

## CLASSIFICATION OF TEXT DATA FROM THE SOCIAL NETWORK TWITTER

I.A. Rytsarev, A.V. Blagov

Samara National Research University, Samara, Russia

**Abstract.** Social networks play an important role in the modern world, and it is important to define the important and popular topics discussed. This article deals with data collection from the social network Twitter, and further clustering and classification of the collected data.

**Keywords:** bigdata, data processing, data analysis, clustering, classification, tf-idf, latent dirichlet allocation.

**Citation:** Rytsarev IA, Blagov AV. Classification of Text Data from the Social Network Twitter. CEUR Workshop Proceedings, 2016; 1638: 851-856. DOI: 10.18287/1613-0073-2016-1638-851-856

### Introduction

The extra-large volumes of data in information technology - data sets the size of which is beyond the capabilities of typical database (DB) for collecting, storage, management and analysis of information [1]. There are many series of approaches, tools and methods of processing such extra large volumes of structured and unstructured data [1-4].

The concept of big data means working with the vast volume of information and varied composition, very frequently updated and located in different sources in order to increase efficiency and create new ones.

At the moment, social networks are at the peak of their popularity, millions of people are already use Facebook and Twitter. Many companies need to analyze the data collected from social networks to assess the relationship of users to their products [5]. Also the analysis of this area is used in solving security issues [6]. Having collected and clustered text data from the social network, it is possible to identify the main themes and events discussed by social network users in different cities and countries.



## 2 Classification of the text information based on the approach with machine learning

Also, there is another approach to solve the problem of the classification - the classification of the information through the machine learning.

Machine learning - the process through which the machine (computer) is able to show the behavior which was not explicitly programmed. There are two types of learning: inductive and deductive.

In the work of researchers involved in cluster analysis of the text information in various search engines, the inductive measure Word2vec is frequently used [10-11].

The principle of the measures is to find connections between the context, the word according to the assumption that the words that are in similar contexts, tend to mean similar things, i.e., be semantically close. More formally, the task is: to maximize the cosine proximity between vectors of the words (scalar product of the vectors) that appear next to each other, and minimizing the cosine proximity between the vectors of words that do not appear next to each other. "Next to each other" in this case means "in close contexts".

Word2vec analyzes the use of words contexts and concludes that they are or are not close in meaning. Since word2vec making such conclusions based on large amounts of text, the conclusions are quite adequate. The algorithms on which word2vec is based are described in detail in [12-13].

The examples of vector distances obtained by word2vec are in Table 1.

**Table 1.** The vector distances between the word «France» and other words by word2vec measure

The word	The vector distance
Paris	0.978443
Spain	0.665923
Belgium	0.665923
Netherlands	0.652428
Italy	0.633130
Portugal	0.577154
Russia	0.571507
Germany	0.563291

One type of a deductive approach can be considered is the Latent Dirichlet Allocation (LDA). This generative model that allows to explain the results of observations with the help of implicit groups, that allows to receive an explanation of why some of the parts of data are similar. Typically, when using this approach, you identify a limited number of topics and further states that each document is a mixture of a small number of topics [15].

For more detailed analysis, it is best to combine different approaches and techniques depending on the amount of processed data.

### 3 The method of collecting information

To research the work of TF-IDF algorithm software tool that allows to collect data directly from the Twitter servers has been developed. The implementation is based on an open Twitter API 2.0 interface. As the object of the study the tweets from Samara region were taken, for this as a selection criterion of the messages geolocation was set to the Samara region (including all settlements of the region). All tweets collected in such a way need to be clustered using TF-IDF metric for short messages of 140 characters in length and the algorithm k-means.

To carry out the data collection from Twitter server a request containing a consumer key and consumersecret key was sent. In reply we have oauth.accessToken, oauth.accessTokenSecret which give us the ability to retrieve data from servers.

The second step is to send the query-request, in response to which a set of tweets returns.

Next, the third step is the counting the TF-IDF metric values for each message.

### 4 Results

The data was collected for 24 hours, by the query-request, which characterizes the Samara region. As the result over 6000 messages has been collected. By applying metrics TF-IDF and the k-means algorithm 22 clusters were obtained. On example, one of the obtained clusters (Figure 2) shows that the messages are similar in meaning, but among them there are messages with "foreign" theme.

```

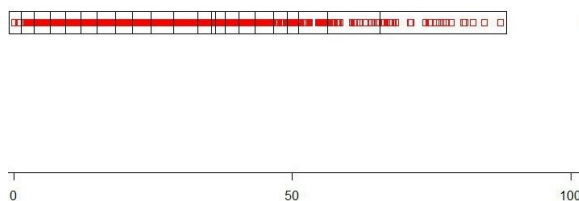
17.417364700059597:   пятьсот тридцать два потому что люблю бгс lovebts
17.412488867909754:   когда же будет тепло lovebts https t со 5wgenfbspt
17.377890466760327:   lovebts пожалуйста пожалуйста прими меня к себе
17.365013020434603:   just posted a photo https t со htavu34orf
17.33785524961285:   why did i wake up 4am
17.26244280531484:   сердечно in смэрэнь самарская https t со xmp1orway7
17.256147680915472:   ну не жожу я думать https t со xmp1orway7
17.245745182632252:   lovebts но все же если однажды увидимся улыбнись
17.242810395943955:   пятьсот двадцать два моя фантазия закончилась lovebts
17.18060503437332:   kiradream это я так у мамы хорошо получал
17.165552773189763:   lovebts я пою в одиночестве всё ту же песню
17.151261464268252:   четьреста пятьдесят четыре глаза улыбки чины lovebts
17.143256397595636:   обожаю такие вечера спасибо вам https t со k13y9ptnx1
17.138387528934295:   старый мост река самара https t со lkv9ocbfnu
17.122425825760548:   г n samara russia https t со zvn1656yup
17.100052153754355:   только версия мультиар lovebts с капитаном дэффи
17.087077012033237:   sakura65651 как время будет заеду в ателье
17.052100143709136:   samara in и сколько уголовных дел планируются
17.02051064681348:   моё утро после пьянки https t со so9ci1p7vbf
17.006630468592117:   r bar terrace https t со li4637xyu1
16.9393484458442:   обожаю людей которые спазывают https t со lndmizayn
16.937749179184305:   lovebts смотри я беспристрастен ко всем кроме тебя
16.935328573424083:   какой же день ужасный а мол быть крутым
16.91188651734057:   ahmythank ego second box 2
16.911526256447324:   ведь есть еще на свете джентльмены спасибо
16.88121521655293:   пятьсот тридцать еще немного твитов перед сном lovebts
16.87353404614601:   так не хочется никуда идти но надо lovebts
16.869827024667654:   lovebts милая просто скажи что хочешь расстаться
16.25709517099365:   уже просто не вывожу ситуацию которая происходит
15.84083916617625:   anastasya2614 блин я думала на 4 приехать

```

Fig. 2. An example of one of the clusters obtained

Apparently, such is not quite accurate result was obtained due to the fact that studied Twits have a 140 character limit. For this reason, for more accurate clustering and further classification it is necessary to modernize TF-IDF measure, introducing additional weighting coefficients corresponding to the number of symbols (words) in the message.

Moreover, high density of the clusters (Figure 3) shows the need for revision of the metric.



**Fig. 3.** The distribution of TF-IDF metric values of processed data on the number line

TF-IDF metric works for short messages with relative accuracy. For this reason it is necessary to optimize this metric: adding weight coefficients, the coefficients associated with the hashtag, the introduction of the normalization coefficient associated with the length of the message and the number of words in it, etc.

## Conclusion

Issues related to clustering and further classification of text data are relevant in relation to the enormous spread of social networks and online services worldwide. Approaches and techniques presented in the article are planned for testing on text data collected from Twitter social network in the Russian segment. Collecting the necessary data is being produced by means of the developed software system, based on time zones and geolocation. It is planned to develop the subject in the direction of the output and optimization of parallel clustering algorithms.

## References

1. Dean J, Ghemawat S. MapReduce: simplified data processing on large clusters. *Communications of the ACM*, 2008; 51: 107-113.
2. Vossen G. Big Data as the new enabler in business and other intelligence. *Vietnam Journal of Computer Science*, 2014; 1: 3-14.
3. Tamhane DS, Sayyad SN. Big Data Analysis Using Hace Theorem. *International Journal of Advanced Research in Computer Engineering & Technology (IJARCET)*, 2015; 4: 18-23.
4. Kazanskiy NL, Protsenko VI, Serafimovich PG. Comparison of system performance for streaming data analysis in image processing tasks by sliding window. *Computer Optics*, 2014; 38(4): 804-810.

5. Tan W, Blake M, Saleh I, Dustdar S. Social-network-sourced big data analytics. *IEEE Internet Computing*, 2013; 5: 62-69.
6. How “Big Data” help to improve security. URL: <http://www.computerra.ru/108760/security-n-big-data/>.
7. Data Mining tasks. Classification and cauterization [In Russian]. URL: <http://www.intuit.ru/studies/courses/6/6/lecture/166>.
8. Blagov A, Rytcarev I, Strelkov K, Khotilin M. Big Data Instruments for Social Media Analysis. *Proceedings of the 5th International Workshop on Computer Science and Engineering*, 2015; 179-184.
9. Ramos J. Using tf-idf to determine word relevance in document queries. URL: <https://www.cs.rutgers.edu/~mlittman/courses/ml03/iCML03/papers/ramos.pdf>.
10. Wang H. Introduction to Word2vec and its application to find predominant word senses. URL: <http://compling.hss.ntu.edu.sg/courses/hg7017/pdf/word2vec%20and%20its%20application%20to%20wsd.pdf>.
11. Yu M, Dredze M. Improving lexical embeddings with semantic knowledge. *Association for Computational Linguistics (ACL)*, 2013; 2: 545-550.
12. Mikolov T, Chen K, Corrado G, Dean J. Efficient Estimation of Word Representations in Vector Space. URL: <http://arxiv.org/pdf/1301.3781.pdf>.
13. Mikolov T, Sutskever I, Chen K, Corrado G, Dean J. Distributed Representations of Words and Phrases and their Compositionality. *Advances in neural information processing systems*, 2013; 3111-3119.
14. MacQueen J. Some Methods for Classification and Analysis of Multivariate Observations. In *Proc. of the 5th Berkeley Symposium on Mathematical Statistics and Probability*, 1967; 281-297.
15. Blei DM, Ng AY, Jordan MI. Latent dirichlet allocation. *The Journal of machine Learning research*, 2003; 3: 993-1022.