

# LEVERAGING MACHINE LEARNING AND LEXICON BASED SOCIAL NETWORKING NOTION INVESTIGATION TO DEVELOP A HYBRID METHODOLOGY FOR AN EFFICACIOUS SENTIMENT ANALYSIS, OPINION MINING, POLARITY DETECTION IN SOCIAL MEDIA SITES

Jaskaran Singh Kohli

*Bachelor's in Technology , CSE*

*( SRM Institute of Science and Technology ,Kattankulathur, Tamil Nadu ,India)*

## ABSTRACT

*Sentiment analysis and opinion mining are firmly combined. Broad research work is being completed in these regions by utilizing unique strategies. These approaches distinguish assessments in a given content as positive, negative or unbiased. Tweets, Facebook posts, client remarks about specific subjects and audits concerning the item, programming and films can be the great wellspring of data. Notion Analysis strategies can be utilized on such information by organizations officials for future arranging and estimating. As the information is gotten from numerous sources and it depends specifically on the client which can be from any piece of the world, so the uproar in information is a typical issue, for example, botch in spellings, syntactic mistakes, and wrong accentuation. Diverse methodologies are accessible for notion investigation which can consequently sort and classify the information. These methodologies are fundamentally ordered as Machine Learning based, Lexicon based and Hybrid. A half and half methodology is the blend of machine learning and dictionary-based methodology for the ideal outcomes, this methodology by and large yields better outcomes. In this exploration work, diverse half breed methods and devices have been talked about and investigated from various angles.*

## I INTRODUCTION

The mix of the dictionary based methodology and machine learning approach have improved the grouping execution contrasted with machine learning and vocabulary approach alone. Because of fast increment and globalization of web, a huge number of clients come online day by day and the measure of client created data and information is expanding with a similar pace. The web has turned into the requirement for a few administrations and organizations in our day by day lives. A great deal of printed information is produced by individuals utilizing social sites, for example, facebook and twitter as posts and tweets. A portion of the sites and online journals from these locales to get the estimations of the clients about a specific theme or the criticism about new item or programming and so on. Extraction of estimations from such information can yield important data about a specific theme. A few devices and methods are accessible now days to separate and

characterize the slants from the gave information as either positive, negative and vocabularies as the real wellspring of query for slant classification.

SR#	Author Name	Year	Referen ce	Features	Accuracy
1	Mudinas	2012	[18]	<ul style="list-style-type: none"> <li>- Concept Level Sentiment Analysis System</li> <li>- integrates lexicon and learning based methods</li> <li>- Achieves significantly higher accuracy in Sentiment Polarity Classification</li> <li>- Offers more structured and readable results with aspect-oriented explanation and justification</li> </ul>	89.64%
2	Zhang	2011	[19]	<ul style="list-style-type: none"> <li>- Hybrid Approach, combines lexicon and learning based sentiment analysis classifiers</li> <li>- Unsupervised method except for the initial lexicon</li> <li>- More desirable and practical method</li> <li>- Adaptive to new fashion in language, neologisms and trends</li> </ul>	85.40%
3	Malandrakis	2013	[15]	<ul style="list-style-type: none"> <li>- Based on the Affective Lexicon and Part of speech tag information</li> <li>- Combination of constrained model with Maximum Entropy model trained on external data</li> </ul>	85.80%
4	Balage	2013	[14]	<ul style="list-style-type: none"> <li>- Hybrid Approach, combines lexicon and learning based classifiers</li> <li>- Expression Level and Message Level Classification</li> <li>- Positive, Negative &amp; Neutral Classification</li> </ul>	65.39%
5	Sommar	2015	[16]	<ul style="list-style-type: none"> <li>- Performs better than lexicon based classifier</li> <li>- Effortless Setup</li> <li>- Prospective Performance</li> <li>- Appealing approach for binary classification</li> </ul>	79.67%
6	A. Shoukry	2015	[21]	<ul style="list-style-type: none"> <li>- Hybrid Approach</li> <li>- used specifically for arabic language and egyptian dialect tweets</li> <li>- Combines Lexicon and learning based approaches for classification</li> </ul>	80.90%
7	Serrano-Guerrero	2015	[24]	<ul style="list-style-type: none"> <li>- 15 Different sentiment analysis test</li> <li>- Comparison of different web services</li> </ul>	-
8	Joseph	2013	[34]	<ul style="list-style-type: none"> <li>- Document level, sentence level, Entity level</li> <li>- Language, Datanews and vision part of IBM watson cloud</li> </ul>	73.60%
9	Tang	2014	[37]	<ul style="list-style-type: none"> <li>- Large Scale twitter specific lexicon</li> <li>- uses seed expansion algorithm for expansion of small list of seeds</li> <li>- uses urban dictionary</li> </ul>	85.65%
10	Akshi	2012	[38]	<ul style="list-style-type: none"> <li>- Hybrid Approach, uses corpus and dictionary based approaches</li> <li>- uses combination of adjectives along with verbs and adverbs</li> </ul>	
11	Syed-Ali	2013	[39]	<ul style="list-style-type: none"> <li>- Hybrid approach, uses sentiment lexicon and linear SVM</li> <li>- effective than the unigram baseline models</li> </ul>	89.13%
12	Mullen	2004	[40]	<ul style="list-style-type: none"> <li>- uses SVM as base</li> <li>- favorability measures for phrases and adjectives</li> </ul>	86.00%
13	Xiang	2014	[42]	<ul style="list-style-type: none"> <li>- based on topic based sentiment mixture model</li> <li>- un-supervised approach</li> </ul>	69.7 F-score
14	Xianghua	2013	[43]	<ul style="list-style-type: none"> <li>- used for automatic discovery of aspects being discussed in chinese social reviews</li> </ul>	91.23%

today contain the area of client's remarks or input so important data can likewise be taken These vocabularies have predefined semantic introductions that are later contrasted and the information informational collection for nonpartisan. Gadgets and strategies from Lexicon based strategy uses space unequivocal word reference ouping as cleared up by Machine learning set up together philosophy concerning the following hand seek after the directed learning computations, for instance, Naive Bayes and Support Vector Machine to make the readiness instructive gathering . By then dependent on this readied dataset the wellsprings of information are broke down and named either positive, negative or some other appraisal. The Hybrid strategy which uses the mix of both vocabulary based system and machine learning approach. The real focus of this blend is to yield the best and immaculate outcomes utilizing the viable once-over of capacities of both vocabulary and machine learning based frameworks, and to beat the inadequacies and deterrents of the two approaches. Different specialists have cemented specific vocabulary and machine learning based procedure to improve and astounding crossbreed contraptions. In this examination, we will consider, dissect and analyze diverse half and half apparatuses and systems for feeling characterization and will talk about various capabilities and correctnesses of the contemplated methodologies.

## HYBRID TOOLS AND TECHNIQUES

### A. pSenti

pSenti is an idea level opinion investigation apparatus that was displayed by , it joins dictionary and learning based supposition order strategies. When contrasted with the unadulterated dictionary based techniques pSenti accomplished more noteworthy precision in opinion quality identification and extremity arrangement. Then again, when the device was thought about against unadulterated machine learning based techniques it yielded somewhat bring down precision. Broad analyses on two diverse datasets i.e., CNet Software Reviews Dataset and IMDB Movie Reviews Dataset for the assessment of the proposed methodology were performed. Learning based methodology utilized in the proposed technique isn't in charge of modest undertakings like change of notion esteems or opinion words recognition however it is additionally in charge of assessment of all parts of notion framework.

The primary part of the framework estimates the given obstinate content and gives the yield as far as aggregate assumption, for example, client input. The last outcomes are appeared with a genuine esteemed score between - 1 and +1 that can be changed as either positive/negative or into a score between 1-5 stars in a last stage. Favorable circumstances of the proposed methodology are that the framework can be stretched out by including new etymological principles or assumption vocabulary can be extended at any occurrence/level. The proposed framework isn't touchy to the adjustments in the subject. It works superior to SentiStrength [5] and dictionary just also yet its precision is somewhat lower than adapting as it were.

### B. Combining Lexicon Based and Learning Based systems for Twitter Sentiment Analysis

For substance level supposition examination, utilized an all-encompassing word reference based system. In the first place, they got extra decided marker, for example words and pictures, by applying Chi-square test on results assembled from the dictionary based technique. Extra troublesome tweets were related to the assistance of new unyielding markers. For segments in the starting late seen tweets, a supposition gathering figuring is utilized to circulate thought uttermost point scores. The result of the vocabulary framework is fundamentally the availability information for the classifier and the entire technique has no manual checking next to test set. This examination utilized five datasets dependent on the demand parts Obama, Harry Potter, Tangled, iPad and Packer. Proposed system accomplished 85.4% precision on the five datasets utilized in this examination. In the proposed framework (LMS) a relative enhancement over the vocabulary based methodology was seen. In any case, it performed dynamically terrible on the other hand with the unadulterated learning-based framework yet having extraordinary position that it doesn't require pre-checked information. In this manner, the proposed framework is direct in execution at any rate cost some execution.

### C. SAIL

Another mixture system was produced by. This examination proposed a framework for twitter

and SMS opinion investigation dependent on various leveled display, full of feeling vocabulary and a dialect demonstrating approach. It is seen that dialect demonstrate was bad alone but rather an enhanced execution was seen when utilizing with vocabulary based model. The various leveled demonstrate demonstrated exceptionally effective notwithstanding utilizing the n-grams, full of feeling evaluations and grammatical feature. The proposed instrument utilizes a full of feeling vocabulary that was precipitously produced from huge corpora of crude web information. Words and bigrams are utilized for full of feeling evaluations computations and insights. To the extent the unconstrained information is concerned the vocabulary models were joined with a learning classifier that depends on the Max-Ent dialect models that are basically educated on an enormous outer dataset. These two order techniques for estimation examination are then consolidated to define the last outcomes. The mix of the two turned out to be full of feeling and yielded better outcomes.

#### **D . NILC\_USP**

The specialists in depicted NILC\_USP framework in SemEval-2013 and proposed a trio order process that consolidates three characterization approaches, for example, the standard based methodology, the vocabulary based methodology, and the machine learning based methodology. The proposed calculation has five stages.

**Standardization:** The initial step is standardization of the given info dataset, it can likewise allude as pre-preparing, it essentially cleans and standardizes the information content, and this progression performs the following tasks.

- Hashtags, URLs, and notification are figured in the solid plan of codes
- Emoticons are ordered by their physical appearance as either merry, lamentable, chuckle, etc and doled out with explicit codes
- Exaltation signals are recognized and stepped, for instance, different signs of objection
- Misspelled words are changed
- Part-of-discourse naming is performed

**The Learning Based Model:** At this stage, it utilizes the starting late referenced SciKit Learn framework, that gives a pipeline structure and draws in a few changes to be connected with the data and plan it as required, making the last model that sorts out the data. By displacing the exhibiting some piece of the pipeline structure it might be attempted with different classifiers to evaluate and find which classifier yields the best and impeccable results. Following three classifiers were tried by the researchers Multinomial Naive Bayes, Bernoulli perfect Bayes, and SVM.

**Dictionary-based Classifier:** In the proposed structure the vocabulary given by SentiStrength was used. This lexicon gives a vocabulary of sentiments, an emoticons once-over, invalidation and boosting words list. The semantic presentation of each word in the given substance is resolved in the proposed computation. The furthest point of the word is reduced if the words are invalidated, in like way the limit is extended when the words are expanded, the classifier names the substance as positive, negative or unprejudiced.

**Machine Learning Classifier:** Labeled models are used by the Machine learning classifiers to learn and gather the given substance, SVM count given by CLiPS configuration was used. In the proposed model, the pack of words, linguistic shape sets and the nearness of negation in the sentences were used as the rundown of capacities by the classifier.

The consequences of this examination demonstrated that the half breed classifier approach could enhance results dependent on the benefit of numerous notion investigation strategies over guideline-based, dictionary-based and machine learning techniques.

**E.Consolidating Lexicon based and Learning based methodologies for enhanced execution and comfort in feeling grouping proposed a crossover way to deal with enhance the execution of feeling investigation process.** The programming dialect picked for the usage of this calculation was Python. The proposed calculation is made out of three stages after pre-handling, the initial segment alludes to the dictionary based model, and it manages to find the ideal parameters for the classifier. While the second part alludes to the learning based model and manages the investigation of the model that performs better. In conclusion, the third part alludes to the mixture show that dissects and chooses the ideal MID proportion.

**The Learning Based Model:** At this stage, it utilizes the starting late referenced SciKit Learn framework, that gives a pipeline structure and draws in a few changes to be connected with the data and plan it as required, making the last model that sorts out the data. By displacing the exhibiting some piece of the pipeline structure it might be attempted with different classifiers to evaluate and find which classifier yields the best and impeccable results. Following three classifiers were tried by the researchers Multinomial Naive Bayes, Bernoulli perfect Bayes, and SVM.LDA exhibit for a seeing subject. Test outcomes showed that the proposed model increment perfect point allotting results, just as assistants in the improvement of end examination accuracy.