

**Sentiment Analysis of Tweets using Apache Flume and Spark**Mayur More<sup>1</sup>, Monika Darekar<sup>2</sup>, Komal Deshmukh<sup>3</sup>

<sup>1</sup> Computer Engineering, AISSMS-  
IOIT, mmore1315@gmail.com <sup>2</sup> Computer Engineering, AISSMS-  
IOIT, darekarmonika24@gmail.com <sup>3</sup> Computer Engineering, AISSMS-  
IOIT, komaldeshmukh1995@gmail.com

**Abstract–**

Through social media the users can share their thoughts with friends, family, and colleagues, and it also gives the user a platform to talk and communicate on their favorite topics. This “unstructured” conversation can give businesses valuable insight into how consumers perceive their brand, and allow them to actively make business decisions to maintain their image. With a rapid increasing of data of sentiments in social media on web has lead the researchers into increased interests regarding opinion mining and sentiment analysis. However, Sentiment Analysis is now considered as a Big Data task due to the large amount of social media available on the web.

To find a technique such that it can efficiently perform sentiment analysis on big data sets was the main focus of this research. In this paper, Hadoop Apache ecosystem’s data ingestion tool was used to perform Sentiment Analysis on the large sets of data consisting tweets and stream processing with Spark. Using this technique the experimental result shows very good efficiency in handling big data sets of sentiment.

**Keywords-** Sentiment Analysis, Opinion Mining, Apache Spark, Apache Flume, Machine Learning

**I. INTRODUCTION**

Sentiment analysis is the field which deals with computational treatment of sentiment, opinion and subjectivity in text. Sentiments can be described as ideas, opinions, judgments or emotions prompted by emotions. World Wide Web has completely changed the way of expressing people’s views. Now a day’s people are expressing their views and thoughts through online blogs, discussion forms and also some Online applications like Facebook, Twitter, etc. If we take Twitter as an example nearly 1 TB of text data is generated within a week in the form of tweets. So, by this it is understood clearly how this Internet is changing the way of living and style of people. Among these tweets can be categorized by the hash value tags for which they are commenting and posting their tweets. So, now many companies and also the survey companies are using this for doing some analytics such that they can predict the success rate of their product.

To calculate their views is very difficult in a normal way by taking these heavy data that are going to generate day by day. If we consider getting the data from Twitter one should use any one programming language to crawl the data from their database or from their web pages. Coming to this problem here we are collecting this data by using Apache Ecosystem’s Data ingestion Tool known as Flume and then processing same data using Spark.

**Levels of sentiment analysis**

**Word level:** Determines whether an expression is neutral or polar and then disambiguates the polarity of the polar expressions.

**Sentence level:** Categorization attempts to classify positive and negative sentiments for each or whether a sentence is subjective or objective.

**Document level:** There are two kinds of approaches

1. Term-counting approaches and
2. Machine learning approaches.

Term-counting approaches usually involve deriving a sentiment measure by calculating the total number of negative and positive terms.

Machine learning approaches recast the sentiment classification problem as a statistical classification task.

## II. PROBLEM DEFINITION

### A. EXISTING SYSTEM

As per the older ways of crawling data and performing the sentiment analysis on those data was done with the help of Java and RDBMS. People were using some coding techniques for crawling the data from the Twitter where they could extract the data from the Twitter web pages by using some code that might be written either in JAVA, Python etc. For those they were using libraries that were provided by the Twitter organization. By using this they were crawling the data that they wanted particularly.

After getting it the raw data was filtered by using some old techniques and also they found out the positive, negative and moderate words from the list of collected words in a text file. All these words would be collected to filter out or do some sentiment analysis on the filtered data. These words can be called as a dictionary set by which they perform sentiment analysis. Also, after performing all these things they wanted to store these results in a database and coming to here they used RDBMS where they were having limitations in creating tables and also accessing the tables effectively.

### B. APPROACHES

**Lexicon** based approach can be used with an assumption work that the collective polarity of a document or sentence is the sum of polarities of the individual words or phrases.

**Keyword spotting** approach can be used to classify the text by affect categories based on the presence of unambiguous affect words such as sad, bored, happy, afraid.

**Bagging**, In this approach, relationships between the individual words are not considered and a document is represented as a mere collection of words. To determine the overall sentiment, sentiments of every word is determined and those values are combined with some aggregation functions.

**Statistical classification** approach is used on elements from machine learning such as support vector machine, latent semantic analysis, "bags of words" and Semantic Orientation-point wise mutual information.

**Lexical affinity** not only detects obvious affect words, It also assigns arbitrary words a probable "affinity" to particular emotions.

