

# Evaluation of Tweet Segmentation using Different Techniques

Patil Umesh A<sup>1</sup>, Prof. Manwade K.B<sup>2</sup>

ME (CSE) Student, Ashokrao Mane Group of Institutions, Vathar Kolhapur, Maharashtra, India<sup>1</sup>

Associate Professor Ashokrao Mane Groups of Institutions, Vathar Kolhapur, Maharashtra, India<sup>2</sup>

**Abstract:** Many private and/or public organizations have been reported to Create and monitor targeted Twitter streams to collect and understand users' opinions about the organizations. Targeted Twitter stream is usually constructed by filtering tweets with user defined selection criteria (e.g., tweets published by users from a selected region, or tweets that match one or more predefined keywords). Targeted Twitter Stream is then monitored to collect and understand users' opinions about the organizations. There is an emerging need for early crisis detection and response with such target stream. Such applications require good named entity recognition (NER) system for Twitter, which is able to automatically discover emerging named entities that are potentially linked to the crisis. In this paper, we present an over-step unsupervised NER system for targeted Twitter stream, called Novel-NER. In the first 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, Novel-NER constructs a random walk model to exploit the gregarious property in the local context derived from the Twitter stream. The highly-ranked segments have a higher chance of being true named entities. We evaluated Novel-NER on two sets of real life tweets simulating two targeted streams. Evaluated using labeled ground truth, Novel-NER achieves comparable performance as with conventional approaches in both streams. Various settings of Novel-NER have also been examined to verify our global context +local context combo idea. As well as we are using the Wikipedia & Microsoft N gram with Association & Correlation.

**Keywords:** Novel-NER constructs a random walk model, global context +local context combo idea

## I. INTRODUCTION

Twitter, as a new type of social media, has seen tremendous growth in recent years. It has attracted great interests from both Industry and academia. Many private and or public organizations have been reported to monitor Twitter stream to collect and understand users' opinions about the organizations. Nevertheless, due to the extremely large volume of tweets published every day, it is practically infeasible and unnecessary to listen and monitor the whole Twitter stream. Therefore, targeted Twitter streams are usually monitored instead; each such stream contains tweets that potentially satisfy some information needs of them on their organization. Targeted Twitter stream is usually constructed by filtering tweets with user defined selection criteria depends on the information needs. For example, the criteria could be a region so that users' opinions from that particular region are collected and monitored; it could also be one or more predefined key words so that Opinions about some particular events/topics/products/services can be monitored. There is also an emerging need for early crisis detection and response with such target stream. For example, a cosmetic company is interested in automatically discovering any new named entities (e.g. person names, competitor names, or location names) in a targeted stream it creates for the company and its products, which may link to a potential PR crisis.

By doing this, the company is able to acquire first hand information about the crisis and make early response. Such applications require good named entity recognition (NER) system for Twitter, which is the focus of this paper. Nevertheless, the nature of tweets brings new challenges. Traditional NER methods on well-formatted documents heavily depend on a phrase's local linguistic features such as capitalization, part of speech (POS) tags of previous words, etc. However, tweets are usually informal in nature and short (up to 140 characters).

They often contain grammatical errors, misspellings, and unreliable capitalizations. These unreliable linguistic features cause traditional methods to perform poorly on tweets. We use the real examples below to illustrate the challenges when applying traditional NER on tweets.

To address the above challenges caused by tweets' error-prone and short nature, this paper presents an overall unsupervised NER system for targeted tweet streams, called Novel-NER. Based on the gregarious property of named entities in targeted tweet stream, Novel-NER recognizes named entities collectively from a batch of tweets in unsupervised manner. More formally, let T be the collection of tweets in question. Novel-NER receives tweets from T in a batch manner. A batch is the set of tweets posted in the targeted Twitter stream. The idea is to

segment an individual tweet into as sequence of consecutive phrases, each of which appears “more than chance” gives an example. More formally, given a tweet off our words  $w_1, w_2, w_3, w_4$ , we segment it as  $w_1 w_2 | w_3 w_4$  rather than  $w_1 | w_2 w_3 w_4$ , if  $C(w_1 w_2) + C(w_3 w_4) > C(w_1) + C(w_2 w_3 w_4)$ , where  $C(\cdot)$  basically captures the probability Being a valid phrase of a segment. A straight forward idea of computing  $Pr(\cdot)$  is to count a segment’s appearance in a very large corpus. The ideal case is that we use the entire collection of tweets published in Twitter to compute the  $Pr(\cdot)$  for all possible segments. Unfortunately, to the best of our Knowledge, such corpus never exists. Instead, we turn to Microsoft

Within one fixed time interval (e.g. second). So,  
 $T = \{T_1, T_2, \dots, T_n\}$

Web N-Gram corpus This N-Gram corpus is based on all the and  $T_i$  is the batch of tweets posted in the  $i$ th interval. Novel-NER the  $n$  recognizes all possible named entities in  $T_i$  regardless of their types. It is noted that currently Novel-NER does not categorize the type of named entity (e.g., person, location). As conventional NER methods fail to address the new challenges posed by emerging media like Twitter, it is more pressing to be able to discover the presence of named entities in targeted Twitter stream before we could categorize their types. Furthermore, even without categorizing the types of named entities, Novel-NER already enable us to make early crisis response. For example, a cosmetic company may Be interested in discovering any new named entity which may directly/indirectly link to the company and subsequently causes a crisis, be it a person name, product name, or company name. Moreover, as a targeted Twitter stream is constructed for a particular information need, we assume that the user who constructs the stream has the background knowledge in interpreting the named entities detected. In the following subsections, overview of Novel-NER.

## II. IMPLEMENTATION STEPS

The performance of proposed system will be evaluated using programming language JAVA software tools and the following flow chart. Tweet segmentation is used to extract the named entity candidates from tweets, or in other words, to identify the correct boundary of potential named entities in tweets.

- i) **Input Tweets:**-Taking targeted tweets stream & applying for further preprocessor.
- ii) **Preprocessor:** - It takes those tweets which are useful for further discussion, for this it uses framework called as HybridSeg with downstream application.
- iii) **Segmentation:** - This process is doing by using the global context with Microsoft N-gram & Wikipedia.
- iv) **Selection of Segments:** - By using Novel-Named Entity Recognition algorithm (Novel-NER) & random walk algorithm it selects the segments.
- v) **Establish Correlation:** - By using genetic algorithm selected segments showing more accuracy on real & large dataset.

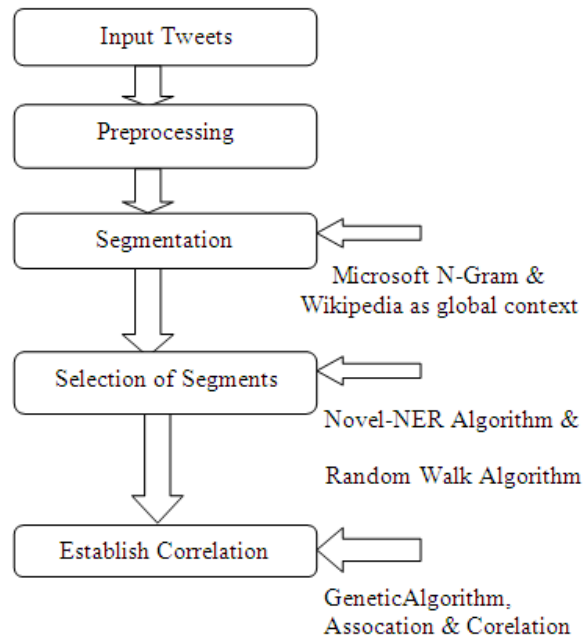


Fig.01 Implementation Steps.

In Novel-NER, information in tweets’ local context and global context are aggregated to calculate the probability that a phrase is a named entity. By doing so, Novel-NER is able to recognize new named entities which may not appear in Wikipedia.

To the best of our knowledge, it is the first to exploit both the local context (in tweets) and the global context (from World Wide Web) together for NER task in twitter.

## III. TWEET SEGMENTATION

In this section, we detail our solution for given an individual tweet  $t \in T_i$ , the problem of tweet segmentation is to split  $t$  into  $m$  consecutive segments,  $t = s_1 s_2 \dots s_m$ ; each segment Contains one or more words.

To obtain the optimal segmentation, we use the following objective function, where  $C$  is the function that measures the stickiness of segment or a tweet defined based on word collocation:

$$\arg \max_{s_1, s_2, \dots, s_m} C(t) = \sum_{i=1}^m C(s_i)$$

A high stickiness score of segments indicates that it is not suitable to further split segments, 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  $l$ , there exist  $2^{l-1}$  possible segmentations. It is inefficient to iterate all of them and compute their Stickiness. We there for design a dynamic programming algorithm to tackle the problem, which is presented in the following

#### IV DYNAMIC PROGRAMMING

Outlines our dynamic programming algorithm for tweet segmentation. The basic idea is to recursively conduct binary segmentations and then evaluates the stickiness of the resultant segments. More formally, given any segments from  $t$  (scan is  $t$  itself or a part of  $t$ ) and supposes  $w_1 w_2 \dots w_n$ , our solution is to conduct a binary segmentation by splitting it in to two adjacent segments  $s_1 = w_1 \dots w_j$  and  $s_2 = w_{j+1} \dots w_n$  by satisfying:

$$\text{Arg max } C(s) = C(s^1) + C(s^2).$$

It gives the lead to a complexity of  $O(l)$ .

#### V SEGMENT STICKINESS FUNCTION

A high stickiness score of segments indicates that further splitting segments would break the correct word collocation. There are a number of collocation measurements. However, all these measures we predefined for two arguments. That is, they were designed to measure the collocation of the  $n$ -grams with the particular binary partition. A variety of studies have been conducted to extend these binary collocation measures to the  $n$ -grams case We defined the stickiness functions by using the generalization frame work proposed in Specifically, the generalized collocation measures of

Point Mutual Information (PMI)

Symmetric Conditional Probability (SCP)

$$\text{PMI Based Stickiness } (w_1 w_2) = \log \frac{\text{Pr}(w_1 w_2)}{\text{Pr}(w_1) \text{Pr}(w_2)}$$

$$\text{SCP Based Stickiness } (w_1 w_2) = \log \frac{\text{Pr}(w_1 w_2)^2}{\text{Pr}(w_1) * \text{Pr}(w_2)}$$

The stickiness of the given segments may be enhanced by using the World Wide Web, the stickiness function is enhanced by

$$C'(s) = C(s) * e^{Q(s)}$$

$C(s)$  is calculated by using the global context from Microsoft Web N- Gram

#### VI. LENGTH NORMALIZATION

Generally longer valid segments have higher chance of being named entity than shorter ones. If we treat all tweet segments equally then it will not give the better result hence we are using the length normalization in which we can

#### VII. EVALUATION

##### 1) TWEETS DATA & PERFORMANCE METRICS

Collections of tweets are used in the experiments to simulate targeted twitter streams. We are collected one type of targeted tweet stream for some specific period. The data collection from period may be any recent

events like Olympics games, P.V. Sindhu, Prime Minister etc on which you want to retrieve the data.

##### 2) COMPARISON

In this section, we compare Novel-NER with two conventional NER systems trained on tweets. Specifically, we train Wikipedia and Microsoft Web N Gram with the labeled tweet data and evaluate their performance. Moreover; we also compare with a tweet septic NER system proposed on the tweet collection.

**Wikipedia:** This system based on the regularized averaged perception approach which uses gazetteers extracted from Wikipedia, word class models derived from unlabeled text, and expressive nonlocal features.

**Microsoft Web N Gram:** This system based on Microsoft Corpus incorporates long-distance information & achieves good Performance consistently across Local & Global Context.

**Novel-NER:** a supervised NER system uses Wikipedia and Microsoft Web N Gram along with correlation & Association, Random model for learning and inference. A set of widely- used effective features are used in Novel-NER, including orthographic, contextual, and dictionary features.

Sr. No.	Method	% Result
1	Wikipedia	73.8
2	Random Walk +Wikipedia	75.6
3	RWW+ Wiki+ Corelation+ Association	82.5

Table 01: Ranking of Segments

##### 3) PERFORMANCE OF TWEET SEGMENTATION

Tweet segmentation is used to extract the named entity candidates from tweets, or in other words, to identify the correct boundary of potential named entities in tweets. It is a critical component because the performance of Novel-NER is heavily affected by the effectiveness of tweet segmentation. Two stickiness functions are defined by using two collocation measures, PMI and SCP or tweet segmentation.

The tweet Segmentation algorithm is also incorporates an external knowledge base Wikipedia. Further, we normalize the segment length to favor long named entities. In this section, we study the impact of the collocation measures (PMI or SCP), the Wikipedia dictionary (Wiki), and the length normalization (Norm) we use tweet segmentation with only PMI or SCP measures as the baseline SCP significantly out performs PMI for tweet segmentation. We believe this is because PMI returns disproportionately high values for frequent items. This property makes PMI prefer longer segments, which is confirmed by manual investigation of the segmentation result.

Length normalization (Norm) is effective and improves the accuracy of tweet segmentation in the two collections for SCP-based stickiness. For example, given a long segment Since PMI prefers longer segments; further preference introduced by Norm does aggravate the problem.

Wikipedia's broad coverage and high quality knowledge help reclaim incorrect decisions made by the lexical statistics or reinforce the correct decisions. The combination of Wikipedia dictionary and length normalization further boosts up the performance of tweet segmentation of SCP-based stickiness. This indicates that Wiki and Norm are complementary to each other.

#### 4) IMPACT OF RANDOM WALK ON SEGMENT RANKING

A random walk model is applied to exploit the gregarious property of named entities in tweets. The final segment ranking output is an aggregation from the stationary probability of the random walk model (local context) and the segment's Wikipedia-based teleportation Wiki probability (global context). We analyze their impact on the performance of segment ranking in this section. Specifically, we investigate the following schemes for segment ranking:

- 1) MFS: A naïve method that ranks the segments based on their frequency in the collection. That is, the most frequent segments are ranked higher.
- 2) Wiki: A naïve method that ranks the segments based on their
- 3) Wikipedia- Based teleportation probability.
- 4) RW: A simple random walk with uniform teleportation. The segments are then ranked based on the stationary probability  $\pi(s)$
- 5) RWW: A random walk with Wikipedia-based teleportation. The segments are then ranked based on the stationary probability  $\pi(s)$
- 6) RWW +Wiki +Correlation + Association: A random walk with Wikipedia based teleportation, while the segments are ranked based on Equation:  $y(s) = e_s \cdot \pi(s)$  & it uses the correlation & association on two different named entities.

### VIII. CONCLUSION

The error prone and short nature of Twitter has brought new challenges to named entity recognition. In this paper, we present a NER system for targeted Twitter stream, called Novel-NER, to address this challenge. Unlike traditional methods, Novel-NER is unsupervised. It does not depend on the unreliable local linguistics features. Instead, it aggregates information garnered from the Worldwide Web to build robust local context and global con-text for tweets. Experimental results show promising results of Novel-NER. Despite its promising results, there is still space for improvement. We plan to study Novel-NER performance in a larger scale. We plan to study the strategy to identify suitable K value. As discussed earlier,

this is because we feel this problem is not as pressing as the problems to correctly locate and recognize presence of named entities in tweets, which exist in methods largely fail. Extension of Novel-NER for entity type classification is also planned for future work

### REFERENCES

- [1] 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.
- [2] C. Li, A. Sun, and A. Datta, "Twevent segment-based event detection from tweets," in CIKM, 2012, pp. 155–164.
- [3] Z. Luo, M. Osborne, and T. Wang, "Opinion retrieval in twitter," in ICWSM, 2012.
- [4] X. Meng, F. Wei, X. Liu, M. Zhou, S. Li, and H. Wang, "Entity centric topic-oriented opinion summarization in twitter," in KDD, 2012, pp. 379–387.
- [5] 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 CIKM, 2011, pp. 1031–1040.
- [6] A. Ritter, Mausam, O. Etzioni, and S. Clark, "Open domain event extraction from twitter," in KDD, 2012, pp. 1104–1112.
- [7] X. Wang, A. McCallum, and X. Wei, "Topical n-grams: Phrase and topic discovery, with an application to information retrieval," in ICDM, 2007, pp. 697–702.
- [8] J. R. Finkel, T. Grenager, and C. Manning, "Incorporating nonlocal information into information extraction systems by gibbs sampling," in ACL, 2005, pp. 363–370.
- [9] K. Nishida, T. Hoshida, and K. Fujimura, "Improving tweet stream classification by detecting changes in word probability," in