive rate was not addressed.

Based on the aforementioned issues such as lack of detection in classification accuracy and failure in detecting the specified event in customer reviews, Associative Regression Decision Rule Mining (ARDRM) technique is presented. The ARDRM technique helps the service provider for improving the hotel customer satisfactory level at different cities. The detailed explanation is presented in forthcoming section.

3. DESIGN OF ASSOCIATIVE REGRESSION DECISION RULE MINING

Our technique Associative Regression based Decision Rule Mining is been done in two step process. First, the Machine Learning Bayes Sentiment Classification use a base classifier where the class labels for each product/service reviews is classified. Then the opinion words and class labels are used to obtain the regressive factor using various probabilistic rules to produce a final decision on improving the customer satisfaction referred to as the Associative Regression Decision Rule model. These two steps are now discussed in detail.

Figure 1 shows the workflow of Associative Regression Decision Rule Mining technique. Given a domain-dependent review comments (i.e. opinion words) extracted from OpinRank dataset that includes the reviews of hotels in 15 different cities, we first extract a list of class labels from the Machine Learning Bayes Sentiment Classification via semantic equivalence of sentiments classification.

For each extracted class labels, we estimate its regression factor which represents the statistical association between opinion words and class labels. The resultant regressive sequence inferred is

then applied with probabilistic rules to arrive at specific set of services preferred by the customers.

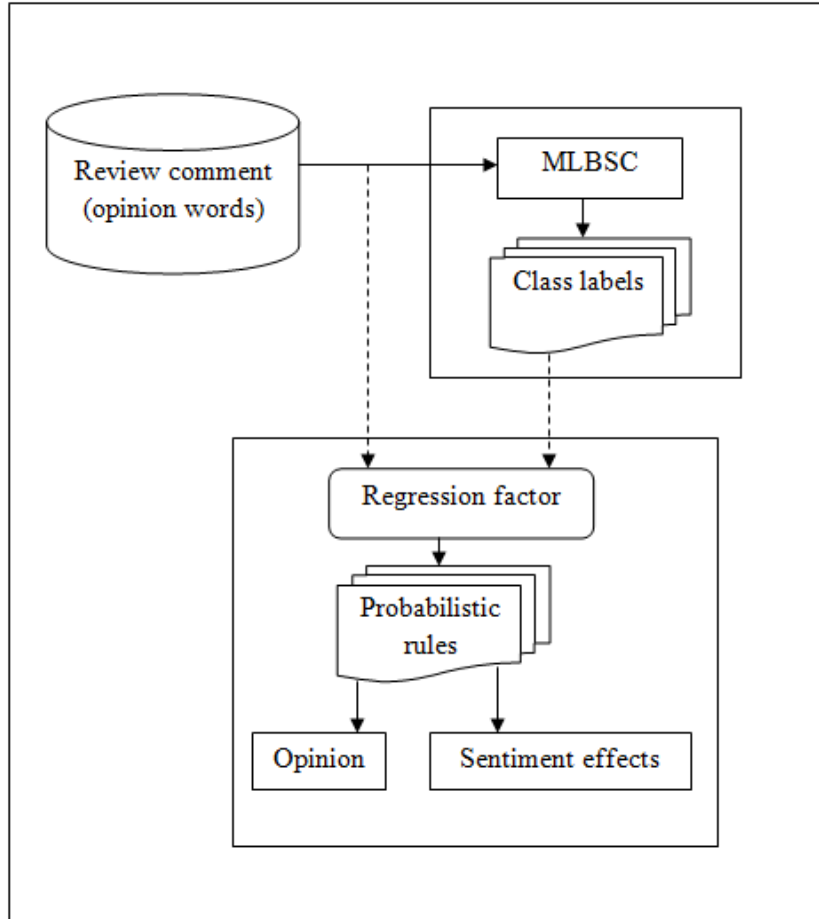


Figure 1. Workflow of Associative Regression Decision Rule Mining technique

3.1. Design of Regressive Sequencing Model

The first step in the design of ARDRM technique is to obtain the class labels generated from Machine Learning Bayes Sentiment Classification (MLBSC) techniques. Here the sentiment class labels are extracted using Probabilistic Bayes Classifier. In MLBSC, Probabilistic Bayes Classifier is applied on the semantic opinion words to evaluate sentiment class label using the maximum likelihood estimates (MLE). The MLE of a training list (i.e. bag of words extracted from OpinRank dataset) belonging to a specific class are mathematically evaluated as given below.

$$MLE \left(\frac{B_i}{C} \right) = \frac{\text{Count of } B_i \text{ in semantic opinion words of Class } C}{\text{Total number of words in semantic opinion words of Class } C} \quad (1)$$

From (1), the maximum likelihood estimates is the ratio of count of semantic opinion words of class 'C' to the total number of words. Followed by this, the class labels generated from MLBSC are subjected to regressive sequencing to infer the sentiments reflected in the customer reviews. The regressive sequencing in ARDRM technique is produced with the aid of support and confidence value.

Let us assume that ' $I = i_1, i_2, \dots, i_n$ ' represents a binary set consisting of opinion words with ' i_1, i_2, \dots, i_n ' referred to as items. Let us further assume that Transaction ' T ' (i.e. review

comments) is the itemset with ' $T \in I$ '. Let ' P ' be the set containing items in ' I ' and transaction ' T ' contains ' P ' if ' $P \in T$ ', then the support denotes the probability of frequent itemsets' occurrence. Smaller value of minsup results in larger number of rules whereas larger value of minsup results in smaller number of rules.

The support of rule ' $P \rightarrow Q$ ' in the transaction database ' TD ' is the ratio between the transaction number including ' P ' and ' Q ' in the transaction sets and all the transaction number, which is then written as ' $SUP(P \rightarrow Q)$ '.

$$SUP(P \rightarrow Q) = \frac{Prob(PQ)}{N} \quad (2)$$

The confidence of the rule ' $P \rightarrow Q$ ' in the transaction sets is the ratio between the transaction number including ' P ' and ' Q ' and those including ' P ', which is written as ' $CONF(P \rightarrow Q)$ '. Therefore,

$$CONF(P \rightarrow Q) = \frac{Prob(PQ)}{Prob(P)} \quad (3)$$

From (2) and (3), the sentiments reflected in customer reviews are obtained by using support and confidence value. This aids in achieving the true positive rate of customer reviews in an extensive manner. Once, the support and confidence value for customer reviews are generated, the regressive sequencing is designed. In MLBSC, the regressive sequencing model uses two variables ' y_1 ' and ' y_2 ' where ' y_1 ' represents '*minsup*' and ' y_2 ' represents '*minconf*' to infer the sentiments reflected in the customer reviews. The mathematical formulates for ' y_1 ' and ' y_2 ' is as given below.

$$x = \delta_0 + \delta_1 \left(\frac{1}{y_1^2}\right) + \delta_2 \left(\frac{1}{y_2^2}\right) \quad (4)$$

$$x = \delta_0 + \delta_1 \left(\frac{1}{y_1}\right) + \delta_2 \left(\frac{1}{y_2}\right) \quad (5)$$

$$x = \delta_0 + \delta_1 (y_1) + \delta_2 (y_2) \quad (6)$$

The regressive factor (i.e. ' y_1 ' and ' y_2 ') of the opinion words with class labels are checked for the association between the words. This in turn improves the associative regression factor in a significant manner. Figure 2 shows the algorithmic description of Regressive Sequencing algorithm.

Input: opinion words ' $I = i_1, i_2, \dots, i_n$ ', Transaction ' T ',	
Output: Optimized true positive rate	
Step 1: Begin	
Step 2:	For each Transaction ' T ' with opinion words ' I '
Step 3:	Measure the value for support using (2)
Step 4:	Measure the value for confidence using (3)
Step 5:	Measure ' <i>minsup</i> ' and ' <i>minconf</i> ' using (6)
Step 6:	End for
Step 7: End	

Figure 2. Regressive Sequencing algorithm

From the above figure 2, the Regressive Sequencing algorithm performs three steps. For each transaction, customer reviews obtained from OpinRank dataset that includes hotel reviews is given as input. The first step evaluates the support value, followed by the measure of confidence value in order to identify the sentiments reflected in customer reviews. Finally, with the objective

of improving the associative regression factor, sentiments reflected in customer reviews ‘*minsup*’ and ‘*minconf*’ are evaluated to check the association between the words.

3.2. Design of Associative Regressive Decision Rule

The second step in the design of ARDRM technique is to construct Associative Regressive Decision Rule. Various probabilistic rules are generated for the class objects in the corresponding classes with more similar patterns together. Based on the probabilistic rules, the opinion and sentiments effect on customer reviews are analyzed to arrive at specific set of services preferred by the customers with their review comments. The Associative Regressive Decision Rule helps the service providers to take decision on how to improve the hotel customer satisfactory level. Next, the frequent itemset generation algorithm is designed to the regressive sequenced dataset which is obtained through regressive sequencing model. Redundant regressive rules generated are eliminated using redundant regressive decision rule testing.

3.2.1. Redundant Regressive Decision Rule Testing

Redundant regressive decision rule testing is performed in ARDRM technique aiming at minimizing the regressive decision rule generation time and removes the redundancy involved. This is performed through elimination of redundant decision rule through regressive model.

$$P_a, P_b, P_c, \dots, P_n \rightarrow \sum_{i=1}^n (y_i, \mu_i, \sigma_i) \quad (7)$$

$$\text{First set of variance } (P_{ab}) = P_a - P_b \quad (8)$$

$$\text{Second set of variance } (P_{bc}) = P_b - P_c \quad (9)$$

From (7), ‘ μ_i ’ symbolizes the mean of the target review for the class objects and ‘ σ_i ’ symbolizes the variance of the target review for the class objects. In (7), the mean and variance are evaluated. The variance of the association rules are calculated using (8), (9) where P_a, P_b, P_c are the target reviews for class object. If the variance (first set) of the association rule is lower than the variance (second set) of the association rule, then redundancy is said to be occurred in the first set. On contrary, if the variance (first set) of the association rule is greater than the variance (second set) of the association rule, then redundancy is said to be occurred in the second set. By using specified threshold value, the redundant rules are eliminated. If the identified redundancy value is obtained within the threshold value, the redundant rules are eliminated. The Redundancy value is possibly occurred within the thresholding value. This in turn minimizes the regressive decision rule generation time.

3.2.2. Associative Regressive Decision Model

Once the redundant rules are eliminated using redundant regressive decision rule testing then, finally associative regressive decision model is designed to arrive at specific set of service preferred by the customers. Building an associative regressive decision model requires selection of a smaller, representative set of rules in order to provide an accurate representation of the training data.

The frequent itemset generation algorithm is shown in figure 3 to select the rule in an efficient manner by first sorting the rule, and then remove the occurrences covered by the rule. As shown in the algorithm, for each itemset, the algorithm starts with the elimination of redundant rule. Followed by this rule redundant removal, the occurrence of redundancy is observed and removed in specified threshold value. Then rule sorting is performed based on the pair of rules. Finally, Associative Regressive Decision model is applied to the generated rules that help the service provider to customer satisfactory level.

Let us consider a pair of rules, ‘ $Rule_1$ ’ and ‘ $Rule_2$ ’ where ‘ $Rule_1 \gg Rule_2$ ’. This implies that ‘ $Rule_1$ ’ has higher preference over ‘ $Rule_2$ ’ and is formulated as given below.

$$if(Rule_1 \gg Rule_2) \rightarrow ARD = Rule_1, Rule_2 \quad (8)$$

$$if(Rule_2 \gg Rule_1) \rightarrow ARD = Rule_2, Rule_1 \quad (9)$$

Input: mean of the target review for the class objects ' μ_i ', variance of the target review for the class objects ' σ_i ', first set P_{ab} , second set P_{bc} ,	
Output: Improved customer satisfactory level	
Step 1:	Begin
Step 2:	For each set P_i
Step 3:	Perform redundant rule elimination through
Step 4:	If $\sigma_i(P_{ab}) < \sigma_i(P_{bc})$
Step 5:	Redundancy is found in said to be occurred in (P_{ab})
Step 6:	End if
Step 7:	If $\sigma_i(P_{ab}) > \sigma_i(P_{bc})$
Step 8:	Redundancy is found in said to be occurred in (P_{bc})
Step 9 :	End if
Step 10:	If ($T_h \geq redundancy\ value$)
Step 11:	The redundant rule is eliminated
Step 12:	End if
Step 13:	End for
Step 15:	Perform rule sorting
Step 16:	If ($Rule_1 \gg Rule_2$)
Step 17:	$ARD = Rule_1, Rule_2$
Step 18:	End if
Step 19:	If ($Rule_2 \gg Rule_1$)
Step 20:	$ARD = Rule_2, Rule_1$
Step 21:	End if
Step 22:	Perform Associative regressive decision model using (10)
Step 23:	End

Figure 3. Associative regressive decision-based frequent itemset generation algorithm

On contrary, if ' $Rule_2 \gg Rule_1$ ', then ' $Rule_2$ ' has higher preference over ' $Rule_1$ '. Once the sorted rules are obtained, the final step is to design Associative Regressive Decision model. The Associative Regressive Decision model is designed in such a way that, the rule has higher support value and has lower variance when ' $Rule_1$ ' and ' $Rule_2$ ' are applied. Then, the mathematical formulates

$$(Rule_1, Rule_2) \rightarrow (MaxSup (Rule_1, Rule_2), MinVar (Rule_1, Rule_2)) \quad (10)$$

Based on (10), several probabilistic rules are generated for the class objects with more similar patterns together and also the rule the service preferred by the customers with their review comments. This in turn helps the service providers to take decision on improving the customer

satisfactory level, thereby improving the review detection accuracy based on the review comments of the customers.

4. EXPERIMENTAL SETTINGS

Associative Regression Decision Rule Mining (ARDRM) technique uses JAVA platform with WEKA tool to predict a predictive pattern for service owner to improve their customer satisfaction based on their review comments. This method is widely used to perform efficient predictive pattern mining model with the tests and training samples. Hotel Customer Service Reviews (eg: OpinRank Dataset - Reviews from TripAdvisor) is taken to perform the experimental work. The training model for OpinRank dataset includes entire hotel reviews situated in 10 different cities (Dubai, Beijing, London, New York, New Delhi, San Francisco, Shanghai, Montreal, Last Vegas and Chicago) with the aid of Java platform and with WEKA tool. This dataset has been chosen because it gives a clear picture that helps in analyzing the comments made by tourists about hotel rooms and food provided. The total number of reviews included in OpinRank dataset is 250,000. For experimental purpose, we reviewed using 350 and the extracted field includes date of review, review title and full review made by the tourists.

The performance of Associative Regression Decision Rule Mining (ARDRM) technique is compared with Taxonomy-Aware Catalog Integration (TACI) [1], and Tweet Analysis for Real-Time Event Detection and Earthquake (TA-RTED) [2]. The tests on OpinRank dataset were conducted to evaluate four parameters: true positive rate, associative regression factor, regressive decision rule generation time and review detection accuracy of similar pattern.

5. DISCUSSION

The Associative Regression Decision Rule Mining (ARDRM) technique is compared against the existing Taxonomy-Aware Catalog Integration (TACI) [1] and Tweet Analysis for Real-Time Event Detection and Earthquake (TA-RTED) [2]. The experimental results using JAVA platform with WEKA are compared and analyzed with the aid of graph form as given below.

5.1. Impact of True positive rate

The true positive rate is the sentiments correctly identified as belonging to a specific class in customer reviews words correctly identified as belonging to a specific class. The true positive rate is mathematically formulated as given below.

$$TPR = \left(\frac{\text{sentiments correctly identified as belonging to a class}}{c} \right) * 100 \quad (11)$$

From (11), the true positive rate 'TPR' is obtained using the class 'C' where each class consists of different number of sentiments extracted from user review. It is measured in terms of percentage (%). The convergence plot for 7 classes is depicted in table 1 and figure 4. From the figure 4 we can note that the proposed ARDRM technique achieved maximum true positive rate on sentiments being correctly identified as belonging to a specific class when compared to other methods.

Table 1. Tabulation for true positive rate

Class (C)	True Positive Rate (%)		
	ARDRM	TACI	TA-RTED
Class 1	84.32	73.15	68.21
Class 2	88.15	77.12	71.08
Class 3	91.35	80.32	74.28
Class 4	85.21	74.18	68.10
Class 5	87.57	76.54	70.46
Class 6	89.32	78.29	72.22
Class 7	92.14	81.11	75.04

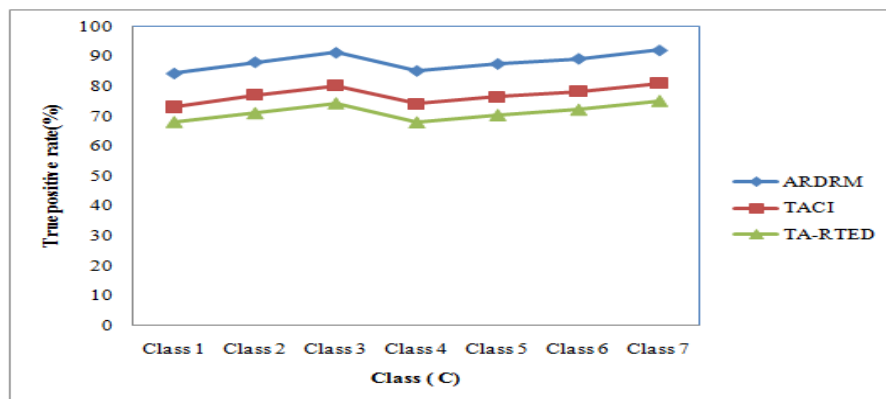


Figure 4 Measure of true positive rate

Figure 4 shows the true positive rate on sentiments being correctly identified as belonging to a specific class is increased with the application of maximum likelihood estimates when compared to the existing methods. The maximum likelihood estimates in ARDRM technique effectively constructs sentiment class label for the testing and training data extracted from OpinRank dataset. Therefore, the true positive rate is improved by 12.52% compared to TACI [1]. Moreover, by evaluating the support and confidence value, probability of frequent itemsets occurrence are made in a significant manner. As a result, the true positive rate is increased by 19.21% compared to TA-RTED [2].

5.2. Impact of associative regression factor

Table 2 shows the associative regression factor using the three methods, ARDRM, TACI [1] and TA-RTED [2] respectively. The associative regression factor in table 2 was measured with the aid of 7 classes and 35 rules generated from 350 customer review words extracted from the OpinRank dataset.

Table 2 Tabulation for associative regression factor

METHODS	ASSOCIATIVE REGRESSION FACTOR (%)
ARDRM	81.35
TACI	74.19
TA-RTED	65.18

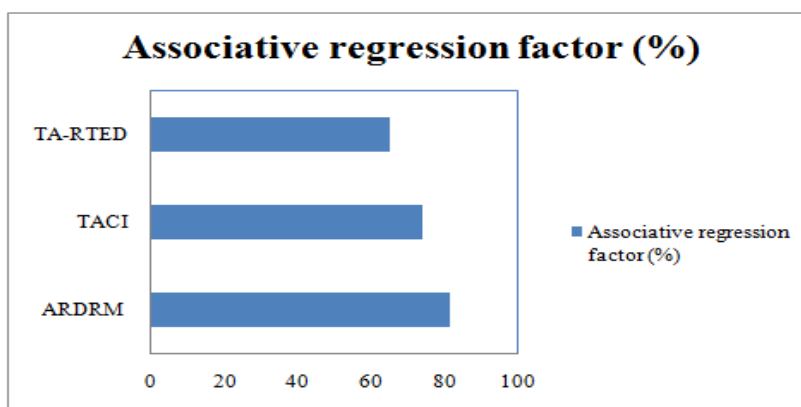


Figure 5. Measure of associative regression factor

Figure 5 shows the measure of associative regression factor with respect to 350 customer review words obtained from OpinRank dataset. The associative regression factor using ARDRM is improved when compared to two other methods [1] and [2]. This is due to the application of Regressive Sequencing algorithm. By applying Regressive Sequencing algorithm, the support and confidence value are evaluated according to the sentiments reflected in the customer review. This in turn improves the associative regression factor using ARDRM by 8.80% compared to TACI and 12.14% compared to TA-RTED respectively.

5.3. Impact of Regressive decision rule generation time

The regressive decision rule generation time is measured using the number of rules and the time to extract single rule. The mathematical formulation for regressive decision rule generation time is given as below.

$$DRGT = \sum_{i=1}^n Rule_i * Time (Rule_i) \quad (12)$$

From (12), the execution time ‘DRGT’ is measured using the number of rules ‘Rule_i’ and measured in terms of milliseconds. Lower the regressive decision rule generation time more efficient the method is said to be. Convergence characteristics for the measure of Time to extract opinions from customer reviews for 35 rules extracted from different customers are considered and compared with two other methods and are shown in table 3.

Table 3. Tabulation for time to extract opinions from reviews

Number of rules	Time to extract opinions from customer reviews (ms)		
	ARDRM	TACI	TA-RTED
5	1.31	1.68	1.85
10	2.51	2.81	3.05
15	3.79	4.02	4.22
20	4.96	5.26	5.48
25	5.85	6.15	6.35
30	6.3	6.60	6.80
35	8.32	8.62	8.82

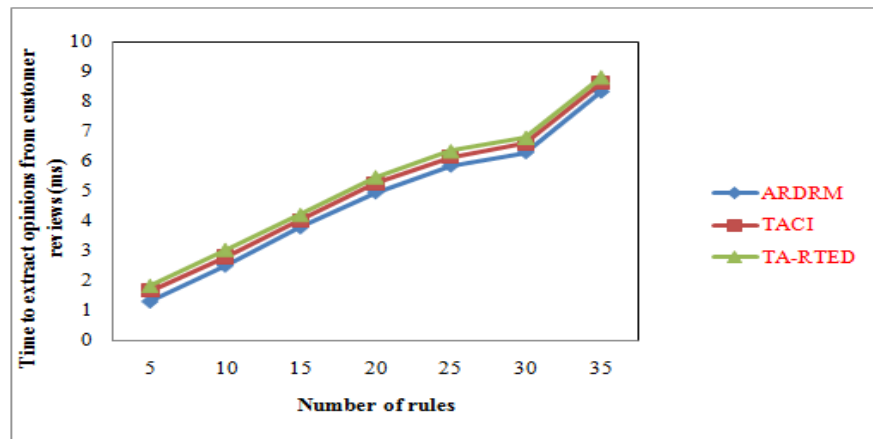


Figure 6 Measure of Regressive decision rule generation time

The targeting results of Regressive decision rule generation time for extracting predictive pattern using ARDRM technique is compared with two state-of-the-art methods [1], [2] in figure 6 is presented for visual comparison based on the number of rules. Our method differs from the FM-TACI [1] and TA-RTED [2] in that we have incorporated associative regressive decision rule. The associative regressive decision rule applies probabilistic rules using the mean and variance

value for performing rule generation. As a result, the Regressive decision rule generation time generating decision rules using ARDRM technique is increased by 9.40 to TACI. Furthermore, by eliminating the redundant rule, further reduces the time for obtaining the regressive decision rule generation by 15.29% compared to TA-RTED.

5.4. Impact of Review detection accuracy of similar pattern

The review detection accuracy of similar pattern is the ratio of number of correct review patterns to the total number of test cases made. The mathematical formulation of review detection accuracy of similar pattern is formulated as given below.

$$A = \left(\frac{\text{No.of correct review patterns}}{\text{Total no.of test cases}} \right) * 100 \quad (13)$$

From (13), the detection accuracy 'A' is measured in a significant manner in terms of percentage (%). Higher the detection accuracy more efficient the method is said to be.

Table 4 Tabulation for review detection accuracy

Customer review words	Review detection accuracy (%)		
	ARDRM	TACI	TA-RTED
50	87.53	74.11	68.21
100	89.31	77.27	71.25
150	92.14	80.10	74.08
200	85.14	73.10	67.07
250	88.21	76.17	70.15
300	88.15	86.11	80.11
350	91.35	79.31	73.21

The comparison of customer review detection accuracy is presented in table 4 with respect to different customer review words. Depending on the customer review words, the customer review detection accuracy either increases or decreases but found to be improved using the proposed ARDRM technique.

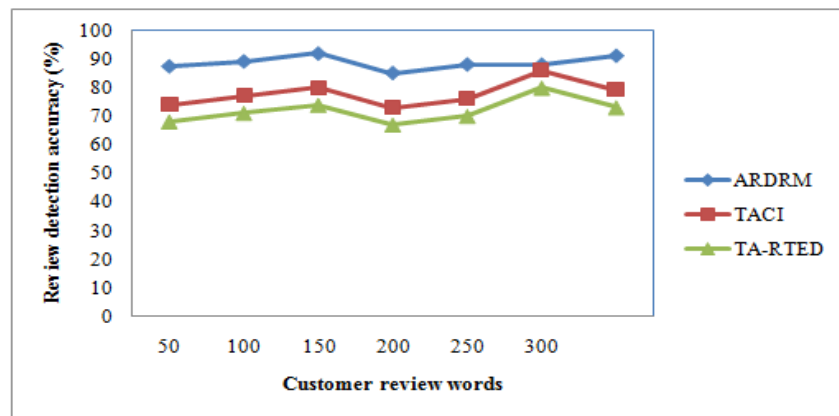


Figure 7. Measure for review detection accuracy

To ascertain the performance of customer review detection accuracy, comparison is made with two other existing works Taxonomy-Aware Catalog Integration (TACI) [1], and Tweet Analysis for Real-Time Event Detection and Earthquake (TA-RTED) [2].

In figure 7, the customer review words varied between 50 and 300. From the figure it is illustrative that the customer review detection accuracy is improved using the proposed ARDRM technique when compared to two other existing works. This is because with the application of

Associative regressive decision-based frequent itemset generation algorithm, the ARDRM technique chooses the rule in a greedy manner by first sorting the rule and detaches the occurrences covered by the rule.

In this way, the customer review detection accuracy is improved using ARDRM by 12.16% when compared to TACI [1]. Furthermore, by applying associative regressive decision model when applied to the generated rules, with higher support value and lower variance improves the customer satisfactory level, therefore improving the review detection accuracy based on their review comments of the customers by 18.93% than when compared to TA-RTED [2].

5.5. Performance analysis of customer review classification accuracy using proposed ARDRM and Extreme Learning Machine

The result analysis of the proposed Associative Regression Decision Rule Mining (ARDRM) method is compared with existing Multi-class Sentiment Classification using Extreme Learning Machine (MSC- ELM) [21].

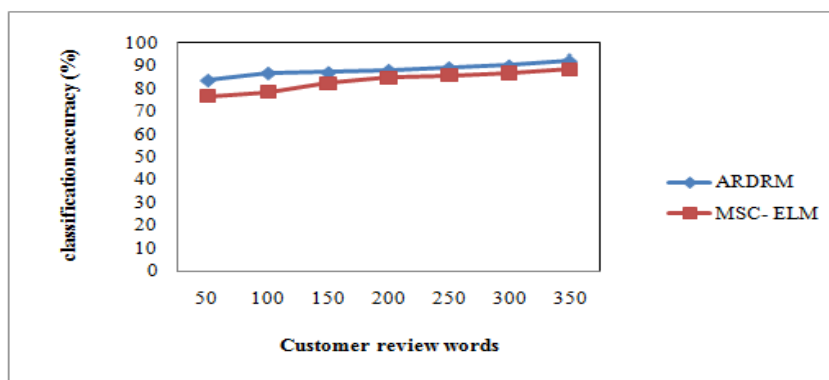


Figure 8. Measure of review classification accuracy

Figure 8 illustrates the customer review classification accuracy of proposed ARDRM and existing Multi-class Sentiment Classification using Extreme Learning Machine (MSC- ELM). From the figure, the customer review classification accuracy is increased in ARDRM. This is because, the Machine Learning Bayes Sentiment Classifier (MLBSC) is applied in ARDRM to classify the class labels for each service reviews. Therefore, the classification accuracy is effectively increased by 6% in ARDRM method compared to existing Multi-class Sentiment Classification using Extreme Learning Machine (MSC- ELM) [21].

6. CONCLUSION

In this work, an effective technique called Associative Regression Decision Rule Mining (ARDRM) is presented. The technique improves the review detection accuracy that in turn improves the customer satisfaction based on their review comments and associative regression factor. The goal of Associative Regression Decision Rule Mining is to improve the true positive rate with sentiments correctly identified as belonging to a specific class and therefore to improve the associative regression factor using the customer review words extracted from OpinRank dataset which significantly contribute to the relevance. To do this, we first designed a Machine Learning Bayes Sentiment Classification technique that measures the sentiment class labels based on the Maximum Likelihood estimates for OpinRank dataset this helps to increase the classification accuracy. Then, based on this measure, we proposed a Regressive Sequencing algorithm for improving the association regression factor in an extensive manner. In addition the associative regressive decision rule with frequent itemset generation algorithm eliminates the redundant rule and therefore reduces the time of extract opinions reviews and therefore true positive rate. Finally, the associative regressive decision model improves the customer review

detection accuracy. Extensive experiments were carried out using JAVA and compared with existing methods. The results show that ARDRM technique offers better performance with an improvement of review detection accuracy by 15.55% and reduces the time taken to extract opinions from reviewers by 12.34% compared to TACI and TA-RTED respectively.

ACKNOWLEDGEMENTS

We thank each one of those who are directly or indirectly helpful in the preparation of this paper. And also we thank our family members for their continuous support and encouragement.

REFERENCES

- [1] Panagiotis Papadimitriou, Panayiotis Tsaparas, Ariel Fuxman, and Lise Getoor, "TACI: Taxonomy-Aware Catalog Integration", *IEEE Transactions on Knowledge and Data Engineering*, Volume 25, Issue 7, July 2013, Pages 1643-1655.
- [2] Takeshi Sakaki, Makoto Okazaki, and Yutaka Matsuo, "Tweet Analysis for Real-Time Event Detection and Earthquake Reporting System Development", *IEEE Transactions on Knowledge and Data Engineering*, Volume 25, Issue 4, April 2013, Pages 919-931.
- [3] Oana Frunza, Diana Inkpen, and Thomas Tran, "A Machine Learning Approach for Identifying Disease-Treatment Relations in Short Texts", *IEEE Transactions on Knowledge and Data Engineering*, Volume 23, Issue 6, June 2011, Pages 801-814.
- [4] Sara Hajian and Josep Domingo-Ferrer, "A Methodology for Direct and Indirect Discrimination Prevention in Data Mining", *IEEE Transactions on Knowledge and Data Engineering*, Volume 25, Issue 7, July 2013, Pages 1445-1459.
- [5] Massimiliano Albanese, Cristian Molinaro, Fabio Persia, Antonio Picariello, and V.S. Subrahmanian, "Discovering the Top-k Unexplained Sequences in Time-Stamped Observation Data", *IEEE Transactions on Knowledge and Data Engineering*, Volume 26, Issue 3, March 2014, Pages 577-594.
- [6] Wenjing Zhang, and Xin Feng, "Event Characterization and Prediction Based on Temporal Patterns in Dynamic Data System", *IEEE Transactions on Knowledge and Data Engineering*, Volume 26, Issue 1, January 2014, Pages 144-156.
- [7] Zhen Hai, Kuiyu Chang, Jung-Jae Kim, and Christopher C. Yang, "Identifying Features in Opinion Mining via Intrinsic and Extrinsic Domain Relevance", *IEEE Transactions on Knowledge and Data Engineering*, Volume 26, Issue 3, March 2014, Pages 623-634.
- [8] Dominik Fisch, Edgar Kalkowski, and Bernhard Sick, "Knowledge Fusion for Probabilistic Generative Classifiers with Data Mining Applications", *IEEE Transactions on Knowledge and Data Engineering*, Volume 26, Issue 3, March 2014, Pages 652-666.
- [9] Jae-Gil Lee, Jiawei Han, Xiaolei Li, and Hong Cheng, "Mining Discriminative Patterns for Classifying Trajectories on Road Networks", *IEEE Transactions on Knowledge and Data Engineering*, Volume 23, Issue 5, May 2011, Pages 713-726.
- [10] Lei Zou, and Lei Chen, "Pareto-Based Dominant Graph: An Efficient Indexing Structure to Answer Top-K Queries", *IEEE Transactions on Knowledge and Data Engineering*, Volume 23, Issue 5, May 2011, Pages 727-741.
- [11] Duen-Yian Yeh, Ching-Hsue Cheng, Yen-Wen Chen, "A predictive model for cerebro vascular disease using data mining", Elsevier, *Expert Systems with Applications*, Volume 38, Issue 7, July 2011, Pages 8970-8977.
- [12] Nima Asadi, Jimmy Lin, Arjen P. de Vries, "Runtime Optimizations for Prediction with Tree-Based Models", *IEEE Transactions on Knowledge and Data Engineering*, Volume 26, Issue 9, September 2014, Pages 2281-2292.
- [13] Tamir Tassa, "Secure Mining of Association Rules in Horizontally Distributed Databases", *IEEE Transactions on Knowledge and Data Engineering*, Volume 26, Issue 4, April 2014, Pages 970-983.

- [14] Farhan Hassan Khan, Saba Bashir , Usman Qamar, "TOM: Twitter opinion mining framework using hybrid classification scheme", Elsevier, Decision Support Systems, Volume 57, January 2014, Pages 245–257.
- [15] Sajid Yousuf Bhat and Muhammad Abulaish, "HOCTracker: Tracking the Evolution of Hierarchical and Overlapping Communities in Dynamic Social Networks", IEEE Transactions on Knowledge and Data engineering, Volume 27, Issue 4, April 2014, Pages 1019-1032.
- [16] Neha Sharma and Hari Om, "Usage of Probabilistic and General Regression Neural Network for Early Detection and Prevention of Oral Cancer", Hindawi Publishing Corporation, The Scientific World Journal, Volume 2015(2015), May 2015, Pages 1-12.
- [17] Yu Wang, "An Incremental Classification Algorithm for Mining Data with Feature Space Heterogeneity", Hindawi Publishing Corporation, Mathematical Problems in Engineering, Volume 2014 (2014), February 2014, Pages 1-10.
- [18] LiangWang, Yu Wang, and Yan Li, "Mining Experiential Patterns from Game-Logs of Board Game", Hindawi Publishing Corporation, International Journal of Computer Games Technology, Volume 2015, December 2014 , Pages 1-21.
- [19] Ramzi A. Haraty, Mohamad Dimishkieh and MehediMasud, "An Enhanced k-Means Clustering Algorithm for Pattern Discovery in Healthcare Data", Hindawi Publishing Corporation, International Journal of Distributed Sensor Networks, Volume 2015, December 2014, Pages 1-12.
- [20] Mareike Lie, Martin Hitziger, and Bernd Huwe, "The Sloping Mire Soil-Landscape of Southern Ecuador: Influence of Predictor Resolution and Model Tuning on Random Forest Predictions", Hindawi Publishing Corporation, Applied and Environmental Soil Science, Volume 2014, February 2014, Pages 1-11.
- [21] Wang, Zhaoxia, and Yogesh Parth "Extreme Learning Machine for Multi-class Sentiment Classification of Tweets." Proceedings of ELM-2015 Volume 1, Springer International Publishing, 2016. 1-11