

000  
001  
002  
003  
004  
005  
006  
007  
008  
009  
010  
011  
012  
013  
014  
015  
016  
017  
018  
019  
020  
021  
022  
023  
024  
025  
026  
027  
028  
029  
030  
031  
032  
033  
034  
035  
036  
037  
038  
039  
040  
041  
042  
043  
044  
045  
046  
047  
048  
049  
050  
051  
052  
053

---

# Topic Modeling via Nonnegative Matrix Factorization on Probability Simplex

---

Anonymous Author(s)

Affiliation

Address

email

## Abstract

One important goal of document modeling is to extract a set of informative topics from a text corpus and produce a reduced representation of each document. In this paper, we propose a novel algorithm for this task based on nonnegative matrix factorization on a probability simplex. We further extend our algorithm by removing global and generic information to produce more diverse and specific topics. In contrast to other matrix factorization methods, such as latent semantic indexing by singular value decomposition, our model has a solid statistical foundation and is based on a generative model for text corpus. In contrast to purely probabilistic approach, such as probabilistic latent semantic indexing (pLSI) solved by Expectation-Maximization, our method is based on efficient block coordinate descent optimization. Experiments demonstrate that the new method generates more meaningful and diverse topics compared with pLSI and LDA with faster convergence behavior.

## 1 Introduction

There has been increasing interests in topic modeling to analyze a large corpus of documents and distill a set of meaningful topics. Current methods can be roughly divided into two categories, 1) matrix decomposition methods and 2) probabilistic methods. The most prominent examples in the first category is latent semantic indexing (LSI) [4] and nonnegative matrix factorization (NMF) based methods [9, 1]. However, they often ignore probabilistic constraints imposed by the data generating process and lack solid statistical foundations.

The most well-known probabilistic methods are probabilistic latent semantic indexing (pLSI) [6, 7] and latent Dirichlet allocation (LDA) [3]. The parameters can be estimated with variational inference or sampling algorithms [3, 5, 2], which suffers from run time efficiency problem.

In this paper, we propose a fast nonnegative matrix factorization algorithm that combines the best of both worlds. It is efficient and respects the constraints imposed by the probabilistic generative models as in pLSI. In addition, we extend our algorithm to remove global and generic information in a principled way and generate more diverse topics in cases of small number of topics.

## 2 NMF Formulation with Probability Constraints

Given  $N$  documents with a  $M$ -sized vocabulary, we create the empirical conditional word distribution matrix  $\mathbf{A} = \hat{P}(w|d)$ , with each of its column summing to one. We approximate the nonnegative matrix  $\mathbf{A} \in \mathbb{R}^{M \times N}$  by a product of two lower-rank nonnegative matrices,  $\mathbf{W} \in \mathbb{R}^{M \times K}$  and  $\mathbf{H} \in \mathbb{R}^{K \times N}$ , where  $K$  is the number of topics. We interpret  $\mathbf{W}$  as topic distributions and  $\mathbf{H}$  as per document mixture proportions. Thus each column is required to have unit  $L_1$  norm, making it a probability distribution.

054 Due to difficulties in optimizing with constraints, we transform the objective function into a regular-  
 055 ized form

$$056 \min_{\mathbf{W} \geq 0, \mathbf{H} \geq 0} \|\mathbf{A} - \mathbf{W}\mathbf{H}\|_F^2 + \alpha \|\mathbf{1}_M^\top \mathbf{W} - \mathbf{1}_K^\top\|_F^2 + \beta \|\mathbf{1}_K^\top \mathbf{H} - \mathbf{1}_N^\top\|_F^2, \quad (1)$$

057 where  $\mathbf{1}_K$  is an all-ones vector of size  $K$ , and  $\alpha$  and  $\beta$  are regularization parameters.

058 Our formulation is closely related to pLSI. To see this, recall that the columns of  $\mathbf{W}$  are the topic  
 059 vectors  $P(w|z)$ , and the columns of  $\mathbf{H}$  are the document specific mixture proportions  $P(z|d)$ . The  
 060 matrix product  $\mathbf{W}\mathbf{H}$  is the conditional distribution of words per document  $P(w|d)$ . By minimiz-  
 061 ing the Frobenius norm  $\|\mathbf{A} - \mathbf{W}\mathbf{H}\|_F^2$ , we obtain factors  $\mathbf{W}$  and  $\mathbf{H}$  that approximate the empirical  
 062 conditional word distribution  $\mathbf{A}$ . In contrast, pLSI maximizes the likelihood of data under the model,  
 063 which is equivalent to minimizing the Kullback-Leibler (KL) divergence of the empirical joint dis-  
 064 tribution  $\hat{P}(w, d)$  and the model  $P(w, d)$ .

065 Compared with pLSI, our method minimizes  $L_2$  distance instead of KL divergence for ease of op-  
 066 timization. Also, we approximate the conditional word distribution instead of the joint distribution.  
 067 In fact, the extra factor  $P(d)$  in pLSI is not useful in discovering topics [3].

068 One way to solve (1) is to use the block coordinate descent framework [8, 10], alternating between  
 069 solving for  $\mathbf{H}$  (with fixed  $\mathbf{W}$ ) and solving for  $\mathbf{W}$  (with fixed  $\mathbf{H}$ ). However, due to difficulty in  
 070 solving for  $\mathbf{W}$  we introduce an auxiliary variable  $\mathbf{Z}$ , which serves as a proxy of  $\mathbf{W}$  and hence  
 071 decouples the non-negativity and the unit  $L_1$  norm constraints. Finally, we solve the following  
 072 optimization problem

$$073 \min_{\mathbf{W} \geq 0, \mathbf{H} \geq 0} \|\mathbf{A} - \mathbf{W}\mathbf{H}\|_F^2 + \gamma \|\mathbf{W} - \mathbf{Z}\|_F^2 + \alpha \|\mathbf{1}_M^\top \mathbf{Z} - \mathbf{1}_K^\top\|_F^2 + \beta \|\mathbf{1}_K^\top \mathbf{H} - \mathbf{1}_N^\top\|_F^2, \quad (2)$$

074 where  $\gamma$  is a regularization parameter.

### 075 3 Three-Block Coordinate Descent Algorithm

076 The optimization problem (2) can be solved by alternating updates for  $\mathbf{H}$ ,  $\mathbf{W}$  and  $\mathbf{Z}$  through the  
 077 following three subproblems:

$$078 \mathbf{H} \leftarrow \operatorname{argmin}_{\mathbf{H} \geq 0} \left\| \begin{pmatrix} \mathbf{A} \\ \sqrt{\beta} \mathbf{1}_N^\top \end{pmatrix} - \begin{pmatrix} \mathbf{W} \\ \sqrt{\beta} \mathbf{1}_K^\top \end{pmatrix} \mathbf{H} \right\|_F^2 \quad (3)$$

$$079 \mathbf{W} \leftarrow \operatorname{argmin}_{\mathbf{W} \geq 0} \left\| \begin{pmatrix} \mathbf{A}^\top \\ \sqrt{\gamma} \mathbf{Z}^\top \end{pmatrix} - \begin{pmatrix} \mathbf{H}^\top \\ \sqrt{\gamma} \mathbf{I}_K^\top \end{pmatrix} \mathbf{W}^\top \right\|_F^2 \quad (4)$$

$$080 \mathbf{Z} \leftarrow \operatorname{argmin} \left\| \begin{pmatrix} \sqrt{\gamma} \mathbf{W} \\ \sqrt{\alpha} \mathbf{1}_K^\top \end{pmatrix} - \begin{pmatrix} \sqrt{\gamma} \mathbf{I}_M \\ \sqrt{\alpha} \mathbf{1}_M^\top \end{pmatrix} \mathbf{Z} \right\|_F^2 \quad (5)$$

081 We use ANLS/BPP [8] to solve the subproblems (3) and (4).

082 The least squares problem (5) has a special structure, and we can apply the Sherman-Morrison  
 083 formula to obtain a direct solution:  $\mathbf{Z} \leftarrow \mathbf{W} - \frac{\alpha}{\gamma + M\alpha} \mathbf{1} \mathbf{1}^\top \mathbf{W} + \left( \frac{\alpha}{\gamma} - \frac{M\alpha^2}{\gamma(\gamma + M\alpha)} \right) \mathbf{1} \mathbf{1}^\top$ , where the  
 084 main cost of computing  $\mathbf{1} \mathbf{1}^\top \mathbf{W}$  can be carried out very efficiently.

085 Theoretically, we need to use large regularization parameters so that the the unit  $L_1$  constraints on the  
 086 columns of  $\mathbf{W}$  and  $\mathbf{H}$  are satisfied. However, large regularization parameters put too much emphasis  
 087 on the constraints during the early stage of iterations. Therefore, we start with reasonably small  
 088 values and then adaptively increase them as iteration progresses [11]. This approach is reminiscent  
 089 of simulated annealing where one starts with high temperature parameters and gradually decreases  
 090 the temperature to zero. The overall algorithm is summarized in Algorithm 1, and we name the  
 091 algorithm t-NMF for topic NMF.

### 092 4 Shifted Non-negative Matrix Factorization

093 Algorithm 1 can also be extended for hierarchical topic modeling, where at each level, we only find  
 094 a small number of topics for a partition of the corpus, and recursively partition the documents and

108  
109  
110  
111  
112  
113  
114  
115  
116  
117  
118  
119  
120  
121  
122  
123  
124  
125  
126  
127  
128  
129  
130  
131  
132  
133  
134  
135  
136  
137  
138  
139  
140  
141  
142  
143  
144  
145  
146  
147  
148  
149  
150  
151  
152  
153  
154  
155  
156  
157  
158  
159  
160  
161

---

**Algorithm 1** tNMF: Nonnegative Matrix Factorization on Probability Simplex

---

1: **input:** Empirical conditional distribution  $\mathbf{A} \in \mathbb{R}_+^{M \times N}$ , the number of topics  $K$ , regularization parameter  $\alpha > 0, \beta > 0, \gamma > 0$ .  
2: **output:** Topics  $\mathbf{W} \in \mathbb{R}_+^{M \times K}$  and mixing proportions  $\mathbf{H} \in \mathbb{R}_+^{K \times N}$ .  
3: **repeat**  
4:   **if** Difference in objective from one iteration to the next is below some threshold and  $(\alpha, \beta, \gamma)$  has not been increased for a certain number of iterations **then**  
5:     Increase the values for  $\alpha, \beta$ , and  $\gamma$ .  
6:   **end if**  
7:    $\mathbf{H} \leftarrow \operatorname{argmin}_{\mathbf{H} \geq 0} \left\| \begin{pmatrix} \mathbf{A} \\ \sqrt{\beta} \mathbf{1}_N^\top \end{pmatrix} - \begin{pmatrix} \mathbf{W} \\ \sqrt{\beta} \mathbf{1}_K^\top \end{pmatrix} \mathbf{H} \right\|_F^2$   
8:    $\mathbf{W} \leftarrow \operatorname{argmin}_{\mathbf{W} \geq 0} \left\| \begin{pmatrix} \mathbf{A}^\top \\ \sqrt{\gamma} \mathbf{Z}^\top \end{pmatrix} - \begin{pmatrix} \mathbf{H}^\top \\ \sqrt{\gamma} \mathbf{1}_K^\top \end{pmatrix} \mathbf{W}^\top \right\|_F^2$   
9:    $\mathbf{Z} \leftarrow \mathbf{W} - \frac{\alpha}{\gamma + M\alpha} \mathbf{1}\mathbf{1}^\top \mathbf{W} + \left( \frac{\alpha}{\gamma} - \frac{M\alpha^2}{\gamma(\gamma + M\alpha)} \right) \mathbf{1}\mathbf{1}^\top$   
10: **until** stopping criterion is reached

---

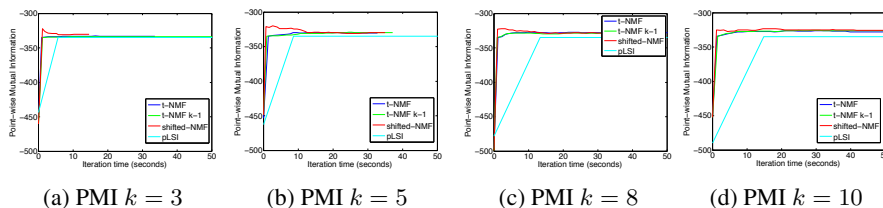


Figure 1: Convergence speed comparison using Point-wise Mutual Information (PMI).

apply NMF. The problem, however, is that the obtained topics are similar to each other when  $K$  is small, often consisting of repeated terms that are global to the corpus. In light of this situation, we propose shifted NMF by explicitly picking out a global topic in hope for more diverse topics.

Specifically, we fix the first topic vector as the average of  $\mathbf{A}$ , *i.e.*,  $\mathbf{W}_{:,1} = \frac{1}{N} \sum_i \mathbf{A}_{:,i}$ . The corresponding mixture proportion, *i.e.*, the first row of  $\mathbf{H}$ , is allowed to change so that it optimizes the objective function. The intuition is that by fixing a global topic, the remaining topics are tilted to explain more specific contents in the corpus, thus arriving at more diverse topics. In addition, the proportion of the global topic is not the same for every document, and it is optimized together with other mixture proportions. The algorithm is similar to t-NMF except that  $\mathbf{W}$  and  $\mathbf{Z}$  are now both one dimension smaller than their counterparts in t-NMF.

## 5 Experiments

We did experiments on the NIPS dataset, which contains 1739 documents from 2000 to 2012 proceedings, with vocabulary size 13648. We used Point-wise Mutual Information (PMI) [1] to evaluate topic quality:  $PMI = \frac{1}{K} \sum_i \sum_{s,t \in \mathcal{T}_i, s < t} \log \frac{D_{st} + \epsilon}{D_s D_t}$ , where  $D_{st}$  is the number of documents in which keywords  $s$  and  $t$  co-occur,  $D_s$  is the number of documents keyword  $s$  occurs, and  $\epsilon$  is a small constant for smoothing. We also computed Average KL divergence (AKL) to measure the distinctiveness of the topics:  $AKL = \frac{2}{K(K-1)} \sum_{i < j} \left( \sum_t W_{ti} \log \frac{W_{ti}}{W_{tj}} \right)$ .

We compared with a Matlab implementation of pLSI in terms of convergence speed and present the comparison in Figures 1 and 2. Our NMF-based algorithms converge quickly in terms of both PMI and average KL divergence. The shifted-variant achieves the fastest convergence when  $k$  is small.

We summarize the topic quality in Table 1 and list the top 10 keywords from each topic in Table 2. Shifted NMF is able to produce more diverse topics by allowing a global topic to explain the generic content in the corpus. t-NMF and LDA discover topics with more repetition of global keywords such as “data” and “model”. In contrast, pLSI performs poorly by selecting non-descriptive tokens, and it is not robust to noise.

## References

[1] S. Arora, R. Ge, and A. Moitra. Learning topic models—going beyond svd. In *Foundations of Computer Science (FOCS), 2012 IEEE 53rd Annual Symposium on*, pages 1–10. IEEE, 2012.

162  
163  
164  
165  
166  
167  
168  
169  
170  
171  
172  
173  
174  
175  
176  
177  
178  
179  
180  
181  
182  
183  
184  
185  
186  
187  
188  
189  
190  
191  
192  
193  
194  
195  
196  
197  
198  
199  
200  
201  
202  
203  
204  
205  
206  
207  
208  
209  
210  
211  
212  
213  
214  
215

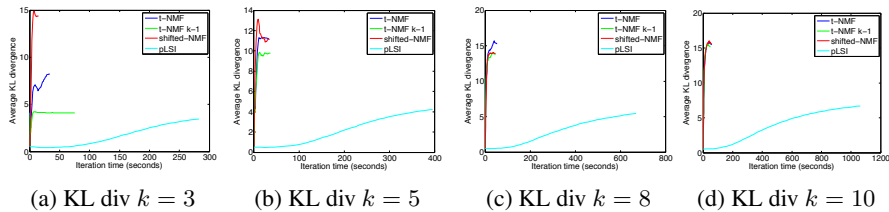


Figure 2: Convergence speed comparison using average KL divergence.

Table 1: Performance comparison of shifted NMF, tNMF, LDA and pLSI on NIPS dataset. AKL is Average KL divergence between pair-wise topic vectors, and PMI is Point-wise Mutual Information.

Method	$K$	AKL	PMI	$K$	AKL	PMI
shifted NMF	3	<b>15.32</b>	<b>-330.99</b>	8	14.54	-328.82
t-NMF	3	7.99	-332.07	8	<b>14.61</b>	<b>-327.14</b>
LDA	3	4.26	-332.36	8	5.72	-329.65
pLSI	3	4.01	-333.32	8	5.87	-330.62
shifted NMF	5	11.30	<b>-329.62</b>	10	15.18	-328.02
tt-NMF	5	<b>11.33</b>	-331.49	10	<b>16.23</b>	<b>-327.78</b>
LDA	5	4.93	-331.03	10	5.90	-327.85
pLSI	5	4.88	-332.71	10	6.72	-328.56

[2] D. M. Blei and M. I. Jordan. Variational inference for dirichlet process mixtures. *Bayesian analysis*, 2006.

[3] D. M. Blei, A. Ng, and M. Jordan. Latent dirichlet allocation. *JMLR*, 3:993–1022, 2003.

[4] S. C. Deerwester, S. T. Dumais, T. K. Landauer, G. W. Furnas, and R. A. Harshman. Indexing by latent semantic analysis. *JASIS*, 1990.

[5] T. L. Griffiths and M. Steyvers. Finding scientific topics. *PNAS*, 101(suppl. 1):5228–5235, 2004.

[6] T. Hofmann. Probabilistic latent semantic indexing. In *Proceedings of the 22nd annual international ACM SIGIR conference on Research and development in information retrieval*, pages 50–57. ACM, 1999.

[7] T. Hofmann. Unsupervised learning by probabilistic latent semantic analysis. *Mach. Learn.*, 2001.

[8] J. Kim and H. Park. Fast nonnegative matrix factorization: An active-set-like method and comparisons. *SIAM J. Sci. Comput.*, 33(6):3261–3281, Nov. 2011.

[9] D. Kuang, C. Ding, and H. Park. Symmetric nonnegative matrix factorization for graph clustering. In *SIAM International Conference on Data Mining (SDM '12)*, pages 106–117, 2012.

[10] L. Li, G. Lebanon, and H. Park. Fast bregman divergence nmf using taylor expansion and coordinate descent. In *KDD 2012*, page 307315. ACM, 2012.

[11] J. Nocedal and S. Wright. *Numerical Optimization*. Springer, 1999.

Table 2: Top key words discovered by shifted NMF, tNMF, LDA and pLSI for different  $K$  on NIPS dataset.

Method	$K$	Topics
shifted NMF	3	network, learning, model, neural, input, function, figure, data, time, networks data, model, algorithm, set, function, learning, models, training, distribution, number model, neurons, figure, input, time, cells, neuron, cell, neural, visual
t-NMF	3	model, figure, time, system, data, neurons, models, cells, input, visual learning, function, algorithm, data, set, error, training, state, problem, number network, neural, networks, input, output, units, training, hidden, layer, weights
LDA	3	model, figure, time, neurons, input, neuron, system, neural, visual, cells network, training, neural, input, networks, units, output, set, learning, hidden learning, function, model, data, algorithm, state, set, error, linear, problem
pLSI	3	kwon, unreliable, cart, shades, finite, minimising, awi, stark, proliferation, thalamus narrowing, martin, cropped, rts, englewood, stochasticity, aij, automatically, instructive, rgb closeness, concentric, adds, sooner, pairwise, dispersed, measurable, medicine, ile, sea