

Nonnegative Matrix Factorization and Probabilistic Latent Semantic Indexing: Equivalence, Chi-square Statistic, and a Hybrid Method

Chris Ding^a, Tao Li^b and Wei Peng^b

^a Lawrence Berkeley National Laboratory, Berkeley, CA 94720

^b School of Computer Science, Florida International University, Miami, FL 33199

Abstract

Non-negative Matrix Factorization (NMF) and Probabilistic Latent Semantic Indexing (PLSI) have been successfully applied to document clustering recently. In this paper, we show that PLSI and NMF optimize the same objective function, although PLSI and NMF are different algorithms as verified by experiments. This provides a theoretical basis for a new hybrid method that runs PLSI and NMF alternatively, each jumping out of local minima of the other method successively, thus achieving better final solution. Extensive experiments on 5 real-life datasets show relations between NMF and PLSI, and indicate the hybrid method lead to significant improvements over NMF-only or PLSI-only methods. We also show that at first order approximation, NMF is identical to χ^2 -statistic.

Introduction

Document clustering has been widely used as a fundamental and effective tool for efficient document organization, summarization, navigation and retrieval of large amount of documents. Generally document clustering problems are determined by the three basic tightly-coupled components: a) the (physical) representation of the given data set; b) The criterion/objective function which the clustering solutions should aim to optimize; c) The optimization procedure (Li 2005).

Among clustering methods, the K-means algorithm has been the most popularly used. A recent development is the Probabilistic Latent Semantic Indexing (PLSI). PLSI is a unsupervised learning method based on statistical latent class models and has been successfully applied to document clustering (Hofmann 1999). (PLSI is further developed into a more comprehensive Latent Dirichlet Allocation model (Blei, Ng, & Jordan 2003).)

Nonnegative Matrix Factorization (NMF) is another recent development for document clustering. Initial work on NMF (Lee & Seung 1999; 2001) emphasizes the contain coherent parts of the original data (images). Later work (Xu, Liu, & Gong 2003; Pauca *et al.* 2004) show the usefulness of NMF for clustering with in experiments on documents

Copyright © 2006, American Association for Artificial Intelligence (www.aaai.org). All rights reserved.

collections, and a recent theoretical analysis (Ding, He, & Simon 2005) shows the equivalence between NMF and K -means / spectral clustering.

Despite significant research on both NMF and PLSI, few attempts have been made to establish the connections between them while highlighting their differences in the clustering framework. Gaussier and Goutte (Gaussier & Goutte 2005) made the first connection between NMF and PLSI, by showing that the local fixed point solutions of the iterative procedures of PLSI and NMF are the same. Their proof is, however, incorrect. NMF and PLSI are different algorithms. They converge to different solutions while starting from the same initial condition, as verified by experiments (see later sections).

In this paper, we first show that both NMF and PLSI optimize the same objective function. This fundamental fact and the L_1 normalization NMF ensures that NMF and PLSI are equivalent.

Second, we show, by an example and extensive experiments, that NMF and PLSI are different algorithms and they converge to different local minima. This leads to a new insight: NMF and PLSI are different algorithms for optimizing the same objective function.

Third, we give a detailed analysis about the NMF and PLSI solutions. They are local minima of the same landscape in a very high dimensional space. We show that PLSI can jump out of the local minima where NMF converges to and vice versa. Based on this, we further propose a hybrid algorithm to run NMF and PLSI alternatively to jump out a series of local minima and finally reach to a much better minimum. Extensive experiments show this hybrid algorithm improves significantly over the standard NMF-only or PLSI-only algorithms.

Data Representations of NMF and PLSI

Suppose we have n documents and m words (terms). Let $F = (F_{ij})$ be the word-to-document matrix: $F_{ij} = F(w_i, d_j)$ is the frequency of word w_i in document d_j .

In this paper, we re-scale the term frequency F_{ij} by $F_{ij} \leftarrow$

F_{ij}/T_w , where $T_w = \sum_{ij} F_{ij}$ is the total number of words. With this stochastic normalization, $\sum_{ij} F_{ij} = 1$. The joint occurrence probability $p(w_i, d_j) = F_{ij}$.

The general form of NMF is

$$F = CH^T, \quad (1)$$

where the matrices $C = (C_{ik}), H = (H_{jk})$ are nonnegative matrices. They are determined by minimizing

$$J_{\text{NMF}} = \sum_{i=1}^m \sum_{j=1}^n F_{ij} \log \frac{F_{ij}}{(CH^T)_{ij}} - F_{ij} + (CH^T)_{ij} \quad (2)$$

PLSI maximize the likelihood

$$\max J_{\text{PLSI}}, \quad J_{\text{PLSI}} = \sum_{i=1}^m \sum_{j=1}^n F_{ij} \log P(w_i, d_j) \quad (3)$$

where $P(w_i, d_j)$ is the factorized (i.e., parameterized or approximated) joint occurrence probability

$$P(w_i, d_j) = \sum_k p(w_i|z_k)p(z_k)p(d_j|z_k), \quad (4)$$

where the probability factors follow the normalization of probabilities

$$\sum_{i=1}^m p(w_i|z_k) = 1, \sum_{j=1}^n p(d_j|z_k) = 1, \sum_{k=1}^K p(z_k) = 1. \quad (5)$$

Equivalence of NMF and PLSI

In this section, we present our main results:

Theorem 1. PLSI and NMF are equivalent.

The proof is better described by the following

Proposition 1. The objective function of PLSI is identical to the objective function of NMF, i.e.,

$$\max J_{\text{PLSI}} \iff \min J_{\text{NMF}} \quad (6)$$

Proposition 2. Column normalized NMF of Eq.(1) is equivalent to the probability factorization of Eq.(4), i.e., $(CH^T)_{ij} = P(w_i, d_j)$.

Proof of Theorem 1: By Proposition 2, NMF (with L_1 -normalization, see §3) is identical to PLSI factorization. By Proposition 1, they minimize the same objective function. Therefore, NMF is identical to PLSI. \square

We proceed to prove Proposition 1 in this section.

Proof of Proposition 1:

First, we note that the PLSI optimization Eq.(3) can be written as $\min \sum_{i=1}^m \sum_{j=1}^n -F_{ij} \log P(w_i, d_j)$. Adding a constant, $\sum_{i=1}^m \sum_{j=1}^n F_{ij} \log F_{ij}$, PLSI is equivalent to solve

$$\min \sum_{i=1}^m \sum_{j=1}^n F_{ij} \log \frac{F_{ij}}{P(w_i, d_j)}.$$

Now since

$$\sum_{i=1}^m \sum_{j=1}^n [P(w_i, d_j) - F_{ij}] = [1 - 1] = 0,$$

we can add this constant to the summation; PLSI is equivalent to minimize

$$\sum_{i=1}^m \sum_{j=1}^n F_{ij} \log \frac{F_{ij}}{P(w_i, d_j)} - F_{ij} + P(w_i, d_j) \quad (7)$$

This is precisely the objective function for NMF. \square

NMF and χ^2 -statistic.

J_{NMF} of Eq.(2) has a somewhat complicated expression. It is related to the Kullback-Leibler divergence. We give a better understanding by relating it to the familiar χ^2 test in statistics. Assume $\frac{|(CH^T)_{ij} - F_{ij}|}{F_{ij}}$ is small. We can write

$$J_{\text{NMF}} = \sum_{i=1}^m \sum_{j=1}^n \frac{[(CH^T)_{ij} - F_{ij}]^2}{2F_{ij}} - \frac{[(CH^T)_{ij} - F_{ij}]^3}{3F_{ij}^2} + \dots \quad (8)$$

This is obtained by setting $\delta_{ij} = (CH^T)_{ij} - F_{ij}$, $z = \delta_{ij}/F_{ij}$, and $\log(1+z) = z - z^2/2 + z^3/3 \dots$; then the ij -th term in J_{NMF} becomes

$$\delta_{ij} - F_{ij} \log \left(1 + \frac{\delta_{ij}}{F_{ij}} \right) = \frac{1}{2} \frac{\delta_{ij}^2}{F_{ij}} - \frac{1}{3} \frac{\delta_{ij}^3}{F_{ij}^2} + \dots$$

Clearly, the first term in J_{NMF} is the χ^2 statistic,

$$\chi^2 = \sum_{i=1}^m \sum_{j=1}^n \frac{[(CH^T)_{ij} - F_{ij}]^2}{F_{ij}}, \quad (9)$$

since F_{ij} is the data and $(CH^T)_{ij}$ is the model fit to it. Therefore, to first order approximation, NMF objective function is a χ^2 statistic. As a consequence, we can associate a confidence to NMF factorization.

The χ^2 form of NMF naturally relates to another NMF cost function, i.e., the sum of squared errors

$$J'_{\text{NMF}} = \sum_{i=1}^m \sum_{j=1}^n [(CH^T)_{ij} - F_{ij}]^2. \quad (10)$$

A comprehensive comparison among J_{NMF}, χ^2 and J'_{NMF} forms of NMF would be useful, but goes beyond the scope of this paper.

Normalizations of NMF

For any given NMF solution (C, H) , there exist a large number of matrices (A, B) such that $AB^T = I$, $CA \geq 0$, $HB \geq 0$. Thus (CA, HB) is also a solution with the same cost function value. Normalization is a way to eliminate this uncertainty. We mostly consider the normalization of columns of C, H . Specifically, let the columns be expressed explicitly, $C = (\mathbf{c}_1, \dots, \mathbf{c}_k)$, $H = (\mathbf{h}_1, \dots, \mathbf{h}_k)$.

¹In this column form, for clustering interpretation (Ding, He, & Simon 2005), \mathbf{c}_k is the centroid for k -th cluster, while \mathbf{h}_k is the posterior probability for k -th cluster. For hard clustering, on each row of H , set the largest element to 1 and the rest to 0.

We consider column normalization. Let the normalized columns be $\tilde{C} = (\tilde{\mathbf{c}}_1, \dots, \tilde{\mathbf{c}}_k)$, $\tilde{H} = (\tilde{\mathbf{h}}_1, \dots, \tilde{\mathbf{h}}_k)$. With this normalization, we can write

$$CH^T = \tilde{C}\tilde{S}\tilde{H}^T, \quad (11)$$

where

$$\tilde{C} = CD_C^{-1}, \quad \tilde{H} = HD_H^{-1}, \quad S = D_C D_H. \quad (12)$$

D_C, D_H are diagonal matrices. Depending on the normalizations in the Hilbert space, the L_p -normalization, the diagonal elements are given by

$$(D_C)_{kk} = \|\tilde{\mathbf{c}}_k\|_p, \quad (D_H)_{kk} = \|\tilde{\mathbf{h}}_k\|_p.$$

For the standard Euclidean distance normalization, i.e., the L_2 -norm

$$\|\tilde{\mathbf{c}}_k\|_2 = 1, \quad \|\tilde{\mathbf{h}}_k\|_2 = 1, \quad (13)$$

This is the same as in singular value decomposition where the non-negativity constraint is ignored.

For probabilistic formulations, such as PLSI, we use the L_1 norm.

$$\|\tilde{\mathbf{c}}_k\|_1 = 1, \quad \|\tilde{\mathbf{h}}_k\|_1 = 1, \quad (14)$$

Due to the non-negativity, these are just the condition that columns sums to 1. D_C contains the column sums of C and D_H contains the column sums of H .

With these clarification, we now prove Proposition 2.

Proof of Proposition 2:

Using L_1 -norm, we obviously have

$$\sum_{i=1}^m \tilde{C}_{ik} = 1, \quad \sum_{j=1}^n \tilde{H}_{jk} = 1, \quad \sum_{k=1}^K S_{kk} = 1,$$

where the last equality is proved as

$$1 = \sum_{ij} F_{ij} = \sum_{i=1}^m \sum_{k=1}^K \sum_{j=1}^n \tilde{C}_{ik} S_{kk} \tilde{H}_{jk} = \sum_{k=1}^K S_{kk}.$$

These can be seen as equivalent to the normalization of probabilities of Eq.(5). Therefore, $\tilde{C}_{ik} = p(w_i|z_k)$, $\tilde{H}_{jk} = p(d_j|z_k)$ and $S_{kk} = p(z_k)$. Thus $F = CH^T = \tilde{C}\tilde{S}\tilde{H}^T$ factorization with L_1 -normalization is identical to PLSI factorization \square

An Illustration of NMF/PLSI Difference

In (Gaussier & Goutte 2005), the authors gave a proof that NMF and PLSI converge to the same solution (fixed point). We believe their proof is not accurate. In their proof, they make use of an relation $P(c|w_i, d_j)^{(t)} = P(w_i, c)^{(t)} P(d_j|c)^{(t)} / P(w_i, d_j)^{(t)}$ which we were unable to reproduce.

Although they optimize the same objective function as shown above, NMF and PLSI are different computational algorithms. This fact is obvious from experiments. In all of our extensive experiments, starting with the same initial

starting C_0, H_0 , NMF and PLSI always converge to different solutions. Here we give an illustration. The input data matrix is

$$X = \begin{pmatrix} .048 & .042 & .047 & .024 & .029 & .026 \\ .035 & .040 & .045 & .016 & .023 & .029 \\ .031 & .019 & .031 & .040 & .045 & .042 \\ .027 & .023 & .031 & .032 & .039 & .045 \\ .047 & .043 & .035 & .026 & .021 & .019 \end{pmatrix}$$

The initial C_0, S_0, H_0 are

$$C_0 S_0 H_0^T = \begin{pmatrix} .24 & .20 \\ .02 & .27 \\ .31 & .16 \\ .07 & .26 \\ .36 & .11 \end{pmatrix} \begin{pmatrix} .34 & 0 \\ 0 & .66 \end{pmatrix} \begin{pmatrix} .18 & .19 \\ .15 & .18 \\ .15 & .21 \\ .18 & .12 \\ .18 & .14 \\ .16 & .16 \end{pmatrix}^T$$

Running the NMF algorithm, the converged solution are

$$\tilde{C}\tilde{S}\tilde{H}^T = \begin{pmatrix} .33 & .14 \\ .29 & .12 \\ .02 & .33 \\ .05 & .29 \\ .32 & .11 \end{pmatrix} \begin{pmatrix} .39 & 0 \\ 0 & .61 \end{pmatrix} \begin{pmatrix} .27 & .14 \\ .28 & .09 \\ .25 & .15 \\ .07 & .18 \\ .06 & .22 \\ .06 & .23 \end{pmatrix}^T$$

Running the PLSI algorithm, the converged solution are

$$C S H^T = \begin{pmatrix} .12 & .31 \\ .10 & .28 \\ .38 & .04 \\ .33 & .07 \\ .08 & .31 \end{pmatrix} \begin{pmatrix} .50 & 0 \\ 0 & .50 \end{pmatrix} \begin{pmatrix} .13 & .25 \\ .09 & .25 \\ .14 & .24 \\ .19 & .09 \\ .22 & .09 \\ .23 & .09 \end{pmatrix}^T$$

One can see NMF solution differs from PLSI solution significantly. Our example shows that starting at the same point in the multi-dimensional space, NMF and PLSI converge to *different* local minima.

However, it is interesting and important to note that the clustering results embedded in the solutions of NMF and PLSI are identical by an examination of H (see footnote 2): the first 3 data points (columns) belong to one cluster, and the rest 3 points belong to another cluster. This result is the same as the K-means clustering. More generally, we introduce a clustering matrix $R = (r_{ij})$, where $r_{ij} = 1$ if $\mathbf{x}_i, \mathbf{x}_j$ belong to the same cluster; $r_{ij} = 0$ otherwise. Thus the clustering results can be expressed as

$$R_{\text{NMF}} = R_{\text{PLSI}} = \begin{pmatrix} 1 & 1 & 1 & 0 & 0 & 0 \\ 1 & 1 & 1 & 0 & 0 & 0 \\ 1 & 1 & 1 & 0 & 0 & 0 \\ 0 & 0 & 0 & 1 & 1 & 1 \\ 0 & 0 & 0 & 1 & 1 & 1 \\ 0 & 0 & 0 & 1 & 1 & 1 \end{pmatrix} \quad (15)$$

Comparison between NMF and PLSI

In this section, we compare the clustering performance of each methods on 5 real-life datasets.

Datasets	# documents	# class
CSTR	476	4
WebKB	4199	4
Log	1367	9
Reuters	2900	10
WebAce	2340	20

Table 1: Dataset Descriptions.

Datasets

We use 5 datasets in our experiments, most of which are frequently used in the information retrieval research. Table 1 summarizes the characteristics of the datasets.

CSTR The dataset contains the abstracts of technical reports (TRs) published in the Computer Science Department at the University of Rochester between 1991 and 2002. The dataset has 476 abstracts, which are divided into four research areas: Natural Language Processing(NLP), Robotics/Vision, Systems, and Theory.

WebKB The dataset contains webpages gathered from university computer science departments. There are about 4199 documents and they are divided into 4 categories: student, faculty, course, project.

Log This dataset contains 1367 log text messages which are grouped into 9 categories, i.e., *configuration, connection, create, dependency, other, report, request, start, and stop*.

Reuters The Reuters-21578 Text Categorization Test collection contains documents collected from the Reuters newswire in 1987. In our experiments, we use a subset of the data collection which includes the 10 most frequent categories among the 135 topics and has about 2900 documents.

WebAce The dataset is from WebACE project (Han *et al.* 1998). It contains 2340 documents consisting news articles from Reuters new service via the Web in October 1997. These documents are divided into 20 classes.

To pre-process the datasets, we remove the stop words using a standard stop list. All HTML tags are skipped and all header fields except subject and organization of the posted articles are ignored. In all our experiments, we first select the top 1000 words by mutual information with class labels.

Evaluation Measures

The above document datasets are standard labeled corpora widely used in the information retrieval literature. We view the labels of the datasets as the objective knowledge on the structure of the datasets. To measure the clustering performance, we use accuracy, entropy, purity and Adjusted Rand Index (ARI) as our performance measures. We expect these measures would provide us with good insights.

Accuracy discovers the one-to-one relationship between clusters and classes and measures the extent to which each cluster contained data points from the corresponding class. It sums up the whole matching degree between all pair class-clusters. Accuracy can be represented as:

$$Accuracy = Max(\sum_{C_k, L_m} T(C_k, L_m))/N, \quad (16)$$

where C_k denotes the k -th cluster, and L_m is the m -th class. $T(C_k, L_m)$ is the number of entities which belong to class m are assigned to cluster k . Accuracy computes the maximum sum of $T(C_k, L_m)$ for all pairs of clusters and classes, and these pairs have no overlaps. The greater accuracy means the better clustering performance.

Purity measures the extent to which each cluster contained data points from primarily one class. In general, the larger the values of purity, the better the clustering solution is. Entropy measures how classes distributed on various clusters. Generally, the smaller the entropy value, the better the clustering quality is. More details on the purity and entropy measures can be found in (Zhao & Karypis 2004).

The Rand Index is defined as the number of pairs of objects which are both located in the same cluster and the same class, or both in different clusters and different classes, divided by the total number of objects (WM 1971). Adjusted Rand Index which adjusts Rand Index is set between $[0, 1]$ (GW & MC 1986). The higher the Adjusted Rand Index, the more resemblance between the clustering results and the labels.

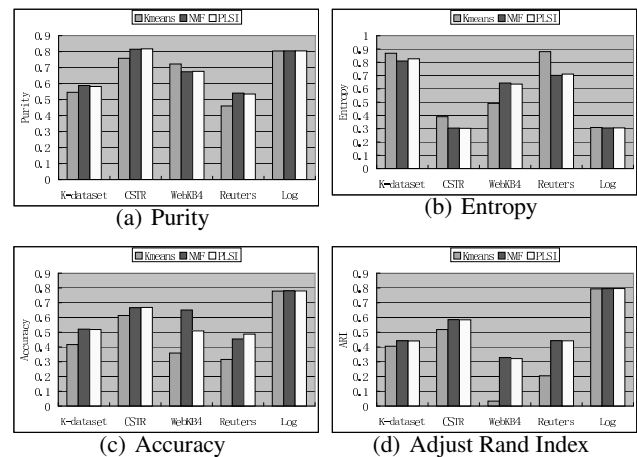


Figure 1: Performance Comparison of NMF and PLSI

Performance Comparison

For each of the five datasets we first run K-means clustering. This serves as a comparison and also initialization. From the K-means solution, H_0 is constructed from the cluster assignments and C_0 is simple the cluster centroids (see footnote 2). The H_0 obtained this way is discrete (0 and 1) and is very sparse (mostly zeroes). This is generally poor for multiplicative updation algorithms. Thus we smooth H_0 by adding 0.2 to every elements of H_0 . We then do necessary normalization on C_0, H_0 . Starting from this smoothed K-means solution, we run NMF or PLSI. From the NMF or PLSI solution, we harden the posterior H (see footnote 2) to obtain a discrete H (containing 0 and 1). From here, the performance measures are computed. We typically run 10 runs and obtain the average.

The clustering solutions of NMF and PLSI are compared based on accuracy, entropy, purity, and Adjusted Rand Index as shown in Figure 1. From these figures, we observe that

NMF and PLSI lead to similar clustering results. For example, as shown in Figure 1(a), in terms of purity value, the differences between the clustering solutions obtained by NMF and PLSI are less than 0.02 in all the datasets. We can observe similar behavior for other performance measures as well ².

Agreements Between NMF and PLSI

However, the closeness of NMF and PLSI on these four performance measures merely indicates the *level* of agreement between the NMF clustering solution and the known class label information is close to the *level* of agreement between PLSI and known class labels.

To understanding the difference between NMF and PLSI, we compare NMF and PLSI solutions *directly*: We measure the number of difference in clustering of data pairs using the clustering matrix R in Eq.(15). To normalize the difference so that datasets of different sizes can be compared with, we measure the relative difference:

$$\delta = \|R_{\text{NMF}} - R_{\text{PLSI}}\|_F / \sqrt{\|R_{\text{NMF}}\|_F^2/2 + \|R_{\text{PLSI}}\|_F^2/2}$$

The computed results, the average of 10 different runs, are listed in line A of Table 2. The results show that the differences between NMF and PLSI are quite substantial for WebKB (24%), and ranges between 1% to 8% in general cases.

	WebAce	CSTR	WebKB	Reuters	Log
A	0.083	0.072	0.239	0.070	0.010
B	0.029	0.025	0.056	0.051	0.010
C	0.022	0.013	0.052	0.040	0.012

Table 2: Dis-agreements between NMF and PLSI. All 3 type experiments begin with the same smoothed K-means. (A) Smoothed K-means to NMF. Smoothed K-means to PLSI. (B) Smoothed K-means to NMF to PLSI. (C) Smoothed K-means to PLSI to NMF.

Function J_{NMF} defines a surface in the multi-dimensional space. Because this global objective function is not a convex function, there are in general a very large number of local minima in the high p -dimensional space. Our experimental results suggest that starting with same initial smoothed K-means solution, NMF and PLSI converge to different local minima. In many cases, NMF and PLSI converge to *nearby* local minima; In other cases they converge to *not-so-nearby* local minima.

A Hybrid NMF-PLSI Algorithm

We have seen that NMF and PLSI optimize the same objective function, but their different detailed algorithms converge

²One thing we need to point out is that, in terms of accuracy, NMF and PLSI have a large difference of about 0.2 on *WebKB* dataset. This is because *WebKB* contain a lot of confusing web-pages that can be assigned to one or more clusters and the accuracy measure takes into account the entire distribution of the documents in a particular cluster and not just the largest class as in the computation of the purity.

to different local minima. An interesting question arises. Starting from a local minimum of NMF, could we jump out the local minimum by running the PLSI algorithm? Strictly speaking, if an algorithm makes an infinitesimal step, it will not jump out of a local minimum (we ignore the situation that the minimum could be saddle points). But PLSI algorithm is a finite-step algorithm, so it is possible to jump out of a local minimum reached by NMF. Vice versa, NMF is also a finite-step algorithm.

Interestingly, experiments indicate we can jump out of local minima this way. The results are shown in Table 2 Lines B & C. In Line B, we start from the K -means solution with smoothing and converge to a local minimum using NMF. Starting from the same local minimum, we run PLSI till convergence. The solution changed and the difference is given in Line B. This change indicates that we jump out of the local minimum. The changes in the solutions are smaller than Line A, as expected. In Line C, we start from the K -means solution with smoothing and then run PLSI to converge to a local minimum; we then jump out of this local minimum by running NMF. The difference of the solutions is given in Line C. The changes in the solutions are smaller than line A, as expected. The changes are also smaller than line B, indicating the local minimum reached by PLSI is perhaps slightly deeper than the local minima reached by NMF.

Based on the ability of NMF for jumping out of local minima of PLSI and vice versa, we propose a hybrid algorithm that alternatively runs NMF and PLSI, with the goal of successive jumping out local minima and therefore converging to a better minimum. The hybrid algorithm consists of 2 steps (1) K-means and smooth. (2) Iterate till converge: (2a) Run NMF to converge. and (2b) Run PLSI to converge. We run the hybrid algorithm on all 5 datasets. The results are listed in Table 3. We observe that: (1) NMF and PLSI always improve upon K-means. (2) Hybrid always improve upon NMF and PLSI; the improvements are significant on 3 out of 5 datasets.

	Reuters	WebKB	CSTR	WebAce	Log
A	0.316	0.410	0.617	0.416	0.775
B	0.454	0.619	0.666	0.520	0.778
C	0.487	0.510	0.668	0.519	0.779
D	0.521	0.644	0.878	0.523	0.781

Table 3: Clustering Accuracy. (A) K-means. (B) NMF-only. (C) PLSI-only. (D) Hybrid.

Summary

In this paper we show that NMF and PLSI optimize the same objective function. Based on this analysis, we propose a hybrid algorithm which alternatively runs NMF and PLSI. Extensive experiments on 5 datasets show the significant improvement of the hybrid method over PLSI or NMF.

Acknowledgment. Chris Ding is partially supported by the US Dept of Energy, Office of Science. Tao Li is partially supported by a 2005 IBM Faculty Award, a 2005 IBM Shared University Research (SUR) Award, and the NSF

grant IIS-0546280. Wei Peng is supported by a Florida International University Presidential Graduate Fellowship.

References

- Blei, D.; Ng, A.; and Jordan, M. 2003. Latent dirichlet allocation. *Journal of Machine Learning Research* 3:993-1022.
- Ding, C.; He, X.; and Simon, H. 2005. On the equivalence of nonnegative matrix factorization and spectral clustering. *Proc. SIAM Data Mining Conf.*
- Gaussier, E., and Goutte, C. 2005. Relation between pLSA and NMF and implications. In *SIGIR '05: Proceedings of the 28th annual international ACM SIGIR conference on Research and development in information retrieval*, 601–602. New York, NY, USA: ACM Press.
- GW, M., and MC, C. 1986. A study of the comparability of external criteria for hierarchical cluster analysis. *Multivar Behav Res* 21:846–850.
- Han, E.-H.; Boley, D.; Gini, M.; Gross, R.; Hastings, K.; Karypis, G.; Kumar, V.; Mobasher, B.; and Moore, J. 1998. WebACE: A web agent for document categorization and exploration. In *Proceedings of the 2nd International Conference on Autonomous Agents (Agents'98)*. ACM Press.
- Hofmann, T. 1999. Probabilistic latent semantic analysis. In *Proceedings of the 15th Annual Conference on Uncertainty in Artificial Intelligence (UAI-99)*, 289–296. San Francisco, CA: Morgan Kaufmann Publishers.
- Lee, D., and Seung, H. S. 1999. Learning the parts of objects by non-negative matrix factorization. *Nature* 401:788–791.
- Lee, D., and Seung, H. S. 2001. Algorithms for non-negative matrix factorization. In Dietterich, T. G., and Tresp, V., eds., *Advances in Neural Information Processing Systems*, volume 13. The MIT Press.
- Li, T. 2005. A general model for clustering binary data. In *KDD*, 188–197.
- Pauca, V. P.; Shahnaz, F.; Berry, M.; and Plemmons, R. 2004. Text mining using non-negative matrix factorization. In *Proc. SIAM Int'l conf on Data Mining (SDM 2004)*, 452–456.
- WM, R. 1971. Objective criteria for the evaluation of clustering methods. *J Am Stat Assoc* 66:846–850.
- Xu, W.; Liu, X.; and Gong, Y. 2003. Document clustering based on non-negative matrix factorization. In *Proc. ACM Conf. Research development in IR (SIGIR)*, 267–273.
- Zhao, Y., and Karypis, G. 2004. Empirical and theoretical comparisons of selected criterion functions for document clustering. *Machine Learning* 55(3):311–331.