## Fast Monte-Carlo Algorithms for finding Low-Rank Approximations

Alan Frieze\* Ravi Kannan<sup>†</sup> Santosh Vempala<sup>‡</sup>

December 9, 2003

#### Abstract

We consider the problem of approximating a given  $m \times n$  matrix  $\mathbf{A}$  by another matrix of specified rank k, which is much smaller than m and n. The Singular Value Decomposition (SVD) can be used to find the "best" such approximation. However, it takes time polynomial in m, n which is prohibitive for some modern applications. In this paper, we develop an algorithm which is qualitatively faster, provided we may sample the entries of the matrix according to a natural probability distribution. In many applications such sampling can be done efficiently. Our **main result** is a randomized algorithm to find the description of a matrix  $\mathbf{D}^*$  of rank at most k so that

$$||\mathbf{A} - \mathbf{D}^*||_F^2 \le \min_{\mathbf{D}, \operatorname{rank}(\mathbf{D}) \le k} ||\mathbf{A} - \mathbf{D}||_F^2 + \varepsilon ||\mathbf{A}||_F^2$$

holds with probability at least  $1 - \delta$  (for any matrix  $\mathbf{M}$ ,  $||\mathbf{M}||_F^2$  denotes the sum of the squares of all the entries of  $\mathbf{M}$ ). The algorithm takes time polynomial in  $k, 1/\varepsilon, \log(1/\delta)$  only and is *independent of m and n*. In particular, this implies that in "constant" time, it can be determined if a given matrix has a good low-rank approximation.

<sup>\*</sup>Supported in part by NSF grant CCR-9530974. Department of Mathematical Sciences, Carnegie Mellon University, Pittsburgh, PA 15213. Email: af1p@andrew.cmu.edu.

<sup>&</sup>lt;sup>†</sup>Computer Science Department, Yale University, New Haven, CT 06511. Email: kannan@cs.yale.edu.

<sup>&</sup>lt;sup>‡</sup>Department of Mathematics, M.I.T., Cambridge, MA 02139. Email: vempala@math.mit.edu

### 1 Introduction

Real-world data often has a large a number of attributes (features/dimensions). A natural question is whether in fact it is generated by a small model, i.e., with a much smaller number of parameters than the number of attributes. One way to formalize this question is the problem of finding a low-rank approximation, i.e., given an  $m \times n$  matrix  $\mathbf{A}$ , find a matrix  $\mathbf{D}$  of rank at most k so that  $||\mathbf{A} - \mathbf{D}||_F$  is as small as possible. Alternatively, if we view the rows of  $\mathbf{A}$  as points in  $\mathbf{R}^n$ , then it is the problem of finding a k-dimensional linear subspace that minimizes the sum of squared distances to the points. This problem arises in many contexts, partly because some matrix algorithms are more efficient for low-rank matrices.

The traditional Singular Value Decomposition (SVD) can be used to solve the problem in time  $O(\min\{mn^2, nm^2\})$ . For many applications motivated by information retrieval and the web, this is too slow and one needs a linear or sublinear algorithm. To speed up SVD-based low-rank approximation, [15] suggested random projection as a pre-processing step, i.e., project the rows of **A** to an  $O(\log n)$ -dimensional subspace and then find the SVD in that subspace. This reduces the worst-case complexity to  $O(mn \log n)$  for a small loss in approximation quality. This is still too high.

How fast can the problem be solved? At first sight, it seems that  $\Omega(mn)$  is a lower bound — if **A** has only a single non-zero entry, one has to examine all its entries to find a good approximation. But suppose we can sample the entries with probability proportional to their magnitudes. Then a constant-sized sample would suffice in this case!

In this paper, we show that a rank k approximation can be found in time polynomial in k and  $1/\varepsilon$  where  $\varepsilon$  is an error parameter, provided we can sample the entries of  $\mathbf{A}$  from a natural probability distribution. The sampling assumptions will be made explicit shortly and also discussed in the context of some applications. Our main result is the following.

**Theorem 1** Given an  $m \times n$  matrix  $\mathbf{A}$ , and  $k, \varepsilon, \delta$ , there is a randomized algorithm which finds the description of a matrix  $\mathbf{D}^*$  of rank at most k so that

$$||\mathbf{A} - \mathbf{D}^*||_F^2 \leq \min_{\mathbf{D}, rank(\mathbf{D}) \leq k} ||\mathbf{A} - \mathbf{D}||_F^2 + \varepsilon ||\mathbf{A}||_F^2$$

holds with probability at least  $1-\delta$ . The algorithm takes time polynomial in  $k, 1/\varepsilon, \log(1/\delta)$  only, independent of m, n. The most complex computational task is to find the first k singular values of a randomly chosen  $s \times s$  submatrix where  $s = O(k^4\varepsilon^{-3})$ . The matrix  $\mathbf{D}^*$  can be explicitly constructed from its description in O(kmn) time.

As a consequence, in  $poly(k, \frac{1}{\varepsilon})$  time, we can determine if **A** has a rank k approximation with error at most  $\varepsilon ||\mathbf{A}||_F$ .

The central idea of our approach is described as follows: We pick p rows of A independently at random, each according to a probability distribution satisfying A1. Suppose these rows form a  $p \times m$  matrix S'. The rows will be scaled to form a matrix S (step 1 of the Algorithm in section 4). It will be relatively easy (Lemma 2) to show that  $S^T S \approx A^T A$ . The intuition

for this is that the (i,j)th entry of  $A^TA$  is the dot product of the ith and jth rows of A and indeed since S has a random sample of rows of A,  $(S^TS)_{i,j}$  estimates this; the scaling is done to make this estimate unbiased. Now from standard LInear Algebra, we can get the SVD of A from the spectral decomposition (SD) of  $A^TA$ , i.e., approximately from the SD of  $S^TS$ . But repeating this, the SD of  $S^TS$  can be read off from the SVD of S which in turn can be read off from the SD of  $SS^T$ . But since  $SS^T$  is just a  $p \times p$  matrix, the problem is reduced to doing the SVD of a constant sized matrix! This still leaves the computation of  $SS^T$ . For this, we apply the sampling trick again — if we pick a sample of p columns of S, to form a  $p \times p$  matrix S, (step 2 of the algorithm), then S0 of S1. Now the SD of S3 of the algorithm indeed the SVD of S4 suffices. This then is the central computation done in step 3 of the algorithm. Besides Lemma 2, the key step in the analysis is showing that we can go from approximate left singular vectors of S5 to approximate right singular vectors with only a small loss.

We present the algorithm in Section 4. It is conceptually quite simple: we select a small random submatrix of  $\mathbf{A}$  and then compute its top k singular values and vectors. These are used to describe a rank k approximation of  $\mathbf{A}$ . Further, it is easy to construct  $\mathbf{D}^*$  from its description. It also follows that the first k singular values of  $\mathbf{A}$  can be computed to within a cumulative additive error of  $\varepsilon||\mathbf{A}||_F$ .

Clearly, the algorithm and Theorem 1 rely on the existence of a good low-rank approximation in a small sample. This observation is the main insight of the paper and we state it formally below. Let the constant c be defined as in Assumption 1 (below).

**Theorem 2** Let  $\mathbf{A}$  be an  $m \times n$  matrix. and S be a sample of s rows of  $\mathbf{A}$  from a distribution satisfying Assumption 1. Let V be the vector space spanned by S. Then with probability at least 9/10, there exist vectors  $\mathbf{y}^{(1)}, \mathbf{y}^{(2)}, \dots \mathbf{y}^{(k)}$  in V such that

$$||\mathbf{A} - \mathbf{A}(\sum_{j=1}^{k} \mathbf{y}^{(j)} \mathbf{y}^{(j)^{T}})||_{F}^{2} \le \min_{\mathbf{D}, rank(\mathbf{D}) \le k} ||\mathbf{A} - \mathbf{D}||_{F}^{2} + \frac{10k}{cs} ||\mathbf{A}||_{F}^{2}.$$
(1)

By elementary linear algebra, the matrix  $\mathbf{A} \sum_{j=1}^{k} \mathbf{y}^{(j)} \mathbf{y}^{(j)^T}$  has its rows in the span of the vectors  $\mathbf{y}^{(1)}, \dots, \mathbf{y}^{(k)}$  and therefore its rank is at most k. In the next section, we give the proof of this existence theorem. The theorem directly gives an  $O(mnk/\varepsilon + \text{poly}(k/\varepsilon))$  algorithm and suggests an algorithm whose running time is linear in m and a small polynomial in k and  $1/\varepsilon$ . Such an algorithm was developed following this paper in [4]. This and other subsequent developments are discussed in Section 6.

Throughout the paper,  $\mathbf{M}^{(i)}$  denotes the *i*th row of the matrix  $\mathbf{M}$ ,  $\mathbf{M}_{(j)}$  denotes the *j*th column and  $M_{i,j}$  is the entry in the *i*th row and *j*th column. Also, for any positive integer r, [r] denotes the set  $\{1, 2, \ldots, r\}$ .

### 1.1 Sampling assumptions

We now state in detail the sampling assumptions we make.

<sup>&</sup>lt;sup>1</sup>The right singular vectors of A are precisely the eigenvectors of  $A^T A$ 

**Assumption 1.** We can sample the rows of **A** so that row i is chosen with probability  $P_i$  satisfying

$$P_i \ge c \frac{|\mathbf{A}^{(i)}|^2}{||\mathbf{A}||_F^2}$$

for some constant  $c \leq 1$  (independent of m, n). The  $P_i$ 's are known to us (if c = 1, then we we don't need to know the  $P_i$ ).

**Assumption 2.** From any row i we can sample entry j with probability  $Q_{j|i}$  satisfying

$$Q_{j|i} \geq c \frac{P_{i,j}}{P_i}, \quad ext{where} \quad P_{i,j} = \frac{\mathbf{A}_{i,j}^2}{||\mathbf{A}||_F^2}.$$

The  $Q_{j|i}$  are known to us (again if c=1, we don't need to know the values).

If the matrix **A** has a known sparsity structure, then we might be able to set up the sampling with very little preprocessing. In particular, if the matrix **A** is dense, i.e., if for all i, j,

$$\mathbf{A}_{ij}^2 \le \frac{c'}{mn} ||\mathbf{A}||_F^2$$

for some constant c', then we we can take  $P_i = 1/m$  and  $Q_{j|i} = 1/n$  for all i, j and c = 1/c'.

In general, for any matrix, by making one pass through the entire matrix, we can set up data structures that let us sample the entries fast from then onwards — O(1) time per sample — so as to satisfy Assumptions 1 and 2. During the pass, we do several things. Suppose M is such that for all i, j

$$\mathbf{A}_{ij}^2 = 0$$
 OR  $\frac{1}{M} \le |\mathbf{A}_{ij}|^2 \le M$ .

We create  $O(\log M)$  bins; during the pass, we put into the lth bin all the entries (i,j) such that  $\frac{2^{l-1}}{M} \leq |\mathbf{A}_{ij}|^2 \leq \frac{2^l}{M}$ . We also keep track of the number of entries in each bin. After this, we pretend all entries in a bin are of equal absolute value and then it is easy to set up a sampler - the details of the data structures are straightforward. In the pass, we also set up similar data structures for each row.

## 1.2 Applications

In this section we discuss our algorithm in the context of applications that rely on low-rank approximation. We show that in several situations we can satisfy the sampling assumptions of our algorithm and thus obtain the SVD approximation more efficiently. Applications that we do not discuss include face recognition and picture compression.

### 1.2.1 Low-Rank Approximations and the Regularity Lemma

The fundamental Regularity Lemma of Szemerédi's in graph theory gives a partition of the vertex set of any graph so that "most" pairs of parts are "nearly regular". (We do not

give details here.) This lemma has a host of applications (see [14]) in graph theory. The lemma was non-constructive in that it only asserted the existence of the partition (but did not give an algorithm to find it.) Alon, Duke, Lefmann, Rödl and Yuster were finally able to give an algorithm to find such a partition in polynomial time [2]. In [9, 10], low-rank approximations of the adjacency matrix of the graph were related to regular partitions. Szemerédi's Lemma and an algorithm for constructing the partition were derived from this connection. While this is not directly relevant to our results, we point it out here as one more case where low-rank approximations are very useful. A more direct application of eigenvector computation and Szemerédi's partition is given in [11].

### 1.2.2 Latent Semantic Indexing

This is a general technique for processing a collection of "documents". We give a very cursory description of this broad area here and discuss its relation to our main problem (see [5, 6, 7, 8] for details and empirical results).

Suppose there are m documents and n "terms" which occur in the documents (terms may be all the words that occur in the documents or key words that occur in them). The model hypothesizes that there are a small number (k) of unknown "topics" which the documents are about. A topic is modelled as a probability distribution on the n terms, i.e., an n-vector of non-negative reals summing to 1. With this model on hand, it is shown (with additional assumptions) that the subspace spanned by the k best topics is close to the span of the top k singular vectors of the so-called "document-term" matrix [15]. The latter is an  $m \times n$  matrix  $\mathbf{A}$  with  $\mathbf{A}_{ij}$  being the frequency of the jth term in the ith document. Alternatively, one can define  $\mathbf{A}_{ij}$  as 0 or 1 depending upon whether the jth term occurs in the ith document.

Here we argue that, in practice, the assumptions of our algorithm are satisfied and it can be used in place of the full SVD. algorithm. It is easy to see that if we are allowed one pass through each document, we can set up data structures for sampling (ideally, the creator of a document could supply a vector of squared term frequencies). Otherwise, if no frequency is too large (this is not unreasonable since words that occur too often, so-called "buzz words", are removed from the analysis), all we need to precompute is the length  $(L_i = \sum_j \mathbf{A}_{ij})$ , of each document. This is typically available (as say "file size") and we pick a document with probability proportional to its length. This is easily seen to satisfy Assumption 1, but without the squares (i.e. we sample the *i*th entry with probability  $L_i/\sum_j L_j$ ). The assumption with the squares is satisfied if all the frequencies are small. Assumption 2 is similarly implemented — given a document, we pick a word uniformly at random from it, i.e.,  $Q_{j|i} = \mathbf{A}_{ij}/L_i$ .

### 1.2.3 Web Search model

Kleinberg [13] proposed an algorithm for the problem of finding the most "important" documents from the set of documents returned by a web search that works by analyzing the  $m \times m$  hyperlink matrix. This matrix **A** has entries  $\mathbf{A}_{ij}$  equal to 1 or 0 depending upon whether the *i*'th document points to the *j*'th. The algorithm sets out to find two unit-length m-vectors  $\mathbf{x}$ ,  $\mathbf{y}$  such that  $\mathbf{x}^T \mathbf{A} \mathbf{y}$  is maximized. This is of course the problem of

finding the singular vectors of  $\mathbf{A}$ . When the keyword has multiple meanings, not only the top, but some of the other singular vectors (with large singular values) are interesting. So, it is of interest to find the largest k singular vectors for some small k.

It is worthwhile to consider our assumptions in this case. For Assumption 1, it is sufficient to sample the documents (roughly) according to the number of hypertext links from them. For Assumption 2, it is sufficient to be able to follow a random link from a document.

## 2 The Singular Value Decomposition

The matrix A can be expressed

$$\mathbf{A} = \sum_{t=1}^{r} \sigma_t \mathbf{u}^{(t)} \mathbf{v}^{(t)^T}$$

where  $\sigma_1 \geq \sigma_2 \geq \ldots \geq \sigma_r \geq 0$  and the  $\mathbf{u}^{(t)}$  form an orthonormal set of vectors as do the  $\mathbf{v}^{(t)}$ . Also  $\mathbf{u}^{(t)^T} \mathbf{A} = \sigma_t \mathbf{v}^{(t)^T}$  and  $\mathbf{A} \mathbf{v}^{(t)} = \sigma_t \mathbf{u}^{(t)}$  for  $1 \leq t \leq r$ . This is called the *singular value decomposition* of  $\mathbf{A}$ . Here r is the rank of A.

By a theorem of Eckart and Young [12], the matrix  $\mathbf{D}_k$  that minimizes of  $||\mathbf{A} - \mathbf{D}||_F$  among all matrices  $\mathbf{D}$  of rank k or less is given by

$$\mathbf{D}_k = \sum_{t=1}^k \mathbf{A} \mathbf{v}^{(t)} \mathbf{v}^{(t)^T}.$$

This implies that

$$||\mathbf{D}_k||_F^2 = \sum_{t=1}^k \sigma_t^2$$
 and  $||\mathbf{A} - \mathbf{D}_k||_F^2 = \sum_{t=k+1}^r \sigma_t^2$ .

We use this notation throughout the paper.

# 3 A small sample contains a good approximation

The goal of this section is to prove Theorem 2, namely the subspace spanned by a sample of rows chosen according to Assumption 1 contains an approximation to **A** that is nearly the best possible.

Let S be a random sample of s rows chosen from a distribution that satisfies Assumption 1. For t = 1, 2, ..., r, we define the vector-valued random variable

$$\mathbf{w}^{(t)} = \frac{1}{s} \sum_{i \in S} \frac{\mathbf{u}_i^{(t)}}{P_i} \mathbf{A}^{(i)}.$$

Then the vectors  $\mathbf{w}^{(t)}$  are clearly in the subspace generated by S. Moreover,

$$\mathbf{E}(\mathbf{w}^{(t)}) = \mathbf{A}^T \mathbf{u}^{(t)} = \sigma_t \mathbf{v}^{(t)},$$

and since  $P_i \geq c \frac{|\mathbf{A}^{(i)}|^2}{\|\mathbf{A}\|_F^2}$ , we have

$$\mathbf{E}(|\mathbf{w}^{(t)} - \sigma_t \mathbf{v}^{(t)}|^2) = \frac{1}{s} \left( \sum_{i=1}^m \frac{|\mathbf{u}_i^{(t)}|^2 |\mathbf{A}^{(i)}|^2}{P_i} \right) - \frac{\sigma_t^2}{s} \le \frac{1}{sc} ||\mathbf{A}||_F^2.$$
 (2)

If we had  $\mathbf{w}^{(t)}$  exactly equal to  $\sigma_t \mathbf{v}^{(t)}$  (instead of just in expectation), then it is easy to see that

$$\mathbf{A} \sum_{t=1}^{k} \mathbf{v} t \mathbf{v}^{(t)^{T}} = \mathbf{A} \sum_{t=1}^{k} \frac{1}{\sigma_{t}^{2}} \mathbf{w} t \mathbf{w}^{(t)^{T}}$$

and this would be sufficient to prove the theorem. We wish to carry this out approximately.

To this end, define  $\hat{\mathbf{y}}^{(t)} = \frac{1}{\sigma_t} \mathbf{w}^{(t)}$  for t = 1, 2, ..., r and let  $V_1 = \operatorname{span}(\hat{\mathbf{y}}^{(1)}, \hat{\mathbf{y}}^{(2)}, ..., \hat{\mathbf{y}}^{(k)}) \subseteq V$ . Let  $\mathbf{y}^{(1)}, \mathbf{y}^{(2)}, ..., \mathbf{y}^{(n)}$  be an orthonormal basis of  $\mathbf{R}^n$  with  $V_1 = \operatorname{span}(\mathbf{y}^{(1)}, \mathbf{y}^{(2)}, ..., \mathbf{y}^{(l)})$  for some  $l \leq k$ . Let

$$\mathbf{F} = \sum_{t=1}^{l} \mathbf{A} \mathbf{y}^{(t)} \mathbf{y}^{(t)^T}$$
 and  $\hat{\mathbf{F}} = \sum_{t=1}^{k} \mathbf{A} \mathbf{v}^{(t)} \hat{\mathbf{y}}^{(t)^T}$ 

The matrix **F** will be our candidate approximation to **A** in the span of *S*. We will bound its error using  $\hat{\mathbf{F}}$ . Note that for any  $i \leq k$  and j > k,  $\hat{\mathbf{y}}^{(i)^T)}\mathbf{y}^{(j)} = 0$ . Thus,

$$||\mathbf{A} - \mathbf{F}||_F^2 = \sum_{i=1}^n |(\mathbf{A} - \mathbf{F})\mathbf{y}^{(i)}|^2$$

$$= \sum_{i=l+1}^n |\mathbf{A}\mathbf{y}^{(i)}|^2$$

$$= \sum_{i=l+1}^n |(\mathbf{A} - \hat{\mathbf{F}})\mathbf{y}^{(i)}|^2$$

$$\leq ||\mathbf{A} - \hat{\mathbf{F}}||_F^2.$$
(3)

Also,

$$\begin{aligned} ||\mathbf{A} - \hat{\mathbf{F}}||_F^2 &= \sum_{i=1}^n |\mathbf{u}^{(i)^T} (\mathbf{A} - \hat{\mathbf{F}})|^2 \\ &= \sum_{i=1}^k |\sigma_i \mathbf{v}^{(i)} - \mathbf{w}^{(i)}|^2 + \sum_{i=k+1}^n \sigma_i^2. \end{aligned}$$

Taking expectations and using (2) we get

$$\mathbf{E}(||\mathbf{A} - \hat{\mathbf{F}}||_F^2) \le \sum_{i=k+1}^n \sigma_i^2 + \frac{k}{sc} ||\mathbf{A}||_F^2.$$
 (4)

Since  $\hat{\mathbf{F}}$  is of rank at most k we have

$$||\mathbf{A} - \hat{\mathbf{F}}||_F^2 \ge ||\mathbf{A} - \mathbf{D}_k||_F^2 = \sum_{i=k+1}^n \sigma_i^2.$$

Thus  $||\mathbf{A} - \hat{\mathbf{F}}||_F^2 - ||\mathbf{A} - \mathbf{D}_k||_F^2$  is a non-negative random variable and (4) implies

$$\mathbf{Pr}(||\mathbf{A} - \hat{\mathbf{F}}||_F^2 - ||\mathbf{A} - \mathbf{D}_k||_F^2 \ge \frac{10k}{sc}||\mathbf{A}||_F^2) \le \frac{1}{10}.$$

From (3) it follows that

$$|\mathbf{Pr}(||\mathbf{A} - \mathbf{F}||_F^2 - ||\mathbf{A} - \mathbf{D}_k||_F^2 \ge \frac{10k}{sc}||\mathbf{A}||_F^2) \le \frac{1}{10}$$

as required.

We next observe that a good low-rank approximation with respect to Frobenius norm implies a good low-rank approximation with respect to the norm  $||\mathbf{M}|| = \max_{|\mathbf{x}|=1} |\mathbf{M}\mathbf{x}|$ .

### Theorem 3 If

$$||\mathbf{A} - \mathbf{A} \sum_{t=1}^k \mathbf{y}^{(t)} \mathbf{y}^{(t)^T}||_F^2 \le ||\mathbf{A} - \mathbf{D}_k||_F^2 + \varepsilon ||\mathbf{A}||_F^2.$$

Then

$$||\mathbf{A} - \mathbf{A} \sum_{t=1}^k \mathbf{y}^{(t)} \mathbf{y}^{(t)^T}||^2 \le (\frac{1}{k+1} + \varepsilon)||\mathbf{A}||_F^2.$$

**Proof.** Let  $\mathbf{B} = \mathbf{A} - \mathbf{A} \sum_{t=1}^{k} \mathbf{y}^{(t)} \mathbf{y}^{(t)^{T}}$ . Suppose that  $\mathbf{B}$  has a unit eigenvector  $\mathbf{x}$  with eigenvalue  $\lambda$  such that

$$\lambda^2 > (\frac{1}{k+1} + \varepsilon) ||\mathbf{A}||_F^2$$

Then we see that

$$||\mathbf{B} - \mathbf{B} \mathbf{x} \mathbf{x}^{T}||_{F}^{2} = ||\mathbf{B}||_{F}^{2} - \lambda^{2} < \sum_{i=k+1}^{n} \sigma_{i}^{2} - \frac{1}{k+1} ||\mathbf{A}||_{F}^{2}.$$
 (5)

The rank of the matrix  $\mathbf{A} \sum_{t=1}^{k} \mathbf{y}^{(t)} \mathbf{y}^{(t)^T} + \mathbf{B} \mathbf{x} \mathbf{x}^T$  is at most k+1, and so

$$\begin{aligned} ||\mathbf{A} - \mathbf{A} \sum_{t=1}^{k} \mathbf{y}^{(t)} \mathbf{y}^{(t)^{T}} - \mathbf{B} \mathbf{x} \mathbf{x}^{T}||_{F}^{2} & \geq ||\mathbf{A} - \mathbf{D}_{k+1}||_{F}^{2} \\ &= \sum_{t=k+2}^{n} \sigma_{t}^{2} \\ &\geq ||\mathbf{A} - \mathbf{D}_{k}||_{F}^{2} - \frac{1}{k+1} ||\mathbf{A}||_{F}^{2}, \end{aligned}$$

which contradicts (5).

# 4 Sampling Algorithm

In this section we present the main "constant time" algorithm to produce the approximation of Theorem 1. What we do below is to first pick a set of p rows of A. We form a matrix S

from these rows after scaling them. We then pick again p columns of S from a probability distribution satisfying a condition of the type stated in Assumption 1 and scale the columns to get a  $p \times p$  matrix W. We find the singular vectors of this matrix and from those, show how to get a good low-rank approximation to A.

### Algorithm

Input: Matrix **A** and error parameter  $\varepsilon > 0$ . Set  $p = \frac{10^7 k^4}{c^3 \varepsilon^3}$ .

1. (Sample rows) Independently choose (rows)  $i_1, i_2, \ldots, i_p$  according to distribution  $P = (P_1, P_2, \ldots, P_m)$  which satisfies Assumption 1, i.e.,

$$P_i \ge c \frac{|\mathbf{A}^{(i)}|^2}{||\mathbf{A}||_F^2}.$$

Let **S** be the  $p \times n$  matrix with rows  $\mathbf{A}^{(i_t)}/\sqrt{pP_{i_t}}$  for  $t = 1, 2, \dots, p$ . Note that if c = 1, the this scaling amounts to normalizing all rows to be of the same length.

2. (**Sample columns**) Independently choose (columns)  $j_1, j_2, \ldots, j_p$  (of **S**) according to a distribution  $P' = (P'_1, P'_2, \ldots, P'_n)$  which satisfies

$$P'_j \ge \frac{c}{2} \frac{|\mathbf{S}_{(j)}|^2}{||\mathbf{S}||_F^2}.$$

(we show below how to do this using Assumption 2.)

Let **W** be the  $p \times p$  matrix with columns  $\mathbf{S}^{(j_t)} / \sqrt{p P'_{j_t}}$  for  $t = 1, 2, \dots, p$ .

3. (Compute SVD) Compute the top k singular vectors  $\mathbf{u}^{(1)}, \mathbf{u}^{(2)}, \dots, \mathbf{u}^{(k)}$  in the column space of  $\mathbf{W}$ .

4. (Filter) Let

$$T = \{t : |\mathbf{W}^T \mathbf{u}^{(t)}|^2 \ge \gamma ||\mathbf{W}||_F^2\},$$

where

$$\gamma = \frac{c\varepsilon}{8k}.$$

For  $t \in T$  let

$$\mathbf{v}^{(t)} = rac{\mathbf{S}^T \mathbf{u}^{(t)}}{|\mathbf{W}^T \mathbf{u}^{(t)}|}$$

Output  $\mathbf{v}^{(t)}$  for  $t \in T$  (the low rank approximation to  $\mathbf{A}$  is  $\mathbf{A} \sum_{t \in T} \mathbf{v}^{(t)} \mathbf{v}^{(t)^T}$ ).

We explore some issues related to the implementation of the above algorithm.

First of all, how do we carry out Step 2? We first pick a row of **S**, each row with probability 1/p; suppose the chosen row is the *i*th row of **A**. Then pick  $j \in \{1, 2, ..., n\}$  with probabilities  $Q_{j|i}$ . This defines the probabilities  $P'_{j}$ . We then have (with  $I = i : \mathbf{A}^{(i)}$  is a row of **S**}),

$$P_j' = \sum_{i \in I} \frac{Q_{j|i}}{p} \geq \sum_{i \in I} \frac{cP_{i,j}}{pP_i} = \sum_{i \in I} \frac{c\mathbf{A}_{i,j}^2}{pP_i||\mathbf{A}||_F^2} = \frac{c}{||\mathbf{A}||_F^2} \sum_{i \in I} \frac{\mathbf{A}_{i,j}^2}{pP_i} = c\frac{|\mathbf{S}_{(j)}|^2}{||\mathbf{A}||_F^2} \geq \frac{c}{2} \frac{|\mathbf{S}_{(j)}|^2}{||\mathbf{S}||_F^2}$$

where the last step is implied by the next lemma.

**Lemma 1** For **S** chosen as in the algorithm, with high probability,

$$||\mathbf{S}||_F^2 \ge \frac{1}{2}||\mathbf{A}||_F^2 \text{ and } ||\mathbf{W}||_F^2 \ge \frac{1}{2}||\mathbf{S}||_F^2.$$

*Proof.* By a routine calculation,

$$\mathbf{E}(||\mathbf{S}||_F^2) = ||\mathbf{A}||_F^2 \quad \text{ and } \quad \mathbf{Var}(||\mathbf{S}||_F^2) \leq \frac{1}{c^2 p} ||\mathbf{A}||_F^4.$$

Next, observe that for any row i of S,

$$||\mathbf{S}^{(i)}||_F^2 \le \frac{||\mathbf{A}||_F^2}{cp}$$

The random variable  $||\mathbf{S}||_F^2$  is a sum of p independent random variables. Therefore,

$$\mathbf{Var}(||\mathbf{S}||_F^2||) = p\mathbf{Var}(||\mathbf{S}^{(i)}||_F^2) \le p\mathbf{E}(||\mathbf{S}^{(i)}||_F^4) \le \frac{1}{c^2p}||\mathbf{A}||_F^4.$$

The lemma now follows using Chebychef's inequality.

## 5 Analysis

The next lemma asserts that a sample N of rows of a matrix M provides a good approximation to M in the sense that  $N^TN$  is close to  $M^TM$ . This will be a key tool in the analysis.

**Lemma 2** Let **M** be an  $a \times b$  matrix and let  $Q = Q_1, Q_2, \ldots, Q_a$  be a probability distribution on  $\{1, 2, \ldots, a\}$  such that

$$Q_i \ge \alpha \frac{|\mathbf{M}^{(i)}|^2}{||\mathbf{M}||_F^2}, \qquad i = 1, 2, \dots, a$$

for some  $0 < \alpha < 1$ . Let  $\sigma = (i_1, i_2, \dots, i_p)$  be a sequence of p independent samples from [a], each chosen according to distribution Q. Let  $\mathbf{N}$  be the  $p \times b$  matrix with

$$\mathbf{N}^{(t)} = rac{\mathbf{M}^{(i_t)}}{\sqrt{pQ_{i_t}}} \qquad t = 1, 2, \dots, p.$$

Then for all  $\theta > 0$ ,

$$|\mathbf{Pr}(||\mathbf{M}^T\mathbf{M} - \mathbf{N}^T\mathbf{N}||_F \ge \theta ||\mathbf{M}||_F^2) \le \frac{1}{\theta^2 \alpha p}.$$

Proof.

$$\begin{aligned} ||\mathbf{M}^T \mathbf{M} - \mathbf{N}^T \mathbf{N}||_F^2 &= \sum_{r,s=1}^b |\mathbf{M}_{(r)}^T \mathbf{M}_{(s)} - \mathbf{N}_{(r)}^T \mathbf{N}_{(s)}|^2 \\ \mathbf{E}(\mathbf{N}_{(r)}^T \mathbf{N}_{(s)}) &= \sum_{t=1}^p \mathbf{E}(\mathbf{N}_{i_t,r} \mathbf{N}_{i_t,s}) \\ &= \sum_{t=1}^p \sum_{i=1}^a Q_i \frac{\mathbf{M}_{i,r} \mathbf{M}_{i,s}}{pQ_i} \\ &= \mathbf{M}_{(r)}^T \mathbf{M}_{(s)} \end{aligned}$$

$$\begin{split} \mathbf{E}(|\mathbf{N}_{(r)}^{T}\mathbf{N}_{(s)} - \mathbf{M}_{(r)}^{T}\mathbf{M}_{(s)}|^{2}) & \leq \sum_{t=1}^{p} \mathbf{E}((\mathbf{N}_{i_{t},r}\mathbf{N}_{i_{t},s})^{2}) \\ & = \sum_{t=1}^{p} \sum_{i=1}^{a} Q_{i} \frac{\mathbf{M}_{i,r}^{2}\mathbf{M}_{i,s}^{2}}{p^{2}Q_{i}^{2}} \\ & \leq \frac{||\mathbf{M}||_{F}^{2}}{\alpha p^{2}} \sum_{t=1}^{p} \sum_{i=1}^{a} \frac{\mathbf{M}_{i,r}^{2}\mathbf{M}_{i,s}^{2}}{|\mathbf{M}^{(i)}|^{2}} \\ & = \frac{||\mathbf{M}||_{F}^{2}}{\alpha p} \sum_{i=1}^{a} \frac{\mathbf{M}_{i,r}^{2}\mathbf{M}_{i,s}^{2}}{|\mathbf{M}^{(i)}|^{2}}. \end{split}$$

Thus,

$$\begin{split} \mathbf{E}(||\mathbf{M}^T\mathbf{M} - \mathbf{N}^T\mathbf{N}||_F^2) &= \sum_{r,s=1}^b \mathbf{E}(|\mathbf{N}_{(r)}^T\mathbf{N}_{(s)} - \mathbf{M}_{(r)}^T\mathbf{M}_{(s)}|^2) \\ &\leq \frac{||\mathbf{M}||_F^2}{\alpha p} \sum_{i=1}^a \frac{1}{|\mathbf{M}^{(i)}|^2} \sum_{r,s=1}^b (\mathbf{M}_{i,r}\mathbf{M}_{i,s})^2 \\ &= \frac{||\mathbf{M}||_F^4}{\alpha p}. \end{split}$$

The result follows from Markov's inequality.

We introduce some notation for the rest of the proof. For a matrix  $\mathbf{M}$  and vectors  $\mathbf{x}^{(i)}, i \in I$  we define

$$\Delta(\mathbf{M}; \mathbf{x}^{(i)}, i \in I) = ||\mathbf{M}||_F^2 - ||\mathbf{M} - \mathbf{M} \sum_{i \in I} \mathbf{x}^{(i)} \mathbf{x}^{(i)^T}||_F^2$$

When the  $\mathbf{x}^{(i)}$  are orthogonal unit vectors, this represents the norm of the projection of  $\mathbf{M}$  to the subspace spanned by the  $\mathbf{x}^{(i)}$ :

$$\Delta(\mathbf{M}; \mathbf{x}^{(i)}, i \in I) = \sum_{i \in I} \mathbf{x}^{(i)^T} \mathbf{M}^T \mathbf{M} \mathbf{x}^{(i)}.$$

Thus, if  $\mathbf{x}^{(t)}, t \in [k]$  are the top k singular vectors of **A**, then

$$\Delta(\mathbf{A}; \mathbf{x}^{(t)}, t \in [k]) = \sum_{t=1}^{k} \sigma_t^2.$$

Lemma 3 Let A, S be matrices with the same number of columns, and

$$||\mathbf{A}^T\mathbf{A} - \mathbf{S}^T\mathbf{S}|| < \theta ||\mathbf{A}||_F^2$$

1. For any pair of unit vectors  $\mathbf{z}, \mathbf{z}'$  in the row space of  $\mathbf{A}$ ,

$$|\mathbf{z}^T \mathbf{A}^T \mathbf{A} \mathbf{z}' - \mathbf{z}^T \mathbf{S}^T \mathbf{S} \mathbf{z}'| \leq \theta ||\mathbf{A}||_F^2$$

2. For any set of unit vectors  $\mathbf{z}^{(1)}, \mathbf{z}^{(2)}, \dots, \mathbf{z}^{(\ell)}, \ell \leq k$  in the row space of  $\mathbf{A}$ ,

$$|\Delta(\mathbf{A}; \mathbf{z}^{(i)}, i \in [\ell]) - \Delta(\mathbf{S}; \mathbf{z}^{(i)}, i \in [\ell])| \le k^2 \theta ||\mathbf{A}||_F^2$$

**Proof.** The first part of the lemma is easy. For the second, using the fact that  $||\mathbf{N}||_F^2 = \text{Tr}(\mathbf{N}\mathbf{N}^T)$  for any matrix  $\mathbf{N}$ , we see that  $\Delta(\mathbf{M}; \mathbf{x}^{(i)}, i \in I)$  equals

$$\begin{split} &2\sum_{i\in I}\mathrm{Tr}(\mathbf{M}\mathbf{x}^{(i)}\mathbf{x}^{(i)^T}\mathbf{M}^T) - \sum_{i,i'\in I}(\mathbf{x}^{(i)^T}\mathbf{x}(i')^T)\mathrm{Tr}(\mathbf{M}\mathbf{x}^{(i)}\mathbf{x}^{(i')^T}\mathbf{M}^T\\ &= &2\sum_{i\in I}\mathbf{x}^{(i)^T}\mathbf{M}^T\mathbf{M}\mathbf{x}^{(i)} - \sum_{i\in I}(|\mathbf{x}^{(i)}|^2)\mathbf{x}^{(i)^T}\mathbf{M}^T\mathbf{M}\mathbf{x}^{(i)} - \sum_{i\neq i'\in I}(\mathbf{x}^{(i)^T}\mathbf{x}^{(i')^T})\mathbf{x}^{(i')^T}\mathbf{M}^T\mathbf{M}\mathbf{x}^{(i)}. \end{split}$$

From this the second part follows.

We are now ready to prove the main theorem.

**Proof of Theorem 1.** We will apply Lemma 2 twice, once to the row sample and once to the induced column sample. It follows from the above lemma and the value of p that with probability at least 9/10 both of the following events hold:

$$||\mathbf{A}^T \mathbf{A} - \mathbf{S}^T \mathbf{S}||_F \le \theta ||\mathbf{A}||_F^2 \quad \text{and } ||\mathbf{S}\mathbf{S}^T - \mathbf{W}\mathbf{W}^T||_F \le \theta ||\mathbf{S}||_F^2$$
 (6)

where

$$\theta = \sqrt{\frac{40}{cp}} = \frac{\varepsilon^{3/2}c}{500k^2}.$$

Assume from now on that these events occur.

It follows from Theorem 2 that with probability at least 9/10 there are unit vectors  $\mathbf{x}^{(t)}, t \in [k]$  in the row space of **S** such that

$$\Delta(\mathbf{A}; \mathbf{x}^{(t)}, t \in [k]) \ge ||\mathbf{A}||_F^2 - ||\mathbf{A} - \mathbf{D}_k||_F^2 - \frac{10k}{cp} \varepsilon ||\mathbf{A}||_F^2 \ge ||\mathbf{D}_k||_F^2 - \frac{\varepsilon}{8} ||\mathbf{A}||_F^2.$$

Applying the second part of Lemma 3 to  $\mathbf{A}, \mathbf{S}$  and the vectors  $\mathbf{x}^{(i)}$ , we see that

$$\Delta(\mathbf{S}; \mathbf{x}^{(t)}, t \in [k]) \ge \Delta(\mathbf{A}; \mathbf{x}^{(t)}, t \in [k]) - k^2 \theta ||\mathbf{A}||_F^2 \ge ||\mathbf{D}_k||_F^2 - \frac{\varepsilon}{4} ||\mathbf{A}||_F^2.$$

Now, **S** and **S**<sup>T</sup> have the same singular values and so there exist unit vectors  $\mathbf{y}^{(t)}$ ,  $t \in [k]$  in the column space of **S** such that

$$\Delta(\mathbf{S}^T; \mathbf{y}^{(t)}, t \in [k]) \ge ||\mathbf{D}_k||_F^2 - \frac{1}{4}\varepsilon||\mathbf{A}||_F^2.$$

Applying Theorem 2 to  $\mathbf{S}^T$  and  $\mathbf{W}^T$ , we see that with probability at least 9/10 there are unit vectors  $\mathbf{z}^{(t)}$ ,  $t \in [k]$  in the column space of  $\mathbf{W}$  such that

$$\Delta(\mathbf{S}^T; \mathbf{z}^{(t)}, t \in [k]) \ge \Delta(\mathbf{S}^T; \mathbf{y}^{(t)}, t \in [k]) - \frac{10k}{cp} ||\mathbf{S}||_F^2 \ge ||\mathbf{D}_k||_F^2 - \frac{3\varepsilon}{8} ||\mathbf{A}||_F^2$$

Applying the second part of Lemma 3 to  $\mathbf{S}^T$ ,  $\mathbf{W}^T$  and the vectors  $\mathbf{z}^{(t)}$ , we see that

$$\Delta(\mathbf{W}^T; \mathbf{z}^{(t)}, t \in [k]) \ge \Delta(\mathbf{S}^T; \mathbf{z}^{(t)}, t \in [k]) - k^2 \theta ||\mathbf{S}||_F^2 \ge ||\mathbf{D}_k||_F^2 - \frac{\varepsilon}{2} ||\mathbf{A}||_F^2.$$

Therefore, the vectors  $\mathbf{u}^{(t)}, t \in [k]$  computed by the algorithm satisfy

$$\Delta(\mathbf{W}^T; \mathbf{u}^{(t)}, t \in [k]) \ge ||\mathbf{D}_k||_F^2 - \frac{\varepsilon}{2} ||\mathbf{A}||_F^2.$$

Note that the highest possible value of  $\Delta(\mathbf{A}; \mathbf{x}^{(t)}, t \in [k])$  is  $||\mathbf{D}_k||_F^2$ . All that remains to show is that in fact  $\Delta(\mathbf{W}^T, .)$  being large implies that  $\Delta(\mathbf{A}, .)$  is large. For this, we construct a suitable set of vectors (as in the algorithm).

Since  $\mathbf{u}^{(t)}, t \in [k]$  are singular vectors,

$$\Delta(\mathbf{W}^T; \mathbf{u}^{(t)}, t \in T) \ge \Delta(\mathbf{W}^T; \mathbf{u}^{(t)}, t \in [k]) - k\gamma ||\mathbf{W}||_F^2 \ge ||\mathbf{D}_k||_F^2 - \frac{5}{8}\varepsilon ||\mathbf{A}||_F^2.$$

Applying Lemma 3 again, this time to  $\mathbf{S}^T$ ,  $\mathbf{W}^T$  and the vectors  $\mathbf{u}^{(t)}$ ,  $t \in T$ , it follows that

$$\Delta(\mathbf{S}^T; \mathbf{u}^{(t)}, t \in T) \ge \Delta(\mathbf{W}^T; \mathbf{u}^{(t)}, t \in T) - k^2 \theta ||\mathbf{S}||_F^2 \ge ||\mathbf{D}_k||_F^2 - \frac{3}{4} \varepsilon ||\mathbf{A}||_F^2.$$

The next and crucial step, is to switch from  $\mathbf{u}^{(t)}$  in the column space of  $\mathbf{S}$  to  $\mathbf{v}^{(t)}$  in the row space of  $\mathbf{S}$ . This is achieved by the following claims whose proof we defer to Section 5.1. For  $t \in T$ ,

Claim 1. 
$$\Delta(\mathbf{S}; \mathbf{v}^{(t)}, t \in T) \ge \Delta(\mathbf{S}^T; \mathbf{u}^{(t)}, t \in T) - \frac{1}{8}\varepsilon ||\mathbf{A}||_F^2$$

Claim 2. 
$$|{\bf v}^{(t)}|^2 \le 1 + \frac{\varepsilon}{16}$$
.

It follows from Lemma 3 that

$$\Delta(\mathbf{A}; \mathbf{v}^{(t)}, t \in T) \ge \Delta(\mathbf{S}; \mathbf{v}^{(t)}, t \in T) - (1 + \frac{\varepsilon}{16})k^2\theta ||\mathbf{A}||_F^2 \ge ||\mathbf{D}_k||_F^2 - \varepsilon ||\mathbf{A}||_F^2$$

(assuming  $\varepsilon < 16$ ). Thus,

$$||\mathbf{A}||_F^2 - ||\mathbf{A} - \mathbf{A} \sum_{t \in T} \mathbf{v}^{(t)} \mathbf{v}^{(t)^T}||_F^2 \ge ||\mathbf{D}_k||_F^2 - \varepsilon ||\mathbf{A}||_F^2.$$

Rearranging terms, we get the conclusion of the theorem.

$$||\mathbf{A} - \mathbf{A} \sum_{t \in T} \mathbf{v}^{(t)} \mathbf{v}^{(t)^T}||_F^2 \le ||\mathbf{A} - \mathbf{D}_k||_F^2 - \varepsilon ||\mathbf{A}||_F^2.$$

### 5.1 Proof of the claims

Observe first that

$$||\mathbf{S}\mathbf{S}^{T}\mathbf{S}\mathbf{S}^{T} - \mathbf{W}\mathbf{W}^{T}\mathbf{W}\mathbf{W}^{T}||_{F} \leq ||\mathbf{S}\mathbf{S}^{T}(\mathbf{S}\mathbf{S}^{T} - \mathbf{W}\mathbf{W}^{T})||_{F} + ||(\mathbf{S}\mathbf{S}^{T} - \mathbf{W}\mathbf{W}^{T})\mathbf{W}\mathbf{W}^{T}||_{F}$$

$$\leq \theta||\mathbf{S}||_{F}^{2}(||\mathbf{S}||_{F}^{2} + ||\mathbf{W}||_{F}^{2}), \tag{7}$$

and that for  $t \neq t' \in T$ ,

$$\mathbf{u}^{(t)^T} \mathbf{W} \mathbf{W}^T \mathbf{u}^{(t')} = \mathbf{u}^{(t)^T} \mathbf{W} \mathbf{W}^T \mathbf{W} \mathbf{W}^T \mathbf{u}^{(t')} = 0.$$

Now consider  $t \neq t' \in T$ . Then

$$(\mathbf{v}^{(t)^T}\mathbf{v}^{(t')})(\mathbf{v}^{(t)^T}\mathbf{S}^T\mathbf{S}\mathbf{v}^{(t')}) = \frac{(\mathbf{u}^{(t)^T}\mathbf{S}\mathbf{S}^T\mathbf{u}^{(t')})(\mathbf{u}^{(t)^T}\mathbf{S}\mathbf{S}^T\mathbf{S}\mathbf{S}^T\mathbf{u}^{(t')})}{|\mathbf{W}^T\mathbf{u}^{(t)}|^2|\mathbf{W}^T\mathbf{u}^{(t')}|^2}.$$

Furthermore,

$$|\mathbf{u}^{(t)^T}\mathbf{S}\mathbf{S}^T\mathbf{u}^{(t')}| = |\mathbf{u}^{(t)^T}(\mathbf{S}\mathbf{S}^T - \mathbf{W}\mathbf{W}^T)\mathbf{u}^{(t')}| \leq \theta||\mathbf{S}||_F^2.$$

Similarly, using (7),

$$|\mathbf{u}^{(t)^T} \mathbf{S} \mathbf{S}^T \mathbf{S} \mathbf{S}^T \mathbf{u}^{(t')}| \le \theta ||\mathbf{S}||_F^2 (||\mathbf{S}||_F^2 + ||\mathbf{W}||_F^2).$$

Hence,

$$|(\mathbf{v}^{(t)^T}\mathbf{v}^{(t')})(\mathbf{v}^{(t)^T}\mathbf{S}^T\mathbf{S}\mathbf{v}^{(t')})| \leq \frac{\theta^2||\mathbf{S}||_F^2(||\mathbf{S}||_F^2 + ||\mathbf{W}||_F^2)}{\gamma^2||\mathbf{W}||_F^4} \leq \frac{12\theta^2}{\gamma^2c^2}||\mathbf{A}||_F^2.$$

Next, for any vector  $\mathbf{u}$  and any matrix  $\mathbf{S}$ 

$$\frac{|\mathbf{S}\mathbf{S}^T\mathbf{u}|}{|\mathbf{S}^T\mathbf{u}|} \ge \frac{|\mathbf{S}^T\mathbf{u}|}{|\mathbf{u}|}.$$

So for  $t \in T$ 

$$\mathbf{v}^{(t)^T}\mathbf{S}^T\mathbf{S}\mathbf{v}^{(t)} = rac{\mathbf{u}^{(t)^T}\mathbf{S}\mathbf{S}^T\mathbf{S}\mathbf{S}^T\mathbf{u}^{(t)}}{|\mathbf{W}^T\mathbf{u}^{(t)}|^2} \geq rac{|\mathbf{S}^T\mathbf{u}^{(t)}|^4}{|\mathbf{W}^T\mathbf{u}^{(t)}|^2}.$$

Observe that the first part of Lemma 3 implies

$$|\mathbf{S}^T \mathbf{u}^{(t)}|^2 - |\mathbf{W}^T \mathbf{u}^{(t)}|^2 \le \theta ||\mathbf{S}||_F^2$$

So,

$$\left| \frac{|\mathbf{S}^T \mathbf{u}^{(t)}|^2}{|\mathbf{W}^T \mathbf{u}^{(t)}|^2} - 1 \right| \le \frac{2\theta}{\gamma} \le \frac{\varepsilon}{16}.$$
 (8)

Claim 2 follows immediately.

We then have

$$\sum_{t \in T} \mathbf{v}^{(t)^T} \mathbf{S}^T \mathbf{S} \mathbf{v}^{(t)} \ge (1 - \frac{\varepsilon}{16}) \sum_{t \in T} \mathbf{u}^{(t)^T} \mathbf{S} \mathbf{S}^T \mathbf{u}^{(t)} = (1 - \frac{\varepsilon}{16}) \Delta(\mathbf{S}^T; \mathbf{u}^{(t)}, t \in T).$$

So,

$$\Delta(\mathbf{S}; \mathbf{v}^{(t)}, t \in T) \ge (1 - \frac{\varepsilon}{16}) \Delta(\mathbf{S}^T; \mathbf{u}^{(t)}, t \in T) - \frac{12k^2\theta^2}{\gamma^2c^2} ||\mathbf{A}||_F^2,$$

which completes the proof of Claim 1.

## 6 Recent work

There have been several developments on the problem of low-rank approximation since a preliminary version of this paper appeared. Drineas et al. [4] give an algorithm whose running time is  $O(mr^2+r^3)$  where  $r=O(k/\varepsilon^2)$ . Although this is theoretically much slower (due to the dependence on m), in practice, the better dependence on k and  $1/\varepsilon$  might make it more practical. An alternative sampling based algorithm was given in [1] with similar bounds on the complexity (i.e., linear in m, polynomial in  $k, 1/\varepsilon$ , and the Frobenius norm approximation (for the 2-norm, they get better bounds). In [3], a lower bound for low-rank approximation is given, which essentially matches the bounds of [4]. It is also shown there that an algorithm with this complexity is not possible with uniform sampling.

## References

- [1] D. Achlioptas and F. McSherry, "Fast Computation of Low Rank Approximations" Proceedings of the 33rd Annual Symposium on Theory of Computing, 2001.
- [2] N. Alon, R. A. Duke, H Lefmann, V. Rödl and R. Yuster, "The algorithmic aspects of the Regularity Lemma," Journal of Algorithms 16 (1994) 80-109.
- [3] Z. Bar-Yossef, "Sampling Lower Bounds via Information Theory", Proceedings of the 35th Annual Symposium on Theory of Computing, 2003.
- [4] P. Drineas, A. Frieze, R. Kannan, S. Vempala, and V. Vinay, *Clustering in large graphs and matrices*, Proc. of the Symposium on Discrete Algorithms, 291–299, 1999.
- [5] M. W. Berry, S. T. Dumais, and G. W. O'Brien. "Using linear algebra for intelligent information retrieval", SIAM Review, 37(4), 1995, 573-595, 1995.
- [6] S. Deerwester, S. T. Dumais, T.K. Landauer, G.W. Furnas, and R.A. Harshman. "Indexing by latent semantic analysis," Journal of the Society for Information Science, 41(6), 391-407, 1990.
- [7] S.T. Dumais, G.W. Furnas, T.K. Landauer, and S. Deerwester, "Using latent semantic analysis to improve information retrieval," In Proceedings of CHI'88: Conference on Human Factors in Computing, New York: ACM, 281-285, 1988.
- [8] S.T. Dumais, "Improving the retrieval of information from external sources", Behavior Research Methods, Instruments and Computers, 23(2), 229-236, 1991.
- [9] A.M.Frieze and R. Kannan, "The Regularity Lemma and approximation schemes for dense problems", Proceedings of the 37th Annual IEEE Symposium on Foundations of Computing, (1996) 12-20.
- [10] A.M.Frieze and R. Kannan, "Quick approximations to matrices and applications," Combinatorica, 19 (1999) 175-220.

- [11] A.M.Frieze and R. Kannan, "A simple algorithm for constructing Szemeredi's Regularity Partition", Electronic Journal of Combinatorics, 6(1) (1999) R17. http://www.math.cmu.edu/~aflp/papers.html.
- [12] G. H. Golub and C. F. Van Loan, Matrix Computations, Johns Hopkins University Press, London, 1989.
- [13] J. Kleinberg, "Authoritative sources in a hyperlinked environment," Proc. 9th ACM-SIAM Symposium on Discrete Algorithms, 1998.
- [14] J. Komlós and M. Simonovits, "Szemerédi's Regularity Lemma and its applications in graph theory", to appear.
- [15] C. Papadimitriou, P. Raghavan, H. Tamaki and S. Vempala, Latent Semantic Indexing: A Probabilistic Analysis, JCSS 61, 217–235, 2000.