

Mining Implicit and Explicit Rules for Customer Data Using Natural Language Processing and Apriori Algorithm

T. Velmurugan¹, B.Hemalatha²

¹Associate Professor, PG and Research Department of Computer Science, D.G. Vaishnav College, Chennai, India.

²Assistant Professor, Department of Computer Applications, Margregorios College, University of Madras, Chennai.

velmurugan_dgvc@yahoo.co.in¹, hemalathamgc11@gmail.com²

Abstract:

The Internet is a massive storehouse of organized and amorphous information. It is not an easy task to examine all these information to draw secret public views and emotions. It's extremely complicated to identify connections with large databases. There are several obsolete and incomplete information files in these archives that are not necessary to retrieve the laws. Therefore, the reliability of the organization rules is significantly affected by these irrelevant data and there is a need to pre-process these documents. In a rapidly growing technology world, there are many customers express their views through online, so organizations are highly dependent on the user's opinion. Natural Language Processing is the branch of machine learning which is about analyzing any text and used to handle predictive analysis. It can assist and interprets to profound the sentence structure in its meaning. A transaction is just a set of items that a customer purchases in a basket. To examine the connection between the items sold in a supermarket, Apriori algorithm is used to identify frequent sets of items that are explicitly bought together. Implicit relationship is ignored in Apriori algorithm which can be identified using sentiment analysis. Sentiment Analysis helps to identify the object and topic from the text to which the feeling is guided. This proposed analysis may assist to identify the implicit product in order to improve the sales by providing offers for respective implicit products.

Keywords: *Natural Language Processing(NLP), Apriori Algorithm, Sentiment Analysis, Implicit products.*

1. Introduction

Data mining is becoming increasingly essential with huge volume of information stored in databases, records and additional sources. If not essential for extracting the interesting information to create strong means of analyzing and possibly interpreting such information might assist in decision making [1]. The process of Knowledge Discovery in Databases (KDD) is the fact component of data mining [2]. In the database, process of Knowledge Discovery includes a few steps to obtain some sort of new knowledge from raw information collection[3]. In data mining, association is well-known approach. Association implies that associated objects which are grouped from a collection. An easy instance for the purpose of discovering association rules is to analyze a big supermarket transaction database [4] known as evaluation of market baskets or mining rules of association. Depending on an association of a specific item in the same transaction on other items, a pattern is revealed in association. For instance, in analysis of market basket, the association technique is used to recognize the products buy together often by the customers. It is vital to find bigger baskets in the trade professional, which can deal with 1000's of items [5]. Depending on this data, enterprises can advertise more products which can have a marketing campaign equivalent to generate a growth in income. This kind of finding helps companies to make definite verdicts namely cross-marketing, catalog design and behavioral analysis of customer shopping.

NLP is the branch of Machine learning and extensively characterized as the programmed control of natural language which is concerned with connection between human language and computers in the areas of computer science and artificial intelligence. NLP is used to analyze massive amount of text and handle predictive analysis. The technique of NLP includes chunking data, stemming and removal of stop words. NLP is useful in review rating for dividing the words, sentences, noun and paragraphs. It determines whether the sentence is Positive and Negative or neutral. NLP may also be useful as translator in case translation of one language to required language. This may produce less noise and leads to a robust data. NLP can be broadly classified into Natural Language Understanding and Natural Language Generation. Techniques used in NLP understanding are Syntax and Semantics. Natural Language Generation includes utilizing repositories to extract and translate linguistic desires into human language. In analyzing online customer reviews, NLP plays an important role for extracting the actual meaning of the sentence as well as customer view about the product.

To discover meaningful pattern and rules there are many datamining techniques and algorithms given by Saurkar et al.,[6] . The main goal of association is to establish the relationship between items which exist in the market. The rule of Association Mining is to determine association rules that fulfill from a specified database by pre-determined minimum number of support and confidence. Usually the issue is broken down into two more sub-problems. First step is to determine those set of items which exceed a predefined limit whose events in the catalog are known as large set of items or frequent set of items. Second step is to create association rules from those large set of items with minimum confidence limitations. The fundamental job for association based rules in mining is to define the connection between transactional database items. Apriori algorithm is used to find the all frequent set of items and produce powerful association laws i.e. for explicit products. Our proposed work concentrates on implicit products as well for better business opportunities.

The remainder of paper is organized as given below: Section 2 presents current knowledge including theoretical and methodological contributions from existing work to proposed work. Section 3 outlines Natural Language Processing implementation details. Section 4 presents Apriori algorithm and Association Rule mining for products. Section 5 presents NLP Vectorization of Reviews. Section 6 concludes the proposed work.

2. Literature Review

A number of articles utilized for the market basket analysis and sentiment analysis. Some of these are discussed here. Manpreet Kaur et al.,[7] proposed that Finding out the patterns due to changes in data is in itself an interesting area to be explored in market basket analysis and for other areas. Spirtes et al.,presented a constraint-based approach to causal discovery, which relies on the conditional independence relationships in the data in [8]. Peska et al.,[9] presented a novel approach to use a specific user behavior pattern as implicit feedback, forming binary relations between objects i.e. analyzing rule relations though it is not insignificant to capture the implicit relations and even to make the actionable well-known rules. Chen et al., [10] propose a comprehensive semantic similarity measureprovides the implicit relations which do not co-occur frequently by referring to the connections among various items but having a maximum probability occurred with identical items of third party. Jinturkar et al.,[11] proposed aframework of map-reduce by the classification of customer reviews for products through which results can be presented in a convenient visual form for the nontechnical user for making better decisions. An application of data mining techniques to a selected business organization with special reference to buying behavior is given by Hilage et al.,[12]. The results of three techniques namely association rule mining technique, rule induction technique and apriori algorithm were combined and efforts were made to understand the correct buying behavior of the customer.

Abulei et al.,[13] has utilized NLP techniques to generate some rules that helps to understand customer opinions and reviews (textual comments) written in the Arabic language. NLP is used for the purpose of understanding each one of them and then converting them to a structured data. Bhargav et

al.,[14] consumer's sentiments are reflected in the form of 'opinion dataset' on internet using sentiment analysis for hotel reviews. Bhatt et al.,[15] proposed a system by finding sentiment of the reviews that performs the classification of customer reviews. Mohan et al., [16] proposed text mining techniques as well as sentiment analysis by analyzing the customer reviews on the restaurant domain. Moreover, priority-based algorithm is used for predicting the reviews polarity which created the rule base to the classifier. Hamano and sato[17] proposed to analyze the targeting competitors and customers by the framework to mine the indirect ARM. Many researchers have proposed indirect association mining and realizing the implicit significance relationship among items [18].Therefore, ARM [19] has only accomplished to the frequent items whereas the infrequent items get ignored.

The concept based on sentiment analysis for developing a system by Zhang et al. [20] has identified the product weakness. Therefore, the producers of product required to improve the product quality and at each situation the system attempted to identify the features of explicit and implicit which gets accurate sentiment words, the sentiment analysis based on sentences. The proposed model has been extracted based on the concept of product customer review whereas both noun and noun phrases have extracted from every review sentences [21]. The rule based technique has deal with discovery of sentiment words and orientation of text [22]. There have been many changes to the Apriori algorithm to maximize its performance and efficiency [23]. Wu et.al suggested another Improved Apriori Algorithm (IAA) to reduce the number of information scans and repetitive operations thus generating regular item sets and association laws. Ayman E. Khedr et al.,[24] proposed a model to improve the prediction accuracy for the future trend of stock market, by considering different types of daily news with different values of numeric attributes during a day using sentiment analysis. The process of analyzing about specific products by customer opinion and their characteristics is known as sentiment analysis [25].

3. Natural Language Processing for Review Data

Syntax relates to the structure of words in a phrase to be grammatical. In NLP, analysis based on syntactic is used to evaluate the alignment of the NL with the grammar guidelines. It can be used to apply and derive significance from linguistic regulations to a group of words. Some syntax methods are Lemmatization, Morphological division, Word segmentation, Part-of-speech tagging, Parsing, Sentence breaking and stemming. Semantics refers to a text's expressed context. The interpretation is one of NLP difficult aspects that have not yet been entirely determined. This contains the use of computer algorithms to comprehend the meaning and description of terms in the form of sentences. Semantic analysis methods are: Named entity recognition (NER), Word sense disambiguation. In this proposed work Natural Language Tool Kit (NLTK) is used for preprocessing of review data and it is implemented in Python.

Tokenizing: For consistency in storage, evaluation and identification, it is important to break down the data that is crawled on the internet. Every analysis of the review is collected as a passage. This is broken into individual sentences and is contained in tables. Therefore, in order to infer the significance, possibly, it is necessary to break each sentence into words also to derive the characteristics of the item that are spoken. In this proposed work,onlinecustomer feedback is given as input for the NLP system. A main component in the NLP system Tokenizer splits the unstructured feedback into tokens. Tokens are the part of a sequence of characters that are combined together in a text. These tokens are the valuable semantic units for processing. Thetokenized includes words, punctuation marks, symbols etc., which can convert a sentence into word level tokens.After importing customer feedback andextracting the required term from tokenizing, generation of required relationship using NLP is done. In this process, the following algorithm may used to splitting the review into sentences.

```
# Import and Segment data
a = []
b = []
f = open ('review_text.text', 'rt');
```

```
for line in f
line = (line.strip())
b.append (line)
sents = tokenize.sent_tokenize(line)
#print (sents)
a.append (sents)
data []
fori in a;
for j in I;
#print (j)
Data.append(j)
# for j in a;
# print (j)
# data.
```

Initially part of speech information is utilized in all NLP tasks to identify adjectives, nouns and root for each word in the text. In the following, we demonstrate how our approach converts customer feedback reviews from unstructured text to a structured data.

POS Tagging: It is the act of marking or identifying each term in a sentence with its POS as in a verb, adverb, noun, adjectives etc., i.e. it labels each word with its POS.

Sentiment sentences extraction and POS tagging:It is recommended to delete any critical material of examination of emotions. Instead of suppressing factual information, all subjective material for potential review has been removed from our report. The emotional content is made up of all sentences with feelings. A sentence of emotion is the one with at least one positive or negative term. First of all, all sentences are tokenized into individual words in English. That word in a sentence has its syntactic function in deciding how the word is used. Often recognized as the POS are the syntactic functions. Throughout English there are 8 pieces of speech: the verb, noun, conjunction, adjective, pronoun, interjection, adverb and preposition. POS taggers were built in NLP to classify words based on their parts of speech. A POS tagger is very useful for sentiment analysis due to the following two reasons:

- Terms such as nouns and pronouns do not typically have any thoughts. With the aid of a POS tagger, it can strip out such terms;
- A POS tagger can also be used to distinguish phrases that can be used in various parts of the expression. For example, as a noun, "enhanced" may execute a number of emotions as an adjective. The POS tagger used for this work is an established max-entropy POS tagger. The tagger can provide 46 separate tags that suggest that more complex syntactic positions can be defined than just 8. For example, Table 7 is a list of all tags found in the POS tagger for verbs. Using the POS tagger every sentence was then labeled. Given the huge amount of sentences, a parallel running Python program was written to increase the tagging speed.

Table 1:POS tags for verbs

Tag	Definition
VB	base form
VBP	present tense, not 3rd person singular
VBZ	present tense, 3rd person singular
VBD	past tense
VBG	present participle
VBN	past participle

Algorithm for Tokenizing and generating the required relationship

```
# generating tree with relationship
Def process Content ();
    Try;
    For item in data
    Tokenized = nltk.word_tokenize(item)
    Tagged = nltk.pos_tag(tokenized)
    Chunkgram = “”” Chunk: {,RB.?? *<VB.??*<NNP>} “””
    Chunkparser = nltk.RegexpParser(chunkgram)
    Chunked = chunkparser.parse(tagged)
    Print(chunked)
    Chunked.draw()
    Except Exception use:
    Print (str(e))
    Process content()
```

This section finds the sentiment in a customer feedback namely adjectives, nouns and adverbs which are used as a features that depict with higher accuracy. NLP techniques chunking data, stemming and removal of stop words are used to get sentiment words. These sentiment words acts as a major parameter to identify implicit products.

4. Mining Association Rules Using Apriori Algorithm

In this stage, the main goal is to use Apriori algorithm to govern the mining organization, which has been accomplished by using highly effective pruning methods and techniques to make significant progress. The raising, help and belief structure to create rules and even targeted to explicit occurrence as a result which even capture the relationship of both implicit and explicit based on ARM. Each column indicates an item while the customer buy number of products with the combination of product related to product ID which represents an item has been occurred in the related transaction.

Table 2: Instance of the stored data

	0	1	2	3	4	5	6	7	8	9	10	11
0	shrimp	almonds	Avocado	Vegetables mix	Green grapes	Whole wheat flour	yams	Cottage cheese	Energy drink	Tomato juice	Low fat yogurt	Green tea
1	burgers	meatballs	Eggs	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN
2	Chutney	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN

3	Turkey	avocado	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN
4	Mineral water	Milk	Energy bar	Whole wheat rice	Green tea	NaN						

The data preprocessing with the Apriori library which is shown in table 2 is used as dataset to manipulate the lift, support and confidence for the respective product available in the stored dataset. Therefore, the whole database is a large list and each operation in the dataset is an internal list within the broad outer list.

Apriori Algorithm: Apriori is an algorithm used to identify frequent sets of items (in our case pairs of items). This is done using a "bottom-up" method, first categorizing individual items that meet a minimum threshold for occurrence. It then expands the collection of items, inserting one item at a time and testing whether the set of items still exceeds the defined limit. The mining association rules problem can be described as follows: Let $I = \{ i_1, i_2, \dots, i_m \}$ be a set of items. Let $T = (t_1, t_2, \dots, t_n)$ be a set of transactions (database), in which each transaction t_i is a set of items such as I . A law of association is an interpretation of the form, $X \rightarrow Y$, where $X \subset I$, $Y \subset I$, and $X \cap Y = \emptyset$. X (or Y) is an itemset, called an itemset. In our proposed work, the connection between the items sold in a supermarket and the set of items in the grocery store is examined. A transaction is just a set of items that a customer purchases in a basket. Table.2 shows the item description for transaction with Product ID for the stored dataset.

Table 3: Product ID with corresponding item description

Product ID	Item description
B00120NH	Burger
B001E0556	Meatballs
B00291ESH	Egg
B004749TK5	Turkey
B000F542TU	Frozen vegetables
B00H66YWR	Cookies
B0034TYYZ	Whole wheat pasta

Each and every transaction t_i is a collection of products that a consumer has bought in a basket in a shop. The collection of all items sold at the shop. A simplistic view of shopping baskets is the information representation in Table 4 transaction type. For e.g., there is no consideration of the quantity and value of each product in the template.

Table 4: Customer transaction details in Supermarket

Customer Transaction	Item description	Product ID
1	Burgers, meat balls, egg	B00120NH, B001E0556, B00291ESH
2	Turkey, egg	B004749TK5, B00291ESH
3	Burger, meat balls, turkey	B00120NH, B001E0556, B004749TK5
4	Burger, meat balls	B00120NH, B001E0556
5	Burger, turkey	B001E0556, B004749TK5

Table 5 illustrates occurrences i.e. the number of times a particular item occurs in the transactions in the iteration 1.

Table 5: Number of times each item occurs

Item set	Occurrence count
Burger	4
Meat balls	3
Egg	2
Turkey	2

Iteration 2: Build item sets of size 2 using the remaining items from Iteration 1 (ie: Burger, Meat balls). Only {Burger, Meat balls} remains and the algorithm stops since there are no more items to add as shown in Table.6.

Table 6: Build item sets of size

Item Set	occurrence count
{ Burger, Meat balls }	3

If we had more orders and items, we can continue to iterate, building item sets consisting of more than 2 elements. For the problem we are trying to solve (i.e.: finding relationships between pairs of items), it suffices to implement Apriori to get to item sets of size 2.

The literature has documented a large number of ARM algorithms that have specific mining efficiencies. Nonetheless, the resultant sets of rules are all the same based on the definition of the principles of association. That is, the collection of relationship rules remaining in T is calculated uniquely given a transaction data set T, a minimum help and a minimum trust. While their computing efficiencies and space needs may be specific, every algorithm will use the same set of rules. The Apriori Algorithm suggested is the best known mining algorithm.

Association Rules Mining: Once the item sets have been generated using Apriori, association rules can be mined. Given that we are only looking at item sets of size 2, the association rules we will generate will be of the form {A} -> {B}. One common application of these rules is in the domain of recommender systems, where customers who purchased item A are recommended item B. Here are three key metrics to remember when determining the laws of association:

Support: This is the proportion of orders comprising the product collection. In the above case, there are a maximum of 5 orders and in 3 of them {Burger, Meat balls} exists, so

$$\text{Support } \{ \text{Burger, Meat balls} \} = 3/5 \text{ or } 60\%$$

The minimum support threshold required by Apriori can be set based on knowledge of your domain. In this grocery dataset for example, since there could be thousands of distinct items and an order can contain only a small fraction of these items, setting the support threshold to 0.01% may be reasonable.

Confidence: Because of two items A and B, confidence tests the number of occasions that item B has been purchased when item A has been purchased. This is represented as:

$$\text{Confidence } \{ A \rightarrow B \} = \text{support } \{ A, B \} / \text{support } \{ A \}$$

Confidence values range from 0 to 1, where 0 indicates that B is never purchased when A is purchased, and 1 indicates that B is always purchased when A is purchased. Remember that the calculation of confidence is directional. It implies that, when item B is purchased, we can also measure the number of periods when item A is bought:

$$\text{Confidence } \{B \rightarrow A\} = \text{support } \{A, B\} / \text{support } \{B\}$$

In our case, when Burger has been purchased, the percentage of times the Meatballs purchased is:

$$\begin{aligned} \text{Confidence } \{ \text{Burger} \rightarrow \text{Meat balls} \} &= \text{support } \{ \text{Burger, Meat balls} \} / \text{support } \{ \text{Burger} \} \\ &= (3/5) / (4/5) \\ &= 0.75 \text{ or } 75\% \end{aligned}$$

A confidence value of 0.75 implies that out of all orders that contain Burger, 75% of them also contain Meat balls. Now, the confidence measure in the opposite direction (i.e. Meat balls->Burger) can be measured.

$$\begin{aligned} \text{Confidence } \{ \text{Meat balls} \rightarrow \text{Burger} \} &= \text{support} \{ \text{Burger, Meat balls} \} / \text{support} \{ \text{Meat balls} \} \\ &= (3/5) / (3/5) \\ &= 1 \text{ or } 100\% \end{aligned}$$

We can observe that all the Meat balls-containing orders often include Burger. Nevertheless, does this imply that there is a connection between these two objects, or in the same orders they exist together.

Lift: In the case of two objects, A and B, lift shows whether there is a correlation between A and B, or whether the two items purely by accident appear together in the same orders (i.e. randomly). Unlike the trust measure, the meaning of which may change depending on the direction (e.g. confidence {A->B} that varies from confidence {B->A}), there is no direction of rising. It ensures the {A, B} lift is always the same as the {B, A} lift.

$$\text{Lift } \{A, B\} = \text{Lift} \{B, A\} = \text{support} \{A, B\} / (\text{support} \{A\} * \text{support} \{B\})$$

In our proposed work, we compute lift as follows:

$$\begin{aligned} \text{Lift } \{ \text{Burger, Meatballs} \} &= \text{support } \{ \text{Burger, Meatballs} \} / (\text{support } \{ \text{Burger} \} * \text{support} \{ \text{Meat balls} \}) \\ &= (3/5) / (4/5 * 3/5) \\ &= 1.25 \end{aligned}$$

One way to understand lift is to think of the denominator as the likelihood that A and B will appear in the same order if there was no relationship between them. In the example above, if Burger occurred in 80% of the orders and Meat balls occurred in 60% of the orders, if there were no partnership between them, we would anticipate both of them to show up together in the same order 48% of the time (i.e.: 80% * 60%).

In order to create a dictionary with item as a key, the list consisting of support, confidence and lift as value. For instance: Relation Record(items = frozenset ({'Burger', 'Meat balls'}, support = 0.01, ordered_statistics = [ordered statistic (items_base = frozenset({'Meat balls'}), items_add= frozenset ({'Burger'})), confidence= 0.75, lift=1.25.

Table.7: The lift, support and confidence for the products in stored dataset

Item1	Item2	Support	Confidence	Lift
Burger	Meat balls	0.01	0.75	1.25
Shrimp	Pasta	0.005	0.322	4.5
Whole Wheat Pasta	Olive oil	0.007	0.27	4.12

Here the lift value is higher, which implies association is stronger. This becomes better rule in predicting something than erratically guessing.

5. Sentiment Analysis for Review Data

Analysis of emotion is a sort of assessment of text which classifies texts depending on the emotional inclination of the opinions they comprise. SA of product reviews in text mining and computational linguistics study has recently become very common. In this proposed work, SA is used to analyze the product feedback given by the customer in the website.

Calculation of word count: In order to calculate the word count, extract all the tokenized review words of customer feedback associated with the product ID which has been arrived through Apriori algorithm for training using deep learning. In order to train the review words, Document Frequency (DF) or Inverse Document Frequency (IDF) can be utilized for determining the word count.

Algorithm for tokenized review words of customer feedback

```
df['reviewText'] = df['reviewText'].fillna("")
#countVectorizer() converts a collection of textdocument
Vectorize = CountVectorizer()
#assign a shorter name for the analyze
#which tokenizes the string
Analyzer = vectorizer.build_analyzer()
defwordcounts(s)
    c = { }
#tokenize the string and continue, if it is not empty
If analyzer(s):
    d = { }
#find counts of the vocabularies and transform to array
w = vectorize.fit_transform([s]).toarray()
#vocabulary and index(index of w)
vc = vectorize.vocabulary_
#items() transforms the dictionary's(word, index) tuple pairs
For k,v in vc.items():
    d[v] = k #d ->index:word
for index, I in enumerate(w[0]):
    c[d[index]] = I # c ->word:count
return c
# add new column to the dataframe
df['word counts'] = df['reviewText'].apply(wordcounts)
df.head()
```

First, all the words in the reviews should be tokenized by fitting Tokenizer class on the data set. These reviewed word count has been vectorized in order to form a selected sentimental words to ensure consistency of the data. After vectorization mapping list of words (tokens) to a list of unique integers using texts_to_sequences class is done. The Frequently used sentimentwords are grouped and given in Table.8.

Table.8: Selected words with Vectorized word count

Selected words	Vectorized word count
Great	46425
Love	32971
Bad	10579
Wonderful	6710
Happy	5406

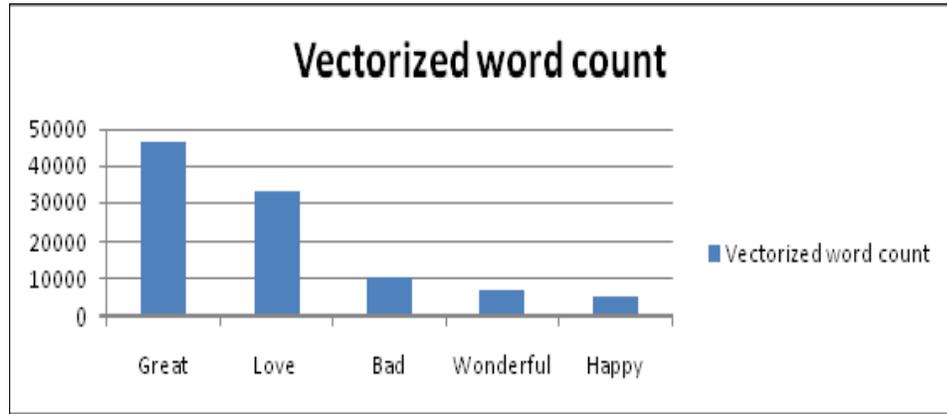


Figure 1: Vectorized Word Count

Figure 1 describes the vectorized word count of selected words in the reviews corresponding to the product_id. Using these vectorized word count it is possible to map customer satisfaction about the product purchased. Figure 2 describes the vectorized words, product_id and frequently used words in the reviews.

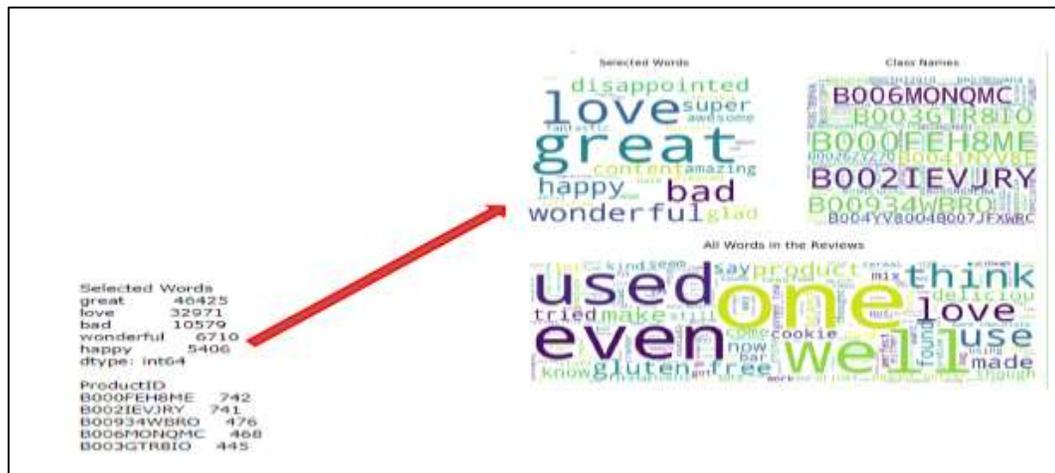


Figure 2: Vectorized Word Count

The above figure describes great is used frequently in more reviews i.e., 46425 times which shows customer satisfaction in buying the products. It also shows Product_id B000FEH8ME occurs 742 times i.e. frequently bought product by the customer.

These are the sentiment words which describes the polarity of the review for our stored dataset. Frequency rating can be done using these words along with product_id. This kind of rating suggests the quality of the product rather than probability or guessing randomly. Depending upon the polarity of the sentence we can also provide different polarity flavors by identifying whether positive or negative associated with a particular feeling, such as anger, sadness, or worry (i.e. negative feelings) or happiness,

love, or enthusiasm (i.e. positive feelings). In this research work, the frequency of sentimental words with their corresponding rating scale is rated from 1 to 5. The following categories of rating levels namely very positive, positive, neutral, negative and very negative. For example, mapped onto 5 star rating in a review very positive = 5 star and very negative = 1 star.

Table 9: Classifying the reviews based on SA

Product_ID	Average Rating
B00HBBYWNW	4
B0047479TA	5
B000F5429A	4
B00B9FWE4A	4
B008UKITG	3
B0078XBN6	3
B00291ES6W	3
9742356831	5
B00EDG3LS	2
B00SARK04	2
B003D41YYY	3
B009F3SC8	1

Frequency review rating helps the online customer to buy the products quickly rather than spending time in reading all the reviews related to that product. It also helps vendors in identifying the popularly purchased product by the customers.

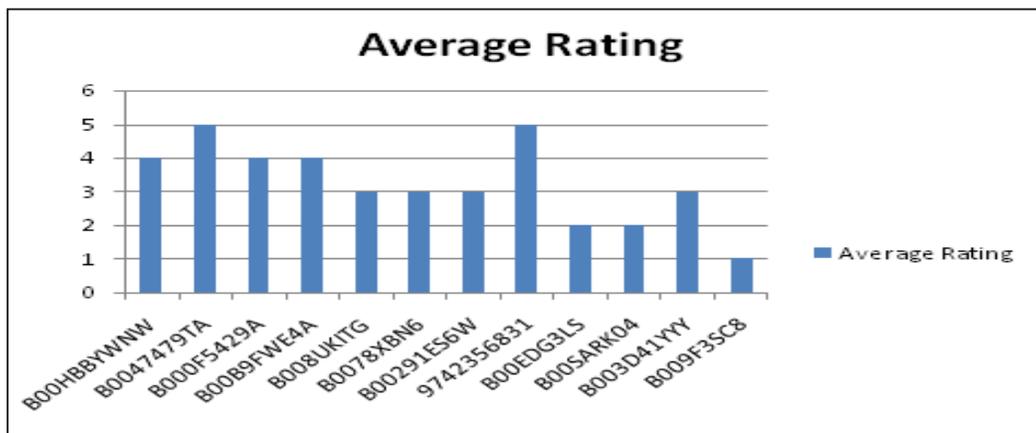


Figure 3: Frequency Review Rating

Figure 3 describes the frequency review ratings for the products from the reviews. Result shows that two frequently bought products completely satisfy the online customers with the rating of 5. Frequency rating

obtained from sentiment analysis extracts the true information about the product which may help explicitly for the customers in buying and implicitly the vendors for their business expansions or profits.

6. Conclusion

Many businesses claim that their business success depends solely on customer satisfaction. Therefore, scientists and educators were encouraged to find better association among purchased products to improve business. The customer before paying the money always prefers to read reviews of the service provider but it is not possible to read all the feedback from the website given by the customers. Every time new information will be provided by every review of the product or feature of the product, there is probability of missing any important feedback given by the customer. In order to overcome the above difficulty there is a need to identify the frequency of review rating. This proposed work concentrates on frequency review rating of products which can determine both implicit and explicit relationship among products. Using NLP, customer reviews can be structured in the form of tokens with pos tagging, word count can be calculated. Association rules can be formed using Apriori Algorithm for products. These rules and word count can be vectorized as sentiment words, which can be used for frequency rating. Our work determines not only frequent itemset in the basket but also implicit product which is less frequent satisfies customer better, which is more important in ecommerce for business expansion. Thus, this proposed work can be recommended for implicit and explicit relationship rule mining among products which helps customer to make quick decision in choosing products as well as in improving business needs.

References

1. Usama Fayyad, Gregory Piatetsky, Shapiro and Padhraic Smyth, “The KDD Process for Extracting Useful Knowledge from Volumes of Data”, *Communications of the ACM*, Vol. 39, No. 11, pp. 27-34, 1996.
2. Gurjit Kaur and Lolita Singh, “Data Mining: An Overview”, *International Journal of Computer Science and Telecommunications*, Vol. 2, No. 2, pp. 336-339, 2011.
3. Reena Hooda and Nasib S. Gill, “Applications and Issues of Data Mining”, *International Journal of Research in IT & Management*, Vol. 2, No. 3, pp. 11-17, 2012.
4. R. Surendiran, K.P. Rajan and M. Sathish Kumar, “Study on the Customer targeting using Association Rule Mining”, *International Journal on Computer Science and Engineering*, Vol. 2, No. 7, pp. 2483-2484, 2010.
5. Luis Caviq, “A Scalable Algorithm for the Market Basket Analysis”, *Journal of Retailing and Consumer Services*, Vol. 14, No. 6, pp. 400-407, 2007.
6. Saurkar Anand V, Bhujade V, Bhagat P, Khaparde A, “A Review Paper on various Data Mining Techniques”, *International Journal of Advanced Research in Computer Science and Software Engineering*, Vol. 4, No. 4, pp. 98-101, 2014.
7. Manpreet Kaura, Shivani Kanga, “Market Basket Analysis: Identify the changing trends of market data using association rule mining”, *International Conference on Computational Modeling and Security*, *Procedia Computer Science*, Vol. 85, pp. 78 – 85, 2016.
8. P. Spirtes and K. Zhang, “Causal discovery and inference: Concepts and recent methodological advances,” *Applied Informatics*, Springer, Vol. 3, No. 1, pp. 1–28, 2016.
9. L. Peska and P. Vojtas, “Using implicit preference relations to improve recommender systems,” *Journal on Data Semantics*, Vol. 6, No. 1, pp. 15–30, 2017.
10. Q. Chen, L. Hu, J. Xu, W. Liu, and L. Cao, “Document similarity analysis via involving both explicit and implicit semantic couplings”, *International Conference of Data Science and Advanced Analysis*, pp. 1–10, 2015.

11. Mugdha Jinturkar and Pradnya Gotmare, “Sentiment Analysis of Customer Review Data using Big Data: A Survey”, *International Journal of Computer Applications (0975 – 8887) Emerging Trends In Computing*, 2016.
12. Hilage Tejaswini A, Kulkarni RV, “ Application of data mining techniques to a selected business organization with special reference to buying behaviour ”, *International Journal of Database Management Systems*, Vol.3, No 4, pp.169-181, 2011.
13. Saleem Abuleil and Khalid Alsamara, “Using Nlp Approach For Analyzing Customer Reviews” , pp. 117– 124, 2017. © Cs & It-Cscp 2017 Doi : 10.5121/Csit.2017.70112.
14. P. Sanjay Bhargav, G. Nagarjuna Reddy, R.V. Ravi Chand, K.Pujitha, Anjali Mathur, “SA for Hotel Rating using Machine Learning Algorithms”, *International Journal of Innovative Technology and Exploring Engineering .ISSN: 2278-3075*, Vol.8 , No.6, 2019.
15. Aashutosh Bhatt, Ankit Patel, Harsh ChhedaandKiranGawande, “Amazon Review Classification and Sentiment Analysis”, *International Journal of Computer Science and Information Technologies*, Vol. 6, No.6, pp. 5107-5110, 2015.
16. Aishwarya Mohan, Manisha.R, Vijayaa.B, Naren.J , “An Approach to Perform Aspect level Sentiment Analysis on Customer Reviews using SentiscoreAlgorithm and Priority Based Classification”, *International Journal of Computer Science and Information Technologies*, Vol. 5, No.3,pp. 4145-4148, 2014.
17. S. Hamano and M. Sato, “Mining indirect association rules”, in *Proc.Ind.Conf. Data Min.*, pp. 106–116, 2004.
18. T. Herawan, A. Noraziah, Z. Abdullah, M. M. Deris, and J. H. Abawajy, “IPMA: Indirect patterns mining algorithm”, *Advanced Methods for Computational Collective Intelligence*, Springer , pp. 187–196, 2013.
19. Q. Wan and A. An, “Efficient mining of indirect associations usinghi-mine”, in *Proc. Conf. Can. Soc. Comput. Stud. Intell*,pp. 206–221, 2003.
20. W. Zhang, H. Xu, W. Wan, “Weakness Finder: Find product weakness from Chinese reviews by using aspects based SA,” *Expert Systems with Applications*, Elsevier, vol. 39, pp. 10283-10291, 2012.
21. A.Jeyapriya, C.S.KanimozhiSelvi, "Extracting Aspects And Mining Opinions In Product Reviews Using Supervised Learning Algorithm", *IEEE*, 2015.
22. Poria S, Cambria E, Winterstein G, Huang G.B, “Sentic patterns: Dependency-based rules for concept-level SA”, *Knowledge- Based Systems*, Vol. 69, pp. 45–63, 2014.
23. M. J. Zaki, M. Ogihara, S. Parthasarathy and W.Li, “Parallel Data Mining For Association Rules On Shared-Memory Multiprocessors”, In *Supercomputing, Proceedings of the ACM/IEEE Conference* , pp. 43-43, 1996.
24. Ayman E. Khedr, S.E.Salama, Nagwa Yaseen, “ Predicting Stock Market Behavior using Data Mining Technique and News Sentiment Analysis”, *International Journal of Intelligent Systems and Applications*, Vol. 7, pp.22-30, 2017.
25. B.Hemalatha, T. Velmurugan, “Direct-Indirect Association Rule Mining for Online Shopping Customer Data Using Natural Language Processing”, *International Journal of Recent Technology and Engineering* , Vo.8 , No.4, pp.11099-11106, 2019.