

# Finding Hierarchy of Topics from Twitter Data

Nghia Duong-Trung, Nicolas Schilling, and Lars Schmidt-Thieme

Information Systems and Machine Learning Lab (ISMLL)  
Universitätsplatz 1, 31141 Hildesheim, Germany  
{duongn,schilling,schmidt-thieme}@ismll.uni-hildesheim.de  
<http://www.ismll.uni-hildesheim.de>

**Abstract.** Topic modeling of text collections is rapidly gaining importance for a wide variety of applications including information retrieval and automatic multimedia indexing. Our motivation is to exploit a hierarchical topic selection via nonnegative matrix factorization to capture the nature content of text posted on Twitter. This paper explores the use of an effective framework to automatically discover hidden topics and their sub-topics. As input, the framework uses textual data. The output is then the discovered structure of topics. We introduce a conceptual topic modeling based on the idea of stability analysis to detect a hierarchy of topics given a text source. In this process, we apply stability measurement in conjunction with nonnegative matrix factorization and WordNet to excavate hidden topics by the scores of conceptual similarity. To demonstrate the effectiveness and generalization, we apply the approach to a large-scale Twitter dataset to investigate the content topics. We also address the problems of several state-of-the-art topic modeling approaches that are unable to handle a large dataset.

**Keywords:** Unsupervised Learning, Semantics in Text Mining, Conceptual Stability, Hierarchy of Topics

## 1 Introduction

Nonnegative matrix factorization (NMF) with nonnegativity constraints has been considered as an efficient representation and an emerged technique for text mining and document clustering [22,2,17,23,9,11]. For any desired low-rank  $K$ , the NMF algorithm groups the data into clusters. The key issue is whether a given low-rank  $K$  helps to decompose the data into appropriate separated clusters. Therefore, the problem we study in this paper is how to effectively and efficiently discover the most appropriate structure of topics giving a text corpus by exploiting the semantic meaning and the conceptual stability. In general, the stability of a clustering model refers to its ability to consistently replicate similar solutions on data randomly generated from the same source. In practice, this involves a repeated re-sampling of data, applying a topic selection model, and evaluating the results by a stability metric which measures the level of discrimination between the resulting clusterings.

We start with the previous work on the idea of random sub-sampling and stability analysis via consensus clustering to discover the number of clusters that best describes the data [16,4,13]. The basic assumption of stability in the context of consensus clustering, in general, is very intuitive: for particular observed data, if we perturb it into different random variabilities, and if they produce the same cluster composition, or *consensus*, without radical difference, we would confidently consider that these clusters represent real structure. Consensus clustering purely captures this procedure. Further work investigated by [4] improved the consensus clustering technique by adding a quantitative evaluation for robustness of the decomposition. They adopted a measure based on the cophenetic correlation coefficient which indicates the dispersion of the consensus matrix. The coefficient is calculated as the Pearson correlation of two distance matrices: the consensus matrix captured the distance between data samples and the average connectivity matrix over many clustering runs. Subsequently, [12,13] formulate the idea of consensus matrix in the latent space learned by NMF.

However, the computation of consensus matrix,  $R^{n \times n}$  matrix where  $n$  is the number of tweets/documents, seems very costly, e.g. large amount of RAM is required. For instance, if we apply the previous method on our experimented Twitter dataset that we describe later in the paper, then 1400GB of RAM is required to store the consensus matrix during model’s computation. Hence, the method provided by [12,13] is insufficient or even impossible for large datasets. To overcome the drawbacks of the construction of consensus matrix, we propose a topic selection approach, called the conceptual stability analysis, to smoothly integrate with NMF that can be applied on large datasets effectively.

Moreover, we also evaluate several state-of-the-art topic modeling approaches via Latent Dirichlet Allocation (LDA) [3]. The first baseline is the topic selection method implemented by [1]. The second baseline is proposed by [6]. We implement the baselines by using [10,18]. However, these methods threw exception due to large dataset during computation. An upper bound of RAM required for each approach is 65GB until an exception occurs.

With these limitation in mind, we introduce an unsupervised topic selection method that enhances the accuracy and effectiveness of NMF-based models in the context of document clustering and topic modeling. We show that our proposed method can work effectively on large dataset within acceptable computing resources such as RAM required and time of computation.

## 2 Theoretical Aspects and Proposed Framework

### 2.1 Nonnegative Matrix Factorization

Consider a dataset  $X \in \mathbb{R}^{n \times m}$  containing a set of  $n$  documents where each document is described by  $m$  many features. The document features are mapped from a dictionary that comprises all words/terms/tokens in the dataset. Each positive entry  $X_{ij}$  is either a raw term frequency or a term frequency - inverse document frequency (TFIDF) score. By  $r$  and  $\tau$ , we denote the sampling rate

and the number of subsets generated from  $X$  respectively. Then each subset  $X_\tau \in \mathbb{R}^{n' \times m}$  is a sample without replacement of  $X$ .

Giving a desired number of topics  $k$ , the NMF algorithm iteratively computes an approximation:

$$X \approx WH, \quad (1)$$

where  $W \in \mathbb{R}^{n \times k}$  and  $H \in \mathbb{R}^{k \times m}$  are nonnegative matrices. The conventional technique to approximate  $W$  and  $H$  is by minimizing the difference between  $X$  and  $WH$  such that:

$$\min_{W \geq 0, H \geq 0} f(W, H) = \frac{1}{2} \sum_{i=1}^n \sum_{j=1}^m (X_{ij} - (WH)_{ij})^2 + \phi(W) + \theta(H), \quad (2)$$

where  $\phi(\cdot)$  and  $\theta(\cdot)$  are regularization terms that are set as follows:

$$\phi(W) = \alpha \|W\|_F^2 \quad \text{and} \quad \theta(H) = \beta \sum_{i=1}^m \|H(:, i)\|_1^2, \quad (3)$$

where  $H(:, i)$  indicates the  $i$ -th column of  $H$ . The  $L_1$  norm term of  $\theta(H)$  promotes sparsity on the rows of  $H$  while the Frobenius norm term of  $\phi(W)$  prevents  $W$  from growing too large. Scalar parameters  $\alpha$  and  $\beta$  are used to control the strength of regularization. The matrices  $W$  and  $H$  are found by minimizing Equation (2) via estimating  $W$  and  $H$  in an alternating fashion using projected gradients or coordinate descent [5].

**Table 1.** Summary of topics discovered.

N	Topic	Number of sub-topics	Number of documents	Share (%)
1	Student Life and Relationship	5	76,415	17.78
2	Information and Networking	8	17,649	4.11
3	Business and Current Affairs	2	72,534	16.88
4	Routine Activities	4	21,960	5.11
5	Leisure and Entertainment	2	81,469	18.96
6	Sport and Games	3	31,812	7.40
7	Pessimism and Negativity	4	78,329	18.23
8	Wishes and Gratitude	5	36,618	8.52
9	Transport and Travel	2	12,887	3.01
	Total	35	429,673	100.00

## 2.2 Conceptual Stability Computation

Now we start discussing our approach of computing stability based on the usage of the *WordNet* hypernym hierarchy [8,15]. Given tokens  $c_p$  and  $c_q$ , then

**Table 2.** Labels of all 9 topics and 35 sub-topics.

1	Student Life and Relationship	2	Information and Networking	4	Routine Activities	7	Pessimism and Negativity
1.1	Friends and Relationship	2.1	News	4.1	Feelings	7.1	Hate and Anger
1.2	Study Life	2.2	Life	4.2	Sleep	7.2	Daily Problems and Complains
1.3	Worry and Confusion	2.3	Mood and Reflections	4.3	Work and People	7.3	School Routines
1.4	Conversations	2.4	Greetings	4.4	Social Media	7.4	Life and Changes
1.5	Social Media and Connections	2.5	Current Regional Events	5	Leisure and Entertainment	9	Transport and Traveling
8	Wishes and Gratitude	2.6	Mates	5.1	Television and Cinema	9.1	Landmarks
8.1	Friends	2.7	Informal Chat	5.2	Reactions	9.2	Journeys
8.2	People	2.8	Religion	6	Sport and Games		
8.3	Anticipation	3	Business and Current Affairs	6.1	Sport Opinions and Discussion		
8.4	Thanks and Affection	3.1	Working-day Activities	6.2	American Football		
8.5	Celebrations	3.2	Events and Socializing	6.3	TV Sport Programs		

$wup(c_p, c_q)$  is the Wu-Palmer similarity [21], which is a scoring method based on how similar the token senses are and where they occur relative to each other in the WordNet hierarchy. Then, the Wu-Palmer similarity is calculated by:

$$wup(c_p, c_q) = \frac{2d}{d_1 + d_2 + 2d} \quad (4)$$

where  $d_1$  and  $d_2$  are the distances that separates the concept  $c_p$  and  $c_q$  from their closest common ancestor and  $d$  is the distance which separates the closest common ancestor of  $c_p$  and  $c_q$  from the root node.

Each row of the low-rank matrix  $H$  represents one of the  $k$  topics and consists of scores for each term. However, we only consider the top  $t \ll m$  terms as they contribute most to the semantic meaning of a topic. In practice, the contribution of each token to topic  $i$  is represented by the scores in the  $i$ -th row in matrix  $H$  generated by NMF. By sorting each row of  $H$ , we can assess the top  $t$  terms for each topic. The set of top  $t$  tokens for all topics of a given  $H$  will be denoted by  $\mathcal{S} = \{R_1, \dots, R_k\}$  such that  $R_i \in \mathbb{R}^t$  is the topic  $i$ -th represented by top  $t$  tokens. Within a topic, we calculate the conceptual stability score as follows:

$$\text{sim}(R_v) = \frac{2}{t(t-1)} \sum_{i=0}^{t-1} \sum_{j=i+1}^t wup(R_{v_i}, R_{v_j}) \quad (5)$$

Similarly, the conceptual stability score between two topics  $R_u$  and  $R_v$  is calculated in the same fashion.

$$\text{sim}(R_u, R_v) = \frac{1}{t^2} \sum_{i=1}^t \sum_{j=1}^t wup(R_{u_i}, R_{v_j}) \quad (6)$$

Finally, we consider the problem of measuring the conceptual stability between two different  $K$ -way topic clusterings  $S_w$  and  $S_l$ . Each ranked list contains top  $t$  tokens that contribute most semantic meaning to the  $i$ -th topic. Then, the conceptual stability between  $S_w$  and  $S_l$  is calculated by:

$$\text{con}(S_w, S_l) = \frac{1}{K} \sum_{k=1}^K \text{sim}(R_{wk}, \pi(R_{lk})), \quad (7)$$

where  $\pi(R_{wi})$  denotes the ranked list  $R_{lj}$  matched to the ranked list  $R_{wi}$  by the permutation  $\pi$ . The optimal permutation  $\pi$  is found by solving the minimal weight bipartite matching problem using the Hungarian method [14].

Moreover, the problem of measuring the conceptual stability within the  $K$ -way topic clustering  $S_w$  itself is also considered. The conceptual stability is then calculated as follows:

$$\text{con}(S_w) = \frac{1}{K} \sum_{k=1}^K \text{sim}(R_{wk}) \quad (8)$$

We now consider the conceptual stability at a particular number of topics  $k$ . At first we apply the NMF on the complete dataset  $X$  to get the factor matrices  $H$  that we consider as the reference ranked lists. Let us define  $S_X$  as the reference  $K$ -way topic clustering, containing  $K$  ranked lists  $S_X = \{R_{X1}, \dots, R_{Xk}\}$ .

Subsequently, we randomly resample  $\tau$  times the documents of the original  $X$  with the sampling rate  $r$  to obtain a random subset of  $X$  which we denote by  $X_\tau$ . We then apply NMF on each  $X_\tau$  to get the factor matrix  $H_\tau$ . This results in  $\tau$  many sets  $\{S_1, \dots, S_\tau\}$  where each set contains  $k$  ranked lists  $S_j = \{R_{j1}, \dots, R_{jk}\}$ . Finally, we calculate the overall semantically conceptual stability at  $k$  as following:

$$\text{stability}(k) = \frac{1}{\tau} \frac{|\sum_{i=1}^{\tau} \text{con}(S_X, S_i) - \sum_{i=1}^{\tau} \text{con}(S_i)|}{\max(\sum_{i=1}^{\tau} \text{con}(S_X, S_i), \sum_{i=1}^{\tau} \text{con}(S_i))} \quad (9)$$

The maximum stability score is achieved if and only if the top  $t$  tokens appear in only one topic  $k$ . Otherwise, the minimum stability score, is obtained if top  $t$  tokens overpoweringly appear in every topic  $k$ .

This process is repeated for a range of topics  $k$ . The most appropriate value of  $k$  is identified by the highest value of  $\text{stability}(k)$  score. However, the scores also reveal the possible range of  $k$  for further investigation. With the  $k$  topic classification finalized at the first level, the dataset is split into sub-datasets where documents are assigned to topic with the highest score, e.g. through the  $W$  matrix.

$$\hat{k}_{X_i} = \underset{k}{\text{argmax}}(W_{ik}) \quad (10)$$

Then, the process is repeated to discover sub-topics in each sub-dataset. Generally, we can expand the procedure deeper in the hierarchy. First, we calculate the most appropriate number of topics  $L$  in the whole dataset  $X$ . Then, a subset of  $X$  is drawn based on each value in range  $k \in l, \dots, L$ . The  $\text{stability}(k)$  score is, in turn, calculated for each subset to find the best number of sub-topics.

---

**Algorithm 1** The conceptual stability analysis approach with 2-level of hierarchy

---

**Input:** Dataset  $X \in \mathbb{R}^{n \times m}$ , range of number of topics  $[K', \dots, K'']$ , number of top tokens  $t$ , sampling rate  $r$ , number of subsets  $\tau$

```

1: /* find  $k$  at the first level in hierarchy */
2: for  $k \in K', \dots, K''$  do
3:   find  $W \in \mathbb{R}^{n \times k}, H \in \mathbb{R}^{k \times m}$  with  $X \approx WH$ 
4:   get  $S_X \in \mathbb{R}^{k \times t}$  from  $H$ 
5:   for  $\tau \in 1, \dots, \tau$  do
6:     draw  $X_\tau \in \mathbb{R}^{n' \times m}$  from  $X$ 
7:     find  $W_\tau \in \mathbb{R}^{n' \times k}, H_\tau \in \mathbb{R}^{k \times m}$  with  $X_\tau \approx W_\tau H_\tau$ 
8:     get  $S_{X_\tau} \in \mathbb{R}^{k \times t}$  from  $H_\tau$ 
9:   end for
10:  calculate  $stability(k)$ , Equation (9)
11: end for
12:  $L = \underset{k}{\operatorname{argmax}} stability(k)$ 
13: /* find  $k$  at the second level in hierarchy */
14: for  $h \in l, \dots, L$  do
15:    $X^h = \emptyset$ 
16:   for  $X_i$  s.t.  $h == \underset{h}{\operatorname{argmax}}(W_{ih})$  do
17:      $X^h = X^h \cup \{X_i\}$ 
18:   end for
19:   for  $k \in K', \dots, K''$  do
20:     find  $W^h \in \mathbb{R}^{p \times k}, H^h \in \mathbb{R}^{k \times m}$  with  $X^h \approx W^h H^h$ 
21:     get  $S_{X^h}^h \in \mathbb{R}^{k \times t}$  from  $H^h$ 
22:     for  $\tau \in 1, \dots, \tau$  do
23:       draw  $X_\tau^h \in \mathbb{R}^{p' \times m}$  from  $X^h$ 
24:       find  $W_\tau^h \in \mathbb{R}^{p' \times k}, H_\tau^h \in \mathbb{R}^{k \times m}$  with  $X_\tau^h \approx W_\tau^h H_\tau^h$ 
25:       get  $S_{X_\tau^h}^h \in \mathbb{R}^{k \times t}$  from  $H_\tau^h$ 
26:     end for
27:     calculate  $stability(k)$  as Equation (9)
28:   end for
29:    $L_h = \underset{k}{\operatorname{argmax}} stability(k)$ 
30: end for

```

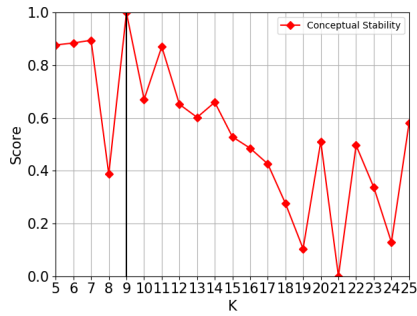
---

For the ease of interpretation, we conduct the experiments within 2-level of hierarchy. An overview of the whole procedure can be seen in Algorithm (1).

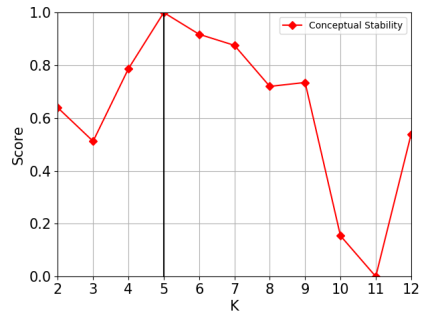
### 3 Empirical Results

#### 3.1 Datasets, Experiment Setup, and Baselines

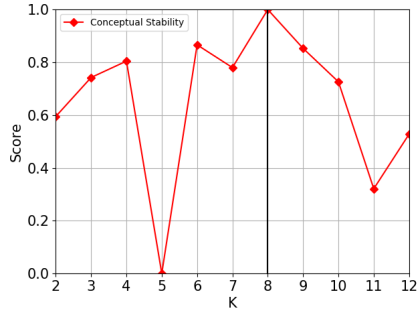
The North America dataset is a large dataset of tweets that was originally used for the geolocation prediction problem [19,20,7]. A document in this dataset is



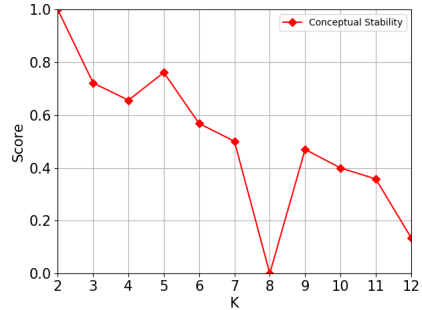
(a) Experiment result on the tweets dataset at the first level of hierarchy



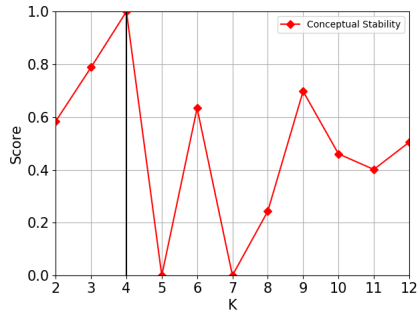
(b) Student Life and Relationship



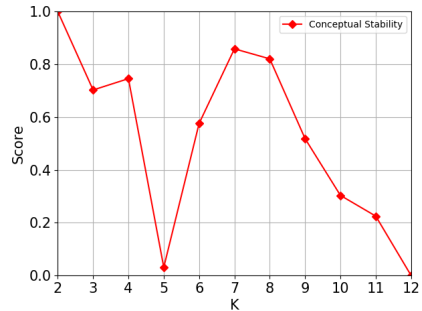
(c) Information and Networking



(d) Business and Current Affairs



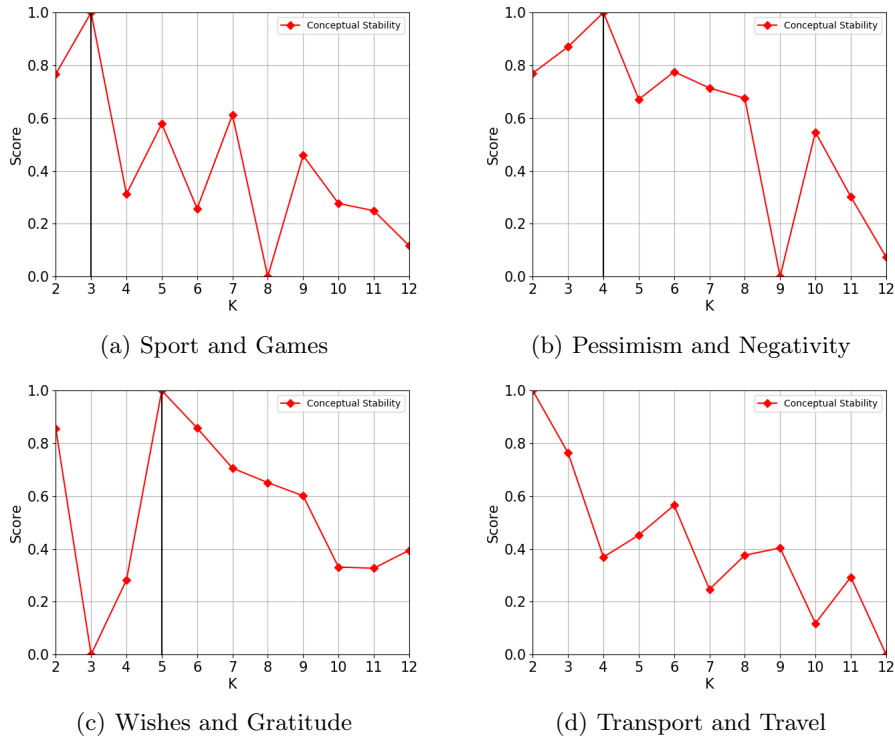
(e) Routine Activities



(f) Leisure and Entertainment

**Fig. 1.** Experiment results on the Tweets dataset. Figure 1(a) shows discovered topics at the first level. Similarly, the other figures present discovered topics at the second level. The appropriate number of topics  $k$  are identified by **peaks** in the plots. The vertical lines represent the highest peaks  $k$ .

the concatenation of all tweets by a single user. There were total 38 million tweets tweeted by 430k users. The tweets were inside a bounding box covering the continuous United States, a part of Canada and a part of Mexico. The final



**Fig. 2.** Experiment results on discovered topics at the second level. The appropriate number of topics  $k$  are identified by **peaks** in the plots. The vertical lines represent the highest peaks  $k$ .

dataset after preprocessing is a very sparse text file that requires 2.2GB to store and contains 430,000 rows, e.g. the number of documents, and 59,368 columns, e.g. the vocabulary size.

In our experiment, we set required model parameters as follows. The range of exploring topics at the first level is  $\{5, \dots, 25\}$ . We expect the range of sub-topics is smaller in the second level so that the range of exploring sub-topics is  $\{2, \dots, 12\}$ . The number of top tokens that characterize a specific topic is set to  $t = 20$ . The sampling rate is set to  $r = 0.8$  and the number of subsets is set to  $\tau = 25$  to cancel out random effects. Our experiments were conducted on a Xeon E5-2670v2 with 2.5GHz clock speed and 128GB of RAM. However, an upper bound of RAM required for our model is 5GB and it takes 4 days to complete.

As we already mentioned in the introduction section, during our experiment, we also compare our method with several state-of-the-art NMF-based and LDA-based topic modeling approaches [13,1,6]. However, all these models either cannot handle a large dataset or throw resource exception during computation.



### 3.2 Topics Discovery

The framework identifies distinctive topics and their sub-topics of documents based on the output of stability scores. In theory, the deepest hierarchy of topics is where documents are recursively classified until one topic only contains one document. We do not specify the exact number of topics beforehand but rather the range of desired topics and the model will figure out the most appropriate values itself. In other words, the model takes (1) a very large textual dataset, (2) a desired range of expectant number of topics, and (3) a desired level of hierarchy. Then, the hierarchy of topics is discovered by considering conceptual stability scores.

Figure (1,2) present the potential number of topics and sub-topics at the first and second levels respectively. Table (1) summarizes topics and their sub-topics explored. As we can see in Table (1), topics at the first level can be divided into two groups based on the % share. People are concerned the most about *Pessimism and Negativity*, *Leisure and Entertainment*, *Student Life and Relationship* and *Business and Current Affairs*. We now describe all the topics and their sub-topics discovery in more detail.

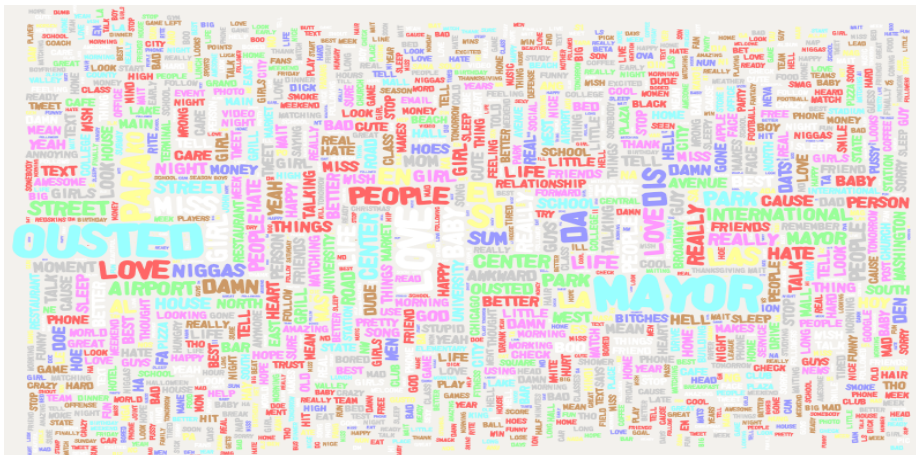
At the first level, the highest peak is found at  $k = 9$  which means that the most distinctive number of topics given North America tweet dataset is 9. However, we also see potential high peaks at  $k = 11$  and  $k = 7$  if we need manually to expand or condense the clustering results respectively. Consequently, the whole dataset at the first level is then divided into 9 sub-datasets that the model continues discovering sub-topics within them.

Next we consider sub-topics. Figure (1(b)) presents that the highest peak is at  $k = 5$  where we clearly see a  $\Lambda$  shape. Similarly, we see the same  $\Lambda$  shape in the 4<sup>th</sup>, 6<sup>th</sup>, 7<sup>th</sup> and 8<sup>th</sup> topics which are presented in Figure (1(e), 2(a), 2(b) and 2(c)) respectively. The peaks made by the  $\Lambda$  shape is the number of sub-topics discovered by the model. Interestingly, the 2<sup>nd</sup> topic, Figure (1(c)), contains only 4.11% of the documents but can be divided into  $k = 8$  distinctive sub-topics. The 3<sup>rd</sup> and 9<sup>th</sup> topics, Figure (1(d), 2(d)) respectively, show an obvious peak that the most suitable number of sub-topics is  $k = 2$ , the left most bound of the experimented range. The 5<sup>th</sup> topic, Figure (1(f)), presents two candidates with high magnitude peaks at  $k = 2$  and  $k = 7$ . Although the highest peak is selected, e.g.  $k = 2$ , as the output for sub-topics consideration, user can manually choose the other peak as the desired output.

### 3.3 Topics Labeling

Having exploited the hierarchical topics structure, we next present their associated labels. Table (2) summarize our labeling schemes. All topics and sub-topics were subjectively labeled to ease the understanding and interpretation in successive spatial distribution analysis. The labels were validated and assigned based on the meaning of top tokens that characterize a specific topic or sub-topic.

More generally, questions of accuracy can be raised about the representativeness of labels as a source for topics demonstration. In each discovered topic

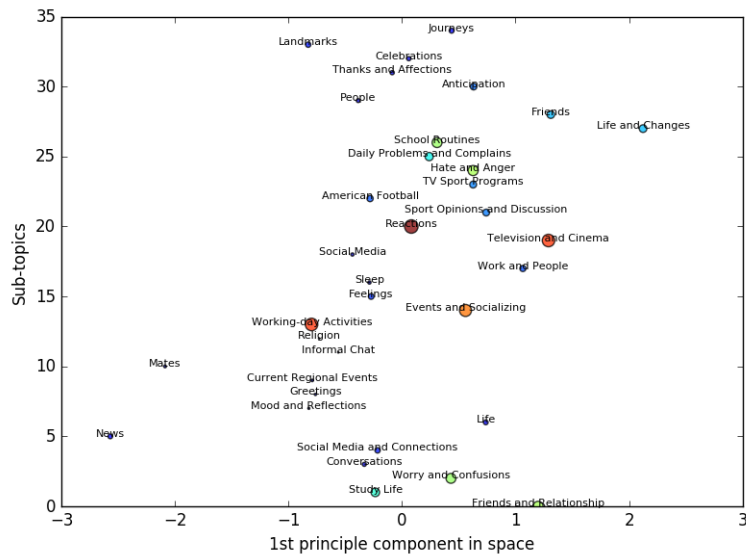


**Fig. 3.** Word cloud of 9 discovered topics. The size of the word depends on the score calculated by the topic model. Offensive language has been removed from the word cloud. The colors for both circles and words are red, yellow, green, blue, purple, brown, gray, white and turquoise for *Student Life and Relationship*, *Information and Networking*, *Business and Current Affairs*, *Routine Activities*, *Leisure and Entertainment*, *Sport and Games*, *Pessimism and Negativity*, *Wishes and Gratitude*, and *Transport and Travel* topic respectively.

and sub-topic, after collecting top tokens based on their meaning contribution, a wide number of heuristic labeling schemes is considered to render each topic representative and distinctive. After the labels are generated, a random selected documents are reviewed and the labels are re-validated if needed. The loop is required to ensure the assigned labels are acceptably appropriate. It is important to consider that the labeling results from this paper reflect Twitter users’ opinions at the time the data was collected, not the population at large. The revealed Twitter topics also were visualized using a comparison word cloud of the top tokens in all topics and sub-topics, e.g. Figure (3). We report the principal component analysis to inspect the subjective distinctiveness of topics in Figure (4).

## 4 Conclusion

In this paper, we propose a topic selection approach to smoothly integrate with NMF that can be applied on large datasets effectively. The model automatically discovers the most distinctive topics and sub-topics in many levels of desired hierarchy by considering conceptual stability scores. The conceptual analysis helps guide the selection of the appropriate number of topics and their sub-topics. The main strength of our approach is that it is entirely unsupervised and does not require any training step. We also demonstrate the practicability of our framework to get a better understanding of textual source. Starting from



**Fig. 4.** The bubble plot of 35 discovered sub-topics. The size of the bubbles corresponds with the % share assigned to that sub-topic. For the ease of interpretation, we report the bubble plot as 1-component principle component analysis.

addressing the drawbacks of consensus matrix models that exist more than a decade, we have provided an effective and powerful framework for large-scale text mining and document clustering via NMF. We also present several state-of-the-art LDA-based topic modeling approaches that are unable to handle large dataset.

## References

1. Arun, R., Suresh, V., Madhavan, C.V., Murthy, M.N.: On finding the natural number of topics with latent dirichlet allocation: Some observations. In: Pacific-Asia Conference on Knowledge Discovery and Data Mining. pp. 391–402. Springer (2010)
2. Berry, M.W., Browne, M.: Email surveillance using non-negative matrix factorization. *Computational & Mathematical Organization Theory* 11(3), 249–264 (2005)
3. Blei, D.M., Ng, A.Y., Jordan, M.I.: Latent dirichlet allocation. *Journal of machine Learning research* 3(Jan), 993–1022 (2003)
4. Brunet, J.P., Tamayo, P., Golub, T.R., Mesirov, J.P.: Metagenes and molecular pattern discovery using matrix factorization. *Proceedings of the national academy of sciences* 101(12), 4164–4169 (2004)
5. Cichocki, A., Anh-Huy, P.: Fast local algorithms for large scale nonnegative matrix and tensor factorizations. *IEICE transactions on fundamentals of electronics, communications and computer sciences* 92(3), 708–721 (2009)

6. Deveaud, R., SanJuan, E., Bellot, P.: Accurate and effective latent concept modeling for ad hoc information retrieval. *Document numérique* 17(1), 61–84 (2014)
7. Duong-Trung, N., Schilling, N., Schmidt-Thieme, L.: Near real-time geolocation prediction in twitter streams via matrix factorization based regression. In: *Proceedings of the 25th ACM International on Conference on Information and Knowledge Management*. pp. 1973–1976. ACM (2016)
8. Fellbaum, C.: *WordNet*. Wiley Online Library (1998)
9. Gillis, N.: The why and how of nonnegative matrix factorization. *Regularization, Optimization, Kernels, and Support Vector Machines* 12(257) (2014)
10. Hornik, K., Grün, B.: topicmodels: An r package for fitting topic models. *Journal of Statistical Software* 40(13), 1–30 (2011)
11. Kim, H., Choo, J., Kim, J., Reddy, C.K., Park, H.: Simultaneous discovery of common and discriminative topics via joint nonnegative matrix factorization. In: *Proceedings of the 21th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*. pp. 567–576. ACM (2015)
12. Kim, J., Park, H.: Sparse nonnegative matrix factorization for clustering (2008)
13. Kuang, D., Choo, J., Park, H.: Nonnegative matrix factorization for interactive topic modeling and document clustering. In: *Partitional Clustering Algorithms*, pp. 215–243. Springer (2015)
14. Kuhn, H.W.: The hungarian method for the assignment problem. *Naval research logistics quarterly* 2(1-2), 83–97 (1955)
15. Miller, G.A.: Wordnet: a lexical database for english. *Communications of the ACM* 38(11), 39–41 (1995)
16. Monti, S., Tamayo, P., Mesirov, J., Golub, T.: Consensus clustering: a resampling-based method for class discovery and visualization of gene expression microarray data. *Machine learning* 52(1-2), 91–118 (2003)
17. Pauca, V.P., Shahnaz, F., Berry, M.W., Plemmons, R.J.: Text mining using non-negative matrix factorizations. In: *SDM*. vol. 4, pp. 452–456. SIAM (2004)
18. Porteous, I., Newman, D., Ihler, A., Asuncion, A., Smyth, P., Welling, M.: Fast collapsed gibbs sampling for latent dirichlet allocation. In: *Proceedings of the 14th ACM SIGKDD international conference on Knowledge discovery and data mining*. pp. 569–577. ACM (2008)
19. Roller, S., Speriosu, M., Rallapalli, S., Wing, B., Baldridge, J.: Supervised text-based geolocation using language models on an adaptive grid. In: *Proceedings of the 2012 Joint Conference on Empirical Methods in Natural Language Processing and Computational Natural Language Learning*. pp. 1500–1510. Association for Computational Linguistics (2012)
20. Wing, B., Baldridge, J.: Hierarchical discriminative classification for text-based geolocation. In: *EMNLP*. pp. 336–348 (2014)
21. Wu, Z., Palmer, M.: Verbs semantics and lexical selection. In: *Proceedings of the 32nd annual meeting on Association for Computational Linguistics*. pp. 133–138. Association for Computational Linguistics (1994)
22. Xie, B., Song, L., Park, H.: Topic modeling via nonnegative matrix factorization on probability simplex. In: *NIPS workshop on topic models: computation, application, and evaluation* (2013)
23. Xu, W., Liu, X., Gong, Y.: Document clustering based on non-negative matrix factorization. In: *Proceedings of the 26th annual international ACM SIGIR conference on Research and development in informaion retrieval*. pp. 267–273. ACM (2003)