Efficient Spectrum-Revealing CUR Matrix Decomposition

Cheng Chen

Ming Gu

Zhihua ZhangPeking University

Shanghai Jiao Tong University

University of California, Berkeley

Yong Yu

Shanghai Jiao Tong University

Weinan Zhang

Shanghai Jiao Tong University

Abstract

The CUR matrix decomposition is an important tool for low-rank matrix approximation. It approximates a data matrix though selecting a small number of columns and rows of the matrix. Those CUR algorithms with gapdependent approximation bounds can obtain high approximation quality for matrices with good singular value spectrum decay, but they have impractically high time complexities. In this paper, we propose a novel CUR algorithm based on truncated LU factorization with an efficient variant of complete pivoting. Our algorithm has gap-dependent approximation bounds on both spectral and Frobenius norms while maintaining high efficiency. Numerical experiments demonstrate the effectiveness of our algorithm and verify our theoretical guarantees.

1 Introduction

CUR matrix decomposition is a well known method for low-rank approximation (Mahoney and Drineas, 2009; Drineas et al., 2008). It approximates the data matrix $\bf A$ as the product of three matrices: $\bf A \approx CUR$, where $\bf C$ and $\bf R$ are composed of sampled columns and sampled rows of matrix $\bf A$ respectively. Compared with other low-rank approximation methods such as truncated Singular Value Decomposition (SVD), CUR decomposition can better preserve the sparsity of the input matrix and facilitate the interpretation of the computed results. Because of these advantages, CUR decomposition is more attractive than SVD in text

Proceedings of the 23rdInternational Conference on Artificial Intelligence and Statistics (AISTATS) 2020, Palermo, Italy. PMLR: Volume 108. Copyright 2020 by the author(s).

mining, collaborative filtering, bioinformatics, image and video processing (Drineas et al., 2008; Xu et al., 2015; Wang and Zhang, 2012, 2013; Mahoney et al., 2008; Mackey et al., 2011).

In recent years, many variants of the CUR decomposition have been developed (Drineas et al., 2008; Wang and Zhang, 2012, 2013; Boutsidis and Woodruff, 2017; Bien et al., 2010; Anderson et al., 2015; Drineas et al., 2006; Wang et al., 2016; Song et al., 2019). Many randomized CUR algorithms adopt random sampling on columns and rows of data matrices. They aim to achieve $1+\varepsilon$ relative error bounds on Frobenius norm with a failure probability δ (Drineas et al., 2008; Wang and Zhang, 2013; Boutsidis and Woodruff, 2017), i.e., $\|\mathbf{A} - \mathbf{CUR}\|_F^2 \le (1+\varepsilon)\|\mathbf{A} - \mathbf{A}_k\|_F^2$, where k is the target rank of approximation. These bounds are tight theoretically, but are less useful in practical applications because they require $O(k/\varepsilon)$ columns and rows. For example, Algorithm 3 of (Boutsidis and Woodruff, 2017) requires to select $4k+4820k/\varepsilon$ columns and rows, but only gets a success probability of 0.16. If we choose k = 10 and $\varepsilon = 0.5$ (very common setting), the number of selected columns and rows is close to 10⁶. For real-world data sets, it is impractical to select so many columns and rows. Anderson et al. (2015) proposed a deterministic column selection method for CUR decomposition with gap-dependent approximation bounds. These bounds allow their method only to choose $\ell = k + O(1)$ columns and rows when the given matrix has a rapidly decaying spectrum. However, their algorithm is highly time-consuming and requires $O(\ell mn(m+n))$ time, where ℓ is the number of selected columns and rows.

A natural question now arises: can we develop an efficient CUR algorithm which is as efficient as previous randomized CUR algorithms while maintaining gap-dependent approximation bounds? In this paper, we address this question via presenting a novel CUR algorithm, namely Spectrum-Revealing CUR (SRCUR) decomposition. Our method borrows the idea from a

classic numerical linear algebra method, i.e., pivoted LU factorization (Trefethen and Bau III, 1997). Unlike existing CUR algorithms which sample columns and rows separately (Drineas et al., 2008; Wang and Zhang, 2013; Boutsidis and Woodruff, 2017; Anderson et al., 2015), the pivoting procedure of LU factorization selects columns and rows simultaneously and considers the inter-connectivity of selected columns and rows. However, the vanilla pivoting procedure is costly and does not bound the relative error. To make use of the advantages and bypass the disadvantages of pivoted LU, we propose Fast Spectrum-Revealing LU (FSRLU) factorization, which is efficient and has theoretical guarantees. Our SRCUR algorithm first adopts FSRLU to select columns and rows and then computes the intersection matrix. The main contributions of this paper are summarized as follows:

- We propose FSRLU algorithm which can efficiently compute truncated LU factorization with complete pivoting. We also analyze the stability and time complexity of FSRLU.
- We propose SRCUR algorithm based on FSRLU. We provide error analysis for SRCUR algorithm in both spectral and Frobenius norms. The proofs show that our method has gap dependency error bounds with regard to the quadratic of the spectral gap $\sigma_{k+1}/\sigma_{p+1}$. These error bounds indicate that our method only needs to select $\ell = k + O(1)$ columns and rows for data matrices with rapidly decaying singular values.
- SRCUR is much faster than existing gapdependent CUR algorithms. The time complexity of our method is $O(\text{nnz}(\mathbf{A})\log n) + \tilde{O}(\ell^2(m+n))$, which is comparable to the state-of-the-art randomized CUR algorithms.
- We validate our method with numerical experiments on four real-world datasets, which demonstrate that our algorithms are substantially faster than existing CUR algorithms while maintaining good performance.

The remainder of the paper is organized as follows. In Section 2, we describe notation and preliminaries used in this paper. In Section 3 we introduce related work. Then we present our algorithms in Section 4 and show theoretical analysis in Section 5. Experimental results are given in Section 6. Conclusions and discussions are provided in Section 7.

2 Notation and Preliminaries

We use \mathbf{I}_m to denote $m \times m$ identity matrix, and $\mathbf{0}$ to denote a zero vector or matrix of appropriate size.

For a vector $\mathbf{x} = (x_1, \dots, x_n)^{\top}$, let $\|\mathbf{x}\|_2 = \sqrt{\sum_i x_i^2}$ be the ℓ_2 -norm of \mathbf{x} . For a matrix $\mathbf{A} = [a_{ij}] \in \mathbb{R}^{m \times n}$, let $\|\mathbf{A}\|_2$ denote the spectral norm and $\|\mathbf{A}\|_F$ denote the Frobenius norm. We use $\operatorname{nnz}(\mathbf{A})$ to denote the number of nonzero elements of \mathbf{A} . Let $\det(\mathbf{A})$ be the determinant of \mathbf{A} and $\operatorname{adj}(\mathbf{A})$ be the adjugate matrix of \mathbf{A} . We use \tilde{O} to hide logarithmic factors.

Let ρ be the rank of matrix **A**. The reduced SVD of **A** is defined as

$$\mathbf{A} = \mathbf{U} \mathbf{\Sigma} \mathbf{V}^{\top} = \sum_{i=1}^{\rho} \sigma_i \mathbf{u}_i \mathbf{v}_i^{\top},$$

where σ_i are positive singular values in the descending order. That is, $\sigma_i(\mathbf{A})$ is the i-th largest singular value of \mathbf{A} . Let $\mathbf{A}_k = \sum_{i=1}^k \sigma_i \mathbf{u}_i \mathbf{v}_i^{\top}$ denote the best rank-k approximation of \mathbf{A} and $\mathbf{A}^{\dagger} = \mathbf{V} \mathbf{\Sigma}^{-1} \mathbf{U}^{\top}$ denote the Moore-Penrose pseudo-inverse of \mathbf{A} . $\kappa(\mathbf{A}) = \sigma_{\max}/\sigma_{\min}$ to denote the condition number of matrix \mathbf{A} . Here σ_{\max} is the maximum singular value and σ_{\min} is the minimum non-zero singular value of \mathbf{A} .of \mathbf{A} when \mathbf{A} is square.

We use the Matlab colon notation for a block matrix. Let $\mathbf{A}_{i,:}$ be the *i*-th row of \mathbf{A} and $\mathbf{A}_{:,j}$ be the *j*-th column of \mathbf{A} . We use $\mathbf{A}_{i:j,:}$ to denote $\begin{bmatrix} \mathbf{A}_{:,p}^{\top}, \mathbf{A}_{i+1,:}^{\top}, \dots, \mathbf{A}_{j,:}^{\top} \end{bmatrix}^{\top}$ and $\mathbf{A}_{:,p:q}$ to denote $\begin{bmatrix} \mathbf{A}_{:,p}, \mathbf{A}_{:,p+1}, \dots, \mathbf{A}_{:,q} \end{bmatrix}$. Let $\mathbf{A}_{i:j,p:q}$ denote the block matrix $\begin{pmatrix} a_{ip} & \dots & a_{iq} \\ \vdots & \ddots & \vdots \\ a_{ip} & \dots & a_{iq} \end{pmatrix}$.

Johnson-Lindenstrauss (JL) transform (Johnson and Lindenstrauss, 1984) is a powerful tool for dimensionality reduction. It has been proved to preserve the vector norm within an ε -error ball. We present the Johnson-Lindenstrauss Lemma as follows:

Lemma 1. For any vector $\mathbf{x} \in \mathbb{R}^d$, $0 < \varepsilon, \delta < 1/2$, there exists a JL transform matrix $\mathbf{S} \in \mathbb{R}^{p \times d}$, with $p = \Theta(\log(1/\delta)\varepsilon^{-2})$, for which satisfies

$$(1 - \varepsilon) \|\mathbf{x}\|_2^2 \le \|\mathbf{S}\mathbf{x}\|_2^2 \le (1 + \varepsilon) \|\mathbf{x}\|_2^2,$$

with probability $1 - \delta$.

3 Related Work

In Section 3.1, we review some representative methods of CUR matrix decomposition. In Section 3.2, we introduce truncated LU factorization with complete pivoting.

3.1 CUR Matrix Decomposition

CUR decomposition approximates the data matrix with actual columns and rows. Given a matrix $\mathbf{A} \in \mathbb{R}^{m \times n}$,

CUR decomposition selects a subset of columns $\mathbf{C} \in \mathbb{R}^{m \times c}$ and a subset of rows $\mathbf{R} \in \mathbb{R}^{r \times n}$ and computes an intersection matrix $\mathbf{U} = \mathbf{C}^{\dagger} \mathbf{A} \mathbf{R}^{\dagger}$. Then, $\widehat{\mathbf{A}} = \mathbf{C} \mathbf{U} \mathbf{R}$ is a low-rank approximation of matrix \mathbf{A} . Some works try to compute an approximation $\widehat{\mathbf{A}}_k$ with exact rank k (Boutsidis and Woodruff, 2017; Anderson et al., 2015). The rank-constrained approximation is more comparable to the truncated SVD decomposition.

The essential of CUR decomposition is how to select columns and rows effectively and efficiently. The traditional approaches select columns according to deterministic pivoting rules (Stewart, 1999; Chan, 1987; Gu and Eisenstat, 1996; Berry et al., 2005), but these methods do not have good approximation errors. Recently, randomized algorithms for CUR decomposition aims to obtain $1+\varepsilon$ relative error bounds in Frobenius norm (Mahoney and Drineas, 2009; Drineas et al., 2008; Wang and Zhang, 2013; Boutsidis and Woodruff, 2017; Mahoney et al., 2011):

$$\|\mathbf{A} - \mathbf{CUR}\|_F^2 \le (1+\varepsilon)\|\mathbf{A} - \mathbf{A}_k\|_F^2.$$

Drineas et al. (2008) randomly sampled columns and rows according to their leverage scores. They achieve $1+\varepsilon$ relative error bounds in Frobenius norm. Wang and Zhang (2013) proposed an adaptive CUR algorithm with better theoretical results, and their results were further improved by (Boutsidis and Woodruff, 2017; Song et al., 2019). However, these methods have to choose $O(k/\varepsilon)$ rows and columns to achieve the $1+\varepsilon$ error bound, which is usually impractical for real-world data set.

Anderson et al. (2015) proposed unweighted column selection for CUR decomposition which provides spectral gap error bounds as the following form:

$$\|\mathbf{A} - \mathbf{CUR}\|_{\varepsilon}^2 \le (1 + O(\omega^2)) \|\mathbf{A} - \mathbf{CUR}\|_{\varepsilon}^2$$

for $\xi \in \{2, F\}$. Here ω is a number that relies on the rate of spectrum decay of \mathbf{A} and the number of oversampling l-k. Their method only needs to choose $\ell = k+O(1)$ columns and rows to obtain a good rank-k approximation when the data matrix has rapid spectrum decay. However, their algorithm is very expensive and requires $O(\ell mn(m+n))$ time to obtain a rank-k approximation.

There have been several investigations of ways to construct the approximation. Wang et al. (2016) proposed an efficient method to compute \mathbf{U} by matrix sketching. Anderson et al. (2015) proposed a numerically stable algorithm to construct $\widehat{\mathbf{A}}$ and $\widehat{\mathbf{A}}_k$. This method avoids computing the pseudo-inverse of matrices by performing QR factorization on \mathbf{C} and \mathbf{R} . We present their algorithm in Appendix A.

Some other variants of CUR algorithms including tensor CUR decomposition (Song et al., 2019) and CUR with

 ℓ_1 -norm (Song et al., 2017) are less relevant to our work.

3.2 LU factorization with Complete Pivoting

LU factorization factors a matrix as a product of a unit lower triangular matrix and an upper triangular matrix with row permutation and/or column permutation. Partial pivoting is widely used in numerical algebra area because it is more efficient. However, it may fail on rank deficient matrices. The complete pivoting, even though very expensive, is necessary for those matrices to guarantee that factorized matrices are well-defined. Recently, Melgaard and Gu (2015) proposed a randomized LU factorization with complete pivoting by leverage random projection.

Rank-revealing LU algorithms (Pan, 2000; Miranian and Gu, 2003) are good quality approximations. They guarantee that the approximation capture the rank of the data matrix within a low polynomial factor in k, m and n, which could be very large. Spectrum-revealing LU (Anderson and Gu, 2017) performs extra swaps to ensure spectral gap error bounds, but it also introduces extra computational costs.

4 Our Approaches

In this section, we propose our main algorithms. We first propose the EELM algorithm which quickly estimates the element with largest magnitude of a given matrix. Then we propose the FSRLU algorithm to efficiently computes the spectrum-revealing LU factorization. Finally, we propose our SRCUR algorithm based on the FSRLU.

4.1 Estimating the Element with Largest Magnitude

Finding the element with largest magnitude in a matrix has many applications (Higham and Relton, 2016). Usually computing $\|\mathbf{A}\|_{\max} \in \mathbb{R}^{m \times n}$ from its definition has O(mn) cost. However, if we can only access matrix \mathbf{A} implicitly, such as through a product $\mathbf{A} = \mathbf{BC}$ or an inverse $\mathbf{A} = \mathbf{B}^{-1}$, the computational cost can be prohibitive if we need to calculate all elements of matrix \mathbf{A} . The method in (Higham and Relton, 2016) uses matrix-vector product to estimate the position of the element with the largest magnitude. But their method does not have theoretical guarantees on the estimated value.

We propose an efficient algorithm to estimate the element with the largest magnitude, called EELM. Given matrix **A** and projected matrix $\mathbf{R} = \Omega \mathbf{A}$ ($\Omega \in \mathbb{R}^{p \times m}$ is a JL transform matrix), EELM first finds the col-

Algorithm 1 Estimating the Element with Largest Magnitude (EELM)

Input: Input matrix $\mathbf{A} \in \mathbb{R}^{m \times n}$, $\mathbf{R} \in \mathbb{R}^{p \times n}$.

Output: row index r, column index c, estimated value x.

- 1: $c = \underset{i \le j \le n}{\operatorname{arg max}} \| \mathbf{R}_{:,j} \|_2$.
- $i \le j \le n$ 2: Compute the *c*-th column of **A**.
- 3: $r = \underset{i \le j \le m}{\operatorname{arg\,max}} |\mathbf{A}_{j,c}|.$
- 4: $x = |\bar{\mathbf{A}}_{r,c}|$.

umn with the maximum ℓ_2 -norm in the sketched matrix \mathbf{R} and then chooses the largest element in the corresponding column of the original matrix \mathbf{A} . We describe EELM in Algorithm 1. Since JL transform preserves the vector norm well, the estimated value can be bounded as the following theorem:

Theorem 1. The estimated value x given by Algorithm 1 satisfies $x \ge \sqrt{\frac{1-\varepsilon}{mn(1+\varepsilon)}} \|\mathbf{A}\|_F$ with probability $1-\delta$.

As shown in the next section, ε only affect the time complexity of our algorithm but have little influence on the error bounds. Thus, in the rest of our paper, we assume that $\varepsilon=0.5$ and $\delta=0.01$ for convenience.

4.2 Fast Spectrum-Revealing LU Factorization

In this section, we describe our FSRLU algorithm which consists of two parts: Fast Pivoting for Truncated LU factorization (FPTLU) and Fast Spectrum-Revealing Pivoting (FSRP). FPTLU computes a rough pivoting for truncated LU, and FSRP performs extra operations on the result of FPTLU in order to bound the error.

4.2.1 Fast Pivoting for Truncated LU Factorization

Truncated LU factorization aims to compute $\Pi_1 A \Pi_2 = \widehat{\mathbf{L}}\widehat{\mathbf{U}} + \begin{pmatrix} \mathbf{0} & \mathbf{0} \\ \mathbf{0} & \mathbf{S} \end{pmatrix}$, where Π_1 and Π_2 are permutation matrices, \mathbf{S} is the Schur complement. In each iteration, the vanilla complete pivoting of LU factorization choose the largest magnitude element in the Schur complement as the pivot. This procedure is very expensive. To reduce the computational cost, we adopt EELM algorithm to estimate pivots of LU. We present FPTLU in Algorithm 2. Note that line 5-7 of Algorithm 2 is standard LU update. Note that the sketched matrix \mathbf{B} can be updated efficiently as in (Anderson and Gu, 2017) (see Appendix B).

Algorithm 2 Fast Pivoting for Truncated LU factorization

Input: data matrix $\mathbf{A} \in \mathbb{R}^{m \times n}$, target number of columns and rows ℓ .

Output: $\widehat{\mathbf{L}} \in \mathbb{R}^{m \times \ell}$, $\widehat{\mathbf{U}} \in \mathbb{R}^{\ell \times n}$, $\mathbf{\Pi}_1$, $\mathbf{\Pi}_2$.

- 1: Generate JL transform matrix $\Omega \in \mathbb{R}^{p_1 \times m}$, compute $\mathbf{B} = \Omega \mathbf{A}$, $\Pi_1 = \mathbf{I}_m$, $\Pi_2 = \mathbf{I}_n$.
- 2: **for** $i = 1, ..., \ell$ **do**
- 3: $[c, r, x] = \text{EELM}(\mathbf{A}, \mathbf{B}).$
- 4: Swap *i*-th and *c*-th column of \mathbf{A} , update $\mathbf{\Pi}_2$, Swap *i*-th and *r*-th row of \mathbf{A} , update $\mathbf{\Pi}_1$.
- 5: $\mathbf{A}_{i:m,i} = \mathbf{A}_{i:m,i} \mathbf{A}_{i:m,1:i-1} \mathbf{A}_{1:i-1,i}, \ \mathbf{A}_{i+1:m,i} = \mathbf{A}_{i+1:m,i} / \mathbf{A}_{i,i}.$
- 6: $\mathbf{A}_{i,i+1:n} = \mathbf{A}_{i,i+1:n} \mathbf{A}_{i,1:i-1} \mathbf{A}_{1:i-1,i+1:n}$
- 7: Update $\mathbf{B}_{::i+1:n}$.
- 8: end for
- 9: Let $\widehat{\mathbf{L}}$ be the lower triangular part of $\mathbf{A}_{:,1:\ell}$ with unit diagonal.
- 10: Let $\hat{\mathbf{U}}$ be the upper triangular part of $\mathbf{A}_{1:\ell::}$.

4.2.2 Fast Spectrum-Revealing Pivoting

Given a truncated LU factorization in the following form.

$$\mathbf{\Pi}_1 \mathbf{A} \mathbf{\Pi}_2 = \begin{pmatrix} \mathbf{L}_{11} & & \\ \mathbf{l}_{21}^\top & 1 & \\ \mathbf{L}_{31} & & \mathbf{I} \end{pmatrix} \begin{pmatrix} \mathbf{U}_{11} & \mathbf{u}_{12} & \mathbf{U}_{13} \\ & \alpha & \mathbf{s}_{12}^\top \\ & \mathbf{s}_{21} & \mathbf{S}_{22} \end{pmatrix},$$

where
$$\widehat{\mathbf{L}} = \begin{pmatrix} \mathbf{L}_{11} \\ \mathbf{l}_{21} \\ \mathbf{L}_{31} \end{pmatrix}$$
, $\widehat{\mathbf{U}} = \begin{pmatrix} \mathbf{U}_{11} & \mathbf{u}_{12} & \mathbf{U}_{13} \end{pmatrix}$ and $\mathbf{S} =$

 $\begin{pmatrix} \alpha & \mathbf{s}_{12}^{\top} \\ \mathbf{s}_{21} & \mathbf{S}_{22} \end{pmatrix}$, FSRP aims to find the element α with the largest magnitude in \mathbf{S} and the largest element β in matrix $|\overline{\mathbf{A}}_{11}^{-\top}|$ in each iteration. Here $\overline{\mathbf{A}}_{11} = \begin{pmatrix} \mathbf{L}_{11} \\ \mathbf{l}_{21}^{\top} & 1 \end{pmatrix} \begin{pmatrix} \mathbf{U}_{11} & \mathbf{u}_{12} \\ \alpha \end{pmatrix}$. Since computing \mathbf{S} is very expensive, we use EELM to estimate α and β rather than compute their accurate value. We present the FSRP in Algorithm 3. In each iteration, the sketched matrix \mathbf{R} and \mathbf{W} only needs a rank-2 update after each swap. We leverage techniques in (Gondzio, 1992) to turn \mathbf{L} and \mathbf{U} back into trapezoidal form after the swaps are performed. We show details of updating $\hat{\mathbf{L}}$ and $\hat{\mathbf{U}}$ in Appendix \mathbf{C} .

4.3 Spectrum-Revealing CUR Decomposition

The permutation matrices produced by FSRP algorithm indicate which columns and rows we should choose for CUR decomposition. Namely, we compute matrix ${\bf C}$ and ${\bf R}$ as

$$\mathbf{C} = \boldsymbol{\Pi}_1^{\top} \widehat{\mathbf{L}} \widehat{\mathbf{U}}_{:,1:\ell} \quad \text{ and } \quad \mathbf{R} = \widehat{\mathbf{L}}_{1:\ell,:} \widehat{\mathbf{U}} \boldsymbol{\Pi}_2^{\top}.$$

Algorithm 3 Fast Spectrum-revealing Pivoting

Input: Truncated LU factorization $\Pi_1 \mathbf{A} \Pi_2 \approx \widehat{\mathbf{L}} \widehat{\mathbf{U}}$ and tolerance f.

Output: $\widehat{\mathbf{L}} \in \mathbb{R}^{m \times \ell}$, $\widehat{\mathbf{U}} \in \mathbb{R}^{\ell \times n}$, $\mathbf{\Pi}_1$, $\mathbf{\Pi}_2$.

- 1: Generate JL transform matrix $\Omega_1 \in \mathbb{R}^{p_2 \times (m-\ell)}$, $\Omega_2 \in \mathbb{R}^{p_3 \times (\ell+1)}$.
- 2: $\mathbf{R} = \mathbf{\Omega}_1 \mathbf{S}$, $[s_r, s_c, \alpha] = \text{EELM}(\mathbf{S}, \mathbf{R})$.
- 3: Swap the $(\ell + 1)$ -th row and the $(\ell + s_r)$ -th row of $\widehat{\mathbf{L}}$ and update $\mathbf{\Pi}_1$ correspondingly.
- 4: Swap the $(\ell + 1)$ -th column and the $(\ell + s_c)$ -th column of $\widehat{\mathbf{U}}$ and update $\mathbf{\Pi}_2$ correspondingly.
- 5: $\mathbf{W} = \mathbf{\Omega}_2 \overline{\mathbf{A}}_{11}^{-\top}, [a_r, a_c, \beta] = \overline{\mathrm{EELM}}(\overline{\mathbf{A}}_{11}^{-\top}, \mathbf{W}).$
- 6: while $\alpha\beta > f$ do
- 7: Exchange the a_r -th row and the $(\ell + 1)$ -th row of **A** and update $\hat{\mathbf{L}}$.
- 8: Exchange the a_c -th column and the $(\ell + 1)$ -th column of **A** and update $\widehat{\mathbf{U}}$.
- 9: Update Π_1 , Π_2 , S, R and perform step 3,4,5.
- 10: Update **W**, let $[a_r, a_c, \beta] = \text{EELM}(\overline{\mathbf{A}}_{11}^{-1}, \mathbf{W})$.
- 11: end while

Algorithm 4 Spectrum-Revealing CUR decomposition

Input: Data matrix $\mathbf{A} \in \mathbb{R}^{m \times n}$, target rank k, target column and row number ℓ , f > 1.

Output: $\widehat{\mathbf{A}} \in \mathbb{R}^{m \times n}$ and $\widehat{\widehat{\mathbf{A}}}_k \in \mathbb{R}^{m \times n}$.

- 1: Perform Algorithm 2 to compute a truncated LU factorization $\Pi_1 \mathbf{A} \Pi_2 \approx \widehat{\mathbf{L}} \widehat{\mathbf{U}}$.
- 2: Perfrom Algorithm 3 to update Π_1 , Π_2 , $\widehat{\mathbf{L}}$ and $\widehat{\mathbf{U}}$.
- 3: Compute $\mathbf{C} = \mathbf{\Pi}_1^{\top} \widehat{\mathbf{L}} \widehat{\mathbf{U}}_{:,1:\ell}$ and $\mathbf{R} = \widehat{\mathbf{L}}_{1:\ell,:} \widehat{\mathbf{U}} \mathbf{\Pi}_2^{\top}$.
- 4: Perform Stable CUR Algorithm to obtain $\widehat{\mathbf{A}}$ and $\widehat{\mathbf{A}}_k.$

To make our algorithm more stable, we construct the approximation matrices by StableCUR (Anderson et al., 2015), which is shown in Appendix A. We show our SRCUR algorithm in Algorithm 4.

5 Theoretical Analysis

In this section, we analyze time complexity of our algorithm in Section 5.1 and provide spectral gap error bounds on singular values, 2-norm and Frobenius norm in Section 5.2.

5.1 Time complexity

It is obvious that the time cost of EELM is O(m+pn). Then we discuss the time complexity of FPTLU. If we choose Gaussian random projection matrix as the JL transform matrix, we only require $O(p_1 \text{nnz}(\mathbf{A}))$ to initialize matrix \mathbf{B} . In each iteration, FPTLU

require $O(p_1n)$ time to update **B** and $O(\ell^2(m+n))$ time to perform standard LU update(line 5-7). Since $p_1 = O(\log n)$, FPTLU can compute a truncated LU factorization in $O(\operatorname{nnz}(\mathbf{A})\log n + \ell^2(m+n))$ time.

To analyze the complexity of the FSRP, we first give the following theorem which shows the upper bound of number of iterations.

Theorem 2. If the input matrix $\hat{\mathbf{L}}$ and $\hat{\mathbf{U}}$ are obtained by Algorithm 2, Algorithm 3 performs at most $O(\ell \log(mn))$ number of iterations with probability 0.99.

We require $O(\operatorname{nnz}(\mathbf{A})\log n + \ell(m+n)\log n)$ to initialize matrix \mathbf{R} and $O(\ell^2\log\ell)$ to initialize matrix \mathbf{W} . In each iteration, the FSRP costs $O(pn+\ell(m+n)+\ell^2\log\ell)$. Using Theorem 2, we know that the time complexity of FSRP is $O(\operatorname{nnz}(\mathbf{A})\log\ell + \ell pn\log(mn) + \ell^2(m+n)\log(mn) + \ell^3\log(mn)\log\ell)$. The StableCUR require $O(\ell^2(m+n))$ time to construct the CUR approximation. Therefore, our SRCUR algorithm needs $O(\operatorname{nnz}(\mathbf{A})\log n) + \tilde{O}(\ell^2(m+n))$ time in total.

5.2 Spectral Gap Error Bounds

We present the norm error bounds in Theorem 3 and the singular value bound in Theorem 4. The detailed proofs are presented in the supplementary material due to page limit.

Theorem 3. For $\gamma = O(f\ell\sqrt{mn})$, Algorithm 4 satisfies:

$$\|\mathbf{A} - \widehat{\mathbf{A}}\|_{2} \leq \|\mathbf{A} - \widehat{\mathbf{A}}\|_{F} \leq \gamma \sigma_{\ell+1}(\mathbf{A}), \tag{1}$$
$$\|\mathbf{A} - \widehat{\mathbf{A}}_{k}\|_{F}^{2} \leq \left(1 + \frac{2\gamma^{2}\sigma_{\ell+1}^{2}(\mathbf{A})}{\sum_{i=k+1}^{\operatorname{rank}(\mathbf{A})} \sigma_{i}^{2}(\mathbf{A})}\right) \|\mathbf{A} - \mathbf{A}_{k}\|_{F}^{2}, \tag{2}$$

$$\|\mathbf{A} - \widehat{\mathbf{A}}_k\|_2^2 \le \left(1 + 2\gamma^2 \frac{\sigma_{\ell+1}^2(\mathbf{A})}{\sigma_{k+1}^2(\mathbf{A})}\right) \|\mathbf{A} - \mathbf{A}_k\|_2^2, \quad (3)$$

with probability 0.98.

Theorem 4. For $1 \leq j \leq \ell$ and $\gamma = O(f\ell\sqrt{mn})$, Algorithm 4 satisfies:

$$\sigma_j(\mathbf{A}) \ge \sigma_j(\widehat{\mathbf{A}}) \ge \sqrt{1 - 2\gamma^2 \left(\frac{\sigma_{\ell+1}(\mathbf{A})}{\sigma_j(\mathbf{A})}\right)^2} \sigma_j(\mathbf{A})$$

 $with\ probability\ 0.98.$

5.3 Discussion

We would like to emphasize the differences between our work and those methods which achieve $1+\varepsilon$ relative error bounds (Drineas et al., 2008; Wang and Zhang, 2013; Boutsidis and Woodruff, 2017). First, $1+\varepsilon$ relative error bounds only bound Frobenius norm, while

we provide bounds on spectral norm, Frobenius norm and singular values. In addition, our bounds are data dependent, while $1+\varepsilon$ relative error bounds depend on the universal constant ε .

Now we compare our bounds with DetUCS algorithm (Anderson et al., 2015). First, both our bounds and those of DetUCS are related to the quadratic of the spectral gap $\sigma_{\ell+1}^2(\mathbf{A})/\sigma_{k+1}^2(\mathbf{A})$. Thus, our bounds have the same order as that of DetUCS. However, DetUCS costs $O(\ell mn(m+n))$ time, which is much higher than our costs.

Finally, we discuss the differences between our method and SRLU (Anderson and Gu, 2017). First, our method and SRLU are designed for different tasks. The goal of SRLU is to obtain a high quality LU factorization, while we aim to compute CUR decomposition. Also, SRLU does not have time guarantees on the extra swap procedure. Thus, our FSRLU is faster than SRLU while maintaining the same theoretical guarantee. Though SRLU provided bounds on $\|\mathbf{\Pi}_1 \mathbf{A} \mathbf{\Pi}_2 - \widehat{\mathbf{L}} \mathbf{M} \widehat{\mathbf{U}}\|_{\xi}$, where $\mathbf{M} = \widehat{\mathbf{L}}^{\dagger} \mathbf{A} \widehat{\mathbf{U}}^{\dagger}$ and $\xi \in \{2, F\}$. Their approximation is not rank-constrained and worse than our result.

6 Experiments

In this section, we empirically evaluate the performance of our SRCUR algorithm. We choose Gaussian random projection matrix as the JL transform matrix in SRCUR. The baseline methods include uniform sampling (Williams and Seeger, 2001), subspace sampling (Drineas et al., 2008), near-optimal sampling (Wang and Zhang, 2013; Boutsidis and Woodruff, 2017) and deterministic unweighted column selection (Anderson et al., 2015).

In our empirical evaluation, we consider following six measures: $\|\mathbf{A} - \widehat{\mathbf{A}}\|_F / \|\mathbf{A} - \mathbf{A}_k\|_F$, $\|\mathbf{A} - \widehat{\mathbf{A}}\|_2 / \|\mathbf{A} - \mathbf{A}_k\|_2$, $\|\mathbf{A} - \widehat{\mathbf{A}}_k\|_F / \|\mathbf{A} - \mathbf{A}_k\|_F$, $\|\mathbf{A} - \widehat{\mathbf{A}}_k\|_2 / \|\mathbf{A} - \mathbf{A}_k\|_2$, $\sigma_k(\widehat{\mathbf{A}})/\sigma_k(\mathbf{A})$ and running time. We compare CUR algorithms on the following four datasets: The Extended Yale Face Database B (Georghiades et al., 2001), Dexter (Guyon et al., 2005), Sido¹ and Reuters². We summarize the information of these data sets in Table 1 and show the singular value distribution in Figure 1. We choose $\ell = k+1, 2k, \cdots, 10k$ for all datasets, where $k = \lceil \min(m,n)/200 \rceil$. We report our results in Figures 2-5. All results of randomized algorithm are averaged over 5 runs. The experiments are all performed in Matlab R2017b.

The figures show that our algorithm is much more efficient than state-of-the-art CUR algorithms. The

Table 1: Summary of datasets for CUR matrix decomposition

Dataset	m	n	$\%\mathrm{nnz}$
Extend Yale Face B	2414	1024	97.84
Dexter	2600	20000	4.8
Sido	4932	12678	9.84
Reuters	8293	18933	0.25

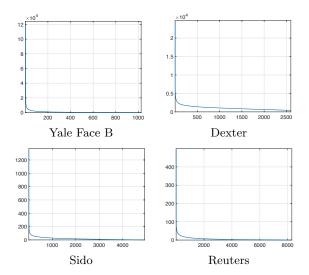


Figure 1: Singular value distribution

main reason is that state-of-the-art CUR algorithms all require to perform SVD or randomized SVD while SR-CUR is SVD-free. The randomized SVD, though has good time complexity Clarkson and Woodruff (2017), is much less efficient than LU update because the constant of SVD is extreme large. Furthermore, the approximation errors of SRCUR algorithm are better than or comparable with other algorithms. Especially, SRCUR performs better than other methods when ℓ is close to k. This result matches our theoretical analysis in Section 5.

7 Conclusions

In this paper, we have proposed Spectrum-Revealing CUR decomposition algorithm. We have built the bridge between CUR decomposition and LU factorization, and leveraged truncated LU factorization to select columns and rows for CUR decomposition. We have develop spectral gap error bound on spectral norm and Frobenius norm and singular values. These bounds confirm high quality approximations computed by our methods for matrices with fast enough singular value decay, a property which is typically observed in practical data matrices. We have also conducted experiments to empirically demonstrate the effectiveness and reliability of our algorithms on four machine learning datasets.

¹http://www.causality.inf.ethz.ch/data/SIDO.html

²http://www.zjucadcg.cn/dengcai/Data/TextData.html

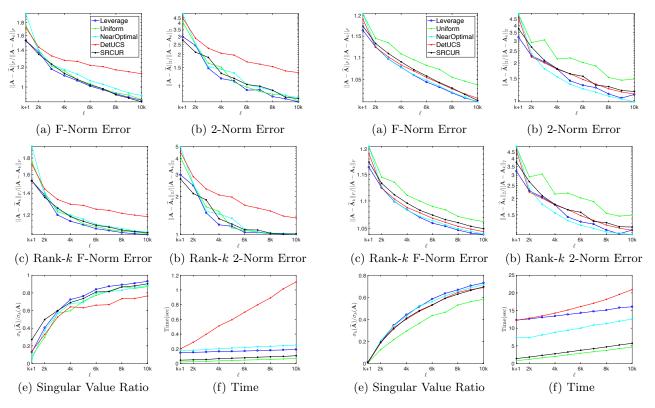


Figure 2: Results of CUR decomposition algorithms comparison on Yale Face B data matrix.

Figure 4: Results of CUR decomposition algorithms comparison on Sido data matrix.

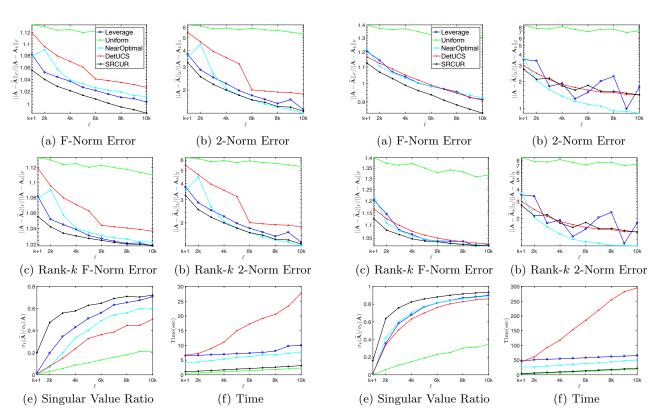


Figure 3: Results of CUR decomposition algorithms comparison on Dexter data matrix.

Figure 5: Results of CUR decomposition algorithms comparison on Reuters data matrix.

Acknowledgements

The work is supported by the National Natural Science Foundation of China (61702327, 61772333, 61632017). Ming Gu is supported by the NSF (1760316). Zhihua Zhang is supported by the National Natural Science Foundation of China (No. 11771002) and Beijing Natural Science Foundation (Z190001). Yong Yu is the corresponding author.

References

- Anderson, D. and Gu, M. (2017). An efficient, sparsity-preserving, online algorithm for low-rank approximation. In *International Conference on Machine Learning*, pages 156–165.
- Anderson, D. G., Du, S., Mahoney, M., Melgaard, C., Wu, K., and Gu, M. (2015). Spectral gap error bounds for improving cur matrix decomposition and the nyström method. In AISTATS.
- Berry, M. W., Pulatova, S. A., and Stewart, G. (2005). Algorithm 844: Computing sparse reduced-rank approximations to sparse matrices. ACM Transactions on Mathematical Software (TOMS), 31(2):252–269.
- Bien, J., Xu, Y., and Mahoney, M. W. (2010). Cur from a sparse optimization viewpoint. In Advances in Neural Information Processing Systems, pages 217–225.
- Boutsidis, C. and Woodruff, D. P. (2017). Optimal cur matrix decompositions. SIAM Journal on Computing, 46(2):543–589.
- Chan, T. F. (1987). Rank revealing qr factorizations. Linear algebra and its applications, 88:67–82.
- Clarkson, K. L. and Woodruff, D. P. (2017). Low-rank approximation and regression in input sparsity time. *Journal of the ACM (JACM)*, 63(6):54.
- Drineas, P., Kannan, R., and Mahoney, M. W. (2006). Fast monte carlo algorithms for matrices iii: Computing a compressed approximate matrix decomposition. *SIAM Journal on Computing*, 36(1):184–206.
- Drineas, P., Mahoney, M. W., and Muthukrishnan, S. (2008). Relative-error cur matrix decompositions. SIAM Journal on Matrix Analysis and Applications, 30(2):844–881.
- Georghiades, A. S., Belhumeur, P. N., and Kriegman, D. J. (2001). From few to many: Illumination cone models for face recognition under variable lighting and pose. *IEEE transactions on pattern analysis and machine intelligence*, 23(6):643–660.
- Gondzio, J. (1992). Stable algorithm for updating dense lu factorization after row or column exchange and row and column addition or deletion. *Optimization*, 23(1):7–26.

- Gu, M. and Eisenstat, S. C. (1996). Efficient algorithms for computing a strong rank-revealing qr factorization. SIAM Journal on Scientific Computing, 17(4):848–869.
- Guyon, I., Gunn, S., Ben-Hur, A., and Dror, G. (2005). Result analysis of the nips 2003 feature selection challenge. In Advances in neural information processing systems, pages 545–552.
- Higham, N. J. and Relton, S. D. (2016). Estimating the largest elements of a matrix. *SIAM Journal on Scientific Computing*, 38(5):C584–C601.
- Johnson, W. B. and Lindenstrauss, J. (1984). Extensions of lipschitz mappings into a hilbert space. *Contemporary mathematics*, 26(189-206):1.
- Mackey, L. W., Jordan, M. I., and Talwalkar, A. (2011). Divide-and-conquer matrix factorization. In Advances in Neural Information Processing Systems, pages 1134–1142.
- Mahoney, M. W. and Drineas, P. (2009). Cur matrix decompositions for improved data analysis. *Proceedings* of the National Academy of Sciences, 106(3):697–702.
- Mahoney, M. W. et al. (2011). Randomized algorithms for matrices and data. Foundations and Trends® in Machine Learning, 3(2):123–224.
- Mahoney, M. W., Maggioni, M., and Drineas, P. (2008). Tensor-cur decompositions for tensor-based data. SIAM Journal on Matrix Analysis and Applications, 30(3):957–987.
- Melgaard, C. and Gu, M. (2015). Gaussian elimination with randomized complete pivoting. arXiv preprint arXiv:1511.08528.
- Miranian, L. and Gu, M. (2003). Strong rank revealing lu factorizations. *Linear algebra and its applications*, 367:1–16.
- Pan, C.-T. (2000). On the existence and computation of rank-revealing lu factorizations. *Linear Algebra and its Applications*, 316(1-3):199–222.
- Song, Z., Woodruff, D. P., and Zhong, P. (2017). Low rank approximation with entrywise l1-norm error. In *Proceedings of the 49th Annual ACM SIGACT Symposium on Theory of Computing*, pages 688–701. ACM.
- Song, Z., Woodruff, D. P., and Zhong, P. (2019). Relative error tensor low rank approximation. In *Proceedings of the Thirtieth Annual ACM-SIAM Symposium on Discrete Algorithms*, pages 2772–2789. Society for Industrial and Applied Mathematics.
- Stewart, G. (1999). Four algorithms for the efficient computation of truncated pivoted qr approximations to a sparse matrix. *Numerische Mathematik*, 83(2):313–323.

- Trefethen, L. N. and Bau III, D. (1997). *Numerical linear algebra*, volume 50. Siam.
- Wang, S. and Zhang, Z. (2012). A scalable cur matrix decomposition algorithm: Lower time complexity and tighter bound. In Advances in Neural Information Processing Systems, pages 647–655.
- Wang, S. and Zhang, Z. (2013). Improving cur matrix decomposition and the nyström approximation via adaptive sampling. *Journal of Machine Learning Research*, 14(1):2729–2769.
- Wang, S., Zhang, Z., and Zhang, T. (2016). Towards more efficient spsd matrix approximation and cur matrix decomposition. *Journal of Machine Learning Research*, 17(210):1–49.
- Williams, C. K. and Seeger, M. (2001). Using the nyström method to speed up kernel machines. In *Advances in neural information processing systems*, pages 682–688.
- Xu, M., Jin, R., and Zhou, Z.-H. (2015). Cur algorithm for partially observed matrices. In *ICML*, pages 1412–1421.