



# A note on element-wise matrix sparsification via a matrix-valued Bernstein inequality

Petros Drineas<sup>a,\*</sup>, Anastasios Zouzias<sup>b</sup>

<sup>a</sup> Department of Computer Science, Rensselaer Polytechnic Institute, United States

<sup>b</sup> Department of Computer Science, University of Toronto, Canada

## ARTICLE INFO

### Article history:

Received 2 June 2010

Received in revised form 3 January 2011

Accepted 11 January 2011

Available online 18 January 2011

Communicated by B. Doerr

### Keywords:

Matrix-valued Bernstein bounds

Matrix sparsification

Algorithms

## ABSTRACT

Given a matrix  $A \in \mathbb{R}^{n \times n}$ , we present a simple, element-wise sparsification algorithm that zeroes out all sufficiently small elements of  $A$  and then retains some of the remaining elements with probabilities proportional to the square of their magnitudes. We analyze the approximation accuracy of the proposed algorithm using a recent, elegant non-commutative Bernstein inequality, and compare our bounds with all existing (to the best of our knowledge) element-wise matrix sparsification algorithms.

© 2011 Elsevier B.V. All rights reserved.

## 1. Introduction

Element-wise matrix sparsification was pioneered by Achlioptas and McSherry [1,2], who described sampling-based algorithms to select a small number of elements from an input matrix  $A \in \mathbb{R}^{n \times n}$  in order to construct a sparse sketch  $\tilde{A} \in \mathbb{R}^{n \times n}$ , which is close to  $A$  in the operator norm. Such sketches were used in approximate eigenvector computations [1,4,2], semi-definite programming solvers [3,7], and matrix completion problems [5,6]. Motivated by their work, we present a simple matrix sparsification algorithm that achieves the best known upper bounds for element-wise matrix sparsification.

Our main algorithm (Algorithm 1) zeroes out “small” elements of  $A$  and randomly samples the remaining elements of  $A$  with respect to a probability distribution that favors “larger” entries.

In Algorithm 1, we let  $e_1, e_2, \dots, e_n \in \mathbb{R}^n$  denote the standard basis vectors for  $\mathbb{R}^n$  (see Section 3.1 for more notation). Our sampling procedure selects  $s$  entries from  $A$  (note that  $\hat{A}$  from the description of Algorithm 1 is sim-

### Algorithm 1 Matrix sparsification algorithm.

1: **Input:**  $A \in \mathbb{R}^{n \times n}$ , accuracy parameter  $\epsilon > 0$ .

2: **Let**  $\hat{A} = A$  and **zero-out** all entries of  $\hat{A}$  that are smaller (in absolute value) than  $\epsilon/2n$ .

3: **Set**  $s$  as in Eq. (1).

4: **For**  $t = 1 \dots s$  (i.i.d. trials with replacement) **randomly sample** indices  $(i_t, j_t)$  (entries of  $\hat{A}$ ), with

$$\mathbb{P}((i_t, j_t) = (i, j)) = p_{ij},$$

$$\text{where } p_{ij} := \hat{A}_{ij}^2 / \|\hat{A}\|_F^2 \text{ for all } (i, j) \in [n] \times [n].$$

5: **Output:**

$$\tilde{A} = \frac{1}{s} \sum_{t=1}^s \frac{\hat{A}_{i_t j_t}}{p_{i_t j_t}} e_{i_t} e_{j_t}^T \in \mathbb{R}^{n \times n}.$$

ply  $A$ , but with elements less than or equal to  $\epsilon/(2n)$  zeroed out) in  $s$  independent, identically distributed (i.i.d.) trials with replacement. In each trial, elements of  $A$  are retained with probability proportional to their squared magnitude. Note that the same element of  $A$  could be selected multiple times and that  $\tilde{A}$  contains at most  $s$  non-zero entries. Theorem 1 is our main quality-of-approximation result for Algorithm 1 and achieves sparsity bounds proportional to  $\|A\|_F^2$ .

\* Corresponding author.

E-mail addresses: drinep@cs.rpi.edu (P. Drineas),  
zouzias@cs.toronto.edu (A. Zouzias).

**Table 1**

Summary of prior work in matrix sparsification results. Given a matrix  $A \in \mathbb{R}^{n \times n}$  and an accuracy parameter  $\epsilon > 0$ , we seek a sparse  $\tilde{A} \in \mathbb{R}^{n \times n}$  such that  $\|A - \tilde{A}\|_2 \leq \epsilon$ . The first column indicates the number of non-zero entries in  $\tilde{A}$ , whereas the second column indicates whether this number is exact or simply holds in expectation. In terms of notation, we let  $b$  denote the  $\max_{i,j} |A_{ij}|$  and  $R$  denote the  $\max_{ij} |A_{ij}| / \min_{A_{ij} \neq 0} |A_{ij}|$ . Finally,  $c_1, c_2, c_3, c_4$  denote unspecified constants.

Comparison with prior results				
Sparsity of $\tilde{A}$		Failure probability	Citation	Comments
$16n \ A\ _F^2 / \epsilon^2 + 8^4 n \log^4 n$	Expected	$e^{-19 \log^4 n}$	[2]	$\epsilon > 4\sqrt{n} \cdot b$ $n \geq 700 \cdot 10^6$
$R \cdot b \cdot n \ A\ _F^2 / \epsilon^2$	Expected	$e^{-\Omega(R \cdot n)}$	[9]	$\epsilon > c_1 \sqrt{n} \cdot R \cdot b, n \geq 1$
$c_2 n \log^2(\frac{n}{\log^2 n}) \log n \ A\ _F^2 / \epsilon^2$	Expected	$1/n$	[12]	$\epsilon > 0, n \geq 300,$ $c_2 \leq 45^2$
$c_3 n \log^3 n \ A\ _F^2 / \epsilon^2$	Expected	$1/n$	[13]	$\epsilon > 0, n \geq 300$ Extends to tensors
$c_4 \sqrt{n} \sum_{ij}  A_{ij}  / \epsilon$	Exact	$e^{-\Omega(m)}$	[4]	$\epsilon > 0, n \geq 1$
$28n \ln(\sqrt{2n}) \ A\ _F^2 / \epsilon^2$	Exact	$1/n$	Theorem 1	$\epsilon > 0, n \geq 1$

**Theorem 1.** Let  $A \in \mathbb{R}^{n \times n}$  be any matrix, let  $\epsilon > 0$  be an accuracy parameter, and let  $\tilde{A}$  be the sparse sketch of  $A$  constructed via Algorithm 1. If

$$s = \frac{28n \ln(\sqrt{2n})}{\epsilon^2} \|A\|_F^2, \tag{1}$$

then, with probability at least  $1 - n^{-1}$ ,

$$\|A - \tilde{A}\|_2 \leq \epsilon.$$

$\tilde{A}$  has at most  $s$  non-zero entries and the construction of  $\tilde{A}$  can be implemented in one pass over the input matrix  $A$  (see Section 3.2).

We conclude this section with Corollary 1, which is a re-statement of Theorem 1 involving the *stable rank* of  $A$ , denoted by  $\mathbf{sr}(A)$  (recall that the stable rank of any matrix  $A$  is defined as the ratio  $\mathbf{sr}(A) := \|A\|_F^2 / \|A\|_2^2$ , which is upper bounded by the rank of  $A$ ). The corollary guarantees relative error approximations for matrices of – say – constant stable rank, such as the ones that arise in [14,6].

**Corollary 1.** Let  $A \in \mathbb{R}^{n \times n}$  be any matrix and let  $\epsilon > 0$  be an accuracy parameter. Let  $\tilde{A}$  be the sparse sketch of  $A$  constructed via Algorithm 1 (with  $\epsilon = \epsilon \|A\|_2$ ). If  $s = 28n \mathbf{sr}(A) \ln(\sqrt{2n}) / \epsilon^2$ , then, with probability at least  $1 - n^{-1}$ ,

$$\|A - \tilde{A}\|_2 \leq \epsilon \|A\|_2.$$

It is worth noting that the sampling algorithm implied by Corollary 1 cannot be implemented in one pass, since we would need a priori knowledge of the spectral norm of  $A$  in order to implement Step 2 of Algorithm 1.

**2. Related work**

In this section (as well as in Table 1), we present a head-to-head comparison of our result with all existing (to the best of our knowledge) bounds on matrix sparsification. In [1,2] the authors presented a sampling method that requires in *expectation*  $16n \|A\|_F^2 / \epsilon^2 + 8^4 n \log^4 n$  non-zero entries in  $\tilde{A}$  in order to achieve an accuracy guarantee  $\epsilon$  with a failure probability of at most  $e^{-19 \log^4 n}$ .

Compared with our result, their bound holds only when  $\epsilon > 4\sqrt{n} \cdot \max_{i,j} |A_{ij}|$  and, in this range, our bounds are superior when  $\|A\|_F^2 / (\max_{i,j} |A_{ij}|)^2 = o(n \log^3 n)$ . It is worth mentioning that the constant involved in [1,2] is two orders of magnitude larger than ours and, more importantly, that the results of [1,2] hold only when  $n \geq 700 \cdot 10^6$ .

In [9], the authors study the  $\|\cdot\|_{\infty \rightarrow 2}$  and  $\|\cdot\|_{\infty \rightarrow 1}$  norms in the matrix sparsification context and they also present a sampling scheme analogous to ours. They achieve (in expectation) a sparsity bound of  $Rn \|A\|_F^2 \max_{i,j} |A_{ij}| / \epsilon^2$  when  $\epsilon \geq \sqrt{nR} \max_{i,j} |A_{ij}|$ ; here  $R = \max_{ij} |A_{ij}| / \min_{A_{ij} \neq 0} |A_{ij}|$ . Thus, our results are superior (in the above range of  $\epsilon$ ) when  $R \cdot \max_{i,j} |A_{ij}| = \omega(\log n)$ .

It is harder to compare our method to the work of [4], which depends on the  $\sum_{i,j=1}^n |A_{ij}|$ . The latter quantity is, in general, upper bounded only by  $n \|A\|_F$ , in which case the sampling complexity of [4] is much worse, namely  $O(n^{3/2} \|A\|_F^2 / \epsilon)$ . Finally, the recent bounds on matrix sparsification via the non-commutative Khintchine’s inequality in [12] are inferior compared to ours in terms of sparsity guarantees by at least  $O(\ln^2(n / \ln^2 n))$ . However, we should mention that the bounds of [12] can be extended to multi-dimensional matrices (tensors), whereas our result does not generalize to this setting; see [13] for details.

**3. Background**

**3.1. Notation**

We let  $[n]$  denote the set  $\{1, 2, \dots, n\}$ . We will use the notation  $\mathbb{P}(\cdot)$  to denote the probability of the event in the parentheses and  $\mathbb{E}(X)$  to denote the expectation of a random variable  $X$ . When  $X$  is a matrix,  $\mathbb{E}(X)$  denotes the element-wise expectation of each entry of  $X$ . For a matrix  $X \in \mathbb{R}^{n \times n}$ ,  $X^{(j)}$  will denote the  $j$ -th column of  $X$  as a column vector and, similarly,  $X_{(i)}$  will denote the  $i$ -th row of  $X$  as a row vector (for any  $i$  or  $j$  in  $[n]$ ). The Frobenius norm  $\|X\|_F$  of the matrix  $X$  is defined as  $\|X\|_F^2 = \sum_{i,j=1}^n X_{ij}^2$ , and the spectral norm  $\|X\|_2$  of the matrix  $X$  is defined as  $\|X\|_2 = \max_{\|y\|_2=1} \|Xy\|_2$ . For two symmetric matrices  $X, Y$  we say that  $Y \succcurlyeq X$  if and only if  $Y - X$  is a positive semi-definite matrix. Finally,  $I_n$  denotes the identity matrix of size  $n$  and  $\ln x$  denotes the natural logarithm of  $x$ .

---

**Algorithm 2** One-pass SELECT algorithm.

---

- 1: **Input:**  $A_{ij}$  for all  $(i, j) \in [n] \times [n]$ , arbitrarily ordered and  $\epsilon > 0$ .
  - 2:  $N = 0$ .
  - 3: **For all**  $(i, j) \in [n] \times [n]$  **such that**  $A_{ij}^2 > \frac{\epsilon^2}{4n^2}$ 
    - $N = N + A_{ij}^2$ .
    - **Set**  $(I, J) = (i, j)$  and  $S = A_{ij}$  **with probability**  $\frac{A_{ij}^2}{N}$ .
  - 4: **Output:** Return  $(I, J)$ ,  $S$  and  $N$ .
- 

3.2. Implementing the sampling in one pass over the input matrix

We now discuss the implementation of Algorithm 1 in one pass over the input matrix  $A$ . Towards that end, we will leverage (a slightly modified version of) Algorithm SELECT (p. 137 of [8]).

We note that Step 3 essentially operates on  $\hat{A}$ . Clearly, in a single pass over the data we can run in parallel  $s$  copies of the SELECT Algorithm (using a total of  $O(s)$  memory) to effectively return  $s$  independent samples from  $\hat{A}$ . Lemma 1 (p. 136 of [8], note that the sequence of the  $A_{ij}^2$ 's is all-positive) guarantees that each of the  $s$  copies of SELECT returns a sample satisfying:

$$\begin{aligned} \mathbb{P}((i_t, j_t) = (i, j)) &= \frac{\hat{A}_{ij}^2}{\sum_{i,j=1}^n \hat{A}_{ij}^2} = \frac{\hat{A}_{ij}^2}{\|\hat{A}\|_F^2}, \quad \text{for all } t = 1, \dots, s. \end{aligned}$$

Finally, in the parlance of Step 5 of Algorithm 1,  $(i_t, j_t)$  is set to  $(I, J)$  and  $p_{i_t j_t}$  is set to  $S^2/N$  for all  $t \in [s]$ .

**4. Proof of Theorem 1**

The proof of Theorem 1 will combine Lemmas 1 and 4 in order to bound  $\|A - \tilde{A}\|_2$  as follows:

$$\begin{aligned} \|A - \tilde{A}\|_2 &= \|A - \hat{A} + \hat{A} - \tilde{A}\|_2 \\ &\leq \|A - \hat{A}\|_2 + \|\hat{A} - \tilde{A}\|_2 \leq \epsilon/2 + \epsilon/2 = \epsilon. \end{aligned}$$

The failure probability of Theorem 1 emerges from Lemma 4, which fails with probability at most  $n^{-1}$  for the choice of  $s$  in Eq. (1). The proof of Lemma 4 will involve an elegant matrix-valued Bernstein bound proven in [14]. See also [10] or [15, Theorem 2.10] for similar bounds.

4.1. Bounding  $\|A - \hat{A}\|_2$

**Lemma 1.** Using the notation of Algorithm 1,  $\|A - \hat{A}\|_2 \leq \epsilon/2$ .

**Proof.** Recall that the entries of  $\hat{A}$  are either equal to the corresponding entries of  $A$  or they are set to zero if the corresponding entry of  $A$  is (in absolute value) smaller than  $\epsilon/(2n)$ . Thus,

$$\begin{aligned} \|A - \hat{A}\|_2^2 &\leq \|A - \hat{A}\|_F^2 = \sum_{i,j=1}^n (A - \hat{A})_{ij}^2 \leq \sum_{i,j=1}^n \frac{\epsilon^2}{4n^2} \\ &\leq \frac{\epsilon^2}{4}. \quad \square \end{aligned}$$

4.2. Bounding  $\|\hat{A} - \tilde{A}\|_2$

In order to prove our main result in this section (Lemma 4) we will leverage a powerful matrix-valued Bernstein bound originally proven in [14] (Theorem 3.2). We restate this theorem, slightly rephrased to better suit our notation.

**Theorem 2.** (See Theorem 3.2 of [14].) Let  $M_1, M_2, \dots, M_s$  be independent, zero-mean random matrices in  $\mathbb{R}^{n \times n}$ . Suppose  $\max_{t \in [s]} \{\|\mathbb{E}(M_t M_t^T)\|_2, \|\mathbb{E}(M_t^T M_t)\|_2\} \leq \rho^2$  and  $\|M_t\|_2 \leq \gamma$  for all  $t \in [s]$ . Then, for any  $\tau > 0$ ,

$$\left\| \frac{1}{s} \sum_{t=1}^s M_t \right\|_2 \leq \tau$$

holds, subject to a failure probability of at most

$$2n \exp\left(-\frac{s\tau^2/2}{\rho^2 + \gamma\tau/3}\right).$$

In order to apply the above theorem, using the notation of Algorithm 1, we set  $M_t = \frac{\hat{A}_{i_t j_t}}{p_{i_t j_t}} e_{i_t} e_{j_t}^T - \hat{A}$  for all  $t \in [s]$  to obtain

$$\frac{1}{s} \sum_{t=1}^s M_t = \frac{1}{s} \sum_{t=1}^s \left[ \frac{\hat{A}_{i_t j_t}}{p_{i_t j_t}} e_{i_t} e_{j_t}^T - \hat{A} \right] = \tilde{A} - \hat{A}. \quad (2)$$

Let  $0_{n \times n}$  denote the  $n \times n$  matrix of all-zeros. It is easy to argue that  $\mathbb{E}(M_t) = 0_{n \times n}$  for all  $t \in [s]$ . Indeed, if we consider that  $\sum_{i,j=1}^n p_{ij} = 1$  and  $\hat{A} = \sum_{i,j=1}^n \hat{A}_{ij} e_i e_j^T$  we obtain

$$\begin{aligned} \mathbb{E}(M_t) &= \sum_{i,j=1}^n p_{ij} \left( \frac{\hat{A}_{ij}}{p_{ij}} e_i e_j^T - \hat{A} \right) \\ &= \sum_{i,j=1}^n \hat{A}_{ij} e_i e_j^T - \sum_{i,j=1}^n p_{ij} \hat{A} = 0_{n \times n}. \end{aligned}$$

Our next lemma bounds  $\|M_t\|_2$  for all  $t \in [s]$ .

**Lemma 2.** Using our notation,  $\|M_t\|_2 \leq 4n\epsilon^{-1} \|\hat{A}\|_F^2$  for all  $t \in [s]$ .

**Proof.** First, using the definition of  $M_t$  and the fact that  $p_{i_t j_t} = \hat{A}_{i_t j_t}^2 / \|\hat{A}\|_F^2$ ,

$$\begin{aligned} \|M_t\|_2 &= \left\| \frac{\hat{A}_{i_t j_t}}{p_{i_t j_t}} e_{i_t} e_{j_t}^T - \hat{A} \right\|_2 \leq \frac{\|\hat{A}\|_F^2}{|\hat{A}_{i_t j_t}|} + \|\hat{A}\|_2 \\ &\leq \frac{2n\|\hat{A}\|_F^2}{\epsilon} + \|\hat{A}\|_F. \end{aligned}$$

The last inequality follows since all entries of  $\hat{A}$  are at least  $\epsilon/(2n)$  and the fact that  $\|\hat{A}\|_2 \leq \|\hat{A}\|_F$ . We can now assume that

$$\|\hat{A}\|_F \leq \frac{2n\|\hat{A}\|_F^2}{\epsilon} \quad (3)$$

to conclude the proof of the lemma. To justify our assumption in Eq. (3), we note that if it is violated, then it must

be the case that  $\|\widehat{A}\|_F < \epsilon/(2n)$ . If that were true, then all entries of  $\widehat{A}$  would be equal to zero. (Recall that all entries of  $\widehat{A}$  are either zero or, in absolute value, larger than  $\epsilon/(2n)$ .) Also, if  $\widehat{A}$  were identically zero, then (i)  $\widetilde{A}$  would also be identically zero and, (ii) all entries of  $A$  would be at most  $\epsilon/(2n)$ . Thus,

$$\|A - \widetilde{A}\|_2 = \|A\|_2 \leq \|A\|_F \leq \sqrt{n^2 \frac{\epsilon^2}{4n^2}} = \frac{\epsilon}{2}.$$

Thus, if the assumption of Eq. (3) is not satisfied, the resulting all-zeros  $\widetilde{A}$  still satisfies Theorem 1.  $\square$

Our next step towards applying Theorem 2 involves bounding the spectral norm of the expectation of  $M_t M_t^T$ . The spectral norm of the expectation of  $M_t^T M_t$  admits a similar analysis and the same bound and is omitted.

**Lemma 3.** *Using our notation,  $\|\mathbb{E}(M_t M_t^T)\|_2 \leq n \|\widehat{A}\|_F^2$  for any  $t \in [s]$ .*

**Proof.** We start by evaluating  $\mathbb{E}(M_t M_t^T)$ ; recall that  $p_{ij} = \widehat{A}_{ij}^2 / \|\widehat{A}\|_F^2$ :

$$\begin{aligned} & \mathbb{E}(M_t M_t^T) \\ &= \mathbb{E}\left(\left(\frac{\widehat{A}_{i_t j_t}}{p_{i_t j_t}} e_{i_t} e_{j_t}^T - \widehat{A}\right)\left(\frac{\widehat{A}_{i_t j_t}}{p_{i_t j_t}} e_{j_t} e_{i_t}^T - \widehat{A}^T\right)\right) \\ &= \sum_{i,j=1}^n p_{ij} \left(\frac{\widehat{A}_{ij}}{p_{ij}} e_i e_j^T - \widehat{A}\right)\left(\frac{\widehat{A}_{ij}}{p_{ij}} e_j e_i^T - \widehat{A}^T\right) \\ &= \sum_{i,j=1}^n \left(\frac{\widehat{A}_{ij}^2}{p_{ij}} e_i e_i^T - \widehat{A}_{ij} \widehat{A} e_j e_i^T - \widehat{A}_{ij} e_i e_j^T \widehat{A}^T + p_{ij} \widehat{A} \widehat{A}^T\right) \\ &= \|\widehat{A}\|_F^2 \sum_{i=1}^n m_i \cdot e_i e_i^T - \sum_{j=1}^n \widehat{A} e_j \sum_{i=1}^n \widehat{A}_{ij} e_i^T \\ &\quad - \sum_{j=1}^n \left(\sum_{i=1}^n \widehat{A}_{ij} e_i\right) (\widehat{A} e_j)^T + \sum_{i,j=1}^n p_{ij} \widehat{A} \widehat{A}^T, \end{aligned}$$

where  $m_i$  is the number of non-zeroes of the  $i$ -th row of  $\widehat{A}$ . We now simplify the above result using a few simple observations:  $\sum_{i,j=1}^n p_{ij} = 1$ ,  $\widehat{A} e_j = \widehat{A}^{(j)}$ ,  $\sum_{i=1}^n \widehat{A}_{ij} e_i = \widehat{A}^{(j)}$ , and  $\sum_{j=1}^n \widehat{A}^{(j)} (\widehat{A}^{(j)})^T = \widehat{A} \widehat{A}^T$ . Thus, we get

$$\begin{aligned} \mathbb{E}(M_t M_t^T) &= \|\widehat{A}\|_F^2 \sum_{i=1}^n m_i \cdot e_i e_i^T - \sum_{j=1}^n \widehat{A}^{(j)} (\widehat{A}^{(j)})^T \\ &\quad - \sum_{j=1}^n \widehat{A}^{(j)} (\widehat{A}^{(j)})^T + \widehat{A} \widehat{A}^T \\ &= \|\widehat{A}\|_F^2 \sum_{i=1}^n m_i \cdot e_i e_i^T - \widehat{A} \widehat{A}^T. \end{aligned}$$

Since  $0 \leq m_i \leq n$  and using Weyl's inequality (Theorem 4.3.1 of [11]), which states that by adding a positive semi-definite matrix to a symmetric matrix all its eigenvalues will increase, we get that

$$-\widehat{A} \widehat{A}^T \preceq \mathbb{E}(M_t M_t^T) \preceq n \|\widehat{A}\|_F^2 I_n.$$

Consequently  $\|\mathbb{E}(M_t M_t^T)\|_2 = \max\{\|\widehat{A}\|_2^2, n \|\widehat{A}\|_F^2\} = n \|\widehat{A}\|_F^2$ .  $\square$

We can now apply Theorem 2 on Eq. (2) with  $\tau = \epsilon/2$ ,  $\gamma = 4n\epsilon^{-1} \|\widehat{A}\|_F^2$  (Lemma 2), and  $\rho^2 = n \|\widehat{A}\|_F^2$  (Lemma 3). Thus, we get that  $\|\widehat{A} - \widetilde{A}\|_2 \leq \epsilon/2$  holds, subject to a failure probability of at most

$$2n \exp\left(-\frac{\epsilon^2 s/8}{(1 + 4/6)n \|\widehat{A}\|_F^2}\right).$$

Bounding the failure probability by  $\delta$  and solving for  $s$ , we get that

$$s \geq \frac{14}{\epsilon^2} n \|\widehat{A}\|_F^2 \ln\left(\frac{2n}{\delta}\right).$$

Using  $\|\widehat{A}\|_F \leq \|A\|_F$  (by construction) concludes the proof of the following lemma, which is the main result of this section.

**Lemma 4.** *Using the notation of Algorithm 1, if  $s \geq 14n\epsilon^{-2} \times \|A\|_F^2 \ln(2n/\delta)$ , then, with probability at least  $1 - \delta$ ,*

$$\|\widehat{A} - \widetilde{A}\|_2 \leq \epsilon/2.$$

### Acknowledgements

We would like to thank the anonymous reviewers for numerous comments that significantly improved the presentation of our work. This research has been supported by the National Science Foundation through NSF CCF 1016501, NSF DMS 1008983, and NSF CCF 545538 awards to Petros Drineas.

### References

- [1] D. Achlioptas, F. McSherry, Fast computation of low rank matrix approximations, in: Proceedings of the Symposium on Theory of Computing (STOC), 2001, pp. 611–618.
- [2] D. Achlioptas, F. McSherry, Fast computation of low rank matrix approximations, Journal of the ACM 54 (2) (2007).
- [3] S. Arora, E. Hazan, S. Kale, Fast algorithms for approximate semidefinite programming using the multiplicative weights update method, in: IEEE Symposium on Foundations of Computer Science (FOCS), 2005, pp. 339–348.
- [4] S. Arora, E. Hazan, S. Kale, A fast random sampling algorithm for sparsifying matrices, in: Proceedings of the International Workshop on Randomization and Approximation Techniques (RANDOM), 2006, pp. 272–279.
- [5] E.J. Candès, B. Recht, Exact matrix completion via convex optimization, Foundations of Computational Mathematics 9 (3) (2009) 717–772.
- [6] E.J. Candès, T. Tao, The power of convex relaxation: Near-optimal matrix completion, IEEE Trans. Inf. Theor. 56 (5) (2010) 2053–2080.
- [7] A. d'Aspremont, Subsampling algorithms for semidefinite programming, available at arXiv:0803.1990v5, November 2009.
- [8] P. Drineas, R. Kannan, M.W. Mahoney, Fast Monte Carlo algorithms for matrices I: Approximating matrix multiplication, SIAM J. Comput. 36 (1) (2006) 132–157.
- [9] A. Gittens, J.A. Tropp, Error bounds for random matrix approximation schemes, available at arXiv:0911.4108, November 2009.

- [10] D. Gross, Recovering low-rank matrices from few coefficients in any basis, available at arXiv:0910.1879, December 2009.
- [11] R.A. Horn, C.R. Johnson, Matrix Analysis, Cambridge University Press, 1990.
- [12] N.H. Nguyen, P. Drineas, T.D. Tran, Matrix sparsification via the Khintchine inequality, manuscript, 2009.
- [13] N.H. Nguyen, P. Drineas, T.D. Tran, Tensor sparsification via a bound on the spectral norm of random tensors, available at arXiv:1005.4732, May 2010.
- [14] B. Recht, A simpler approach to matrix completion, available at arXiv:0910.0651, October 2009.
- [15] J.A. Tropp, User-friendly tail bounds for sums of random matrices, available at arXiv:1004.4389, April 2010.