
Provable Alternating Gradient Descent for Non-negative Matrix Factorization with Strong Correlations

Yuanzhi Li¹ Yingyu Liang¹

Abstract

Non-negative matrix factorization is a basic tool for decomposing data into the feature and weight matrices under non-negativity constraints, and in practice is often solved in the alternating minimization framework. However, it is unclear whether such algorithms can recover the ground-truth feature matrix when the weights for different features are highly correlated, which is common in applications. This paper proposes a simple and natural alternating gradient descent based algorithm, and shows that with a mild initialization it provably recovers the ground-truth in the presence of strong correlations. In most interesting cases, the correlation can be in the same order as the highest possible. Our analysis also reveals its several favorable features including robustness to noise. We complement our theoretical results with empirical studies on semi-synthetic datasets, demonstrating its advantage over several popular methods in recovering the ground-truth.

1. Introduction

Non-negative matrix factorization (NMF) is an important tool in data analysis and is widely used in image processing, text mining, and hyperspectral imaging (e.g., (Lee & Seung, 1997; Blei et al., 2003; Yang & Leskovec, 2013)). Given a set of observations $\mathbf{Y} = \{y^{(1)}, y^{(2)}, \dots, y^{(n)}\}$, the goal of NMF is to find a feature matrix $\mathbf{A} = \{a_1, a_2, \dots, a_D\}$ and a non-negative weight matrix $\mathbf{X} = \{x^{(1)}, x^{(2)}, \dots, x^{(n)}\}$ such that $y^{(i)} \approx \mathbf{A}x^{(i)}$ for any i , or $\mathbf{Y} \approx \mathbf{A}\mathbf{X}$ for short. The intuition of NMF is to write each data point as a *non-negative* combination of the features.

Authors listed in alphabetic order. ¹Princeton University, Princeton, NJ, USA. Correspondence to: Yuanzhi Li <yuanzhil@cs.princeton.edu>, Yingyu Liang <yingyul@cs.princeton.edu>.

Proceedings of the 34th International Conference on Machine Learning, Sydney, Australia, PMLR 70, 2017. Copyright 2017 by the author(s).

By doing so, one can avoid cancellation of different features and improve interpretability by thinking of each $x^{(i)}$ as a (unnormalized) probability distribution over the features. It is also observed empirically that the non-negativity constraint on the coefficients can lead to better features and improved downstream performance of the learned features.

Unlike the counterpart which factorizes $\mathbf{Y} \approx \mathbf{A}\mathbf{X}$ without assuming non-negativity of \mathbf{X} , NMF is usually much harder to solve, and can even be NP-hard in the worse case (Arora et al., 2012b). This explains why, despite all the practical success, NMF largely remains a mystery in theory. Moreover, many of the theoretical results for NMF were based on very technical tools such as algebraic geometry (e.g., (Arora et al., 2012b)) or tensor decomposition (e.g. (Anandkumar et al., 2012)), which undermine their applicability in practice. Arguably, the most widely used algorithms for NMF use the *alternative minimization scheme*: In each iteration, the algorithm alternatively keeps \mathbf{A} or \mathbf{X} as fixed and tries to minimize some distance between \mathbf{Y} and $\mathbf{A}\mathbf{X}$. Algorithms in this framework, such as multiplicative update (Lee & Seung, 2001) and alternative non-negative least square (Kim & Park, 2008), usually perform well on real world data. However, alternative minimization algorithms are usually notoriously difficult to analyze. This problem is poorly understood, with only a few provable guarantees known (Awasthi & Risteski, 2015; Li et al., 2016). Most importantly, these results are only for the case when the coordinates of the weights are from essentially independent distributions, while in practice they are known to be correlated, for example, in correlated topic models (Blei & Lafferty, 2006). As far as we know, there exists no rigorous analysis of practical algorithms for the case with strong correlations.

In this paper, we provide a theoretical analysis of a natural algorithm AND (Alternative Non-negative gradient Descent) that belongs to the practical framework, and show that it probably recovers the ground-truth given a mild initialization. It works under general conditions on the feature matrix and the weights, in particular, allowing strong correlations. It also has multiple favorable features that are unique to its success. We further complement our theoretical analysis by experiments on semi-synthetic data, demon-

strating that the algorithm converges faster to the ground-truth than several existing practical algorithms, and providing positive support for some of the unique features of our algorithm. Our contributions are detailed below.

1.1. Contributions

In this paper, we assume a generative model of the data points, given the ground-truth feature matrix \mathbf{A}^* . In each round, we are given $y = \mathbf{A}^*x$,¹ where x is sampled i.i.d. from some *unknown* distribution μ and the goal is to recover the ground-truth feature matrix \mathbf{A}^* . We give an algorithm named AND that starts from a mild initialization matrix and provably converges to \mathbf{A}^* in polynomial time. We also justify the convergence through a sequence of experiments. Our algorithm has the following favorable characteristics.

1.1.1. SIMPLE GRADIENT DESCENT ALGORITHM

The algorithm AND runs in stages and keeps a working matrix $\mathbf{A}^{(t)}$ in each stage. At the t -th iteration in a stage, after getting one sample y , it performs the following:

$$\begin{aligned} \text{(Decode)} \quad z &= \phi_\alpha \left((\mathbf{A}^{(0)})^\dagger y \right), \\ \text{(Update)} \quad \mathbf{A}^{(t+1)} &= \mathbf{A}^{(t)} + \eta \left(yz^\top - \mathbf{A}^{(t)} z z^\top \right), \end{aligned}$$

where α is a threshold parameter,

$$\phi_\alpha(x) = \begin{cases} x & \text{if } x \geq \alpha, \\ 0 & \text{otherwise,} \end{cases}$$

$(\mathbf{A}^{(0)})^\dagger$ is the Moore-Penrose pseudo-inverse of $\mathbf{A}^{(0)}$, and η is the update step size. The decode step aims at recovering the corresponding weight for the data point, and the update step uses the decoded weight to update the feature matrix. The final working matrix at one stage will be used as the $\mathbf{A}^{(0)}$ in the next stage. See Algorithm 1 for the details.

At a high level, our update step to the feature matrix can be thought of as a gradient descent version of alternative non-negative least square (Kim & Park, 2008), which at each iteration alternatively minimizes $L(\mathbf{A}, \mathbf{Z}) = \|\mathbf{Y} - \mathbf{AZ}\|_F^2$ by fixing \mathbf{A} or \mathbf{Z} . Our algorithm, instead of performing a complete minimization, performs only a stochastic *gradient descent* step on the feature matrix. To see this, consider one data point y and consider minimizing $L(\mathbf{A}, z) = \|y - \mathbf{A}z\|_F^2$ with z fixed. Then the gradient of \mathbf{A} is just $-\nabla L(\mathbf{A}) = (y - \mathbf{A}z)z^\top$, which is exactly the update of our feature matrix in each iteration.

As to the decode step, when $\alpha = 0$, our decoding can be regarded as a one-shot approach minimizing $\|\mathbf{Y} - \mathbf{AZ}\|_F^2$

restricted to $\mathbf{Z} \geq 0$. Indeed, if for example projected gradient descent is used to minimize $\|\mathbf{Y} - \mathbf{AZ}\|_F^2$, then the projection step is exactly applying ϕ_α to \mathbf{Z} with $\alpha = 0$. A key ingredient of our algorithm is choosing α to be larger than zero and then decreasing it, which allows us to outperform the standard algorithms.

Perhaps worth noting, our decoding only uses $\mathbf{A}^{(0)}$. Ideally, we would like to use $(\mathbf{A}^{(t)})^\dagger$ as the decoding matrix in each iteration. However, such decoding method requires computing the pseudo-inverse of $\mathbf{A}^{(t)}$ at every step, which is extremely slow. Instead, we divide the algorithm into stages and in each stage, we only use the starting matrix in the decoding, thus the pseudo-inverse only needs to be computed once per stage and can be used across all iterations inside. We can show that our algorithm converges in polylogarithmic many stages, thus gives us to a much better running time. These are made clear when we formally present the algorithm in Section 4 and the theorems in Section 5 and 6.

1.1.2. HANDLING STRONG CORRELATIONS

The most notable property of AND is that it can provably deal with *highly* correlated distribution μ on the weight x , meaning that the coordinates of x can have very strong correlations with each other. This is important since such correlated x naturally shows up in practice. For example, when a document contains the topic “machine learning”, it is more likely to contain the topic “computer science” than “geography” (Blei & Lafferty, 2006).

Most of the previous theoretical approaches for analyzing alternating between decoding and encoding, such as (Awasthi & Risteski, 2015; Li et al., 2016; Arora et al., 2015), require the coordinates of x to be pairwise-independent, or almost pairwise-independent (meaning $\mathbb{E}_\mu[x_i x_j] \approx \mathbb{E}_\mu[x_i] \mathbb{E}_\mu[x_j]$). In this paper, we show that algorithm AND can recover \mathbf{A}^* even when the coordinates are highly correlated. As one implication of our result, when the sparsity of x is $O(1)$ and each entry of x is in $\{0, 1\}$, AND can recover \mathbf{A}^* even if each $\mathbb{E}_\mu[x_i x_j] = \Omega(\min\{\mathbb{E}_\mu[x_i], \mathbb{E}_\mu[x_j]\})$, matching (up to constant) the highest correlation possible. Moreover, we do not assume any prior knowledge about the distribution μ , and the result also extends to general sparsities as well.

1.1.3. PSEUDO-INVERSE DECODING

One of the feature of our algorithm is to use Moore-Penrose pseudo-inverse in decoding. Inverse decoding was also used in (Li et al., 2016; Arora et al., 2015; 2016). However, their algorithms require carefully finding an inverse such that certain norm is minimized, which is not as efficient as the vanilla Moore-Penrose pseudo-inverse. It was also observed in (Arora et al., 2016) that Moore-Penrose

¹We also consider the noisy case; see 1.1.5.

pesudo-inverse works equally well in practice, but the experiment was done only when $\mathbf{A} = \mathbf{A}^*$. In this paper, we show that Moore-Penrose pseudo-inverse also works well when $\mathbf{A} \neq \mathbf{A}^*$, both theoretically and empirically.

1.1.4. THRESHOLDING AT DIFFERENT α

Thresholding at a value $\alpha > 0$ is a common trick used in many algorithms. However, many of them still only consider a fixed α throughout the entire algorithm. Our contribution is a new method of thresholding that first sets α to be high, and gradually decreases α as the algorithm goes. Our analysis naturally provides the explicit rate at which we decrease α , and shows that our algorithm, following this scheme, can provably converge to the ground-truth \mathbf{A}^* in polynomial time. Moreover, we also provide experimental support for these choices.

1.1.5. ROBUSTNESS TO NOISE

We further show that the algorithm is robust to noise. In particular, we consider the model $y = \mathbf{A}^*x + \zeta$, where ζ is the noise. The algorithm can tolerate a general family of noise with bounded moments; we present in the main body the result for a simplified case with Gaussian noise and provide the general result in the appendix. The algorithm can recover the ground-truth matrix up to a small blow-up factor times the noise level in *each example*, when the ground-truth has a good condition number. This robustness is also supported by our experiments.

2. Related Work

Practical algorithms. Non-negative matrix factorization has a rich empirical history, starting with the practical algorithms of (Lee & Seung, 1997; 1999; 2001). It has been widely used in applications and there exist various methods for NMF, e.g., (Kim & Park, 2008; Lee & Seung, 2001; Cichocki et al., 2007; Ding et al., 2013; 2014). However, they do not have provable recovery guarantees.

Theoretical analysis. For theoretical analysis, (Arora et al., 2012b) provided a fixed-parameter tractable algorithm for NMF using algebraic equations. They also provided matching hardness results: namely they show there is no algorithm running in time $(mW)^{o(D)}$ unless there is a sub-exponential running time algorithm for 3-SAT. (Arora et al., 2012b) also studied NMF under separability assumptions about the features, and (Bhattacharyya et al., 2016) studied NMF under related assumptions. The most related work is (Li et al., 2016), which analyzed an alternating minimization type algorithm. However, the result only holds with strong assumptions about the distribution of the weight x , in particular, with the assumption that the coordinates of x are independent.

Topic modeling. Topic modeling is a popular generative model for text data (Blei et al., 2003; Blei, 2012). Usually, the model results in NMF type optimization problems with $\|x\|_1 = 1$, and a popular heuristic is *variational inference*, which can be regarded as alternating minimization in KL-divergence. Recently, there is a line of theoretical work analyzing tensor decomposition (Arora et al., 2012a; 2013; Anandkumar et al., 2013) or combinatorial methods (Awasthi & Risteski, 2015). These either need strong structural assumptions on the word-topic matrix \mathbf{A}^* , or need to know the distribution of the weight x , which is usually infeasible in applications.

3. Problem and Definitions

We use $\|\mathbf{M}\|_2$ to denote the 2-norm of a matrix \mathbf{M} . $\|x\|_1$ is the 1-norm of a vector x . We use $[\mathbf{M}]_i$ to denote the i -th row and $[\mathbf{M}]^i$ to denote the i -th column of a matrix \mathbf{M} . $\sigma_{\max}(\mathbf{M})$ ($\sigma_{\min}(\mathbf{M})$) stands for the maximum (minimal) singular value of \mathbf{M} , respectively. We consider a generative model for non-negative matrix factorization, where the data y is generated from²

$$y = \mathbf{A}^*x, \quad \mathbf{A}^* \in \mathbb{R}^{W \times D}$$

where \mathbf{A}^* is the ground-truth feature matrix, and x is a non-negative random vector drawn from an unknown distribution μ . The goal is to recover the ground-truth \mathbf{A}^* from i.i.d. samples of the observation y .

Since the general non-negative matrix factorization is NP-hard (Arora et al., 2012b), some assumptions on the distribution of x need to be made. In this paper, we would like to allow distributions as general as possible, especially those with strong correlations. Therefore, we introduce the following notion called (r, k, m, λ) -general correlation conditions (GCC) for the distribution of x .

Definition 1 (General Correlation Conditions, GCC). *Let $\Delta := \mathbb{E}[xx^\top]$ denote the second moment matrix.*

1. $\|x\|_1 \leq r$ and $x_i \in [0, 1], \forall i \in [D]$.
2. $\Delta_{i,i} \leq \frac{2k}{D}, \forall i \in [D]$.
3. $\Delta_{i,j} \leq \frac{m}{D^2}, \forall i \neq j \in [D]$.
4. $\Delta \succeq \frac{k}{D}\lambda \mathbf{I}$.

The first condition regularizes the sparsity of x .³ The second condition regularizes each coordinate of x_i so that there is no x_i being large too often. The third condition

²Section 6.2 considers the noisy case.

³Throughout this paper, the sparsity of x refers to the ℓ_1 norm, which is much weaker than the ℓ_0 norm (the support sparsity). For example, in LDA, the ℓ_1 norm of x is always 1.

regularizes the maximum pairwise correlation between x_i and x_j . The fourth condition always holds for $\lambda = 0$ since $\mathbb{E}[xx^\top]$ is a PSD matrix. Later we will assume this condition holds for some $\lambda > 0$ to avoid degenerate cases. Note that we put the weight k/D before λ such that λ defined in this way will be a positive constant in many interesting examples discussed below.

To get a sense of what are the ranges of k , m , and λ given sparsity r , we consider the following most commonly studied non-negative random variables.

Proposition 1 (Examples of GCC).

1. If x is chosen uniformly over s -sparse random vectors with $\{0, 1\}$ entries, then $k = r = s$, $m = s^2$ and $\lambda = 1 - \frac{1}{s}$.
2. If x is uniformly chosen from Dirichlet distribution with parameter $\alpha_i = \frac{s}{D}$, then $r = k = 1$ and $m = \frac{1}{sD}$ with $\lambda = 1 - \frac{1}{s}$.

For these examples, the result in this paper shows that we can recover \mathbf{A}^* for aforementioned random variables x as long as $s = O(D^{1/6})$. In general, there is a wide range of parameters (r, k, m, λ) such that learning \mathbf{A}^* is doable with polynomially many samples of y and in polynomial time.

However, just the GCC condition is not enough for recovering \mathbf{A}^* . We will also need a mild initialization.

Definition 2 (ℓ -initialization). *The initial matrix \mathbf{A}_0 satisfies for some $\ell \in [0, 1)$,*

1. $\mathbf{A}_0 = \mathbf{A}^*(\mathbf{\Sigma} + \mathbf{E})$, for some diagonal matrix $\mathbf{\Sigma}$ and off-diagonal matrix \mathbf{E} .
2. $\|\mathbf{E}\|_2 \leq \ell$, $\|\mathbf{\Sigma} - \mathbf{I}\|_2 \leq \frac{1}{4}$.

The condition means that the initialization is not too far away from the ground-truth \mathbf{A}^* . For any $i \in [D]$, the i -th column $[\mathbf{A}_0]^i = \mathbf{\Sigma}_{i,i}[\mathbf{A}^*]^i + \sum_{j \neq i} \mathbf{E}_{j,i}[\mathbf{A}^*]^j$. So the condition means that each feature $[\mathbf{A}_0]^i$ has a large fraction of the ground-truth feature $[\mathbf{A}^*]^i$ and a small fraction of the other features. $\mathbf{\Sigma}$ can be regarded as the magnitude of the component from the ground-truth in the initialization, while \mathbf{E} can be regarded as the magnitude of the error terms. In particular, when $\mathbf{\Sigma} = \mathbf{I}$ and $\mathbf{E} = 0$, we have $\mathbf{A}_0 = \mathbf{A}^*$. The initialization allows $\mathbf{\Sigma}$ to be a constant away from \mathbf{I} , and the error term \mathbf{E} to be ℓ (in our theorems ℓ can be as large as a constant).

In practice, such an initialization is typically achieved by setting the columns of \mathbf{A}_0 to reasonable ‘‘pure’’ data points that contain one major feature and a small fraction of some others (e.g. (Iida, 2016; Awasthi & Risteski, 2015)).

Algorithm 1 Alternating Non-negative gradient Descent (AND)

Input: Threshold values $\{\alpha_0, \alpha_1, \dots, \alpha_s\}$, T , \mathbf{A}_0

- 1: $\mathbf{A}^{(0)} \leftarrow \mathbf{A}_0$
- 2: **for** $j = 0, 1, \dots, s$ **do**
- 3: **for** $t = 0, 1, \dots, T$ **do**
- 4: On getting sample $y^{(t)}$, do:
- 5: $z^{(t)} \leftarrow \phi_{\alpha_j}((\mathbf{A}^{(0)})^\dagger y^{(t)})$
- 6: $\mathbf{A}^{(t+1)} \leftarrow \mathbf{A}^{(t)} + \eta (y^{(t)} - \mathbf{A}^{(t)} z^{(t)}) (z^{(t)})^\top$
- 7: **end for**
- 8: $\mathbf{A}^{(0)} \leftarrow \mathbf{A}^{(T+1)}$
- 9: **end for**

Output: $\mathbf{A} \leftarrow \mathbf{A}^{(T+1)}$

4. Algorithm

The algorithm is formally describe in Algorithm 1. It runs in s stages, and in the j -th stage, uses the same threshold α_j and the same matrix $\mathbf{A}^{(0)}$ for decoding, where $\mathbf{A}^{(0)}$ is either the input initialization matrix or the working matrix obtained at the end of the last stage. Each stage consists of T iterations, and each iteration decodes one data point and uses the decoded result to update the working matrix. It can use a batch of data points instead of one data point, and our analysis still holds.

By running in stages, we save most of the cost of computing $(\mathbf{A}^{(0)})^\dagger$, as our results show that only polylogarithmic stages are needed. For the simple case where $x \in \{0, 1\}^D$, the algorithm can use the same threshold value $\alpha = 1/4$ for all stages (see Theorem 1), while for the general case, it needs decreasing threshold values across the stages (see Theorem 4). Our analysis provides the hint for setting the threshold; see the discussion after Theorem 4, and Section 7 for how to set the threshold in practice.

5. Result for A Simplified Case

In this section, we consider the following simplified case:

$$y = \mathbf{A}^* x, \quad x \in \{0, 1\}^D. \quad (5.1)$$

That is, the weight coordinates x_i ’s are binary.

Theorem 1 (Main, binary). *For the generative model (5.1), there exists $\ell = \Omega(1)$ such that for every (r, k, m, λ) -GCC x and every $\epsilon > 0$, Algorithm AND with $T = \text{poly}(D, \frac{1}{\epsilon})$, $\eta = \frac{1}{\text{poly}(D, \frac{1}{\epsilon})}$, $\{\alpha_i\}_{i=1}^s = \{\frac{1}{4}\}_{i=1}^s$ for $s = \text{polylog}(D, \frac{1}{\epsilon})$ and an ℓ initialization matrix \mathbf{A}_0 , outputs a matrix \mathbf{A} such that there exists a diagonal matrix $\mathbf{\Sigma} \succeq \frac{1}{2}\mathbf{I}$ with $\|\mathbf{A} - \mathbf{A}^* \mathbf{\Sigma}\|_2 \leq \epsilon$ using $\text{poly}(D, \frac{1}{\epsilon})$ samples and iterations, as long as*

$$m = O\left(\frac{kD\lambda^4}{r^5}\right).$$

Therefore, our algorithm recovers the ground-truth \mathbf{A}^* up to scaling. The scaling is unavoidable since there is no assumption on \mathbf{A}^* , so we cannot, for example, distinguish \mathbf{A}^* from $2\mathbf{A}^*$. Indeed, if we in addition assume each column of \mathbf{A}^* has norm 1 as typical in applications, then we can recover \mathbf{A}^* directly. In particular, by normalizing each column of \mathbf{A} to have norm 1, we can guarantee that $\|\mathbf{A} - \mathbf{A}^*\|_2 \leq O(\epsilon)$.

In many interesting applications (for example, those in Proposition 1), k, r, λ are constants. The theorem implies that the algorithm can recover \mathbf{A}^* even when $m = O(D)$. In this case, $\mathbb{E}_\mu[x_i x_j]$ can be as large as $O(1/D)$, the same order as $\min\{\mathbb{E}_\mu[x_i], \mathbb{E}_\mu[x_j]\}$, which is the highest possible correlation.

5.1. Intuition

The intuition comes from assuming that we have the ‘‘correct decoding’’, that is, suppose magically for every $y^{(t)}$, our decoding $z^{(t)} = \phi_{\alpha_j}(\mathbf{A}^\dagger y^{(t)}) = x^{(t)}$. Here and in this subsection, \mathbf{A} is a shorthand for $\mathbf{A}^{(0)}$. The gradient descent is then $\mathbf{A}^{(t+1)} = \mathbf{A}^{(t)} + \eta(y^{(t)} - \mathbf{A}^{(t)}x^{(t)})(x^{(t)})^\top$. Subtracting \mathbf{A}^* on both side, we will get

$$(\mathbf{A}^{(t+1)} - \mathbf{A}^*) = (\mathbf{A}^{(t)} - \mathbf{A}^*)(\mathbf{I} - \eta x^{(t)}(x^{(t)})^\top)$$

Since $x^{(t)}(x^{(t)})^\top$ is positive semidefinite, as long as $\mathbb{E}[x^{(t)}(x^{(t)})^\top] \succ 0$ and η is sufficiently small, $\mathbf{A}^{(t)}$ will converge to \mathbf{A}^* eventually.

However, this simple argument does not work when $\mathbf{A} \neq \mathbf{A}^*$ and thus we do not have the correct decoding. For example, if we just let the decoding be $\tilde{z}^{(t)} = \mathbf{A}^\dagger y^{(t)}$, we will have $y^{(t)} - \mathbf{A}\tilde{z}^{(t)} = y^{(t)} - \mathbf{A}^\dagger \mathbf{A} y^{(t)} = (\mathbf{I} - \mathbf{A}^\dagger \mathbf{A})\mathbf{A}^* x^{(t)}$. Thus, using this decoding, the algorithm can never make any progress once \mathbf{A} and \mathbf{A}^* are in the same subspace.

The most important piece of our proof is to show that after *thresholding*, $z^{(t)} = \phi_\alpha(\mathbf{A}^\dagger y^{(t)})$ is much closer to $x^{(t)}$ than $\tilde{z}^{(t)}$. Since \mathbf{A} and \mathbf{A}^* are in the same subspace, inspired by (Li et al., 2016) we can write \mathbf{A}^* as $\mathbf{A}(\Sigma + \mathbf{E})$ for a diagonal matrix Σ and an off-diagonal matrix \mathbf{E} , and thus the decoding becomes $z^{(t)} = \phi_\alpha(\Sigma x^{(t)} + \mathbf{E}x^{(t)})$. Let us focus on one coordinate of $z^{(t)}$, that is, $z_i^{(t)} = \phi_\alpha(\Sigma_{i,i}x_i^{(t)} + \mathbf{E}_i x^{(t)})$, where \mathbf{E}_i is the i -th row of \mathbf{E} . The term $\Sigma_{i,i}x_i^{(t)}$ is a nice term since it is just a rescaling of $x_i^{(t)}$, while $\mathbf{E}_i x^{(t)}$ mixes different coordinates of $x^{(t)}$. For simplicity, we just assume for now that $x_i^{(t)} \in \{0, 1\}$ and $\Sigma_{i,i} = 1$. In our proof, we will show that the threshold will remove a large fraction of $\mathbf{E}_i x^{(t)}$ when $x_i^{(t)} = 0$, and keep a large fraction of $\Sigma_{i,i}x_i^{(t)}$ when $x_i^{(t)} = 1$. Thus, our decoding is much more accurate than without thresholding. To show this, we maintain a crucial property that for our decoding matrix, we always have $\|\mathbf{E}_i\|_2 = O(1)$. Assuming

this, we first consider two extreme cases of \mathbf{E}_i .

1. Ultra dense: all coordinates of \mathbf{E}_i are in the order of $\frac{1}{\sqrt{d}}$. Since the sparsity of $x^{(t)}$ is r , as long as $r = o(\sqrt{d})\alpha$, $\mathbf{E}_i x^{(t)}$ will not pass α and thus $z_i^{(t)}$ will be decoded to zero when $x_i^{(t)} = 0$.
2. Ultra sparse: \mathbf{E}_i only has few coordinate equal to $\Omega(1)$ and the rest are zero. Unless $x^{(t)}$ has those exact coordinates equal to 1 (which happens not so often), then $z_i^{(t)}$ will still be zero when $x_i^{(t)} = 0$.

Of course, the real \mathbf{E}_i can be anywhere in between these two extremes, and thus we need more delicate decoding lemmas, as shown in the complete proof.

Furthermore, more complication arises when each $x_i^{(t)}$ is not just in $\{0, 1\}$ but can take fractional values. To handle this case, we will set our threshold α to be large at the beginning and then keep shrinking after each stage. The intuition here is that we first decode the coordinates that we are most confident in, so we do not decode $z_i^{(t)}$ to be non-zero when $x_i^{(t)} = 0$. Thus, we will still be able to remove a large fraction of error caused by $\mathbf{E}_i x^{(t)}$. However, by setting the threshold α so high, we may introduce more errors to the nice term $\Sigma_{i,i}x_i^{(t)}$ as well, since $\Sigma_{i,i}x_i^{(t)}$ might not be larger than α when $x_i^{(t)} \neq 0$. Our main contribution is to show that there is a nice trade-off between the errors in \mathbf{E}_i terms and those in $\Sigma_{i,i}$ terms such that as we gradually decreases α , the algorithm can converge to the ground-truth.

5.2. Proof Sketch

For simplicity, we only focus on one stage and the expected update. The expected update of $\mathbf{A}^{(t)}$ is given by

$$\mathbf{A}^{(t+1)} = \mathbf{A}^{(t)} + \eta(\mathbb{E}[yz^\top] - \mathbf{A}^{(t)}\mathbb{E}[zz^\top]).$$

Let us write $\mathbf{A}^{(0)} = \mathbf{A}^*(\Sigma_0 + \mathbf{E}_0)$ where Σ_0 is diagonal and \mathbf{E}_0 is off-diagonal. Then the decoding is given by

$$z = \phi_\alpha((\mathbf{A}^{(0)})^\dagger y) = \phi_\alpha((\Sigma_0 + \mathbf{E}_0)^{-1}x).$$

Let Σ, \mathbf{E} be the diagonal part and the off-diagonal part of $(\Sigma_0 + \mathbf{E}_0)^{-1}$.

The key lemma for decoding says that under suitable conditions, z will be close to Σx in the following sense.

Lemma 2 (Decoding, informal). *Suppose \mathbf{E} is small and $\Sigma \approx \mathbf{I}$. Then with a proper threshold value α , we have*

$$\mathbb{E}[\Sigma x x^\top] \approx \mathbb{E}[z z^\top], \quad \mathbb{E}[\Sigma x z^\top] \approx \mathbb{E}[z z^\top].$$

Now, let us write $\mathbf{A}^{(t)} = \mathbf{A}^*(\Sigma_t + \mathbf{E}_t)$. Then applying the above decoding lemma, the expected update of $\Sigma_t + \mathbf{E}_t$ is

$$\Sigma_{t+1} + \mathbf{E}_{t+1} = (\Sigma_t + \mathbf{E}_t)(\mathbf{I} - \Sigma \Delta \Sigma) + \Sigma^{-1}(\Sigma \Delta \Sigma) + \mathbf{R}_t$$

where $\Delta = \mathbb{E}[xx^\top]$ and \mathbf{R}_t is a small error term.

Our second key lemma is about this update.

Lemma 3 (Update, informal). *Suppose the update rule is*

$$\Sigma_{t+1} + \mathbf{E}_{t+1} = (\Sigma_t + \mathbf{E}_t)(1 - \eta\Lambda) + \eta\mathbf{Q}\Lambda + \eta\mathbf{R}_t$$

for some PSD matrix Λ and $\|\mathbf{R}_t\|_2 \leq C''$. Then

$$\begin{aligned} \|\Sigma_t + \mathbf{E}_t - \mathbf{Q}\|_2 &\leq \|\Sigma_0 + \mathbf{E}_0 - \mathbf{Q}\|_2(1 - \eta\lambda_{\min}(\Lambda))^t \\ &\quad + \frac{C''}{\lambda_{\min}(\Lambda)}. \end{aligned}$$

Applying this on our update rule with $\mathbf{Q} = \Sigma^{-1}$ and $\Lambda = \Sigma\Delta\Sigma$, we know that when the error term is sufficiently small, we can make progress on $\|\Sigma_t + \mathbf{E}_t - \Sigma^{-1}\|_2$. Furthermore, by using the fact that $\Sigma_0 \approx \mathbf{I}$ and \mathbf{E}_0 is small, and the fact that Σ is the diagonal part of $(\Sigma_0 + \mathbf{E}_0)^{-1}$, we can show that after sufficiently many iterations, $\|\Sigma_t - \mathbf{I}\|_2$ blows up slightly, while $\|\mathbf{E}_t\|_2$ is reduced significantly. Repeating this for multiple stages completes the proof.

We note that most technical details are hidden, especially for the proofs of the decoding lemma, which need to show that the error term \mathbf{R}_t is small. This crucially relies on the choice of α , and relies on bounding the effect of the correlation. These then give the setting of α and the bound on the parameter m in the final theorem.

6. More General Results

6.1. Result for General x

This subsection considers the general case where $x \in [0, 1]^D$. Then the GCC condition is not enough for recovery, even for $k, r, m = O(1)$ and $\lambda = \Omega(1)$. For example, GCC does not rule out the case that x is drawn uniformly over $(r - 1)$ -sparse random vectors with $\{\frac{1}{D}, 1\}$ entries, when one cannot recover even a reasonable approximation of \mathbf{A}^* since a common vector $\frac{1}{D} \sum_i [\mathbf{A}^*]^i$ shows up in all the samples. This example shows that the difficulty arises if each x_i constantly shows up with a small value. To avoid this, a general and natural way is to assume that each x_i , once being non-zero, has to take a large value with sufficient probability. This is formalized as follows.

Definition 3 (Decay condition). *A distribution of x satisfies the order- q decay condition for some constant $q \geq 1$, if for all $i \in [D]$, x_i satisfies that for every $\alpha > 0$,*

$$\Pr[x_i \leq \alpha \mid x_i \neq 0] \leq \alpha^q.$$

When $q = 1$, each x_i , once being non-zero, is uniformly distributed in the interval $[0, 1]$. When q gets larger, each x_i , once being non-zero, will be more likely to take larger

values. We will show that our algorithm has a better guarantee for larger q . In the extreme case when $q = \infty$, x_i will only take $\{0, 1\}$ values, which reduces to the binary case.

In this paper, we show that this simple decay condition, combined with the GCC conditions and an initialization with constant error, is sufficient for recovering \mathbf{A}^* .

Theorem 4 (Main). *There exists $\ell = \Omega(1)$ such that for every (r, k, m, λ) -GCC x satisfying the order- q condition, every $\epsilon > 0$, there exists T, η and a sequence of $\{\alpha_i\}^4$ such that Algorithm AND, with ℓ -initialization matrix \mathbf{A}_0 , outputs a matrix \mathbf{A} such that there exists a diagonal matrix $\Sigma \succeq \frac{1}{2}\mathbf{I}$ with $\|\mathbf{A} - \mathbf{A}^*\Sigma\|_2 \leq \epsilon$ with $\text{poly}(D, \frac{1}{\epsilon})$ samples and iterations, as long as*

$$m = O\left(\frac{kD^{1-\frac{1}{q}}\lambda^{4+\frac{4}{q}}}{r^{5+\frac{6}{q+1}}}\right).$$

As mentioned, in many interesting applications, $k = r = \lambda = \Theta(1)$, where our algorithm can recover \mathbf{A}^* as long as $m = O(D^{1-\frac{1}{q+1}})$. This means $\mathbb{E}_\mu[x_i x_j] = O(D^{-1-\frac{1}{q+1}})$, a factor of $D^{-\frac{1}{q+1}}$ away from the highest possible correlation $\min\{\mathbb{E}_\mu[x_i], \mathbb{E}_\mu[x_j]\} = O(1/D)$. Then, the larger q , the higher correlation it can tolerate. As q goes to infinity, we recover the result for the case $x \in \{0, 1\}^D$, allowing the highest order correlation.

The analysis also shows that the decoding threshold should be $\alpha = \left(\frac{\lambda\|\mathbf{E}_0\|_2}{r}\right)^{\frac{2}{q+1}}$ where \mathbf{E}_0 is the error matrix at the beginning of the stage. Since the error decreases exponentially with stages, this suggests to decrease α exponentially with stages. This is crucial for AND to recover the ground-truth; see Section 7 for the experimental results.

6.2. Robustness to Noise

We now consider the case when the data is generated from $y = \mathbf{A}^*x + \zeta$, where ζ is the noise. For the sake of demonstration, we will just focus on the case when $x_i \in \{0, 1\}$ and ζ is random Gaussian noise $\zeta \sim \gamma\mathcal{N}(0, \frac{1}{W}\mathbf{I})$.⁵ A more general theorem can be found in the appendix.

Definition 4 ((ℓ, ρ) -initialization). *The initial matrix \mathbf{A}_0 satisfies for some $\ell, \rho \in [0, 1]$,*

1. $\mathbf{A}_0 = \mathbf{A}^*(\Sigma + \mathbf{E}) + \mathbf{N}$, for some diagonal matrix Σ and off-diagonal matrix \mathbf{E} .
2. $\|\mathbf{E}\|_2 \leq \ell$, $\|\Sigma - \mathbf{I}\|_2 \leq \frac{1}{4}$, $\|\mathbf{N}\|_2 \leq \rho$.

Theorem 5 (Noise, binary). *Suppose each $x_i \in \{0, 1\}$. There exists $\ell = \Omega(1)$ such that for every (r, k, m, λ) -GCC x , every $\epsilon > 0$, Algorithm AND with $T = \text{poly}(D, \frac{1}{\epsilon})$, $\eta =$*

⁴In fact, we will make the choice explicit in the proof.

⁵we make this scaling so $\|\zeta\|_2 \approx \gamma$.

$\frac{1}{\text{poly}(D, \frac{1}{\epsilon})}$, $\{\alpha_i\}_{i=1}^s = \{\frac{1}{4}\}_{i=1}^4$ and an (ℓ, ρ) -initialization \mathbf{A}_0 for $\rho = O(\sigma_{\min}(\mathbf{A}^*))$, outputs \mathbf{A} such that there exists a diagonal matrix $\Sigma \succeq \frac{1}{2}\mathbf{I}$ with

$$\|\mathbf{A} - \mathbf{A}^*\Sigma\|_2 \leq O\left(\epsilon + r \frac{\sigma_{\max}(\mathbf{A}^*)}{\sigma_{\min}(\mathbf{A}^*)\lambda}\gamma\right)$$

using $\text{poly}(D, \frac{1}{\epsilon})$ iterations, as long as $m = O\left(\frac{kD\lambda^4}{r^5}\right)$.

The theorem implies that the algorithm can recover the ground-truth up to $r \frac{\sigma_{\max}(\mathbf{A}^*)}{\sigma_{\min}(\mathbf{A}^*)\lambda}$ times γ , the noise level in each sample. Although stated here for Gaussian noise for simplicity, the analysis applies to a much larger class of noises, including adversarial ones. In particular, we only need to the noise ζ have sufficiently bounded $\|\mathbb{E}[\zeta\zeta^\top]\|_2$; see the appendix for the details. For the special case of Gaussian noise, by exploiting its properties, it is possible to improve the error term with a more careful calculation, though not done here.

7. Experiments

To demonstrate the advantage of AND, we complement the theoretical analysis with empirical study on semi-synthetic datasets, where we have ground-truth feature matrices and can thus verify the convergence. We then provide support for the benefit of using decreasing thresholds, and test its robustness to noise. In the appendix, we further test its robust to initialization and sparsity of x , and provide qualitative results in some real world applications.⁶

Setup. Our work focuses on convergence of the solution to the ground-truth feature matrix. However, real-world datasets in general do not have ground-truth. So we construct semi-synthetic datasets in topic modeling: first take the word-topic matrix learned by some topic modeling method as the ground-truth \mathbf{A}^* , and then draw x from some specific distribution μ . For fair comparison, we use one not learned by any algorithm evaluated here. In particular, we used the matrix with 100 topics computed by the algorithm in (Arora et al., 2013) on the NIPS papers dataset (about 1500 documents, average length about 1000). Based on this we build two semi-synthetic datasets:

1. DIR. Construct a 100×5000 matrix \mathbf{X} , whose columns are from a Dirichlet distribution with parameters $(0.05, 0.05, \dots, 0.05)$. Then the dataset is $\mathbf{Y} = \mathbf{A}^*\mathbf{X}$.
2. CTM. The matrix \mathbf{X} is of the same size as above, while each column is drawn from the logistic normal prior in the correlated topic model (Blei & Lafferty, 2006). This leads to a dataset with strong correlations.

⁶The code is public on <https://github.com/PrincetonML/AND4NMF>.

Note that the word-topic matrix is non-negative. While some competitor algorithms require a non-negative feature matrix, AND does not need such a condition. To demonstrate this, we generate the following synthetic data:

3. NEG. The entries of the matrix \mathbf{A}^* are i.i.d. samples from the uniform distribution on $[-0.5, 0.5]$. The matrix \mathbf{X} is the same as in CTM.

Finally, the following dataset is for testing the robustness of AND to the noise:

4. NOISE. \mathbf{A}^* and \mathbf{X} are the same as in CTM, but $\mathbf{Y} = \mathbf{A}^*\mathbf{X} + \mathbf{N}$ where \mathbf{N} is the noise matrix with columns drawn from $\gamma\mathcal{N}(0, \frac{1}{W}\mathbf{I})$ with the noise level γ .

Competitors. We compare the algorithm AND to the following popular methods: Alternating Non-negative Least Square (ANLS (Kim & Park, 2008)), multiplicative update (MU (Lee & Seung, 2001)), LDA (online version (Hoffman et al., 2010)),⁷ and Hierarchical Alternating Least Square (HALS (Cichocki et al., 2007)).

Evaluation criterion. Given the output matrix \mathbf{A} and the ground truth matrix \mathbf{A}^* , the *correlation error* of the i -th column is given by

$$\varepsilon_i(\mathbf{A}, \mathbf{A}^*) = \min_{j \in [D], \sigma \in \mathbb{R}} \{ \|\sigma[\mathbf{A}^*]^i - [\mathbf{A}]^j\|_2 \}.$$

Thus, the error measures how well the i -th column of \mathbf{A}^* is covered by the best column of \mathbf{A} up to scaling. We find the best column since in some competitor algorithms, the columns of the solution \mathbf{A} may only correspond to a permutation of the columns of \mathbf{A}^* .⁸

We also define the *total correlation error* as

$$\varepsilon(\mathbf{A}, \mathbf{A}^*) = \sum_{i=1}^D \varepsilon_i(\mathbf{A}, \mathbf{A}^*).$$

We report the total correlation error in all the experiments.

Initialization. In all the experiments, the initialization matrix \mathbf{A}_0 is set to $\mathbf{A}_0 = \mathbf{A}^*(\mathbf{I} + \mathbf{U})$ where \mathbf{I} is the identity matrix and \mathbf{U} is a matrix whose entries are i.i.d. samples from the uniform distribution on $[-0.05, 0.05]$. Note that this is a very weak initialization, since $[\mathbf{A}_0]^i = (1 + \mathbf{U}_{i,i})[\mathbf{A}^*]^i + \sum_{j \neq i} \mathbf{U}_{j,i}[\mathbf{A}^*]^j$ and the magnitude of the noise component $\sum_{j \neq i} \mathbf{U}_{j,i}[\mathbf{A}^*]^j$ can be larger than the signal part $(1 + \mathbf{U}_{i,i})[\mathbf{A}^*]^i$.

⁷We use the implementation in the sklearn package (<http://scikit-learn.org/>)

⁸In the Algorithm AND, the columns of \mathbf{A} correspond to the columns of \mathbf{A}^* without permutation.

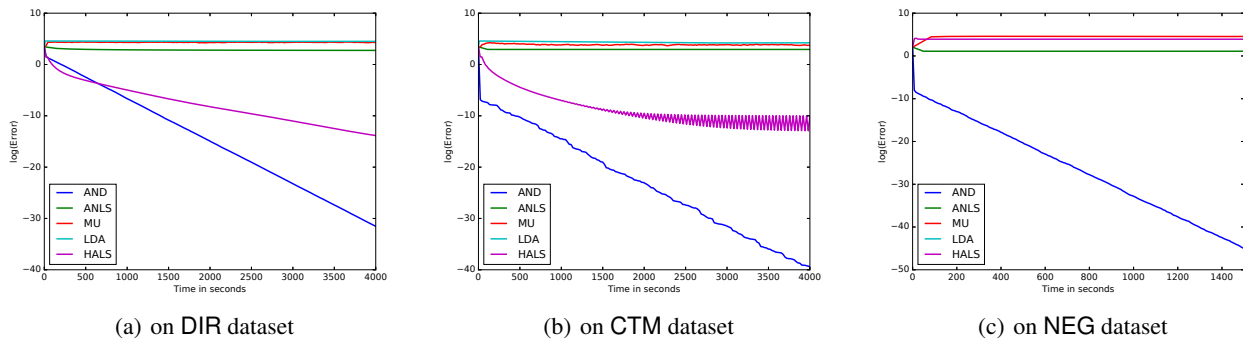


Figure 1. The performance of different algorithms on the three datasets. The x -axis is the running time (in seconds), the y -axis is the logarithm of the total correlation error.

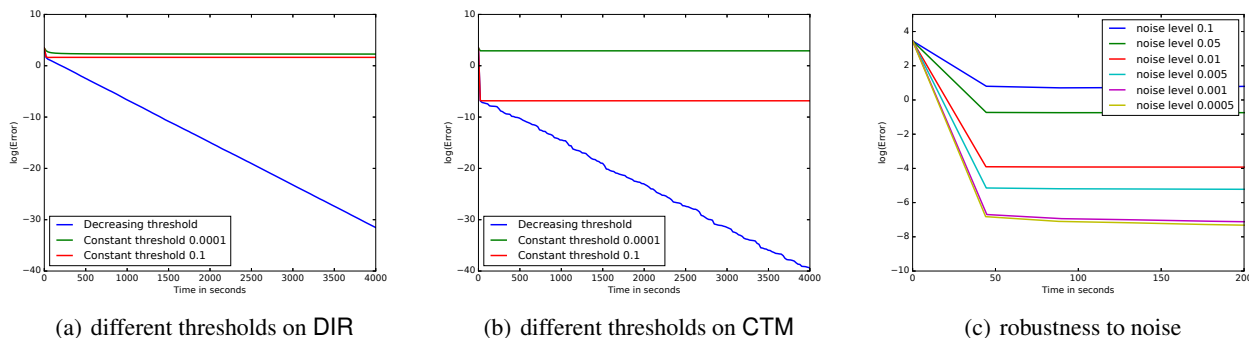


Figure 2. The performance of the algorithm AND with different thresholding schemes, and its robustness to noise. The x -axis is the running time (in seconds), the y -axis is the logarithm of the total correlation error. (a)(b) Using different thresholding schemes on the DIR/CTM dataset. “Decreasing threshold” refers to the scheme used in the original AND, “Constant threshold c ” refers to using the threshold value c throughout all iterations. (c) The performance in the presence of noises of various levels.

Hyperparameters and Implementations. For most experiments of AND, we used $T = 50$ iterations for each stage, and thresholds $\alpha_i = 0.1/(1.1)^{i-1}$. For experiments on the robustness to noise, we found $T = 100$ leads to better performance. Furthermore, for all the experiments, instead of using one data point at each step, we used the whole dataset for update.

7.1. Convergence to the Ground-Truth

Figure 1 shows the convergence rate of the algorithms on the three datasets. AND converges in linear rate on all three datasets (note that the y -axis is in log-scale). HALS converges on the DIR and CTM datasets, but the convergence is in slower rates. Also, on CTM, the error oscillates. Furthermore, it doesn’t converge on NEG where the ground-truth matrix has negative entries. ANLS converges on DIR and CTM at a very slow speed due to the non-negative least square computation in each iteration.⁹ All the other algo-

⁹We also note that even the thresholding of HALS and ANLS designed for non-negative feature matrices is removed, they still

do not converge to the ground-truth, suggesting that they do not have recovery guarantees.

7.2. The Threshold Schemes

Figure 2(a) shows the results of using different thresholding schemes on DIR, while Figure 2(b) shows that those on CTM. When using a constant threshold for all iterations, the error only decreases for the first few steps and then stop decreasing. This aligns with our analysis and is in strong contrast to the case with decreasing thresholds.

7.3. Robustness to Noise

Figure 2(c) shows the performance of AND on the NOISE dataset with various noise levels γ . The error drops at the first few steps, but then stabilizes around a constant related to the noise level, as predicted by our analysis. This shows that it can recover the ground-truth to good accuracy, even when the data have a significant amount of noise.

do not converge on NEG.

Acknowledgements

This work was supported in part by NSF grants CCF-1527371, DMS-1317308, Simons Investigator Award, Simons Collaboration Grant, and ONR-N00014-16-1-2329. This work was done when Yingyu Liang was visiting the Simons Institute.

References

- Lda-c software. <https://github.com/blei-lab/lda-c/blob/master/readme.txt>, 2016. Accessed: 2016-05-19.
- Anandkumar, A., Kakade, S., Foster, D., Liu, Y., and Hsu, D. Two svds suffice: Spectral decompositions for probabilistic topic modeling and latent dirichlet allocation. Technical report, 2012.
- Anandkumar, A., Hsu, D., Javanmard, A., and Kakade, S. Learning latent bayesian networks and topic models under expansion constraints. In *ICML*, 2013.
- Arora, S., Ge, R., and Moitra, A. Learning topic models – going beyond svd. In *FOCS*, 2012a.
- Arora, S., Ge, R., Halpern, Y., Mimno, D., Moitra, A., Sontag, D., Wu, Y., and Zhu, M. A practical algorithm for topic modeling with provable guarantees. In *ICML*, 2013.
- Arora, S., Ge, R., Ma, T., and Moitra, A. Simple, efficient, and neural algorithms for sparse coding. In *COLT*, 2015.
- Arora, Sanjeev, Ge, Rong, Kannan, Ravindran, and Moitra, Ankur. Computing a nonnegative matrix factorization–provably. In *STOC*, pp. 145–162. ACM, 2012b.
- Arora, Sanjeev, Ge, Rong, Koehler, Frederic, Ma, Tengyu, and Moitra, Ankur. Provable algorithms for inference in topic models. In *Proceedings of The 33rd International Conference on Machine Learning*, pp. 2859–2867, 2016.
- Awasthi, Pranjali and Risteski, Andrej. On some provably correct cases of variational inference for topic models. In *NIPS*, pp. 2089–2097, 2015.
- Bhattacharyya, Chiranjib, Goyal, Navin, Kannan, Ravindran, and Pani, Jagdeep. Non-negative matrix factorization under heavy noise. In *Proceedings of the 33rd International Conference on Machine Learning*, 2016.
- Blei, David and Lafferty, John. Correlated topic models. *Advances in neural information processing systems*, 18: 147, 2006.
- Blei, David M. Probabilistic topic models. *Communications of the ACM*, 2012.
- Blei, David M, Ng, Andrew Y, and Jordan, Michael I. Latent dirichlet allocation. *JMLR*, 3:993–1022, 2003.
- Cichocki, Andrzej, Zdunek, Rafal, and Amari, Shun-ichi. Hierarchical als algorithms for nonnegative matrix and 3d tensor factorization. In *International Conference on Independent Component Analysis and Signal Separation*, pp. 169–176. Springer, 2007.
- Ding, W., Rohban, M.H., Ishwar, P., and Saligrama, V. Topic discovery through data dependent and random projections. *arXiv preprint arXiv:1303.3664*, 2013.
- Ding, W., Rohban, M.H., Ishwar, P., and Saligrama, V. Efficient distributed topic modeling with provable guarantees. In *AISTAT*, pp. 167–175, 2014.
- Hoffman, Matthew, Bach, Francis R, and Blei, David M. Online learning for latent dirichlet allocation. In *advances in neural information processing systems*, pp. 856–864, 2010.
- Kim, Hyunsoo and Park, Haesun. Nonnegative matrix factorization based on alternating nonnegativity constrained least squares and active set method. *SIAM journal on matrix analysis and applications*, 30(2):713–730, 2008.
- Lee, Daniel D and Seung, H Sebastian. Unsupervised learning by convex and conic coding. *NIPS*, pp. 515–521, 1997.
- Lee, Daniel D and Seung, H Sebastian. Learning the parts of objects by non-negative matrix factorization. *Nature*, 401(6755):788–791, 1999.
- Lee, Daniel D and Seung, H Sebastian. Algorithms for non-negative matrix factorization. In *NIPS*, pp. 556–562, 2001.
- Li, Yuanzhi, Liang, Yingyu, and Risteski, Andrej. Recovery guarantee of non-negative matrix factorization via alternating updates. *Advances in neural information processing systems*, 2016.
- Yang, Jaewon and Leskovec, Jure. Overlapping community detection at scale: a nonnegative matrix factorization approach. In *Proceedings of the sixth ACM international conference on Web search and data mining*, pp. 587–596. ACM, 2013.