# Feature Extraction and Machine learning Approach (SMO) in Sentiment Analysis and Opinion Mining

Yadwinder Singh[1], Er. Rohini Sharma[2]
[1]M.Tech(Scholar) , [2]Assistant Professor
Department of Computer Science and Engineering, Gurukul Vidyapeeth Institute of Engineering & Technology

*Abstract*—Social Media of the improvement with online documents such as the comments of novel blog and article have conventional nice devotion; and the sentiment analysis via online papers has become 1 famous investigation area. Investigation of judgment orientation is purpose to find the useful the valuable orientation information, it becomes a study focus in the nature language processing, mainly in micro blog based on current semantic match, this paper presents a sentence verification method taking benefits of an enhanced algorithm for intention like Facebook comments term semantic value. Firstly this paper described the technique sentimental analysis. The proposed work has implemented create the database in two ways i.e. first single comment which is time consume & second one multiple data upload 1 at a time which is time saving approach. Apply the feature extraction technique means to find which comment is unique & create the Ascii code in every alphabet cause of comment has safe. Three categories we have described i.e. positive, negative and neutral. The classification technique classify the data in the 2 modules i.e. training module & testing module is defined the percentage of the each category.

*Keywords*— Sentimental, Manual, Processing, Keyword Processing, Vector Method, SMO Technique.

## I.    INTRODUCTION

Sentiment is a view, feeling, opinion or assessment of [1] a person for some product, event or service .Sentiment Analysis or Opinion Mining is a stimulating Text Mining & Natural Language Processing problematic for automatic extraction, organization and summarization of opinions and emotions expressed in online text .Sentiment study is replacing traditional & web based reviews conducted by businesses for finding public opinion about objects like products & facilities. Sentiment Inquiry also assists entities & organizations interested in knowing what other persons comment about a specific product, service topic, issue and event to find an optimal choice for which they are looking for.

Sentiment analysis is of great value for business intelligence applications [2], where business analysts can analyses public sentiments about products, services, and policies. Sentiment Analysis in the context of Administration Intelligence aims at removing public views on administration strategies and decisions to infer possible public reaction on implementation of certain policies.

The sentiment found within comments, feedback or critiques provide useful indicators for many dissimilar purposes. These opinions can be categorized either into two categories: positive and negative; or into an n-point gauge, e.g., very decent, good, acceptable, bad, very bad. In this respect, a sentiment analysis task can be inferred as a classification assignment where each group represents a sentiment. Sentiment analysis provides companies with capitals to estimation the extent of product receipt and to determine strategies to improve product quality. It also simplifies policy makers or legislators to analyses public sentiments with respect to policies, public services or political issues [3].

## II.    TYPES OF SENTIMENT ANALYSIS

Firstly we need to understand the techniques that communal media dealer's usage to determine sentimentality. As I mention above, there are many types of sentiment analysis. However, for the resolutions of this article we will distillate on 3 :

### A.    Manual Processing

Human understanding of sentiment is definitely the most mature and accurate judge of sentiment. However, it still isn't 100% accurate.  Very rare retailers still use this process without the other use of a tool. This is due to the prolific growth of social media.  Allowing to Seth Grimes, communal is the fastest growing source of enterprise analytical data.

### B.    Keyword Processing

Keyword dispensation algorithms allocate a degree of positivity or negativity to an individual word, then it gives and overall proportion score to the sustenance. For example, progressive words, great, like, love or undesirable words: terrible, dislike. The recompenses of this technique are that it is very fast, foreseeable & cheap to implement & run.

### C.    Natural Language Processing

(NLP also called: script analytics, data mining, computational linguistics) NLP refers to computer systems that procedure human language in positions of its denotation. NLP understands that numerous words make a phrase, several phrases create a sentence &, ultimately, verdicts convey ideas. NLP works by analyzing language for its sense.  NLP structures are used for in a amount of areas such as changing speech to text, language translation and grammar checks.

Table 1 Advantages and Disadvantages of sentiment analysis

| Advantages | Disadvantages |
|---|---|
| The ability to adapt and create trained models for specific purposes and contexts | The low applicability to new data because it is necessary the availability of labelled data that could be costly or even prohibitive. |
| Wider term coverage | Finite digit of words in the dictionaries & the meeting of a fixed sentiment orientation and score to words |
| Lexicon/learning symbiosis, the detection and dimension of sentiment at the impression level & the lesser sensitivity to changes in topic domain | Noisy reviews[4] |

## III. RELATED WORK

M. S. Usha et.,al(2013)[5] Sentimental analysis, a sub regulation within data mining and computational linguistics, mentions to the computational procedure for mining, sympathetic and assess the opinions expressed in many view rich resources like blogs, argument forums etc. The objective of sentiment analysis is to recognize emotional states in online text. Utmost of the time classifiers qualified in one area do not do well in another area. Also the difficulty in existing methods is not to detect sentimentality & subjects. Sentiments may vary with topics. This paper proposes a new model called Combined Sentiment Topic model to detect sentiments and topics at the same time from text. This replica is based on Gibbs specimen process. Besides, unlike overseen approach to view mining which often fails to create good presentation when unstable to other domains, the unverified natural world of CST makes it highly moveable to other domains. CST model achieves better compared to present semi- controlled approaches. Xinzhi Wang1 et.,al(2013)[6] With the development of social media, online ID such as the comments of information articles, blogs and micro blogs have conventional great attention; and the sentiment psychoanalysis via online ID has become one well-liked investigate area. This paper focuses on establishing user sentimental room obtained from online ID to analyses user's personalized sentiments, which aims to identify user's emotional feature. Dislike, sentiment and attributes of user are firstly employed to build user's modified sentimental universe. Then, the general obliges of user sentiments space are proposed to calculate user's behavior. And finally we seek out maudlin leaders who paly pivotal role in the leading public opinions. Our works can stretch some suggestions for conclusion makers when urgent event happen. Bin Wen1 et .,al(2012)[7] Research of judgment orientation is aim to find the useful alignment information, it becomes a research focus in the nature language dispensation, especially in Micro-blog. Based on the occurred How Net semantic similarity, this paper presents a sentence location identification technique taking benefit of an improved algorithm for calculating Chinese term semantic alignment value. Firstly, this paper familiarized the progressive method of term semantic orientation computation; then, combining adverbs

& unifications points a sentence documentation algorithm. And the experimental results show that the proposed approach was suitable for judging sentences' sentimental orientation. Ye Wu et.,al(2011)[8] Recently, research about social systems has attracted wonderful benefits. It can be considered that the links of online social networks describe the relationships between individuals. Analyzing online data from social networks provides opportunities for extracting qualities of sentimental inspiration, which also supports to get over the corner of current research on sentiment analysis. In this paper we project models to study both romantic influencing probabilities and influenced probabilities for consumers of Twitter, 1 of the utmost popular online communal media. They find that there is a high correlation between Twitter users' manipulating likelihoods & influenced probabilities, and the majority of users keep sentimental balance on both.

## IV. PROPOSED ALGORITHM

### A. Vector Method:

The vectorization is the method of converting an algorithm from a scalar implementation, in which does an process on 1 pair of operands at a period, to a vector process, if a particular instruction can mention to a vector (series of together values)

1. The totally effect, it adds a form of parallelism to software in which 1 order or operation is applied to manifold pieces of information. When done on computing systems that support such actions, the advantage is more efficient processing & better application performance. Many general-purpose microprocessors today feature multimedia additions that support SIMD (single-instruction-multiple-information) parallelism. And when the hardware is coupled with FORTRAN compilers that provision it, developers of scientific & engineering applications have an easier time delivering extra efficient, better performance software

2. Performance or efficiency benefits from vectorization depend on the code arrangement. But, in general, the automatic & near automatic techniques introduced below are most productive in delivering better performance or efficiency. The methods offering the most control require greater application knowledge & skill in knowing where they should be implementing. But these more intrusive techniques, such as those that may involve compiler directives or different source code alterations, can yield potentially greater performance and efficiency benefit when properly used. Feature engineering is an extremely basic and essential task for Sentiment Analysis. Changing a piece of text to a feature vector is the basic phase in any data driven approach to SA. In the following section we will see some commonly used features used in Sentiment Analysis & their critique.

• Term Presence vs. Term Frequency
Term frequency has always been considered essential in traditional Information Retrieval and Text Classification tasks.

But [9] found that term presence is more important to Sentiment analysis than term frequency. That is, binary-valued feature vectors in which the entries merely specify whether a term occurs (value 1) or not (value 0). This is no counter-intuitive as in the numerous examples we saw before that the presence of even a single string sentiment bearing words can opposite the polarity of the whole sentence. It has also been seen that the occurrence of rare words contain more information than frequently occurring words, a phenomenon so-called Hapax Legomena.

- Term Position

Words appearing in positive positions in the text carry more sentiment or weightage than words appearing elsewhere. This is similar to IR where words looking in topic Titles, Subtitles or Abstracts etc are given extra weightage than those appearing in the body. In the example given in Section 1.3.c, although the text contains positive words through, the presence of a negative sentiment at the end sentence plays the determining role in determining the sentiment. Thus generally words appearing in the 1st few sentences and last few sentences in a text are given more weightage than those appearing elsewhere.

- N-gram Features

N-grams are accomplished of capturing context to some extent & are widely used in Natural Language Processing tasks. Whether higher order n-grams are useful is a matter of debate. [10] reported that unigrams outperform bigrams when classifying movie reviews by sentiment polarity, but [11] found that in some settings, bigrams and trigrams perform better.

*B. SMO Technique*

Classification Method is a modest algorithm that quickly solves the SVM QP problem without any extra medium storage and without raising an iterative numerical routine for each sub-problem. SMO decomposes the overall QP problematic into QP sub-problems comparable to Osuna's method. SMO chooses to solve the smallest possible optimization difficult at every step. For the standard SVM QP problematic, the smallest possible optimization problem involves two Lagrange multipliers since the Lagrange multipliers must obey a linear equality constraint[12]. At every step, SMO chooses two Lagrange multipliers to jointly optimize, discoveries the optimal values for these multipliers, and updates the SVM to reflect the new optimum values. The benefit of SMO lies in the fact that solving for two Lagrange multipliers can be done logically. Thus, a complete inner iteration due to numerical QP optimization is avoided.

The inner loop of the algorithm can be expressed in a minor amount of C code, slightly than invoking an entire iterative QP library repetitive.
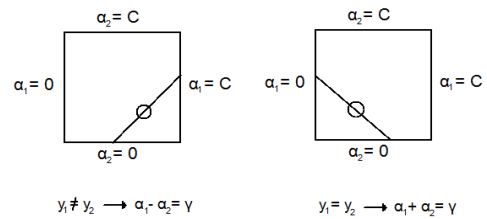


Fig. 1: SMO Techniques

Even though more optimization sub-problems are resolved in the course of the algorithm, each sub-problem is so fast that the complete QP problem can be answered fast. In addition, SMO does not require extra matrix storage (ignoring the minor quantities of memory compulsory to store any 2x2 matrices required by SMO). Thus, very large SVM training problems can't inside of the memory of an ordinary individual computer or workstation. Because operation of large matrices is avoided, SMO may be fewer susceptible to numerical precision problems.
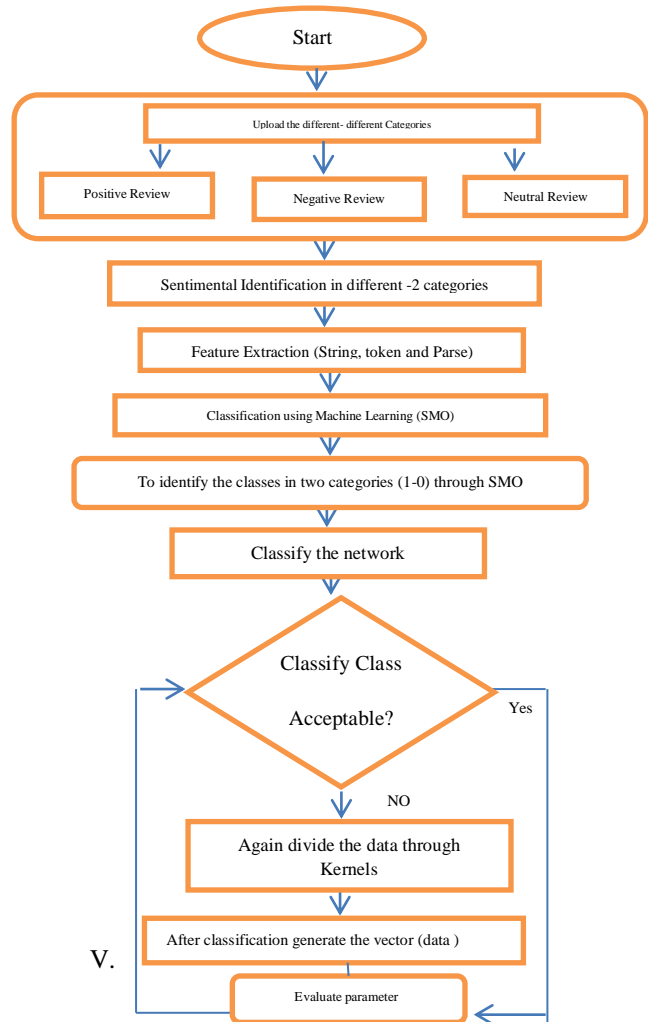


Fig. 2: Proposed Flow chart

## VI. SIMULATION MODEL

This Process shows that which sentiment detected from a input with the help of SMO and system's knowledge base. It compare the features of data and datasets and generate best matching with the knowledge base generate as system's output in this section. This section shows the matching possibilities of three different categories. One sentence having multiple sentiments. So we need to analyze whole sentence for best solution for every input.

All the parameters show the efficiency of the system. FRR and FAR is used find the rejection rate and acceptance rate that how correctly the system is accepting the right authentication and false authentication accurately and efficiently. The one another parameter ACURACY is used to check the detection accuracy and processing the textual data over different knowledge base. The values of FRR and FAR need to be less and accuracy will be high in best solutions systems.
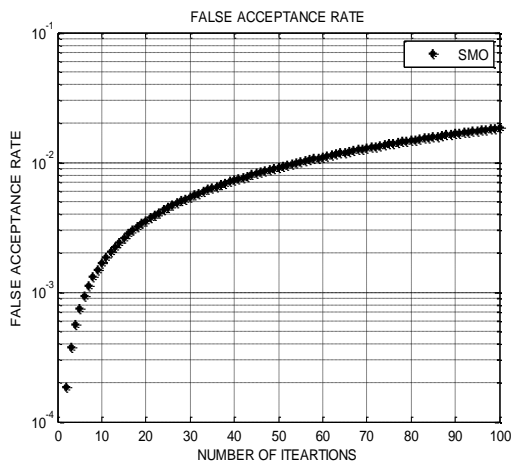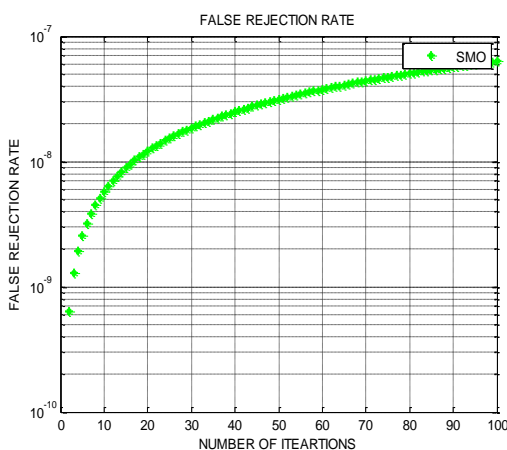


Fig. 3: False Acceptance Rate in SMO

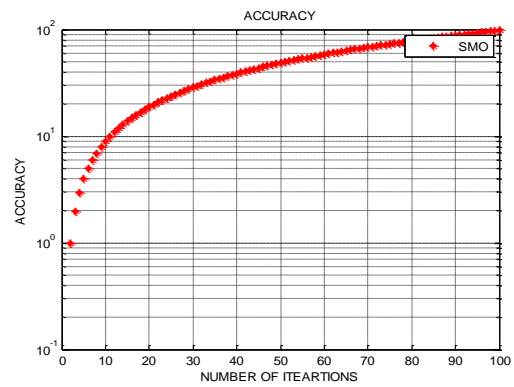

Fig. 4: False Rejection Rate in SMO



Fig. 5: Accuracy



Fig. 6: Comparison between Existing and Proposed Parameters

Table 2.Comparison Matrixes

| Inputs | FRR | FAR | ACCURACY |
|---|---|---|---|
| Happy | 6.781 | .001892 | 99.8107 |
| This is bad. | 5.884 | 0.00425 | 97.568 |
| I am going. | 6.662 | .001758 | 98.995 |

In this fig 6 described the comparison proposed and existing approach. We improve the performance parameters i.e. Accuracy in various sentiments.

## VII. CONCLUSION AND FUTURE SCOPE

This paper describes the investigation problem of studying sentiments in social site via online website which is a significant topic of view knowledge. The basic definition of the comments and sentimental characteristics, we can calculate client's personal features. Then laws are described as the normal constraints of human sentiments based on sentiment area. In this paper we can apply the matrix based technique for feature

extraction means obtained the data percentage of the sentiment and calculate the percentage of the sentiment category. Classify the sentiment using SMO architecture which has shown the performance according to the iterations. In testing form evaluate the parameters like false acceptance and false rejection this is the errors to find the testing part and decrease the error rate because of increase the accuracy and detect the fit category sentiment.

In future, more work is needed on further improving the presentation measures. Sentiment analysis can be applied for novel applications. Although the techniques and algorithms used for sentiment analysis are advancing fast, however, a lot of difficulties in this field of study remain unsolved. The main stimulating aspects exist in use of other languages, dealing with negation expressions; produce a summary of opinions based on product features/attributes, complexity of sentence/ document, handling of implicit product features, etc.

## VIII.    REFERENCES

[1]. Batool, Rabia, et al. "Precise tweet classification and sentiment analysis."Computer and Information Science (ICIS), 2013 IEEE/ACIS 12th International Conference on. IEEE, 2013.

[2]. Fong, Simon, et al. "Sentiment analysis of online news using MALLET."Computational and Business Intelligence (ISCBI), 2013 International Symposium on. IEEE, 2013.

[3]. Prabowo, Rudy, and Mike Thelwall. "Sentiment analysis: A combined approach." Journal of Informetrics 3.2 (2009): 143-157.

[4]. Alessia, D., et al. "Approaches, Tools and Applications for Sentiment Analysis Implementation." International Journal of Computer Applications125.3 (2015).

[5]. Usha, M. S., and M. Indra Devi. "Analysis of sentiments using unsupervised learning techniques." Information Communication and Embedded Systems (ICICES), 2013 International Conference on. IEEE, 2013.

[6]. Wang, Xinzhi, and Xiangfeng Luo. "Sentimental Space Based Analysis of User Personalized Sentiments." Semantics, Knowledge and Grids (SKG), 2013 Ninth International Conference on. IEEE, 2013.

[7]. Wen, Bin, Wenhua Dai, and Junzhe Zhao. "Sentence Sentimental Classification Based on Semantic Comprehension." Computational Intelligence and Design (ISCID), 2012 Fifth International Symposium on. Vol. 2. IEEE, 2012.

[8]. Wu, Ye, and Fuji Ren. "Learning sentimental influence in twitter." Future Computer Sciences and Application (ICFCSA), 2011 International Conference on. IEEE, 2011.

[9]. Vinodhini, G., and R. M. Chandrasekaran. "Sentiment analysis and opinion mining: a survey." International Journal 2, no. 6 (2012).

[10].Taboada, Maite, Julian Brooke, Milan Tofiloski, Kimberly Voll, and Manfred Stede. "Lexicon-based methods for sentiment analysis." Computational linguistics 37, no. 2 (2011): 267-307.

[11].Agarwal, Apoorv, Boyi Xie, Ilia Vovsha, Owen Rambow, and Rebecca Passonneau. "Sentiment analysis of twitter data." In Proceedings of the workshop on languages in social media, pp. 30-38. Association for Computational Linguistics, 2011.

[12].Platt, John C. "12 fast training of support vector machines using sequential minimal optimization." Advances in kernel methods (1999): 185-208.