# Factoring nonnegative matrices with linear programs

Victor Bittorf[*], Benjamin Recht[*], Christopher Ré[*], and Joel A. Tropp[†]

[*] Computer Sciences Department, University of Wisconsin-Madison

[†] Computational and Mathematical Sciences, California Institute of Technology

## Abstract

This paper describes a new approach, based on linear programming, for computing nonnegative matrix factorizations (NMFs). The key idea is a data-driven model for the factorization where the most salient features in the data are used to express the remaining features. More precisely, given a data matrix $X$, the algorithm identifies a matrix $C$ that satisfies $X \approx CX$ and some linear constraints. The constraints are chosen to ensure that the matrix $C$ selects features; these features can then be used to find a low-rank NMF of $X$. A theoretical analysis demonstrates that this approach has guarantees similar to those of the recent NMF algorithm of Arora et al. (2012). In contrast with this earlier work, the proposed method extends to more general noise models and leads to efficient, scalable algorithms. Experiments with synthetic and real datasets provide evidence that the new approach is also superior in practice. An optimized C++ implementation can factor a multigigabyte matrix in a matter of minutes.

**Keywords.** Nonnegative Matrix Factorization, Linear Programming, Stochastic gradient descent, Machine learning, Parallel computing, Multicore.

## 1 Introduction

Nonnegative matrix factorization (NMF) is a popular approach for selecting features in data [15–17, 22]. Many machine-learning and data-mining software packages (including Matlab [3], R [11], and Oracle Data Mining [1]) now include heuristic computational methods for NMF. Nevertheless, we still have limited theoretical understanding of when these heuristics are correct.

The difficulty in developing rigorous methods for NMF stems from the fact that the problem is computationally challenging. Indeed, Vavasis has shown that NMF is NP-Hard [26]; see [4] for further worst-case hardness results. As a consequence, we must instate additional assumptions on the data if we hope to compute nonnegative matrix factorizations in practice.

In this spirit, Arora, Ge, Kannan, and Moitra (AGKM) have exhibited a polynomial-time algorithm for NMF that is provably correct—provided that the data is drawn from an appropriate model, based on ideas from [7]. The AGKM result describes one circumstance where we can be sure that NMF algorithms are capable of producing meaningful answers. This work has the potential to make an impact in machine learning because proper feature selection is an important preprocessing step for many other techniques. Even so, the actual impact is damped by the fact that the AGKM algorithm is too computationally expensive for large-scale problems and is not tolerant to departures from the modeling assumptions. Thus, for NMF, there remains a gap between the theoretical exercise and the actual practice of machine learning.

The present work presents a scalable, robust algorithm that can successfully solve the NMF problem under appropriate hypotheses. Our first contribution is a new formulation of the nonnegative feature selection problem that only requires the solution of a single linear program. Second, we provide a theoretical analysis of this algorithm. This argument shows that our method succeeds under the same modeling assumptions as the AGKM algorithm with an additional *margin constraint* that is common in machine learning. We prove that if there exists a unique, well-defined model, then we can recover this model accurately; our error bound improves substantially on the error bound for the AGKM algorithm in the high SNR regime. One may argue that NMF only "makes sense" (i.e., is well posed) when a unique solution exists, and so we believe our result has independent interest. Furthermore, our algorithm can be adapted for a wide class of noise models.

In addition to these theoretical contributions, our work also includes a major algorithmic and experimental component. Our formulation of NMF allows us to exploit methods from operations research and database systems to design solvers that scale to extremely large datasets. We develop an efficient stochastic gradient descent (SGD) algorithm that is (at least) two orders of magnitude faster than the approach of AGKM when both are implemented in Matlab. We describe a parallel implementation of our SGD algorithm that can robustly factor matrices with $10^5$ features and $10^6$ examples in a few minutes on a multicore workstation.

Our formulation of NMF uses a data-driven modeling approach to simplify the factorization problem. More precisely, we search for a small collection of rows from the data matrix that can be used to express the other rows. This type of approach appears in a number of other factorization problems, including rank-revealing QR [14], interpolative decomposition [19], subspace clustering [9, 23], dictionary learning [10], and others. Our computational techniques can be adapted to address large-scale instances of these problems as well.

## 2   Separable Nonnegative Matrix Factorizations and Hott Topics

**Notation.** For a matrix $\boldsymbol{M}$ and indices $i$ and $j$, we write $\boldsymbol{M}_{i\cdot}$ for the $i$th row of $\boldsymbol{M}$ and $\boldsymbol{M}_{\cdot j}$ for the $j$th column of $\boldsymbol{M}$. We write $M_{ij}$ for the $(i, j)$ entry.

Let $\boldsymbol{Y}$ be a nonnegative $f \times n$ data matrix with columns indexing examples and rows indexing features. Exact NMF seeks a factorization $\boldsymbol{Y} = \boldsymbol{FW}$ where the feature matrix $\boldsymbol{F}$ is $f \times r$, where the weight matrix $\boldsymbol{W}$ is $r \times n$, and both factors are nonnegative. Typically, $r \ll \min\{f, n\}$.

Unless stated otherwise, we assume that each row of the data matrix $\boldsymbol{Y}$ is normalized so it sums to one. Under this hypothesis, we may also assume that each row of $\boldsymbol{F}$ and of $\boldsymbol{W}$ also sums to one [4].

It is notoriously difficult to solve the NMF problem. Vavasis showed that it is NP-complete to decide whether a matrix admits a rank-$r$ nonnegative factorization [26]. AGKM proved that an exact NMF algorithm can be used to solve 3-SAT in subexponential time [4].

The literature contains some mathematical analysis of NMF that can be used to motivate algorithmic development. Thomas [24] developed a necessary and sufficient condition for the existence of a rank-$r$ NMF. More recently, Donoho and Stodden [7] obtained a related sufficient condition for uniqueness. AGKM exhibited an algorithm that can produce a nonnegative matrix factorization under a weaker sufficient condition. To state their results, we need a definition.

**Definition 2.1** *A set of vectors $\{\boldsymbol{v}_1, \dots, \boldsymbol{v}_r\} \subset \mathbb{R}^d$ is* simplicial *if no vector $\boldsymbol{v}_i$ lies in the convex hull of $\{\boldsymbol{v}_j : j \neq i\}$. The set of vectors is $\alpha$-robust simplicial if, for each $i$, the $\ell_1$ distance from $\boldsymbol{v}_i$ to the convex hull of $\{\boldsymbol{v}_j : j \neq i\}$ is at least $\alpha$. Figure 1 illustrates these concepts.*

Algorithm 1: AGKM: Approximably Separable Nonnegative Matrix Factorization [4]

---

1: Initialize $R = \emptyset$.

2: Compute the $f \times f$ matrix $\boldsymbol{D}$ with $D_{ij} = \|\boldsymbol{X}_{i\cdot} - \boldsymbol{X}_{j\cdot}\|_1$.

3: **for** $k = 1, \ldots f$ **do**

4:     Find the set $\mathcal{N}_k$ of rows that are at least $5\epsilon/\alpha + 2\epsilon$ away from $\boldsymbol{X}_{k\cdot}$.

5:     Compute the distance $\delta_k$ of $\boldsymbol{X}_{k\cdot}$ from $\mathrm{conv}(\{\boldsymbol{X}_{j\cdot} : j \in \mathcal{N}_k\})$.

6:     **if** $\delta_k > 2\epsilon$, add $k$ to the set $R$.

7: **end for**

8: Cluster the rows in $R$ as follows: $j$ and $k$ are in the same cluster if $D_{jk} \leq 10\epsilon/\alpha + 6\epsilon$.

9: Choose one element from each cluster to yield $\boldsymbol{W}$.

10: $\boldsymbol{F} = \arg\min_{\boldsymbol{Z} \in \mathbb{R}^{f \times r}} \|\boldsymbol{X} - \boldsymbol{Z}\boldsymbol{W}\|_{\infty,1}$
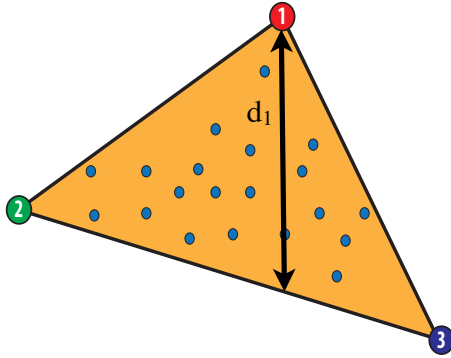
---



Figure 1: Numbered circles are hott topics. Their convex hull (orange) contains the other topics (small circles), so the data admits a separable NMF. The arrow $d_1$ marks the $\ell_1$ distance from hott topic (1) to the convex hull of the other two hott topics; definitions of $d_2$ and $d_3$ are similar. The hott topics are $\alpha$-robustly simplicial when each $d_i \geq \alpha$.

These ideas support the uniqueness results of Donoho and Stodden and the AGKM algorithm. Indeed, we can find an NMF of $\boldsymbol{Y}$ efficiently if $\boldsymbol{Y}$ contains a set of $r$ rows that is simplicial and whose convex hull contains the remaining rows.

**Definition 2.2** *An NMF* $\boldsymbol{Y} = \boldsymbol{F}\boldsymbol{W}$ *is called* separable *if the rows of* $\boldsymbol{W}$ *are simplicial and there is a permutation matrix* $\boldsymbol{\Pi}$ *such that*

$$\boldsymbol{\Pi}\boldsymbol{F} = \left[ \begin{array}{c} \mathbf{I}_r \\ \boldsymbol{M} \end{array} \right]. \tag{1}$$

To compute a separable factorization of $\boldsymbol{Y}$, we must first identify a simplicial set of rows from $\boldsymbol{Y}$. Afterward, we compute weights that express the remaining rows as convex combinations of this distinguished set. We call the simplicial rows *hott* and the corresponding features *hott topics*.

This model allows us to express all the features for a particular instance if we know the values of the instance at the simplicial rows. This assumption can be justified in a variety of applications. For example, in text, knowledge of a few keywords may be sufficient to reconstruct counts of the other words in a document. In vision, localized features can be used to predict gestures. In audio data, a few bins of the spectrogram may allow us to reconstruct the remaining bins.

While a nonnegative matrix one encounters in practice might not admit a separable factorization, it may be *well-approximated* by a nonnegative matrix with separable factorization. AGKM derived an algorithm for nonnegative matrix factorization of a matrix that is well-approximated by a separable factorization. To state their result, we introduce a norm on $f \times n$ matrices:

$$\|\boldsymbol{\Delta}\|_{\infty,1} := \max_{1 \leq i \leq f} \sum_{j=1}^{n} |\Delta_{ij}|.$$

**Theorem 2.3 (AGKM [4])** *Let $\epsilon$ and $\alpha$ be nonnegative constants satisfying $\epsilon \leq \frac{\alpha^2}{20 + 13\alpha}$. Let $\boldsymbol{X}$ be a nonnegative data matrix. Assume $\boldsymbol{X} = \boldsymbol{Y} + \boldsymbol{\Delta}$ where $\boldsymbol{Y}$ is a nonnegative matrix whose rows*

3

*have unit $\ell_1$ norm, where $\boldsymbol{Y} = \boldsymbol{FW}$ is a rank-r separable factorization in which the rows of $\boldsymbol{W}$ are $\alpha$-robust simplicial, and where $\|\boldsymbol{\Delta}\|_{\infty,1} \leq \epsilon$. Then Algorithm 1 finds a rank-r nonnegative factorization $\hat{\boldsymbol{F}}\hat{\boldsymbol{W}}$ that satisfies the error bound $\|\boldsymbol{X} - \hat{\boldsymbol{F}}\hat{\boldsymbol{W}}\|_{\infty,1} \leq 10\epsilon/\alpha + 7\epsilon$.*

In particular, the AGKM algorithm computes the factorization exactly when $\epsilon = 0$. Although this method is guaranteed to run in polynomial time, it has many undesirable features. First, the algorithm requires a priori knowledge of the parameters $\alpha$ and $\epsilon$. It may be possible to calculate $\epsilon$, but we can only estimate $\alpha$ if we know which rows are hott. Second, the algorithm computes all $\ell_1$ distances between rows at a cost of $O(f^2 n)$. Third, for every row in the matrix, we must determine its distance to the convex hull of the rows that lie at a sufficient distance; this step requires us to solve a linear program for each row of the matrix at a cost of $\Omega(fn)$. Finally, this method is intimately linked to the choice of the error norm $\|\cdot\|_{\infty,1}$. It is not obvious how to adapt the algorithm for other noise models. We present a new approach, based on linear programming, that overcomes these drawbacks.

# 3 Main Theoretical Results: NMF by Linear Programming

This paper shows that we can factor an approximately separable nonnegative matrix by solving a linear program. A major advantage of this formulation is that it scales to very large data sets.

Here is the key observation: Suppose that $\boldsymbol{Y}$ is any $f \times n$ nonnegative matrix that admits a rank-r separable factorization $\boldsymbol{Y} = \boldsymbol{FW}$. If we pad $\boldsymbol{F}$ with zeros to form an $f \times f$ matrix, we have

$$\boldsymbol{Y} = \boldsymbol{\Pi}^T \left[ \begin{array}{cc} \mathbf{I}_r & \mathbf{0} \\ \boldsymbol{M} & \mathbf{0} \end{array} \right] \boldsymbol{\Pi}\boldsymbol{Y} =: \boldsymbol{CY}$$

We call the matrix $\boldsymbol{C}$ *factorization localizing*. Note that any factorization localizing matrix $\boldsymbol{C}$ is an element of the polyhedral set

$$\Phi(\boldsymbol{Y}) := \{\boldsymbol{C} \geq \mathbf{0} \; : \; \boldsymbol{CY} = \boldsymbol{Y}, \; \mathrm{Tr}(\boldsymbol{C}) = r, \; C_{jj} \leq 1 \; \forall j, \; C_{ij} \leq C_{jj} \; \forall i,j\}.$$

Thus, to find an exact NMF of $\boldsymbol{Y}$, it suffices to find a feasible element of $\boldsymbol{C} \in \Phi(\boldsymbol{Y})$ whose diagonal is integral. This task can be accomplished by linear programming. Once we have such a $\boldsymbol{C}$, we construct $\boldsymbol{W}$ by extracting the rows of $\boldsymbol{X}$ that correspond to the indices $i$ where $C_{ii} = 1$. We construct the feature matrix $\boldsymbol{F}$ by extracting the nonzero columns of $\boldsymbol{C}$. This approach is summarized in Algorithm 2. In turn, we can prove the following result.

**Theorem 3.1** *Suppose $\boldsymbol{Y}$ is a nonnegative matrix with a rank-r separable factorization $\boldsymbol{Y} = \boldsymbol{FW}$. Then Algorithm 2 constructs a rank-r nonnegative matrix factorization of $\boldsymbol{Y}$.*

As the theorem suggests, we can isolate the rows of $\boldsymbol{Y}$ that yield a simplicial factorization by solving a single linear program. The factor $\boldsymbol{F}$ can be found by extracting columns of $\boldsymbol{C}$.

## 3.1 Robustness to Noise

Suppose we observe a nonnegative matrix $\boldsymbol{X}$ whose rows sum to one. Assume that $\boldsymbol{X} = \boldsymbol{Y} + \boldsymbol{\Delta}$ where $\boldsymbol{Y}$ is a nonnegative matrix whose rows sum to one, which has a rank-r separable factorization

---

**Algorithm 2** Separable Nonnegative Matrix Factorization by Linear Programming

---

**Require:** An $f \times n$ nonnegative matrix $\boldsymbol{Y}$ with a rank-$r$ separable NMF.
**Ensure:** An $f \times r$ matrix $\boldsymbol{F}$ and $r \times n$ matrix $\boldsymbol{W}$ with $\boldsymbol{F} \geq \boldsymbol{0}$, $\boldsymbol{W} \geq \boldsymbol{0}$, and $\boldsymbol{Y} = \boldsymbol{FW}$.
 1: Find the unique $\boldsymbol{C} \in \Phi(\boldsymbol{Y})$ to minimize $\boldsymbol{p}^T \operatorname{diag}(\boldsymbol{C})$ where $\boldsymbol{p}$ is any vector with distinct values.
 2: Let $I = \{i \ : \ C_{ii} = 1\}$ and set $\boldsymbol{W} = \boldsymbol{Y}_{I\cdot}$ and $\boldsymbol{F} = \boldsymbol{C}_{\cdot I}$.

---

---

**Algorithm 3** Approximably Separable Nonnegative Matrix Factorization by Linear Programming

---

**Require:** An $f \times n$ nonnegative matrix $\boldsymbol{X}$ that satisfies the hypotheses of Theorem 3.2.
**Ensure:** An $f \times r$ matrix $\boldsymbol{F}$ and $r \times n$ matrix $\boldsymbol{W}$ with $\boldsymbol{F} \geq \boldsymbol{0}$, $\boldsymbol{W} \geq \boldsymbol{0}$, and $\left\| \boldsymbol{X} - \boldsymbol{FW} \right\|_{\infty,1} \leq 2\epsilon$.
 1: Find $\boldsymbol{C} \in \Phi_{2\epsilon}(\boldsymbol{X})$ that minimizes $\boldsymbol{p}^T \operatorname{diag} \boldsymbol{C}$ where $\boldsymbol{p}$ is any vector with distinct values.
 2: Let $I = \{i \ : \ C_{ii} = 1\}$ and set $\boldsymbol{W} = \boldsymbol{X}_{I\cdot}$.
 3: Set $\boldsymbol{F} = \arg\min_{\boldsymbol{Z} \in \mathbb{R}^{f \times r}} \left\| \boldsymbol{X} - \boldsymbol{ZW} \right\|_{\infty,1}$

---

$\boldsymbol{Y} = \boldsymbol{FW}$ such that the rows of $\boldsymbol{W}$ are $\alpha$-robust simplicial, and where $\left\| \boldsymbol{\Delta} \right\|_{\infty,1} \leq \epsilon$. Define the polyhedral set

$$\Phi_\tau(\boldsymbol{X}) := \left\{ \boldsymbol{C} \geq \boldsymbol{0} \ : \ \left\| \boldsymbol{CX} - \boldsymbol{X} \right\|_{\infty,1} \leq \tau, \operatorname{Tr}(\boldsymbol{C}) = r, C_{jj} \leq 1 \ \forall j, C_{ij} \leq C_{jj} \ \forall i, j \right\}$$

The set $\Phi(\boldsymbol{X})$ consists of matrices $\boldsymbol{C}$ that *approximately* locate a factorization of $\boldsymbol{X}$. We can prove the following result.

**Theorem 3.2** *Suppose that $\boldsymbol{X}$ satisfies the assumptions stated in the previous paragraph. Furthermore, assume that for every row $\boldsymbol{Y}_{j,\cdot}$ that is not hott, we have the margin constraint $\left\| \boldsymbol{Y}_{j,\cdot} - \boldsymbol{Y}_{i,\cdot} \right\| \geq d_0$ for all hott rows $i$. Then we can find a nonnegative factorization satisfying $\left\| \boldsymbol{X} - \hat{\boldsymbol{F}}\hat{\boldsymbol{W}} \right\|_{\infty,1} \leq 2\epsilon$ provided that $\epsilon < \frac{\min\{\alpha d_0, \alpha^2\}}{9(r+1)}$. Furthermore, this factorization correctly identifies the hott topics appearing in the separable factorization of $\boldsymbol{Y}$.*

Algorithm 3 requires the solution of two linear programs. The first minimizes a cost vector over $\Phi_{2\epsilon}(\boldsymbol{X})$. This lets us find $\hat{\boldsymbol{W}}$. Afterward, the matrix $\hat{\boldsymbol{F}}$ can be found by setting

$$\hat{\boldsymbol{F}} = \arg \min_{\boldsymbol{Z} \geq \boldsymbol{0}} \ \left\| \boldsymbol{X} - \boldsymbol{Z}\hat{\boldsymbol{W}} \right\|_{\infty,1}. \tag{2}$$

Our robustness result requires a *margin-type* constraint assuming that the original configuration consists either of duplicate hott topics, or topics that are reasonably far away from the hott topics. On the other hand, under such a margin constraint, we can construct a considerably better approximation than that guaranteed by the AGKM algorithm. Moreover, unlike AGKM, our algorithm does not need to know the parameter $\alpha$.

The proofs of Theorems 3.1 and 3.2 can be found in the appendix. The main idea is to show that we can only represent a hott topic efficiently using the hott topic itself. Some earlier versions of this paper contained incomplete arguments, which we have remedied. For a signifcantly stronger robustness analysis of Algorithm 3, see the recent paper [12].

Having established these theoretical guarantees, it now remains to develop an algorithm to solve the LP. Off-the-shelf LP solvers may suffice for moderate-size problems, but for large-scale matrix factorization problems, their running time is prohibitive, as we show in Section 5. In Section 4, we turn to describe how to solve Algorithm 3 efficiently for large data sets.

## 3.2 Related Work

Localizing factorizations via column or row subset selection is a popular alternative to direct factorization methods such as the SVD. Interpolative decomposition such as Rank-Revealing QR [14] and CUR [19] have favorable efficiency properties as compared to factorizations (such as SVD) that are not based on exemplars. Factorization localization has been used in subspace clustering and has been shown to be robust to outliers [9, 23].

In recent work on dictionary learning, Esser et al. and Elhamifar et al. have proposed a factorization localization solution to nonnegative matrix factorization using group sparsity techniques [8, 10]. Esser et al. prove asymptotic exact recovery in a restricted noise model, but this result requires preprocessing to remove duplicate or near-duplicate rows. Elhamifar shows exact representative recovery in the noiseless setting assuming no hott topics are duplicated. Our work here improves upon this work in several aspects, enabling finite sample error bounds, the elimination of any need to preprocess the data, and algorithmic implementations that scale to very large data sets.

## 4 Incremental Gradient Algorithms for NMF

The rudiments of our fast implementation rely on two standard optimization techniques: dual decomposition and incremental gradient descent. Both techniques are described in depth in Chapters 3.4 and 7.8 of Bertsekas and Tstisklis [5].

We aim to minimize $\boldsymbol{p}^T \operatorname{diag}(\boldsymbol{C})$ subject to $\boldsymbol{C} \in \Phi_\tau(\boldsymbol{X})$. To proceed, form the Lagrangian

$$\mathcal{L}(\boldsymbol{C}, \beta, \boldsymbol{w}) = \boldsymbol{p}^T \operatorname{diag}(\boldsymbol{C}) + \beta(\operatorname{Tr}(\boldsymbol{C}) - r) + \sum_{i=1}^{f} w_i \left( \|\boldsymbol{X}_{i\cdot} - [\boldsymbol{C}\boldsymbol{X}]_{i\cdot}\|_1 - \tau \right)$$

with multipliers $\beta$ and $\boldsymbol{w} \geq \boldsymbol{0}$. Note that we do not dualize out all of the constraints. The remaining ones appear in the constraint set $\Phi_0 = \{\boldsymbol{C} \ : \ \boldsymbol{C} \geq \boldsymbol{0}, \ \operatorname{diag}(\boldsymbol{C}) \leq 1, \ \text{and} \ C_{ij} \leq C_{jj} \ \text{for all} \ i, j\}$.

Dual subgradient ascent solves this problem by alternating between minimizing the Lagrangian over the constraint set $\Phi_0$, and then taking a subgradient step with respect to the dual variables

$$w_i \leftarrow w_i + s\left( \|\boldsymbol{X}_{i\cdot} - [\boldsymbol{C}^\star\boldsymbol{X}]_{i\cdot}\|_1 - \tau \right) \quad \text{and} \quad \beta \leftarrow \beta + s(\operatorname{Tr}(\boldsymbol{C}^\star) - r)$$

where $\boldsymbol{C}^\star$ is the minimizer of the Lagrangian over $\Phi_0$. The update of $w_i$ makes very little difference in the solution quality, so we typically only update $\beta$.

We minimize the Lagrangian using projected incremental gradient descent. Note that we can rewrite the Lagrangian as

$$\mathcal{L}(\boldsymbol{C}, \beta, \boldsymbol{w}) = -\tau \boldsymbol{1}^T \boldsymbol{w} - \beta r + \sum_{k=1}^{n} \left( \sum_{j \in \operatorname{supp}(\boldsymbol{X}_{\cdot k})} w_j \|\boldsymbol{X}_{jk} - [\boldsymbol{C}\boldsymbol{X}]_{jk}\|_1 + \mu_j(p_j + \beta)C_{jj} \right) .$$

Here, $\operatorname{supp}(\boldsymbol{x})$ is the set indexing the entries where $\boldsymbol{x}$ is nonzero, and $\mu_j$ is the number of nonzeros in row $j$ divided by $n$. The incremental gradient method chooses one of the $n$ summands at random and follows its subgradient. We then project the iterate onto the constraint set $\Phi_0$. The projection onto $\Phi_0$ can be performed in the time required to sort the individual columns of $\boldsymbol{C}$ plus a linear-time operation. The full procedure is described in Appendix B. In the case where we expect a unique solution, we can drop the constraint $C_{ij} \leq C_{jj}$, resulting in a simple clipping procedure: set all

---
**Algorithm 4** HOTTOPIXX: Approximate Separable NMF by Incremental Gradient Descent
---
**Require:** An $f \times n$ nonnegative matrix $\boldsymbol{X}$. Primal and dual stepsizes $s_p$ and $s_d$.
**Ensure:** An $f \times r$ matrix $\boldsymbol{F}$ and $r \times n$ matrix $\boldsymbol{W}$ with $\boldsymbol{F} \geq \boldsymbol{0}$, $\boldsymbol{W} \geq \boldsymbol{0}$, and $\left\| \boldsymbol{X} - \boldsymbol{FW} \right\|_{\infty,1} \leq 2\epsilon$.
  1: Pick a cost $\boldsymbol{p}$ with distinct entries.
  2: Initialize $\boldsymbol{C} = \boldsymbol{0}$, $\beta = 0$
  3: **for** $t = 1, \ldots, N_{epochs}$ **do**
  4:    **for** $i = 1, \ldots n$ **do**
  5:       Choose $k$ uniformly at random from $[n]$.
  6:       $\boldsymbol{C} \leftarrow \boldsymbol{C} + s_p \cdot \text{sign}(\boldsymbol{X}_{\cdot k} - \boldsymbol{C}\boldsymbol{X}_{\cdot k})\boldsymbol{X}_{\cdot k}^T - s_p \, \text{diag}(\boldsymbol{\mu} \circ (\beta \boldsymbol{1} - \boldsymbol{p}))$.
  7:    **end for**
  8:    Project $\boldsymbol{C}$ onto $\Phi_0$.
  9:    $\beta \leftarrow \beta + s_d(\text{Tr}(\boldsymbol{C}) - r)$
10: **end for**
11: Let $I = \{i \ : \ C_{ii} = 1\}$ and set $\boldsymbol{W} = \boldsymbol{X}_{I\cdot}$.
12: Set $\boldsymbol{F} = \arg\min_{\boldsymbol{Z} \in \mathbb{R}^{f \times r}} \left\| \boldsymbol{X} - \boldsymbol{Z}\boldsymbol{W} \right\|_{\infty,1}$
---

negative items to zero and set any diagonal entry exceeding one to one. In practice, we perform a tradeoff. Since the constraint $C_{ij} \leq C_{jj}$ is used solely for symmetry breaking, we have found empirically that we only need to project onto $\Phi_0$ every $n$ iterations or so.

This incremental iteration is repeated $n$ times in a phase called an *epoch*. After each epoch, we update the dual variables and quit after we believe we have identified the large elements of the diagonal of $\boldsymbol{C}$. Just as before, once we have identified the hott rows, we can form $\boldsymbol{W}$ by selecting these rows of $\boldsymbol{X}$. We can find $\boldsymbol{F}$ just as before, by solving (2). Note that this minimization can also be computed by incremental subgradient descent. The full procedure, called HOTTOPIXX, is described in Algorithm 4.

## 4.1   Sparsity and Computational Enhancements for Large Scale.

For small-scale problems, HOTTOPIXX can be implemented in a few lines of Matlab code. But for the very large data sets studied in Section 5, we take advantage of natural parallelism and a host of low-level optimizations that are also enabled by our formulation. As in any numerical program, memory layout and cache behavior can be critical factors for performance. We use standard techniques: in-memory clustering to increase prefetching opportunities, padded data structures for better cache alignment, and compiler directives to allow the Intel compiler to apply vectorization.

Note that the incremental gradient step (step 6 in Algorithm 4) only modifies the entries of $\boldsymbol{C}$ where $\boldsymbol{X}_{\cdot k}$ is nonzero. Thus, we can parallelize the algorithm with respect to updating either the rows or the columns of $\boldsymbol{C}$. We store $\boldsymbol{X}$ in large contiguous blocks of memory to encourage hardware prefetching. In contrast, we choose a dense representation of our localizing matrix $\boldsymbol{C}$; this choice trades space for runtime performance.

Each worker thread is assigned a number of rows of $\boldsymbol{C}$ so that all rows fit in the shared L3 cache. Then, each worker thread repeatedly scans $\boldsymbol{X}$ while marking updates to multiple rows of $\boldsymbol{C}$. We repeat this process until all rows of $\boldsymbol{C}$ are scanned, similar to the classical block-nested loop join in relational databases [21].

# 5 Experiments

Except for the speedup curves, all of the experiments were run on an identical configuration: a dual Xeon X650 (6 cores each) machine with 128GB of RAM. The kernel is Linux 2.6.32-131.

In small-scale, synthetic experiments, we compared HOTTOPIXX to the AGKM algorithm and the linear programming formulation of Algorithm 3 implemented in Matlab. Both AGKM and Algorithm 3 were run using CVX [13] coupled to the SDPT3 solver [25]. We ran HOTTOPIXX for 50 epochs with primal stepsize 1e-1 and dual stepsize 1e-2. Once the hott topics were identified, we fit $F$ using two cleaning epochs of incremental gradient descent for all three algorithms.

To generate our instances, we sampled $r$ hott topics uniformly from the unit simplex in $\mathbb{R}^n$. These topics were duplicated $d$ times. We generated the remaining $f - r(d+1)$ rows to be random convex combinations of the hott topics, with the combinations selected uniformly at random. We then added noise with $(\infty, 1)$-norm error bounded by $\eta \cdot \frac{\alpha^2}{20+13\alpha}$. Recall that AGKM algorithm is only guaranteed to work for $\eta < 1$. We ran with $f \in \{40, 80, 160\}$, $n \in \{400, 800, 1600\}$, $r \in \{3, 5, 10\}$, $d \in \{0, 1, 2\}$, and $\eta \in \{0.25, 0.95, 4, 10, 100\}$. Each experiment was repeated 5 times.

Because we ran over 2000 experiments with 405 different parameter settings, it is convenient to use the *performance profiles* to compare the performance of the different algorithms [6]. Let $\mathcal{P}$ be the set of experiments and $\mathcal{A}$ denote the set of different algorithms we are comparing. Let $Q_a(p)$ be the value of some performance metric of the experiment $p \in \mathcal{P}$ for algorithm $a \in \mathcal{A}$. Then the performance profile at $\tau$ for a particular algorithm is the fraction of the experiments where the value of $Q_a(p)$ lies within a factor of $\tau$ of the minimal value of $\min_{b \in \mathcal{A}} Q_b(p)$. That is,

$$P_a(\tau) = \frac{\#\{p \in \mathcal{P} \ : \ Q_a(p) \leq \tau \min_{a' \in \mathcal{A}} Q_{a'}(p)\}}{\#(\mathcal{P})} .$$

In a performance profile, the higher a curve corresponding to an algorithm, the more often it outperforms the other algorithms. This gives a convenient way to contrast algorithms visually.

Our performance profiles are shown in Figure 2. The first two figures correspond to experiments with $f = 40$ and $n = 400$. The third figure is for the synthetic experiments with all other values of $f$ and $n$. In terms of $(\infty, 1)$-norm error, the linear programming solver typically achieves the lowest error. However, using SDPT3, it is prohibitively slow to factor larger matrices. On the other hand, HOTTOPIXX achieves better noise performance than the AGKM algorithm in much less time. Moreover, the AGKM algorithm must be fed the values of $\epsilon$ and $\alpha$ in order to run. HOTTOPIXX does not require this information and still achieves about the same error performance.

We also display a graph for running only four epochs (hott (fast)). This algorithm is by far the fastest algorithm, but does not achieve as optimal a noise performance. For very high levels of noise, however, it achieves a lower reconstruction error than the AGKM algorithm, whose performance degrades once $\eta$ approaches or exceeds 1 (Figure 2(f)). We also provide performance profiles for the root-mean-square error of the nonnegative matrix factorizations (Figure 2 (d) and (e)). The performance is qualitatively similar to that for the $(\infty, 1)$-norm.

We also coded HOTTOPIXX in C++, using the design principles described in Section 4.1, and ran on three large data sets. We generated a large synthetic example (jumbo) as above with $r = 100$. We generated a co-occurrence matrix of people and places from the ClueWeb09 Dataset [2], normalized by TFIDF. We also used HOTTOPIXX to select features from the RCV1 data set to recognize the class CCAT [18]. The statistics for these data sets can be found in Table 1.

In Figure 3 (left), we plot the speed-up over a serial implementation. In contrast to other parallel methods that exhibit memory contention [20], we see superlinear speed-ups for up to 20 threads
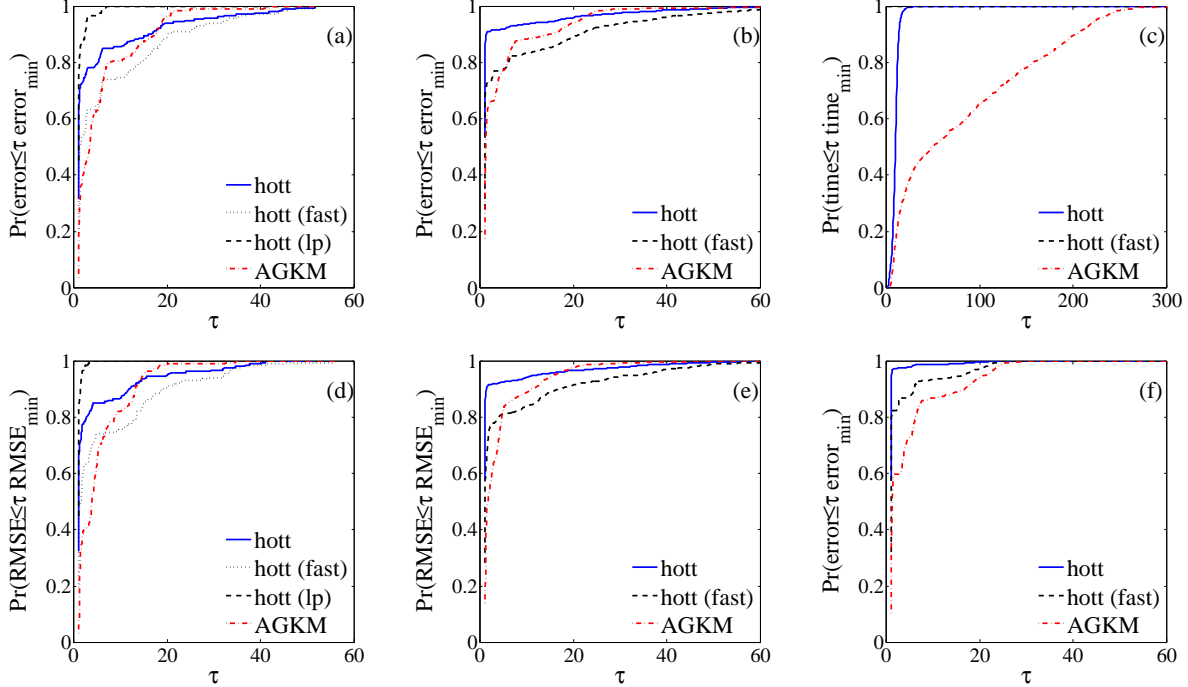
Figure 2: Performance profiles for synthetic data. (a) $(\infty, 1)$-norm error for $40 \times 400$ sized instances and (b) all instances. (c) is the performance profile for running time on all instances. RMSE performance profiles for the (d) small scale and (e) medium scale experiments. (f) $(\infty, 1)$-norm error for the $\eta \geq 1$. In the noisy examples, even 4 epochs of HOTTOPIXX is sufficient to obtain competitive reconstruction error.

| data set | features | documents | nonzeros | size (GB) | time (s) |
|----------|----------|-----------|----------|-----------|----------|
| jumbo    | 1600     | 64000     | 1.02e8   | 2.7       | 338      |
| clueweb  | 44739    | 351849    | 1.94e7   | 0.27      | 478      |
| RCV1     | 47153    | 781265    | 5.92e7   | 1.14      | 430      |

Table 1: Description of the large data sets. Time is to find 100 hott topics on the 12 core machines.

due to hardware prefetching and cache effects. All three of our large data sets can be trained in minutes, showing that we can scale HOTTOPIXX on both synthetic and real data. Our algorithm is able to correctly identify the hott topics on the jumbo set. For clueweb, we plot the RMSE Figure 3 (middle). This curve rolls off quickly for the first few hundred topics, demonstrating that our algorithm may be useful for dimensionality reduction in Natural Language Processing applications. For RCV1, we trained an SVM on the set of features extracted by HOTTOPIXX and plot the misclassification error versus the number of topics in Figure 3 (right). With 1500 hott topics, we achieve 7% misclassification error as compared to 5.5% with the entire set of features.

# 6   Discussion

This paper provides an algorithmic and theoretical framework for analyzing and deploying any factorization problem that can be posed as a linear (or convex) factorization localizing program.
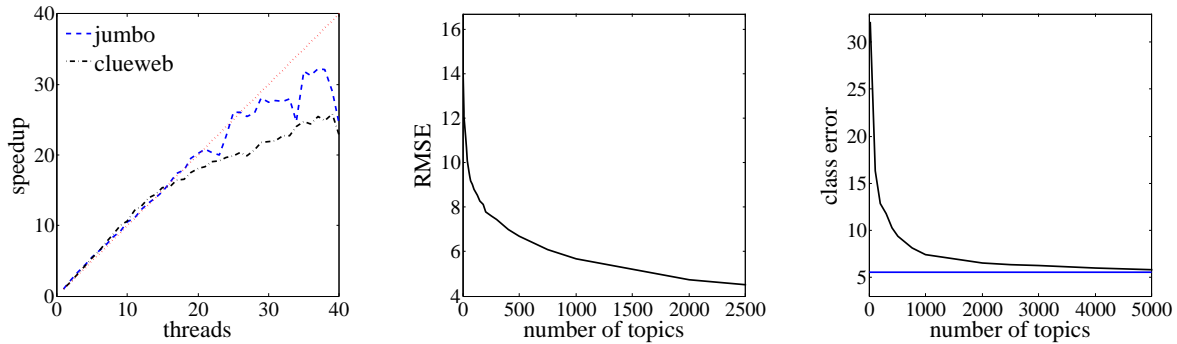
9

Figure 3: (left) The speedup over a serial implementation for HOTTOPIXX on the jumbo and clueweb data sets. Note the superlinear speedup for up to 20 threads. (middle) The RMSE for the clueweb data set. (right) The test error on RCV1 CCAT class versus the number of hott topics. The horizontal line indicates the test error achieved using all of the features.

Future work should investigate the applicability of HOTTOPIXX to other factorization localizing algorithms, such as subspace clustering, and should revisit earlier theoretical bounds on such prior art.

### Acknowledgments

# References

[1] docs.oracle.com/cd/B28359_01/datamine.111/b28129/algo_nmf.htm.

[2] lemurproject.org/clueweb09/.

[3] www.mathworks.com/help/toolbox/stats/nnmf.html.

[4] S. Arora, R. Ge, R. Kannan, and A. Moitra. Computing a nonnegative matrix factorization – provably. To appear in STOC 2012. Preprint available at \arxiv.org/abs/1111.0952, 2011.

[5] D. P. Bertsekas and J. N. Tsitsiklis. *Parallel and Distributed Computation: Numerical Methods*. Athena Scientific, Belmont, MA, 1997.

[6] E. D. Dolan and J. J. Moré. Benchmarking optimization software with performance profiles. *Mathematical Programming, Series A*, 91:201–213, 2002.

[7] D. Donoho and V. Stodden. When does non-negative matrix factorization give a correct decomposition into parts? In *Advances in Neural Information Processing Systems*, 2003.

[8] E. Elhamifar, G. Sapiro, and R. Vidal. See all by looking at a few: Sparse modeling for finding representative objects. In *Proceedings of CVPR*, 2012.

[9] E. Elhamifar and R. Vidal. Sparse subspace clustering. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2009.

[10] E. Esser, M. Möller, S. Osher, G. Sapiro, and J. Xin. A convex model for non-negative matrix factorization and dimensionality reduction on physical space. *IEEE Transactions on Image Processing*, 2012. To appear. Preprint available at `arxiv.org/abs/1102.0844`.

[11] R. Gaujoux and C. Seoighe. NMF: A flexible R package for nonnegative matrix factorization. *BMC Bioinformatics*, 11:367, 2010. `doi:10.1186/1471-2105-11-367`.

[12] N. Gillis. Robustness analysis of hotttopixx, a linear programming model for factoring nonnegative matrices. `arxiv.org/1211.6687`, 2012.

[13] M. Grant and S. Boyd. CVX: Matlab software for disciplined convex programming, version 1.21. `http://cvxr.com/cvx`, May 2010.

[14] M. Gu and S. C. Eisenstat. Efficient algorithms for computing a strong rank-revealing QR factorization. *SIAM Journal on Scientific Computing*, 17:848–869, 1996.

[15] T. Hofmann. Probabilistic latent semantic indexing. In *Proceedings of the 22nd Annual International SIGIR Conference on Research and Development in Information Retrieval*, 1999.

[16] D. D. Lee and H. S. Seung. Learning the parts of objects by non-negative matrix factorization. *Nature*, 401:788–791, 1999.

[17] D. D. Lee and H. S. Seung. Algorithms for non-negative matrix factorization. In *Advances in Neural Information Processing Systems*, 2001.

[18] D. Lewis, Y. Yang, T. Rose, and F. Li. RCV1: A new benchmark collection for text categorization research. *Journal of Machine Learning Research*, 5:361–397, 2004.

[19] M. W. Mahoney and P. Drineas. CUR matrix decompositions for improved data analysis. *Proceedings of the National Academy of Sciences*, 106:697–702, 2009.

[20] F. Niu, B. Recht, C. Ré, and S. J. Wright. HOGWILD!: A lock-free approach to parallelizing stochastic gradient descent. In *Advances in Neural Information Processing Systems*, 2011.

[21] L. D. Shapiro. Join processing in database systems with large main memories. *ACM Transactions on Database Systems*, 11(3):239–264, 1986.

[22] P. Smaragdis. Non-negative matrix factorization for polyphonic music transcription. In *IEEE Workshop on Applications of Signal Processing to Audio and Acoustics*, pages 177–180, 2003.

[23] M. Soltanolkotabi and E. J. Candès. A geometric analysis of subspace clustering with outliers. Preprint available at `arxiv.org/abs/1112.4258`, 2011.

[24] L. B. Thomas. Problem 73-14, rank factorization of nonnegative matrices. *SIAM Review*, 16(3):393–394, 1974.

[25] K. C. Toh, M. Todd, and R. H. Tütüncü. *SDPT3: A MATLAB software package for semidefinite-quadratic-linear programming*. Available from `http://www.math.nus.edu.sg/~mattohkc/sdpt3.html`.

[26] S. A. Vavasis. On the complexity of nonnegative matrix factorization. *SIAM Joural on Optimization*, 20(3):1364–1377, 2009.

# A   Proofs

Let $\boldsymbol{Y}$ be a nonnegative matrix whose rows sum to one. Assume that $\boldsymbol{Y}$ admits an exact separable factorization of rank $r$. In other words, we can write $\boldsymbol{Y} = \boldsymbol{FW}$ where the rows of $\boldsymbol{W}$ are $\alpha$-robust simplicial and

$$\boldsymbol{\Pi F} = \begin{bmatrix} \mathbf{I}_r \\ \boldsymbol{M} \end{bmatrix}$$

for some permutation $\boldsymbol{\Pi}$. Let $I$ denote the indices of the rows in $\boldsymbol{Y}$ that correspond with the identity matrix in the factorization, which we have called the *hott rows*. Then we can write each row $j$ that is not hott as a convex combination of the hott rows:

$$\boldsymbol{Y}_{j\cdot} = \sum_{k \in I} M_{jk} \boldsymbol{Y}_{k\cdot}. \quad \text{for each } j \notin I.$$

As we have discussed, we may assume that $\sum_k M_{jk} = 1$ for each $j \notin I$ because each row of $\boldsymbol{Y}$ sums to one.

The first lemma offers a stronger bound on the coefficients $M_{jk}$ in terms of the distance between row $j$ and the hott rows.

**Lemma A.1** *For an index $\ell$, suppose that the row $\boldsymbol{Y}_{\ell\cdot}$ has distance greater than $\delta$ from a hott topic $\boldsymbol{Y}_{i\cdot}$ with $i \in I$. Then $M_{\ell i} \leq 1 - \delta/2$.*

**Proof** We can express the $\ell$th row as a convex combination of hott rows: $\boldsymbol{Y}_{\ell\cdot} = \sum_{k \in I} M_{\ell k} \boldsymbol{Y}_{k\cdot}$. For each $i \in I$, we can bound $M_{\ell i}$ as follows.

$$
\begin{aligned}
\delta \leq \|\boldsymbol{Y}_{i\cdot} - \boldsymbol{Y}_{\ell\cdot}\|_1 &= \left\| \boldsymbol{Y}_{i\cdot} - \sum_{k \in I} M_{\ell k} \boldsymbol{Y}_{k\cdot} \right\|_1 \\
&= \left\| (1 - M_{\ell i}) \boldsymbol{Y}_i - \sum_{k \in I \setminus \{i\}} M_{\ell k} \boldsymbol{Y}_{k\cdot} \right\|_1 \\
&\leq \|(1 - M_{\ell i}) \boldsymbol{Y}_i\|_1 + \sum_{k \in I \setminus \{i\}} M_{\ell k} \|\boldsymbol{Y}_{k\cdot}\|_1 \\
&= 2(1 - M_{\ell i}).
\end{aligned}
$$

The inequality is the triangle inequality. To reach the last line, we use the fact that each row of $\boldsymbol{Y}$ has $\ell_1$ norm equal to one. Furthermore, $\sum_{k \in I \setminus \{i\}} M_{\ell k} = 1 - M_{\ell i} \geq 0$ because each row of $\boldsymbol{M}$ consists of nonnegative numbers that sum to one. Rearrange to complete the argument.  ∎

The next lemma is the central tool in our proofs. It tells us that any representation of a hott row has to involve rows that are close in $\ell_1$ norm to a hott row. To state the result, we define for each hott row $i$

$$\mathcal{B}_\delta(i) = \{j : \|\boldsymbol{Y}_{i\cdot} - \boldsymbol{Y}_{j\cdot}\|_1 \leq \delta\}.$$

In other words, $\mathcal{B}_\delta(i)$ contains the indices of all rows with $\ell_1$ distance no greater than $\delta$ from the hott topic $\boldsymbol{Y}_{i\cdot}$.

**Lemma A.2** *Let $\boldsymbol{c} \in \mathbb{R}^f$ be a nonnegative vector whose entries sum to one. For some hott row $i \in I$, suppose that $\|\boldsymbol{c}^T \boldsymbol{Y} - \boldsymbol{Y}_{i\cdot}\|_1 \leq \tau$. Then*

$$\sum_{j \in \mathcal{B}_\delta(i)} c_j \geq 1 - \frac{2\tau}{\min\{\alpha\delta, \alpha^2\}}. \tag{3}$$

**Proof** Let us introduce notation for the quantity of interest: $w_i = w_i(\boldsymbol{c}) = \sum_{j \in \mathcal{B}_\delta(\boldsymbol{X}_{i\cdot})} c_j$. We may assume that $w_i < 1$, or else the result holds trivially. Since the entries of $\boldsymbol{c}$ sum to one, we have

$$0 < 1 - w_i = \sum_j c_j - \sum_{j \in \mathcal{B}_\delta(i)} c_j = \sum_{j \notin \mathcal{B}_\delta(i)} c_j.$$

Next, we introduce the extra assumption that $\delta < \alpha$. It is clear that $w_i$ increases monotonically with $\delta$, so any lower bound on $w_i$ that we establish in this case extends to a bound that holds for larger $\delta$. Since the hott topics are $\alpha$-robust simplicial, all the other hott topics are at least $\delta$ away from $\boldsymbol{Y}_{i\cdot}$ in $\ell_1$ norm. Therefore, the hott row $i$ is the unique hott row listed in $\mathcal{B}_\delta(i)$.

To establish the result, we may as well assume that $w_i(\boldsymbol{c})$ achieves its minimum possible value subject to the constraints that the value of $\boldsymbol{c}^T \boldsymbol{Y}$ is fixed and that $\boldsymbol{c}$ is a nonnegative vector that sums to one. We claim that this minimum such $w_i$ occurs if and only if $c_j = 0$ for all $j \in \mathcal{B}_\delta(i) \setminus \{i\}$. We complete the proof under this additional surmise.

The assertion in the last paragraph follows from an argument by contradiction. Suppose that $w_i(\boldsymbol{c})$ were minimized at a vector $\boldsymbol{c}$ where $c_j > 0$ for some $j \in \mathcal{B}_\delta(i) \setminus \{i\}$. Then we can construct another set of coefficients $\tilde{\boldsymbol{c}}$ that satisfies the constraints and leads to a smaller value of $w_i$. We have the representation $\boldsymbol{Y}_{j\cdot} = \sum_{k \in I} M_{jk} \boldsymbol{Y}_{k\cdot}$. Set $\tilde{c}_j = 0$; set $\tilde{c}_k = c_k + c_j M_{jk}$ for each $k \in I$, and set $\tilde{c}_k = c_k$ for all remaining $k \notin I \cup \{j\}$. It is easy to verify that $\tilde{\boldsymbol{c}}^T \boldsymbol{Y} = \boldsymbol{c}^T \boldsymbol{Y}$ and that $\tilde{\boldsymbol{c}}$ is a nonnegative vector whose entries sum to one. But the value of $w_i$ is strictly smaller with the coefficients $\tilde{\boldsymbol{c}}$:

$$w_i(\tilde{\boldsymbol{c}}) = \sum_{k \in \mathcal{B}_\delta(i)} \tilde{c}_k < \sum_{k \in \mathcal{B}_\delta(i)} c_k = w_i(\boldsymbol{c})$$

In this relation, all the summands cancel, except for the one with index $j$. But $\tilde{c}_j = 0 < c_j$. It follows that the minimum value of $w_i$ cannot occur when $c_j > 0$. Compactness of the constraint set assures us that there is some vector $\boldsymbol{c}$ of coefficients that minimizes $w_i(\boldsymbol{c})$, so we must conclude that the minimizer $\boldsymbol{c}$ satisfies $c_j = 0$ for $j \in \mathcal{B}_\delta(i) \setminus \{i\}$.

Let us continue. Owing to the assumption that $\boldsymbol{Y}_{i\cdot}$ is no farther than $\tau$ from $\boldsymbol{c}^T \boldsymbol{Y}$, we have

$$\begin{aligned}
\tau &\geq \left\| \boldsymbol{Y}_{i\cdot} - \sum_j c_j \boldsymbol{Y}_{j\cdot} \right\|_1 \\
&= \left\| (1 - w_i)\boldsymbol{Y}_{i\cdot} - \sum_{j \notin \mathcal{B}_\delta(i)} c_j \boldsymbol{Y}_{j\cdot} \right\|_1 \\
&= (1 - w_i) \left\| \boldsymbol{Y}_{i\cdot} - \frac{1}{1 - w_i} \sum_{j \notin \mathcal{B}_\delta(i)} \sum_{k \in I} c_j M_{jk} \boldsymbol{Y}_{k\cdot} \right\|_1.
\end{aligned} \tag{4}$$

The first line follows when we split the sum over $j$ based on whether or not the components fall in $\mathcal{B}_\delta(i)$. Then we apply the property that $c_j = 0$ for $j \in \mathcal{B}_\delta(i) \setminus \{i\}$, and we identify the quantity $w_i$. In the last line, we factored out $1 - w_i$, and we introduced the separable factorization of $\boldsymbol{Y}$.

13

Next, for each $k \in I$, define

$$\pi_k := \frac{1}{1 - w_i} \sum_{j \notin \mathcal{B}_\delta(i)} c_j M_{jk},$$

and note that $\pi_k \geq 0$. Furthermore,

$$\sum_{k \in I} \pi_k = \frac{1}{1 - w_i} \sum_{j \notin \mathcal{B}_\delta(i)} c_j \sum_{k \in I} M_{jk} = \frac{1}{1 - w_i} \sum_{j \notin \mathcal{B}_\delta(i)} c_j = 1$$

because the rows of $\boldsymbol{M}$ sum to one and because of the definition of $w_i$. Lemma A.1 implies that $\pi_i$ satisfies the bound

$$\pi_i = \frac{1}{1 - w_i} \sum_{j \notin \mathcal{B}_\delta(i)} c_j M_{ji} \leq \frac{1 - \delta/2}{1 - w_i} \sum_{j \notin \mathcal{B}_\delta(i)} c_j = 1 - \delta/2. \tag{5}$$

Indeed, the lemma is valid because $\boldsymbol{Y}_{j\cdot}$ is at least a distance of $\delta$ away from $\boldsymbol{Y}_{i\cdot}$ for every $j \notin \mathcal{B}_\delta(i)$.

With these observations, we can continue our calculation from (4):

$$\tau \geq (1 - w_i) \left\| \boldsymbol{Y}_{i\cdot} - \sum_{k \in I} \pi_k \boldsymbol{Y}_{k\cdot} \right\|_1$$

$$= (1 - w_i)(1 - \pi_i) \left\| \boldsymbol{Y}_{i\cdot} - \sum_{k \in I \setminus \{i\}} \frac{\pi_k}{1 - \pi_i} \boldsymbol{Y}_{k\cdot} \right\|_1$$

$$\geq (1 - w_i)(1 - \pi_i)\alpha$$

$$\geq (1 - w_i)(\delta/2)\alpha.$$

The first identity follows when we combine the $i$th term in the sum with $\boldsymbol{Y}_{i\cdot}$. The inequality depends on the assumption that $\boldsymbol{W}$ is $\alpha$-robust simplicial; any convex combination of $\{\boldsymbol{Y}_{k\cdot} : k \notin I\}$ is at least $\alpha$ away from $\boldsymbol{Y}_{i\cdot}$ in $\ell_1$ norm. Afterward, we use the bound (5). Rearrange the final expression to complete the argument. ∎

## A.1   Proof of Theorem 3.1

This result is almost obvious when there are no duplicated rows. Indeed, since the hott topics form a simplicial set and the matrix $\boldsymbol{Y}$ admits a separable factorization, the only way we can represent all $r$ hott topics exactly is to have $C_{ii} = 1$ for every hott row $i$. This exhausts the trace constraint, and we see that every other diagonal entry $C_{kk} = 0$ for every not hott row $k$. The only matrices that are feasible identify the hott rows on the diagonal. They must represent the remaining rows using linear combinations of the hott topics because of the constraints $\boldsymbol{CY} = \boldsymbol{Y}$ and $C_{ij} \leq C_{jj}$. It follows that the only feasible matrices are factorization localizing matrices.

When there are duplicated rows, the analysis is slightly more delicate. By the same argument as above, all the weight on the diagonal must be concentrated on hott rows. But the objective $\boldsymbol{p}^T \text{diag}(\boldsymbol{C})$ ensures that, out of any set of duplicates of a given topic, we always pick the duplicate row $j$ where $p_j$ is smallest; otherwise, we could reduce the objective further. Therefore, the diagonal of $\boldsymbol{C}$ identifies all $r$ distinct hott topics, and we select each one duplicate of each topic. As before, the other constraints ensure that the remaining rows are represented with this distinguished choice of hott topic exemplars. Therefore, the only minimizers are factorization localizing matrices that identify each hott topic exactly once.

## A.2    Proof of Theorem 3.2

Let $X = Y + \Delta$. The matrix $X$ is the observed data, with rows scaled to have unit sum, and the perturbation matrix $\Delta$ satisfies $\lVert\Delta\rVert_{\infty,1} \leq \epsilon$. We assume that $Y$ is a nonnegative matrix whose rows sum to one, and we posit that it admits a rank-$r$ separable NMF $Y = FW$ where $W$ is $\alpha$-robust simplicial. We write $I$ for the set of rows corresponding to hott topics in $Y$.

Suppose that $C_0$ is a factorization localizing matrix for the underlying matrix $Y$. That is, $C_0 Y = Y$ and each row of $C_0$ sums to one. It follows that

$$\lVert C_0\Delta - \Delta\rVert_{\infty,1} \leq (\lVert C_0\rVert_{\infty,1} + \lVert I\rVert_{\infty,1})\lVert\Delta\rVert_{\infty,1} \leq 2\epsilon.$$

Using our decomposition $X = Y + \Delta$, we quickly verify that

$$\lVert C_0 X - X\rVert_{\infty,1} \leq \lVert C_0 Y - Y\rVert_{\infty,1} + \lVert C_0\Delta - \Delta\rVert_{\infty,1} \leq 2\epsilon.$$

The point here is that a factorization localizing matrix for $Y$ serves as an approximate factorization localizing matrix for $X$.

Our approach for constructing an approximate factorization of $X$ requires us to minimize a cost function $t^T \operatorname{diag}(C)$ over the constraint set

$$\Phi_{2\epsilon}(X) = \left\{ C \geq 0 : \lVert CX - X\rVert_{\infty,1} \leq 2\epsilon, \operatorname{Tr}(C) = r, C_{jj} \leq 1\ \forall j,\ C_{ij} \leq C_{jj}\ \forall i, j \right\}. \qquad (6)$$

Note that the factorization localizing matrix $C_0$ for $Y$ is a member of this set, so the optimization problem we solve in Theorem 3.2 is feasible.

Suppose that $C \in \Phi_{2\epsilon}(X)$ is arbitrary. Let us check that the row sums of $C$ are not much larger than one. To that end, note that

$$C\mathbf{1} = CX\mathbf{1} = X\mathbf{1} + (CX - X)\mathbf{1} = \mathbf{1} + (CX - X)\mathbf{1}.$$

We have twice used the fact that every row of $X$ sums to one. For any row $c$ of the matrix $C$, this formula yields $c^T\mathbf{1} \leq 1 + 2\epsilon$ since $\lVert CX - X\rVert_{\infty,1} \leq 2\epsilon$. As a consequence,

$$\lVert C\Delta - \Delta\rVert_{\infty,1} \leq (\lVert C\rVert_{\infty,1} + \lVert I\rVert_{\infty,1})\lVert\Delta\rVert_{\infty,1} \leq (1 + 2\epsilon + 1)\epsilon = 2\epsilon + 2\epsilon^2.$$

We may conclude that

$$\lVert CY - Y\rVert_{\infty,1} \leq \lVert CX - X\rVert_{\infty,1} + \lVert C\Delta - \Delta\rVert_{\infty,1} \leq 4\epsilon + 2\epsilon^2.$$

The margin assumption states that $\lVert Y_{\ell\cdot} - Y_{i\cdot}\rVert > d_0$ for every hott topic $i \in I$ and every row $\ell \notin I$. For any $i \in I$, Lemma A.2 ensures that any approximate representation $c^T Y$ of the $i$th row $Y_{i\cdot}$ with error at most $4\epsilon + 2\epsilon^2$ satisfies

$$c_i = \sum_{j \in \mathcal{B}_{d_0}(i)} c_j \geq 1 - \frac{8\epsilon + 4\epsilon^2}{\min\{\alpha d_0, \alpha^2\}}.$$

In particular, every matrix $C$ in the set $\Phi_{2\epsilon}(X)$ has $C_{ii} \geq 1 - (8\epsilon + 4\epsilon^2)/\min\{\alpha d_0, \alpha^2\}$ for each hott topic $i$. To ensure that hott topic $i$ has weight $C_{ii}$ greater than $1 - 1/(r+1)$ for each $i$, we need

$$\epsilon < \sqrt{1 + \frac{\min\{\alpha d_0, \alpha^2\}}{4(r+1)}} - 1 < \frac{\min\{\alpha d_0, \alpha^2\}}{9(r+1)}$$

15

Since there are $r$ hott rows, they carry total weight greater than $r(1 - 1/(r + 1))$. Given the trace constraint, that leaves less than $1 - 1/(r + 1)$ for the remaining rows. We see that each of the $r$ hott rows must carry more weight than every row that is not hott, so we can easily identify them.

Once we have identified the set $I$ of hott topics, we simply solve the second linear program

$$\underset{B}{\text{minimize}} \left\| X - \begin{bmatrix} I \\ B \end{bmatrix} X_I \right\|_{\infty,1} \tag{7}$$

to find a $2\epsilon$-accurate factorization.

# B    Projection onto $\Phi_0$

To project onto the set $\Phi_0$, note that we can compute the projection one column at a time. Moreover, the projection for each individual column amounts to (after permuting the entries of the column),

$$\{x \in \mathbb{R}^f \ : \ 0 \le x_i \le x_1 \ \forall i, x_1 \le 1\}.$$

Assume, again without loss of generality, that we want to project a vector $z$ with $z_2 \ge z_3 \ge \ldots \ge z_n$. Then we need to solve the quadratic program

$$\begin{array}{ll} \text{minimize} & \frac{1}{2}\|z - x\|^2 \\ \text{subject to} & 0 \le x_i \le x_1 \ \forall i, x_1 \le 1 \end{array} \tag{8}$$

The optimal solution can be found as follows. Let $k_c$ be the largest $k \in \{2, \ldots, f\}$ such that

$$z_{k_c+1} \le \Pi_{[0,1]} \left( \sum_{k=1}^{k_c} z_k \right) =: \mu$$

where $\Pi_{[0,1]}$ denotes the projection onto the interval $[0, 1]$. Set

$$\hat{x}_i = \begin{cases} \mu & i \le k_c \\ (z_i)_+ & i > k_c \end{cases}.$$

Then $\hat{x}$ is the optimal solution. A linear time algorithm for computing $\hat{x}$ is given by Algorithm 5

To prove that $\hat{x}$ is optimal, define

$$y_i = \begin{cases} z_i - \mu & i \le k_c \\ \min(z_i, 0) & i > k_c \end{cases}.$$

$y_i$ is the gradient of $\frac{1}{2}\|x - z\|^2$ at $\hat{x}$. Consider the LP

$$\begin{array}{ll} \text{minimize} & -y^T x \\ \text{subject to} & 0 \le x_i \le x_1 \ \forall i, x_1 \le 1 \end{array}.$$

$\hat{x}$ is an optimal solution for this LP because the cost is negative on the negative entries, 0 on the nonnegative entries that are larger than $k_c$, positive for $2 \le k \le k_c$, and nonpositive for $k = 1$. Hence, by the minimum principle, $\hat{x}$ is also a solution of (8).

**Algorithm 5** Column Squishing

---

**Require:** A vector $\boldsymbol{z} \in \mathbb{R}^f$ with $z_2 \geq z_3 \geq \ldots \geq z_n$.

**Ensure:** The projection of $\boldsymbol{z}$ onto $\{\boldsymbol{x} \in \mathbb{R}^f \ : \ 0 \leq x_i \leq x_1 \ \forall i \, , x_1 \leq 1\}$.

1:  $\mu \leftarrow z_1$.
2:  **for** $k = 2, \ldots, f$ **do**
3:      **if** $z_k \leq \Pi_{[0,1]}(\mu)$, Set $k_c = k - 1$ and **break**
4:      **else** set $\mu = \frac{k-1}{k}\mu + \frac{1}{k}z_k$.
5:  **end for**
6:  $x_1 \leftarrow \Pi_{[0,1]}(\mu)$
7:  **for** $k = 2, \ldots, k_c$ **set** $x_k = \Pi_{[0,1]}(\mu)$.
8:  **for** $k = (k_c + 1), \ldots, f$ **set** $x_k = (z_i)_+$.
9:  **return** $\boldsymbol{x}$

---