CS388 NLP Project Report

# Text Mining for Hidden Relations and Trending

C. Vic Hu
vic@cvhu.org

Ali Unwala
aliunwala@gmail.com

*Abstract*—The main contribution of this paper has two folds. First, we formalized an illustration method to visualize topic trending on a bidirectional spanning tree, which gives a meaningful intuition to discover how hidden thematic structures in large archive of text documents change, merge and split over time. Second, we proposed an algorithm, Thematic Particle Clustering, that combines probabilistic sampling, clustering and gradient descent methods to predict upcoming topics based on a sequence of history topics. The effectiveness of our methods is demonstrated through a collection of 10,000 patent data in the field of robotics spanning over 30 years.

## I. Introduction

## II. Background

Given a collection of text-based patent documents, one intuitive idea to find out trending patterns is to examine the underlying thematic structures hidden in the text. Based on the vocabulary distribution, we want to know what the intrinsic topics are implied in the given context. One common way to do exactly that is the Probabilistic Topic Models formalized by Blei et al. [1]

### A. Probabilistic Topic Models

The main objective of topic modeling is to automatically discover the unobserved hidden structures–the topics, per-document topic distributions, and the per-document per-word topic assignments, while a collection of text documents is the only observable variables.



Fig. 1: A sample patent document (partial)

For example in Fig. 1, we have annotated a selection of words, with topics distinguished by colors. For the orange topic, we get words like flange, surface and extending, which could be interpreted as the attachment of hardware components. Similarly, the blue and green topics could be translated into topics about installation and mounting respectively. By looking at the text, most human being with common comprehensibility could easily tell what a patent data like Fig. 1 is about, and accordingly highlight the relevant keywords that compose such topics.

Nonetheless, the efficiency and accuracy of human labors don't scale up easily when the size or complexity of these patent documents increases. The objective of probabilistic topic modeling is to automate this inference process and to provide hidden insights and meaningful intelligence of big data. If we are able to successfully construct a probable thematic structure from a large archive of text data for each time slice in a sequence, we could presumably infer how these topics inherit or inspire each other, and most excitingly, predict the most likely topics in the future.

### B. Latent Dirichlet Allocation

Latent Dirichlet allocation, or LDA, is the simplest topic model [2] that assigns each word in the documents a distribution over a fixed number of topics. Instead of having a hard boundary between topic collections, LDA provides a distribution of topics per document, giving the likelihood of a mixed proportion of topic assignments. Namely, all text documents share the same set of topic collection but with different proportions to each topic. For instance in Fig. 2, although there are $K = 100$ topics overall, only a few topics were actually activated.
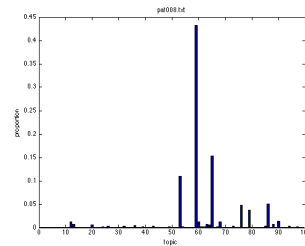


Fig. 2: A sample topic proportion of a patent

To build the generative probabilistic model, we compute the joint distribution and use it to estimate the posterior

| | |
|---|---|
| $\beta_k$, $k = 1 \cdots K$ | The K topics, represented by a distribution over words. |
| $\theta_d$, $d = 1 \cdots D$ | Topic proportions for document d, where $\theta_{d,k}$ is the topic proportion of topic k for document d. |
| $z_d$, $d = 1 \cdots D$ | Topic assignments for document d, where $z_{d,n}$ is the topic assignment for the n-th word in document d. |
| $w_d$, $d = 1 \cdots D$ | The observed words for document d, where $w_{d,n}$ is the n-th word in document d. |

TABLE I: Topic modeling notations

probability. With the notation specified in Table. I, the LDA generative process can be formalized as the following joint probability of both hidden and observed random variables:

$$p(\beta_{1:K}, \theta_{1:D}, z_{1:D}, w_{1:D})$$
$$= \prod_{i=1}^{K} p(\beta_i) \prod_{d=1}^{D} p(\theta_d) \left( \prod_{n=1}^{N} p(z_{d,n}|\theta_d) p(w_{d,n}|\beta_{1:K}, z_{d,n}) \right)$$

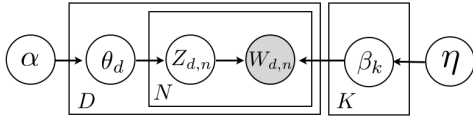which can also be expressed as a graphical model:



Fig. 3: LDA graphical model. Nodes represent variables, while edges indicate the dependency relations. The shaded node is the only observed variable (document words), and all others are the hidden variables. The $D$ plate denotes the replicated variables product over $D$ documents, while the $N$ plate denotes replication over $N$ words in each document.

Note that there are several conditional dependencies implied in the graphical models, which reflects the main principles of how LDA "think" the documents are generated:

1) Randomly pick a distribution $\theta_d$ over topics.
2) For each word in the document
   a) Randomly choose a topic from the previously-chosen distribution $\theta_{d,n}$.
   b) Randomly choose a word from the corresponding distribution $Z_{d,n}$.

Assuming this generative process is how our documents are created, now LDA uses the graphical model in Fig. 3 to infer the posterior probability of the hidden structures given our observable:

$$p(\beta_{1:K}, \theta_{1:D}, z_{1:D}|w_{1:D}) = \frac{p(\beta_{1:K}, \theta_{1:D}, z_{1:D}, w_{1:D})}{p(w_{1:D})}$$

The computation of possible topic structures is often intractable and the posterior distribution can only be approximated in most cases. To form an approximation algorithm, topic modeling can generally be categorized as sampling-based algorithms and variational algorithms. The most popular sampling method for topic modeling is Gibbs sampling, which introduces a sequence of random variables to construct a Markov chain and collects samples from the limiting distribution to estimate the posterior. Instead of using samples to approximate the posterior, variational methods find the closest parameterized distribution candidate by solving optimization problems [2] [3].



Fig. 4: The top 3 topics of a sample patent

C. Limitations & Potential Improvements

Although LDA provides a powerful perspective to browsing and interpreting the implicit topic structures in our patent corpus, there are a few limitations it imposes against further discoveries. An extensive amount of research has been focused on relaxing some of the assumptions made by LDA to make it more flexible and suitable for various adaptations in more sophisticated context.

LDA is essentially a bag-of-words probabilistic model. Namely, it constructs a word-frequency vector for each document but disregards the word ordering and the neighboring context. Although this assumption looses the syntactic information and sometimes seems unrealistic when processing natural language, it is usually good enough when capturing the document semantics and simplifying hidden structural inferences. Nonetheless, for more sophisticated tasks such as

language generation or writing style modeling, the bag-of-words assumption is apparently insufficient and needs to be relaxed. In these cases, there are variants of topic models that generate topic words conditioned on the previous word [4], or switches between LDA and hidden Markov models (HMM) [5].

The LDA graphical model in Fig. 3 is invariant to the ordering of our patent documents, which could be inappropriate if the hidden thematic structure is actually dependent on sequential information such as years published, which is typical in document collections spanning years, decades or centuries. To discover how the topics change over time, the dynamic topic model [6] treats topics as a sequence of distributions over words and tracks how they change over time.

In either LDA or more sophisticated dynamic topic models [6], the number of topics $\beta_{1:K}$ is determined manually and assumed to be fixed. One elegant approach provided by the Bayesian nonparametric topic model [7] is to find a hierarchical tree of topics, in which new documents can now imply previously undiscovered topics.

To include additional attribute information associated with the documents such as authorships, titles, geolocation, citations and many others, an active branch of research has been performed to incorporate meta-data in topic models. The author-topic model [8] associates author similarity based on their topic proportions, the relational topic model [9] assumes document links are dependent on their topic proportion distances, and more general purpose methods such as Dirichlet-multinomial regression models [10] and supervised topic models [11].

Many other extensions of LDA are available, including the correlated topic model [12], pachinko allocation machine, [13], spherical topic model [14], sparse topic models [15] and bursty topic models [16].

## III. Problem Definition and Algorithm

### A. Task Definition

### B. Algorithm Definition

1) Assign K topics to N particles uniformly
2) Add Gaussian noise to particles
3) Cluster particles into K groups (TF-IDF weights with cosine/Euclidean distances)
4) Compare the clusters with topics from the next year, apply discounts to current weights, and adjust to new weights
5) repeat

## IV. Experimental Evaluation

### A. Methodology

### B. Results

### C. Discussion

## V. Related Work

## VI. Future Work

## VII. Conclusion

## References

[1] Blei, D. Introduction to Probabilistic Topic Models. Princeton University. 2011.

[2] Blei, D., Ng, A. and Jordan, M. Latent Dirichlet allocation. Journal of Machine Learning Research, 3:993-1022, January 2003.

[3] Hoffman, M., Blei, D. and Bach, F. On-line learning for latent Dirichlet allocation. In Neural Information Processing Systems, 2010.

[4] Wallach, H. Topic modeling: Beyond bag of words. In Proceedings of the 23rd International Conference on Machine Learning, 2006.

[5] Griffiths, T., Steyvers, M., Blei, D. and Tenenbaum, J. Integrating topics and syntax. In L. K. Saul, Y. Weiss, and L. Bottou, editors, Advances in Neural Information Processing Systems 17, pages 537-544, Cambridge, MA, 2005. MIT Press.

[6] Blei, D. and Lafferty, J. Dynamic topic models. In International Conference on Machine Learning, pages 113-120, New York, NY, USA, 2006. ACM

[7] Teh, Y., Jordan, M., Beal, M. and Blei, D. Hierarchical Dirichlet process. Journal of the American Statistical Association, 101(476):1566-1581, 2006.

[8] Rosen-Zvi, M., Griffiths, T., Steyvers, M. and Smith, P. The author-topic model for authors and documents. In Proceedings of the 20th Conference on Uncertainty in Artificial Intelligence, pages 487-494. AUAI Press, 2004.

[9] Chang, J. and Blei, D. Hierarchical relational models for document networks. Annals of APplied Statistics, 4(1), 2010.

[10] Mimno, D. and McCallum, A. Topic models conditioned on arbitrary features with Dirichlet-multinomial regression. In Uncertainty in Artificial Intelligence, 2008.

[11] Blei, D. and McAuliffe, J. Supervised topic models. In Neural Information Processing Systems, 2007.

[12] Blei, D. and Lafferty, J. A correlated topic model of Science. Annals of Applied Statistics, 1(1):17-35, 2007.

[13] Li, W. and McCallum, A. Pachinko allocation: DAG-structured mixture models of topic correlations. In International Conference on Machine Learning, pages 577-584, 2006.

[14] Reisinger, J., Waters, A. ,Silverthorn, B. and Mooney, R. Spherical topic models. In International Conference on Machine Learning, 2010.

[15] Wang, C. and Blei, D. Decoupling sparsity and smoothness in the discrete hierarchical dirichlet process. In Y. Bengio, D. Schuurmans, J. Lafferty, C. K. I. Williams, and A. Culotta, editors, Advances in Neural Information Processing Systems 22, pages 1982-1989. 2009.

[16] Doyle, G. and Elkan, C. Accounting for burstiness in topic models. In International Conference on Machine Learning, pages 281-288. ACM, 2009.