# Opinion Mining of Twitter Data for Hotel Review Analysis

## SRIKANTH M S, KAMAD BHATT, HARSH RAJ, MD. ADNAN AHMED

Assistant Professor, Dept. of Computer science and Engineering, Sapthagiri college of Engineering

(VTU)Bengaluru, India, srikanthms@sapthagiri.edu.in

Student, Dept. of Computer science and Engineering, Sapthagiri college of Engineering

(VTU), Bengaluru, India, kamadbhatt96@gmail.com

Student, Dept. of Computer science and Engineering, Sapthagiri college of Engineering

(VTU), Bengaluru,India, rajharsh1997@gmail.com

Student, Dept. of Computer science and Engineering, Sapthagiri college of Engineering

(VTU), Bengaluru, India, adnan.workoffice@gmail.com

**Abstract**— *The rapid increase in mountains of unstructured textual data accompanied by proliferation of tools to analyse them has opened up great opportunities and challenges for text mining research. The automatic labelling of text data is hard because people often express opinions in complex ways that are sometimes difficult to comprehend. The labelling process involves huge amount of efforts and mislabelled datasets usually lead to incorrect decisions. In this paper, we design a frame work for sentiment analysis with opinion mining for the case of hotel customer feedback. Most available datasets of hotel reviews are not labelled which presents a lot of works for researchers as fares text data pre-processing task is concerned. Moreover, sentiment datasets are often highly domain sensitive and hard to create because sentiments are feelings such as emotions, attitudes and opinions that are commonly rife with idioms, onomatopoeias, homophones, phonemes, alliterations and acronyms. The proposed framework is termed sentiment polarity that automatically prepares a sentiment dataset for training and testing to extract unbiased opinions of hotel services from reviews to discover a suitable machine learning algorithm for the classification component of the framework.*

***Keywords-sentiment analysis, text mining, association rule, bag of words, opinion mining***

## I. INTRODUCTION

In recent years, the world has experienced a tremendous rise in the volume of textual data especially for the unstructured data generated from people who express opinions through various web and social media platforms for different reasons. Mountains of these textual data, initially could be equated to garbage which would need to be disposed from time to time. However, with the advancement in storage capacity accompanied by the increasing sophistication in data mining tools, opportunities and challenges have been created for analysing and deriving useful insights from these mountains of data.
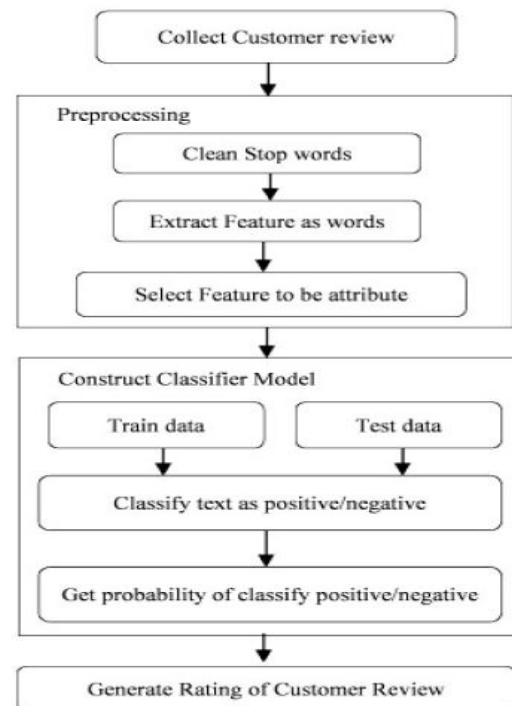
In this paper, we have chosen textual data in the form of hotel reviews for sentiment analysis with opinion mining from customer perspectives. Sentiment analysis uses the techniques of natural language processing and computational linguistics to automate the classification of sentiments generated from reviews. Hotels provide satisfaction, security, comfort, luxury and lodging services for travellers and people on vacation. Mining hotel reviews is desirable to gain deeper knowledge of

customer expectations and support effective management of customer relationships. It would enable the hotel managers to have a good understanding of customer needs, discover areas for further improvement and improve service quality. The hotel reviews are provided exclusively by customers who have made reservations at a particular hotel. Customers post feedback about hotels which include hygiene, quality of food, location, customer service quality and hospitality exhibited by hotel staff. Moreover, sentiment analysis of hotel reviews is crucial to understand hidden patterns generated by data that would help to effectively improve performance.

## II. RELATED WORKS

The opinion mining has become one of popular research area. The challenge is in process of opinion mining or sentiment analysis that is unstructured and noisy data on website. A part of opinion mining refers using of natural language processing (NLP) by proposed different method of dictionary for sentiment analysis of text as corpus, lexicon and specific language dictionary. They tried to extract word from sentences for removal stop word or unnecessary word automatically. In addition, various dictionaries are solved by machine learning methods, which try to rank scoring of various dictionaries. For example, the paper in used fuzzy logic algorithm to collect the ranking of different dictionary into rule for classify the opinion. After word segmentation process is removal stop words by dictionary checking. The group of researches focuses on the

calculating polarity of words to trend in positive or negative in a cluster of interest's customers that are extracted from texts and compared the word occurrence of whole sentence. If the word extractions have weight from dictionary of emotional words, it is calculated to answer the comment as positive or negative.

However, the customer review has different behaviour with the product. The proposed classifier model is presented using association rule in. The popular classifier model is nave Bayes compared with other model, which there are different sources such as social media and web site. From these researches are used classifier models that are the same objective to classified opinion. Our approach is different from them, this paper uses the advantage of classifier model to generate the rating value from classifier which is not only shown classify opinion as positive and negative and also factors analysis to impact the customer who posted or commented to positive and negative.



**Fig.1** Proposed Methodology for generating score of customer review using opinion mining
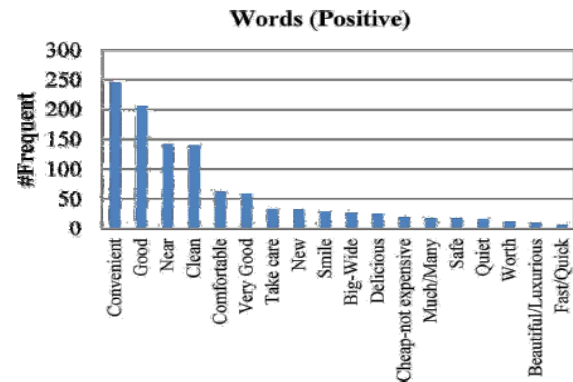
## III. PROPOSED METHODOLOGY

The proposed methodology used Thai customer review's hotels from a website of hotel agent service, which service in hotel reservation directly. The target of classify customer review from this website because the comment is posted from customer who is serviced checked-in and checked-out from hotel. The system has cleaned the promotion of hotel's comment which has

only existed customer review given comment and rating. The numbers of open opinion texts are collected 400 customer reviews that are used service to checked-in/out the hotels in Bangkok, Thailand. The processes started from collected data and pre-processing is cleaned data by removal stop words and using the high frequency of word which will be selected into attribute for using classifier model. The classifier model will be solving the text of customer review that is positive of negative from training data and test data which are train from behaviour posting from customer of hotel service group. The proposed methodology is detailed as follows, Pre-processing.
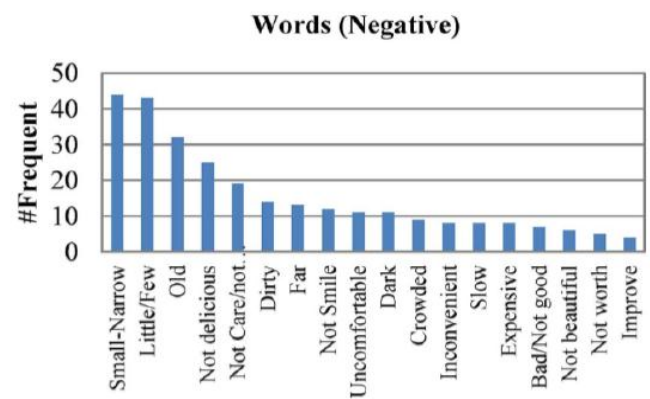
The feature selection is to be attributes in classifier that will be extracted words from these customer reviews as words occurred frequently to 36 words. There are positive and negative in Table I, which are ordered by descending frequent.

**TABLE I.** FEATUE SELECTION FROM REQUENT WORDS

| No. | Words (Positive) | #Frequent | Words (Negative) | #Frequent |
|-----|------------------|-----------|------------------|-----------|
| 1. | Convenient | 245 | Small/N arrow | 44 |
| 2. | Good | 206 | Little/Few | 43 |
| 3. | Near | 142 | Old | 32 |
| 4. | Clean | 140 | Not delicious | 25 |
| 5. | Comfortable | 62 | Not Care/not impression | 19 |
| 6. | Very Good | 59 | Dirty | 14 |
| 7. | Take care | 33 | Far | 13 |
| 8. | New | 32 | Not Smile | 12 |
| 9. | Smile | 29 | Uncomfortable | 11 |
| 10. | Big/Wide | 26 | Dark | 11 |
| 11. | Delicious | 25 | Crowded | 9 |
| 12. | Cheap/not expensive | 19 | Inconvenient | 8 |
| 13. | Much/Many | 17 | Slow | 8 |
| 14. | Safe | 17 | Expensive | 8 |
| 15. | Quiet | 16 | Bad/Not good | 7 |
| 16. | Worth | 12 | Not beautiful | 6 |
| 17. | Beautiful/Luxurious | 9 | Not worth | 5 |
| 18. | Fast/Quick | 6 | Improve | 4 |



**Fig. 2**. Frequent word of positive opinions



**Fig. 3.** Frequent word of negative opinions

The frequent words of positive are analysed for attribute transformation individual text of customer review. The training and test data are separated into 3 sets: set 1 is composed 5 positive and 5 negative words; set is 2 is composed of 10 positive and 10 negative words and set 3 is composed of all positive and negative words in Table II as follows,

**TABLE II.** DATA SETS FOR CLASSIFIER MODELS

| Datasets | Words |
|----------|-------|
| Set(10words) | Positive: convenient, good, near, clean, comfortable |
| | Negative: small/narrow, little/few, old, not delicious, not care/not impression |
| Set2(20words) | Positive: convenient, good, near, clean, comfortable, very good, take care, new, smile, big/wide |
| | Negative: small/narrow, little/few, old, not delicious, not care/not impression, dirty, far, not smile, uncomfortable, dark |

| Data sets | Words |
|---|---|
| Set3(36 words) | Positive: convenient, good, near, clean, comfortable, very good, take care, new, smile, big/wide, delicious, cheap/not expensive, much/many, safe, quiet, worth, beautiful/luxurious, fast/quick |
|  | Negative: small/narrow, little/few, old, not delicious, not care/not impression, dirty, far, not smile, uncomfortable, dark, crowded, inconvenient, slow, expensive, bad, not good, not beautiful, not worth, improve |

B. Model Construction

From data sets lead to model construction. The classifier models are used 2 models which are decision Tree (C4.5) and naive Bayes to classify texts as class labels: positive or negative. Each data set is trained to model and test model that given predicted class labels follows probability trending of classifier model. The classifier models are described as bellows,

• Decision Tree (C4.5)

The decision tree learning was proposed as a model of data classification for a class label, which called ID3 and developed to C4.5. In addition, decision tree is clearly represented through a tree diagram. It starts from the first node is a root node. The root node selects an attribute as words in opinion from the best value of measurement. Each attribute has its own values i.e. true/false, which are separated by branch links composed of original attributes. At the end, the data reveals a class which represents a leaf node (i.e. positive/negative).

The advantage of the decision tree is for ordering attributes that are the best measurement as After the distinguished information of attribute is calculated, the entropy value is also calculated to define the summary of each branch needed be clearly separated from attribute. The highest gained value of the attribute A results in the best attribute to classify data set which is calculated and range between 0 and 1.

• Naive Bayes

Naive Bayes is an algorithm of probability based on Bayes theorem of learning. It aims to create a model in the form of probability. The advantage of naive Bayes is an effective method which is easy processing. The probability of the classification data with prior knowledge is denoted by P (ail Vj), where ai refers to the attribute I and Vj refers to class label therefore, the classification has been calculated for this

probability. The highest probability of ai is depended on Vj foreach class is trend to answer of classification. The range of probability is between 0 and 1.

## IV. EXPERIMENTAL RESULTS

The experimental results retested with open opinion texts from 400 customer reviews from a website of hotel agent service. The results are compared percentage of accuracy between decision tree model (C4.5) and naive Bayes [16] and difference the number of features are extracted as 10,20 and 36 words respectively. The accuracy of naive Bayes is given values that are higher than decision tree all of data sets. Moreover, the highest of accuracy value is 94.37% with 20 words and also average of naive Bayes is higher than decision tree to 93.61% in table.
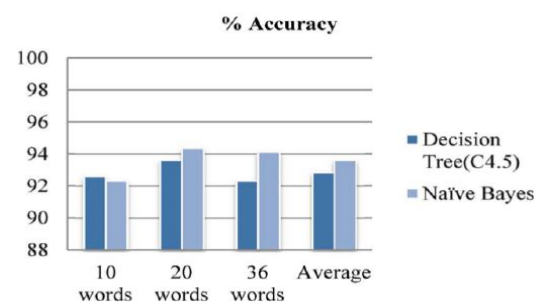


Fig. 4. Comparison of decision tree and naive Bayes

However, the advantage of decision trees model shown structure of words has related and priority following the entropy values. For example, of decision tree with 10 words training data, the hotels should take care of customer, location is near tourist attraction, convenient in room are ready, therefore, the review is trend to good

and positive. The words relationships able to translate to IFTHEN rules as follows,

Rule1: IF not care= false THEN Positive

Rule2: IF not care= true and near= true THEN Positive

Rule3: IF not care= true and near= false and convenient=true THEN Positive

Rule4: IF not care = true and near = false and convenient=false and good= true THEN Positive

Rule5: IF not care = true and near = false and convenient=false and good= false THEN Negative

The decision tree with 20-word training data shown the first service is dirty and lower level are far which is related to clean and good, moreover, not care word is related to smile and far again. These keywords are translated into rules, for example, Rule3: If customer complains far but have other good convenient, customer still gives positive score. And Rule4: If customer complain far but do not have any good convenient, customer still give negative score. Moreover, Rule 5-8: dirty room is first factor to decide of negative score. All relationship of word has IFTHEN rules as follows,

Rule 1: IF dirty = false and far = false THEN Positive

Rule2: IF dirty = false and far = true and clean = true THEN Positive

Rule3: IF dirty = false and far = true and clean = false and good= true THEN Positive

Rule4: IF dirty = false and far = true and clean = false and good= false THEN Negative

Rule5: IF dirty = true and not care= true THEN Negative

Rule6: IF dirty = true and not care= false and smile = true then Negative

Rule7: IF dirty = true and not care= false and smile = false and far = true then Negative

Rule8: IF dirty = true and not care= false and smile = false and far = false and uncomfortable= true then Negative

Rule9: IF dirty = true and not care = false and smile = false and far = false and uncomfortable= false then Positive.

The decision tree with all word training data shows word relationships such as dirty, far, not care, not good, many/much, smile, expensive, good, far, near, uncomfortable in form of tree. In these experimental results show that some words have affected to class label. For example, Rule 5, even if customer review in text as

expensive and near, customer still has opinion as positive, whereas, Rule6 has expensive and good word in customer review, customer given negative rating to service. All relationship of word has IF-THEN rules as follows,

Rule 1: IF dirty = false and far = true and many/much= false THEN Positive

Rule2: IF dirty = false and far = true and many/much= true THEN Negative

Rule3: IF dirty = false and far = false and not good= true and good= true THEN Positive

Rule3: IF dirty = false and far = false and not good= true and good= false THEN Negative

Rule4: IF dirty = false and far = false and not good= false and expensive= false THEN Positive

Rule5: IF dirty = false and far = false and not good= false and expensive= true and near= true THEN Positive

Rule6: IF dirty = false and far = false and not good= false and expensive= true and near = false and good = true THEN Negative

Rule7: IF dirty = false and far = false and not good= false and expensive= true and near = false and good = false THEN Positive

Rule8: IF dirty = true and not care= true THEN Negative

Rule9: IF dirty = true and not care = false and smile = true THEN Negative

Rule10: IF dirty = true and not care= false and smile = false and far = true THEN Negative

Rule11: IF dirty = true and not care= false and smile = falls and far = false and uncomfortable= true THEN Negative

Rule12: IF dirty = true and not care = false and smile = false and far = false and uncomfortable = false THEN Positive
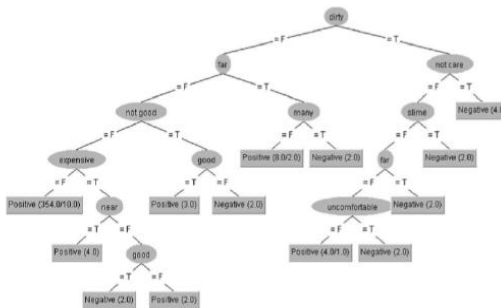


**Fig. 5**. Decision tree from training data (36 words)

However, the rating generating is testing by naive Bayes by probability trend to predict class label in Table IV. The table IV shows RMSE of different data sets. The lowest of RMSE IS 36 words testing data that give rating that are similar to actual score from customer review to 0.2326. The rating of 20 words and 10 words are slightly higher value than 30 words to 0.2390 and 0.3669 respectively. The average of naive Bayes model generates rating value that is similar actual rating as 0.2792 and median as 0.2390.

TABLE II. ROOT MEAN SQUARE ERROR OF NAIVE BAYES

| Attributes | Root Mean Square Error (RMSE) |
|---|---|
| 10 words | 0.3660 |
| 20 words | 0.2390 |
| 36 words | 0.2326 |
| Average | 0.2792 |

## V. CONCLUSION

The sentiment analysis with opinion mining framework reported in this paper can be incorporated into a hotel technology system that can help improve customer relationship management. What good is a system that predicts the polarity of sentiments if it works with the wrongly labelled data? From the sentiment polarity exercise that we did, we found out that some comments may be wrongly viewed as neutral while they will be either positive or negative. The following example was viewed as a neutral comment. "That hotel is surely a HELLTEL!" This comment is truly negative and sarcastic, but because the word HELLTEL does not exist in the English vocabulary it was classified under the neutral class. However, most comments were labelled with a much better accuracy. We believe that a lot of research can be done in this area especially in fine tuning the feature extraction algorithm of the framework so that classification error is minimized. The system is expected to determine sentiments the way human beings do and labelled datasets are normally used for the system to learn automatically. The proposed framework tries to make sure that sentences are correctly labelled, the proposed framework in this paper helps in automatic labelling of sentiment datasets. The primary constraint with this point of view level execution is that it is space specific. Regardless, simply little changes (in context vectors) would be required to use this neural Associations Rule in another region result administering most hoisted review on the hotel name with exactness and we are using machine figuring it is separating the unlabelled data and indicating unlabelled data with outline.

## REFERENCES

[1] S. 1. Wu, R.D. Chiang and Z.H. Ji, Development of a Chinese opinion mining system for application to Internet online forum, The Journal of Supercomputing, Springer US[Online], 2016.

[2] Z. Li, L.Liu and C.Li, Analysis of customer satisfaction from Chinese reviews using opinion mining, Proceeding of the 6th IEEE International Conference on Software Engineering and Service Science(ICSESS).
2015, pp.95-99.

[3] Q.Su, X.Xu, H.Guo, Z.Guo, X. Wu, X. Zhang and B.Swen. Hidden Sentiment association in Chinese web opinion mining. Proceeding of the 17th International Conference on World Wide Web, 2008, pp.959-968.

[4] S.Atia and K. Shaalan, Increasing the accuracy of opinion mining in Arabic. Proceeding of the 1st

International conference on Arabic computing linguistics, 2015, pp.l 06-113.

[5] R.M. Duwairi and I. Qarqaz, Arabic Sentiment Analysis using Supervised Classification. Proceeding of 2014 International Conference on Future Internet of Things and Cloud. 2014, pp. 579-583.

[6] H.S. Le, T.V. Le and T.V. Pham, Aspect Analysis for Opinion Mining of Vietnamese Text. Proceeding of International Conference on Advance Computing and Application, 2015, pp.118-123.

[7] T. Chumwatana, Using sentiment analysis technique for analyzing Thai customer satisfaction from social media. Proceeding of the 5th International Conference on Computing and Informatics, 2015, pp.659-664.

[8] S.Ahmed and A.Danti, A novel Approach for sentimental analysis and opinion mining based on sentiwordnet using web data. Proceeding of International Conference on Trends in Automation, Communications and Computing Technology, 2015, pp.1-5.

[9] R.K. Bakshi, N. Kaur, R. Kaur and G.Kaur, Opinion mining and sentiment analaysis, Proceeding of the 3rd International Conference on Computing for Sustainable Global Development, 2016, pp. 452-455.

[10] P.Barnaghim, 1.G. Breslin and P. Ghaffari, Opinion mining and sentiment polarity on Twiiter and correlation between events and sentiment, Proceeding of the 2nd International Conference on Big Data Computing Service and Application, 2016, pp. 52-57.

[11] N. Kumari and S. N. Singh, Sentiment analysis on E-commerce application by using opinion mining, Proceeding of the 6th International Conference-Cloud System and Big Data Engineering (Confluence),
2016, pp. 320-325.

[12] V.B. Raut and D.D. Londhe, "Survey on opinion mining and summarization of user review on web", International Journal of Computer Science and Information Technology, Vol. 5(2), 2014, pp.1026-1030.

[13] 1. Fiaidhi, O. Mohammed, S. Mohammed, S. Fong, and T.H, Kim, Opinion Mining over twiiterspace: Classifying tweets programmatically using the R approach. Proceeding of the 7th International Conference on Digital Information Management, 2012, pp. 313-319.

[14] M. R. Islam and Minhaz F. Zibran, "Exploration and Exploitation of Developers' Sentimental Variations in Software Engineering", Internation Journal of Software Innovation, Vo1.4(4), 2016, pp.35-55.

[15] Y.Yokoyama, T. Hochin and H. Nomiya, "Estimation of Factor Scores of Impressions of Question and Answer Statements", International Journal of Software Innovation, Vol. 1(4), 2013, pp.53-66.

[16] L. Lin, 1. Li, R. Zhang, W. Yu and C. Sun, Opinion mining and sentiment analysis in social networks: A retweeting structure-aware approach. Proceeding of the 7th International Conference on Utility and Cloud Computing, 2014, pp.890-895.

[17] A.H. Al-hamaami and S. H. Shahrour, Development of an opinion blog mining system, Proceeding of the 4th International Conference