

Natural Language Processing of the Hybrid Context for Tweet Segmentation

Miss. Varpe Kanchan Nanasaheb¹ Prof. Sandip A. Kahate²

²Assistant Professor

^{1,2}SPCOE, Otur, Pune, Maharashtra, India

Abstract— A time-line based framework for topic summarization is in the Twitter. Summarization on the topics by sub-topics along time line to fully capture rapid topic evolution can be done on Twitter. An HybridSeg can be in corporate to local context knowledge with global knowledge bases for better tweet segmentation on the social network communication. HybridSeg can having of two phases: first step, the existing NER algorithms are applied to on the batches of tweets. The NERs are then employed to guide the tweet segmentation process. Second step, Hybrid-Seg can be manage the tweet segmentation results iteratively by exploiting all kind of segments in the batch of tweets. Experiments on two tweet datasets show that HybridSeg significantly improves tweet segmentation quality compared with the state of the art algorithm.

Key words: Twitter, Segmentation, Hybrid Seg, NER

I. INTRODUCTION

Recently, Twitter has become one of the most popular social networking sites. The people can freely post short messages (called tweets) up to 140 characters. Twitter has rapidly gained worldwide popularity, with over 140 million active users generating over 340 million tweets daily in March 2012. The rapid proliferation of Twitter posts presents a big information. A user to get an overview of important topics on Twitter by reading all tweet record everyday by day. The information redundancy and the informal writing style, it is time consuming to find useful information about a topic from a large number of tweets. Requirements of topic explanation, places, and search from hundreds of thousands of tweets on the media. Specifically, a summary that provides representative information of topics with no redundancy[2].

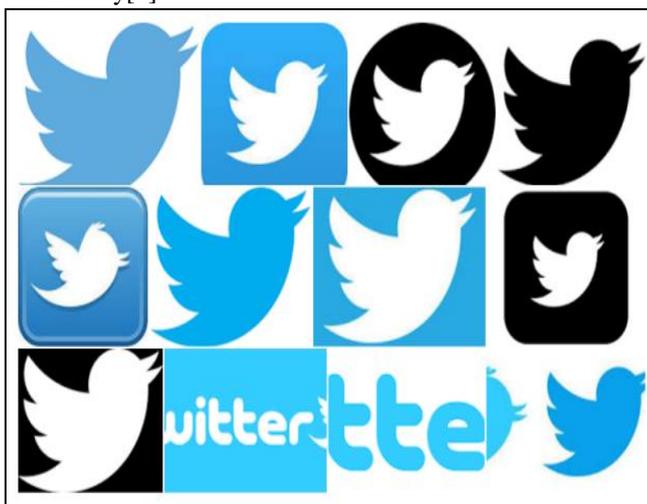


Fig. 1: TWEETER Logo

Twitter is an online social networking service that enables users to send and read short 140-character messages called "tweets". Registered users can read and post tweets, but those who are unregistered can only read the tweets.

Users access Twitter through the website interface, SMS or mobile device app. Twitter is based in San Francisco and has more than 25 offices around all over the world.

Twitter was created in March 2006 by Jack Dorsey, Evan Williams, Biz Stone, and Noah Glass and launched in July 2006. The service rapidly gained worldwide popularity, with more than 100 million users posting 340 million tweets a day in 2012. The service also handled 1.6 billion search queries per day.[14][15][16] In 2013, it was one of the ten most-visited websites and has been described as "the SMS of the Internet" As of March 2016, Twitter has more than 310 million monthly active users[9].

Tweets are created on through the social networks. Twitter has attracted millions of users to share most up-to-date information, resulting in large volumes of data produced everyday. However, many applications in Information Retrieval (IR) and Natural Language Processing (NLP) suffer severely from the noisy and short nature of tweets. An propose a novel framework for tweet segmentation in a batch model, called HybridSeg. An splitting tweets into meaningful segments, the semantic or context information is well preserved and easily extracted. HybridSeg 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 (i.e., Global context) and the probability of a segment being a phrase within the batch of tweets (i.e., Local context).

II. RELATED WORK

Global context: Tweets are posted on the Social Media, Communication. The NER phrases are postes in tweets. The global context are derived in the from Web pages e.g., Microsoft Web N-Gram corpusor the Wikipedia therefore helps identifying the meaningful segments in tweets.

Local context: The method utilizing local language features is denoted by HybridSeg NER. It obtains the confident segments are based on the voting results of multiple NER tools. HybridSegNGram segments tweets can be estimate the term-dependency factor within a batches of tweets[3].

The tweet segmentation and NER are considered important subtasks in NLP. Many existing NLP techniques heavily an features, such as POS tags of the surrounding words, word capitalization, trigger words (e.g., Mr., Dr.). These glossary features, together with effective learning algorithms (e.g., hidden markov model (HMM) and conditional random field (CRF)), There have been a lot of attempts to incorporate tweet's unique characteristics into the conventional NLP. To improve POS tagging on tweets, train a POS tagger by using CRF model with tweet-specific features[1].

Length normalization: As the tweet segmentation is to extraction of the meaningful phrases into the segment forms, longer phrases can be preferred for preserving more topically specific meanings of that segments.

Presence in Wikipedia: In framework of Wikipedia serves as an external dictionary of valid or invalid names or phrases or the segments. Each text in Wikipedia refers to a Wikipedia entry even if the entry has not been created.

Tweet Datasets: There are two tweet datasets in our segmentation: SIN and SGE. The two datasets were used for simulating two targeted Twitter streams that should be divided in to the batches or POS. Randomly select 5; 000 tweets are published on one day in each tweet collection. After discarding retweets and tweets with inconsistent annotations, 4; 422 tweets from SIN and 3; 328 tweets from SGE are used for evaluation. The agreement of annotation on tweet level is 81% and 62% for SIN and SGE respectively[8].

III. HYBRIDSEG FRAMEWORK

The proposed Hybrid Seg framework segments in batch mode or POS. A tweets from a Twitter stream are grouped into batches or POS by their publishing time using a fixed time interval day by day invention of the tweets. Each batch of tweets are segmented by the Hybrid Segment.

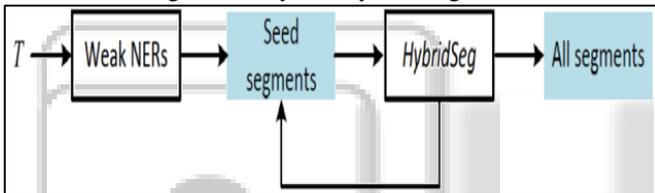


Fig. 2: Process for Hybrid Segmentation

Named Entity Recognition: Here select NER as a downstream application to demonstrate the benefit of tweet segmentation to investigation of segment based NER algorithms. The Named Entities from a pool of segments for exploiting the co-occurrences of named entities of the sements.

Evaluate the accuracy of named entity recognition based on segments or the tweets. The two NER methods, Random Walk-based (RW-based) and POS-based (Part-Of-Speech) of the NER NLP. Through out the result answer two questions: (i) which methods is more effective?, and (ii) does better segmentation can be lead to better NER accuracy?

So evaluate five variations of the two methods namely Global Seg RW, Hybrid Seg RW, Hybrid Seg POS, Global Seg POS, Unigram POS.

Here GlobalSeg denotes 1) HybridSegWeb since it only uses global context, And 2) HybridSeg refers to HybridSegIter, the best method.

IV. TWITTER ARCHITECTURE

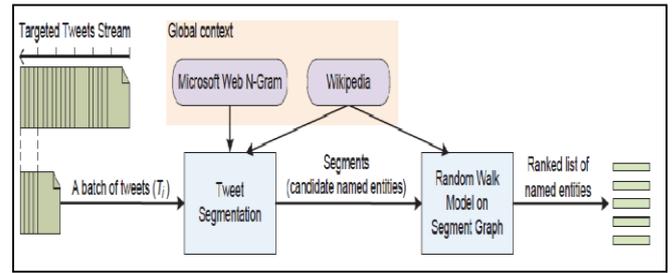


Fig. 3: System Architecture for Tweet NER

A. Tweet Segmentation

Given a tweet t from batch T , the problem of tweet segmentation is to split the words in $t = w_1 w_2 : : w_n$ into $m \leq 1$ consecutive segments, $t = s_1, s_2, : : : : : s_m$, where each segment of s_i contains one or more than words for the segment splitting.

See the detail solution for tweet segmentation. Given an individual tweet $t \in T_i$, the problem of tweet segmentation is to split t into m segments like RW & POS, the $t = s_1, s_2, : : : : : s_m$; each tweet segment contains one or more words. To obtain the optimal segmentation, we use the following objective function, where C is the function that can be measures the stickiness of a segment or a tweet defined based on word collocation of the tweets:

$$\text{Argmax} C(t) = \sum_{i=1}^n C(s_i)$$

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 , there exists $2^l - 1$ possible segmentations. It should be inefficient to iterate all of them and compute their stickiness[3].

B. Observations for Tweet Segmentation

Tweets are considered noisy with lots of informal abbreviations and grammatical errors. An tweets can be posts mainly information sharing and communication on the social media among many purposes.

1) Observation 1:

Word collocations of NLP and common phrases in English are well preserved in Tweets. Many entities and common phrases are in tweets for information sharing and dissemination. In this sense, $Pr(s)$ can be estimated by counting a segment's appearances in a very large English corpus (i.e., global context). In our implementation, we turn to Microsoft Web N-Grams.

2) Observation 2:

Many tweets can be having linguistic features. Although many tweets contain unreliable linguistic features like misspellings and unreliable capitalizations there exist tweets composed in proper English. The features in these tweets can be enable with NER relatively high accuracy.

3) Observation 3:

Tweets in a batches are not topically independent with each other within a time window. Many tweets published within a short time period talk about the same theme. These similar tweets largely share the same segments. For example, similar tweets have been grouped together to collectively detect events, and an event can be represented by the common discriminative segments across tweets [1].

Tags	Definition	Example
N	N common noun (NN, NNS)	books; someone
^	proper noun (NNP, NNPS)	lebron; usa; iPad
\$	numeral (CD)	2010; four; 9:30

Table 1:

Three POS tags as the indicator of a segment being a noun phrase, reproduced.

V. RESULTS

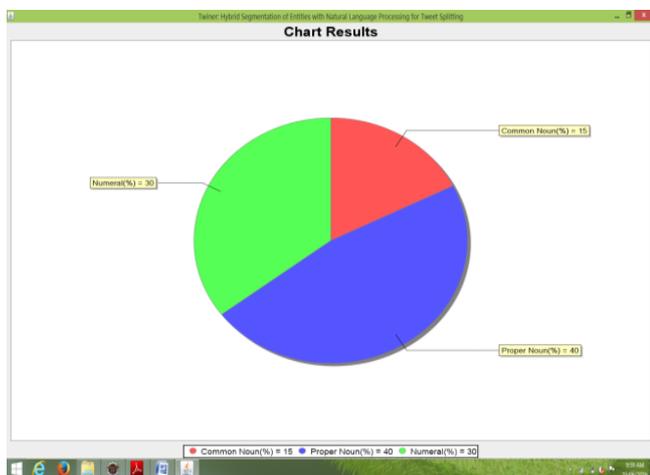


Fig. 3:

A Result can be Shows the Contents Proper Noun(40%), Common Noun(15%), Numerical Noun(30%).

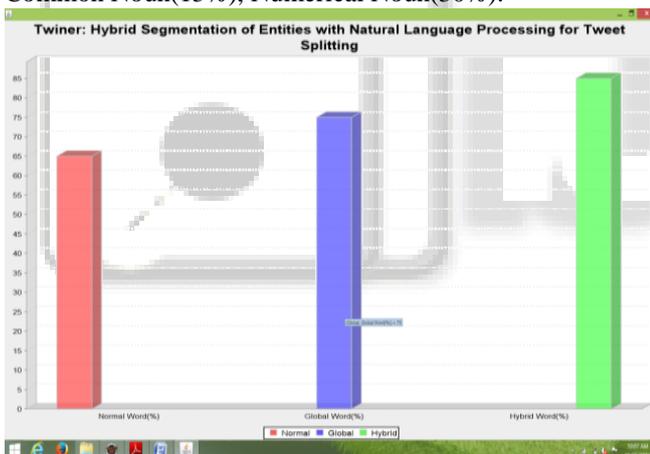


Fig. 4:

VI. CONCLUSION

As Present the HybridSeg framework which segments tweets into meaningful phrases called segments using both the global context and local context. Through our framework, demonstrate that local features are more reliable than term dependency in guiding the segmentation process.

ACKNOWLEDGMENT

We would like to thank Our Principal Dr.G.U.Kharat for valuable guidance at all steps while framing this paper. We are extremely thankful to P. G. Coordinator Prof. S. A.Kahate for guidance and review of this paper. I would also like to thanks the all faculty members of "Sharadchandra Pawar College of Engineering, Otur(M.S.), India".

REFERENCES

- [1] 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, SUBMISSION 2013.
- [2] 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, VOL. 27, NO. 2, FEBRUARY 2015.
- [3] 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.
- [4] S. Guo, M.-W. Chang, and E. Kiciman, "To link or not to link? A study on end-to-end tweet entity linking," in Proc. Conf. North Amer. Chapter Assoc. Comput. Linguistics: Human Language Technol., 2013, pp. 1020–1030.
- [5] A. Sil and A. Yates, "Re-ranking for joint named-entity recognition and linking," in Proc. 22nd ACM Int. Conf. Inf. Knowl. Manage.2013, pp. 2369–2374.
- [6] J. Gao, M. Li, C. Huang, and A. Wu, "Chinese word segmentation and named entity recognition: A pragmatic approach," in Comput. Linguist., vol. 31, pp. 531–574, 2005.
- [7] Y. Zhang and S. Clark, "A fast decoder for joint word segmentation andpos-tagging using a single discriminative model," in Proc. Conf. Empirical Methods Natural Language Process., 2010,pp. 843–852.
- [8] W. Jiang, L. Huang, and Q. Liu, "Automatic adaption of annotation standards: Chinese word segmentation and pos tagging—A case study," in Proc. Joint Conf. 47th Annu. Meeting ACL 4th Int. Joint Conf. Natural Language Process. AFNLP, 2009, pp. 522–530.
- [9] X. Zeng, D. F. Wong, L. S. Chao, and I. Trancoso, "Graph based semi supervised model for joint chinese word segmentation and part-of-speech tagging," in Proc. Annu. Meeting Assoc.Comput. Linguistics, 2013, pp. 770–779.
- [10] W.Jiang, M.Sun "Discriminative learning with natural annotations: Word segmentation as a case study,"in Proc. Annu. Meeting Assoc. Comput. Linguistics, 2013, pp.761–769.
- [11]R. Mihalcea and A. Csomai, "Wikify!:. linking documents to encyclopedic knowledge," in Proc. 16th ACM Conf. Inf. Knowl. Manage.,2007, pp. 233–242.
- [12]W. Wu, H. Li, H. Wang, and K. Q. Zhu, "Probase: A probabilistic taxonomy for text understanding," in Proc. ACM SIGMOD Int.Conf. Manage. Data, 2012, pp. 481–492.
- [13]K. Wang, C. Thrasher, E. Viegas, X. Li, and P. Hsu, "An overviewof microsoft web n-gram corpus and applications," in Proc. HLT-NAACL Demonstration Session, 2010, pp. 45–48.
- [14]F. A. Smadja, "Retrieving collocations from text: Xtract," Comput. Linguist., vol. 19, no. 1, pp. 143–177, 1993.
- [15]H. Ney, U. Essen, and R. Kneser, "On structuring probabilistic dependences in stochastic language

- modelling,” *Comput. Speech Language*, vol. 8, pp. 1–38, 1994.
- [16] S. F. Chen and J. Goodman, “An empirical study of smoothing techniques for language modeling,” in *Proc. 34th Annu. Meeting Assoc. Comput. Linguistics*, 1996, pp. 310–318.
- [17] Miss. Varpe Kanchan Nanasahab B.E. in Computer Engg. SPPU (2014). Diploma in Computer MSBTE (2011). B.Com. YCMOU (2011). Post Graduate Student, Computer Engg, Sharadchandra Pawar College Of Engineering, Dumbarwadi, Otur, Junner. Pune. (M.S.) India. Savitribai Phule Pune University. Pune-412409. +91-976199606 varpe.kanchan8@gmail.com
- [18] Prof. Sandip A. Kahate B.E. in computer science and engineering from Amravati university, M.E. in Wireless Communication and Computing, from Nagpur University and preparation for Ph. D. registration. He is currently working as an Assistant Professor in Computer Engineering Department, Sharadchandra Pawar College of Engineering, At. Post-Otur, Dist-Pune-412409 (M.S.), India. +91-7028068147 sandip.kahate@gmail.com
- [19] He has 10 years of teaching experience. He is author of 1 research paper, with around 10 papers in international journal and 3 in international conference in India and abroad. His areas of interest are Wireless Communication and computing, network security and Ad-Hoc Network.

