
Algorithms for ℓ_p Low-Rank Approximation

Flavio Chierichetti¹ Sreenivas Gollapudi² Ravi Kumar² Silvio Lattanzi³ Rina Panigrahy²
David P. Woodruff⁴

Abstract

We consider the problem of approximating a given matrix by a low-rank matrix so as to minimize the entry-wise ℓ_p -approximation error, for any $p \geq 1$; the case $p = 2$ is the classical SVD problem. We obtain the first provably good approximation algorithms for this version of low-rank approximation that work for every value of $p \geq 1$, including $p = \infty$. Our algorithms are simple, easy to implement, work well in practice, and illustrate interesting tradeoffs between the approximation quality, the running time, and the rank of the approximating matrix.

1. Introduction

The problem of low-rank approximation of a matrix is usually studied as approximating a given matrix by a matrix of low rank so that the Frobenius norm of the error in the approximation is minimized. The Frobenius norm of a matrix is obtained by taking the sum of the squares of the entries in the matrix. Under this objective, the optimal solution is obtained using the singular value decomposition (SVD) of the given matrix. Low-rank approximation is useful in large data analysis, especially in predicting missing entries of a matrix by projecting the row and column entities (e.g., users and movies) into a low-dimensional space. In this work we consider the low-rank approximation problem, but under the general entry-wise ℓ_p norm, for any $p \in [1, \infty]$.

There are several reasons for considering the ℓ_p version of

^{*}Equal contribution ¹Sapienza University, Rome, Italy. Work done in part while visiting Google. Supported in part by a Google Focused Research Award, by the ERC Starting Grant DMAP 680153, and by the SIR Grant RBSI14Q743. ²Google, Mountain View, CA ³Google, Zurich, Switzerland ⁴IBM Almaden, San Jose, CA. Correspondence to: Flavio Chierichetti <flavio@di.uniroma1.it>, Sreenivas Gollapudi <sgollapu@yahoo.com>, Ravi Kumar <ravi.k53@gmail.com>, Silvio Lattanzi <silviol@google.com>, Rina Panigrahy <rinapy@gmail.com>, David P. Woodruff <dpwoodru@us.ibm.com>.

low-rank approximation instead of the usually studied ℓ_2 (i.e., Frobenius) version. For example, it is widely acknowledged that the ℓ_1 version is more robust to noise and outliers than the ℓ_2 version (Candès et al., 2011; Huber, 1981; Xu & Yuille, 1995). Several data mining and computer vision-related applications exploit this insight and resort to finding a low-rank approximation to minimize the ℓ_1 error (Lu et al., 2014; Meng & Torre, 2013; Wang & Yeung, 2013; Xiong et al., 2011). Furthermore, the ℓ_1 error is typically used as a proxy for capturing sparsity in many applications including robust versions of PCA, sparse recovery, and matrix completion; see, for example (Candès et al., 2011; Xu et al., 2012). For these reasons the problem has already received attention (Gillis & Vavasis, 2015) and was suggested as an open question in a survey on sketching techniques for linear algebra (Woodruff, 2014). Likewise, the ℓ_∞ version (dubbed also as the Chebyshev norm) has been studied for the past many years (Goreinov & Tyrtshnikov, 2001; 2011), though to the best of our knowledge, no result with theoretical guarantees was known for ℓ_∞ before our work. Our algorithm is quite general, and works for every $p \geq 1$.

Working with ℓ_p error, however, poses many technical challenges. First of all, unlike ℓ_2 , the general ℓ_p space is not amenable to spectral techniques. Secondly, the ℓ_p space is not as nicely behaved as the ℓ_2 space, for example, it lacks the notion of orthogonality. Thirdly, the ℓ_p version quickly runs into computational complexity barriers: for example, even the rank-1 approximation in ℓ_1 has been shown to be NP-hard by Gillis and Vavasis (Gillis & Vavasis, 2015). However, there has been no dearth in terms of heuristics for the ℓ_p low-rank approximation problem, in particular for $p = 1$ and $p = \infty$: this includes alternating convex (and, in fact, linear) minimization (Ke & Kanade, 2005), methods based on expectation-maximization (Wang et al., 2012), minimization with augmented Lagrange multipliers (Zheng et al., 2012), hyperplanes projections and linear programming (Brooks et al., 2013), and generalizations of the Wiberg algorithm (Eriksson & van den Hengel, 2012). These heuristics, unfortunately, do not come with any performance guarantees. While theoretical approximation guarantees have been given for the rank-1 version for the GF(2) and the Boolean cases (Dan et al., 2015), to the best of our knowledge there have been no provably good

(approximation) algorithms for general matrices, or for rank more than one, or for general ℓ_p .

1.1. Our contributions

In this paper we obtain the first provably good algorithms for the ℓ_p rank- k approximation problem for every $p \geq 1$. Let $n \times m$ be the dimensions of the input matrix. From an algorithmic viewpoint, there are three quantities of interest: the running time of the algorithm, the approximation factor guaranteed by the algorithm, and the actual number of vectors in the low-rank approximation that is output by the algorithm (even though we only desire k).

Given this setting, we show three main algorithmic results intended for the case when k is not too large. First, we show that one can obtain a $(k + 1)$ -approximation to the rank- k problem in time $m^k \text{poly}(n, m)$; note that this running time is not polynomial once k is larger than a constant. To address this, next we show that one can get an $O(k)$ -approximation to the best k -factorization in time $O(\text{poly}(nm))$; however, the algorithm returns $O(k \log m)$ columns, which is more than the desired k (this is referred to as a bi-criteria approximation). Next, we combine these two algorithms. We first show that the output of the second algorithm can further be refined to output exactly k vectors, with an approximation factor of $\text{poly}(k)$ and a running time of $O(\text{poly}(n, m)(k \log n)^k)$. The running time now is polynomial as long as $k = O(\log n / \log \log n)$. Finally, we show that for any constant $p \geq 1$, we can obtain an approximation factor of $(k \log m)^{O(p)}$ and a running time of $\text{poly}(n, m)$ for every value of k .

Our first algorithm is existential in nature: it shows that there are k columns in the given matrix that can be used, along with an appropriate convex program, to obtain a $(k + 1)$ -approximation. Realizing this as an algorithm would therefore naïvely incur a factor m^k in the running time. Our second algorithm works by sampling columns and iteratively “covering” the columns of the given matrix, for an appropriate notion of covering. In each round of sampling our algorithm uniformly samples from a remaining set of columns; we note here that it is critical that our algorithm is adaptive as otherwise uniform sampling would not work. While this is computationally efficient and maintains an $O(k)$ -approximation to the best rank- k approximation, it can end up with more than k columns, in fact $O(k \log m)$. Our third algorithm fixes this issue by combining the first algorithm with the notion of a near-isoperimetric transformation for the ℓ_p -space, which lets us transform a given matrix into another matrix spanning the same subspace but with small ℓ_p distortion. Our fourth algorithm uses ℓ_p leverage scores to overcome the sampling step; this improves the running time while mildly worsening the approximation factor.

A useful feature of our algorithms is that they are uniform with respect to all values of p . We test the performance of our algorithms, for $p = 1$ and $p = \infty$, on real and synthetic data and show that they produce low-rank approximations that are substantially better than what the SVD (i.e., $p = 2$) would obtain.

1.2. Related work

Very recently, Song, Woodruff, and Zhong (Song et al., 2017), obtained a low-rank approximation that holds for every $p \in [1, 2]$, using sketching-based techniques. In particular, their main result is an $(O(\log m) \text{poly}(k))$ -approximation in $\text{nnz}(A) + (n + m) \text{poly}(k)$ time, for every k , where $\text{nnz}(A)$ is the number of non-zero entries in A . In our work, we also obtain such a result for $p \in [1, 2]$ but via very different sampling-based methods. In addition, we obtain an algorithm with a $\text{poly}(k)$ -approximation factor that is independent of m and n , though this latter algorithm requires $k = O(\log n / \log \log n)$ in order to be polynomial time. Another result in (Song et al., 2017) shows how to achieve a $(k \text{poly}(\log k))$ -approximation, in $n^{O(k)}$ time for $p \in [1, 2]$. For k larger than a constant, this is larger than polynomial time, whereas our algorithm with $\text{poly}(k)$ -approximation is polynomial time for k as large as $\Theta(\log n / \log \log n)$. Importantly, our results hold for every $p \in [1, \infty]$, rather than only when $p \in [1, 2]$.

In addition we note that there exist papers solving problems that, at first blush, might seem similar to ours. For instance, (Deshpande et al., 2011) study a convex relation, and a rounding algorithm to solve the subspace approximation problem (an ℓ_p generalization of the least squares fit), which is related to but different from our problem. Also, (Feldman et al., 2007) offer a bi-criteria solution for another related problem of approximating a set of points by a collection of flats; they use convex relaxations to solve their problem and are limited to bi-criteria solutions, unlike ours. Finally, in some special settings robust PCA can be used to solve ℓ_1 low-rank approximation (Candès et al., 2011). However, robust PCA and ℓ_1 low-rank approximation have some apparent similarities but they have key differences. Firstly, ℓ_1 low-rank approximation allows to recover an approximating matrix of any chosen rank, whereas robust PCA returns some matrix of some unknown (possibly full) rank. While variants of robust PCA have been proposed to force the output rank to be a given value (Netrapalli et al., 2014; Yi et al., 2016), these variants make additional noise model and incoherence assumptions on the input matrix, whereas our results hold for every input matrix. Secondly, in terms of approximation quality, it is unclear if near-optimal solutions of robust PCA provide near-optimal solutions for ℓ_1 low-rank approximation.

Finally, we mention example matrices for which the SVD

gives a poor approximation factor for ℓ_p -approximation error. First, suppose $p < 2$ and $k = 1$. Consider the following $n \times n$ block diagonal matrix composed of two blocks: a 1×1 matrix with value n and an $(n-1) \times (n-1)$ matrix with all 1s. The SVD returns as a solution the first column, and therefore incurs polynomial in n error for $p = 2 - \Omega(1)$. Now suppose $p > 2$ and $k = 1$. Consider the following $n \times n$ block diagonal matrix composed of two blocks: a 1×1 matrix with value $n-2$ and an $(n-1) \times (n-1)$ matrix with all 1s. The SVD returns as a solution the matrix spanned by the bottom block, and so also incurs an error polynomial in n for $p = 2 + \Omega(1)$.

2. Background

For a matrix M , let $M_{i,j}$ denote the entry in its i th row and j th column and let M_i denote its i th column. Let M^T denote its transpose and let $|M|_p = \left(\sum_{i,j} |M_{i,j}|^p\right)^{1/p}$ denote its entry-wise p norm. Given a set $S = \{i_1, \dots, i_t\}$ of column indices, let $M_S = M_{i_1, \dots, i_t}$ be the matrix composed of the columns of M with the indices in S .

Given a matrix M with m columns, we will use $\text{span } M = \{\sum_{i=1}^m \alpha_i M_i \mid \alpha_i \in \mathbb{R}\}$ to denote the vectors spanned by its columns. If M is a matrix and v is a vector, we let $d_p(v, M)$ denote the minimum ℓ_p distance between v and a vector in $\text{span } M$:

$$d_p(v, M) = \inf_{w \in \text{span } M} |v - w|_p.$$

Let $A \in \mathbb{R}^{n \times m}$ denote the input matrix and let $k > 0$ denote the target rank. We assume, without loss of generality (wlog.), that $m \leq n$. Our first goal is, given A and k , to find a subset $U \in \mathbb{R}^{n \times k}$ of k columns of A and $V \in \mathbb{R}^{k \times m}$ to minimize the ℓ_p error, $p \geq 1$, given by

$$|A - UV|_p.$$

Our second goal is, given A and k , to find $U \in \mathbb{R}^{n \times k}$, $V \in \mathbb{R}^{k \times m}$ to minimize the ℓ_p error, $p \geq 1$, given by

$$|A - UV|_p.$$

Note that in the second goal, we do not require U be a subset of columns.

We refer to the first problem as the k -columns subset selection problem in the ℓ_p norm, denoted k -CSS $_p$, and to the second problem as the rank- k approximation problem in the ℓ_p norm, denoted k -LRA $_p$.¹ In the paper we often call U, V the k -factorization of A . Note that a solution to k -CSS $_p$ can be used as a solution to k -LRA $_p$, but not necessarily vice versa.

¹ k -LRA $_2$ is the classical SVD problem of finding $U \in \mathbb{R}^{n \times k}$, $V \in \mathbb{R}^{k \times m}$ so as to minimize $|A - UV|_2$.

In this paper we focus on solving the two problems for general p . Let U^*V^* be a k -factorization of A that is optimal in the ℓ_p norm, where $U^* \in \mathbb{R}^{n \times k}$ and $V^* \in \mathbb{R}^{k \times m}$, and let $\text{opt}_{k,p}(A) = |A - U^*V^*|_p$. An algorithm is said to be an α -approximation, for an $\alpha \geq 1$, if it outputs $U \in \mathbb{R}^{n \times k}$, $V \in \mathbb{R}^{k \times m}$ such that

$$|A - UV|_p \leq \alpha \cdot \text{opt}_{k,p}(A).$$

It is often convenient to view the input matrix as $A = U^*V^* + \Delta = A^* + \Delta$, where Δ is some error matrix of minimum ℓ_p -norm. Let $\delta = |\Delta|_p = \text{opt}_{k,p}(A)$.

We will use the following observation.

Lemma 1. *Let $U \in \mathbb{R}^{n \times k}$ and $v \in \mathbb{R}^{n \times 1}$. Suppose that there exists $x \in \mathbb{R}^{k \times 1}$ such that $\delta = |U \cdot x - v|_p$. Then, there exists a polynomial time algorithm that, given U and v , finds $y \in \mathbb{R}^{k \times 1}$ such that $|U \cdot y - v|_p \leq \delta$.*

Proof. This ℓ_p regression problem is a convex program and well-known to be solvable in polynomial time. \square

3. An $(m^k \text{ poly}(nm))$ -time algorithm for k -LRA $_p$

In this section we will present an algorithm that runs in time $m^k \text{ poly}(nm)$ and produces a $(k+1)$ -approximation to k -CSS $_p$ (hence to k -LRA $_p$) of a matrix $A \in \mathbb{R}^{n \times m}$, $m \leq n$, for any $p \in [1, \infty]$. The algorithm simply tries all possible subsets of k columns of A for producing one of the factors, U , and then uses Lemma 1 to find the second factor V .

3.1. The existence of one factor in A

For simplicity, we assume that $|\Delta_i|_p > 0$ for each column i . To satisfy this, we can add an arbitrary small random error to each entry of the matrix. For instance, for any $\gamma > 0$, and to each entry of the matrix, we can add an independent uniform value in $[-\gamma, \gamma]$. This would guarantee that $|\Delta_i|_p > 0$ for each $i \in [m]$.

Recall that $A = A^* + \Delta$ is the perturbed matrix, and we only have access to A , not A^* . Consider Algorithm 1 and its output S . Note that we cannot actually run this algorithm since we do not know A^* . It is a hypothetical algorithm used for the purpose of our proof, i.e., the algorithm serves as a proof that there exists a subset of k columns of A providing a good low rank approximation. In Theorem 3 we prove that the columns in A indexed by the subset S can be used as one factor of a k -factorization of A .

Before proving the main theorem of the section, we show a useful property of the matrix \tilde{A}^* , i.e., the matrix having the vector $A_i^*/|\Delta_i|_p$ as the i th column. Then we will use this property to prove Theorem 3.

Algorithm 1 Enumerating and selecting k columns of A .

Require: A rank k matrix A^* and perturbation matrix Δ
Ensure: k column indices of A^*

- 1: For each column index i , let $\tilde{A}_i^* \leftarrow A_i^*/|\Delta_i|_p$.
- 2: Write $\tilde{A}^* = \tilde{U} \cdot \tilde{V}$, s.t. $\tilde{U} \in \mathbb{R}^{n \times k}$, $\tilde{V} \in \mathbb{R}^{k \times m}$.
- 3: Let S be the subset of k columns of $\tilde{V} \in \mathbb{R}^{k \times m}$ that has maximum determinant in absolute value (note that the subset S indexes a $k \times k$ submatrix).
- 4: Output S .

Lemma 2. For each column \tilde{A}_i^* of \tilde{A}^* , one can write $\tilde{A}_i^* = \sum_{j \in S} M_i(j) \tilde{A}_j^*$, where $|M_i(j)| \leq 1$ for all i, j .

Proof. Fix an $i \in \{1, \dots, m\}$. Consider the equation $\tilde{V}_S M_i = \tilde{V}_i$ for $M_i \in \mathbb{R}^k$. We can assume the columns in \tilde{V}_S are linearly independent, since wlog., \tilde{A}^* has rank k . Hence, there is a unique solution $M_i = (\tilde{V}_S)^{-1} \tilde{V}_i$. By Cramer's rule, the j th coordinate $M_i(j)$ of M_i satisfies $M_i(j) = \frac{\det(\tilde{V}_S^j)}{\det(\tilde{V}_S)}$, where \tilde{V}_S^j is the matrix obtained by replacing the j th column of \tilde{V}_S with \tilde{V}_i . By our choice of S , $|\det(\tilde{V}_S^j)| \leq |\det(\tilde{V}_S)|$, which implies $|M_i(j)| \leq 1$. Multiplying both sides of equation $\tilde{V}_S M_i = \tilde{V}_i$ by \tilde{U} , we have $\tilde{A}_S^* M_i = \tilde{A}_i^*$. \square

Now we prove the main theorem of this section.

Theorem 3. Let $U = A_S$. For $p \in [1, \infty]$, let M_1, \dots, M_m be the vectors whose existence is guaranteed by Lemma 2 and let $V \in \mathbb{R}^{k \times n}$ be the matrix having the vector $|\Delta_i|_p \cdot (M_i(1)/|\Delta_{i_1}|_p, \dots, M_i(k)/|\Delta_{i_k}|_p)^T$ as its i th column. Then, $|A_i - (UV)_i|_p \leq (k+1)|\Delta_i|_p$ and hence $|A - UV|_p \leq (k+1)|\Delta|_p$.

Proof. We consider the generic column $(UV)_i$.

$$\begin{aligned}
 (UV)_i &= |\Delta_i|_p \sum_{j=1}^k \left(\frac{M_i(j)}{|\Delta_{i_j}|_p} A_{i_j} \right) \\
 &= |\Delta_i|_p \sum_{j=1}^k \left(\frac{M_i(j)}{|\Delta_{i_j}|_p} (A_{i_j}^* + \Delta_{i_j}) \right) \\
 &= |\Delta_i|_p \sum_{j=1}^k \left(M_i(j) \tilde{A}_{i_j}^* + M_i(j) \frac{\Delta_{i_j}}{|\Delta_{i_j}|_p} \right) \\
 &= |\Delta_i|_p \tilde{A}_i^* + |\Delta_i|_p \sum_{j=1}^k \left(M_i(j) \frac{\Delta_{i_j}}{|\Delta_{i_j}|_p} \right) \\
 &= A_i^* + \sum_{j=1}^k \left(|\Delta_i|_p \cdot M_i(j) \frac{\Delta_{i_j}}{|\Delta_{i_j}|_p} \right) \triangleq A_i^* + E_i.
 \end{aligned}$$

Observe that E_i is the weighted sum of k vectors, $\frac{\Delta_{i_1}}{|\Delta_{i_1}|_p}, \dots, \frac{\Delta_{i_k}}{|\Delta_{i_k}|_p}$, having unit ℓ_p -norm. Observe fur-

ther that, since the sum of their weights satisfies $|\Delta_i|_p \sum_{j=1}^k |M_i(j)| \leq k|\Delta_i|_p$, we have that the ℓ_p -norm of E_i is not larger than $|E_i|_p \leq k|\Delta_i|_p$. The proof is complete using the triangle inequality:

$$\begin{aligned}
 |A_i - (UV)_i|_p &\leq |A_i^* - A_i|_p + |A_i^* - (UV)_i|_p \\
 &= |\Delta_i|_p + |E_i|_p \\
 &\leq (k+1)|\Delta_i|_p. \quad \square
 \end{aligned}$$

3.2. An m^k poly(nm)-time algorithm

In this section we give an algorithm that returns a $(k+1)$ -approximation to the k -LRA $_p$ problem in time m^k poly(nm).

Algorithm 2 A $(k+1)$ -approximation to k -LRA $_p$.

Require: An integer k and a matrix A

Ensure: $U \in \mathbb{R}^{n \times k}$, $V \in \mathbb{R}^{k \times m}$ s.t. $|A - UV|_p \leq (k+1)\text{opt}_{k,p}(A)$.

- 1: **for all** $I \in \binom{[m]}{k}$ **do**
- 2: Let $U = A_I$
- 3: Use Lemma 1 to compute a matrix V that minimizes the distance $d_I = |A - UV|_p$
- 4: **end for**
- 5: Return U, V that minimizes d_I , for $I \in \binom{[m]}{k}$

The following statement follows directly from the existence of k columns in A that make up a factor U having small ℓ_p error (Theorem 3).

Theorem 4. Algorithm 2 obtains a $(k+1)$ -approximation to k -LRA $_p$ in time m^k poly(nm).

4. A poly(nm)-time bi-criteria algorithm for k -CSS $_p$

We next show an algorithm that runs in time poly(nm) but returns $O(k \log m)$ columns of A that can be used in place of U , with an error $O(k)$ times the error of the best k -factorization. In other words, it obtains more than k columns but achieves a polynomial running time; we will later build upon this algorithm in Section 5 to obtain a faster algorithm for the k -LRA $_p$ problem. We also show a lower bound: there exists a matrix A for which the best possible approximation for the k -CSS $_p$, for $p \in (2, \infty)$, is $k^{\Omega(1)}$.

Definition 5 (Approximate coverage). Let S be a subset of k column indices. We say that column A_i is c_p -approximately covered by S if for $p \in [1, \infty)$ we have $\min_{x \in \mathbb{R}^{k \times 1}} |A_S x - A_i|_p \leq c \frac{100(k+1)^p |\Delta_i|_p}{n}$, and for $p = \infty$, $\min_{x \in \mathbb{R}^{k \times 1}} |A_S x - A_i|_\infty \leq c(k+1)|\Delta|_\infty$. If $c = 1$, we say A_i is covered by S .

We first show that if we select a set R columns of size $2k$

uniformly at random in $\binom{[m]}{2k}$, with constant probability we cover a constant fraction of columns of A .

Lemma 6. *Suppose R is a set of $2k$ uniformly random chosen columns of A . With probability at least $2/9$, R covers at least a $1/10$ fraction of columns of A .*

Proof. Let i be a column index of A selected uniformly at random and not in R . Let $T = R \cup \{i\}$ and let η be the cost of the best ℓ_p rank- k approximation to A_T . Note that T is a uniformly random subset of $2k + 1$ columns of A .

Case: $p < \infty$. Since T is a uniformly random subset of $2k + 1$ columns of A , $\mathbf{E}_T[\eta^p] = \frac{(2k+1)|\Delta|_p^p}{n}$. Let \mathcal{E}_1 denote the event “ $\eta^p \leq \frac{10(2k+1)|\Delta|_p^p}{n}$ ”. By a Markov bound, $\Pr[\mathcal{E}_1] \geq 9/10$.

By Theorem 3, there exists a subset L of k columns of A_T for which $\min_x |A_L x - A_T|_p \leq (k+1)\eta^p$. Since i is itself uniformly random in the set T , it holds that $\mathbf{E}_i[\min_x |A_L x - A_i|_p \leq \frac{(k+1)^p \eta^p}{2k+1}]$. Let \mathcal{E}_2 denote the event “ $\min_x |A_L x - A_i|_p \leq \frac{10(k+1)^p \eta^p}{2k+1}$ ”. By a Markov bound, $\Pr[\mathcal{E}_2] \geq 9/10$.

Let \mathcal{E}_3 denote the event “ $i \notin L$ ”. Since i is uniformly random in the set T , $\Pr[\mathcal{E}_3] \geq \frac{k+1}{2k} > 1/2$.

Clearly $\Pr[\mathcal{E}_1 \wedge \mathcal{E}_2 \wedge \mathcal{E}_3] \geq 3/10$. Conditioned on $\mathcal{E}_1 \wedge \mathcal{E}_2 \wedge \mathcal{E}_3$, we have

$$\begin{aligned} \min_x |A_R x - A_i|_p &\stackrel{\mathcal{E}_3}{\leq} \min_x |A_L x - A_i|_p \\ &\stackrel{\mathcal{E}_2}{\leq} \frac{10(k+1)^p \eta^p}{2k+1} \\ &\stackrel{\mathcal{E}_1}{\leq} \frac{100(k+1)^p |\Delta|_p^p}{n}, \end{aligned}$$

which implies that i is covered by R . Note that the first inequality uses that L is a subset of R given \mathcal{E}_3 , and so the regression cost using A_L cannot be smaller than that of using A_R .

Let Z_i be an indicator variable if i is covered by R and let $Z = \sum_i Z_i$. We have $\mathbf{E}[Z] = \sum_i \mathbf{E}[Z_i] \geq \sum_i \frac{3}{10} = 3m/10$; hence $\mathbf{E}[m - Z] \leq \frac{7m}{10}$. By a Markov bound, $\Pr[m - Z \geq \frac{9m}{10}] \leq \frac{7}{9}$.

Case $p = \infty$. Then $\eta \leq |\Delta|_\infty$ since A_T is a submatrix of A . By Theorem 3, there exists a subset L of k columns of A_T for which $\min_x |A_L x - A_T|_\infty \leq (k+1)\eta$. Defining \mathcal{E}_3 as before and conditioning on it, we have

$$\begin{aligned} \min_x |A_R x - A_i|_\infty &\leq \min_x |A_L x - A_i|_\infty \\ &\leq \min_x |A_L x - A_T|_\infty \\ &\leq (k+1)|\Delta|_\infty, \end{aligned}$$

i.e., i is covered by R . Again defining Z_i to be the event that

i is covered by R , we have $\mathbf{E}[Z_i] \geq \frac{1}{2}$, and so $\mathbf{E}[m - Z] \leq \frac{m}{2}$, which implies $\Pr[m - Z \geq \frac{9m}{10}] \leq \frac{5}{9} < \frac{7}{9}$. \square

We are now ready to introduce Algorithm 3. We can wlog. assume that the algorithm knows a number N for which $|\Delta|_p \leq N \leq 2|\Delta|_p$. Indeed, such a value can be obtained by first computing $|\Delta|_2$ using the SVD. Note that although one does not know Δ , one does know $|\Delta|_2$ since this is the Euclidean norm of all but the top k singular values of A , which one can compute from the SVD of A . Then, note that for $p < 2$, $|\Delta|_2 \leq |\Delta|_p \leq n^{2-p} |\Delta|_2$, while for $p \geq 2$, $|\Delta|_p \leq |\Delta|_2 \leq n^{1-2/p} |\Delta|_p$. Hence, given $|\Delta|_2$, there are only $O(\log n)$ values of N , one of which will satisfy $|\Delta|_p \leq N \leq 2|\Delta|_p$. One can take the best solution found by Algorithm 3 for each of the $O(\log n)$ guesses to N .

Algorithm 3 Selecting $O(k \log m)$ columns of A .

Require: An integer k , and a matrix $A = A^* + \Delta$.

Ensure: $O(k \log m)$ columns of A

SELECTCOLUMNS(k, A)

if number of columns of $A \leq 2k$ **then**

 return all the columns of A

else

repeat

 Let R be uniform at random $2k$ columns of A

until at least $(1/10)$ -fraction columns of A are c_p -approximately covered

 Let $A_{\bar{R}}$ be the columns of A not approximately covered by R

 return $A_R \cup \text{SELECTCOLUMNS}(k, A_{\bar{R}})$

end if

Theorem 7. *With probability at least $9/10$, Algorithm 3 runs in time $\text{poly}(nm)$ and returns $O(k \log m)$ columns that can be used as a factor of the whole matrix inducing ℓ_p error $O(k|\Delta|_p)$.*

Proof. First note that if $|\Delta|_p \leq N \leq 2|\Delta|_p$ and if i is covered by a set R of columns, then i is c_p -approximately covered by R for a constant c_p ; here $c_p = 2^p$ for $p < \infty$ and $c_\infty = 2$. By Lemma 6, the expected number of repetitions of selecting $2k$ columns until $(1/10)$ -fraction of columns of A are covered is $O(1)$. When we recurse on SELECTCOLUMNS on the resulting matrix $A_{\bar{R}}$, each such matrix admits a rank- k factorization of cost at most $|\Delta|_p$. Moreover, the number of recursive calls to SELECTCOLUMNS can be upper bounded by $\log_{10} m$. In expectation there will be $O(\log m)$ total repetitions of selecting $2k$ columns, and so by a Markov bound, with probability $9/10$, the algorithm will choose $O(k \log m)$ columns in total and run in time $\text{poly}(nm)$.

Let S be the union of all columns of A chosen by the algorithm. Then for each column i of A , for $p \in [1, \infty)$,

we have $\min_x |A_S x - A_i|_p^p \leq \frac{100(k+1)^p 2^p |\Delta|_p^p}{n}$, and so $\min_X |A_S X - A|_p^p \leq 100(k+1)^p 2^p |\Delta|_p^p$. For $p = \infty$ we instead have $\min_x |A_S x - A_i|_\infty \leq 2(k+1)|\Delta|_\infty$, and so $\min_X |A_S X - A|_\infty \leq 2(k+1)|\Delta|_\infty$. \square

4.1. A lower bound for k -CSS $_p$

In this section we prove an existential result showing that there exists a matrix for which the best approximation to the k -CSS $_p$ is $k^{\Omega(1)}$.

Lemma 8. *There exists a matrix A such that the best approximation for the k -CSS $_p$ problem, for $p \in (2, \infty)$, is $k^{\Omega(1)}$.*

Proof. Consider $A = (k+1)I_{k+1}$, where I_{k+1} is the $(k+1) \times (k+1)$ identity matrix. And consider the matrix $B = (k+1) \cdot I_{k+1} - E$, where E is the $(k+1) \times (k+1)$ all ones matrix. Note that B has rank at most k , since the sum of its columns is 0.

Case: $2 < p < \infty$. If we choose any k columns of A , then the ℓ_p cost of using them to approximate A is $(k+1)$. On the other hand, $|A - B|_\infty = 1$, which means that ℓ_p cost of B is smaller or equal than $((k+1)^2)^{1/p}$.

Case: $p = \infty$. If we choose any k columns of A , then the ℓ_∞ cost of using them to approximate A is $k+1$. On the other hand, $|A - B|_\infty = 1$, which means that ℓ_∞ cost of B is smaller or equal than 1. \square

Note also that in (Song et al., 2017) the authors show that for $p = 1$ the best possible approximation is $\Omega(\sqrt{k})$ up to $\text{poly}(\log k)$ factors.

5. A $((k \log n)^k \text{poly}(mn))$ -time algorithm for k -LRA $_p$

In the previous section we have shown how to get a rank- $O(k \log m)$, $O(k)$ -approximation in time $\text{poly}(nm)$ to the k -CSS $_p$ and k -LRA $_p$ problems. In this section we first show how to get a rank- k , $\text{poly}(k)$ -approximation efficiently starting from a rank- $O(k \log m)$ approximation. This algorithm runs in polynomial time as long as $k = O\left(\frac{\log n}{\log \log n}\right)$. We then show how to obtain a $(k \log m)^{O(p)}$ -approximation ratio in polynomial time for every k .

Let U be the columns of A selected by Algorithm 3.

5.1. An isoperimetric transformation

The first step of our proof is to show that we can modify the selected columns of A to span the same space but to have small distortion. For this, we need the following notion of isoperimetry.

Definition 9 (Almost isoperimetry). *A matrix $B \in \mathbb{R}^{n \times m}$ is almost- ℓ_p -isoperimetric if for all x , we have*

$$\frac{|x|_p}{2m} \leq |Bx|_p \leq |x|_p.$$

We now show that given a full rank $A \in \mathbb{R}^{n \times m}$, it is possible to construct in polynomial time a matrix $B \in \mathbb{R}^{n \times m}$ such that A and B span the same space and B is almost- ℓ_p -isoperimetric.

Lemma 10. *Given a full (column) rank $A \in \mathbb{R}^{n \times m}$, there is an algorithm that transforms A into a matrix B such that $\text{span } A = \text{span } B$ and B is almost- ℓ_p -isoperimetric. Furthermore the running time of the algorithm is $\text{poly}(nm)$.*

Proof. In (Dasgupta et al., 2009), specifically, Equation (4) in the proof of Theorem 4, the authors show that in polynomial time it is possible to find a matrix B such that $\text{span } B = \text{span } A$ and for all x ,

$$|x|_2 \leq |Bx|_p \leq \sqrt{m}|x|_2,$$

for any $p \geq 1$.

If $p < 2$, their result implies

$$\frac{|x|_p}{\sqrt{m}} \leq |x|_2 \leq |Bx|_p \leq \sqrt{m}|x|_2 \leq \sqrt{m}|x|_p,$$

and so rescaling B by \sqrt{m} makes it almost- ℓ_p -isoperimetric. On the other hand, if $p > 2$, then

$$|x|_p \leq |x|_2 \leq |Bx|_p \leq \sqrt{m}|x|_2 \leq m|x|_p,$$

and rescaling B by m makes it almost- ℓ_p -isoperimetric. \square

Note that the algorithm used in (Dasgupta et al., 2009) relies on the construction of the Löwner–John ellipsoid for a specific set of points. Interestingly, we can also show that there is a more simple and direct algorithm to compute such a matrix B ; this may be of independent interest. We provide the details of our algorithm in the full version (Chierichetti et al., 2017).

5.2. Reducing the rank to k

The main idea for reducing the rank is to first apply the almost- ℓ_p -isoperimetric transformation to the factor U to obtain a new factor Z^0 . For such a Z^0 , the ℓ_p -norm of $Z^0 V$ is at most the ℓ_p -norm of V . Using this fact we show that V has a low-rank approximation and a rank- k approximation of V translates into a good rank- k approximation of UV . But a good rank- k approximation of V can be obtained by exploring all possible k -subsets of rows of V , as in Algorithm 2. More formally, in Algorithm 4 we give the pseudo-code to reduce the rank of our low-rank approximation from $O(k \log m)$ to k . Let $\delta = |\Delta|_p = \text{opt}_{k,p}(A)$.

Algorithm 4 An algorithm that transforms an $O(k \log m)$ -rank matrix decomposition into a k -rank matrix decomposition without inflating the error too much.

Require: $U \in \mathbb{R}^{n \times O(k \log m)}$, $V \in \mathbb{R}^{O(k \log m) \times m}$

Ensure: $W \in \mathbb{R}^{n \times k}$, $Z \in \mathbb{R}^{k \times m}$

- 1: Apply Lemma 10 to U to obtain matrix W^0
- 2: Apply Lemma 1 to obtain matrix Z^0 , s.t. $\forall i$, $|W^0 Z_i^0 - (UV)_i|_p$ is minimized
- 3: Apply Algorithm 2 with input $(Z^0)^T \in \mathbb{R}^{n \times O(k \log m)}$ and k to obtain X and Y
- 4: Set $Z \leftarrow X^T$
- 5: Set $W \leftarrow W^0 Y^T$
- 6: Output W and Z

Theorem 11. Let $A \in \mathbb{R}^{n \times m}$, $U \in \mathbb{R}^{n \times O(k \log m)}$, $V \in \mathbb{R}^{O(k \log m) \times m}$ be such that $|A - UV|_p = O(k\delta)$. Then, Algorithm 4 runs in time $O(k \log m)^k (mn)^{O(1)}$ and outputs $W \in \mathbb{R}^{n \times k}$, $Z \in \mathbb{R}^{k \times m}$ such that $|A - WZ|_p = O((k^4 \log k)\delta)$.

5.3. Improving the running time

Interestingly it is possible to improve the running time to $(mn)^{O(1)}$ for every k and every constant $p \geq 1$, at the cost of a poly($k \log m$)-approximation instead of the poly(k)-approximation we had previously. See the full version (Chierichetti et al., 2017) for a proof.

Theorem 12. Let $A \in \mathbb{R}^{n \times m}$, $1 \leq k \leq \min(m, n)$, and $p \geq 1$ be an arbitrary constant. Let $U \in \mathbb{R}^{n \times O(k \log m)}$ and $V \in \mathbb{R}^{O(k \log m) \times m}$ be such that $|A - UV|_p = O(k\delta)$. There is an algorithm that runs in time $(mn)^{O(1)}$ and outputs $W \in \mathbb{R}^{n \times k}$, $Z \in \mathbb{R}^{k \times m}$ such that $|A - WZ|_p = (k \log m)^{O(p)} \delta$.

6. Experiments

In this section, we show the effectiveness of Algorithm 2 compared to the SVD. We run our comparison both on synthetic as well as real data sets. For the real data sets, we use matrices from the FIDAP set² and a word frequency dataset from UC Irvine³. The FIDAP matrix is 27×27 with 279 real asymmetric non-zero entries. The KOS blog entries matrix, representing word frequencies in blogs, is 3430×6906 with 353160 non-zero entries. For the synthetic data sets, we use two matrices. For the first, we use a 20×30 random matrix with 184 non-zero entries—this random matrix was generated as follows: independently, we set each entry to 0 with probability 0.7, and to a uniformly random

value in $[0, 1]$ with probability 0.3. Both matrices are full rank. For the second matrix, we use a random ± 1 20×30 matrix.

In all our experiments, we run a simplified version of Algorithm 2, where instead of running for all possible $\binom{m}{k}$ subsets of k columns (which would be computationally prohibitive), we repeatedly sample k columns, a few thousand times, uniformly at random. We then run the ℓ_p -projection on each sampled set and finally select the solution with the smallest ℓ_p -error. (While this may not guarantee provable approximations, we use this a reasonable heuristic that seems to work well in practice, without much computational overhead.) We focus on $p = 1$ and $p = \infty$.

Figure 1 illustrates the relative performance of Algorithm 2 compared to the SVD for different values of k on the real data sets. In the figure the green line is the ratio of the total error. The ℓ_1 -error for Algorithm 2 is always less than the corresponding error for the SVD and in fact consistently outperforms the SVD by roughly 40% for small values of k on the FIDAP matrix. On the larger KOS matrix, the relative improvement in performance with respect to ℓ_∞ -error is more uniform (around 10%).

We observe similar trends for the synthetic data sets as well. Figures 2 and 3 illustrate the trends. Algorithm 2 performs consistently better than the SVD in the case of ℓ_1 -error for both the matrices. In the case of ℓ_∞ -error, it outperforms SVD by around 10% for higher values of k on the random matrix. Furthermore, it consistently outperforms SVD, between 30% and 50%, for all values of k on the random ± 1 matrix.

To see why our ℓ_∞ error is always 1 for a random ± 1 matrix A , note that by setting our rank- k approximation to be the zero matrix, we achieve an ℓ_∞ error of 1. This is optimal for large values of n and m and small k as can be seen by recalling the notion of the *sign-rank* of a matrix $A \in \{-1, 1\}^{n \times m}$, which is the minimum rank of a matrix B for which the sign of $B_{i,j}$ equals $A_{i,j}$ for all entries i, j . If the sign-rank of A is larger than k , then for any rank- k matrix B , we have $\|A - B\|_\infty \geq 1$ since necessarily there is an entry $A_{i,j}$ for which $|A_{i,j} - B_{i,j}| \geq 1$. It is known that the sign-rank of a random $m \times m$ matrix A , and thus also of a random $n \times m$ matrix A , is $\Omega(\sqrt{m})$ with high probability (Forster, 2002).

7. Conclusions

We studied the problem of low-rank approximation in the entry-wise ℓ_p error norm and obtained the first provably good approximation algorithms for the problem that work for every $p \geq 1$. Our algorithms are extremely simple, which makes them practically appealing. We showed the effectiveness of our algorithms compared with the SVD on

²http://math.nist.gov/MatrixMarket/data/SPARSKIT/_fidap/_fidap005.html

³<https://archive.ics.uci.edu/ml/datasets/Bag+of+Words>

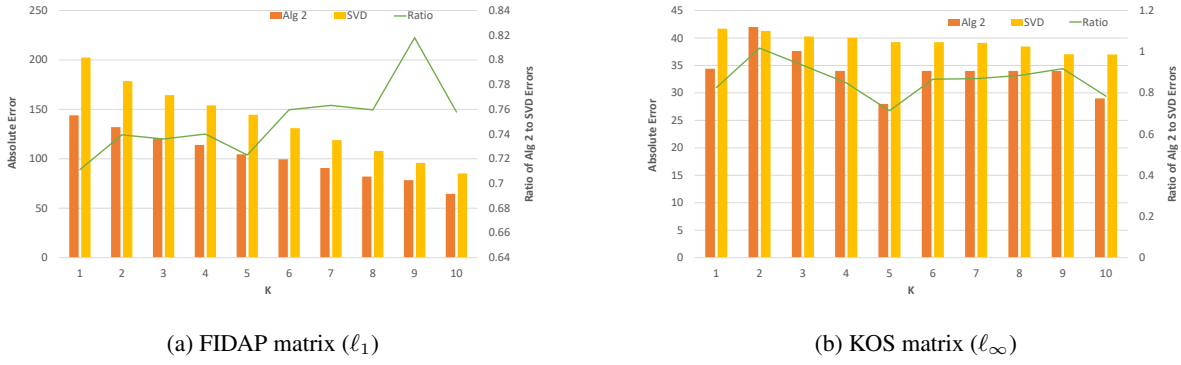


Figure 1: Comparing the performance of Algorithm 2 with SVD on the real data sets.

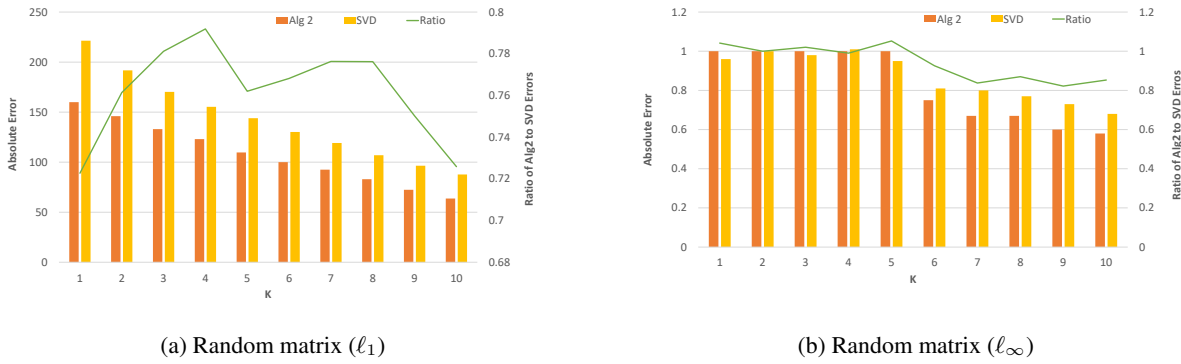


Figure 2: Comparing the performance of Algorithm 2 with SVD on the random matrix.

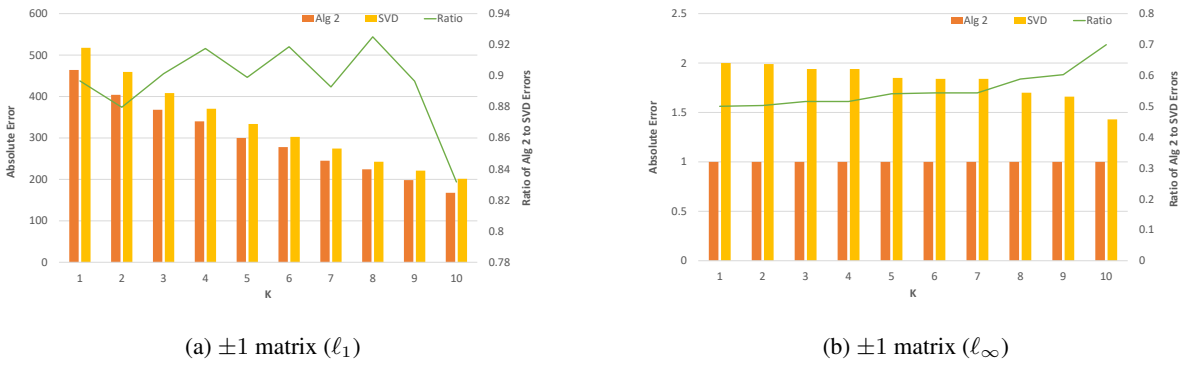


Figure 3: Comparing the performance of Algorithm 2 with SVD on the ± 1 matrix.

real and synthetic data sets. We obtain a $k^{O(1)}$ approximation factor for every p for the column subset selection problem, and we showed an example matrix for this problem for which a $k^{\Omega(1)}$ approximation factor is necessary. It is unclear if better approximation factors are possible by designing algorithms that do not choose a subset of input columns to span the output low rank approximation. Resolving this would be an interesting and important research direction.

References

- Brooks, J. Paul, Dulá, Jose' H., and Boone, Edward L. A pure ℓ_1 -norm principal component analysis. *Computational Statistics & Data Analysis*, 61:83–98, 2013.
- Candès, Emmanuel J., Li, Xiaodong, Ma, Yi, and Wright, John. Robust principal component analysis? *JACM*, 58(3):11:1–11:37, 2011.
- Chierichetti, Flavio, Gollapudi, Sreenivas, Kumar, Ravi, Lattanzi, Silvio, Panigrahy, Rina, and Woodruff, David P. Algorithms for ℓ_p low-rank approximation. Technical Report 1705.06730, arXiv, 2017.
- Dan, Chen, Hansen, Kristoffer A., Jiang, He, Wang, Liwei, and Zhou, Yuchen. On low rank approximation of binary matrices. Technical Report 1511.01699v1, arXiv, 2015.
- Dasgupta, Anirban, Drineas, Petros, Harb, Boulos, Kumar, Ravi, and Mahoney, Michael W. Sampling algorithms and coresets for ℓ_p regression. *SICOMP*, 38(5)(2060–2078), 2009.
- Deshpande, Amit, Tulsiani, Madhur, and Vishnoi, Nisheeth K. Algorithms and hardness for subspace approximation. In *SODA*, pp. 482–496, 2011.
- Eriksson, Anders and van den Hengel, Anton. Efficient computation of robust low-rank matrix approximations using the L_1 norm. *PAMI*, 34(9):1681–1690, 2012.
- Feldman, Dan, Fiat, Amos, Sharir, Micha, and Segev, Danny. Bi-criteria linear-time approximations for generalized k -mean/median/center. In *SoCG*, pp. 19–26, 2007.
- Forster, Jürgen. A linear lower bound on the unbounded error probabilistic communication complexity. *J. Comput. Syst. Sci.*, 65(4):612–625, 2002.
- Gillis, Nicolas and Vavasis, Stephen A. On the complexity of robust PCA and ℓ_1 -norm low-rank matrix approximation. Technical Report 1509.09236, arXiv, 2015.
- Goreinov, Sergei A. and Tyrtshnikov, Eugene E. The maximal-volume concept in approximation by low-rank matrices. *Contemporary Mathematics*, 208:47–51, 2001.
- Goreinov, Sergei A. and Tyrtshnikov, Eugene E. Quasi-optimality of skeleton approximation of a matrix in the Chebyshev norm. *Doklady Mathematics*, 83(3):374–375, 2011.
- Huber, Peter J. *Robust Statistics*. John Wiley & Sons, New York., 1981.
- Ke, Qifa and Kanade, Takeo. Robust L_1 norm factorization in the presence of outliers and missing data by alternative convex programming. In *CVPR*, pp. 739–746, 2005.
- Lu, Cewu, Shi, Jiaping, and Jia, Jiaya. Scalable adaptive robust dictionary learning. *TIP*, 23(2):837–847, 2014.
- Meng, Deyu and Torre, Fernando. D. L. Robust matrix factorization with unknown noise. In *ICCV*, pp. 1337–1344, 2013.
- Netrapalli, Praneeth, Niranjan, U. N., Sanghavi, Sujay, Anandkumar, Animashree, and Jain, Prateek. Non-convex robust PCA. In *NIPS*, pp. 1107–1115, 2014.
- Song, Zhao, Woodruff, David P., and Zhong, Pelin. Low rank approximation with entrywise ℓ_1 -norm error. In *STOC*, 2017.
- Wang, Naiyan and Yeung, Dit-Yan. Bayesian robust matrix factorization for image and video processing. In *ICCV*, pp. 1785–1792, 2013.
- Wang, Naiyan, Yao, Tiansheng, Wang, Jingdong, and Yeung, Dit-Yan. A probabilistic approach to robust matrix factorization. In *ECCV*, pp. 126–139, 2012.
- Woodruff, David P. Sketching as a tool for numerical linear algebra. *Foundations and Trends in Theoretical Computer Science*, 10(1-2):1–157, 2014.
- Xiong, Liang, Chen, Xi, and Schneider, Jeff. Direct robust matrix factorization for anomaly detection. In *ICDM*, pp. 844–853, 2011.
- Xu, Huan, Caramanis, Constantine, and Sanghavi, Sujay. Robust PCA via outlier pursuit. *TOIT*, 58(5):3047–3064, 2012.
- Xu, Lei and Yuille, Alan L. Robust principal component analysis by self-organizing rules based on statistical physics approach. *IEEE Transactions on Neural Networks*, 6(1):131–143, 1995.
- Yi, Xinyang, Park, Dohyung, Chen, Yudong, and Caramanis, Constantine. Fast algorithms for robust pca via gradient descent. In *NIPS*, pp. 4152–4160, 2016.
- Zheng, Yinqiang, Liu, Guangcan, Sugimoto, Shigeki, Yan, Shuicheng, and Okutomi, Masatoshi. Practical low-rank matrix approximation under robust L_1 -norm. In *CVPR*, pp. 1410–1417, 2012.