# Tweet Segmentation and Enhancement of Tweets

Sonam Meshram<sup>1</sup>, Hirendra Hajare<sup>2</sup>

<sup>1</sup>Gondwana University, Ballarpur Institute of Technology, Ballarpur, India

<sup>2</sup>Professor, Gondwana University, Ballarpur Institute of Technology, Chandrapur, India

Abstract: Twitter is a biggest connecting site that includes various types of users. Many users share their data and it is updated sites so data should be maintained properly and accessing in proper way. Hence mining algorithm helps to managing data. Many application such as Information Retrieval and Natural Language Processing contains some errors and short nature of tweets, hence to recover of such type of tweets tweet classification is used. Data mining algorithm used in the classification of tweets hence it is easily access and easy to understand.

Keywords: Twitter, tweet segmentation, named entity recognition, k-means algorithm, support vector machine algorithm.

#### 1. Introduction

Twitter is a type of social connecting media, has been tremendous growth in the recent years. It includes the all type of users and it has attracted great interests from both of industries and academic field. The twitter stream is monitored and to collect then understand users opinions about the organization. It is required to detect and response with such targeted stream, such application requires a good named entity recognition (NER). [1], [2], [9]. Twitter is big source of continuously and instantly updated information. The Social networking sites are updated and most important communication channel with its capability of providing the most up-to-date and news oriented information. The targeted twitter stream to focus the tweet segmentation and its arrangement. Twitter is a micro blogging service that founded in the 2006 and it is one of the most popular and it is fastest broadcasting sites, growing online social networking sites with more than 190 million Twitter accounts. The social networking sites includes various types of peoples and hence data can be share one to another that time data must be safe and it is nothing but the malicious data or message to send another user. Hence the targeted stream which helps to remove such type of spam or messages and it is preserving from the spam.

Twitter is a social networking sites that enables users to send and red short 140-charactes messages called as tweets. Each and every user wants to their data must be safe and prevented from the hackers. Many social communities thought there data must be spam free means that errors free. The error can be grammatical also and the spam data can be affected your system and hence that malicious data harmful to the system and that's why it is detected properly and preserving that such type of spam and hence system must be error free[3]. The targeted twitter proposed system of tweet segmentation helps to removing the errors and protected from the illegal messages. Hence it is used for the improving the quality of tweets. The social networking sites which will be much updated day by day and that's why the data should be effective in nature. The data mining concept very useful in the targeted twitter. Data mining is an interdisciplinary subfield of computer science. It is the computational process of discovering patterns in large data sets involving methods .The overall goal of the data mining process is to extract information from a data set and transform it into an understandable structure for further use. Data mining is a collection of tools and techniques. It is one of several technologies required to support customer-centric enterprise.[8]. It is useful in the tweet segmentation and with the help of data mining algorithm the data must easily maintained and easy to access.

## 2. Related Work

Twitter includes millions of users and data must be up-todate. The novel framework for tweet segmentation called as HybridSeg. The local linguistic features are more reliable for learning local context and high accuracy is achieved named entity recognition by using segment based part-of-speech (POS) tagging [1],[10]. The Chao Yang focuses on the empirical study and new design for twitter spammer's fighter. With the help of machine learning detection techniques features and the goal is to provide the first empirical analysis of the evasion tactics and in-depth analysis of those evasion tactics [3]. Make a comprehensive and empirical analysis of the evasion tactics utilized by Twitter spammers. The online social networking sites such as twitter and Facebook are now part of many people's daily routine and hence it is updated. Spammers have utilized Twitter as a new platform to achieve their malicious goals such as sending spam messages, spreading malware, hosting botnet and control (C&C) channels and performing other illicit activities [3]. The named entity recognition (NER) used in twitter stream for the monitoring and response to the stream. The unsupervised NER system known as TwiNER. First step is that global context obtained from the Wikipedia and partition of tweets by using dynamic algorithm [2]. The TwiNER system is the first to exploit both the local context in tweets and the global context from the World Wide Web together for named entity recognition task in twitter [2]. An experimental study of the named entity recognition in tweets that focuses on the demonstrating the tools for part-of-speech (POS) tagging. Showing that benefits of features generated from T-pos and T-chunk in the segmenting named entities [4]. In corpus linguistics, part-of-speech tagging or POST tagging or wordcategory disambiguation, is the process of marking up a word in a text or corpus as corresponding to a particular part of

speech, based on both its definition and its context. The new approach for twitter user modeling and tweet recommendation by using named entities and its extracted from the tweets [5]. The previous work in that the named entity extraction (NEE) and linking for tweets it is the hybrid approach. The named entity extraction is for locate phrases in the text that represent names of persons. The approaches is that named entity generation and linking then its filtering [6].

### 3. Tweet Segmentation

The tweet segmentation is the task of twitter stream. The goal of work is to classify tweets into section hence it can be understand easily. The previous work of the tweets is that the tokenization hence named entity recognition is used. Both tweet segmentation and named entity recognition are considered the subtask of the Natural Language Processing (NLP) [1]. The segmentation is to split the tweet segmentation is that the tweet is to be split into consecutive segments. Tweet segmentation it is important job of the previous paper. Twitter is a social networking sites and it contains the millions of people interact each other. Hence the data should be maintained properly. Tweets are very high time-sensitive nature so that many phrases like "she eatin" cannot be found in external knowledge bases. Observe that tweets from many official accounts of organizations and advertisers are likely well written. Then the named entity recognition helps with the high accuracy of tweets [1], [5]. Hence the overall study about the twitter and there challenges there is an need to be a segmented manner of data. The property of named entities in the targeted tweet stream and it is a collectively from a batch of tweets in unsupervised manner. Basically, let T be the collection of the tweets that posted in the targeted twitter stream within the one fixed time interval. For example, India is the biggest country. That sentence is to be segmented is that (India) | (is the) | (biggest) | (country). The job of tweet segmentation is that the data is to be splited [1]. The traditional named entity recognition method is the well formatted documents heavily depends on the phrases local linguistic features.

The capitalization and part of speech is the previous work of the tweets [2]. The previous work related to the tweet segmentation is focuses towards by using the algorithms that includes the random walk( RW) and the part-of-speech (POS) . The co-occurrence of names entities in the twitter stream by applying the random walk and the another part-ofspeech tags of the constituents words in segments. That the segment are likely to be a noun phrase are considered as a named entity [1]. To overcoming the some features of the related tweets and hence another features can be applying and tweets are in error free and preserving from the spam. Whenever the tweets can be segmented then some grammatical errors are present in such phrases and hence overcoming in the targeted twitter stream apply algorithm and named entity concept for that the tweet segmentation.

### 4. Tweet Classification

The tweet segmentation is the split the tweets. The related work of the tweets hence it is contained large number of

some features which will be absent hence it is to be implemented hence that features is to be added in this work of tweets. The classification is distribute the term or data. Hence the tweet can be categorizes some manner that should be related to that the particular tweet phrases. Tweet segmentation is the task to divides the tweet in some segmented manner not in the word manner, because the study of that segment based are better than the word based. Using the clustering algorithm to improve the nature of the tweets. Hence this paper to enhance the features of tweet by using Kmeans algorithm. The data mining is handled the large number of data. Data mining is the exploration and analysis of large quantities of data in order to discover meaningful pattern and rules. The goal of data mining is to allow a corporation to improve its marketing, sales and customer support operations through better understanding of its customer. The data mining algorithm is to be implemented for that the commercial application purposes. The techniques is to be borrowed from the statistics, computer science and machine learning research [8.]

The data mining algorithm is used that is k-means algorithm. Basically cluster analysis is one of the major data analysis method and the k-means clustering algorithm is mainly used for the various applications. For generating and the collecting data the growth of database has been large day by day. Hence the practically impossible to extract useful information from them by applying conventional database analysis techniques. That of the effective mining method are essential to extract information from large databases [7]. K-means clustering algorithm which has likely the nearest neighbor that depends on geometric interpretation of metric ideas used in k-means. It brings general topic that related association and distance. K-means not only the algorithm but also automatic cluster detection [8]. The idea is that to classifying the given set of data into the k number of disjoint cluster and then that the value of k is the fixed in advance. The algorithm can be categorize into two phases , the first phase is that defines k centroids one for the each cluster. The another phase is to take each point related to the given data set and it associate it to the nearest centroid [7]. The k-means algorithm is very helpful in the targeted stream because of that the tweet segmented are classified. Hence the term classification means the segmented tweet can be detected and it can be categorize in the specific region. Then by applying the algorithm such as k-means and it is a clustering algorithm and it is used in the detection also. The classification of tweets is that the it can be divided and hence particular tweet is to be section wised distributed. The previous work of the tweets is that the there is no any classification of section wised, that overcoming such type of tweets in this paper. With the help of data mining algorithm the data can be handled properly and hence it can be classified region wised. The many kinds of messages are to be exchanges hence it is need to be security of that the tweets. Spam is the illegal type of massage hence some features can be implement to that type of short term types of tweets. In the social networking sites such as twitter includes various types of users, each and every person can be posted there tweets in any field such as it should be related to sport, entertainment, education, commerce and current event also. The targeted twitter stream that segmented the tweet and then it should be categorized in that the particular section

by using this algorithm effectiveness of tweets is to be improved. In data processing, filtering of all data will be done. The punctuation, symbols, deletion of email ids etc. will be removed which is not important. Topic allocation means allocating data in the form of field like we do in our PC, we allocate movies as per the category i.e. Hollywood movies, Bollywood movies, animated movies etc. So like this there is need of allocating topics category wise. This work will be done in proposed work. Topic detection will be done after topic allocation for that topic K-means will be used. Topic K-means will use for feature extraction [8].

Another algorithm is used in this work is the support vector machine (SVM) is widely used in object detection & recognition, content-based image retrieval, text recognition, biometrics, speech recognition. In machine learning, support vector machines are supervised learning models with associated learning algorithms that analyze data and recognize patterns, used for classification and regression analysis. SVM also support vector networks [11]. Named entity linking (NEL) is the task of that exploring which correct person, place and events is referred to by a mention. The linking approach to determine the particular named entity and the support vector machine to predict which candidates are true positions and which one are not [6]. The idea of the targeted twitter stream is that whenever the data can be segmented hence tweet segmentation can be performed then that tweet are classified next job of this task is that the current event detection mechanism should be performed with the help of the support vector machine algorithm. Social networking sites includes the user interface features that's why the targeted stream also present the user interface characteristics and hence the many users can be interconnected to each other and exchanged there information. The main objective of this system is that To Classification of tweets, it provides to removing the noisy tweets then to identify the spam word and preserve this. It provides Current event detection. The concept of named entity ranking that is research in the previous work and that can be named entity play important role in all of the tweet segmentation [1], [2], [4]. Hence by applying the mining rules the accessing data easily and improves the efficiency of targeted stream. With the help of tweet segmentation and its classification that improves the targeted twitter stream. The support vector machine algorithm is very important in this work of tweets the task of the tweet segmentation of tweets is to be segmented means the tweets can be split. The main task of this work is the classify given tweets. Social networking sites include various people and that type of data is to be managed in the specific way. Most of the tweets is related to the some special field hence the another user has been seen in that of respective field. That type of work is to be maintained in this work and hence by using data mining algorithms the features of tweets is to be improve and the tweets enhancement is to be maintained.

### 5. Conclusion

The tweet segmentation and classification helps to preserving the semantic meaning of tweets. This paper proposes a new tweet classification which helps to improve the accuracy and efficiency of tweets and hence the tweet shows in specific region. The segment based tweet it is better than that of another word based. The current event detection is also helpful for the traffic analysis. For future work The graphical analysis and improves again the segmentation analysis. The data can be preserving from the spam and hence the tweets are secured nature.

#### References

- [1] Chenliang Li, Aixin Sun, Jianshu Weng and Qi Hi, "Tweet Segmentation and Its Application to Named Entity Recognition," IEEE, vol. 27, No. 2, February 2015.(conference style).
- [2] Chenliang Li, Jianshu Weng, Qi Hi, Yuxia Yao, Anwitaman Datta, Aixin Sun and Bu-Sung Lee, "TwiNER: Named Entity Recognition in Targeted Twitter Stream," School of Computer Engineering ,Singapore, August 2012.(journal style)
- [3] Chao Yang , Robert Harkreader and Guofei Gu, "Empirical Evluation and New Design for Fighting Evolving Twitter Spammers," Member, IEEE, vol. 8, No. 8, August 2013.(conference style)
- [4] Alian Ritter, Sam Clark, Mausam and Oream Etzioni, "Named Entity Recognition in Tweets: An Experimental Study," Computer Science and Engineering University of Washingtn, USA. (technical report style)
- [5] Deniz Karatay and Pinar Karatay, "User Interest Modeling in Twitter with Named Entity Recognition," Turkey, vol. 1395, 18<sup>th</sup> May 2015.(technical workshop report style)
- [6] Mena B. Habib , Maurice van Keulen and Zhemin Zhu, "Named Entity Extraction and Linking Challenges," University of Twente Microposts , 7<sup>TH</sup> April 2014.(technical workshop report style)
- [7] K. A. Abdul Nazeer and M. P. Sebastian, "Improving the Accuracy and Efficiency of k-means Clustering Algorithm," London, U.K., vol. I, July 2009.(conference style)
- [8] Wiley, "Data Mining Techniques," second edition.(book style)
- [9] David Nadeau and Satoshi Sekine, "A survey of named entity recognition and classification," National Research Council Canada / New York University.(report style)
- [10] Chenliang Li, Aixin Sun, Jianshu Weng, and Qi He"Tweet Segmentation and Its Application to Named Entity Recognition," Ieee Transactions On Knowledge And Data Engineering, 2013.(conference style)
- [11] Hiep-Thun Do, Nguyen-Khang Pham, Thanh-Nghi Do,"A SIMPLE,FAST SUPPORT VECTOR MACHINE ALGORITHM FOR DATA MINING," Fundamentl and Applied IT Reaserch Symposium 2005.(conference style)