

AN EFFICIENT APPROACH FOR ANALYSING THE TWITTER POSTS USING CLASSIFICATION TECHNIQUES

Manikandan V¹, Rajkumar S²

¹Student, K.S.Rangasamy Coll ege of Technology, Tiruchengode, Tamilnadu, India.

²Assistant Professor, K.S.Rangasamy College of Technology, Tiruchengode, Tamilnadu, India.

Abstract: Twitter may be a stimulating platform for the dissemination of stories. The real-time nature and brevity of the tweets are conducive to sharing of knowledge related to important events as they unfold. But, one of the simplest challenges is to hunt out the tweets that characterize as news within the ocean of tweets. The paper propose a totally unique method for detecting and tracking breaking news from Twitter in real-time. Filtering the stream of incoming tweets to urge obviate junk tweets employing a text classification algorithm and also compare the performance of varied supervised SVM text classification algorithms for this task. After classification then cluster similar tweets, so that, tweets within the same cluster relate to an equivalent real-life event and should be termed as a breaking news. Finally, rank the news employing a dynamic scoring system which also allows us to trace the news over a period of some time.

Keywords: SVM, Twitter.

1. INTRODUCTION:

The real-time nature and shortness of the tweets encourages user to speak real-time events using slightest of text. Sakaki et al. used Twitter for early detection of earthquakes within the hope of sending word about them before they even hit. In fact, thanks to this real-time nature, Twitter are often used as a sensor to collect up-to-date information about the state of the planet. The goal is to style a system to be used for detecting and tracking breaking news in real-time on Twitter.

An approach to detect and track breaking news in presence of noisy data stream without counting on traditional news publishers. We evaluate different algorithms which classify tweets as either news or junk. Also show how a standard density based clustering algorithm are often used for detecting clusters during a stream of streaming data and propose a singular technique to parallelize classification of tweets using RabbitMQ. Finally, the paper also proposes a completely unique dynamic rating system for ranking and tracking news.

Twitter:

Twitter is an online news and social networking service where users post and interact with messages, known as "tweets." These messages were originally restricted to 140 characters, but on November 7, 2017, the limit was doubled to 280 characters for all languages except Japanese, Korean and Chinese. Registered users can post tweets, but those who are unregistered can only read them. Users access Twitter

through its website interface, Short Message Service (SMS) or mobile device application software Twitter, Inc. is based in San Francisco, California, United States, and has more than 25 offices around the world.

Social networks witness a rapid growth development and recommendation on social networks is becoming essential. Recommendation of quality and useful information to the users is an issue here. How to identify user's interest and provide customized recommendation for each user becomes a challenging problem. In this project recommendation of news articles to the users of twitter is proposed. This project aims to investigate a framework to combine tag correlation and user social interest for recommendation. News related tweet contents are extracted and tags for the contents are generated using context inference and user preference analysis. A user tag retrieval strategy is developed to assign tags for users and a user-tag matrix is created to provide the initial weights for user's tags. Every user will have different or similar tags based on his/her social interest. News articles relating to the users tag are recommended for high quality and interest. A user-tag matrix if formed for every user. The tags with higher ranking is considered for recommendation. N-gram extraction is used to find the occurrence of similar tags. RankSVM algorithm is used to rank the different tags of the same user and Rank Aggregation algorithm is used for ranking the overall tags hence a customized recommendation of news articles is provided to the users.

News Detection and Tracking:

Twitter has been used as one of the communication channels for spreading breaking news. We propose a method to collect, group, rank and track breaking news in Twitter. Since

short length messages make similarity comparison difficult, we boost scores on proper nouns to improve the grouping results. Each group is ranked based on popularity and reliability factors. Current detection method is limited to facts part of messages. We developed an application called “Hotstream” based on the proposed method. Users can discover breaking news from the Twitter timeline. Each story is provided with the information of message originator, story development and activity chart. This provides a convenient way for people to follow breaking news and stay informed with real-time updates.

Breaking news is defined by Wiktionary as “news that has either just happened or is currently happening. Breaking news may contain incomplete information, factual error or poor editing because of rush.” With this definition Twitter can fit the needs of breaking news delivery.

However, news posted in Twitter requires an effort to discover it. Firstly, users often have problems of deciding which users to follow. That is, to find users with interesting tweets. Secondly, users need to read through status updates and follow links to obtain further information. To ease these problems and to deliver breaking news effectively, we propose a method to collect, group, rank and track breaking news in Twitter. This work is a contribution to the area of Topic Detection and Tracking (TDT). The tasks we focus are first story detection, cluster detection, and tracking.

Online Clustering:

Clustering is the grouping of a particular set of objects based on their characteristics, aggregating them according to their similarities. Regarding to data mining, this methodology partitions the data implementing a specific join algorithm, most suitable for the desired information analysis.

This clustering analysis allows an object not to be part of a cluster, or strictly belong to it, calling this type of grouping hard partitioning. In the other hand, soft partitioning states that every object belongs to a cluster in a determined degree. More specific divisions can be possible to create like objects belonging to multiple clusters, to force an object to participate in only one cluster or even construct hierarchical trees on group relationships.

There are several different ways to implement this partitioning, based on distinct models. Distinct algorithms are applied to each model, differentiating its properties and results. These models are distinguished by their organization and type of relationship between them.

Classification of Tweets:

Millions of users share opinions on various topics using micro-blogging every day. Twitter is a very popular

microblogging site where users are allowed a limit of 140 characters; this kind of restriction makes the users is concise as well as expressive at the same time. For that reason, it becomes a rich source for sentiment analysis and belief mining. The aim of this paper is to develop such a functional classifier which can correctly and automatically classify the sentiment of an unknown tweet. In our work, we propose techniques to classify the sentiment label accurately. We introduce two methods: one of the methods is known as sentiment classification algorithm (SCA) based on k-nearest neighbor (KNN) and the other one is based on support vector machine (SVM). We also evaluate their performance based on real tweets.

These days social networks, blogs, and other media produce a huge amount of data on the World Wide Web. This huge amount of data contains crucial opinion related information that can be used to benefit for businesses and other aspects of commercial and scientific industries. Manual tracking and extracting this useful information from this massive amount of data is almost impossible. Sentiment analysis of user posts is required to help taking business decisions. It is a process which extracts sentiments or opinions from reviews which are given by users over a particular subject, area or product in online.

Categorize the sentiment into two types: 1) positive or 2) negative that determine the general attitude of the people to a particular topic. Our principal goal is to correctly detect sentiment of tweets as more as possible. This paper has two main parts: the first one is to classify sentiment of tweets by using some feature and in the second one we use machine learning algorithm SVM. In both the cases, we use five-fold cross validation method to determine the accuracy. We propose two approaches for sentiment analysis. One of the technique facilitates KNN and the other uses SVM.

2. EXISTING SYSTEM:

In NB (Naïve Bayes), documents are projected into a low dimensional topic space by assigning each word with a latent topic. It employs an extra generative process on the topic proportion of each document and models the whole corpus via a hierarchical Bayesian framework. The BoW representation disregards the linguistic structures between the words. It the consumer expectation not predicted clearly. Less accuracy prediction on opinion analysis. User review based word alignment is cumbersome. High in latency to analyze the datasets. We can find the topic distribution for each of the document and compare them for similarity. As these are probability distributions, we make use of a modified KL-divergence method. Querying makes use of similarity ranking to find the documents which are most similar to a given a query.

Thus system can retrieve the data properly from database and also get keyword ratings explicitly from the users. In tlets based collaborative filtering technique, tlets similarity computation and prediction computation modules have been implemented. Recommended lists are generated on non-purchased keywords of login user. So I have computed system predicted ratings for all non-purchased keywords of login user. To calculate system predicted rating of target keyword, first I have obtained 5 most similar tlets and then used lighted sum approach for rating prediction computation.

REFERENCES:

1. Burnap, P., Rana, O., Williams, M., Housley, W., Edwards, A., Morgan, J., Sloan, L., Conejero, J., 2015. COSMOS: towards an integrated and scalable service for analysing social media on demand. *Int. J. Parallel Emergent Distrib. Proceedings of the IEEE Computer Vision and Pattern Recognition*.
2. Campbell, H., Engelbrecht, I., 2018. The Baboon Spider Atlas – using citizen science and the ‘fear factor’ to map baboon spider (Araneae: Theraphosidae) diversity and distributions in Southern Africa. *International Journal of Computer Applications*, Vol. 120, no. 15, 2015.
3. Carter, S., Iterkamp, W., Tsagkias, M., 2016. Microblog language identification: overcoming the limitations of short, unedited and idiomatic text. *Lang. Resour. Eval.* 47,195215. *Proceedings of the IEEE mag. on Computer Vision and Pattern Recognition*
4. Venter, O., Sanderson, E.W., Magrach, A., Allan, J.R., Behr, J., Jones, K.R., Possingham, H.P., Laurance, W.F., Wood, P., Fekete, B.M., Levy, M.A., Watson, J.E.M., *Proceedings of the IEEE mag. on global implications for biodiversity conservation. Nat. Commun.* 7, 12558
5. Cong, L., Wu, B., Morrison, A.M., Shu, H., Wang, M., 2016. Analysis of wildlife tourism experiences with endangered species: an exploratory study of encounters with giant pandas in Chengdu, China. *International Journal of Science, Engineering and Technology Research (IJSETR)*, vol 4, no. 4, pp. 816-820, 2015.
6. Crampton, J.W., Graham, M., Poorthuis, A., Shelton, T., Stephens, M., Wilson, M.W., Zook, M., 2016. Beyond the geotag: situating ‘big data’ and leveraging the potential of the geotag. In *2016 International Journal on Advanced Computer Science Applications and Technologies (ACSAT)*, pp. 379-384. IEEE, 2016.
7. van Zanten, B.T., Van Berkel, D.B., Meentemeyer, R.K., Smith, J.W., Tieskens, K.F., Verburg, P.H., 2016. Continental-scale quantification of landscape values using social media data. *Proc. Natl. Acad. Sci.* 113, 12974–12979. 2015. *International Journal of Emerging Research in Management & Technology*, Vol 4, no. 11, pp. 110-118, 2015.
8. Croitoru, A., Wayant, N., Crooks, A., Radzikowski, J., Stefanidis, A., 2016. Linking cyber and physical spaces through community detection and clustering in social media feeds. *Comput. Environ. Urban. Syst.* 53, 47–64. *Proceedings of the fourth International Journal on weblogs and social media.* 2017.
9. Gliozzo, G., Petteorelli, N., Muki Haklay, M., Using crowdsourced imagery to detect cultural ecosystem services: a case study in South Wales, UK. *Ecol. Soc.* 21, art6. *Journal paper on Visual Computing* on 2016.
10. Huifang Ma, MeihuiziJia, Di Zhang, Xianghong Lin, “Combining tag correlation and user social relation for microblog recommendation”. *IEEE Paper on Computer Vision and Pattern Recognition*(2016)
11. Grabowicz, P.A., Ramasco, J.J., Moro, E., Pujol, J.M., Eguiluz, V.M., 2017. Social features of online networks: the strength of intermediary ties in online social media. *PLoS One* 7, e29358. *IEEE Paper on Advances in Neural Information Processing Systems (NIPS)*
12. Greer, K., Day, K., McCutcheon, S., 2017. Efficacy and perception of trail use enforcement in an urban natural reserve in San Diego, California. *J. Outdoor Recreat. Tour.* 18,56–64. 2016 *European paper on Computer Vision (ECCV)*.
13. Hampton, S., Strasser, C., Tewksbury, J., Gram, W., Budden, A., Batcheller, A., Duke, C., Porter, J.H., 2018. Big data and the future of ecology. *Front. Ecol. Environ.* 11, 156–162. in *Proc. 2013 IEEE paper on Applications of Computer Vision*.
14. Hausmann, A., Toivonen, T., Heikinheimo, V., Tenkanen, H., Slotow, R., Di Minin, E., 2017. Social media reveal that charismatic species are not the main attractor of ecotourists to sub-Saharan protected areas. *Sci. Rep.* 7, 763. 2015 *IEEE paper on Computer Science. Springer Verlag, Berlin, Germany.*
15. Hausmann, A., Toivonen, T., Slotow, R., Tenkanen, H., Moilanen, A., Heikinheimo, V., Di Minin, E., 2018. Social media data can be used to understand tourists' preferences for nature-based experiences in protected areas. *Conserv. Lett.* 11. In *Proceedings of the ACL 2014 IEEE paper on Language Technologies and Computational Social Science*
16. Thelwall, M., 2009. Social network sites. Users and uses. In: *Zelkowitz, M. (Ed.), Advances in Computers. Elsevier*, pp. 19–73. *International Journal of Linguistic Annotation (LAW IX)*, pages 112–123
17. Tieskens, K.F., Schulp, C.J.E., Levers, C., Lieskovský, J., Kuemmerle, T., Plieninger, T., Verburg, P.H., 2017 In *Proceedings IEEE paper of Characterizing European cultural landscapes: accounting for structure, management intensity and value of agricultural and forest landscapes. Land Use Policy* 62, 29–39.
18. Van Berkel, D.B., Tabrizian, P., Dorning, M.A., Smart, L., Newcomb, D., Mehaffey, M., Neale, A., Meentemeyer, R.K., 2018 In *Proceedings IEEE paper of Quantifying the visual-sensory landscape qualities*

that contribute to cultural ecosystem services using
social media and LiDAR. *Ecosyst. Serv.* 31, 326–335.