# LSA-MSC: LEXICAL SEMANTIC ANALYZE BASED MULTI-LEVEL SEMANTIC CLUSTERING FOR USER OPINION PREDICTION IN WEB MINING RESOURCE

[1]Murugannatham A, [2]Victor S P

[1]Assoicate Professor, [2]Dean & Associate Professor

[1]Department of Computer Science[PG],

[1]Kristu Jayanti College(Autonomous), Bangalore, India.

[2] Department of Computer Science,

[2]St. Xavier's College(Autonomous), Tirunelveli, India.

*Abstract :* The web mining expertise has more crucial information to process user opinions, reviews, talks, responses, feedback and provides productive consumer input. The opinion mining and human-agent interaction communities are currently addressing sentiment analysis from different perspectives that comprise the web mining resources, on the one hand, disparate sentiment-related phenomena and computational representations had more problems, and on the other hand, different extraction and dialog process management methods to present clustering evaluation to optimize the result. This paper is to propose a Lexical semantic analyze based multi-level semantic clustering algorithm to ensemble the evaluation of user opinion. Clustering product feature is the essential task to mine opinions from unstructured online reviews because different customers usually express the same feature with different words or phrases. Cluster ensembles the relative clusters measure that have been applied to accomplish this task using clustering accuracy. The resultant provided has great impact of mining optimistic result with time complexity with best feature case analysis.

*Index Terms* **- web mining, opinion mining, sentiment analysis, clustering, rank prediction**

## 1. Introduction

Sentiment analysis deals with establishment and classification of subjective information present in web mining. This could no longer necessarily be reality-identify based as humans have one-of-a-kind emotions towards the similar product, service, subject matter, occasion or person. Opinion extraction is vital part to identify the user hidden case opinion with a view to target the precise about the product to find where is the real opinion is expressed. Opinion from a person on a specialized subject might not matter except about the direct opinion. For opinion extraction and summarization the numerous entities are necessary.

The rapid growth of Internet technologies has resulted in an increase in online user's content creation. User generated information, usually represents the people's opinions, thoughts, reviews and sentiments. Automatic sensing and analysis of opinions about products, brands, political publications, etc. is a challenging job. Opinions are expressed in different ways in different domains. Constructing and labeling corpus for every domain is a pricey thing. Words from source and target domains are not always similar; hence classifier trained with one domain to another domain leads into poor performance. Therefore the need of domain adaptation algorithms arises to diminish domain's reliance and tagging costs. Using adapted maximum entropy with bipartite clustering in the proposed work, opinionated words are classified in two categories as positive and negative. Results demonstrated that proposed approach performs fairly well compared to baseline method

A main task of opinion mining is that of classifying documents by their polarization, i.e. whether a document is written with a positive feeling or a negative feeling. In opinion mining, a word's polarization is often used as an element in machine learning.

However, the polarity of some words cannot be determined without domain information. Reusability of knowledge of one domain to another domain is a key issue in opinion mining. Transfer learning is nothing but applying previously learned knowledge to solve new problems faster or with better solutions reviewed several trends of transfer learning. Consider the problem, where the task is to automatically classify the reviews of a domain, such as a movie, into positive and negative orientation. For this task, the first collection of reviews from different domains can be done. Then train a classifier on the reviews with positive and negative labels. The largest amount of reviews is needed to maintain good classification performance. Labeling reviews for each domain is time consuming as well as expensive process. Hence, the domain adaptation need arises which can use knowledge of one domain to another one.
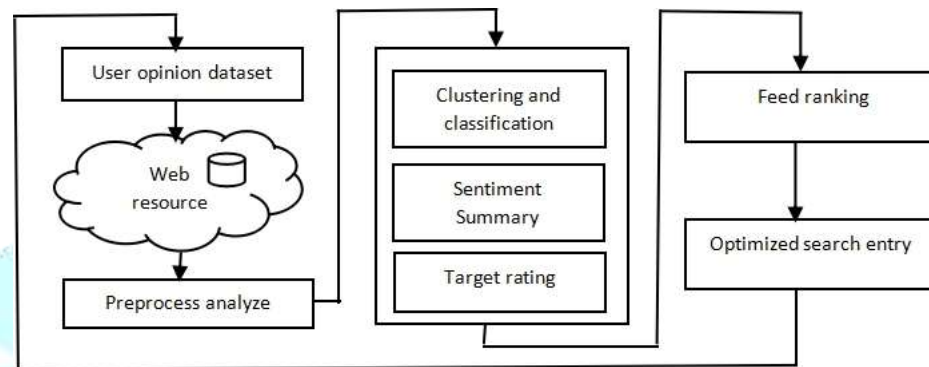


Figure 1: Optimization of opinion mining from web resource

An opinion word need not express the same sentiment everywhere. These words could be domain specific. Like the word "small" may represent a negative orientation in a web domain but if used in consumer product applications it is a positive. Different types of generalized lexicons are available which gives general sentiment polarity which is not applicable for the above scenario. Therefore, lexicons which can give polarity of the same word in different domains using a same lexicon database are needed. A proposed framework addresses major issue in opinion mining, i.e. domain adaptation. Cost domain adaptation is required because of high labeling. The proposed approach attempts to build a domain adaptable lexicon using adapted maximum entropy with bipartite clustering.

It is easy to make score and finalize the outcome whereas unstructured review that may usually include feedback in the form of text and images from various social monitoring tools and online shopping sites. In market each product may be introduced on the basis of some latest features they hold and they can either uplift or downsize the demand of that product. In this proposal we focus on identifying features from unstructured reviews and clustered it manually. Then the polarity of each clustered comments is determined to undergo further classification using supervised method Naive Bayes for prior distribution of each featured class attribute towards positivity and negativity.

This paper presents cluster approaches used in process of opinion mining and opinion summarization. This paper also explains outcomes and comparison of these approaches. Such comparison will be helpful to find methods that best suits for required scenarios in opinion mining.

## 2. Literature review

The opinion mining holds various points to mine useful information is mainly to abstract evaluation objects and evaluation words, and then judge the evaluation words' polarity [1]. After that, for each evaluation object, cluster its evaluation words.

Recommender systems attempt to predict items in which a user might be interested, given some information about the user's and items' profiles. Most existing recommender systems use content-based or collaborative filtering methods or hybrid

methods that combine both techniques [2]. created. To generate product recommendations and to solve determining the polarity of a sentiment bearing expression requires more than a simple bag-of-words approach [3]. In particular, words or constituents within the expression can interact with each other to yield a particular overall polarity.

The ensemble framework is applied to sentiment classification tasks, with the aim of efficiently integrating different feature sets and classification algorithms to synthesize a more accurate classification procedure [5]. Tree kernels encode not only syntactic structure information, but also sentiment related information [6], such as sentiment boundary and sentiment polarity, which are important features to opinion mining.

Selecting features for multi-class classification is critically important task for pattern recognition and machine learning applications. Especially challenging is selecting an optimal subset of features from high-dimensional data from text [8, 9], which typically have many more variables than observations and contain significant noise, missing components, or outliers. Because Opinions are central to almost all human activities and are key influence of our behaviors [13, 14], approaches to opinion feature extraction rely on mining patterns only from a single review corpus [15], ignoring the nontrivial disparities in word distributional characteristics of opinion features across different corpora.

Most Prior work has measured the emotional expressions in users' tweets and then performed various analysis and learning [17]. However, how to utilize those learned knowledge from the observed tweets and the context information to predict users' opinions toward specific topics they had not directly given yet is a novel problem presenting both challenges and opportunities. One of the most challenging problems in aspect-based opinion mining is aspect extraction [19], which aims to identify expressions that describe aspects of products (called aspect expressions) and categorize domain-specific synonymous expressions. Distant supervision that considers emoticons as natural sentiment labels in the micro-blog texts has been widely used in social media sentiment analysis [22]. However, the previous distant supervision works were normally trained based on an isolate set of data, and they were not capable to deal with the scenario where the texts are continuously increasing and the topics are constantly changing. To formalize this, the Fastest Threshold Clustering Algorithm in the domain of Twitter Sentiment Analysis [24], which was never tried by any researchers earlier in this specific area and also a novel method is designed and implemented to analyze the impact of a trendy product using the real data samples gathered from the micro-blog.

## 3.    Implementation of Proposed solution

Mining based on opinions can extract useful information from users' comments. After doing cluster and analysis on the information, users can get a detailed understanding of the commodity, then can determine to buy the commodity or not. In this paper, we extract evaluation objects and evaluation words of summarization from social comments, then cluster the evaluation objects judge the polarity of evaluation words and determine their polarity intensity values, then use multilevel lexical semantic clustering algorithm to cluster the evaluation words. For every kind of target evaluation object and each kind of evaluation word a count on the proportion of indexing direct from ratings are maintained and the results are shown to users in an intuitive way. This paper uses noun phrase pattern to match comments to extract evaluation objects and put forward the semantic words extraction algorithm. On judging the evaluation words' polarity, this paper establishes an emotional seed dictionary for each target object of opinion words. The method establishing dictionary for every attribute can eliminate the influence that less-correlation evaluation words have on the polarity judgment.

The proposed system presents a detailed description of the user opinion prediction about the commercial comments. The system consists of five main phases, Data Collection, Data Pre-processing, lexical pattern Clustering, feature extraction summarization and feed ranking Analysis. Figure.2 shows the architecture of the proposed system.
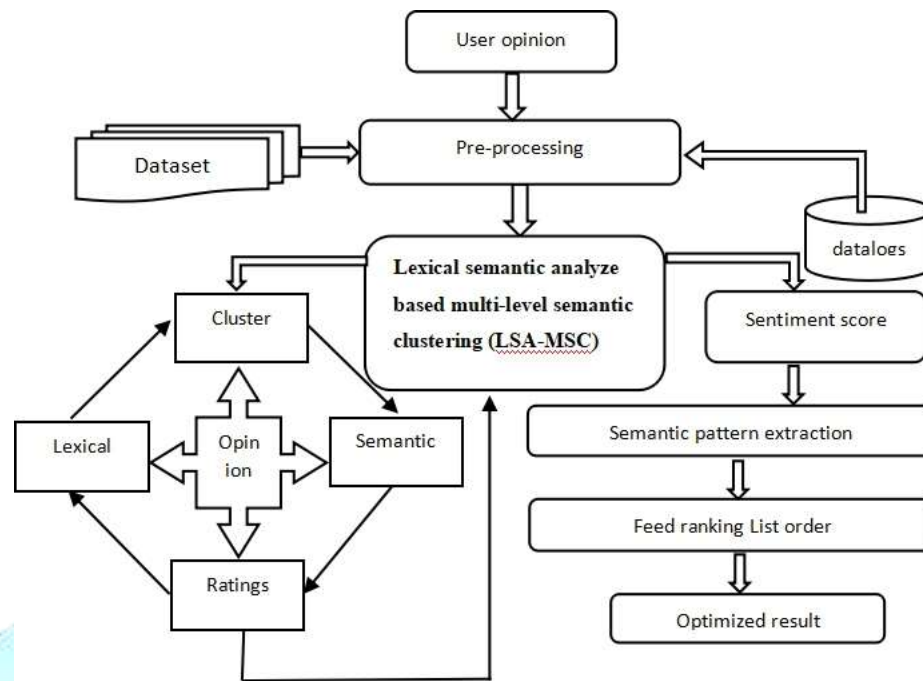
Figure 2 Architecture diagram of proposed implementation

The above figure 2. Shows the multilevel semantic prediction based on the user opinion by following way to analyze the semantic relation from hidden sentiments

In this paper, we mainly focus on short text messages with casual language usage as opinion. The system attempts to create groups of user postage on a particular topic based on content similarity and then providing concise opinion summary for each group. As the different user can ask for the summary at any instant the traditional clustering methods cannot satisfy the real-time necessity of this application. This problem is designed as progressive tweet stream clustering task. The system finds best clusters representing a different group of opinion towards a specific event. Every tweet within a cluster is classified as positive, negative and neutral sentiment. Finally, the important terms are extracted from each cluster to create indexing key-term cloud which enables users to get a brief understanding of opinions corresponding to an event. The contribution begins to demonstrate the pervasiveness of ambiguity in short texts and the limitations of traditional approaches in handling them:

• Achieve better accuracy of short text understanding, using knowledge-intensive approaches based on lexical-semantic analysis

• Improve the efficiency of our approaches to facilitate real-time applications

A) Semantic clustering evaluation objects: The extracted evaluation objects are various, but many nouns express the same meaning and users are concerned with only certain kinds of attributes. So we should cluster such various kinds of evaluation objects to determine the final evaluation object. In this paper, to cluster evaluation objects, we express each of the evaluation objects as a vector. First, we select twenty high-frequency evaluation words then the values of elements in the vector are the times that evaluation object and the twenty high-frequency evaluation words appear together in the comment text. In this way the evaluation object is converted to a twenty dimensional vector, then we cluster these evaluation objects. The clustering algorithm used in this paper is semantic clustering algorithm because this algorithm is high-efficiency and it is appropriate for large-scale data. This algorithm is more accurate than other cluster algorithms. When using multi-level clustering algorithm, the most critical jobs are the selection of initial centers and to determine the value of K. For the selection of initial centers, we should try to ensure that every center does not

belong to the same category. That is to say the distance between each two centers should be as large as possible. So this paper adopts the method of maximum distance to determine initial center.

B) Sentiment analysis is the process of discovering and grouping opinions hidden in the social message. In proposed system which identifies the sentiment on tweets collected from presidential debate performance event. Tweets are rated as belonging to one of four categories: Positive, negative, mixed and other. In addition to this, consider classifying each message along six sentiment labels and generates six-dimensional mood vector for each day in the timeline. The system for predicting optimized result by finding the political sentiment on twitter.

C ) Document summarization, Extractive summarization which involves selecting relevant sentences from the document that belong to summary and Abstractive summarization which creates summaries by synthesizing and rewriting sentences based on contextual understanding. In this paper, we concentrate on extractive summarization. Most of the techniques in Extractive summarization finds a score of each sentence in the document and then select top ranked sentences. Work focuses on maximizing the number of informative content-words by assigning a score to each term in document cluster using frequency, position information and then finding the set of sentences in the document cluster that maximizes the sum of these scores. Multi analyze is a system for a single as well as multi-document, multi-lingual summarization of opinion which performs sentences scoring based on sentence-level and sentence features such as cluster centroids, position, TF-IDF etc.

## 3.1 Pre-processing of nonlinear Input dataset

Data Collection and pre-processing of data implies the dataset contains product reviews taken from redundant nonlinear data from many product domains. This dataset contains three types of opinions positive, negative and neutralized format. These files are extracted using opinions splitter and reviews are written into stored procedure. In the Proposed Framework dataset contains positive opinions, negative opinions, and neutral case opinions for each domain. For the experiment 2000 opinions from each products as 1000 positive and 1000 negative opinions are taken. To remove noise from input dataset pre-processing is needed. Pre-processing is nothing but cleaning of data. Unnecessary words like "a, an, of, the, I, it, you, and, etc." are called stop words in English which overheads text documents. Hence the removal of these words improves effectiveness of information retrieval. To remove stop words from sentences, text file is used which consists of list of English stop words. Stanford parser is used to extract part of speech like a noun, adjective, adverb, etc. Different parts of speech tags such as noun, preposition, verb, adjective and adverbs assigned using parser known as Part-Of-Speech (POS) tagging. This is a very important step in opinion mining because finding opinionated words is not as easy as it gives adjectives, adverbs from the sentence. Adjectives and sometimes adverbs are indicator of opinion.

## 3.2 Lexical semantic pattern model

The external knowledge is indispensable for short text to understand from opinion words, which in turn benefits to many real-world applications that need to handle large amount of short texts from direct ratings. The lexical-semantic relationships between terms (namely words and phrases) are from a well-known probabilistic network and a web corpus. The proposed knowledge is intensive approach to understand short texts effectively and efficiently using lexical term indexing patterns.

Input: Source and target domains pre-processed review documents. For all features extracted from the input dataset.

1.Initialize dataset x={x1, x2, x3,..},n is total number of features.

2.For x= i=1 to n+1

For each text data Di from Ds

    Text T = Extract Text from Di.

$$\text{Sentence set Ss} = \left(\sum_{n=1}^{size(Ds)} Text€\ Di\right) \times Splitby\ (.)$$

Repeat until convergence Calculate the probability of clusters each word as c

Where x= xi +ni calculated the iteration,

F(c,di)= tf-idf value of each word considered as feature

End for

3.Weight assignment by Point-wise Mutual Information

For y=assign term set value process p,i=n

For each term Tk from Ti

If Di€ then

Identify relation it has with other concepts.

Relation Set Rs = ∑(Concepts €Di)+Ci.   (2)

Add relations to the word set Ni.

$$P(x,y) = \log_n \frac{p(x^\wedge y)}{p(x)p(y)}$$

End

End

End for

where x and y are words from documents

4. Common and uncommon words from source and target domains are clustered using weight value out point

Ot -> Ni as neutralized data value.

5. Output: positive and negative classified terms and word clusters.

Our approach exploits external knowledge and infers the best segmentation based on the semantics among the terms, which reduces its dependency on POS tagging. Furthermore, in order to calculate semantic relatedness, the set of terms (namely the segmentation of a short text) should be determined first, which raises the necessity to accomplish segmentation first.

### 3.3 Semantic short opinion weight evaluation

The sentiment score through semantic similarity has been proven to be much more preferable than surface similarity. However, an incorrect cluster of short texts leads to incorrect semantic similarity which creates a neutralized case data. The machine learning approach applicable to sentiment analysis mostly belongs to classification of relevant review score in general and text classification technique. In a machine learning based classification, two sets of opinions reviews are required, training and a test set. A training set is used by an automatic classifier to learn the differentiating characteristics of text reviews, and a test set is used to validate the performance of the automatic classifier.

First, it randomly selects a clusters C = (u, v) with probability proportional to its weight.

Input: G = (V, E); W (E) = {w(C) |C ∈ E} Output: G = (V, E); s (G)

step1: Initialize V = ∅ ; E = ∅ as dominant clusters

step2: while E ∅ do

step3: Randomly select the C cluster centroids to opinion match case

Randomly select clusters C = (u, v) from E with probability proportional to its

weight

Step4: If (i=0; c>v; c++)

V = V ∪ {u, v}; E = E ∪ {C} for positive state;

V = V − {u, v}; E = E − {C} for negative state;

Step 5: Calculate the centroid of all opinions factors in each cluster weight

For each t ∈ V do

if e = (u, t) E or e = (v, t) E then

V = V − {t} ;

Remove edges linked to t from E: E = E − {C = (t, ∗ )};

End if

End for

End if

End while

Step6: calculate average edge weight: $s(G) = \sum Terms(Ts) \nexists O(c)$

Output acquired: (i) Cluster centroids weighting factors; C (ii) Cluster labels of opinion dataset

The opinion clustering method was used for detecting hidden patterns in our unlabeled data samples of 1500 real data-sets about the sentiments from online data-sets. Right here we've initialized 3 because the range of clusters for appearing the analysis. The statistics factors have been randomly chosen because the centroids C. primarily based on the distance degree known as Euclidean distance, opinions points have been assigned to the closest cluster centroid. The imply value of the cluster changed into expected and the newly received centroids were updated therefore the steps were repeated until the similar information factors were consecutively allotted to every cluster. The clusters generated changed into categorized into wonderful, terrible, neutral sentiments of the famous brand to predict the effect of the corresponding product emblem.

### 3.4 Multilevel Feature-based sentiment analysis

Features clustering only unigrams are selected for training an opinion classifier. Features i.e. the words are classified as domain independent words and domain specific words, depending on the frequency of words in particular domains. Dominant words like indexing are constructed between common and uncommon words. Both the sets are connected by edges which gives the relationship between two vertices. Weight of edge is calculated based on the distance between two nodes. The smaller their distance, the larger weight can be assigned to the corresponding edge. Importance of clustering is to reduce the mismatch between domain specific words of source and target domain. Two sets of features are extracted as an output with polarity. Finally, all classified words are summarized in a single table with polarity value, class and domain name.

Since the problem of grouping feature expressions is a clustering task, two common measures for evaluating clustering are use to find the relevance score.. Given a data set PDS, its gold partition is G = {Pd1,…, pd2,…, pdn}, where k is the given number of clusters. The Groups partition DS into k disjoint subsets, Ds1,…, Dsi, …, Dsk.

For each feature expression Fi in PDs

Fi ← all sentences containing Fi in G

For each sentence Pos-fs ∈ Si :

For each positive word Pds in opinion dataset [-Ds, Ds]:

For each sentence Neg-fs ∈ Si :

For each Negative word Pds in opinion dataset [-Ds, Ds]:

If v does not co-occur with vi in any review sentences:

Pds ← Ds, Fi ← words from all ds, j = Ds1,Ds2, …, |Fi|

End

To find positive opinion score in relational text,

$$\text{Pos-Fs} = \sum_p^{pn} \frac{Total\ extracted\ index\ features\ in\ cluster\ (Ds)}{total\ number\ of\ clusters\ (Dsk)}$$

To find negative opinion score in relational text,

$$\text{Neg-Fs} = \sum_p^{pn} \frac{Total\ extracted\ non-index\ features\ in\ cluster\ (Ds)}{total\ number\ of\ clusters\ (Dsk)}$$

Others are considered as neutralize.

The scores of opinion from each term are not some arbitrary scores assigned to them. These scores reflect the positivity, neutrality and negativity of the terms in the related context. Instead of working with high-dimensional vectors work with 3-dimensinal vectors.

## 3.5 Lexical semantic multilevel clustering and prediction

The clustering of the text opinions is performed by computing the semantic bound and semantic closeness measure. Then the method selects the top class according to the semantic closeness measure. From the selected class, the method identifies the list of terms and with the term set identified from the reviews text, the method identifies the non-common elements.

Initialize clusters $C_{new}$ set as ci;

Compute sentiment score measure $C_{new} = \int \frac{cluster\ C(i,j)}{Number of\ Max\ an\ Min\ in\ Cnew}$

For Ci terms assign centroid cluster value to Cn

$C_{new}$ →cluster set Ct!=0;

For each cluster Cj in C

If $C_{new}$ (Ts,Cj)>=closeness edge(Cj)

Add max_weight value opinion to $C_{new}$ →cluster set Ct;

Positive opinion OPj.<-- c;

Else

Add min_weight value opinion to $C_{new}$ →cluster set Ct;

Negative opinion OPi.<-- c;

End if

End for

Select non Match clusters $C_{nm} \rightarrow$ neutral Update $C_n$ in C;

Update OPj in C; Update OPi in C;C$\leftarrow$(OP(i,j));

End for

Using the opinion sentiments related elements identified, the method computes the frequency of the terms towards each class. Based on the frequency computed and the semantic closeness measure, the method computes the truth weight for each class. Based on truth weight computed the opinion is assigned with the selected class as positive or negative.

### 3.6 Feed ranking Dataset and Reviews of User opinion

The feed ranking generalize the opinions to rank the order as per the aspect ratio to computing a single score for each term: starting from a seed set of words with arbitrary scores and propagate them to the other words. Compute a score for each sentence and also for its neighbors. These systems focus on objective parts to extract factual information, they can discard subjective sentences. It feeds the maximum entropy clusters for positive and negative opinion cases, then the recommendation system can recommend items with positive feedback and not recommend items with negative feedback.

### 4. RESULTS AND DISCUSSION

The results are carried out through with Microsoft frame work with real time product dataset rating from UCI repository social summarized opinion dataset. The data sets are carried out opinions from web product reviews. The proposed Lexical semantic analyze based multi-level semantic clustering and opinion intelligence generation algorithm has been implemented and tested for its efficiency. The proposed method has yield the opinions of mind harvest about product features to give efficient results on clustering evaluation and improves the performance also on predicting opinion rates. Parameters are tabulated given below.

**Table 4.1: Details of Data set**

| Parameter | Value |
|---|---|
| Number of opinions | 2000-5000 |
| Number of case thoughts | Positive, negative, neutralized |
| Datasets used | online product reviews |

Above Table 4.1, shows the details of data set being used to evaluate the performance of the proposed Lexical semantic

analyze based multi-level semantic clustering approach. The performance of LSA-MSC is evaluated through clustering accuracy (cs), false classification ratio (Fcr), time complexity (Ts) and frequent occurrence (Fs) The resultant figure given below shows the performance,

$$\text{Clustering accuracy (cs)} = \sum_{k=0}^{k=n} \times \frac{\text{total number of cluster group dataset(Cds)}}{\text{Total originate dataopinions(Tr)}}$$
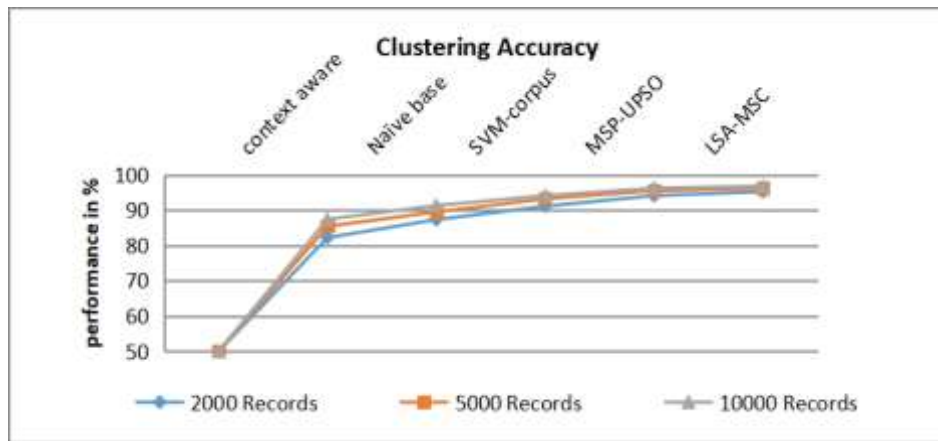
**Figure 4.6: Impact of Clustering Accuracy**

Above Figure 4.6, shows the comparison of clustering accuracy and shows that the proposed method has produced higher clustering accuracy of 96.8 % as well as than other methods.

$$\text{False classification Ratio (Fcr)} = \sum_{k=0}^{k=n} \times \frac{\text{clustering Accuracy(cs)}}{\text{Total no of failted cluster rate(Fr)}}$$

**Table 4.2 Evaluation of clustering accuracy**

| Clustering accuracy produced by dissimilar methods in % | | | | | |
|---|---|---|---|---|---|
| Methods/number of opinions | Context aware | Naïve base | SVM-corpus | MSP-UPSO | LSA-MSC |
| 2000opinions | 82.2 | 87.3 | 91.1 | 94.1 | 95.2 |
| 5000 opinions | 85.4 | 89.5 | 93.2 | 95.5 | 96.2 |
| 10000 opinions | 87.4 | 91.3 | 94.1 | 96.1 | 96.8 |

Above Table 4.2, shows the comparison of clustering accuracy produced 2000 opinions as 95.2%, 5000 opinions as 96.2% and 10000 opinions as 96.8 % shows that the proposed approach has produced higher clustering accuracy.
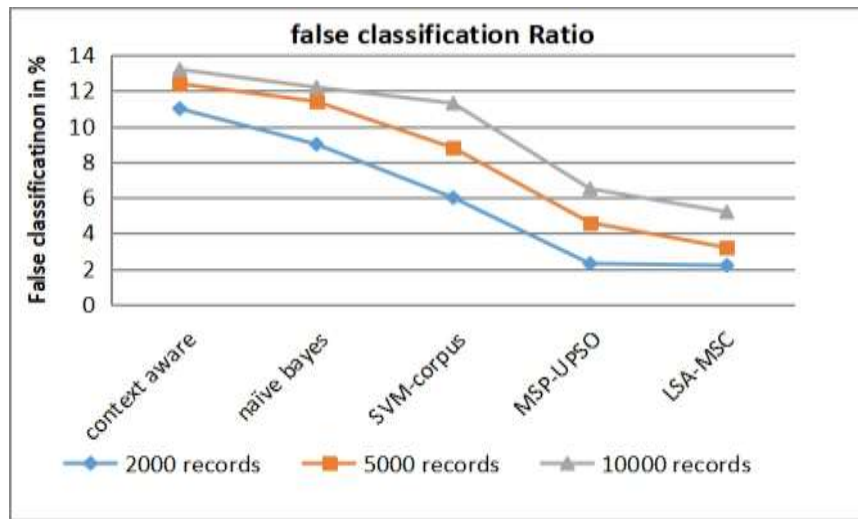
**Figure 4.7: Impact of false classification**

Above Figure 4.7, shows the comparison of false classification ratio produced by different methods and the proposed method has produced less false classification ratio than other methods.

**Table 4.3: Evaluation of false classification**

| False classification produced by dissimilar methods in % | | | | | |
|---|---|---|---|---|---|
| Methods/number of opinions | Context aware | Naïve base | SVM-corpus | MSP-UPSO | LSA-MSC |
| 2000 opinions | 11 | 9 | 6 | 2.3 | 2.2 |
| 5000 opinions | 12.4 | 8.5 | 8.8 | 4.6 | 3.2 |
| 10000 opinions | 13.2 | 12.2 | 11.3 | 6.5 | 5.2 |

Above Table 4.3, shows the comparison of false classification ratio produced 2000 opinions as 2.2%, 5000 opinions as 3.2% and 10000 opinions as 5.2 % shows that the proposed approach produces less false classification ratio.

$$\text{Time complexity (Tc)} = \sum_{k=0}^{k=n} \times \frac{\text{clustering Accuracy(cs)+false classification(Fcr)}}{\text{Time taken(Ts)}}$$
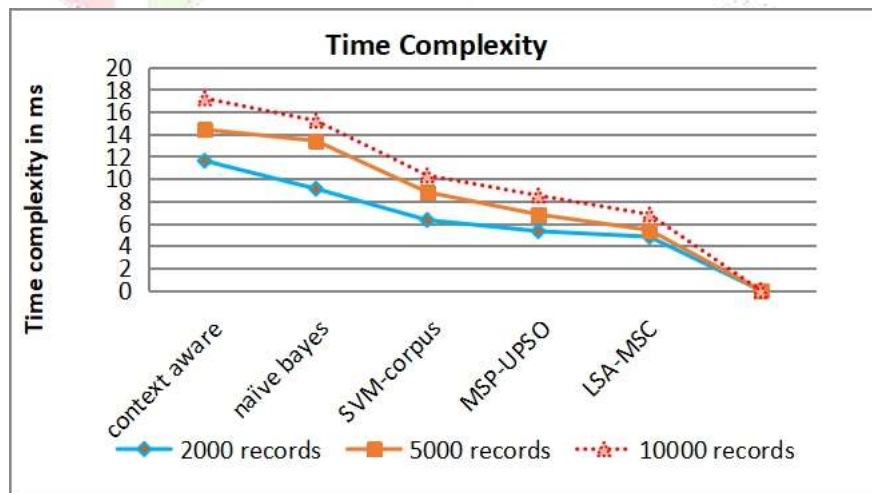


**Figure 4.8: Impact of time complexity**

Above Figure 4.8, shows the comparison of time complexity produced by different methods and shows that the proposed approach has produced less time complexity than other methods.

**Table 4.4: Evaluation of time complexity**

| Time complexity produced by dissimilar methods in milliseconds (ms) | | | | | |
|---|---|---|---|---|---|
| Methods/number of opinions | Context aware | Naïve base | SVM-corpus | MSP-UPSO | LSA-MSC |
| 2000 opinions | 11.6 | 9.1 | 6.3 | 5.3 | 4.8 |
| 5000 opinions | 14.4 | 13.4 | 8.8 | 6.6 | 5.4 |
| 10000 opinions | 17.2 | 15.2 | 10.3 | 8.5 | 6.8 |

Above Table 4.4, shows the comparison of time complexity LSA-MSC produced 2000 opinions as 4.8(ms), 5000 opinions as 5.4(ms) and 10000 opinions as 6.8(ms) shows that the proposed approach has produced less time complexity.

$$\text{Frequent occurrence (Fc)} = \sum_{k=0}^{k=n} \times \frac{\text{repeated clusters(Rs)} + \text{irrelavant clusters(Irc)}}{\text{Total number of clusters}}$$
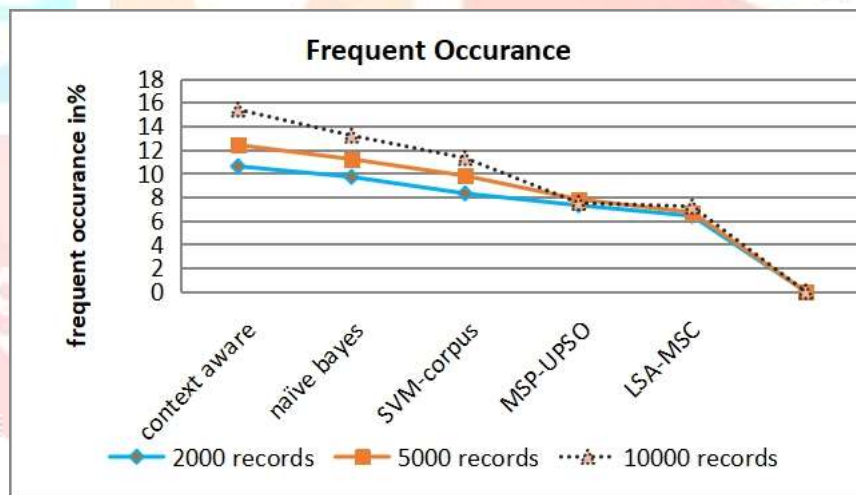


**Figure 4.9: Impact of Frequent occurrence**

Above Figure 4.9, shows the comparison of frequent occurrence produced by different methods and shows that the proposed approach has produced less frequent than other methods.

**Table 4.5: Evaluation of time complexity**

| Frequent occurrence produced by different methods in % | | | | | |
|---|---|---|---|---|---|
| Methods/Number of opinions | Context aware | Naïve base | SVM-corpus | MSP-UPSO | LSA-MSC |
| 2000 opinions | 10.6 | 9.7 | 8.3 | 7.3 | 6.4 |
| 5000 opinions | 12.4 | 11.2 | 9.8 | 7.8 | 6.7 |

| 10000 opinions | 15.4 | 13.2 | 11.3 | 7.5 | 7.2 |
|---|---|---|---|---|---|

Above Table 4.5, shows the comparison of frequent occurrence LSA-MSC produced 2000 opinions as 6.4%, 10000 opinions as 6.7% and 10000opinions as 7.2 % shows that the proposed approach has produced less frequent occurrence.

## 5.  CONCLUSION

The process of user opinion mining on the hidden occurrence from any real world event is very high. This capability of sentiment analysis can be utilized using the event summarization system. The progressive multilevel clustering proves to be useful for real time product review summarization problem as well the lexical semantic patterns. It updates the clusters progressively for newly arriving opinions from relevant score of product ratings. Then these clusters are summarized using lexical indexing term which gives brief understanding to the users about the event at any instance. Sentiment analysis is performed on each cluster. Then the clusters are classified as positive, negative and neutral sentiment. Sentiment Analysis gives more depth understanding about views of people related to the particular event. The proposed system produces higher efficient result on good performance - 96.8 % compared to the other eventual methods.

**REFERENCES**

[1] B. Gao, T.-Y. Liu, G. Feng, T. Qin, Q.-S. Cheng, and W.-Y. Ma, "Hierarchical taxonomy preparation for text categorization using consistent bipartite spectral graph co-partitioning," Knowledge and Data Engineering, IEEE Transactions on, vol. 17, no. 9, pp. 1263–1273, 2005.

[2] Aciar S., Zhang D., Simoff S., and Debenham J, "Informed Recommender: Basing Recommendations on Consumer Product Reviews," In Intelligent Systems, IEEE, Vol.22, Issue 03, pp.39-47.2007

[3] Y. Choi and C. Cardie, "Learning with Compositional Semantics as Structural Inference for Sub sentential Sentiment Analysis," Proceedings of the Conference on Empirical Methods in Natural Language Processing (EMNLP), pp. 793-801, 2008.

[4] S. Na and L. Xumin, "Research on k-means Clustering Algorithm An Improved k-means Clustering Algorithm", Intelligent Information Technology and Security Informatics, pp. 63-67, IEEE, April 2010.

[5] R. Xia, C. Zong, and S. Li, "Ensemble of Feature Sets and Classification Algorithms for Sentiment Classification," Information Sciences, vol. 181, no. 6, pp. 1138-1152, 2011.

[6] P. Jiang, C. Zhang, H. Fu, Z. Niu1 and Q. Yang, "An Approach Based on Tree Kernels for Opinion Mining of Online Product Reviews", Data Mining, pp. 256–265, IEEE, Dec. 2010.

[7] Abbasi, S. France, Z. Zhang, and H. Chen, "Selecting attributes for sentiment classification using feature relation networks," IEEE Transactions on Knowledge and Data Engineering (TKDE), vol. 23, no. 3, pp. 447-462, 2011.

[8] Pan SinnoJialin , Xiaochuan Ni , Jian-Tao Sun , Qiang Yang and Zheng Chen (2010), "Cross-Domain Sentiment Classification via Spectral Feature Alignment", Proceedings of the 19th International World Wide Web Conference, ACM, Raleigh, USA, April 26-30, 2010.

[9] Q. Cheng, H. Zhou, and J. Cheng, "The Fisher-Markov Selector: Fast Selecting Maximally Separable Feature Subset for Multiclass Classification with Applications to High-Dimensional Data," Pattern Analysis and Machine Intelligence, Vol. 33, pp. 1217–1233, IEEE, June. 2011.

[10] C. Lin, Y. He, R. Everson, and S. Ruger, "Weakly supervised joint sentiment-topic detection from text," IEEE Transactions on Knowledge and Data Engineering (TKDE), vol. 24, no. 6, pp. 1134-1145, 2012.

[11] Z. Zhai, B. Liu, H. Xu, and P. Jia, "Constrained lda for grouping product features in opinion mining," in Advances in knowledge discovery and data mining. Springer, pp. 448–459.2011.

[12] X. Yu, Y. Liu, X. Huang, and A. An, "Mining online reviews for predicting sales performance: A case study in the movie domain," IEEE Transactions on Knowledge and Data Engineering (TKDE), vol. 24, no. 4, pp. 720-734, 2012.

[13] B. Liu, "Sentiment Analysis and Opinion Mining," Synthesis Lectures on Human Language Technologies, Morgan & Claypool, pp. vol. 5, no. 1, pp. 1-165, 2012.

[14] Y. Yang, Y. Ma, and H. Lin, "Clustering product features in opinion mining," Journal of Chinese Information Processing, vol. 26, no. 3, pp. 104–108, 2012.

[15] Z. Hai, K. Chang, J. Kim, and C. C. Yang, "Identifying Features in Opinion Mining via Intrinsic and Extrinsic Domain Relevance," IEEE Transactions on Knowledge and Data Engineering (TKDE), vol. 26, no. 3, pp. 447-462, 2014.

[16] T. Wang, Y. Cai, G. Zhang, Y. Liu, J. Chen, and H. Min, Product Feature Summarization by Incorporating Domain Information. Springer Berlin Heidelberg, 2013.

[17] F. Ren and Y. Wu, "Predicting user-topic opinions in twitter with social and topical context," Affective Computing, IEEE Transactions on, vol. 4, no. 4, pp. 412–424, 2013.

[18] L. Zhang and B. Liu, "Aspect and entity extraction for opinion mining," in Data mining and knowledge discovery for big data. Springer,  pp. 1–40. 2014.

[19] T. Wang, Y. Cai, H.-f. Leung, R. Y. Lau, Q. Li, and H. Min, "Product aspect extraction supervised with online domain knowledge," Knowledge Based Systems, vol. 71, pp. 86–100, 2014.

[20] L. N. Ferreira, L. Zhao, " A Time series Clustering Technique based on Community Detection in Networks", Proceeding in Computer Science, Elsevier, pp. 183-190, 2015. International Journal of Advanced Research in CS and Software Engineering, Vol. 5, Issue 8, Aug. 2015.

[21] N. Pooranam, G. Shyamala, "A Statistical Method of Knowledge Extraction on Online Stock Forum Using Subspace Clustering with Outlier Detection", International Journal of Innovative Research in Science, Engineering and Technology, Vol.5, Issue 5, May.2016.

[22] R.Xia , J.Jiang, and H.He, "Distantly Supervised Lifelong Learning for Large-Scale Social Media Sentiment Analysis" IEEE transactions on affective computing, vol. 8, no. 4 ,pp. 751 – 770,2017.

[23] B. Burscher, R. Vliegenthart, H. Claes, "Frames Beyond Words: Applying Cluster and Sentiment Analysis to News Coverage of the Nuclear Power," International Social Science Computer Research, vol 2,Issue 32 pp. 1-16, 2015.

[24] H Suresh, G Raj . "A Novel Cluster-Based Unsupervised Technique for Twitter Sentiment Analysis", International Journal of Innovative Research in Computer and Communication Engineering, Vol. 6, Issue 7,pp.345,362, 2017.

[25] L.Zhu, A. Galstyan, "Tripartite Graph clustering for Dynamic Sentiment Analysis on Social Media", International Social and Information Networks, vol 14, pp.1531-1542, June. 2014.