



Spatial Statistics 2015: Emerging Patterns

An Author-Topic based Approach to Cluster Tweets and Mine their Location

Mohamed Morchid^a, Yonathan Portilla^{a,b}, Didier Josselin^{c,a}, Richard Dufour^a, Eitan Altman^{b,a}, Marc El-Beze^a, Jean-Valère Cossu^a, Georges Linarès^a, Alexandre Reiffers-Masson^{a,b}

^aLaboratoire d'Informatique d'Avignon, LIA, 339 chemin des Meinajariès, Agroparc BP 91228, 84911 Avignon cedex 9, France

^bINRIA, B.P 93, 06902 Sophia Antipollis Cedex, France

^cUMR ESPACE 7300 ; 74 rue Louis Pasteur, 84029 Avignon Cedex, France

Abstract

Social Networks became a major actor in information propagation. Using the Twitter popular platform, mobile users post or relay messages from different locations. The tweet content, meaning and location show how an event-such as the bursty one “JeSuisCharlie” happened in France in January 2015 is comprehended in different countries. This research aims at clustering the tweets according to the co-occurrence of their terms, including the country, and forecasting the probable country of a non located tweet, knowing its content. First, we present the process of collecting a large quantity of data from the Twitter website. We finally have a set of 2.189 located tweets about “Charlie”, from the 7th to the 14th of January. We describe an original method adapted from the Author-Topic (AT) model based on the Latent Dirichlet Allocation method (LDA). We define a homogeneous space containing both lexical content (words) and spatial information (country). During a training process on a part of the sample, we provide a set of clusters (topics) based on statistical relations between lexical and spatial terms. During a clustering task, we evaluate the method effectiveness on the rest of the sample that reaches up to 95% of good assignment.

© 2015 The Authors. Published by Elsevier B.V.

Peer-review under responsibility of Spatial Statistics 2015: Emerging Patterns committee.

Keywords: Author-Topic model, Tweet location

1. Context of the study and state of the art

The exponential growth of available data on the Web enables users to access a large quantity of information. Micro-blogging platforms evolve in the same way, offering an easy way to disseminate ideas, opinions or common facts under the form of short text messages. Depending on the sharing platform used, the size of these messages can

1878-0296 © 2015 The Authors. Published by Elsevier B.V.

Peer-review under responsibility of Spatial Statistics 2015: Emerging Patterns committee.

be limited to a maximum number of words or characters. Although Twitter is a recent information-sharing model, it has been widely studied. Many works have focused on various aspects of Twitter, such as social impact [1] event detection [2], user influence [3], sentiment analysis [4], hash-tag analysis [5] or theme classification [6].

The aim of the proposed approach is to locate a given tweet by using the tweet content (a set of words). Nonetheless, the Twitter service does not allow to send messages whose size exceeds 140 characters. This constraint causes the use of a particular vocabulary that is often unusual, noisy, full of new words, including misspelled or even truncated words [7]. Indeed, the goal of these messages is to include a lot of information with a small number of characters. Thus, it may be difficult to understand the meaning of a short text message (STM) with only the tweet content (words). Several approaches have been proposed to represent the tweet content. The classical bag-of-words approach [8] is usually used for text document representation in the context of keyword extraction. This method estimates the Term Frequency-Inverse Document Frequency (TF-IDF) of the document terms. Although this unsupervised approach is effective for a large collection of documents, it seems unusable in the particular case of short messages since most of the words occur only once (*hapax legomena* [9]).

Other approaches propose to consider the document as a mixture of latent topics to work around segments of errors. These methods build a higher-level representation of the document in a topic space. All these methods are commonly used in the Information Retrieval (IR) field. They consider documents as a bag-of-words without taking account of the words order. Nevertheless, they demonstrated their performance on various tasks. Several approaches were proposed such as Probabilistic LSA (PLSA) [10] or Latent Dirichlet Allocation (LDA) [11]. LDA is a generative model of statistics which considers a document, seen as a bag-of-words, as a mixture probability of latent topics. In opposition to a multinomial mixture model, LDA considers that a theme is associated with each occurrence of a word composing the document, rather than associating a topic with the complete document.

Thereby, a document can belong to different topics from a word to another. However, it is noted that the word occurrences are connected by a latent variable which controls the global respect of the topic distribution in the document. These latent topics are characterized by a distribution of word probabilities which are associated with them. During the LDA learning process, distribution of words into each topic is estimated automatically. Nonetheless, the location associated with the tweet is not directly taken into account in the topic model. As a result, such a system considers separately the tweet content (words), to learn a topic model, and the labels (location) to train a classifier. Thus, the relation between the tweet content and its location (country) is crucial to efficiently locate (unknown) new tweets.

In this paper, we propose to build a topic model, called author-topic (AT) [12,13] that takes into consideration all information contained in a tweet: the content itself (words), the label (country) and the relation between the distribution of words into the tweet and the location, considered as a latent relation. From this model, a vector representation in a continuous space is built for each tweet. Then, a supervised classification approach, based on Support Vector Machines (SVM) [14] is applied. For mathematical and methodological details, see [6, 13].

2. Experimental protocol applied on the tweet “Charlie”

We propose to evaluate the approach on a Twitter corpus. This corpus is composed of tweets from 16 countries. A classification approach based on Support Vector Machines (SVM) is performed to find out the most likely country of a given tweet, in two stages: a training and an assignment [14]. The table 1 lists the corpus of tweets. This data set is split in three parts depending of the tweet emission day in January 2015: 887 tweets between the 7th and the 8th, 471 tweets between the 9th and the 10th and 881 tweets between the 11th and the 14th.

We obtain 1.520 tweets for the training phase of the AT models and 669 for the validation (testing) phase which corresponds to a corpus of 2.189 tweets for the whole 16 countries (roughly 137 tweets for each country). The number of topics contained in the AT model strongly influences the quality of this model. Indeed, an AT model with only few topics will be more general than one with a large number of classes (granularity of the model). For a sake comparison, a set of 100 AT models is learnt (between 5 and 105). As the classification of tweets requires a multi-class classifier, the SVM (one-against-one) method is chosen with a linear kernel. This method gives a better accuracy than the (one-against-rest) one [15].

Table 1. Number of tweets in each time period (January, 2015).

Country name	7 th to 8 th		9 th to 10 th		11 th to 14 th	
	Train	Test	Train	Test	Train	Test
France	287	124	171	74	259	111
United-Kingdom	90	39	58	25	99	43
United-States	77	34	58	26	110	48
Brazil	30	13	11	5	32	15
Italia	28	12	16	7	20	9
Spain	20	9	14	6	14	6
Turkey	14	7			13	6
The Nederland	16	8			7	3
Canada	9	5			7	4
Belgium	9	5				
Mexico	8	4				
Colombia	9	5				
Philippines					9	4
Argentina					8	4
Germany	8	4				
India	9	4				

3. Results of the classification

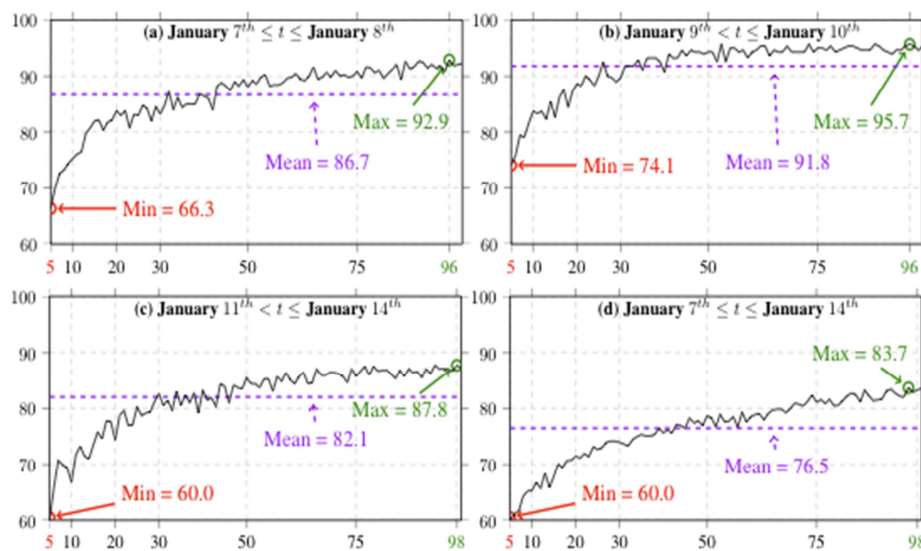


Fig. 1. Country classification accuracy (%) using various author topic-based representations on the test sets with different time periods. The X axis represents the number of classes contained in the topic space (between 5 and 100).

A main event named here “Je suis Charlie” (bursty event on Internet sharing platforms) happened in Paris during the January 7th to 14th 2015 epoch. The figure 1 shows the accuracies (percentage of countries found) obtained during the country location task of this event (tweets) for the different time slots.

The first remark is that the higher the number of topics of the AT model, the better the location accuracy. Indeed, regardless the time period, the best accuracies are reached with an AT model of size 96-98 topics and contrariwise,

the worst accuracies are observed with AT models with a small number of topics (5). This is indeed a well-known scale effect.

The approach used to automatically locate a tweet obtains very good results (more than 95% for the 9th to the 10th of January (b)). In a same way, we notice that the more precise the AT model (high number of topics), the higher the accuracy. A topic model with a thin granularity allows us to better characterize the content of a given message. Finally, we can point out that the best result is reached during the 9th to the 10th with an accuracy of 95.7% which corresponds to the second attack in Paris.

4. Conclusion and further works

In this paper, we present an efficient way to deal with short text messages from Internet micro-blogging platforms which are highly error prone. The approach seeks to map a tweet into a high-level representation using the Author-topic (AT) model that takes into consideration all information contained in a tweet: the content itself (words), the label (country) and the relation between the distribution of words in the tweet and its location, considered as a latent relation. A high-level representation allows us to obtain very promising results during the identification of the country. Experiments conducted on a Twitter corpus showed the effectiveness of the proposed AT model with an accuracy reached of more than 95%. However, these results are based on the peculiar words of significantly different languages and it does not give any semantic information on the way the Charlie event was perceived by the population from these countries. This could be an additional interesting approach to develop, especially if it is linked to a map of tweets, to observe how such a worldwide event makes a buzz in space and time.

References

1. Haewoon Kwak, Changhyun Lee, Hosung Park, and Sue Moon, "What is Twitter, a social network or a news media?," in *WWW*, 2010, p. 591–600.
2. Wayne Xin Zhao, Jing Jiang, Jianshu Weng, Jing He, Ee-Peng Lim, Hongfei Yan, and Xiaoming Li, "Comparing twitter and traditional media using topic models" in *Advances in Information Retrieval*. 2011.
3. Meeyoung Cha, Hamed Haddadi, Fabricio Benevenuto and Krishna P. Gummadi, "Measuring user influence in twitter: The million follower fallacy," in *Intern. Conference on Weblogs and Social Media (ICWSM)*, 2010.
4. Andranik Tumasjan, Timm O. Sprenger, Philipp G. Sandner, and Isabell M. Welpe, "Predicting elections with Twitter: What 140 characters reveal about political sentiment," in *ICWSM*, 2010, p. 178–185.
5. Huang Jeff, Thornton Katherine M, and Efthimiadis Efthimis N, "Conversational tagging in twitter," in *Proceedings of the 21st ACM conference on Hypertext and hypermedia. ACM*, 2010, p. 173–178.
6. Morchid M., Dufour R., and Linarès G., "A LDA-based topic classification approach from highly imperfect automatic transcriptions," in *LREC'14*, 2014.
7. Monojit Choudhury, Rahul Saraf, Vijit Jain, Sudeshna Sarkar, and Anupam Basu, "Investigation and modelling of the structure of texting language," in *IJCAI*, 2007, pp. 63–70.
8. Salton G. and Buckley C., "Term-weighting approaches in automatic text retrieval," *Information processing & management*, vol. 24, no. 5, p. 513–523, 1988.
9. Renoufand Antoinette, Baayen Harald, "Aviating among the hapax legomena: Morphological grammaticalisation in current British newspaper english," *Explorations in corpus linguistics*, no. 23, p. 181-1998.
10. Hofmann T., "Probabilistic latent semantic analysis," in *Proc. of Uncertainty in Artificial Intelligence, UAI'99*, 1999, p. 21.
11. Blei D.M., Ng A.Y., and Jordan M.I., "Latent Dirichlet Allocation," *The Journal of Machine Learning Research*, vol. 3, pp. 993–1022, 2003.
12. Rosen-Zvi M., Griffiths T., Steyvers M., and Smyth P., "The author-topic model for authors and documents," in *UAI'04*, 2004, pp. 487–494.
13. Morchid Mohamed, Richard Dufour, Mohamed Bouallegue and Georges Linarès, "Author-topic based representation of call-center conversations," in *International Spoken Language Technology Workshop (SLT)*, 2014.
14. Vladimir Vapnik and A. Lerner, "Pattern recognition using generalized portrait method," *Automation and Remote Control*, vol. 24, p. 774–780, 1963.
15. Guo-Xun Yuan, C-H Ho, and Chih-Jen Lin, "Recent advances of large-scale linear classification," *Proceedings of the IEEE*, vol. 100, no. 9, p. 2584–2603, 2012.