

Topic Classification using Latent Dirichlet Allocation

Adrian Guthals (aguthals@cs.ucsd.edu),
David Larson (dplarson@ucsd.edu),

CSE 250B: Project #3
University of California, San Diego

February 28, 2013

Abstract

LDA, Gibbs sampling, topic classification of documents, datasets used, results and their meaning, conclusions

1 Introduction

Introduce topic classification of documents. Then transition into formal definitions of LDA and Gibbs Sampling.

Elkan's lecture notes [1]

1.1 Latent Dirichlet Allocation

Latent Dirichlet Allocation (LDA) is a multinomial

$$p(\gamma|\alpha) = \frac{1}{D(\alpha)} \prod_{s=1}^m \gamma_s^{\alpha_s-1} \quad (1)$$

$$D(\alpha) = \int_{\gamma} \prod_{s=1}^m \gamma_s^{\alpha_s-1} \quad (2)$$

$$D(\alpha) = \frac{\prod_{s=1}^m \Gamma(\alpha_s)}{\Gamma(\sum_{s=1}^m \alpha_s)} \quad (3)$$

1.2 Gibbs Sampling

$$p(z_i = j | \bar{z}', \bar{w}) \propto \frac{q'_{jw_i} + \beta_{w_i}}{\sum_t q'_{jt} + \beta_t} \frac{n'_{mj} + \alpha_j}{\sum_k n'_{mk} + \alpha_k} \quad (4)$$

2 Design and Analysis of Algorithms

Discuss how we're implementing LDA and Gibbs Sampling.

Dataset	Documents	Vocabulary	α	β
Classic400	400	6205	0.01	0.1
KOS	400	6906	0.01	0.1

Table 1: Composition of the two datasets used in this study and the chosen Gibbs hyperparameters α and β . Note that the KOS dataset used is a reduced version of the original KOS dataset, which contains 3430 documents.

3 Design of Experiments

3.1 Datasets

Two datasets were classified using LDA: Classic400, a collection of English documents from three research areas (aeronautics, medicine, and library science); and KOS, a collection of English blog posts from dailykos.com (see Table 1 for details on their sizes) [2, 3]. While the Classic400 dataset provided the true labels of its source topics, the KOS dataset provided no such information. Also, to reduce run times, we elected to use a reduced version of the KOS dataset containing only the first 400 of the original 3430 documents.

3.2 Hyperparameters

For both datasets we apply LDA with three topics ($K = 3$). Since α is positively correlated to the number of topics present per document and we are only considering three topics, we chose values of $\alpha \ll 1$, which implies that every document contains only a few topics. Likewise, we chose values of $\beta \ll 1$, which implies that every word in every document corresponds to a small number of topics. Previous studies have looked into the effect of α and β on the clustering results for LDA models, some going so far as to suggest likely values based on the number of topics in a dataset (see [4]). However, for the purposes of this study our selection criteria of α and β is only that they are both small, and we have therefore chosen to omit any sort of optimization of their values (see Table 1).

3.3 Stopping Criteria of Gibbs Sampling

For both datasets we chose the Gibbs Sampling stopping criteria as 100 epochs, although the majority of improvement in the LDA classification occurs during the first 10 epochs.

3.4 Results

3.4.1 Document Topic Clustering

Figure 1 shows the clustering of documents relative to each topic as projected onto a triangular simplex, where the each vertex indicates $\theta = 1$ the denoted topic.

3.4.2 Most Common Words Per Topic

If a LDA model is setup to classify more than 4 topics, then a simplex plot will not be a viable method for checking the clustering of the results. Instead, lists of the most commonly occurred words per topic can be analyzed to gauge the performance of the LDA. Table 2 shows the most commonly occurred words for each topic of the two datasets. The most common words from Classic400 for each topic are all related and correspond well to the true topic labels (topic 1: medicine, topic 2: aeronautics, and topic 3: library science). Meanwhile, the most common words for KOS for each topic seem to overlap and be derived from same general topic, namely

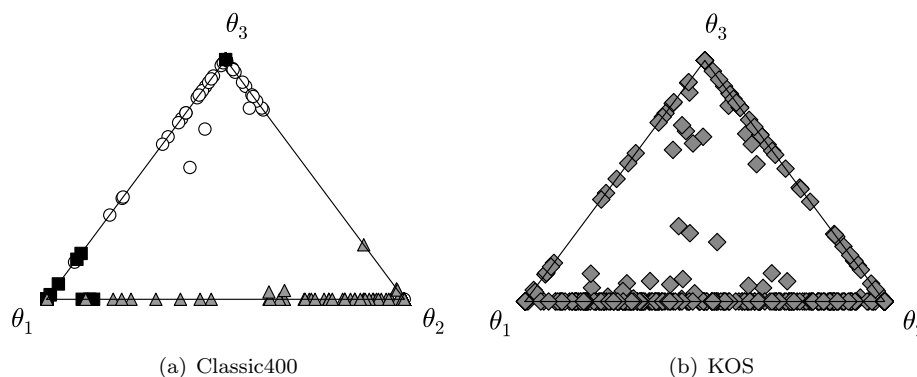


Figure 1: Clustering of documents based on their θ_1 , θ_2 and θ_3 , projected on a triangular simplex. Documents in Classic400 are labeled based on their true labels (topic 1: white circles with black edges, topic 2: black solid squares, topic 3: gray solid triangles) while documents in KOS are have the same marker label as the true labels are unknown.

politics, or more specifically, the 2008 US Presidential campaign between George Bush and John Kerry. As the source of KOS is a political blog site, it is not surprising that politics shows up prominently as a topic. However, the lack of diversity between the most common words of each topic and therefore lack of unique topics is more likely the result of selecting only the first 400 documents from the original dataset rather than an issue with our LDA model implementation.

4 Findings and Lessons Learned

Thoughts on: LDA as a model, Gibbs Sampling as a training method, performance issues, results of the experiments

References

- [1] C. Elkan, “Text mining and topic models,” February 2013.
- [2] “Classic400 dataset.” [Online]. Available: <ftp://ftp.cs.cornell.edu/pub/smart/>
- [3] A. Frank and A. Asuncion, “UCI machine learning repository,” 2010. [Online]. Available: <http://archive.ics.uci.edu/ml>
- [4] T. L. Griffiths and M. Steyvers, “Finding scientific topics,” *Proceedings of the National Academy of Sciences of the United States of America*, vol. 101, no. Suppl 1, pp. 5228–5235, 2004. [Online]. Available: <http://www.pnas.org/content/101/suppl.1/5228.abstract>

Rank	Classic400			KOS		
	Topic 1	Topic 2	Topic 3	Topic 1	Topic 2	Topic 3
1	patients	boundary	system	november	bush	kerry
2	ventricular	layer	retrieval	kerry	campaign	time
3	cases	wing	scientific	vote	republican	bush
4	fatty	mach	research	voting	race	people
5	left	supersonic	language	polls	war	jul
6	acids	wings	systems	house	general	party
7	nickel	velocity	journals	poll	elections	john
8	aortic	ratio	field	bush	percent	democratic
9	blood	shock	methods	senate	state	media
10	glucose	effects	subjects	governor	news	voters

Table 2: Ten most commonly occurred words for each topic classification for the classic400 and KOS datasets. Rank indicates how frequently a word appears in each topic, with 1 being the most occurred.