



## ABSTRACT:

Opinion Mining is the field of study that analyses people's opinions, views and emotions from text. It is one of the most active research areas in natural language processing and is also widely studied in data mining, Web mining, and text mining. The growing importance of sentiment analysis coincides with the growth of social media such as reviews, forum discussions, blogs, micro-blogs, Twitter, and social networks. And with the help of technological advances, we now have a huge volume of opinionated data recorded in digital form for analysis. Opinion Mining systems are being applied in almost every business and social domain because opinions are central to almost all human activities and are key influencers of our behaviors. Our beliefs and perceptions of reality, and the choices we make, are largely conditioned on how others see and evaluate the world. For this reason, when we need to make a decision we often seek out the opinions of others. This is true not only for individuals but also for organizations. Twitter is an active medium to express opinions and views of the public in general. We present a system that will classify the most recent tweets extracted according to a search query into positive and negative classes, thus generating an overall opinion of the public. We will be using classification algorithms in order to classify the tweets into the positive and negative categories.

## INTRODUCTION:

Opinion Mining, commonly referred to as sentiment analysis is used to extract the information from a text with respect to its polarity to the given subject. The goal of such a system is to effectively process subjective information. It identifies opinion bearing words, which are usually adjectives which define the nature and polarity of the text. In general, opinions can be expressed about anything, a product, a service, a topic, an individual, an organization or an event. It allows classification of the text or feature, i.e. whether the polarity of the opinion is positive, negative or neutral. Given a set of tweets, the system classifies them into positive and negative classes. This is clearly a classification learning problem. But it is also different from topic based text classification. In topic based classification, topic related words are important. However, in this case, opinion related words are important for us, e.g. good, bad, worst, excellent, horrible, etc. On the basis of these opinion words and proper weightage to the features, one can obtain the overall sentiment of the tweet related to the topic.

## METHODOLOGY:

In this project, we will be doing Opinion Mining on tweets pertaining to a particular search topic. The most recent tweets related to the topic will be searched on twitter, after which their polarity will be classified as positive or negative. As a result, the current overall sentiment among the public related to that particular topic can be obtained by applying this on a large set of tweets extracted from the public at a particular time. For the classification task, a Twitter Sentiment Corpus containing 1.6 million classified tweets is used as a training set to build our classifier. It is to be used to build a classifier which will classify the tweets searched as positive and negative. On the basis of this classified data on the search dataset, we need to determine the overall public opinion of the search query topic. Given a tweet **T**, searched from the most recent tweets from the search query, classify it into either of the two classes – positive or negative. Repeat the process for the entire set of current tweets and thus get the overall opinion of the search query among the public.

## Literature Survey:

### OPINION MINING TECHNIQUES:

#### 1. PROBABILISTIC CLASSIFIERS:

Probabilistic classifiers use mixture models for classification. The mixture model assumes that each class is a component of the mixture. Each mixture component is a generative model that provides the probability of sampling a particular term for that component. Three of the most famous probabilistic classifiers are discussed here.

##### 1.1 NAÏVE BAYES CLASSIFIER (NB):

The Naïve Bayes Classifier is the simplest and most commonly used classifier. It computes the posterior probability of a class based on the distribution of words in a document.  $P(\text{label})$  is the prior probability of a label or the likelihood that a random feature set the label.  $P(\text{features}|\text{label})$  is the prior probability that a given feature set is being classified as a label.  $P(\text{features})$  is the prior probability that a given feature set is occurred.

##### 1.2 BAYESIAN NETWORK (BN) :

The main assumption of the NB classifier is the independence of the features. The other extreme assumption is to assume that all the features are fully dependent. This leads to the Bayesian Network model which is a directed acyclic graph whose nodes represent random variables, and edges represent conditional dependencies. BN is considered a complete model for the variables and their relationships.

### RESULTS AND DISCUSSIONS:

METHOD	ACCURACY(%)	TRAINING DATA SIZE (number of tweets)	TEST DATA SIZE (number of tweets)
Naïve Bayes <sup>7</sup>	69.97%	20,000	500
Maximum Entropy	73.36%	5,000	500

The training for Maximum Entropy is much more resource consuming than that of Naïve Bayes<sup>7</sup> because of the exponential calculation of Entropies of its features. Hence, training takes a greater time for Maximum Entropy Model. Since it doesn't consider independence of features, unlike Naïve Bayes<sup>7</sup>, it is likely to give better results than that of Naïve Bayes, as evident from the results of the Stanford Test dataset.

TOPIC	POSITIVE SENTIMENT (%) NB	NEGATIVE SENTIMENT (%) NB	POSITIVE SENTIMENT (%) ME	NEGATIVE SENTIMENT (%) ME	NET SENTIMENT
Phil Hughes	42.645	57.354	32.713	67.286	NEGATIVE
Ebola	44.390	55.610	36.569	63.430	NEGATIVE
Interstellar	56.61	43.39	59.72	40.28	POSITIVE

