
A Scalable CUR Matrix Decomposition Algorithm: Lower Time Complexity and Tighter Bound

Shusen Wang and Zhihua Zhang
College of Computer Science & Technology
Zhejiang University
Hangzhou, China 310027
{wss, zhzhang}@zju.edu.cn

Abstract

The CUR matrix decomposition is an important extension of Nyström approximation to a general matrix. It approximates any data matrix in terms of a small number of its columns and rows. In this paper we propose a novel randomized CUR algorithm with an expected relative-error bound. The proposed algorithm has the advantages over the existing relative-error CUR algorithms that it possesses tighter theoretical bound and lower time complexity, and that it can avoid maintaining the whole data matrix in main memory. Finally, experiments on several real-world datasets demonstrate significant improvement over the existing relative-error algorithms.

1 Introduction

Large-scale matrices emerging from stocks, genomes, web documents, web images and videos everyday bring new challenges in modern data analysis. Most efforts have been focused on manipulating, understanding and interpreting large-scale data matrices. In many cases, matrix factorization methods are employed to construct compressed and informative representations to facilitate computation and interpretation. A principled approach is the truncated singular value decomposition (SVD) which finds the best low-rank approximation of a data matrix. Applications of SVD such as eigenface [20, 21] and latent semantic analysis [4] have been illustrated to be very successful.

However, the basis vectors resulting from SVD have little concrete meaning, which makes it very difficult for us to understand and interpret the data in question. An example in [10, 19] has well shown this viewpoint; that is, the vector $[(1/2)\text{age} - (1/\sqrt{2})\text{height} + (1/2)\text{income}]$, the sum of the significant uncorrelated features from a dataset of people's features, is not particularly informative. The authors of [17] have also claimed: "it would be interesting to try to find basis vectors for all experiment vectors, using actual experiment vectors and not artificial bases that offer little insight." Therefore, it is of great interest to represent a data matrix in terms of a small number of actual columns and/or actual rows of the matrix.

The *CUR matrix decomposition* provides such techniques, and it has been shown to be very useful in high dimensional data analysis [19]. Given a matrix \mathbf{A} , the CUR technique selects a subset of columns of \mathbf{A} to construct a matrix \mathbf{C} and a subset of rows of \mathbf{A} to construct a matrix \mathbf{R} , and computes a matrix \mathbf{U} such that $\hat{\mathbf{A}} = \mathbf{CUR}$ best approximates \mathbf{A} . The typical CUR algorithms [7, 8, 10] work in a two-stage manner. Stage 1 is a standard column selection procedure, and Stage 2 does row selection from \mathbf{A} and \mathbf{C} simultaneously. Thus Stage 2 is more complicated than Stage 1.

The CUR matrix decomposition problem is widely studied in the literature [7, 8, 9, 10, 12, 13, 16, 18, 19, 22]. Perhaps the most widely known work on the CUR problem is [10], in which the authors devised a randomized CUR algorithm called the *subspace sampling algorithm*. Particularly, the algorithm has $(1 + \epsilon)$ relative-error ratio with high probability (w.h.p.).

Unfortunately, all the existing CUR algorithms require a large number of columns and rows to be chosen. For example, for an $m \times n$ matrix \mathbf{A} and a target rank $k \leq \min\{m, n\}$, the state-of-the-art CUR algorithm — the subspace sampling algorithm in [10] — requires exactly $\mathcal{O}(k^4 \epsilon^{-6})$ rows or $\mathcal{O}(k \epsilon^{-4} \log^2 k)$ rows in expectation to achieve $(1 + \epsilon)$ relative-error ratio w.h.p. Moreover, the computational cost of this algorithm is at least the cost of the truncated SVD of \mathbf{A} , that is, $\mathcal{O}(\min\{mn^2, nm^2\})$.¹ The algorithms are therefore impractical for large-scale matrices.

In this paper we develop a CUR algorithm which beats the state-of-the-art algorithm in both theory and experiments. In particular, we show in Theorem 5 a novel randomized CUR algorithm with lower time complexity and tighter theoretical bound in comparison with the state-of-the-art CUR algorithm in [10].

The rest of this paper is organized as follows. Section 3 introduces several existing column selection algorithms and the state-of-the-art CUR algorithm. Section 4 describes and analyzes our novel CUR algorithm. Section 5 empirically compares our proposed algorithm with the state-of-the-art algorithm.

2 Notations

For a matrix $\mathbf{A} = [a_{ij}] \in \mathbb{R}^{m \times n}$, let $\mathbf{a}^{(i)}$ be its i -th row and \mathbf{a}_j be its j -th column. Let $\|\mathbf{A}\|_1 = \sum_{i,j} |a_{ij}|$ be the ℓ_1 -norm, $\|\mathbf{A}\|_F = (\sum_{i,j} a_{ij}^2)^{1/2}$ be the Frobenius norm, and $\|\mathbf{A}\|_2$ be the spectral norm. Moreover, let \mathbf{I}_m denote an $m \times m$ identity matrix, and $\mathbf{0}_{mn}$ denotes an $m \times n$ zero matrix.

Let $\mathbf{A} = \mathbf{U}_\mathbf{A} \boldsymbol{\Sigma}_\mathbf{A} \mathbf{V}_\mathbf{A}^T = \sum_{i=0}^{\rho} \sigma_{\mathbf{A},i} \mathbf{u}_{\mathbf{A},i} \mathbf{v}_{\mathbf{A},i}^T = \mathbf{U}_{\mathbf{A},k} \boldsymbol{\Sigma}_{\mathbf{A},k} \mathbf{V}_{\mathbf{A},k}^T + \mathbf{U}_{\mathbf{A},k\perp} \boldsymbol{\Sigma}_{\mathbf{A},k\perp} \mathbf{V}_{\mathbf{A},k\perp}^T$ be the SVD of \mathbf{A} , where $\rho = \text{rank}(\mathbf{A})$, and $\mathbf{U}_{\mathbf{A},k}$, $\boldsymbol{\Sigma}_{\mathbf{A},k}$, and $\mathbf{V}_{\mathbf{A},k}$ correspond to the top k singular values. We denote $\mathbf{A}_k = \mathbf{U}_{\mathbf{A},k} \boldsymbol{\Sigma}_{\mathbf{A},k} \mathbf{V}_{\mathbf{A},k}^T$. Furthermore, let $\mathbf{A}^\dagger = \mathbf{U}_{\mathbf{A},\rho} \boldsymbol{\Sigma}_{\mathbf{A},\rho}^{-1} \mathbf{V}_{\mathbf{A},\rho}^T$ be the Moore-Penrose inverse of \mathbf{A} [1].

3 Related Work

Section 3.1 introduces several relative-error column selection algorithms related to this work. Section 3.2 describes the state-of-the-art CUR algorithm in [10]. Section 3.3 discusses the connection between the column selection problem and the CUR problem.

3.1 Relative-Error Column Selection Algorithms

Given a matrix $\mathbf{A} \in \mathbb{R}^{m \times n}$, column selection is a problem of selecting c columns of \mathbf{A} to construct $\mathbf{C} \in \mathbb{R}^{m \times c}$ to minimize $\|\mathbf{A} - \mathbf{C}\mathbf{C}^\dagger \mathbf{A}\|_F$. Since there are $\binom{n}{c}$ possible choices of constructing \mathbf{C} , so selecting the best subset is a hard problem. In recent years, many polynomial-time approximate algorithms have been proposed, among which we are particularly interested in the algorithms with relative-error bounds; that is, with $c \geq k$ columns selected from \mathbf{A} , there is a constant η such that

$$\|\mathbf{A} - \mathbf{C}\mathbf{C}^\dagger \mathbf{A}\|_F \leq \eta \|\mathbf{A} - \mathbf{A}_k\|_F.$$

We call η the *relative-error ratio*. We now present some recent results related to this work.

We first introduce a recently developed deterministic algorithm called the *dual set sparsification* proposed in [2, 3]. We show their results in Lemma 1. Furthermore, this algorithm is a building block of some more powerful algorithms (e.g., Lemma 2), and our novel CUR algorithm also relies on this algorithm. We attach the algorithm in Appendix A.

Lemma 1 (Column Selection via Dual Set Sparsification Algorithm). *Given a matrix $\mathbf{A} \in \mathbb{R}^{m \times n}$ of rank ρ and a target rank $k (< \rho)$, there exists a deterministic algorithm to select $c (> k)$ columns of \mathbf{A} and form a matrix $\mathbf{C} \in \mathbb{R}^{m \times c}$ such that*

$$\|\mathbf{A} - \mathbf{C}\mathbf{C}^\dagger \mathbf{A}\|_F \leq \sqrt{1 + \frac{1}{(1 - \sqrt{k/c})^2}} \|\mathbf{A} - \mathbf{A}_k\|_F.$$

¹Although some partial SVD algorithms, such as Krylov subspace methods, require only $\mathcal{O}(mnk)$ time, they are all numerical unstable. See [15] for more discussions.

Moreover, the matrix \mathbf{C} can be computed in $T_{\mathbf{V}_{\mathbf{A},k}} + \mathcal{O}(mn + nck^2)$, where $T_{\mathbf{V}_{\mathbf{A},k}}$ is the time needed to compute the top k right singular vectors of \mathbf{A} .

There are also a variety of randomized column selection algorithms achieving relative-error bounds in the literature: [3, 5, 6, 10, 14].

An randomized algorithm in [2] selects only $c = \frac{2k}{\epsilon}(1 + o(1))$ columns to achieve the expected relative-error ratio $(1 + \epsilon)$. The algorithm is based on the approximate SVD via random projection [15], the dual set sparsification algorithm [2], and the adaptive sampling algorithm [6]. Here we present the main results of this algorithm in Lemma 2. Our proposed CUR algorithm is motivated by and relies on this algorithm.

Lemma 2 (Near-Optimal Column Selection Algorithm). *Given a matrix $\mathbf{A} \in \mathbb{R}^{m \times n}$ of rank ρ , a target rank k ($2 \leq k < \rho$), and $0 < \epsilon < 1$, there exists a randomized algorithm to select at most*

$$c = \frac{2k}{\epsilon} (1 + o(1))$$

columns of \mathbf{A} to form a matrix $\mathbf{C} \in \mathbb{R}^{m \times c}$ such that

$$\mathbb{E}^2 \|\mathbf{A} - \mathbf{C}\mathbf{C}^\dagger \mathbf{A}\|_F \leq \mathbb{E} \|\mathbf{A} - \mathbf{C}\mathbf{C}^\dagger \mathbf{A}\|_F^2 \leq (1 + \epsilon) \|\mathbf{A} - \mathbf{A}_k\|_F^2,$$

where the expectations are taken w.r.t. \mathbf{C} . Furthermore, the matrix \mathbf{C} can be computed in $\mathcal{O}((mnk + nk^3)\epsilon^{-2/3})$.

3.2 The Subspace Sampling CUR Algorithm

Drineas *et al.* [10] proposed a two-stage randomized CUR algorithm which has a relative-error bound w.h.p. Given a matrix $\mathbf{A} \in \mathbb{R}^{m \times n}$ and a target rank k , in the first stage the algorithm chooses exactly $c = \mathcal{O}(k^2 \epsilon^{-2} \log \delta^{-1})$ columns (or $c = \mathcal{O}(k \epsilon^{-2} \log k \log \delta^{-1})$ in expectation) of \mathbf{A} to construct $\mathbf{C} \in \mathbb{R}^{m \times c}$; in the second stage it chooses exactly $r = \mathcal{O}(c^2 \epsilon^{-2} \log \delta^{-1})$ rows (or $r = \mathcal{O}(c \epsilon^{-2} \log c \log \delta^{-1})$ in expectation) of \mathbf{A} and \mathbf{C} simultaneously to construct \mathbf{R} and \mathbf{U} . With probability at least $1 - \delta$, the relative-error ratio is $1 + \epsilon$. The computational cost is dominated by the truncated SVD of \mathbf{A} and \mathbf{C} .

Though the algorithm is ϵ -optimal with high probability, it requires too many rows get chosen: at least $r = \mathcal{O}(k \epsilon^{-4} \log^2 k)$ rows in expectation. In this paper we seek to devise an algorithm with mild requirement on column and row numbers.

3.3 Connection between Column Selection and CUR Matrix Decomposition

The CUR problem has a close connection with the column selection problem. As aforementioned, the first stage of existing CUR algorithms is simply a column selection procedure. However, the second stage is more complicated. If the second stage is naïvely solved by a column selection algorithm on \mathbf{A}^T , then the error ratio will be at least $(2 + \epsilon)$.

For a relative-error CUR algorithm, the first stage seeks to bound a construction error ratio of $\frac{\|\mathbf{A} - \mathbf{C}\mathbf{C}^\dagger \mathbf{A}\|_F}{\|\mathbf{A} - \mathbf{A}_k\|_F}$, while the second stage seeks to bound $\frac{\|\mathbf{A} - \mathbf{C}\mathbf{C}^\dagger \mathbf{A}\mathbf{R}^\dagger \mathbf{R}\|_F}{\|\mathbf{A} - \mathbf{C}\mathbf{C}^\dagger \mathbf{A}\|_F}$ given \mathbf{C} . Actually, the first stage is a special case of the second stage where $\mathbf{C} = \mathbf{A}_k$. Given a matrix \mathbf{A} , if an algorithm solving the second stage results in a bound $\frac{\|\mathbf{A} - \mathbf{C}\mathbf{C}^\dagger \mathbf{A}\mathbf{R}^\dagger \mathbf{R}\|_F}{\|\mathbf{A} - \mathbf{C}\mathbf{C}^\dagger \mathbf{A}\|_F} \leq \eta$, then this algorithm also solves the column selection problem for \mathbf{A}^T with an η relative-error ratio. Thus the second stage of CUR is a generalization of the column selection problem.

4 Main Results

In this section we introduce our proposed CUR algorithm. We call it *the fast CUR algorithm* because it has lower time complexity compared with SVD. We describe it in Algorithm 1 and give a theoretical analysis in Theorem 5. Theorem 5 relies on Lemma 2 and Theorem 4, and Theorem 4 relies on Theorem 3. Theorem 3 is a generalization of [6, Theorem 2.1], and Theorem 4 is a generalization of [2, Theorem 5].

Algorithm 1 The Fast CUR Algorithm.

- 1: **Input:** a real matrix $\mathbf{A} \in \mathbb{R}^{m \times n}$, target rank k , $\epsilon \in (0, 1]$, target column number $c = \frac{2k}{\epsilon}(1 + o(1))$, target row number $r = \frac{2c}{\epsilon}(1 + o(1))$;
 - 2: // Stage 1: select c columns of \mathbf{A} to construct $\mathbf{C} \in \mathbb{R}^{m \times c}$
 - 3: Compute approximate truncated SVD via random projection such that $\mathbf{A}_k \approx \tilde{\mathbf{U}}_k \tilde{\Sigma}_k \tilde{\mathbf{V}}_k$;
 - 4: Construct $\mathcal{U}_1 \leftarrow$ columns of $(\mathbf{A} - \tilde{\mathbf{U}}_k \tilde{\Sigma}_k \tilde{\mathbf{V}}_k)$; $\mathcal{V}_1 \leftarrow$ columns of $\tilde{\mathbf{V}}_k^T$;
 - 5: Compute $\mathbf{s}_1 \leftarrow$ Dual Set Spectral-Frobenius Sparsification Algorithm ($\mathcal{U}_1, \mathcal{V}_1, c - 2k/\epsilon$);
 - 6: Construct $\mathbf{C}_1 \leftarrow \mathbf{A} \text{Diag}(\mathbf{s}_1)$, and then delete the all-zero columns;
 - 7: Residual matrix $\mathbf{D} \leftarrow \mathbf{A} - \mathbf{C}_1 \mathbf{C}_1^\dagger \mathbf{A}$;
 - 8: Compute sampling probabilities: $p_i = \|\mathbf{d}_i\|_2^2 / \|\mathbf{D}\|_F^2, i = 1, \dots, n$;
 - 9: Sampling $c_2 = 2k/\epsilon$ columns from \mathbf{A} with probability $\{p_1, \dots, p_n\}$ to construct \mathbf{C}_2 ;
 - 10: // Stage 2: select r rows of \mathbf{A} to construct $\mathbf{R} \in \mathbb{R}^{r \times n}$
 - 11: Construct $\mathcal{U}_2 \leftarrow$ columns of $(\mathbf{A} - \tilde{\mathbf{U}}_k \tilde{\Sigma}_k \tilde{\mathbf{V}}_k)^T$; $\mathcal{V}_2 \leftarrow$ columns of $\tilde{\mathbf{U}}_k^T$;
 - 12: Compute $\mathbf{s}_2 \leftarrow$ Dual Set Spectral-Frobenius Sparsification Algorithm ($\mathcal{U}_2, \mathcal{V}_2, r - 2c/\epsilon$);
 - 13: Construct $\mathbf{R}_1 \leftarrow \text{Diag}(\mathbf{s}_2) \mathbf{A}$, and then delete the all-zero rows;
 - 14: Residual matrix $\mathbf{B} \leftarrow \mathbf{A} - \mathbf{A} \mathbf{R}_1^\dagger \mathbf{R}_1$; Compute $q_j = \|\mathbf{b}^{(j)}\|_2^2 / \|\mathbf{B}\|_F^2, j = 1, \dots, m$;
 - 15: Sampling $r_2 = 2c/\epsilon$ rows from \mathbf{A} with probability $\{q_1, \dots, q_m\}$ to construct \mathbf{R}_2 ;
 - 16: **return** $\mathbf{C} = [\mathbf{C}_1, \mathbf{C}_2], \mathbf{R} = [\mathbf{R}_1^T, \mathbf{R}_2^T]^T$, and $\mathbf{U} = \mathbf{C}^\dagger \mathbf{A} \mathbf{R}^\dagger$.
-

4.1 Adaptive Sampling

The relative-error adaptive sampling algorithm is established in [6, Theorem 2.1]. The algorithm is based on the following idea: after selecting a proportion of columns from \mathbf{A} to form \mathbf{C}_1 by an arbitrary algorithm, the algorithms randomly samples additional c_2 columns according to the residual $\mathbf{A} - \mathbf{C}_1 \mathbf{C}_1^\dagger \mathbf{A}$. Boutsidis *et al.* [2] used the adaptive sampling algorithm to decrease the residual of the dual set sparsification algorithm and obtained an $(1 + \epsilon)$ relative-error bound. Here we prove a new bound for the adaptive sampling algorithm. Interestingly, this new bound is a generalization of the original one in [6, Theorem 2.1]. In other words, Theorem 2.1 of [6] is a direct corollary of our following theorem in which $\mathbf{C} = \mathbf{A}_k$ is set.

Theorem 3 (The Adaptive Sampling Algorithm). *Given a matrix $\mathbf{A} \in \mathbb{R}^{m \times n}$ and a matrix $\mathbf{C} \in \mathbb{R}^{m \times c}$ such that $\text{rank}(\mathbf{C}) = \text{rank}(\mathbf{C} \mathbf{C}^\dagger \mathbf{A}) = \rho$, ($\rho \leq c \leq n$), we let $\mathbf{R}_1 \in \mathbb{R}^{r_1 \times n}$ consist of r_1 rows of \mathbf{A} , and define the residual $\mathbf{B} = \mathbf{A} - \mathbf{A} \mathbf{R}_1^\dagger \mathbf{R}_1$. Additionally, for $i = 1, \dots, m$, we define*

$$p_i = \|\mathbf{b}^{(i)}\|_2^2 / \|\mathbf{B}\|_F^2.$$

We further sample r_2 rows i.i.d. from \mathbf{A} , in each trial of which the i -th row is chosen with probability p_i . Let $\mathbf{R}_2 \in \mathbb{R}^{r_2 \times n}$ contains the r_2 sampled rows and let $\mathbf{R} = [\mathbf{R}_1^T, \mathbf{R}_2^T]^T \in \mathbb{R}^{(r_1+r_2) \times n}$. Then the following inequality holds:

$$\mathbb{E} \|\mathbf{A} - \mathbf{C} \mathbf{C}^\dagger \mathbf{A} \mathbf{R}^\dagger \mathbf{R}\|_F^2 \leq \|\mathbf{A} - \mathbf{C} \mathbf{C}^\dagger \mathbf{A}\|_F^2 + \frac{\rho}{r_2} \|\mathbf{A} - \mathbf{A} \mathbf{R}_1^\dagger \mathbf{R}_1\|_F^2,$$

where the expectation is taken w.r.t. \mathbf{R}_2 .

4.2 The Fast CUR Algorithm

Based on the dual set sparsification algorithm of Lemma 1 and the adaptive sampling algorithm of Theorem 3, we develop a randomized algorithm to solve the second stage of CUR problem. We present the results of the algorithm in Theorem 4. Theorem 5 of [2] is a special case of the following theorem where $\mathbf{C} = \mathbf{A}_k$.

Theorem 4 (The Fast Row Selection Algorithm). *Given a matrix $\mathbf{A} \in \mathbb{R}^{m \times n}$ and a matrix $\mathbf{C} \in \mathbb{R}^{m \times c}$ such that $\text{rank}(\mathbf{C}) = \text{rank}(\mathbf{C} \mathbf{C}^\dagger \mathbf{A}) = \rho$, ($\rho \leq c \leq n$), and a target rank k ($\leq \rho$), the proposed randomized algorithm selects $r = \frac{2\rho}{\epsilon}(1 + o(1))$ rows of \mathbf{A} to construct $\mathbf{R} \in \mathbb{R}^{r \times n}$, such that*

$$\mathbb{E} \|\mathbf{A} - \mathbf{C} \mathbf{C}^\dagger \mathbf{A} \mathbf{R}^\dagger \mathbf{R}\|_F^2 \leq \|\mathbf{A} - \mathbf{C} \mathbf{C}^\dagger \mathbf{A}\|_F^2 + \epsilon \|\mathbf{A} - \mathbf{A}_k\|_F^2,$$

where the expectation is taken w.r.t. \mathbf{R} . Furthermore, the matrix \mathbf{R} can be computed in $\mathcal{O}((mnk + mk^3)\epsilon^{-2/3})$ time.

Based on Lemma 2 and Theorem 4, here we present the main theorem for the fast CUR algorithm.

Table 1: A summary of the datasets.

Dataset	Type	size	Source
Redrock	natural image	18000 × 4000	http://www.agarwala.org/efficient_gdc/
Arcene	biology	10000 × 900	http://archive.ics.uci.edu/ml/datasets/Arcene
Dexter	bag of words	20000 × 2600	http://archive.ics.uci.edu/ml/datasets/Dexter

Theorem 5 (The Fast CUR Algorithm). *Given a matrix $\mathbf{A} \in \mathbb{R}^{m \times n}$ and a positive integer $k \ll \min\{m, n\}$, the fast CUR algorithm (described in Algorithm 1) randomly selects $c = \frac{2k}{\epsilon}(1 + o(1))$ columns of \mathbf{A} to construct $\mathbf{C} \in \mathbb{R}^{m \times c}$ with the near-optimal column selection algorithm of Lemma 2, and then selects $r = \frac{2c}{\epsilon}(1 + o(1))$ rows of \mathbf{A} to construct $\mathbf{R} \in \mathbb{R}^{r \times n}$ with the fast row selection algorithm of Theorem 4. Then we have*

$$\mathbb{E}\|\mathbf{A} - \mathbf{CUR}\|_F = \mathbb{E}\|\mathbf{A} - \mathbf{C}(\mathbf{C}^\dagger \mathbf{A} \mathbf{R}^\dagger)\mathbf{R}\|_F \leq (1 + \epsilon)\|\mathbf{A} - \mathbf{A}_k\|_F.$$

Moreover, the algorithm runs in time $\mathcal{O}\left(mnk\epsilon^{-2/3} + (m + n)k^3\epsilon^{-2/3} + mk^2\epsilon^{-2} + nk^2\epsilon^{-4}\right)$.

Since $k, c, r \ll \min\{m, n\}$ by the assumptions, so the time complexity of the fast CUR algorithm is lower than that of the SVD of \mathbf{A} . This is the main reason why we call it the fast CUR algorithm.

Another advantage of this algorithm is avoiding loading the whole $m \times n$ data matrix \mathbf{A} into main memory. None of three steps — the randomized SVD, the dual set sparsification algorithm, and the adaptive sampling algorithm — requires loading the whole of \mathbf{A} into memory. The most memory-expensive operation throughout the fast CUR Algorithm is computing the Moore-Penrose inverse of \mathbf{C} and \mathbf{R} , which requires maintaining an $m \times c$ matrix or an $r \times n$ matrix in memory. In comparison, the subspace sampling algorithm requires loading the whole matrix into memory to compute its truncated SVD.

5 Empirical Comparisons

In this section we provide empirical comparisons among the relative-error CUR algorithms on several datasets. We report the relative-error ratio and the running time of each algorithm on each data set. The relative-error ratio is defined by

$$\text{Relative-error ratio} = \frac{\|\mathbf{A} - \mathbf{CUR}\|_F}{\|\mathbf{A} - \mathbf{A}_k\|_F},$$

where k is a specified target rank.

We conduct experiments on three datasets, including natural image, biology data, and bags of words. Table 1 briefly summarizes some information of the datasets. Redrock is a large size natural image. Arcene and Dexter are both from the UCI datasets [11]. Arcene is a biology dataset with 900 instances and 10000 attributes. Dexter is a bag of words dataset with a 20000-vocabulary and 2600 documents. Each dataset is actually represented as a data matrix, upon which we apply the CUR algorithms.

We implement all the algorithms in MATLAB 7.10.0. We conduct experiments on a workstation with 12 Intel Xeon 3.47GHz CPUs, 12GB memory, and Ubuntu 10.04 system. According to the analysis in [10] and this paper, k , c , and r should be integers far less than m and n . For each data set and each algorithm, we set $k = 10, 20$, or 50 , and $c = \alpha k$, $r = \alpha c$, where α ranges in each set of experiments. We repeat each set of experiments for 20 times and report the average and the standard deviation of the error ratios. The results are depicted in Figures 1, 2, 3.

The results show that the fast CUR algorithm has much lower relative-error ratio than the subspace sampling algorithm. The experimental results well match our theoretical analyses in Section 4. As for the running time, the fast CUR algorithm is more efficient when c and r are small. When c and r become large, the fast CUR algorithm becomes less efficient. This is because the time complexity of the fast CUR algorithm is linear in ϵ^{-4} and large c and r imply small ϵ . However, the purpose of CUR is to select a small number of columns and rows from the data matrix, that is, $c \ll n$ and $r \ll m$. So we are not interested in the cases where c and r are large compared with n and m , say $k = 20$ and $\alpha = 10$.

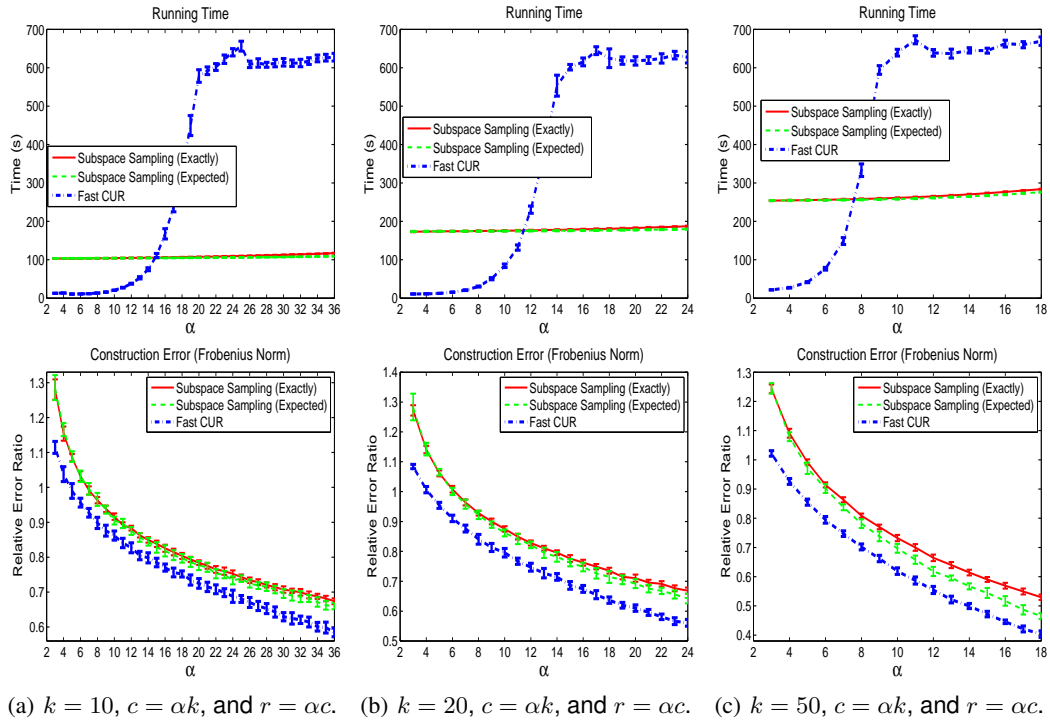


Figure 1: Empirical results on the Redrock data set.

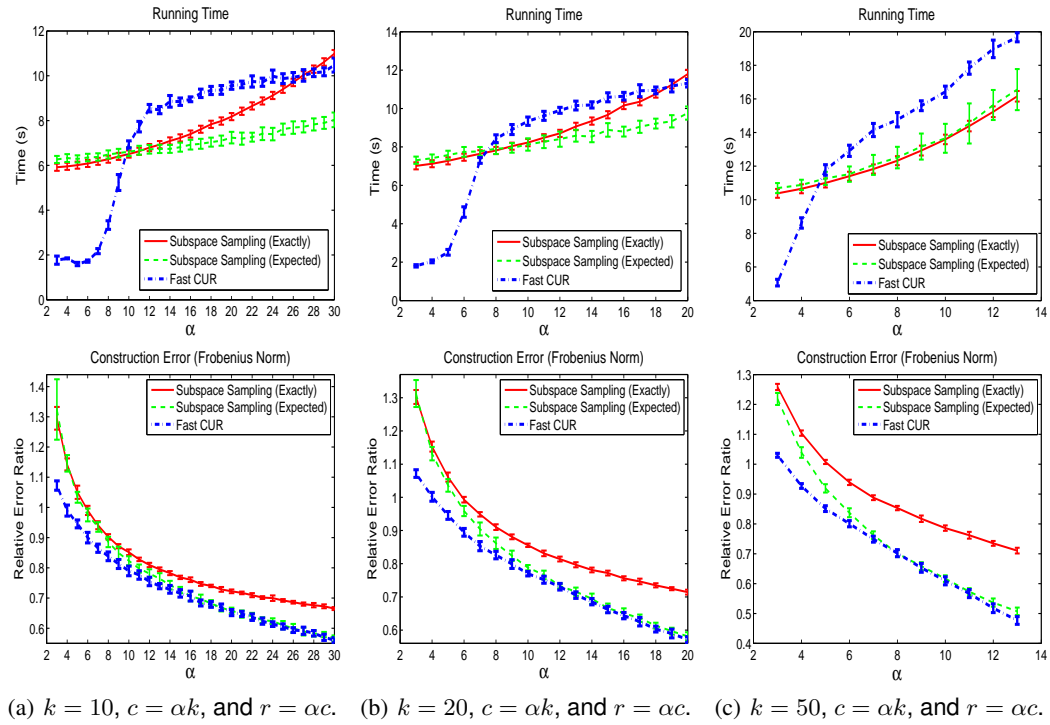
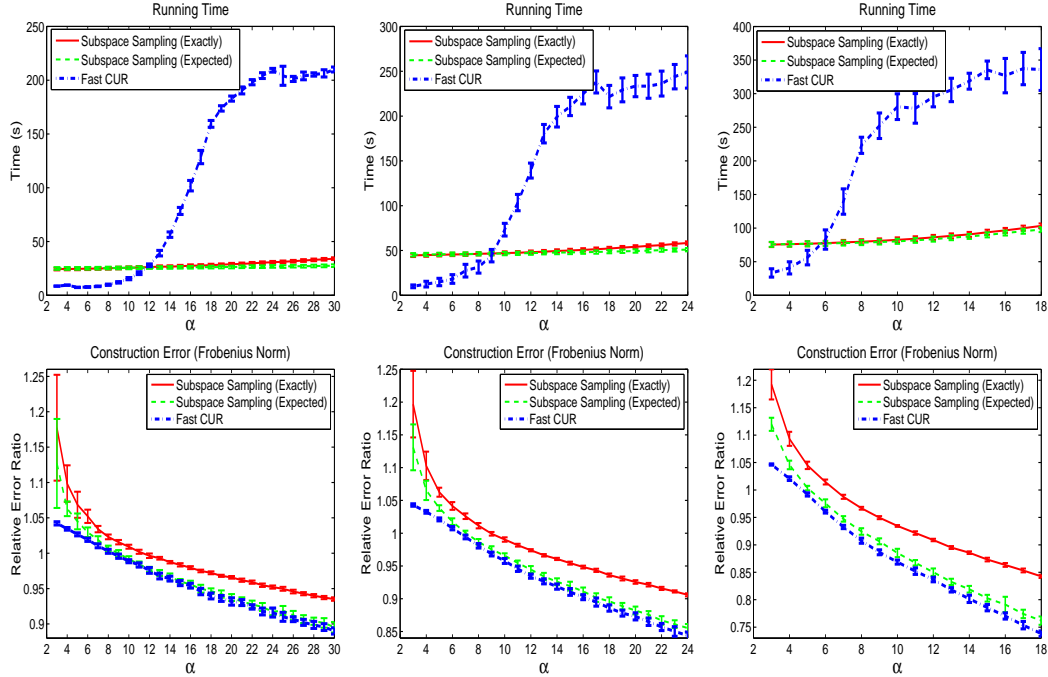


Figure 2: Empirical results on the Arcene data set.



(a) $k = 10$, $c = \alpha k$, and $r = \alpha c$. (b) $k = 20$, $c = \alpha k$, and $r = \alpha c$. (c) $k = 50$, $c = \alpha k$, and $r = \alpha c$.

Figure 3: Empirical results on the Dexter data set.

6 Conclusions

In this paper we have proposed a novel randomized algorithm for the CUR matrix decomposition problem. This algorithm is faster, more scalable, and more accurate than the state-of-the-art algorithm, i.e., the subspace sampling algorithm. Our algorithm requires only $c = 2k\epsilon^{-1}(1 + o(1))$ columns and $r = 2c\epsilon^{-1}(1 + o(1))$ rows to achieve $(1+\epsilon)$ relative-error ratio. To achieve the same relative-error bound, the subspace sampling algorithm requires $c = \mathcal{O}(k\epsilon^{-2} \log k)$ columns and $r = \mathcal{O}(c\epsilon^{-2} \log c)$ rows selected from the original matrix. Our algorithm also beats the subspace sampling algorithm in time-complexity. Our algorithm costs $\mathcal{O}(mnk\epsilon^{-2/3} + (m+n)k^3\epsilon^{-2/3} + mk^2\epsilon^{-2} + nk^2\epsilon^{-4})$ time, which is lower than $\mathcal{O}(\min\{mn^2, m^2n\})$ of the subspace sampling algorithm when k is small. Moreover, our algorithm enjoys another advantage of avoiding loading the whole data matrix into main memory, which also makes our algorithm more scalable. Finally, the empirical comparisons have also demonstrated the effectiveness and efficiency of our algorithm.

A The Dual Set Sparsification Algorithm

For the sake of completeness, we attach the dual set sparsification algorithm here and describe some implementation details. The dual set sparsification algorithms are deterministic algorithms established in [2]. The fast CUR algorithm calls the *dual set spectral-Frobenius sparsification algorithm* [2, Lemma 13] in both stages. We show this algorithm in Algorithm 2 and its bounds in Lemma 6.

Lemma 6 (Dual Set Spectral-Frobenius Sparsification). *Let $\mathcal{U} = \{\mathbf{x}_1, \dots, \mathbf{x}_n\} \subset \mathbb{R}^l$, ($l < n$), contains the columns of an arbitrary matrix $\mathbf{X} \in \mathbb{R}^{l \times n}$. Let $\mathcal{V} = \{\mathbf{v}_1, \dots, \mathbf{v}_n\} \subset \mathbb{R}^k$, ($k < n$), be a decompositions of the identity, i.e. $\sum_{i=1}^n \mathbf{v}_i \mathbf{v}_i^T = \mathbf{I}_k$. Given an integer r with $k < r < n$, Algorithm 2 deterministically computes a set of weights $s_i \geq 0$ ($i = 1, \dots, n$) at most r of which are non-zero, such that*

$$\lambda_k \left(\sum_{i=1}^n s_i \mathbf{v}_i \mathbf{v}_i^T \right) \geq \left(1 - \sqrt{\frac{k}{r}} \right)^2 \quad \text{and} \quad \text{tr} \left(\sum_{i=1}^n s_i \mathbf{x}_i \mathbf{x}_i^T \right) \leq \|\mathbf{X}\|_F^2.$$

Algorithm 2 Deterministic Dual Set Spectral-Frobenius Sparsification Algorithm.

- 1: **Input:** $U = \{\mathbf{x}_i\}_{i=1}^n \subset \mathbb{R}^l$, ($l < n$); $\mathcal{V} = \{\mathbf{v}_i\}_{i=1}^n \subset \mathbb{R}^k$, with $\sum_{i=1}^n \mathbf{v}_i \mathbf{v}_i^T = \mathbf{I}_k$ ($k < n$); $k < r < n$;
- 2: **Initialize:** $\mathbf{s}_0 = \mathbf{0}_{m \times 1}$, $\mathbf{A}_0 = \mathbf{0}_{k \times k}$;
- 3: Compute $\|\mathbf{x}_i\|_2^2$ for $i = 1, \dots, n$, and then compute $\delta_U = \frac{\sum_{i=1}^n \|\mathbf{x}_i\|_2^2}{1 - \sqrt{k/r}}$;
- 4: **for** $\tau = 0$ to $r - 1$ **do**
- 5: Compute the eigenvalue decomposition of \mathbf{A}_τ ;
- 6: Find an index j in $\{1, \dots, n\}$ and compute a weight $t > 0$ such that

$$\delta_U^{-1} \|\mathbf{x}_j\|_2^2 \leq t^{-1} \leq \frac{\mathbf{v}_j^T (\mathbf{A}_\tau - (L_\tau + 1)\mathbf{I}_k)^{-2} \mathbf{v}_j}{\phi(L_\tau + 1, \mathbf{A}_\tau) - \phi(L_\tau, \mathbf{A}_\tau)} - \mathbf{v}_j^T (\mathbf{A}_\tau - (L_\tau + 1)\mathbf{I}_k)^{-1} \mathbf{v}_j;$$

where

$$\phi(L, \mathbf{A}) = \sum_{i=1}^k (\lambda_i(\mathbf{A}) - L)^{-1}, \quad L_\tau = \tau - \sqrt{rk};$$

- 7: Update the j -th component of \mathbf{s}_τ and \mathbf{A}_τ : $\mathbf{s}_{\tau+1}[j] = \mathbf{s}_\tau[j] + t$, $\mathbf{A}_{\tau+1} = \mathbf{A}_\tau + t\mathbf{v}_j \mathbf{v}_j^T$;
 - 8: **end for**
 - 9: **return** $\mathbf{s} = \frac{1 - \sqrt{k/r}}{r} \mathbf{s}_r$.
-

The weights s_i can be computed deterministically in $\mathcal{O}(rnk^2 + nl)$ time.

Here we would like to mention the implementation of Algorithm 2, which is not described in detailed by [2]. In each iteration the algorithm performs once eigenvalue decomposition: $\mathbf{A}_\tau = \mathbf{W}\mathbf{\Lambda}\mathbf{W}^T$. (\mathbf{A}_τ is guaranteed to be positive semi-definite in each iteration). Since

$$(\mathbf{A}_\tau - \alpha\mathbf{I}_k)^q = \mathbf{W}\text{Diag}\left((\lambda_1 - \alpha)^q, \dots, (\lambda_k - \alpha)^q\right)\mathbf{W}^T,$$

we can efficiently compute $(\mathbf{A}_\tau - (L_\tau + 1)\mathbf{I}_k)^q$ based on the eigenvalue decomposition of \mathbf{A}_τ . With the eigenvalues at hand, $\phi(L, \mathbf{A}_\tau)$ can also be computed directly.

Acknowledgments

This work has been supported in part by the Natural Science Foundations of China (No. 61070239), the Google visiting faculty program, and the Scholarship Award for Excellent Doctoral Student granted by Ministry of Education.

References

- [1] Adi Ben-Israel and Thomas N.E. Greville. *Generalized Inverses: Theory and Applications. Second Edition*. Springer, 2003.
- [2] Christos Boutsidis, Petros Drineas, and Malik Magdon-Ismail. Near-optimal column-based matrix reconstruction. *CoRR*, abs/1103.0995, 2011.
- [3] Christos Boutsidis, Petros Drineas, and Malik Magdon-Ismail. Near optimal column-based matrix reconstruction. In *Proceedings of the 2011 IEEE 52nd Annual Symposium on Foundations of Computer Science, FOCS '11*, pages 305–314, 2011.
- [4] Scott Deerwester, Susan T. Dumais, George W. Furnas, Thomas K. Landauer, and Richard Harshman. Indexing by latent semantic analysis. *Journal of The American Society for Information Science*, 41(6):391–407, 1990.
- [5] Amit Deshpande and Luis Rademacher. Efficient volume sampling for row/column subset selection. In *Proceedings of the 2010 IEEE 51st Annual Symposium on Foundations of Computer Science, FOCS '10*, pages 329–338, 2010.
- [6] Amit Deshpande, Luis Rademacher, Santosh Vempala, and Grant Wang. Matrix approximation and projective clustering via volume sampling. *Theory of Computing*, 2(2006):225–247, 2006.
- [7] Petros Drineas. Pass-efficient algorithms for approximating large matrices. In *In Proceeding of the 14th Annual ACM-SIAM Symposium on Discrete Algorithms*, pages 223–232, 2003.

- [8] Petros Drineas, Ravi Kannan, and Michael W. Mahoney. Fast monte carlo algorithms for matrices iii: Computing a compressed approximate matrix decomposition. *SIAM Journal on Computing*, 36(1):184–206, 2006.
- [9] Petros Drineas and Michael W. Mahoney. On the Nyström method for approximating a gram matrix for improved kernel-based learning. *Journal of Machine Learning Research*, 6:2153–2175, 2005.
- [10] Petros Drineas, Michael W. Mahoney, and S. Muthukrishnan. Relative-error CUR matrix decompositions. *SIAM Journal on Matrix Analysis and Applications*, 30(2):844–881, September 2008.
- [11] A. Frank and A. Asuncion. UCI machine learning repository, 2010.
- [12] S. A. Goreinov, E. E. Tyrtyshnikov, and N. L. Zamarashkin. A theory of pseudoskeleton approximations. *Linear Algebra and Its Applications*, 261:1–21, 1997.
- [13] S. A. Goreinov, N. L. Zamarashkin, and E. E. Tyrtyshnikov. Pseudo-skeleton approximations by matrices of maximal volume. *Mathematical Notes*, 62(4):619–623, 1997.
- [14] Venkatesan Guruswami and Ali Kemal Sinop. Optimal column-based low-rank matrix reconstruction. In *Proceedings of the Twenty-Third Annual ACM-SIAM Symposium on Discrete Algorithms*, SODA '12, pages 1207–1214. SIAM, 2012.
- [15] Nathan Halko, Per-Gunnar Martinsson, and Joel A. Tropp. Finding structure with randomness: Probabilistic algorithms for constructing approximate matrix decompositions. *SIAM Review*, 53(2):217–288, 2011.
- [16] John Hopcroft and Ravi Kannan. *Computer Science Theory for the Information Age*. 2012.
- [17] Finny G. Kuruvilla, Peter J. Park, and Stuart L. Schreiber. Vector algebra in the analysis of genome-wide expression data. *Genome Biology*, 3:research0011–research0011.1, 2002.
- [18] Lester Mackey, Ameet Talwalkar, and Michael I. Jordan. Divide-and-conquer matrix factorization. In *Advances in Neural Information Processing Systems 24*. 2011.
- [19] Michael W. Mahoney and Petros Drineas. CUR matrix decompositions for improved data analysis. *Proceedings of the National Academy of Sciences*, 106(3):697–702, 2009.
- [20] L. Sirovich and M. Kirby. Low-dimensional procedure for the characterization of human faces. *Journal of the Optical Society of America A*, 4(3):519–524, Mar 1987.
- [21] Matthew Turk and Alex Pentland. Eigenfaces for recognition. *Journal of Cognitive Neuroscience*, 3(1):71–86, 1991.
- [22] Eugene E. Tyrtyshnikov. Incomplete cross approximation in the mosaic-skeleton method. *Computing*, 64:367–380, 2000.

B Proofs

B.1 Proof of Theorem 3

Theorem 3 can be equivalently expressed in Theorem 7. In order to stick to the column space convention of [2], we prove Theorem 7 instead of Theorem 3.

Theorem 7 (Adaptive Sampling Algorithm). *Given a matrix $\mathbf{A} \in \mathbb{R}^{m \times n}$ and a matrix $\mathbf{R} \in \mathbb{R}^{r \times n}$ such that $\text{rank}(\mathbf{R}) = \text{rank}(\mathbf{A}\mathbf{R}^\dagger\mathbf{R}) = \rho$, ($\rho \leq r \leq m$), let $\mathbf{C}_1 \in \mathbb{R}^{m \times c_1}$ consist of c_1 columns of \mathbf{A} , and define the residual $\mathbf{B} = \mathbf{A} - \mathbf{C}_1\mathbf{C}_1^\dagger\mathbf{A}$. For $i = 1, \dots, n$, let*

$$p_i = \|\mathbf{b}_i\|_2^2 / \|\mathbf{B}\|_F^2,$$

where \mathbf{b}_i is the i -th column of the matrix \mathbf{B} . Sample a further c_2 columns from \mathbf{A} in c_2 i.i.d. trials, where in each trial the i -th column is chosen with probability p_i . Let $\mathbf{C}_2 \in \mathbb{R}^{m \times c_2}$ contains the c_2 sampled columns and let $\mathbf{C} = [\mathbf{C}_1, \mathbf{C}_2] \in \mathbb{R}^{m \times (c_1+c_2)}$ contain the columns of both \mathbf{C}_1 and \mathbf{C}_2 , all of which are columns of \mathbf{A} . Then the following inequality holds:

$$\mathbb{E}\|\mathbf{A} - \mathbf{C}\mathbf{C}^\dagger\mathbf{A}\mathbf{R}^\dagger\mathbf{R}\|_F^2 \leq \|\mathbf{A} - \mathbf{A}\mathbf{R}^\dagger\mathbf{R}\|_F^2 + \frac{\rho}{c_2} \|\mathbf{A} - \mathbf{C}_1\mathbf{C}_1^\dagger\mathbf{A}\|_F^2.$$

where the expectation is taken w.r.t. \mathbf{C}_2 .

Proof. With a little abuse of symbols, we use bold uppercase letters to denote matrix random variables and bold lowercase to denote vector random variables, without distinguishing between matrix/vector random variables and constant matrices/vectors.

We denote the j -th column of $\mathbf{V}_{\mathbf{A}\mathbf{R}^\dagger\mathbf{R}, \rho} \in \mathbb{R}^{n \times \rho}$ as \mathbf{v}_j , and the (i, j) -th entry of $\mathbf{V}_{\mathbf{A}\mathbf{R}^\dagger\mathbf{R}, \rho}$ as v_{ij} . Define vector random variables $\mathbf{x}_{j,(l)} \in \mathbb{R}^m$ such that for $j = 1, \dots, n$ and $l = 1, \dots, c_2$,

$$\mathbf{x}_{j,(l)} = \frac{v_{ij}}{p_i} \mathbf{b}_i = \frac{v_{ij}}{p_i} (\mathbf{a}_i - \mathbf{C}_1\mathbf{C}_1^\dagger\mathbf{a}_i) \quad \text{with probability } p_i, \quad \text{for } i = 1, \dots, n,$$

Note that $\mathbf{x}_{j,(l)}$ is a linear function of a column of \mathbf{A} sampled from the above defined distribution. We have that

$$\begin{aligned} \mathbb{E}[\mathbf{x}_{j,(l)}] &= \sum_{i=1}^n p_i \frac{v_{ij}}{p_i} \mathbf{b}_i = \mathbf{B}\mathbf{v}_j, \\ \mathbb{E}\|\mathbf{x}_{j,(l)}\|_2^2 &= \sum_{i=1}^n p_i \frac{v_{ij}^2}{p_i^2} \|\mathbf{b}_i\|_2^2 = \sum_{i=1}^n \frac{v_{ij}^2}{\|\mathbf{b}_i\|_2^2 / \|\mathbf{B}\|_F^2} \|\mathbf{b}_i\|_2^2 = \|\mathbf{B}\|_F^2. \end{aligned}$$

Then we let $\mathbf{x}_j = \frac{1}{c_2} \sum_{l=1}^{c_2} \mathbf{x}_{j,(l)}$, we have

$$\begin{aligned} \mathbb{E}[\mathbf{x}_j] &= \mathbb{E}[\mathbf{x}_{j,(l)}] = \mathbf{B}\mathbf{v}_j, \\ \mathbb{E}\|\mathbf{x}_j - \mathbf{B}\mathbf{v}_j\|_2^2 &= \mathbb{E}\left\| \mathbf{x}_j - \mathbb{E}[\mathbf{x}_j] \right\|_2^2 = \frac{1}{c_2} \mathbb{E}\left\| \mathbf{x}_{j,(l)} - \mathbb{E}[\mathbf{x}_{j,(l)}] \right\|_2^2 = \frac{1}{c_2} \mathbb{E}\|\mathbf{x}_{j,(l)} - \mathbf{B}\mathbf{v}_j\|_2^2. \end{aligned}$$

According to the construction of $\mathbf{x}_1, \dots, \mathbf{x}_\rho$, we define the c_2 columns of \mathbf{A} to be $\mathbf{C}_2 \in \mathbb{R}^{m \times c_2}$. Note that all the random variables $\mathbf{x}_1, \dots, \mathbf{x}_\rho$ lie in the subspace $\text{span}(\mathbf{C}_1) + \text{span}(\mathbf{C}_2)$. We define random variables

$$\mathbf{w}_j = \mathbf{C}_1\mathbf{C}_1^\dagger\mathbf{A}\mathbf{R}^\dagger\mathbf{R}\mathbf{v}_j + \mathbf{x}_j = \mathbf{C}_1\mathbf{C}_1^\dagger\mathbf{A}\mathbf{v}_j + \mathbf{x}_j, \quad \text{for } j = 1, \dots, \rho,$$

where the second equality follows from Lemma 8 that $\mathbf{A}\mathbf{R}^\dagger\mathbf{R}\mathbf{v}_j = \mathbf{A}\mathbf{v}_j$ if \mathbf{v}_j is one of the top ρ right singular vectors of $\mathbf{A}\mathbf{R}^\dagger\mathbf{R}$. Then we have that any set of random variables $\{\mathbf{w}_1, \dots, \mathbf{w}_\rho\}$ lies in $\text{span}(\mathbf{C}) = \text{span}(\mathbf{C}_1) + \text{span}(\mathbf{C}_2)$. Let $\mathbf{W} = [\mathbf{w}_1, \dots, \mathbf{w}_\rho]$ be a matrix random variable, we have that $\text{span}(\mathbf{W}) \subset \text{span}(\mathbf{C})$. The expectation of \mathbf{w}_j is

$$\mathbb{E}[\mathbf{w}_j] = \mathbf{C}_1\mathbf{C}_1^\dagger\mathbf{A}\mathbf{v}_j + \mathbb{E}[\mathbf{x}_j] = \mathbf{C}_1\mathbf{C}_1^\dagger\mathbf{A}\mathbf{v}_j + \mathbf{B}\mathbf{v}_j = \mathbf{A}\mathbf{v}_j,$$

therefore we have that

$$\mathbf{w}_j - \mathbf{A}\mathbf{v}_j = \mathbf{x}_j - \mathbf{B}\mathbf{v}_j.$$

The expectation of $\|\mathbf{w}_j - \mathbf{A}\mathbf{v}_j\|_2^2$ is

$$\begin{aligned}
\mathbb{E}\|\mathbf{w}_j - \mathbf{A}\mathbf{v}_j\|_2^2 &= \mathbb{E}\|\mathbf{x}_j - \mathbf{B}\mathbf{v}_j\|_2^2 = \frac{1}{c_2}\mathbb{E}\|\mathbf{x}_{j,(l)} - \mathbf{B}\mathbf{v}_j\|_2^2 \\
&= \frac{1}{c_2}\mathbb{E}\|\mathbf{x}_{j,(l)}\|_2^2 - \frac{2}{c_2}(\mathbf{B}\mathbf{v}_j)^T\mathbb{E}[\mathbf{x}_{j,(l)}] + \frac{1}{c_2}\|\mathbf{B}\mathbf{v}_j\|_2^2 \\
&= \frac{1}{c_2}\mathbb{E}\|\mathbf{x}_{j,(l)}\|_2^2 - \frac{1}{c_2}\|\mathbf{B}\mathbf{v}_j\|_2^2 = \frac{1}{c_2}\|\mathbf{B}\|_F^2 - \frac{1}{c_2}\|\mathbf{B}\mathbf{v}_j\|_2^2 \\
&\leq \frac{1}{c_2}\|\mathbf{B}\|_F^2
\end{aligned} \tag{1}$$

To complete the proof, we let the matrix variable

$$\mathbf{F} = \left(\sum_{q=1}^{\rho} \sigma_q^{-1} \mathbf{w}_q \mathbf{u}_q^T\right) \mathbf{A} \mathbf{R}^\dagger \mathbf{R},$$

where σ_q is the q -th largest singular value of $\mathbf{A} \mathbf{R}^\dagger \mathbf{R}$ and \mathbf{u}_q is the corresponding left singular vector of $\mathbf{A} \mathbf{R}^\dagger \mathbf{R}$. The column space of \mathbf{F} is contained in $\text{span}(\mathbf{W}) \subset \text{span}(\mathbf{C})$, and thus

$$\|\mathbf{A} \mathbf{R}^\dagger \mathbf{R} - \mathbf{C} \mathbf{C}^\dagger \mathbf{A} \mathbf{R}^\dagger \mathbf{R}\|_F^2 \leq \|\mathbf{A} \mathbf{R}^\dagger \mathbf{R} - \mathbf{W} \mathbf{W}^\dagger \mathbf{A} \mathbf{R}^\dagger \mathbf{R}\|_F^2 \leq \|\mathbf{A} \mathbf{R}^\dagger \mathbf{R} - \mathbf{F}\|_F^2.$$

We use \mathbf{F} to bound the error $\|\mathbf{A} \mathbf{R}^\dagger \mathbf{R} - \mathbf{C} \mathbf{C}^\dagger \mathbf{A} \mathbf{R}^\dagger \mathbf{R}\|_F^2$:

$$\begin{aligned}
\mathbb{E}\|\mathbf{A} - \mathbf{C} \mathbf{C}^\dagger \mathbf{A} \mathbf{R}^\dagger \mathbf{R}\|_F^2 &= \mathbb{E}\|\mathbf{A} - \mathbf{A} \mathbf{R}^\dagger \mathbf{R} + \mathbf{A} \mathbf{R}^\dagger \mathbf{R} - \mathbf{C} \mathbf{C}^\dagger \mathbf{A} \mathbf{R}^\dagger \mathbf{R}\|_F^2 \\
&= \mathbb{E}\left[\|\mathbf{A} - \mathbf{A} \mathbf{R}^\dagger \mathbf{R}\|_F^2 + \|\mathbf{A} \mathbf{R}^\dagger \mathbf{R} - \mathbf{C} \mathbf{C}^\dagger \mathbf{A} \mathbf{R}^\dagger \mathbf{R}\|_F^2\right] \\
&\leq \|\mathbf{A} - \mathbf{A} \mathbf{R}^\dagger \mathbf{R}\|_F^2 + \mathbb{E}\|\mathbf{A} \mathbf{R}^\dagger \mathbf{R} - \mathbf{F}\|_F^2,
\end{aligned} \tag{2}$$

where (2) follows from that $\mathbf{A}(\mathbf{I} - \mathbf{R}^\dagger \mathbf{R})$ is orthogonal to $(\mathbf{I} - \mathbf{C} \mathbf{C}^\dagger) \mathbf{A} \mathbf{R}^\dagger \mathbf{R}$. Since $\mathbf{A} \mathbf{R}^\dagger \mathbf{R}$ and \mathbf{F} both lies on the space spanned by the right singular vectors of $\mathbf{A} \mathbf{R}^\dagger \mathbf{R}$, i.e. $\{\mathbf{v}_j\}_{j=1}^{\rho}$, we decompose $\mathbf{A} \mathbf{R}^\dagger \mathbf{R} - \mathbf{F}$ along $\{\mathbf{v}_j\}_{j=1}^{\rho}$:

$$\begin{aligned}
\mathbb{E}\|\mathbf{A} - \mathbf{C} \mathbf{C}^\dagger \mathbf{A} \mathbf{R}^\dagger \mathbf{R}\|_F^2 &\leq \|\mathbf{A} - \mathbf{A} \mathbf{R}^\dagger \mathbf{R}\|_F^2 + \mathbb{E}\|\mathbf{A} \mathbf{R}^\dagger \mathbf{R} - \mathbf{F}\|_F^2, \\
&= \|\mathbf{A} - \mathbf{A} \mathbf{R}^\dagger \mathbf{R}\|_F^2 + \sum_{j=1}^{\rho} \mathbb{E}\left\|\left(\mathbf{A} \mathbf{R}^\dagger \mathbf{R} - \mathbf{F}\right) \mathbf{v}_j\right\|_2^2 \\
&= \|\mathbf{A} - \mathbf{A} \mathbf{R}^\dagger \mathbf{R}\|_F^2 + \sum_{j=1}^{\rho} \mathbb{E}\left\|\mathbf{A} \mathbf{R}^\dagger \mathbf{R} \mathbf{v}_j - \left(\sum_{q=1}^{\rho} \sigma_q^{-1} \mathbf{w}_q \mathbf{u}_q^T\right) \sigma_j \mathbf{u}_j\right\|_2^2 \\
&= \|\mathbf{A} - \mathbf{A} \mathbf{R}^\dagger \mathbf{R}\|_F^2 + \sum_{j=1}^{\rho} \mathbb{E}\left\|\mathbf{A} \mathbf{R}^\dagger \mathbf{R} \mathbf{v}_j - \mathbf{w}_j\right\|_2^2 \\
&= \|\mathbf{A} - \mathbf{A} \mathbf{R}^\dagger \mathbf{R}\|_F^2 + \sum_{j=1}^{\rho} \mathbb{E}\|\mathbf{A} \mathbf{v}_j - \mathbf{w}_j\|_2^2
\end{aligned} \tag{3}$$

$$\leq \|\mathbf{A} - \mathbf{A} \mathbf{R}^\dagger \mathbf{R}\|_F^2 + \frac{\rho}{c_2} \|\mathbf{B}\|_F^2, \tag{4}$$

where (3) follows from Lemma 8 and (4) follows from (1). \square

Lemma 8. Given a matrix $\mathbf{A} \in \mathbb{R}^{m \times n}$ and a matrix $\mathbf{R} \in \mathbb{R}^{r \times n}$ such that $\text{rank}(\mathbf{A} \mathbf{R}^\dagger \mathbf{R}) = \text{rank}(\mathbf{R}) = \rho$, ($\rho \leq r \leq m$). Let $\mathbf{v}_j \in \mathbb{R}^n$ be the j -th top right singular vector of $\mathbf{A} \mathbf{R}^\dagger \mathbf{R}$, then we have that

$$\mathbf{A} \mathbf{R}^\dagger \mathbf{R} \mathbf{v}_j = \mathbf{A} \mathbf{v}_j, \quad \text{for } j = 1, \dots, \rho.$$

Proof. First let $\mathbf{V}_{\mathbf{R},\rho} \in \mathbb{R}^{n \times \rho}$ contain the top ρ right singular vectors of \mathbf{R} , then the projection of \mathbf{A} onto the row space of \mathbf{R} is $\mathbf{A} \mathbf{R}^\dagger \mathbf{R} = \mathbf{A} \mathbf{V}_{\mathbf{R},\rho} \mathbf{V}_{\mathbf{R},\rho}^T$. Let the thin SVD of $\mathbf{A} \mathbf{V}_{\mathbf{R},\rho} \in \mathbb{R}^{m \times \rho}$ be $\tilde{\mathbf{U}} \tilde{\Sigma} \tilde{\mathbf{V}}^T$, where $\tilde{\mathbf{V}} \in \mathbb{R}^{\rho \times \rho}$. Then the compact SVD of $\mathbf{A} \mathbf{R}^\dagger \mathbf{R}$ is

$$\mathbf{A} \mathbf{R}^\dagger \mathbf{R} = \mathbf{A} \mathbf{V}_{\mathbf{R},\rho} \mathbf{V}_{\mathbf{R},\rho}^T = \tilde{\mathbf{U}} \tilde{\Sigma} \tilde{\mathbf{V}}^T \mathbf{V}_{\mathbf{R},\rho}^T.$$

According to the definition, \mathbf{v}_j is the j -th column of $(\mathbf{V}_{\mathbf{R},\rho}\tilde{\mathbf{V}}) \in \mathbb{R}^{n \times \rho}$, and thus \mathbf{v}_j lies on the column space of $\mathbf{V}_{\mathbf{R},\rho}$, and \mathbf{v}_j is orthogonal to $\mathbf{V}_{\mathbf{R},\rho^\perp}$. Finally, since $\mathbf{A} - \mathbf{A}\mathbf{R}^\dagger\mathbf{R} = \mathbf{A}\mathbf{V}_{\mathbf{R},\rho^\perp}\mathbf{V}_{\mathbf{R},\rho^\perp}^T$, we have that \mathbf{v}_j is orthogonal to $\mathbf{A} - \mathbf{A}\mathbf{R}^\dagger\mathbf{R}$, which directly proves the lemma. \square

B.2 Proof of Theorem 4

Proof. This randomized algorithm has three steps: approximate SVD via randomized projection [15], deterministic column selection via dual set sparsification algorithm [2] shown in Lemma 1, and the adaptive sampling algorithm of Theorem 3. This algorithm is a generalization of the near-optimal column selection algorithm of Lemma 2.

Given $\mathbf{A} \in \mathbb{R}^{m \times n}$ and a target rank $k < r_1$, Step 1 compute an approximate truncated SVD of \mathbf{A} in $\mathcal{O}(mnk/\epsilon_0)$ time such that $\mathbf{A}_k \approx \tilde{\mathbf{A}}_k = \tilde{\mathbf{U}}_k \tilde{\Sigma}_k \tilde{\mathbf{V}}_k^T$:

$$\mathbb{E}\|\mathbf{A} - \tilde{\mathbf{U}}_k \tilde{\Sigma}_k \tilde{\mathbf{V}}_k^T\|_F^2 \leq (1 + \epsilon_0)\|\mathbf{A} - \mathbf{A}_k\|_F^2.$$

Step 2 selects r_1 rows of \mathbf{A} to construct \mathbf{R}_1 by the dual set sparsification algorithm taking \mathcal{U} and \mathcal{V} as input, where \mathcal{U} contains all the m columns of $(\mathbf{A}^T - \tilde{\mathbf{A}}_k^T) \in \mathbb{R}^{n \times m}$, \mathcal{V} contains all the m columns of $\tilde{\mathbf{U}}_{\mathbf{A},k}^T \in \mathbb{R}^{k \times m}$. Theorem 4 in [2] has shown that

$$\mathbb{E}\|\mathbf{A} - \mathbf{A}\mathbf{R}_1^\dagger\mathbf{R}_1\|_F^2 \leq (1 + \epsilon_0)\left(1 + \frac{1}{(1 - \sqrt{k/r_1})^2}\right)\|\mathbf{A} - \mathbf{A}_k\|_F^2,$$

where the expectation is taken w.r.t. \mathbf{R}_1 . Step 2 costs $\mathcal{O}(mr_1k^2 + mn)$ time.

Step 3 samples additional r_2 rows of \mathbf{A} to construct $\mathbf{R}_2 \in \mathbb{R}^{r_2 \times n}$ by the adaptive sampling algorithm of Theorem 3. Let $\mathbf{R} = [\mathbf{R}_1^T, \mathbf{R}_2^T]^T \in \mathbb{R}^{(r_1+r_2) \times n}$. We apply Theorem 3 and have that

$$\begin{aligned} \mathbb{E}_{\mathbf{R}}\|\mathbf{A} - \mathbf{C}\mathbf{C}^\dagger\mathbf{A}\mathbf{R}^\dagger\mathbf{R}\|_F^2 &= \mathbb{E}_{\mathbf{R}_1}\left[\mathbb{E}_{\mathbf{R}_2}\left[\|\mathbf{A} - \mathbf{C}\mathbf{C}^\dagger\mathbf{A}\mathbf{R}^\dagger\mathbf{R}\|_F^2 \middle| \mathbf{R}_1\right]\right] \\ &\leq \mathbb{E}_{\mathbf{R}_1}\left[\|\mathbf{A} - \mathbf{C}\mathbf{C}^\dagger\mathbf{A}\|_F^2 + \frac{\rho}{r_2}\|\mathbf{A} - \mathbf{A}\mathbf{R}_1^\dagger\mathbf{R}_1\|_F^2\right] \\ &\leq \|\mathbf{A} - \mathbf{C}\mathbf{C}^\dagger\mathbf{A}\|_F^2 + \frac{\rho}{r_2}(1 + \epsilon_0)\left(1 + \frac{1}{(1 - \sqrt{k/r_1})^2}\right)\|\mathbf{A} - \mathbf{A}_k\|_F^2. \end{aligned}$$

By setting $r_1 = \mathcal{O}(k\epsilon^{-2/3})$, $r_2 \approx \frac{2\rho}{\epsilon}$, and $\epsilon_0 = \epsilon^{2/3}$, we conclude that

$$\mathbb{E}\|\mathbf{A} - \mathbf{C}\mathbf{C}^\dagger\mathbf{A}\mathbf{R}^\dagger\mathbf{R}\|_F^2 \leq \|\mathbf{A} - \mathbf{C}\mathbf{C}^\dagger\mathbf{A}\|_F^2 + \epsilon\|\mathbf{A} - \mathbf{A}_k\|_F^2.$$

The total computation time of the three steps is $\mathcal{O}(mnk/\epsilon_0 + mr_1k^2 + mn) = \mathcal{O}((mnk + mk^3)\epsilon^{-2/3})$ \square

B.3 Proof of Theorem 5

Proof. Since \mathbf{C} is constructed by columns of \mathbf{A} , the column space of \mathbf{C} is contained in the column space of \mathbf{A} , so $\text{rank}(\mathbf{C}\mathbf{C}^\dagger\mathbf{A}) = \text{rank}(\mathbf{C}) = \rho \leq c$, and thus the assumptions of Theorem 4 are satisfied. Lemma 2 and Theorem 4 together prove Theorem 5:

$$\begin{aligned} \mathbb{E}^2\|\mathbf{A} - \mathbf{C}\mathbf{U}\mathbf{R}\|_F &\leq \mathbb{E}\|\mathbf{A} - \mathbf{C}\mathbf{U}\mathbf{R}\|_F^2 = \mathbb{E}_{\mathbf{C},\mathbf{R}}\|\mathbf{A} - \mathbf{C}\mathbf{C}^\dagger\mathbf{A}\mathbf{R}^\dagger\mathbf{R}\|_F^2 \\ &= \mathbb{E}_{\mathbf{C}}\left[\mathbb{E}_{\mathbf{R}}\left[\|\mathbf{A} - \mathbf{C}\mathbf{C}^\dagger\mathbf{A}\mathbf{R}^\dagger\mathbf{R}\|_F^2 \middle| \mathbf{C}\right]\right] \\ &\leq \mathbb{E}_{\mathbf{C}}\left[\|\mathbf{A} - \mathbf{C}\mathbf{C}^\dagger\mathbf{A}\|_F^2 + \epsilon\|\mathbf{A} - \mathbf{A}_k\|_F^2\right] \\ &\leq (1 + 2\epsilon)\|\mathbf{A} - \mathbf{A}_k\|_F^2. \end{aligned}$$

Finally we have $\mathbb{E}\|\mathbf{A} - \mathbf{C}\mathbf{U}\mathbf{R}\|_F \leq (1 + \epsilon)\|\mathbf{A} - \mathbf{A}_k\|_F$ because $1 + 2\epsilon \leq (1 + \epsilon)^2$.

The time cost of the fast CUR algorithm is the sum of Stage 1, Stage 2, and the Moore-Penrose inverse of \mathbf{C} and \mathbf{R} , i.e. $\mathcal{O}((mnk + nk^3)\epsilon^{-2/3}) + \mathcal{O}((mnk + mk^3)\epsilon^{-2/3}) + \mathcal{O}(mc^2) + \mathcal{O}(nr^2) = \mathcal{O}(mnk\epsilon^{-2/3} + (m + n)k^3\epsilon^{-2/3} + mk^2\epsilon^{-2} + nk^2\epsilon^{-4})$. \square