
Visualizing the contributions of a user in the Social Web

AUTHORS:

Maria Schmidt
Maria-Stamatoula Karavolia

UNIVERSITY:

University of Fribourg

COURSE NAME:

Web Monitoring and Analysis

SUPERVISOR:

Aleksandar Drobnjak – Research
Assistant

May 21, 2016

Contents

| | | |
|----------|---|-----------|
| 1 | Introduction | 3 |
| 2 | Motivation | 3 |
| 3 | Visualize text corpus | 5 |
| 3.1 | Literature review | 5 |
| 3.2 | Graph based visualization | 5 |
| 4 | Topic models | 6 |
| 4.1 | Literature review | 6 |
| 4.2 | Latent dirichlet allocation | 7 |
| 4.3 | Non-negative Matrix Factorization | 9 |
| 5 | Similarity between short texts | 10 |
| 5.1 | Literature review | 10 |
| 5.2 | Cosine Similarity | 10 |
| 6 | Prototype implementation | 11 |
| 6.1 | Framework | 11 |
| 6.2 | Data Collection | 11 |
| 6.2.1 | Retrieve data from Facebook | 11 |
| 6.2.2 | Retrieve data from Twitter | 11 |
| 6.3 | Data visualization | 12 |
| 6.3.1 | Contribution over time | 12 |
| 6.3.2 | Summary of posts | 12 |
| 6.3.3 | Topic models | 12 |
| 6.3.4 | Impact of posts | 12 |
| 7 | Conclusion | 13 |

List of Figures

| | | |
|---|--|---|
| 1 | An LDA example. | 8 |
| 2 | Illustration of approximate non-negative matrix factorization: the matrix A is represented by the two smaller matrices W and H , which, when multiplied, approximately reconstruct A | 9 |

1 Introduction

In the recent decades, the notion of social networks have attracted the interest of researchers and the curiosity of the social behavioral sciences, since there is a wide spread usage of the internet by people all over the world that interact with each other and exchange web content through numerous online communities such as Facebook¹, Twitter² and many other. An interesting aspect in social networks is the visualization of the contributions of a user posts and their impact across different platforms. Over time the contributions of an user in social networks is growing and this creates the need to have a simple overview of the data by visualizing it. Data visualization offers a quick way to present the data in a way that can reveal valuable hidden insights. Thus, through the visualization, users can easily understand what are the hot posts, that is those posts are able to attract a greater attention or interest.

2 Motivation

klgdfkgdkfgdkljgkldf

¹<http://www.facebook.com>

²<http://www.twitter.com>

contributons visualization

3 Visualize text corpus

3.1 Literature review

3.2 Graph based visualization

4 Topic models

As our collective knowledge continues to be digitized and stored—in the form of news, blogs, web pages, scientific articles, books, images, sound, video, and social networks—it becomes more difficult to find and discover what we are looking for. We need new computational tools to help organize, search and understand these vast amounts of information. To this end, machine learning provides topic models which are a suite of algorithms for discovering themes (or topics) that spread through a collection of documents.

In this section, we provide an overview of existing literature on topics models and then we describe two types of existing topic models: (i) the probabilistic model: Latent Dirichlet allocation and (ii) the non-probabilistic model: Non-negative Matrix Factorization.

4.1 Literature review

Topic modeling is gaining increasingly attention and is applied in a wide range of areas including on social networks such as Twitter and Facebook. From the view of methodology, topic models are separated into two groups: the non-probabilistic and probabilistic approaches.

Most of probabilistic approaches are based on Latent Dirichlet Allocation (LDA) [2], which is the most popular standard tool in topic modeling. As a result, LDA has been used and extended in a variety of ways, and in particular for social networks and social media, so a great number of research papers that deal with LDA have been proposed.

Ramage et al. [13, 14] extended LDA to a supervised form and studied its application in micro-blogging environment. More particularly, in [14] the Labeled LDA, a novel model of multi-labeled corpora that directly addresses the credit assignment problem is introduced. In [13], a scalable implementation of a partially supervised learning model (Labeled LDA) is proposed for discovering topics in microblogs like Twitter. This model maps the content of the Twitter feed into dimensions such as substance, style, status, and social characteristics of posts. Thus, this approach helps to efficiently characterize selected Twitter users along these learned dimensions and indicates that topic models can provide interpretable summaries of users' tweet posts. Phan et al. [12] studied the problem of modeling short text through LDA. In particular, a general framework, based on LDA, for building classifiers with hidden topics discovered from large-scale data collections that can deal successfully with short and sparse text Web segments.

Moreover, Chang et al. [3] proposed a novel probabilistic topic model to analyze text corpora and infer descriptions of the entities and of relationships between those entities on Wikipedia. McCallum et al. [8] proposed a probabilistic generative model to simultaneously discover groups among the entities and topics among the corresponding text. Zhang et al. [18] introduced a model to incorporate LDA into a community detection process.

More specifically, in this paper they designed a hierarchical Bayesian network based approach, namely GWNLDA(Generic-Weighted Network-LDA) which is inspired by LDA for discover probabilistic communities from complex networks. Similar work can be found in [7] and [10].

Standard LDA is often less coherent when applied to microblog content like Twitter because tweets are short. To overcome this difficulty, some previous studies proposed to aggregate all the tweets of a user as a single document. In [1], a topical classification of Twitter users and messages is provided. This paper deals with the problem of using topic models in microblogs by proposing schemes based on LDA and one extended model based on the Author-Topic model. It also presents that topic models aims some classification problems by indicating that topic models learned from aggregated messages by the same user obtaining higher accuracy. Also, a different approach on topic modeling of Tweets is provided in [9]. This paper focus on how to improve clustering metrics and topic coherence with existing algorithms. More specifically, it provides two novel schemes that lead to significantly improved LDA topic models on Twitter content without requiring any modification of the underlying LDA machinery. The first one is about pooling tweets by hashtags that yields a great improvement in all metrics for topic coherence across three diverse Twitter datasets, and the second is an automatic hashtag assignment scheme further improves the hashtag pooling results on a subset of metrics.

Non-probabilistic topic models are also very popular. One of the most known representative model is the Non-Negative Matrix Factorization (NMF) [11, 6]. Yan et.al [17], proposed a novel term weight called Ncut-weighted, which measures term's discriminability according to the words cooccurrences, for short text clustering. More particularly, the experiments show that the clustering performance of NMF is greatly improved with terms weighted by the Ncut-weight. Due to the severe sparsity of short texts, in[16], a different approach on the non-negative matrix factorization is introduced. This approach first learns topics from term correlation data using symmetric non-negative matrix factorization, and then infers the topics of documents. The experimental results on three short text data sets show that this method provides substantially better performance than other baseline methods like LDA.

4.2 Latent dirichlet allocation

The idea behind Latent Dirichlet allocation (LDA) [2, 4, 15], which is an unsupervised machine learning technique, is to model documents as arising from multiple topics, where a topic is defined to be a distribution over a fixed vocabulary of terms. Specifically, we assume that K topics are associated with a collection, and that each document exhibits these topics with different proportions. The interaction between the observed documents and hidden topic structure is manifest in the probabilistic generative process associated

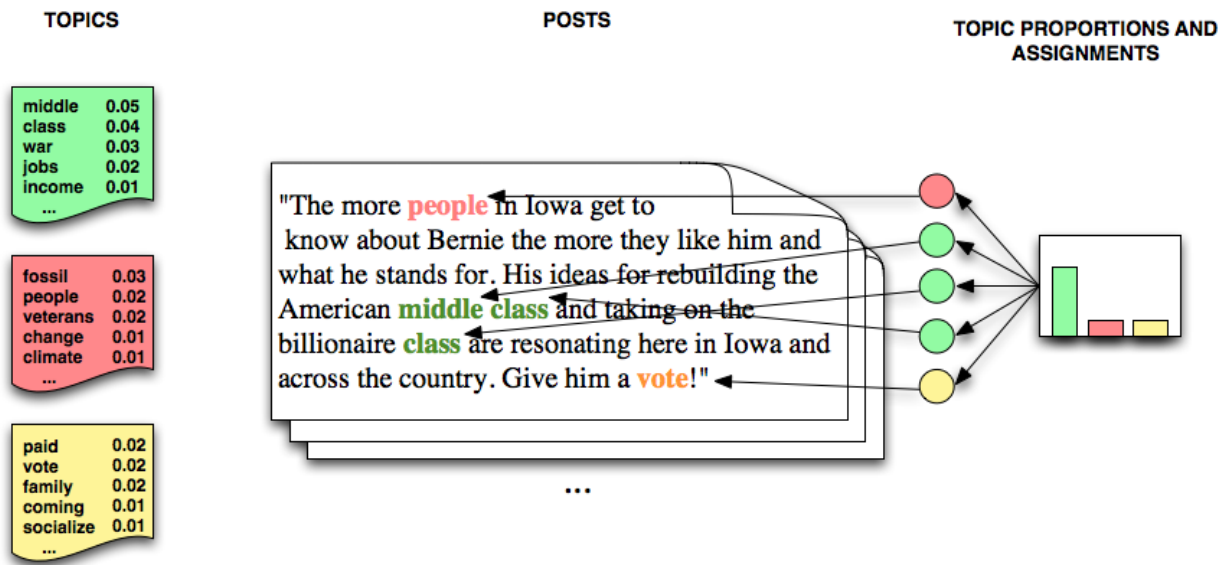


Figure 1: An LDA example.

with LDA. This generative process is as follows:

To generate a document:

1. Randomly choose a distribution over topics.
2. For each word in the document
 - (a) Randomly choose a topic from the distribution over topics in step #1.
 - (b) Randomly choose a word from the corresponding distribution over the vocabulary

So, this statistical model reflects the intuition that documents exhibit multiple topics. Each document exhibits the topics with different proportion (step #1); each word in each document is drawn from one of the topics (step #2b), where the selected topic is chosen from the per-document distribution over topics (step #2a).

In Figure 1, an LDA example is illustrated. We assume that we have three topics, which are distributions over words, exist for the whole collection (far left). Each document is assumed to be generated as follows. First choose a distribution over the topics (the histogram at right) and then, for each word, choose a topic assignment (the colored coins) and choose the word from the corresponding topic. As we can see, in this example the particular document is assigned to the first topic with higher probability compare to the others.

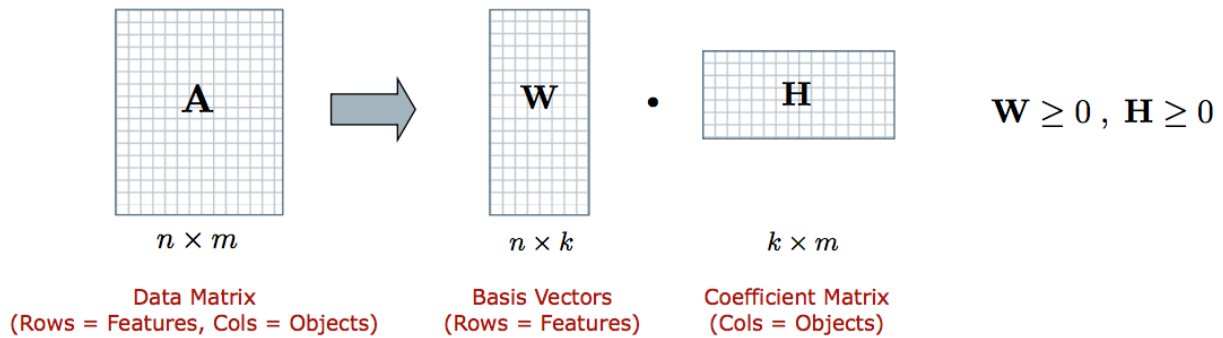


Figure 2: Illustration of approximate non-negative matrix factorization: the matrix \mathbf{A} is represented by the two smaller matrices \mathbf{W} and \mathbf{H} , which, when multiplied, approximately reconstruct \mathbf{A}

4.3 Non-negative Matrix Factorization

Non-negative matrix factorization (NMF) is an unsupervised family of algorithms from linear algebra that simultaneously perform dimension reduction and clustering. NMF was first introduced by Paatero and Tapper [11] as positive matrix factorization and subsequently popularized by Lee and Seung [6].

In Figure 2, NMF takes a non-negative matrix \mathbf{A} as an input, and factorizes it into two smaller non-negative matrices \mathbf{W} and \mathbf{H} , each having k dimensions. When multiplied together, these factors approximate the original matrix \mathbf{A} . The specified parameter k controls the number of topics that will be produced. The rows of the matrix \mathbf{W} provides weights that indicate the strength of association between documents and topics. The columns of the \mathbf{H} that indicate the strength of association between terms and topics. By ordering the values in a given column and selecting the top-ranked terms, we can produce a description of the corresponding topic.

5 Similarity between short texts

5.1 Literature review

5.2 Cosine Similarity

For creating the user timeline, we need to find similar posts that a user share both in Twitter and Facebook. For this purpose, we calculate the cosine similarity metric. The cosine similarity [5] between two vectors (or two documents on the Vector Space) is a measure that calculates the cosine of the angle between them. This metric is a measurement of orientation and not magnitude, it can be seen as a comparison between documents on a normalized space because we're not taking into the consideration only the magnitude of each word count (tf-idf) of each document, but the angle between the documents.

Given two posts t_a and t_b , their cosine similarity is

$$\cos(\mathbf{t_a}, \mathbf{t_b}) = \frac{\mathbf{t_a} \mathbf{t_b}}{\|\mathbf{t_a}\| \|\mathbf{t_b}\|} \quad (1)$$

where t_a and t_b are m -dimensional vectors over the term set $T = t_1, \dots, t_m$. Each dimension represents a term with its weight in the document, which is non-negative. As a result, the cosine similarity is non-negative and bounded between $[0, 1]$.

We predefine a threshold to accept two similar posts to have similarity at least 60%. This portion lets us a great amount of similar posts and also recognize twitter posts that have urls, hashtags and mentions.

6 Prototype implementation

In this section we provide a detailed description of our framework, the data collection and the technologies that used in the front and back ends.

6.1 Framework

Our web application is based on Django³, which is a free and open-source web framework, written in Python, and follows the model–view–controller (MVC) architectural pattern.

6.2 Data Collection

Our data collection consists of real-world data from Facebook and Twitter and focus on 28 public persons such as, politicians and athletes because they tend to post the same content on Twitter and Facebook more than a normal user.

6.2.1 Retrieve data from Facebook

6.2.2 Retrieve data from Twitter

The data was obtained by quering the timeline API of Twitter with the username of each person related to politicians and athletes. For this procedure we used the Tweepy⁴, which is a Python library for accessing the Twitter API. We were able to collect a fixed number of tweets because Twitter only allows access to a users most recent 3240 tweets. The attributes of the data along with their definitions are displayed in Table 1.

Table 1: Description of the attributes of Twitter data

| Attribute | Description |
|-----------|------------------------------------|
| ID | The id of the twitter post |
| date | The date when the tweet was posted |
| text | The text of the tweet |
| likes | The number of likes of a tweet |

³<https://www.djangoproject.com/>

⁴<http://www.tweepy.org/>

6.3 Data visualization

6.3.1 Contribution over time

6.3.2 Summary of posts

6.3.3 Topic models

6.3.4 Impact of posts

7 Conclusion

References

- [1]
- [2] D. M. Blei, A. Y. Ng, and M. I. Jordan. Latent dirichlet allocation. *J. Mach. Learn. Res.*, 3:993–1022, Mar. 2003.
- [3] J. Chang, J. L. Boyd-Graber, and D. M. Blei. Connections between the lines: augmenting social networks with text. In J. F. E. IV, F. Fogelman-Soulié, P. A. Flach, and M. Zaki, editors, *KDD*, pages 169–178. ACM, 2009.
- [4] G. Heinrich. Parameter estimation for text analysis. *Web: <http://www.arbylon.net/publications/text-est.pdf>*, 2005.
- [5] A. Huang. Similarity measures for text document clustering. pages 49–56, 2008.
- [6] D. D. Lee and H. S. Seung. Learning the parts of objects by nonnegative matrix factorization. *Nature*, 401:788–791, 1999.
- [7] Y. Liu, A. Niculescu-Mizil, and W. Gryc. Topic-link lda: Joint models of topic and author community. In *Proceedings of the 26th Annual International Conference on Machine Learning*, ICML ’09, pages 665–672, New York, NY, USA, 2009. ACM.
- [8] A. McCallum, X. Wang, and N. Mohanty. Joint group and topic discovery from relations and text. In *Proceedings of the 2006 Conference on Statistical Network Analysis*, ICML’06, pages 28–44, Berlin, Heidelberg, 2007. Springer-Verlag.
- [9] R. Mehrotra, S. Sanner, W. Buntine, and L. Xie. Improving lda topic models for microblogs via tweet pooling and automatic labeling. In *Proceedings of the 36th international ACM SIGIR conference on Research and development in information retrieval*, pages 889–892. ACM, 2013.
- [10] R. M. Nallapati, A. Ahmed, E. P. Xing, and W. W. Cohen. Joint latent topic models for text and citations. In *Proceedings of the 14th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*, KDD ’08, pages 542–550, New York, NY, USA, 2008. ACM.
- [11] P. Paatero and U. Tapper. Positive matrix factorization: A non-negative factor model with optimal utilization of error estimates of data values. *Environmetrics*, 5(2):111–126, 1994.
- [12] X.-H. Phan, L.-M. Nguyen, and S. Horiguchi. Learning to classify short and sparse text & web with hidden topics from large-scale data collections. In *Proceedings of the 17th International Conference on World Wide Web*, WWW ’08, pages 91–100, New York, NY, USA, 2008. ACM.

-
- [13] D. Ramage, S. Dumais, and D. Liebling. Characterizing microblogs with topic models. In *Proc. ICWSM 2010*. American Association for Artificial Intelligence, May 2010.
 - [14] D. Ramage, D. Hall, R. Nallapati, and C. D. Manning. Labeled lda: A supervised topic model for credit attribution in multi-labeled corpora. In *Proceedings of the 2009 Conference on Empirical Methods in Natural Language Processing: Volume 1 - Volume 1*, EMNLP '09, pages 248–256, Stroudsburg, PA, USA, 2009. Association for Computational Linguistics.
 - [15] M. Steyvers and T. Griffiths. *Probabilistic topic models*, volume 427. Handbook of Latent Semantic Analysis, 2007.
 - [16] X. Yan, J. Guo, S. Liu, X. Cheng, and Y. Wang. Learning topics in short texts by non-negative matrix factorization on term correlation matrix. In *Proceedings of the SIAM International Conference on Data Mining*, 2013.
 - [17] X. Yan, J. Guo, S. Liu, X.-q. Cheng, and Y. Wang. Clustering short text using ncut-weighted non-negative matrix factorization. In *Proceedings of the 21st ACM International Conference on Information and Knowledge Management*, CIKM '12, pages 2259–2262, New York, NY, USA, 2012. ACM.
 - [18] H. Zhang, C. L. Giles, H. C. Foley, and J. Yen. Probabilistic community discovery using hierarchical latent gaussian mixture model. In *Proceedings of the 22Nd National Conference on Artificial Intelligence - Volume 1*, AAAI'07, pages 663–668. AAAI Press, 2007.