



International Journal of Innovative Research in Computer and Communication Engineering

(An ISO 3297: 2007 Certified Organization)

Vol. 4, Issue 3, March 2016

A Survey on Named Entity Recognition with the Use of Tweet Segmentation

Vivek Dhole, Dhanshri Patil

M.E. Student, Dept. of Computer Network, N.M.I.E.T. Talegaon Dabhade, Savitribai Phule Pune University, Pune,
India

Assistant Professor, Dept. of Computer Network, N.M.I.E.T. Talegaon Dabhade, Savitribai Phule Pune University,
Pune, India

ABSTRACT: Twitter is attracting huge number of clients to share most up-to-date information, bringing about extensive volumes of information delivered each day. However, numerous applications are experiencing severely from noisy and short nature of tweets. So Proposed a novel system for tweet segmentation in a batch mode called Hybrid Segment. The semantic and context information is well preserved by dividing tweets into significant fragments, and is effectively extracted by the downstream applications. The optimal segmentation of a tweet is found by maximizing the sum of stickiness scores of its candidate segments. Experiments on tweet data sets show that quality of tweet segmentation is improved by learning the global, local contexts compared with the use of global context alone. Through investigation, recognized that local linguistic features are reliable for learning local context compared with term-dependency. As an application attempted to demonstrate that high accuracy is achieved in NER.

KEYWORDS: Tweet Segmentation, Named Entity Recognition(NER).

I. INTRODUCTION

MICROBLOGGING sites like as the Twitter have reshaped the way individuals discover, share, and disseminate timely information. Numerous associations have reported to create and monitor targeted Twitter streams to collect and understand user's opinions. Targeted Twitter stream is developed by filtering the tweets with predefined choice criteria (e.g., tweets that are published by users from a geographical region, tweets that match one or more predefined keywords). Due to the invaluable business value of timely information from these tweets it is imperative i.e important to understand tweets' language for the large body of downstream applications like as named entity recognition (NER) event detection and summarization [1] opinion mining sentiment analysis and many others.

Given the length of a tweet is constrained (i.e., 140 characters) and no limitations on its written work styles, tweets frequently contain incorrect spellings, grammatical errors, and informal abbreviations. The error-prone and the short nature of tweets often make the word-level language models for tweets are less dependable. For instance, given a tweet 'I call her, there is no answer, her phone in the bag, she dancin,' there is no clue or no idea to guess its true theme by disregarding word order (i.e., bag-of-word model). The circumstance is further described and exacerbated with the limited context provided by the tweet. That is, more than one clarification for this tweet is derived by different readers if the tweet is considered in isolation. On the other hand, due to the noisy nature of tweets, the core semantic information is well preserved in tweets in the form of semantic phrases or named entities. Case in point, the creating expression "she dancin" in the related tweets shows that its a key thought—it portrays the tweet into gathering of tweets which discusses the tune, She Dancin an example subject in the Bay Area.

Overview- First, to obtain tweets on the target event precisely, apply semantic analysis of a tweet. For example, users make tweets which can be such as 'Now it is shaking' or 'Earthquake' for which shaking or earthquake are considered as keywords, but users may also make tweets such as 'Someone is shaking hands with their friends' Or 'I am attending an Earthquake Conference.' After this then try and prepare the training data and then devise a classifier using a Support Vector Machine (SVM) based on features such as the number of words, keywords in a tweet, and the context of target-

International Journal of Innovative Research in Computer and Communication Engineering

(An ISO 3297: 2007 Certified Organization)

Vol. 4, Issue 3, March 2016

event words. After doing this then achieve a probabilistic spatiotemporal model of an event, then make a crucial assumption that each Twitter user is regarded and considered as a sensor and each tweet as sensory information.

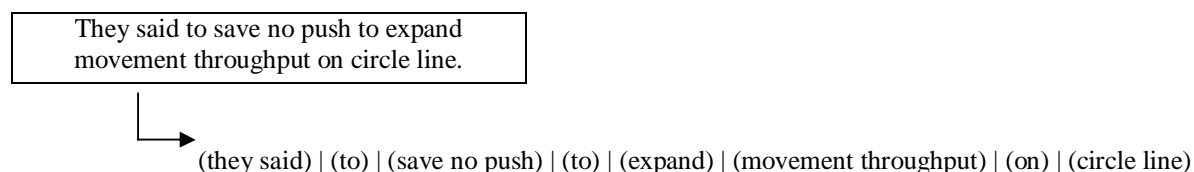


Fig. 1. Sample of tweet segmentation.

Project Idea- The main focus is given on the undertaking of tweet division. The objective of this assignment is to part a tweet into an arrangement of back to back n-grams ($n > 1$), each of which is known as a segment. A segment can be a named entity (e.g., a film title "discovering nemo"), a semantically significant data unit, or whatever other sorts of expressions which appear more than by possibility or which show up more than by chance. In the above example the tweet is split into eight segments and the semantically significant segments in the example are protected. Since these segments protect semantic importance of the tweet more unequivocally than each of its constituent words does, the theme of this tweet can be better caught in the resulting preparing of this tweet.

Then identified two directions for the research. One is to further improve the segmentation quality by considering more local factors and other is to investigate the adequacy of the segmentation-based representation for undertakings such as tweets summarization, search, hashtag recommendation, and so forth.

II. RELATED WORK

In the work [2] Author said, Social events are events that happen between individuals where no less than one individual knows about the other and of the event occurring. Extracting social events can play an important role in a wide range of applications, for example, the development of interpersonal organization. Here author has presented the assignment of social event extraction for tweets, an essential and important source of fresh events. One fundamental test is the absence of data in a single tweet, which is established in the short and noise-prone nature of tweets. Author proposes to all in all concentrate get-togethers from numerous comparable tweets utilizing a novel element diagram, to collect the redundancy in tweets, that is the repeated occurrences of a social event in several tweets. Then evaluate the method on a human annotated data set, and demonstrate that it outperforms all baselines, with an absolute gain.

In the work [3] Author said, the difficulties of Named Entities Recognition (NER) for tweets lie in the lacking data in a tweet and the inaccessibility of preparing information. Author proposed to consolidate a K-Nearest Neighbors classifier with a linear Conditional Random Fields model under a semi-supervised learning framework to tackle these challenges. The K-Nearest Neighbors based classifier conducts pre-labeling to collect global coarse evidence across tweets while the Conditional Random Fields model conducts sequential labeling to capture fine-grained information encoded in a tweet. The semi-administered learning plus the gazetteers ease the absence of preparing information. Broad analyses demonstrate the benefits of the technique over the baselines and in addition the effectiveness of KNN and semi supervised learning.

In the work [4] Author said, Event detection from tweets is an imperative task to understand the current events/topics attracting an expansive number of common users. However, the unique characteristics of tweets (e.g., short and noisy content, diverse and fast changing topics, and substantial information volume) make event detection a challenging assignment. Most existing strategies proposed for well written documents (e.g., news articles) cannot be directly adopted. Author proposed a segment-based event detection framework for tweets, called Twevent. Twevent first recognizes bursty tweet segments as event segments and then bunches the event segments into events considering both their frequency distribution and content similarity. More particularly, each tweet is split into non-overlapping segments. The bursty segments are distinguished within a settled time window based on their frequency patterns, and each bursty segment is described by the set of tweets containing the section distributed inside of that time window. The likeness between a pair of bursty segments is computed using their associated tweets. Author likewise demonstrates that Twevent is proficient and versatile, prompting an alluring answer for event detection from tweets.



International Journal of Innovative Research in Computer and Communication Engineering

(An ISO 3297: 2007 Certified Organization)

Vol. 4, Issue 3, March 2016

III. PROPOSED WORK

Here the focus is given on the task of tweet segmentation. The objective of this assignment is to split and divide a tweet into a sequence of consecutive n-grams, each of which is called a segment. A segment can be a named entity (e.g., a film title "discovering nemo"), a semantically significant data unit, or whatever other sorts of expressions which appear more than by possibility or which show up more than by chance.

Propose a generic tweet segmentation framework to achieve high quality tweet segmentation, named Hybrid Segment. Hybrid Segment learns from the global and local contexts, and has the capacity of gaining from pseudo feedback.

Global context : Tweets are posted which are used for data sharing and correspondence. The named entities and semantic phrases are well preserved in tweets.

Local context. Tweets are highly time-touchy i.e. time-sensitive because of this many emerging phrases such as 'She Dancin' cannot be found in external knowledge bases. Therefore, considering a large number of tweets published within a short time period containing the phrase (e. g., a day), it is not difficult to recognize "She Dancin" as a valid and meaningful segment. Therefore investigate or examine two local contexts, which are the local linguistic features and local collocation.

Pseudo feedback: The segments which are recognized based on the local context with high confidence serve as good feedback to extract or get more meaningful segments. The learning from this pseudo feedback is conducted iteratively and the method which implements the iterative learning is named Hybrid Segment_{iter}.

ADVANTAGES:

-The work is also related to entity linking, it is to identify the mention of a named entity and link it to an entry in a knowledge base like Wikipedia.

-The framework, show that local linguistic features are more reliable than term-dependency in directing the segmentation process. This discovering opens opportunities for tools, devices produced developed for formal text to be applied to tweets which are believed to be much more noisy than formal text.

-Helps in safeguarding and preserving Semantic meaning of tweets.

IV. SYSTEM ARCHITECTURE

The general overview of the system architecture is shown in the diagram. The system mainly involves six phases. The general information on the phases is as follows:

Data Gathering also known as Information Gathering is the process of collecting a Twitter user's data, including user's friends' posts as well as user's own posts. In this phase, user-friend relationship is also extracted and friends' relative ranking is produced as an output.

Knowledge Base Construction is the process of creating a chart based Knowledge base of Turkish Wikipedia article titles and their links to one another, in order to validate named entity candidates generated as an output of Named Entity Recognition phase. Keeping this knowledge base up to date is also included in this phase. Although other phases iteratively follow each other and one's output is the other's input, this phase is independent and conducted in parallel.

Data Preprocessing includes removing unnecessary parts of tweet texts for example, mentions, hashtags, smileys, vocatives, links and so forth. Since informal writing style is commonly adopted in tweets, this phase is also responsible from normalizing the tweet text such as getting rid of unnecessarily repeated characters, slang words, correcting asciification related issues.

Named Entity Recognition is the next phase of data preprocessing phase. In this phase, tweet segmentation on preprocessed tweets is carried out by means of global context and segments as candidate named entities are produced,

International Journal of Innovative Research in Computer and Communication Engineering

(An ISO 3297: 2007 Certified Organization)

Vol. 4, Issue 3, March 2016

generated. Then, these candidates are validated as named entities or ignored by usage of previously constructed knowledge base of Turkish Wikipedia article titles.

User Interest Model Generation phase is a must requirement. In this phase, using named entities extracted from user's and user's friends' tweets and user-friend relationships, a user interest model is produced. In other words, a Twitter user is represented via weighted named entities.

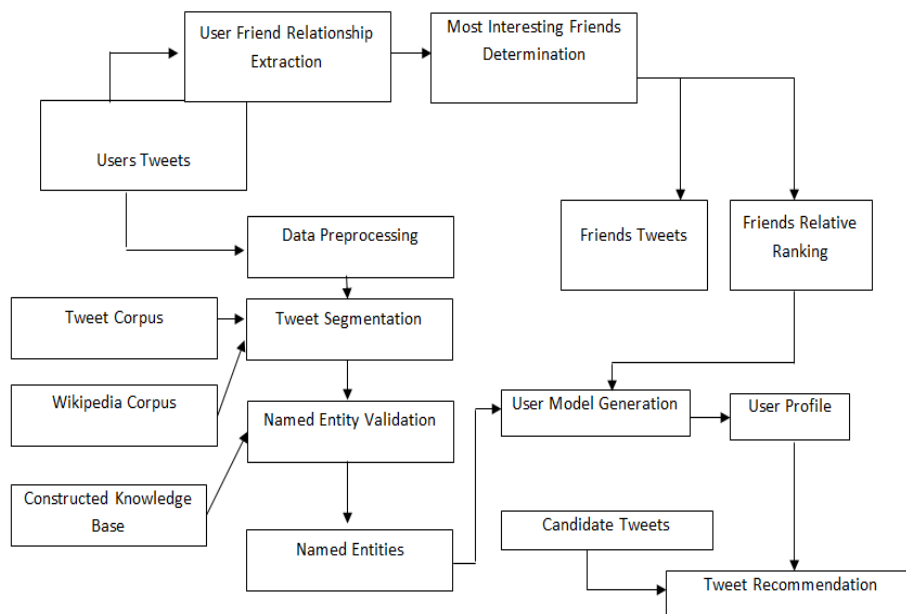


Fig 2 Architecture.

Tweet Recommendation is the last stage, where two sorts of proposal applications connected by contrasting hopeful tweets and the produced client interest model. Tweet classification which is the task of deciding or choosing whether an applicant tweet is interesting for the user or not, and tweet ranking which aims and plan to sort tweets from the most recommendable to the minimum recommendable are performed in this phase.

V. IMPLEMENTATION STRATEGIES

The simulation studies involve the deterministic small network topology with 5 nodes as shown in Fig.1. The proposed energy efficient algorithm is implemented with MATLAB. We transmitted same size of data packets through source node 1 to destination node 5. Proposed algorithm is compared between two metrics Total Transmission Energy and Maximum Number of Hops on the basis of total number of packets transmitted, network lifetime and energy consumed by each node. We considered the simulation time as a network lifetime and network lifetime is a time when no route is available to transmit the packet. Simulation time is calculated through the CPU TIME function of MATLAB. Our results shows that the metric total transmission energy performs better than the maximum number of hops in terms of network lifetime, energy consumption and total number of packets transmitted through the network.

Hybrid Segment Framework

The proposed Hybrid Segment framework segments tweets in batch mode. Tweets from a focused Twitter stream are assembled into clusters by their publication time using a fixed time interval. Every cluster of tweets are then segmented by Hybrid Segment on the whole.

International Journal of Innovative Research in Computer and Communication Engineering

(An ISO 3297: 2007 Certified Organization)

Vol. 4, Issue 3, March 2016

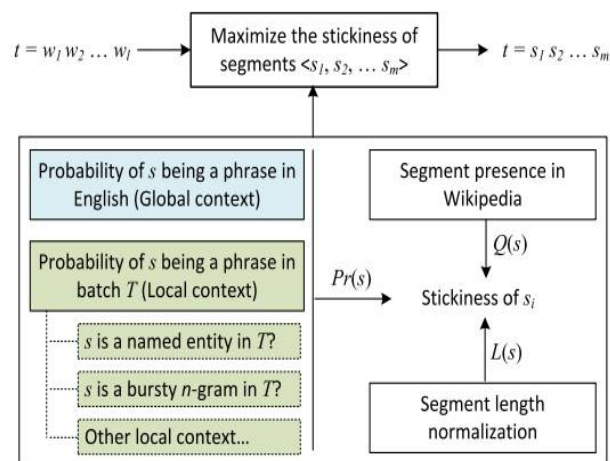


Fig 3 Hybrid Segment framework.

Tweet Segmentation

Tweets are viewed as noisy with lots of informal abbreviations and grammatical errors. However, tweets are posted mostly for data sharing and correspondence among numerous reasons.

Given an individual tweet t belongs to T_i , the issue in tweet segmentation is to split t into n consecutive segments, $t = s_1, s_2, \dots, s_n$; every segment contains one or more than one words. A high stickiness score of segment s demonstrates that it is not suitable to further split segment s , as it breaks the right and correct word collocation. As such, a high stickiness value indicates that a segment can't be further split at any inward position. If the word length of tweet t is L , there exist $2L-1$ conceivable segmentations.

Perceptions for Tweet Segmentation:

Perception 1. Word collocations of named entities and common phrases in English are well preserved in Tweets. Numerous named entities and common phrases are protected in tweets for data sharing and spread.

Perception 2. Numerous tweets contain useful linguistic features. Although numerous tweets contain unreliable linguistic features like misspellings and unreliable capitalizations, there exist tweets formed in proper English.

Perception 3. Tweets in a focused stream are not topically independent to each other within a time window. Numerous tweets distributed inside a short time period talk about the same theme. These similar tweets largely share the same segments.

The latter two perceptions essentially reveal the same phenomenon: local context in a batch of tweets complements global context in segmenting tweets.

Segment-Based Named Entity Recognition

Here select named entity recognition as a downstream application to demonstrate the advantage of tweet segmentation. Investigate two segment-based NER algorithms. The first one identifies named entities from a pool of segments (extracted by Hybrid Segment) by exploiting the co-occurrences of named entities. The second one does as such taking into account the POS tags of the constituent expressions of the segments.

NER by Random Walk:

The principal NER algorithm is based on the perception that a named entity often co-occurs with other named entities in a batch of tweets. Depending on this perception, build a segment graph. A node in this graph is a segment recognized by Hybrid Segment. An edge exists between two hubs if they co-occur in some tweets; and the heaviness of the edge is



International Journal of Innovative Research in Computer and Communication Engineering

(An ISO 3297: 2007 Certified Organization)

Vol. 4, Issue 3, March 2016

measured by Jaccard Coefficient between the two corresponding segments. A random walk model is then applied to the segment graph.

NER by POS Tagger:

Because to the short nature of tweets, the gregarious property may be weak. The second calculation then investigates the part-of speech tags in tweets for NER by considering noun phrases as named elements utilizing segment rather than word as a unit. A segment might show up in various tweets and its constituent words might be appointed diverse POS tags in these tweets. At that point then evaluate the probability of a segment which is a noun phrase by considering the POS tags of its constituent words of all appearances.

VI. CONCLUSION AND FUTURE WORK

Twitter, which is a new type of social media, is picking up significance and has attracted great interests and from both industry and scholarly world. Numerous private and open associations have been reported to monitor Twitter stream to gather and comprehend i.e collect and understand users' opinions about the associations. However it is practically infeasible and unnecessary to monitor and listen the entire Twitter stream, due to its large volume. Therefore, targeted Twitter streams are typically checked. The Hybrid Segment framework which segments tweets into important phrases called segments using both global and local context. Through the framework, it is shown that local linguistic features are much more reliable than the term-dependency in directing the segmentation process. This finding then opens opportunities for devices created for formal content to be applied to tweets which are believed to be much noisier than formal text. Tweet segmentation safeguards and preserves the semantic meaning of tweets, which benefits many downstream applications like e.g., named entity recognition. And a segment-based named entity recognition method achieves much better accuracy than the word-based alternative.

REFERENCES

- [1] Chenliang Li, Aixin Sun, JianshuWeng, and Qi He, "Tweet Segmentation and Its Application to Named Entity Recognition" IEEE TRANSACTIONS ON KNOWLEDGE AND DATA ENGINEERING, VOL. 27, NO. 2, FEBRUARY 2015
- [2] X. Liu, X. Zhou, Z. Fu, F. Wei, and M. Zhou, "Extracting social events for tweets using a factor graph," in Proc. AAAI Conf. Artif. Intell., 2012.
- [3] X. Liu, S. Zhang, F. Wei, and M. Zhou, "Recognizing named entities in tweets," in Proc. 49th Annu. Meeting Assoc. Comput. Linguistics: Human Language Technol., pp. 359-367, 2011
- [4] C. Li, A. Sun, and A. Datta, "Twevent: segment-based event detection from tweets," in Proc. 21st ACM Int. Conf. Inf. Knowl. Manage., pp. 155-164, 2012
- [5] 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
- [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
- [7] H. Ney, U. Essen, and R. Kneser, "On structuring probabilistic dependences in stochastic language modelling," Comput. Speech Language, vol. 8

BIOGRAPHY

Mr. VivekKrishnadeo Dhole (P.G. Student)Comp. Networks, NMVPMs, Nutan Maharashtra Institute of Engineering and Technology, TalegaonDabhade, Pune, India.

Prof. DhanshriPatil(Assistant professor), Comp.Network, NMVPMs, Nutan Maharashtra Institute of Engineering and Technology, TalegaonDabhade, Pune, India.