# Alternatives for engineering and evalution of algorithms for LDA.

Di Wang, Victor Huang, James Cook, Andrew Gambardella, Chenyu Zhao, Satish Rao

## ABSTRACT

In this paper, we empirically study standard algorithms for topic modelling using Latent Dirichlet Allocation. We also develop geometry based algorithms for learning the Latent Dirichlet Allocation (LDA) which are as effective as inference based methods while being more efficient. In particular, the most effective algorithms use a novel combination of dimension reduction, projection, $k$-means, and a scaling procedure.

We argue that the evaluation of topic modelling should include a prediction task that is closely aligned with the model. LDA, in particular, generates the words in a document. Thus, predicting missing words is appropriate for the evaluation of the LDA model, both when applied to a corpus and for genarated data. Moreover, this task is useful for recommendation systems, keyword suggestion, tag prediction, etc. Finally, for practitioners this evaluation task may be more indicative of the usefulness of various methods than internal measures such as perplexity or than indirect methods where other prediction algorithms are applied to learned topic models.

Finally, to empirically study algorithms, one needs to understand properties of datasets. We thus provide an analysis that predicts whether a corpus can be effectively modelled. For generated datasets, the prediction is based on the diversity of topics in a typical document, and the concentration of words in a topic.

## 1. INTRODUCTION

The Latent Dirichlet Allocation topic model has been tremendously influential in the development of methods for analyzing documents, and other types of data. Numerous algorithms for learning such models as well as variations of this model have been proposed; see [4] for a recent survey. The bulk of these methods are based on inference techniques; these methods proceed by using learning the parameters using methods which include Gibbs sampling, EM, variational methods, and modifying the model for easier computation. The evaluation of the methods have proceeded typically along two lines. One is to train on a set of documents and evaluate how well the resulting model predicts a test set, typically using perplexity as a measure. Another is to use the model as a set of features for some classification task, say for example document classification, and apply some learning method, say for example, using a support vector machine.

In this paper, we suggest modifying both the inference approach to learning topic models and suggest a different evaluation method. The first yields an improved algorithm for learning LDA, and the second allows for better understanding of topic learning efficacy.

In terms of inference algorithms, we take the view that the model parameters are geometric objects. That is, topic centers are simply points in the word space when the data consist of documents. This view is quite a traditional view of data, e.g., that it is generated from points in space under some noise model. This view of data has long been associated with algorithms such as nearest neigbors or $k$-means for classification or summarization. Unfortunately, as we show, $k$-means and nearest neighbors algorithms are inferior to existing LDA inference algorithms for LDA generated data. We can, however, combine $k$-means, along with dimension reduction, and a scaling step to produce an effective algorithm for this type of data. For real world data, our methods remain as effective as LDA algorithms and more effective than $k$-means, but fall short of nearest neighbor techniques (as do all topic modelling approaches.) We call refer to our algorithm as the Projector algorithm.

We note here that recent work in [1] use linear algebraic methods, which can be viewed to some extent as geometric algorithms, to give an algorithm that provably learns the LDA parameters given a polynomial number of examples and polynomial time.

The context in which we do our evaluation varies from the standard. We measure performance on the task of predicting a set of dropped out words in a document. Topic models in general and the Latent Direchlet Allocation topic model, in particular, generates documents word by word. Thus, a model's effectiveness for predicting a word from the document is both fair and is easy to understand as a raw score and allows for the comparison to any prediction method whether it is specifies a distribution or not. We also note that this task is interesting in itself; if a document represents a set of products bought by a user, predicting a new product is essentially the recommendation problem.

Before proceeding, we note that this paper proposes an algorithm, argues for more informative evaluations, and does experimental studies existing algorithms. We feel these things go very much together. The proof for a methodology for algorithm development is, after all, an effective algorithm.

We proceed with related work in section 2, a description of the lda model in section 3, the data generation process in section 4, a description of the algorithms we study in section 6, the results of our experiments in section 7, a bit of analysis in section 8, and conclude in 9.

## 2. RELATED WORK

An early version of topic models could be seen in latent semantic indexing which views a document as being about a subset of a larger set of topics. Topics are in turn are about (positive numbers) or not about (negative numbers) certain words [13]. They used this view to provide an algorithm based on linear algebra that learned the structure of their model. They also argued that this model provided some insight into the wide applicability of principal components analysis in data analysis. The algorithmic tool at the core of their work (and principal components) was the singular value decomposition.

Shortly thereafter, a probabilistic framework was provided by Puzicha and Hofmann [7]. They termed their method probabilisitc latent semantic indexing (pLSI), and modelled topics as probability distributions over words (no negative numbers in this description). Documents were then presumed to be generated as a mixture of these topics. The large number of parameters in the PLSI model motivated Blei, Ng and Jordan to provide a generative model for these parameters as well; they provided the enormously influential Latent Dirichlet Allocation(LDA) topic model [6].

The most closely related work in theory are the aforementioned algorithms provided by Arora et.al. [2] which learn pLSI under certain assumptions, and a breakthrough by Anandkumar[1] which gives a provably converging algorithm for learning LDA with polynomial data requirements. There is a long literature in learning mixtures of distributions in statistics and more recently in theoretical computer science. See [9] for a recent breakthrough on learning mixtures of Gaussians and a discussion of this field.

There is ample work describing methods for optimizing LDA [4]. Again, recent a recent examples are described in [11], [12].

There is also significant work on extending LDA and topic models to better model data as well as to augment them with other types of information. These are discussed in [4]. The hierarchical LDA model [5] is perhaps the most relevant. This model is based on the chinese restuarant process where new customers arrive and chose to join a table or create a new one. This process can be used to create a hierarchical structure of documents and topics. The authors argue that this is a better model for real data and we tend to agree though our preliminary experiments do not bear this out. Still, we began by understanding the simpler model which continues to have wide application.

Other examples of variations of the LDA model include adding the notion of a manifold to the model [8], adding partially labeled data in [14], and sparse LDA models [19].

This paper, in addition to providing some insight into the effectiveness of LDA algorithms, seeks to provide an infrastructure to evaluate effective topic modelling methods. We should point out that the machine learning community is made efforts in this regard. See for example, MLComp [10]. We also highlight the java based infrastructure that we used for topic modelling [11] which is part of a larger MALLET java package for machine learning.

## 3. LDA MODEL

LDA was introduced in [6] as a generative process. As a model it is widely applied in various domains such as information retrieval, collaborative filtering, vision, and bioinformatics. In this work we will adopt the language of text collections, and denote entities as 'word', 'document', and 'corpus', since they give intuition in discussing topic models.

The basic idea is that there exist $k$ underlying latent topics. Each document is a mixture of the latent topics, where the topic mixture is drawn from a dirichlet distribution. More precisely, there are $n$ documents, $m$ words, and $k$ topics. The model has a $m \times k$ word-topic matrix $A$, where the $i$-th column $A_i$ specifies a multi-

nomial distribution on the $m$ words for topic $i$. For a document $w$, we first choose $\vec{\theta}$, its distribution over topics, which can take values in the $(k-1)$-simplex, and has the following Dirichlet distribution

$$p(\theta|\vec{\alpha}) = \frac{\Gamma(\sum_{i=1}^{k}\alpha_i)}{\Pi_{i=1}^{k}\Gamma(\alpha_i)}\theta_1^{\alpha_1-1}\cdots\theta_k^{\alpha_k-1}$$

where $\vec{\alpha}$ is parameter of the model. The number of words in document $w$ is sampled from $Poisson(l)$. For each word $w_i$, a topic $z_i \sim Multinomial(\vec{\theta})$ is chosen, then the actual word $w_i \sim Multinomial(A_{z_i})$ is generated. Equivalently, in matrix form, there are the $m \times k$ word-topic matrix $A$, and $k \times n$ topic-document matrix $W$ whose columns are drawn i.i.d from $Dirichlet(\vec{\alpha})$. The product $M = AW$ is the $m \times n$ term-document matrix where column $M_i$ is document $i$'s distribution on words. Document $i$ is generated by sampling words i.i.d from $Multinomial(M_i)$. We are interested in the case where $A$ is of full rank, since if the columns of $A$ are not independent, intuitively it means there exists some document which is covered completely by a set of topics $I$, but at the same time also completely covered by another set of topics $J$ which is disjoint from $I$. In our experiments, the randomly generated $A$ matrices are almost always of full rank.

## 4. DATA

Since our focus is on how effectively the algorithms learn the model, we use synthetic datasets generated from the LDA model for a range of parameters. Our data generator takes in parameters

- $n, m, k$, number of documents, words, and topics respectively

- $\alpha$, the Dirichlet parameter for generating documents' distributions over topics as in the LDA model. In our experiments we work with symmetric Dirichlet distributions, where $\alpha_i = \ldots = \alpha_k = \alpha$

- $\beta$, we generate the columns of word-topic matrix $A$ from a $m$ dimensional Dirichlet distribution with parameter $\vec{\beta}$. Again we work with symmetric Dirichlet where $\beta_1 = \ldots = \beta_m = \beta$.

- $l$, the Poisson parameter controlling the expected number of words in a document.

Intuitively the Dirichlet parameter $\alpha$ is a crude measure of the sparsity of the sampled distribution over topics. When $\alpha = 1$, all points in the $k-1$ simplex have the same probability density. When $\alpha < 1$, distributions that are more concentrated on a few topics are prefered by the Dirichlet. The same applies to $\beta$ and the topic's
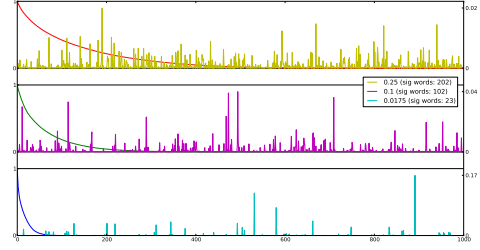


**Figure 1: Plot of distributions on words for various $\beta$. $m = 1000$, each distribution is plotted along with its cdf after sorting the words by popularity. Refer to the y-axis on the right for the scaling of the distributions. In general, larger $\beta$ values yield flatter distributions.**

distribution on words. See figure 1 for typical word distributions sampled from the Dirichlet distribution with various $\beta$'s as parameter.

To help understand the dataset, we compute the values $sig\_topic$ and $sig\_word$. For a document with distribution $\vec{\theta}$ over topics, $sig\_topic$ is the smallest $t$ such that the union of the $t$ heaviest topics in $\vec{\theta}$ has an aggregate probability of at least 0.8. Intuitively, $sig\_topic$ is the number of significant topics of a document. Analogously, for a topic's distribution over words, $sig\_word$ is the smallest number of most popular words with an aggregate probability of at least 0.8. Instead of using $\alpha$ and $\beta$, we use the average $sig\_topic$ and average $sig\_word$ to characterize our datasets.

We also evaluate on a topic distribution which obeys a power law. Instead of using the same $\alpha$ parameter for all the topics, we vary it so that the expected usage.

We also evaluate the methods on some standard real world datasets. We used the Classic-3 datasets (Cran,Med,Cisi) [15], a corpus (AP) from the Associated Press [16] pruning stop words and infrequent words, and a bag of words dataset (NIPS) from UC Irvine[18], MovieLens [17].

## 5. EXPERIMENTS

### 5.1 Prediction task

For a corpus of documents, we randomly divide the documents into the training set and the testing set, where each document is put into the training set with probability $p_t$ independently. For each document in the testing set, we hold out a certain percentage $H$ of the distinct words in the document uniformly at random. The training set is given to the algorithms. After training, each algorithm gets the testing documents, and for each document predicts $s$ terms not in the observed part of the

testing document. We use the precison of the $s$ predicted terms as the score of an algorithm on a specific testing document. The score of an algorithm on the corpus is its average score over all testing documents. In our experiments, we use $p_t = 0.9, H = 30\%, s = 3$ as our parameters.

This prediction task is widely applicable in practice, especially in online settings with a large amount of user-generated data. A familiar example is the collaborative filtering system by which Netflix leverages the preferences of a large user population to recommend new films to a customer after observing his or her viewing history. For our purpose, the prediction task provides a more straightforward algorithmic measure than statistical measures such as perplexity and likelihood, which are commonly adopted in machine learning, but not applicable to algortihms that don't reconstruct a statistical model.

## 5.2 Recovery task

For the algorithms that reconstruct a topic matrix $\hat{A}$ in the process, we also measure how close $\hat{A}$ and $A$ are. For each learned topic $\hat{A}_i$ and real topic $A_j$, we compute $cos(\hat{A}_i, A_j)$, then find a maximum matching between the learned topics and real topics. We evaluate the average cosine similarity between the matched real and learned topics. We also carry out the above computation using total variation distance between distributions, and get same qualitative results between algorithms.

## 6. ALGORITHMS

## 6.1 Previous Methods.

We method produces a word distribution for teach test document and outputs the most frequently unseen words in this distribution.

- **Baseline:** uses the training set word distribution.

- **KNN:** KNN finds the $k$ most similar training documents to a testing document, where similarity is defined as the cosine between the two documents as vectors in the word space. For a testing document, KNN predicts the $s$ unseen words that are most frequent in its $k$ closest training documents. Notice baseline is just KNN with $k$ equal to the number of training documents.

- **LDA:** Davide Blei's implementation [3] based of variational EM uses an "estimation" phase to train a topic model on the training data. Then use the "inference" routine to infer a topic distribution for a test document which than uses the topic definitions to produce a word distribution for the document.

- K-means: Uses $k$-means with cosine similarilty to produce a topic matrix, and uses LDA inference to produce a word distibution for each test document.

- **LDA(MALLET)** Uses the implementation of LDA in Mallet based on Gibbs sampling [11]

- **LDAT, LDAC** For generated LDA datasets, we have two "cheating" algorithms as benchmarks. LDAT knows the real term-topic matrix $A$ of the model and uses the "inference" routine to find a word distribution for each document. LDAC knows the real term-document matrix $M$ for each testing document, and uses the real word distribution for a test document. LDAC is the best we can do given sampling noise.

- **LSI**. Computes the best rank $k$ subspace approximation of the document-word matrix. Then projects test document into the subspace to find a word distribution.

- **Projector**. We have two versions of Projector which is described below which produces a topic word matrix. Then LDA inference is used on each test document to produce a word distribution.

## 6.2 Projector

Projector is our new algorithm that builds upon LSI, and reconstructs a term-topic probability matrix $\hat{A}$. The motivation is that SVD is computationally more efficient than the LDA algorithm, and has a clear geometric interpretation, but doesn't recover the topics as distributions of words. We aim to start from the subspace computed by SVD, and use some straightforward operation to construct the topics. Our algorithm is based on geometric intuition of the documents as points in the high dimensional word space. The algorithm is as follows

**Input** $\hat{M}$: observed distributions of training documents, $k$: number of topics, $\delta$: algorithm parameter

**Shift** Shift the training documents to be centered at the origin.
$center = \frac{1}{n} \sum_{i=1}^{n} \hat{M}_i$
$\hat{M}_i = \hat{M}_i - center \qquad \forall i = 1, \ldots, n$

**SVD** Compute the U, the best rank $(k-1)$-dimension approximation to the column space of $\hat{M}$
Project all $\hat{M}_i$'s to the subspace $U$, denote $V_i$ as the projections.

**Clustering** Use k-means to cluster the $V_i$'s into $k$ clusters, where in the k-means algorithm the distance between two points $x, y$ is defined as $1 - cos(x, y)$. Let $C_1, \ldots, C_k$ be the centers of the $k$ clusters (center in the sense as in euclidean distance).

**Scale** Scale $C_1, \ldots, C_k$ by the smallest common scalar so that $\delta n$ of $V_1, \ldots, V_n$ are contained in the hull with $C_1, \ldots, C_k$ as vertices.

**Whitening** Make all $C_i$ distribution over words: $C_i = C_i + center$, truncate the negative entries in $C_i$ to be 0, normalize $C_i$ so the sum of entries is 1. Return $\hat{A}_i = C_i$ be the recovered topics.

We illustrate in figure 6.2 how our algorithm works using the a visualization on two datasets with $k = 3$ topics. Notice after the *Shift* step, we want to find the best $(k-1)$-dimensional subspace since the columns from the topic-document matrix are from the $(k-1)$-simplex.

We use estimated $\hat{A}$ and the inference procedure of the LDA algorithm to predict words for testing documents. We use the inference procedure of LDA since LDAT also uses it, and then we can attribute the performance difference between LDAT and Projector to the quality of $\hat{A}$ compared to the real topics.

We also experimented with a version of projector which does not use LSI as a first step; it just proceeds with $k$-means, then we project the documents into the subspace that contains the means and scale as above. The results were quite similar to the results above so we do not include them here.

## 7. EXPERIMENT RESULTS

We generated datasets from the LDA model for a range of parameters, and tested the algorithms discussed in previous section on the prediction task. For the LDA algorithm and our Projector algorithm, we also have results for the recovery task. We experimented with $k$ from 3 to 30, and in each case, a set of $\alpha$ and $\beta$ to cover a wide range of *sig_word* and *sig_topic*. See figure 3 for the prediction results of the algorithms on a representative set of datasets.

In table 1, we give results for the various algorithms on the uniform $\alpha$ and $\beta$ datasets. Here, we see that projector is competitive with all other algorithms. In table 2, we give cosine similarity measures for these datasets for mallet lda, k-means, and projector.

For $\alpha$ chosen so that the topic distribution obeys a power law, we present results in table **??**. Here, we see that projector again performs robustly well. We note that the Mallet implementation optimizes over varying values of $\alpha$.

Finally, table 4 contains results for the word prediction task on the standard real world datasets we discussed earlier.

Also see table 4 for prediction results on real datasets. There we see that LDA is no longer dominated by LSI and Projector. Still, nearest neighbor's performance dominates. In the appendix, in table 6 we show the most popular words in a sampling of topic vectors recovered by Projector in the AP dataset.

Typically, we saw the topic matching to correlate with prediction performance. See figure **??** to see this result.

### 7.1 Runtime

Figure 4, compares the runtimes of Projector with Mallet LDA and suggests that Mallet's runtime increases linearly with respect to average document length and the number of documents. Projector appears sublinear in each, which suggests that the our timings are influenced by fixed startup costs. Not showen, is Mallet's better scaling with respect to vocabulary size as Projector uses dense matrices. Projector's dependence remains linear in this case.

## 8. ANALYSIS

In the previous section, we defined the notion of typical number of words in the support of a topic, or sig_words, and the notion of a typical number of topics in a document. Clearly, if both get very large one gets to a trivial topic model where each document is generated by choosing words independently from a single probability distribution.

On the other hand if sig_words and sig_topics are small each document should have relatively small support compared with the corpus. Thus, in such cases we should easily distinguish the data from the trivial topic model. We will proceed by showing that the probability of cooccurrence in documents of two words $i$ and $j$ differs significantly from in the trivial topic model. This can be represented as a matrix which refer to as the co-occurence matrix.

The expected co-occurence matrix can, of course, be calculated precisely from the topic and document distributions. But we give simple, even trivial, calculations that provide insight based on a simplified topic model.

### 8.1 A Uniform Topic Model.

We proceed by calculating the difference in co-occurence matrices of a corpus generated by a nontrivial topic model from the trivial topic model.

Let $k$ be the total number of topics. Let $m$ be the vocabulary size. Let $t$ be the number of topics in a document and assume that each word is chosen uniformly from these $t$ topics. Each topic is a uniform distribution over $w$ words. Let $l$ be the number of words in the document. We examine word cooccurrences for $w_i$ and $w_j$: that is, the probability that two words generated in a document are $w_i$ and $w_j$. Note that in a document of length $l$, there are $\binom{l}{2}$ possible cooccurences between $w_i$ and $w_j$. We compute the probability of a word cooccurrence for data from a topic distribution versus the trivial model with the working vocabulary: the union of significant

| sig_topics | sig_words | Baseline | LSI-15 | kmeans-15 | knn-25 | lda-15 | ldaC | ldaT | malletlda-15 | projector-15 |
|---|---|---|---|---|---|---|---|---|---|---|
| 5.09 | 19.4 | 0.21 | 0.5 | 0.23 | 0.46 | 0.28 | 0.57 | 0.54 | 0.52 | **0.54** |
| 3.03 | 199 | 0.06 | **0.18** | 0.06 | 0.15 | 0.06 | 0.23 | 0.23 | 0.15 | **0.18** |
| 5.01 | 53.6 | 0.11 | 0.3 | 0.11 | 0.27 | 0.12 | 0.36 | 0.34 | 0.31 | **0.32** |
| 1.19 | 19.73 | 0.04 | 0.81 | 0.81 | 0.81 | 0.73 | 0.84 | 0.84 | 0.79 | **0.82** |
| 7.42 | 55 | 0.17 | 0.21 | 0.17 | 0.19 | 0.17 | 0.3 | 0.26 | **0.25** | 0.23 |
| 1.19 | 206.13 | 0.06 | **0.32** | 0.29 | 0.31 | 0.27 | 0.32 | 0.32 | 0.29 | **0.32** |
| 2.46 | 53.2 | 0.09 | 0.53 | 0.48 | 0.51 | 0.49 | 0.56 | 0.56 | 0.52 | **0.54** |
| 7.36 | 20.6 | 0.22 | 0.35 | 0.22 | 0.3 | 0.22 | 0.48 | 0.41 | **0.41** | 0.4 |
| 3.06 | 21.53 | 0.15 | 0.63 | 0.56 | 0.59 | 0.6 | 0.69 | 0.66 | 0.65 | **0.67** |
| 3.07 | 120.4 | 0.07 | 0.25 | 0.1 | 0.24 | 0.08 | 0.3 | 0.28 | 0.23 | **0.26** |
| 2.48 | 119.53 | 0.06 | 0.29 | 0.18 | 0.29 | 0.12 | 0.32 | 0.32 | 0.28 | **0.3** |
| 2.51 | 204.33 | 0.06 | 0.19 | 0.08 | 0.15 | 0.07 | 0.24 | 0.23 | 0.17 | **0.22** |
| 3.06 | 53.53 | 0.11 | 0.48 | 0.33 | 0.46 | 0.39 | 0.52 | 0.5 | 0.48 | **0.5** |
| 5.07 | 205.67 | 0.07 | **0.09** | 0.06 | 0.06 | 0.07 | 0.14 | 0.13 | 0.07 | **0.09** |
| 1.2 | 51.6 | 0.04 | 0.68 | 0.64 | **0.69** | 0.6 | 0.71 | 0.7 | 0.68 | **0.69** |
| 5.05 | 120.87 | 0.08 | 0.18 | 0.08 | 0.14 | 0.08 | 0.23 | 0.21 | 0.16 | **0.2** |
| 7.37 | 116.53 | **0.1** | 0.09 | **0.1** | 0.09 | **0.1** | 0.17 | 0.15 | **0.1** | **0.1** |
| 2.5 | 19.4 | 0.12 | 0.68 | 0.64 | 0.66 | 0.61 | 0.75 | 0.73 | 0.68 | **0.72** |
| 7.39 | 208.13 | **0.08** | 0.05 | **0.08** | 0.05 | **0.08** | 0.11 | 0.1 | 0.07 | 0.06 |
| 1.18 | 118.07 | 0.06 | **0.48** | **0.48** | **0.48** | 0.42 | 0.49 | 0.49 | 0.47 | **0.48** |

Table 1: Bold indicates champion non-cheating method for each row.

| sig_topics | sig_words | lda-15 | malletlda-15 | projector-15 | kmeans-cosine | lda-cosine | mallet-cosine | projector-cosine |
|---|---|---|---|---|---|---|---|---|
| 5.09 | 19.4 | 0.28 | 0.52 | **0.54** | 0.14 | 0.62 | 0.98 | 0.98 |
| 3.03 | 199 | 0.06 | 0.15 | **0.18** | 0.45 | 0.48 | 0.67 | 0.88 |
| 5.01 | 53.6 | 0.12 | 0.31 | **0.32** | 0.23 | 0.41 | 0.94 | 0.95 |
| 1.19 | 19.73 | 0.73 | 0.79 | **0.82** | 0.15 | 0.93 | 0.99 | 1.0 |
| 7.42 | 55 | 0.17 | **0.25** | 0.23 | 0.23 | 0.37 | 0.76 | 0.73 |
| 1.19 | 206.13 | 0.27 | 0.29 | **0.32** | 0.45 | 0.81 | 0.83 | 0.94 |
| 2.46 | 53.2 | 0.49 | 0.52 | **0.54** | 0.22 | 0.84 | 0.97 | 0.98 |
| 7.36 | 20.6 | 0.22 | **0.41** | 0.4 | 0.14 | 0.39 | 0.98 | 0.94 |
| 3.06 | 21.53 | 0.6 | 0.65 | **0.67** | 0.15 | 0.91 | 0.99 | 0.99 |
| 3.07 | 120.4 | 0.08 | 0.23 | **0.26** | 0.34 | 0.49 | 0.87 | 0.94 |
| 2.48 | 119.53 | 0.12 | 0.28 | **0.3** | 0.34 | 0.59 | 0.88 | 0.95 |
| 2.51 | 204.33 | 0.07 | 0.17 | **0.22** | 0.45 | 0.49 | 0.75 | 0.9 |
| 3.06 | 53.53 | 0.39 | 0.48 | **0.5** | 0.22 | 0.75 | 0.96 | 0.97 |
| 5.07 | 205.67 | 0.07 | 0.07 | **0.09** | 0.45 | 0.46 | 0.42 | 0.69 |
| 1.2 | 51.6 | 0.6 | 0.68 | **0.69** | 0.23 | 0.84 | 0.98 | 0.99 |
| 5.05 | 120.87 | 0.08 | 0.16 | **0.2** | 0.34 | 0.41 | 0.78 | 0.86 |
| 7.37 | 116.53 | **0.1** | **0.1** | 0.1 | 0.33 | 0.39 | 0.43 | 0.58 |
| 2.5 | 19.4 | 0.61 | 0.68 | **0.72** | 0.14 | 0.86 | 0.99 | 0.99 |
| 7.39 | 208.13 | **0.08** | 0.07 | 0.06 | 0.46 | 0.46 | 0.23 | 0.41 |
| 1.18 | 118.07 | 0.42 | 0.47 | **0.48** | 0.33 | 0.84 | 0.93 | 0.97 |

Table 2: Bold indicates champion results and we use cosine similiarity where matching to model topics uses maximum weight matching.

| sig_topics | sig_words | Baseline | lda-15 | ldaC | ldaT | malletlda-15 | projector-15 | lda-cosine | mallet-cosine | projector-cosine |
|---|---|---|---|---|---|---|---|---|---|---|
| 1.4 | 98.73 | 0.12 | 0.5 | 0.53 | 0.52 | 0.43 | **0.51** | 0.95 | 0.8 | 0.98 |
| 3.59 | 203 | 0.06 | 0.06 | 0.18 | 0.15 | 0.12 | **0.13** | 0.46 | 0.65 | 0.87 |
| 2.29 | 52.93 | 0.13 | 0.49 | 0.55 | 0.54 | 0.47 | **0.5** | 0.91 | 0.88 | 0.92 |
| 3.61 | 54.8 | 0.19 | 0.19 | 0.44 | 0.45 | 0.35 | **0.43** | 0.46 | 0.79 | 0.98 |
| 1.69 | 100.8 | 0.08 | 0.4 | 0.45 | 0.45 | 0.38 | **0.44** | 0.91 | 0.87 | 0.97 |
| 2.29 | 202 | 0.06 | 0.07 | 0.23 | 0.22 | 0.19 | **0.21** | 0.5 | 0.73 | 0.87 |
| 1.38 | 204.67 | 0.06 | 0.15 | 0.34 | 0.33 | 0.27 | **0.3** | 0.58 | 0.78 | 0.95 |
| 3.58 | 99.53 | 0.16 | 0.16 | 0.35 | 0.35 | 0.27 | **0.33** | 0.44 | 0.81 | 0.9 |
| 1.66 | 54.2 | 0.1 | 0.58 | 0.61 | 0.6 | 0.47 | **0.58** | 0.92 | 0.89 | 0.99 |
| 2.27 | 97.53 | 0.11 | 0.32 | 0.37 | 0.38 | 0.32 | **0.36** | 0.88 | 0.85 | 0.97 |
| 1.72 | 205.2 | 0.06 | 0.08 | 0.25 | 0.23 | 0.19 | **0.23** | 0.52 | 0.7 | 0.94 |
| 1.4 | 52.87 | 0.12 | 0.61 | 0.69 | 0.68 | 0.6 | **0.64** | 0.85 | 0.91 | 0.92 |

Table 3: Pareto distribution for topics. Bold indicates best non-cheating result and cosine similarities are included.

| | Baseline | LSI-15 | knn-15 | lda-15 | malletlda-15 | projector-15 |
|------|----------|--------|--------|--------|--------------|--------------|
| AP | 0.21 | 0.3 | 0.28 | 0.26 | 0.26 | 0.24 |
| Cacm | 0.07 | 0.07 | 0.12 | 0.1 | 0.08 | 0.1 |
| Cisi | 0.13 | 0.13 | 0.16 | 0.15 | 0.15 | 0.15 |
| Cran | 0.18 | 0.27 | 0.29 | 0.25 | 0.25 | 0.24 |
| Med | 0.08 | 0.13 | 0.13 | 0.12 | 0.13 | 0.13 |
| Nips | 0.66 | 0.69 | 0.77 | 0.75 | 0.75 | 0.68 |

Table 4: Experiment results on real datasets. We pick the result of the best among a few parameters for each algorithm. We note that increasing the number of topics to 30 does not change the results.

words in all topics which roughly has size $v = \min(m, kw)$. For two words from the same topic, a document contains their common topic with probabiltiy $t/k$, and the probability that both words are chosen is $(1/tw)^2$. Thus the probability of cooccurence from being in the same topic is $\frac{t}{k}(1/tw)^2$. We assume that the background probability that the two words cooccur in the other case (or in the case that there are no topics) is $1/v^2$. [1]

When the background co-occurrence is much smaller than the topic cooccurrence the topic model should be easy to learn. For example, when $\frac{t}{k}1/(tw)^2 >> 1/v^2$ or when $1 >> (\frac{kw}{v})^2\frac{t}{k}$, then we should have good performance. In figure **??**, we plot this ratio against the performance of the projector algorithm and see that things degrade as this value increases.

We note that technically even under the weaker condition that $tw < v$, an algorithm could information theoretically determine the topics with enough data but the benefit over the baseline algorithm will be small.

## 8.2 Dependence on document length and number.

The basic unit for determining the co-occurence values is the number of word pairs in the corpus. The number of word pairs grows quadratically with number of words in a document and linearly with number of document. This would indicate that the performance should improve more as document length grows as compared with document number. We note that the pairs inside the document are not independent but there is an analysis that shows a nonlinear benefit.

## 8.3 An external measure.

For real world data, there is no access, of course, to the parameters used in the discussion above. We instead use the chi squared measure on the cooccurence matrix, $W$. The analysis above is really just an anlysis of the cooccurence matrix. The Chi Squared measure is defined

as

$$\sum_{i,j} \frac{(W_{i,j} - E_{i,j})^2}{E_{i,j}},$$

where $E_{i,j}$ is the expected number of cooccurrences if documents are generated from a trivial topic model with a word distribution equal to the marginal distribution of the data.

We remove all words for the corpus that occur infrequently and compute the Chi-squared measure on the corresponding co-occurrence matrix. In figure 5, we see that this measure gives us reasonable insight on the generated data into when we get good performance on the prediction task.

## 9. CONCLUSION

In this paper, we made progress toward understanding the performance of algorithms on data generated by the LDA model. Using our prediction task, we were able to find an improved algorithm for this task as well as for learning the topics in a generated topic model. This prediction task itself may be a reasonable candidate for practitioners to use to evaluate algorithms that are trying to find interesting topics.

We provided some rules of thumb for when algorithms can make use of topic structure in generated data. It is important to see if these rules of thumb extend to real world datasets. We are continuing on this aspect currently.
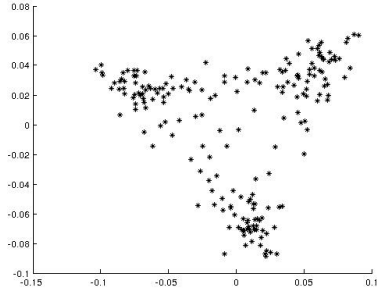
We note that we tested Hierarchical LDA on our task of predicting more words. Admittedly, we had to modify the algorithm for this task by sampling many hLDA hidden states to estimate the word distributions for the test documents. We found it to be substantially inferior to Gibbs LDA on real datasets.

Our framework is currently available on github and we are working to make it user friendly.
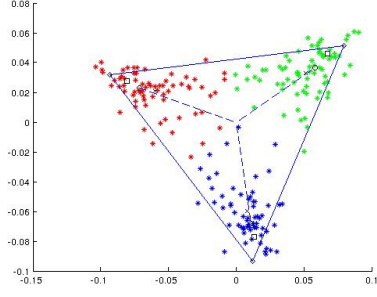
## 10. REFERENCES

[1] Animashree Anandkumar, Dean P. Foster, Daniel Hsu, Sham M. Kakade, and Yi-Kai Liu. Two svds suffice: Spectral decompositions for probabilistic topic modeling and latent dirichlet allocation. *CoRR*, abs/1204.6703, 2012.
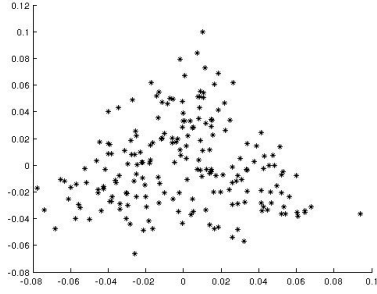
---

[1] We note that it is possible to set up topics so that the additional correlations one recieves in one topic are exactly cancelled out by other topics. This setup corresponds to coding up a parity problem in the set of topics, but it seems unlikely to arise in any reasonable topic model. Still, the "rough" calculations here fail. Indeed, we emphasize that the arguments are heuristic.

[2] Sanjeev Arora, Rong Ge, and Ankur Moitra. Learning Topic Models âĂŤ Going beyond SVD. 2012.

[3] David M. Blei. lda-c, 2003.

[4] David M. Blei. Probabilistic topic models. *Commun. ACM*, 55(4):77–84, 2012.

[5] David M. Blei, Thomas L. Griffiths, Michael I. Jordan, and Joshua B. Tenenbaum. Hierarchical topic models and the nested chinese restaurant process. In Sebastian Thrun, Lawrence K. Saul, and Bernhard Schölkopf, editors, *NIPS*. MIT Press, 2003.

[6] David M Blei, Andrew Y Ng, and Michael I Jordan. Latent Dirichlet Allocation. *The Journal of Machine Learning Research*, 3:993–1022, 2003.

[7] Thomas Hofmann. Latent semantic models for collaborative filtering. *ACM Trans. Inf. Syst.*, 22(1):89–115, 2004.

[8] Seungil Huh and Stephen E. Fienberg. Discriminative topic modeling based on manifold learning. In *Proceedings of the 16th ACM SIGKDD international conference on Knowledge discovery and data mining*, KDD '10, pages 653–662, New York, NY, USA, 2010. ACM.

[9] Adam Tauman Kalai, Ankur Moitra, and Gregory Valiant. Disentangling gaussians. *Commun. ACM*, 55(2):113–120, 2012.

[10] Percy Lang. Mlcomp. http://mlcomp.org, 2010.

[11] Andrew Kachites McCallum. Mallet: A machine learning for language toolkit. http://mallet.cs.umass.edu, 2002.

[12] Indraneel Mukherjee and David M. Blei. Relative performance guarantees for approximate inference in latent dirichlet allocation. In Daphne Koller, Dale Schuurmans, Yoshua Bengio, and Léon Bottou, editors, *NIPS*, pages 1129–1136. Curran Associates, Inc., 2008.

[13] Christos H Papadimitriou, Prabhakar Raghavan, Hisao Tamaki, and Santosh Vempala. Latent Semantic Indexing : A Probabilistic Analysis. In *PODS*, 1997.

[14] Daniel Ramage, Christopher D. Manning, and Susan Dumais. Partially labeled topic models for interpretable text mining. In *Proceedings of the 17th ACM SIGKDD international conference on Knowledge discovery and data mining*, KDD '11, pages 457–465, New York, NY, USA, 2011. ACM.

[15] Med, cran, cisi.

[16] Ap. http://www.cs.princeton.edu/ blei/lda-c/index.html.

[17] Movielens. http://www.grouplens.org/node/73.

[18] Kos, nips. http://archive.ics.uci.edu/ml/ datasets/Bag+of+Words.

[19] Limin Yao, David Mimno, and Andrew McCallum. Efficient methods for topic model inference on streaming document collections. In *Proceedings of the 15th ACM SIGKDD international conference on Knowledge discovery and data mining*, KDD '09, pages 937–946, New York, NY, USA, 2009. ACM.
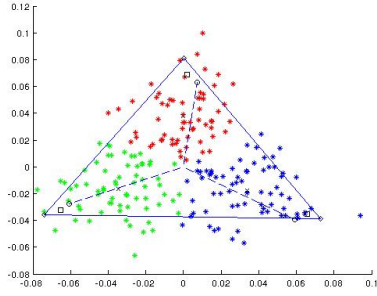
(a) $\alpha = 0.1, \beta = 0.25, k = 3$



(b) Algorithm illustration with $\delta = 0.8$



(c) $\alpha = 0.8, \beta = 0.25, k = 3$



(d) Algorithm illustration with $\delta = 0.8$

**Figure 2: Illustration of Projector. The left figures are the $V_i$'s after the SVD step. In the right figures, the black 'o's at the ends of dotted lines are the real topic, black $\times$ are the $C_i$'s before scaling, black $\diamond$ are $C_i$'s after scaling, and black $\square$ are the recovered $\hat{A}_i$'s. All points in the plot are after shifting and projected on the SVD subspace.**
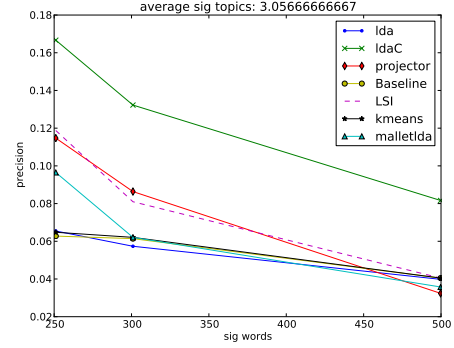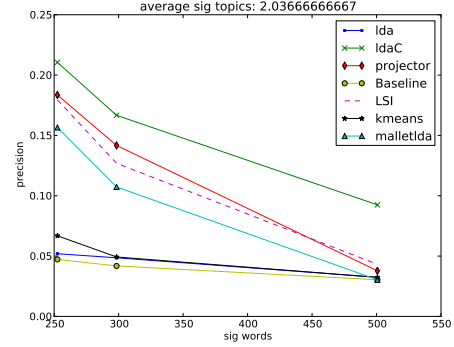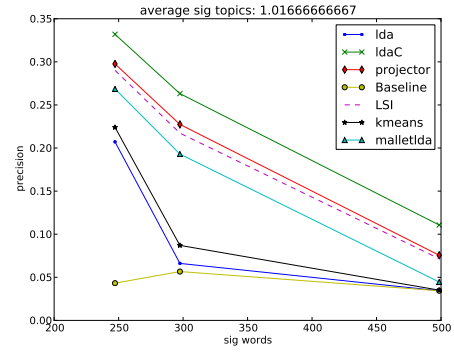








**Figure 3: Results of various algorithms on 16 generated datasets, with $k = 20, n = 1000, m = 1000, l = 75$. Each subfigure has a fixed $sig\_topic$, the plots are result of the prediction task versus $sig\_word$**

Figure 5: The x-axis is the log of the Chi Squared measure on the co-occurrence matrix for the words that occur with more than average frequency. Performance relative to baseline increases with increasing Chi squared value.
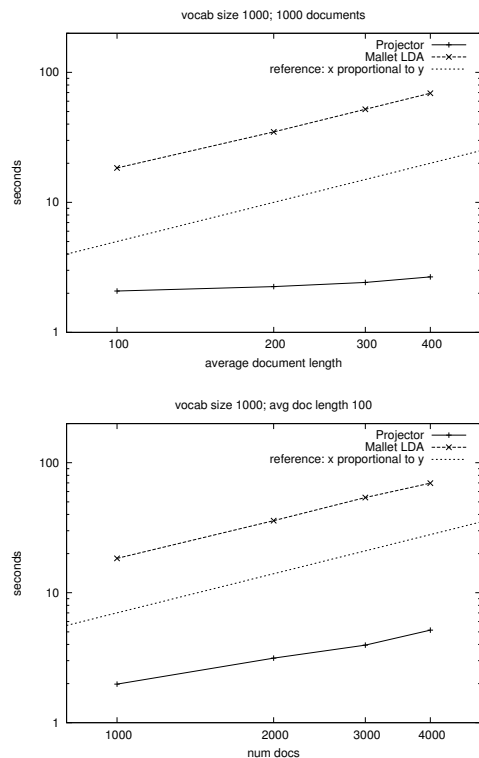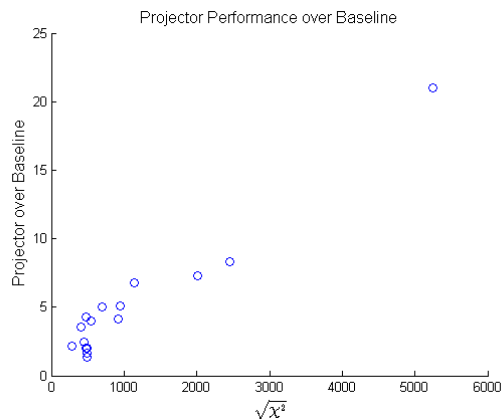


Figure 4: The runtimes of Mallet's LDA implementation and projector a plotted against the number of documents and the average document lengths in the figures above.

| Civil Rights | Economics | Politics | Stock | Legal |
|---|---|---|---|---|
| years | percent | states | stock | court |
| west | rate | government | exchange | wednesday |
| black | year | united | points | federal |
| american | prices | union | closed | judge |
| war | reported | war | tuesday | trial |
| lived | increase | minister | average | convicted |
| jr | compared | countries | tokyo | state |
| chicago | rose | military | nikkei | death |
| story | report | meeting | close | attorney |
| social | higher | soviet | share | police |
| martin | highest | leader | ↓nancial | prison |
| progress | time | world | shares | district |
| series | index | president | index | charges |
| blacks | march | administration | volume | accused |
| side | mortgages | political | timesstock | man |
| neighborhood | billion | forces | times | charged |
| dream | ↓gures | nations | million | ↓led |
| king | in'ation | foreign | percent | case |
| in'uence | retail | prime | prices | years |
| remains | august | capital | investors | ordered |

Figure 6: The top 20 words (ranked by probability) for 5 sparse topics recovered by Projector on the AP dataset. Topic names derived from top words.