# A Review Paper on: The Exploiting Concept of Twitter NER of Segmentation

Miss. Anjum Inamdar [1], Miss. Kartiki Wahatole [2], Miss. Vishakha Shinde [3],
Miss. Harshata Tohake[4] , Prof.G.S.Deokate[5], Prof.B.S.Kurhe [6]

[1]*Student,Computer,SPCOE, Pune,Maharashtra, India,*
[2]*Student,Computer,SPCOE, Pune, Maharashtra, India,*
[3]*Student,Computer,SPCOE, Pune,Maharashtra, India,*
[4]*Student,Computer,SPCOE, Pune,Maharashtra, India,*
[5]*Assi.Professor, Computer,SPCOE, Pune,Maharashtra, India,*
[6]*Assi.Professor, Computer,SPCOE, Pune,Maharashtra, India.*

## ABSTRACT

*The information given on the social media is delivered to each and every person within some fraction of the second. The social networks such as Facebook or Twitter is the platform which is being widely used for posting what is happening?, what are the crimes happened? , what steps had been taken against that crime? , and it also give individual to express each and every emotion on such platform. The opinions changes to person to person and the posts may create a different impact on the individual. So the reaction may be positive or negative. The impact of the negative thought may be so strong that may result into social unrest. The social network some-times used for the planning the social unrest or gather people for such activity for example the candle march after some crimes, is the social unrest which has been mostly happened in India. There must be some tool or applications that can be used for detect such post and predict the Railway issues. Each and every post or action to be analyzed and then prediction is done whether the unrest will happen or not. Due to prediction of unrest before it happens will really help the investigator/police to prepare for that situation or to completely stop such activity. There is a need of such application which will predict the railway issues and this prediction will really help to get the details such as the location or the date.*

**Keyword : -** *Facebook, Activity detection, Social networks, Tweeter.*

## 1.INTRODUCTION

Twitter, is an  new type of social media, It has seen tremendous growth in recent years. It has attracted by both industry and academic. Many private, public organizations have reported to monitor Twitter stream which collect and also understand opinions of the users about the organizations[11]. Due to the very large volume of tweets published each and every day, its practically infeasible and it is unnecessary to listen and also monitor the whole Twitter stream. Therefore, the targeted Twitter streams are  monitored instead; each such stream contains tweets that potentially satisfy some information needs of the monitoring organization. Targeted Twitter stream is usually constructed by filtering[7][13] tweets with user-defined selection criteria depends on the information needs. For example, the criterion could be a region  so that users opinions from that particular region are collected and monitored; it could also be one or more predefined keywords so that opinions about some particular events/topics/products/services can be monitored. The idea is to  segment an individual tweet into a sequence of consecutive phrases, each of which appears more than chance. After removing the stop words, a tweet My shoes are gg to compete in the youth olympic games sailing competition[1][3][5]. It just needs a mast and a rudder is segmented into seven parts. In the solution for tweet segmentation. Given an individual tweet t Ti, the problem of tweet segmentation is to split t into m consecutive segments, t =s1s2...sm; each segment contains one or more words. To obtain the optimal segmentation. A high stickiness score of

segment s indicates that it is not suitable to further split segment s, as it breaks the correct word collocation. In other words, a high stickiness value indicates that a segment cannot be further split at any internal position. If the word length of tweet t is L, possible segmentations. It is insufficient to iterate all of them and compute their stickiness[2][13]. Twitter has become one of the most important channels for people to and, share, and disseminate timely information. As of March Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for components of this work owned by others than ACM [15][10]or the author must be honored. To copy otherwise, or republish to post on servers or to redistribute to lists, requires prior specific permission and/or a fee. There are more than 140 million active Twitter users with over 340 million tweets posted in a 1 day. Due to its large volume of timely information generated by its millions of users, it is imperative to understand tweets language for the tremendous downstream applications like named entity recognition (NER)[5][8], event detection and summarization, opinion mining, sentiment analysis[3].Status Messages posted on Social Media websites such as Face book and Twitter present a new and challenging style of text for language technology due to their noisy and informal nature. Like SMS, tweets are particularly there. Yet tweets provide a unique compilation of information that is more up to-date and inclusive than news articles, due to the low-barrier to tweeting, and the proliferation of mobile devices[4][11]. One main challenge is the lack of information in a single tweet, which is rooted in the short and noise-prone nature of tweets[8]. To collectively extract social events from multiple similar tweets using a novel factor graph, to harvest the redundancy in tweets, i.e., the repeated occurrences of a social event in several tweets[1]. Twitter has several characteristics which present unique challenges and opportunities for the task of open-domain event extraction[9].

## 2. LITERATURE SURVEY:

Many applications in Information Retrieval (IR) and Natural Language Processing (NLP) suffer severely from the noisy and short nature of tweets. In this paper, propose a novel framework for tweet segmentation in a batch mode, called Hybrid Segment. By splitting tweets into meaningful segments, the semantic or context information is well preserved and easily extracted by the downstream applications. Hybrid Segment finds the optimal segmentation of a tweet by maximizing the sum of the stickiness scores of its candidate segments. The stickiness score considers the probability of a segment being a phrase in English and the probability of a segment being a phrase within the batch of tweets .Present a novel 2-step unsupervised NER system for targeted Twitter stream, called TWINER. In the _rst step, it leverages on the global context obtained from Wikipedia and Web N-Gram corpus to partition tweets into valid segments (phrases) using a dynamic programming algorithm. Each such tweet segment is a candidate named entity. It is observed that the named entities in the targeted stream usually exhibit a gregarious property, due to the way the targeted stream is constructed. In the second step, TWINER constructs a random walk model to exploit the gregarious property in the local context derived from the Twitter stream[1]. A novel framework for tweet segmentation in a batch mode, called Hybrid Segment. Hybrid Segment incorporates local context knowledge with global knowledge bases for better tweet segmentation. Hybrid Segment consists of two steps: learning from the shelf weak NERs and learning from pseudo feedback. In the first step, the existing NER tools are applied to a batch of tweets. The named entities recognized by these NERs are then employed to guide the tweet segmentation process. In the second step, Hybrid-Segment adjusts the tweet segmentation results iteratively by exploiting all segments in the batch of tweets in a collective manner. Experiments on two tweet datasets show that Hybrid Segment significantly improves tweet segmentation quality compared with the state of the art algorithm[2].An experimental study, re-building the NLP pipeline beginning with part-of-speech tagging, through chunking, to named-entity recognition. Novel T-NER system doubles F1 score compared with the Stanford NER system. T-NER leverages the redundancy inherent in tweets to achieve this performance, using Labeled LDA to exploit Freebase dictionaries as a source of distant supervision. Labeled LDA out performs co training, increasing F1 by 25 percent over ten common entity types[3].

To combine a K-Nearest Neighbors (KNN) classifier with a linear Conditional Random Fields (CRF) model under a semi-supervised learning framework to tackle these challenges. The KNN based classifier conducts pre-labeling to collect global coarse evidence across tweets while the CRF model conducts sequential labeling to capture fine-grained information encoded in a tweet. The semi-supervised learning plus the gazetteers alleviate the lack of training data[4].The task of social event extraction for tweets, an important source of fresh events. One main challenge is the lack of information in a single tweet, which is rooted in the short and noise-prone nature of tweets.

So propose to collectively extract social events from multiple similar tweets using a novel factor graph, to harvest the redundancy in tweets, i.e., the repeated occurrences of a social event in several tweets[5].This paper describes TWICAL the first open-domain event-extraction and categorization system for Twitter. So demonstrate that accurately extracting an open-domain calendar of significant events from Twitter is indeed feasible. In addition Presenting a novel approach for discovering important event categories and classifying extracted events based on latent variable models. By leveraging large volumes of unlabeled data, approach achieves a 14 percent increase in maximum F1 over a supervised baseline[7].Afterwards, develop a target (i.e. entity) dependent sentiment classification approach to identifying the opinion towards a given target (i.e. entity) of tweets. Finally, the opinion summary is generated through integrating information from dimensions of topic, opinion and insight, as well as other factors (e.g. topic relevancy, redundancy and language styles) in an unified optimization framework. Conduct extensive experiments on a real-life data set to evaluate the performance of individual opinion summarization modules as well as the quality of the produced summary. The promising experiment results show the effectiveness of the proposed framework and algorithms[8]

.Consider the problem of finding opinionated tweets about a given topic. Automatically construct opinionated lexical from sets of tweets matching specific patterns indicative of opinionated messages. When incorporated into a learning to rank approach, results show that automatically opinionated information yields retrieval performance comparable with a manual method[9].A segment-based event detection system for tweets, called Tweet event. Tweet event first detects bursty tweet segments as event segments and then clusters the event segments into events considering both their frequency distribution and content similarity. More specifically, each tweet is split into non-overlapping segments (i.e., phrases possibly refer to named entities or semantically meaningful information units). The bursty segments are identified within a fixed time window based on their frequency patterns ,and each bursty segment is described by the set of tweets containing the segment published within that time window[10].This paper proposes a Hidden Markov Model (HMM) and an HMM-based chunk tagger, from which a named entity (NE) recognition (NER) system is built to recognize and classify names, times and numerical quantities. Through the HMM, system is able to apply and integrate four types of internal and external evidences: 1) simple deterministic internal feature of the words, such as capitalization and digitalization;2) internal semantic feature of important triggers; 3) internal gazetteer feature; 4)external macro context feature[11].Target out-of-vocabulary words in short text messages and propose a method for identifying and normalizing ill-formed words. Method uses a classifier to detect ill-formed words, and generates correction candidates based on morphophonemic similarity. Both word similarity and context are then exploited to select the most probable correction candidate for the word[12].

## 3. OBJECTIVE:

Using the Named Entity Recognition(NER) algorithm developed an online summaries or the historical summaries of tweeter which will also combines data with historical data with new data. Evaluate the queries of user.

## 4. PROBLEM STATEMENT:

The data is then divided into clusters by using a clustering algorithm which will divide the data according to the class it will belong. Clusters will be of crime, politics, and religious and so on. The monitoring/keyword filtering is done on the cluster so that only that data will be left which contains the words related to the riots or civil unrest.

## 5. SYSTEM ARCHITECTURE:

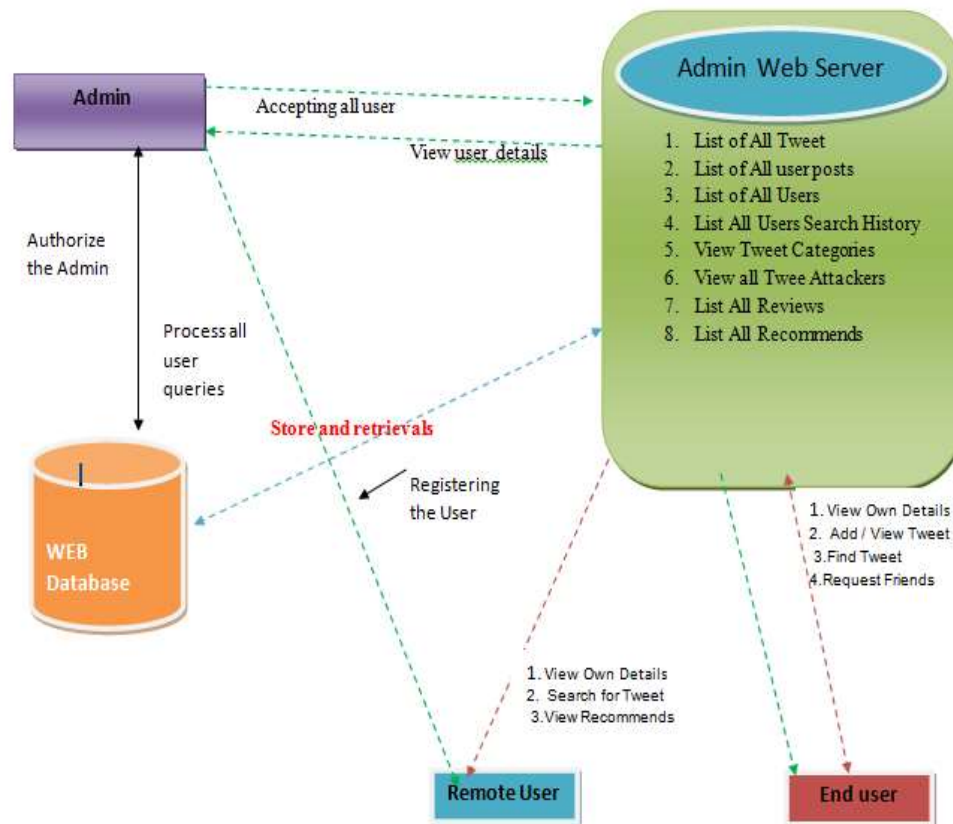Fig. 1. System Architecture

**5.1 Existing System:-**

Many existing NLP[2]techniques heavily rely on linguistic features, such as POS tags of the surrounding words, word capitalization, trigger words (e.g., Mr., Dr.), and gazetteers. These linguistic features, together with effective supervised learning algorithms (e.g., hidden markov model (HMM)[5] and conditional random field (CRF))[5], achieve very good performance on formal text corpus. However, these techniques experience severe performance deterioration on tweets because of the noisy and short nature of the latter.

In Existing System, to improve POS tagging on tweets, Ritter et al. train a POS tagger by using CRF model with conventional and tweet-specific features. Brown clustering is applied in their work to deal with the ill-formed words.[1]

**5.2 Proposed System:-**

To achieve high quality tweet segmentation, propose a generic tweet segmentation framework, named Hybrid Segment. Hybrid Segment learns from both global and local contexts, and has the ability of learning from pseudo feedback.[1]

**5.3 Global context[1]: -**

Tweets are posted for information sharing and communication. The named entities and semantic phrases are well preserved in tweets[3].

**5.4 Local context[1]:-**

Tweets are highly time-sensitive so that many emerging phrases like "She Dancin" cannot be found in external knowledge bases. However, considering a large number of tweets published within a short time period (e.g., a day) containing the phrase, it is not difficult to recognize "She Dancin" as a valid and meaningful segment. Therefore investigate two local contexts, namely local linguistic features and local collocation[1][2].

## 6. ALGORITHM:

Named Entity Recognition Process

The semi-supervised NER algorithm[3]

Step 1: L - a small set of labeled training data
Step 2: U - unlabeled data
Step 3: Loop for k iterations:
Step 4: Train a classifier Ck based on L;
Step 5: Extract new data D based on Ck;
Step 6: Add D to L;

Extract new data D based on Ck

i) Classify kth portion of U and compute confidence scores;
ii) Find high-confidence Named Entity segments and use them to tag other low confidence tokens
iii) Find qualilied O tokens
iv) Extract selected NE and O tokens as well as their neighbors
v) Shuffle part of the NEs in the extracted data
vi) Add extracted data to D

### 6.1 Advantages:-

Our work is also related to entity linking (EL). EL is to identify the mention of a named entity and link it to an entry in a knowledge base like Wikipedia.

Through our framework, demonstrate that local linguistic features are more reliable than term-dependency in guiding the segmentation process. This finding opens opportunities for tools developed for formal text to be applied to tweets which are believed to be much more noisy than formal text[5][7s].

Helps in preserving Semantic meaning of tweets.

### 6.2 Disadvantges:-

Given the limited length of a tweet (i.e., 140 characters) and no restrictions on its writing styles, tweets often contain grammatical errors, misspellings, and informal abbreviations.

The error-prone and short nature of tweets often make the word-level language models for tweets less reliable.

## 7. CONCLUSION

This paper presents an a prototype which supported continuous tweet stream summarization. A tweet stream clustering algorithm to compress tweets into clusters and maintains them in an online fashion. Then, it uses a Rank summarization algorithm for generating online summaries and historical summaries with arbitrary time durations. The topic evolution can be detected automatically, allowing System to produce dynamic timelines for tweet streams by using Local and Global Context[1].

## 8. ACKNOWLEDGEMENT

comments during my research work. I should also like to acknowledge the contribution of my Principal Dr.G.U.Kharat.

## 9. REFERENCES

[1] C. Li, J. Weng, Q. He, Y. Yao, A. Datta, A. Sun, and B.-S. Lee, "Twiner: Named entity recognition in targeted twitter stream," in Proc. 35th Int. ACM SIGIR Conf. Res. Develop. Inf. Retrieval, 2012, pp. 721–730.

[2] C. Li, A. Sun, J. Weng, and Q. He, "Exploiting hybrid contexts for tweet segmentation," in Proc. 36th Int. ACM SIGIR Conf. Res. Develop. Inf. Retrieval, 2013, pp. 523–532.

[3] A. Ritter, S. Clark, Mausam, and O. Etzioni, "Named entity recognition in tweets: An experimental study," in Proc. Conf. Empirical Methods Natural Language Process., 2011, pp. 1524–1534.

[4] X. Liu, S. Zhang, F. Wei, and M. Zhou, "Recognizing named entities in tweets," in Proc. 49th Annu. Meeting Assoc. Compute. Linguistics: Human Language Technol., 2011, pp. 359–367.

[5] X. Liu, X. Zhou, Z. Fu, F. Wei, and M. Zhou, "Exacting social events for tweets using a factor graph," in Proc. AAAI Conf. Artif. Intell., 2012, pp. 1692–1698.

[6] A. Cui, M. Zhang, Y. Liu, S. Ma, and K. Zhang, "Discover breaking events with popular hashtags in twitter," in Proc. 21st ACM Int. Conf. Inf. Knowl. Manage., 2012, pp. 1794–1798.

[7] A. Ritter, Mausam, O. Etzioni, and S. Clark, "Open domain event extraction from twitter," in Proc. 18th ACM SIGKDD Int. Conf. Knowledge Discovery Data Mining, 2012, pp. 1104–1112.

[8] X. Meng, F. Wei, X. Liu, M. Zhou, S. Li, and H. Wang, "Entity centric topic-oriented opinion summarization in twitter," in Proc. 18th ACM SIGKDD Int. Conf. Knowledge Discovery Data Mining, 2012, pp. 379–387.

[9] Z. Luo, M. Osborne, and T. Wang, "Opinion retrieval in twitter," in Proc. Int. AAAI Conf. Weblogs Social Media, 2012, pp. 507–510

[10] X. Wang, F. Wei, X. Liu, M. Zhou, and M. Zhang, "Topic sentiment analysis in twitter: a graph-based hash tag sentiment classification approach," in Proc. 20th ACM Int. Conf. Inf. Knowl. Manage.,2011, pp. 1031–1040.

[11] K.-L. Liu, W.-J. Li, and M. Guo, "Emoticon smoothed language models for twitter sentiment analysis," in Proc. AAAI Conf. Artif ntell., 2012, pp. 1678–1684.

[12] S. Hosseini, S. Unankard, X. Zhou, and S. W. Sadiq, "Location oriented phrase detection in microblogs," in Proc. 19th Int. Conf. Database Syst. Adv. Appl., 2014, pp. 495–509.

[13] C. Li, A. Sun, and A. Datta, "Twevent: segment-based event detection from tweets," in Proc. 21st ACM Int. Conf. Inf. Knowl. Manage.,2012, pp. 155–164.

[14] L. Ratinov and D. Roth, "Design challenges and misconceptions in named entity recognition," in Proc. 13th Conf. Compute .Natural Language Learn., 2009, pp. 147–155.

[15] J. R. Finkel, T. Grenager, and C. Manning, "Incorporating nonlocal information into information extraction systems by Gibbs sampling," in Proc. 43rd. 40th Annu. Meeting Assoc. Comput. Linguistics, 2002, pp. 473–480.

## BIOGRAPHIES

| | |
|---|---|
| | **Miss.Anjum Inamdar** is a student of 7th semester in Department of Computer Science, Sharadchandra Pawar college of Engg, Otur She is working on the project titled The Exploiting Concept of Twitter NER of Segmentation. This paper is the outcome of the application being developed. |
| | **Miss.Kartiki Wahatole** is a student of 7th semester in Department of Computer Science, Sharadchandra Pawar college of Engg,Otur She is working on the project titled The Exploiting Concept of Twitter NER of Segmentation. This paper is the outcome of the application being developed. |
| | **Miss.Vishakha Shinde** is a student of 7th semester in Department of Computer Science,Sharadchandra Pawar college of Engg,Otur She is working on the project titled The Exploiting Concept of Twitter NER of Segmentation. This paper is the outcome of the application being developed. |
| | **Miss.Harshata Tohake** is a student of 7th semester in Department of Computer Science,Sharadchandra Pawar college of Engg,Otur She is working on the project titled The Exploiting Concept of Twitter NER of Segmentation. This paper is the outcome of the application being developed. |