# Search optimization of named entities from twitter streams

To cite this article: K Mohammed Fazeel *et al* 2017 *IOP Conf. Ser.: Mater. Sci. Eng.* **263** 042042

View the article online for updates and enhancements.

## Related content

# Search optimization of named entities from twitter streams

**Mohammed Fazeel K, Simama Hassan Mottur , Jasmine Norman and Mangayarkarasi R**

School of Information Technology and Engineering, VIT University, Vellore-632014, Tamil Nadu, India.

E-mail: jasmine@vit.ac.in

**Abstract**. With Enormous number of tweets, People often face difficulty to get exact information about those tweets. One of the approach followed for getting information about those tweets via Google .There is not any accuracy tool developed for search optimization and  as well as getting information  about those tweets. So, this system contains the search optimization and functionalities for getting information about those tweets. Another problem faced here are the tweets that contains grammatical errors, misspellings, non-standard abbreviations, and meaningless capitalization. So, these problems can be eliminated by the use of this tool. Lot of time can be saved and as well as by the use of efficient search optimization each information about those particular tweets can be obtained.

## 1.Introduction

Twitter is one of the most used social media by celebrities, business tycoons all over the world. Twitter is used by most of the people in the world, a lot of information's, posts, news and interactions are made by celebrities, business tycoons. The administration quickly increases world-wide prominence. Twitter tweet's are limited to only 140 characters strings. Celebrities, business tycoons and fans are being interacted using tweets personally by their respective user names. Twitter's tweets are trending these days as tweets of celebrities and business tycoons are shown in the news now days. Twitter's tweets post by using #tags. If the users of the twitter tweet about the several topics and issues consistently, the tweets become trending all over the world. So that, everyone knows about the most discussed topic across the globe.
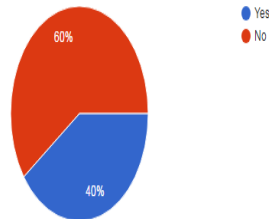
## 2.Literature Survey

Do you use twitter?



**Figure. 1 Survey result of first Question**

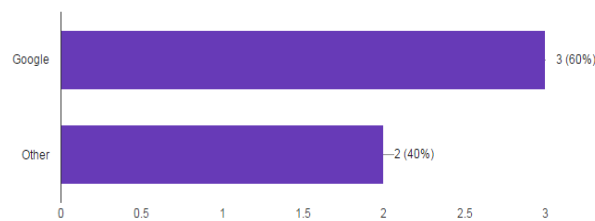How do you get to know about #tag tweets?



**Figure. 2 Survey result of Second Question**
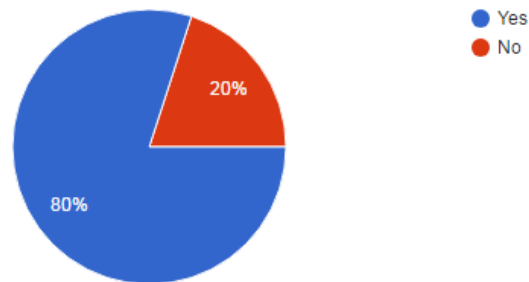
Do you need any tool for search optimization?



**Figure. 3 Survey result of third Question**

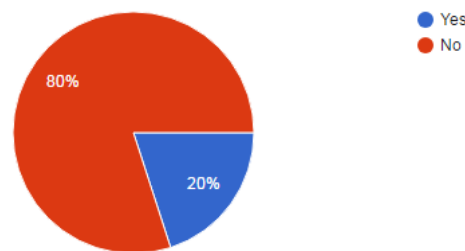Do you get full information about those #tweets?



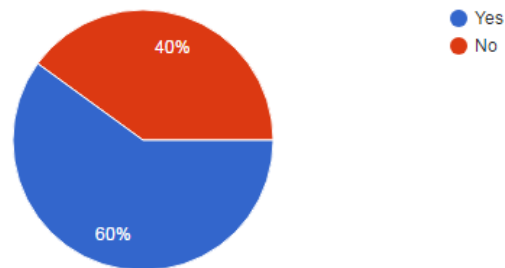**Figure. 4 Survey result of fouth Question**

**Figure. 5 Survey result of fifth Question**

From the referred papers gain the information of  normalization of social media verbal[4] ,conversion of grammatical mistake text to proper English[2] ,able to get the knowledge of  removal of confusion from context, recognition of named entity from a context[6] and way of twitter normalization with the concept of 1 to n[5] ,ways of classifying tweets[1]. In the Existing System , there are several tools where normal text is disambiguated its accuracy is around  90%  because it contains without any grammatical errors,  but for tweets but its accuracy lies between 40-50% because tweets contain Accents, Abbreviation, Misspelled words  and extra special characters which reduces the accuracy. Search Optimization is not efficient.

## 3.Proposed System

The proposed system consists of Search optimization, .In which tweets areconverted into information.By, which a lot of tweets accuracy can be improved. And detailed Information can be acquired with.in a short span of time. Including that, filtration of grammatical errors, Misspelled words, non- standard abbreviations and meaningless capitalizations are done. A lot can time be saved.

Methodologies

In the initial part tweets are retrieved and then normalization of tweets are done.Take each entity  link it with Wikipedia to disambiguate, from there we could get the entire Details about the given nouns in the tweets ,by following this method we are saving the Time and providing an effective search operation for the twitter Streams.
It contains four phases: Data collection, pre-processing of tweets and normalization     Valuation

### 3.1Data Collection
- Create an api in Twitter.
- Generate Tokens
- Make Use of Twitter package to fetch the tweets.
- Retrieve tweets from Twitter database.
- Store in local database for future Usage and Remove Extra Fields.

### 3.2Preprocessing of Tweets And Normalization

[4]By Using the Regular Expression we negate the unnormalised tweets
removal of urls
removal of non-alphanumeric characters
removal of hash tags
removal of retweet tags

*3.3 Use Of NTLK*

With the help of the NTLK package (Python), the normalized tweets are preprocessed
To extract the name entities.
[2]The extracted Entities are stored in a Separate table


**Figure. 6 NTLK**

*3.4 Valuation*

- Take the named Entities.

*3.5 Linking*

- Take each entity

- Link it with Wikipedia to disambiguate.


**Figure. 7 Linking**

**4.Conclusion**

As the data's are very large in twitter doing search optimization plays a major role in order to improve the accuracy rate. By the increase of accuracy rate, Detailed Information can be obtained. This can have a lot of impact in the years to come. Users don't need to go through Search engines every time. By using search optimization tool, accuracy of tweets information can be improved.  As well as. People will be able to access the quality factors like user-friendly, Reliability.   .

**References**

[1]     David Harrison Building a taxonomy of tweet types and automatically classifying  tweets in to
            these types

[2]     Neelmay Desaia and Narvekarb 2015 Normalization of Noisy Text Data *International Conference on Advanced Computing Technologies and Applications  ICACTA*

[3]     DrlimlianTze 2015 A Practical Introduction to natural language Processing

[4]     Eleanor Clarka and kenji Arakia 2014 Text Normalization in social Media:Progress,problems and Applications for a preprocessing System of casual English *Pacific Association for computational Linguistics (PACLING)*

[5]     Yafeng ren , Jiayuan  Deng and Donghong 2016 Twitter Normalization via 1 to N Recovering *Springer International Publishing AG*

[6]     Rizzo G and Trancy R NERD:A framework for evaluating Named Entity Recognition Tools in the web of data ISWC'11

[7]     Gurpreet Singh Khanuja and Sachin Yadav 2013 Normalisation of SMS Text IIT Kanpur Conference of NLP 2013

[8]     Catherine Kobus 2008 Normalizing SMS: are two metaphors better than one