

1 **ACCURACY AND STABILITY OF CUR DECOMPOSITIONS WITH**  
2 **OVERSAMPLING \***

3 TAEJUN PARK<sup>†</sup> AND YUJI NAKATSUKASA<sup>†</sup>

4 **Abstract.** This work investigates the accuracy and numerical stability of CUR decompositions  
5 with oversampling. The CUR decomposition approximates a matrix using a subset of columns and  
6 rows of the matrix. When the number of columns and the rows are the same, the CUR decomposition  
7 can become unstable and less accurate due to the presence of the matrix inverse in the core matrix.  
8 Nevertheless, we demonstrate that the CUR decomposition can be implemented in a numerical stable  
9 manner and illustrate that oversampling, which increases either the number of columns or rows in  
10 the CUR decomposition, can enhance its accuracy and stability. Additionally, this work devises an  
11 algorithm for oversampling motivated by the theory of the CUR decomposition and the cosine-sine  
12 decomposition, whose competitiveness is illustrated through experiments.

13 **Key words.** Low-rank approximation, CUR decomposition, stability analysis, oversampling

14 **MSC codes.** 15A23, 65F55, 65G50

15 **1. Introduction.** The computation of a low-rank approximation to a matrix is  
16 omnipresent in the computational sciences [53]. It has seen increasing popularity due  
17 to its importance in tackling large scale problems. An important aspect of low-rank  
18 approximation is the selection of bases that approximate the span of a matrix’s column  
19 and/or row spaces. In this work, we consider a natural choice of using a subset of rows  
20 and columns of the original matrix as low-rank bases, namely the CUR decomposition.  
21 The CUR decomposition [7, 38, 49], also known as a “matrix skeleton” approximation  
22 [28], is a low-rank approximation that approximates a matrix  $A$  as a product<sup>1</sup>

23 (1.1) 
$$A \underset{m \times n}{\approx} \underset{m \times k}{C} \underset{k \times k}{Z} \underset{k \times n}{R}$$

24 where in MATLAB notation,  $C = A(:, J)$  is a  $k$ -subset of columns of  $A$  with  $J \subseteq$   
25  $\{1, \dots, m\}$  being the column indices and  $R = A(I, :)$  is a  $k$ -subset of the rows of  
26  $A$  with  $I \subseteq \{1, \dots, n\}$  being the row indices. The factors  $C$  and  $R$  are subsets of  
27 the original matrix  $A$  that inherit certain properties of the original matrix, such as  
28 sparsity or nonnegativity, a property that is absent in the truncated SVD. They also  
29 assist with data interpretation by revealing the important rows and columns of  $A$ . The  
30 CUR decomposition is also memory efficient when the entries of  $A$  can be computed  
31 or extracted quickly because we can just store the column and row indices without  
32 explicitly storing  $C$  and  $R$ .

33 There are two common choices for  $Z$ , which we refer to as the core matrix:  $C^\dagger AR^\dagger$   
34 or  $A(I, J)^{-1}$  (or more generally  $A(I, J)^\dagger$ ) [31]. The choice  $Z = C^\dagger AR^\dagger$  minimizes the  
35 Frobenius norm error  $\|A - CZR\|_F$  given the choice of  $C$  and  $R$ . Hence, we refer to  
36 the CUR decomposition with this choice as CURBA (CUR with Best Approximation)  
37 for shorthand in this work. While this choice is more robust, it involves the full ma-  
38 trix  $A$  to form  $C^\dagger AR^\dagger$ , which can be costly, requiring  $O(mnk)$  operations for a dense

---

\*Date: November 20, 2024

**Funding:** TP was supported by the Heilbronn Institute for Mathematical Research. YN is supported by EPSRC grants EP/Y010086/1 and EP/Y030990/1.

<sup>†</sup>Mathematical Institute, University of Oxford, Oxford, OX2 6GG, UK, ([park@maths.ox.ac.uk](mailto:park@maths.ox.ac.uk), [nakatsukasa@maths.ox.ac.uk](mailto:nakatsukasa@maths.ox.ac.uk)).

<sup>1</sup>The notation  $Z$  is temporarily used to denote the core matrix instead of the traditional notation  $U$  as we use  $U$  to denote the intersection of  $C$  and  $R$ , i.e.,  $U = A(I, J)$  later.

39 matrix  $A$ . On the other hand, the choice  $Z = A(I, J)^{-1}$ , which is often called the  
 40 cross approximation [13], is more efficient as we only require the overlapping entries  
 41 of  $C$  and  $R$ , and need not even read the whole matrix  $A$ . We refer to this version  
 42 of the CUR decomposition as CURCA (CUR with Cross Approximation) for short-  
 43 hand. However, this choice has the drawback that  $A(I, J)$  can be (nearly) singular,  
 44 which can lead to poor approximation error. Indeed, the ill-conditioning of  $A(I, J)$   
 45 for CURCA becomes alarming when computing its matrix (pseudo)inverse, a concern  
 46 that has been raised in various papers [20, 40, 49]. For this reason, the related in-  
 47 terpolative decomposition [39] is often advocated to avoid numerical instability with  
 48 CUR, especially CURCA. The aim of this paper is to address the instability of the  
 49 CURCA by demonstrating that it can be implemented in a numerically stable man-  
 50 ner using the  $\epsilon$ -pseudoinverse in the presence of roundoff errors. The  $\epsilon$ -pseudoinverse  
 51 variant takes the  $\epsilon$ -pseudoinverse of  $A(I, J)$  by first truncating the singular values of  
 52  $A(I, J)$  less than  $\epsilon$  before computing its pseudoinverse, which has been explored in a  
 53 related context in [8, 12, 42, 44]. We abbreviate the CURCA with  $\epsilon$ -pseudoinverse as  
 54 the stabilized CURCA (SCURCA for short). Furthermore, we illustrate that oversam-  
 55 pling, which increases either the number of columns or rows in the CURCA, can be  
 56 incorporated to enhance its accuracy and stability. Specifically, we recommend over-  
 57 sampling in proportion to the target rank; see Section 5.3. Throughout this work,  
 58 we concentrate on CURCA and take  $U = A(I, J)$  unless otherwise stated, such that  
 59 CURCA can be written as  $CU^\dagger R$ , but the analysis and the relevant counterparts  
 60 regarding CURBA can be found in Appendix A.<sup>2</sup>

61 *Overview.* In this work, we focus on the CURCA. Section 2 provides background  
 62 information on the CUR decomposition, discusses related and existing works and out-  
 63 lines our contributions. In Section 3, we prove a theoretical bound for the CURCA  
 64 and its  $\epsilon$ -pseudoinverse variant, the SCURCA with oversampling and show that the  
 65 SCURCA can be computed in a numerically stable way in the presence of roundoff  
 66 errors.<sup>3</sup> In Section 4, we study oversampling of indices for the CURCA with the pur-  
 67 pose of improving its accuracy and numerical stability based on the theory developed  
 68 in Section 3. The setting we are interested in is when we have row indices  $I$  and  
 69 column indices  $J$  with  $|I| = |J| = k$ , how should we oversample one of  $I$  or  $J$ . We  
 70 address this question by proposing a deterministic way of oversampling for which the  
 71 rationale for our idea can be explained by the cosine-sine (CS) decomposition. The  
 72 construction could be of use in other contexts where (over)sampling from a subspace  
 73 is desired; a common task that finds use, for example, in model reduction [11] and  
 74 active learning [48]. We conclude with numerical illustrations in Section 5, demon-  
 75 strating the stability of our method in computing the CURCA and the strength of  
 76 oversampling.

77 *Notation.* Throughout, we use  $\|\cdot\|_2$  for the spectral norm or the vector- $\ell_2$  norm,  
 78  $\|\cdot\|_F$  for the Frobenius norm and  $\|\cdot\|$  for any unitarily invariant norms. We use dagger  
 79  $\dagger$  to denote the pseudoinverse of a matrix and  $\llbracket A \rrbracket_k$  to denote the best rank- $k$  approx-  
 80 imation to  $A$  in any unitarily invariant norm, i.e., the approximation derived from  
 81 truncated SVD [35]. Unless specified otherwise,  $\sigma_i(A)$  denotes the  $i$ th largest singular

<sup>2</sup>For the CURBA, there already exists a numerically stable way of computing the CURBA given by the StableCUR algorithm in [1]. Furthermore, the effect of oversampling is less effective as the CURBA,  $A \approx C(C^\dagger A R^\dagger)R$  is a more robust and stable (but expensive) approximation than the CURCA,  $A \approx CU^\dagger R$ ; see Appendix A.

<sup>3</sup>The stability analysis of the plain CURCA,  $CU^\dagger R$  is not covered in this work; in our analysis we require using the  $\epsilon$ -pseudoinverse. However, we observe its numerical stability in practice; see Section 5.1.

82 value of the matrix  $A$ . We use MATLAB style notation for matrices and vectors. For  
 83 example, for the  $k$ th to  $(k + j)$ th columns of a matrix  $A$  we write  $A(:, k : k + j)$ . We  
 84 use  $I$  and  $J$  for the row and the column indices respectively and set  $\Pi_I = I_m(:, I)$   
 85 and  $\Pi_J = I_n(:, J)$  so that  $A(:, J) = A\Pi_J$  and  $A(I, :) = \Pi_I^T A$  for  $A \in \mathbb{R}^{m \times n}$ . Here,  $I_m$   
 86 denotes the  $m \times m$  identity matrix. Lastly, we use  $|I|$  to denote the cardinality of the  
 87 index set  $I$  and define  $[n] := \{1, 2, \dots, n\}$ .

88 **2. Background, Related Work and Contributions.** In the CUR decompo-  
 89 sition, it is important to get a good set of row and column indices as they dictate the  
 90 quality of the low-rank approximation. There are many practical methods, but they  
 91 largely fall into two different categories: pivoting or sampling. For pivoting based  
 92 methods, we use pivoting schemes such as column pivoted QR (CPQR) [27] or LU  
 93 with complete pivoting [51] on  $A$  or the singular vectors of  $A$  to obtain the pivots  
 94 which we then use as row or column indices. For example, these can be used on the  
 95 dominant singular vectors of  $A$  to obtain the pivots and the discrete empirical inter-  
 96 polation method (DEIM) is a popular example [11, 24, 49]. There are also pivoting  
 97 schemes that provide a strong theoretical guarantee [29]. Applying pivoting schemes  
 98 directly on  $A$  can be prohibitive for large matrices, and for this reason a number of  
 99 randomized algorithms based on *sketching* have been proposed [20, 25, 54]. These  
 100 methods “sketch” the original matrix down to a smaller-sized matrix using random-  
 101 ization and perform the pivoting schemes there. See [20] for a comparative study of  
 102 randomized pivoting algorithms. On the other hand, for sampling based methods, we  
 103 sample the column or row indices from some probability distribution obtained from  
 104 certain information about  $A$ . For example, a popular choice is the row norms of the  
 105 dominant singular vectors of  $A$  for leverage scores [23, 38]. There are other sampling  
 106 strategies such as uniform sampling [12], volume sampling [13, 17, 18, 28], DPP sam-  
 107 pling [16], and BSS sampling [5, 6]. In particular, volume sampling leads to a CUR  
 108 approximation that have close-to-optimal error guarantees [13, 56]. Sampling based  
 109 methods can also be prohibitive for large matrices, so a smaller-sized proxy of  $A$  via  
 110 sketching have been proposed [22]. There is also a deterministic sampling method for  
 111 leverage scores [46] and hybrid methods such as L-DEIM [26].

112 The vast majority of CUR decompositions of  $A$  comes with a theoretical guarantee  
 113 that involves the term

114 (2.1) 
$$\|V(J, :)^{-1}\|_2 = \frac{1}{\sigma_{\min}(V(J, :))},$$

115 where  $V \in \mathbb{R}^{n \times k}$  is the  $k$  (approximate) dominant right singular vectors of  $A \in \mathbb{R}^{m \times n}$   
 116 and  $J \subset [n]$  with  $|J| = k$  is the set of column indices; see Theorems 3.3 and A.3. A  
 117 similar term involving the left singular vectors and a row index set is also present. The  
 118 term (2.1) is usually the deciding factor for the accuracy of the CUR decomposition  
 119 and therefore, a majority of the algorithms aim at choosing a set of indices  $J$  that  
 120 would diminish the effect of (2.1). A natural way of improving (2.1) is to *oversample*,  
 121 that is, obtain extra indices  $J_0 \in \mathbb{R}^p$  distinct from  $J$  so that  $V(J \cup J_0, :) \in \mathbb{R}^{(k+p) \times k}$   
 122 becomes a rectangular matrix. By appending more rows to  $V(J, :)$ ,  $V(J \cup J_0, :)$  has  
 123 a larger minimum singular value by the Courant-Fischer min-max theorem, which  
 124 improves the accuracy of the CUR decomposition, as we will see in the later sections.  
 125 The topic of oversampling is not new, in particular, it is shown in [1] that oversampling  
 126 improves the accuracy of CURBA when the singular values decay rapidly. In the  
 127 context of sampling based methods, oversampling has been suggested for theoretical  
 128 guarantees; see e.g. [7, 12, 38]. It is easy to oversample for sampling based methods as

129 we can simply sample more than  $k$  indices from the given probability distribution. In  
 130 contrast, it is often difficult to oversample for pivoting based methods. This is because  
 131 we typically perform pivoting schemes on a smaller-sized surrogate  $X \in \mathbb{R}^{n \times k}$ , for  
 132 example, the sketch of  $A$  and the pivots beyond the first  $k$  indices carry little to no  
 133 information. Therefore, for pivoting based methods, usually a different strategy is  
 134 used for oversampling. In the context of DEIM, various ways of oversampling have  
 135 been suggested [9, 26, 47, 57]. Specifically, Zimmermann and Willcox [57] show that  
 136 oversampling can improve the condition number of oblique projections, and Donello  
 137 et al. [19] proves a bound for DEIM projectors<sup>4</sup> with oversampling, which extends the  
 138 proof without oversampling from [49]. A majority of the aforementioned literature  
 139 focuses on the DEIM projector or the CURBA. However, oversampling can often be  
 140 more effective for the CURCA as we not only improve the bound involving (2.1) (see  
 141 Theorem 3.3), we also make  $A(I, J)$  a rectangular matrix when we oversample either  
 142  $I$  or  $J$  (but not both), which improves the condition number of  $A(I, J)$  making the  
 143 CURCA more accurate when computing the pseudoinverse of  $A(I, J)$ . Therefore, we  
 144 advocate oversampling when possible over the standard choice  $|I| = |J|$ .

145 *Existing work.* The numerical stability of the stabilized CURCA,  $A \approx CU_\epsilon^\dagger R$ ,  
 146 involving the  $\epsilon$ -pseudoinverse has not been studied previously to our knowledge. How-  
 147 ever, there are some related works exploring the idea of  $\epsilon$ -pseudoinverse in the core  
 148 matrix. The  $\epsilon$ -pseudoinverse takes the core matrix  $U = A(I, J)$ , truncates its sin-  
 149 gular values that are less than  $\epsilon$ , and computes the pseudoinverse of the resulting  
 150 matrix. Firstly, Chiu and Demanet [12] study the  $\epsilon$ -pseudoinverse in the context of  
 151 the CURCA with uniform sampling and show that when the matrix is incoherent, the  
 152 algorithm succeeds. However, this paper does not contain stability analysis, and has  
 153 the condition that the matrix needs to be incoherent. The authors in [8, 44] explore  
 154 the  $\epsilon$ -pseudoinverse in the context of the symmetric Nyström method,  $A \approx CU_\epsilon^\dagger C^T$ ,  
 155 applied to symmetric indefinite matrices. However, they show that the  $\epsilon$ -pseudoinverse  
 156 can deteriorate the accuracy of the symmetric Nyström method when applied to sym-  
 157 metric indefinite matrices. Lastly, Nakatsukasa [42] studies the generalized Nyström  
 158 algorithm<sup>5</sup> with  $\epsilon$ -pseudoinverse and oversampling, which provides a numerically sta-  
 159 ble way of computing the generalized Nyström method and proves its stability. The  
 160 paper also demonstrates that oversampling is necessary for a stable approximation in  
 161 the generalized Nyström method.

162 In a related work, Hamm and Huang study the stability of sampling for CUR  
 163 decompositions in [32] and the perturbation bounds of CUR decompositions in [33].  
 164 In [32], they study the problem of determining when a column submatrix of a rank  
 165  $k$  matrix  $A$  also has rank  $k$ . This is important as if the chosen columns or rows of  
 166 a matrix is (numerically) rank-deficient then  $U$  is also (numerically) rank-deficient,  
 167 which may cause numerical issues when computing the pseudoinverse of  $U$ . In [33],  
 168 they derive perturbation bounds for CUR decompositions under the influence of a  
 169 noise matrix. They investigate several variants of the CUR decomposition including  
 170 the  $\epsilon$ -pseudoinverse variant, SCURCA, and provide perturbation bounds in terms of  
 171 the noise matrix. This work is related, but different from our work as they investigate  
 172 the accuracy of the CUR decomposition for the perturbed matrix  $\tilde{A} = A + E$ , while we

<sup>4</sup>DEIM projector is given by  $\mathcal{P}_U = U(\Pi^T U)^\dagger \Pi^T$  where  $U \in \mathbb{R}^{n \times k}$  is an orthonormal matrix and  $\Pi \in \mathbb{R}^{n \times k}$  is a submatrix of the identity matrix that picks  $k$  chosen rows.

<sup>5</sup>The generalized Nyström method is a variant of the CURCA where instead of subsets of rows and columns of  $A$  that approximate the row and column space of  $A$ , we have the sketches,  $Y^T A$  and  $A X$  that approximate the row and column space of  $A$  where  $X$  and  $Y$  are random embeddings. See [42, 52] for more details.

173 investigate the accuracy and stability for a numerical implementation of the CURCA  
 174 in the presence of rounding errors.

175 We review two existing algorithms for oversampling. First, Gidisu and Hochsten-  
 176 bach [26] oversample  $p$  extra indices by choosing the largest  $p$  leverage scores (row  
 177 norms of (approximate) dominant singular vectors) out of the unchosen indices. The  
 178 complexity is  $\mathcal{O}(nk)$  for computing the leverage scores of an orthonormal matrix  
 179  $V \in \mathbb{R}^{n \times k}$ . If we are only given an approximator that is not orthonormal then (ap-  
 180 proximate) orthonormalization needs to be done, usually at a cost of  $\mathcal{O}(nk^2)$ . Second,  
 181 Peherstorfer, Drmač and Gugercin [47], iteratively select  $p$  extra indices for oversam-  
 182 pling to maximize the minimum singular value of  $V(J, :)$  in a greedy fashion. This  
 183 approach, which is called the GappyPOD+E algorithm, uses perturbation bounds on  
 184 the eigenvalues given in [36] to find the next index that maximizes the lower bound  
 185 for the minimum singular value of  $V(J, :)$ . This approach is also a special case of [57].  
 186 The algorithm runs with complexity  $\mathcal{O}((k+p)^2k^2 + nk^2p)$  where  $p$  is the number of  
 187 indices we oversample by. Again, if  $V$  is not orthonormal to begin with, then (ap-  
 188 proximate) orthonormalization needs to be done usually at a cost of  $\mathcal{O}(nk^2)$ . For a  
 189 treatment of approximate orthonormalization, see, for example [3, 4].

190 *Contributions.* Our contribution is twofold: (1) presenting a method for comput-  
 191 ing the CURCA in a numerically stable manner with a theoretical guarantee, and  
 192 (2) advocating the use of oversampling to improve the accuracy and stability of the  
 193 CURCA.

194 Our first and main contribution lies in presenting a method for computing the  
 195 CURCA in a numerically stable manner, accompanied by an analysis that guarantees  
 196 its stability. We show that with the  $\epsilon$ -pseudoinverse in the core matrix, the SCURCA,  
 197  $A \approx CU_\epsilon^\dagger R$ , can be computed in a numerically stable manner by taking the following  
 198 steps. First, we compute each row of  $CU_\epsilon^\dagger$  using a backward stable underdetermined  
 199 linear solver. Then we compute the SCURCA by multiplying  $CU_\epsilon^\dagger$  by  $R$ . See Section  
 200 3.2 for details. In addition to the stability analysis, we also analyze the CURCA and  
 201 its  $\epsilon$ -pseudoinverse variant, SCURCA in exact arithmetic by deriving a relative norm  
 202 bound. While our analysis does not cover the stability of plain CURCA,  $CU^\dagger R$ , we  
 203 observe its stability in practice without the  $\epsilon$ -truncation; see Sections 3.2 and 5.1.

204 Our secondary contribution involves advocating the use of oversampling for the  
 205 CURCA and for providing a deterministic algorithm to oversample row or column in-  
 206 dices. We show that oversampling improves the accuracy and stability of the CURCA  
 207 by providing a theoretical analysis that demonstrates the benefits of oversampling.  
 208 We show that oversampling should be done such that it increases the minimum singu-  
 209 lar value(s) of  $V(J, :)$  where  $V \in \mathbb{R}^{n \times k}$  is the  $k$  (approximate) dominant right singular  
 210 vectors of  $A$  and  $J$  is a set of indices with  $|J| = k$ . Our algorithm is motivated by  
 211 the cosine-sine (CS) decomposition and runs with complexity  $\mathcal{O}(nk^2 + nkp)$  where  $k$   
 212 is the target rank and  $p$  is the oversampling parameter. Note that this complexity  
 213 only refers to the cost of the oversampling process, not the whole CUR process; the  
 214 complexity of the initial CUR process is separate and depends on the method used to  
 215 get the initial set of indices. We show that our algorithm is competitive with existing  
 216 algorithms and in particular, performs similarly to the GappyPOD+E algorithm in [47],  
 217 in which the oversampling process runs with complexity  $\mathcal{O}((k+p)^2k^2 + nk^2p)$ . The  
 218 numerical experiments illustrate that oversampling is recommended.

219 **3. Accuracy and stability of the stabilized CURCA.** In this section, we  
 220 study two topics related to the CURCA,  $A \approx CU^\dagger R$ . We first analyze the accuracy

221 of the CURCA,

$$222 \quad A_{IJ} = A(:, J)A(I, J)^\dagger A(I, :) = A\Pi_J(\Pi_I^T A\Pi_J)^\dagger \Pi_I^T A =: CU^\dagger R$$

223 and its  $\epsilon$ -pseudoinverse variant, the SCURCA,

$$224 \quad A_{IJ}^\epsilon = A(:, J)A(I, J)_\epsilon^\dagger A(I, :) = A\Pi_J(\Pi_I^T A\Pi_J)_\epsilon^\dagger \Pi_I^T A = CU_\epsilon^\dagger R,$$

225 and show that the  $\epsilon$ -truncation in the core matrix compromises the accuracy of the  
 226 CURCA only by  $\epsilon$  times the condition number of the CURCA; see Remark 3.5.  
 227 We then analyze the numerical stability of SCURCA in the presence of roundoff  
 228 errors and show that SCURCA satisfies a similar bound under roundoff errors, mak-  
 229 ing the SCURCA numerically stable as long as the selected rows and columns well-  
 230 approximate the dominant row and column spaces of  $A$ ; see Section 3.2.

231 We begin with some preliminaries: oblique projectors and standard assumptions.  
 232 In the proofs below, we frequently use oblique projectors and their properties. We  
 233 use  $\mathcal{P}_{X,Y} := X(Y^T X)^\dagger Y^T$  where  $X \in \mathbb{R}^{n \times k}$  and  $Y \in \mathbb{R}^{n \times \ell}$  to denote an oblique  
 234 projection onto the column space of  $X$  if  $k \leq \ell$  and  $Y^T X$  has full column rank or  
 235 onto the row space of  $Y^T$  if  $\ell \leq k$  and  $Y^T X$  has full row rank. For example, the CUR  
 236 decomposition  $A_{IJ}$  can be written as

$$237 \quad A_{IJ} = A\Pi_J(\Pi_I^T A\Pi_J)^\dagger \Pi_I^T A = \mathcal{P}_{A\Pi_J, \Pi_I} A = A\mathcal{P}_{\Pi_J, A^T \Pi_I}.$$

238 Some of the important properties of projectors [50] are

- 239 1.  $\mathcal{P}_{X,Y}\mathcal{P}_{X,Y} = \mathcal{P}_{X,Y}$ ,
- 240 2.  $\mathcal{P}_{X,Y}X = X$  if  $Y^T X$  has full column rank,
- 241 3.  $Y^T \mathcal{P}_{X,Y} = Y^T$  if  $Y^T X$  has full row rank,
- 242 4.  $\|\mathcal{P}_{X,Y}\|_2 = \|I - \mathcal{P}_{X,Y}\|_2$  if  $\mathcal{P}_{X,Y} \neq 0, I$ .

243 Lastly, sometimes we can simplify the norm of oblique projectors, which is given by  
 244 the lemma below.

245 LEMMA 3.1. *Let  $\mathcal{P}_{X,Y} \in \mathbb{R}^{n \times n}$  be a projector where  $X \in \mathbb{R}^{n \times k}$ ,  $Y \in \mathbb{R}^{n \times \ell}$  and*  
 246  *$Y^T X \in \mathbb{R}^{\ell \times k}$  all have full column rank (so  $k \leq \ell$ ). Then*

$$247 \quad (3.1) \quad \|\mathcal{P}_{X,Y}\| = \|(Y^T Q_X)^\dagger Y^T\|$$

248 *for any unitarily invariant norm  $\|\cdot\|$  where  $Q_X$  is an orthonormal matrix spanning*  
 249 *the columns of  $X$ .*

250 *Proof.* Let  $X = Q_X R_X$  be the thin QR decomposition of  $X$ . Then

$$251 \quad \|\mathcal{P}_{X,Y}\| = \|X(Y^T X)^\dagger Y^T\| = \|Q_X R_X (Y^T Q_X R_X)^\dagger Y^T\| = \|(Y^T Q_X)^\dagger Y^T\|,$$

252 since  $Y^T Q_X \in \mathbb{R}^{\ell \times k}$  has full column rank and  $R_X \in \mathbb{R}^{k \times k}$  is nonsingular as  $Y^T X$  has  
 253 full column rank.  $\square$

254 Now, we lay out some generic assumptions that hold in our theorems below. The  
 255 assumptions are

256 *Assumption 3.2.*

- 257 1.  $|I| = |J| = k \leq \text{rank}(A)$  where  $k$  is the target rank,
- 258 2.  $A(I, J) \in \mathbb{R}^{k \times k}$  is a non-singular matrix,

259 3.  $X(:, J) \in \mathbb{R}^{k \times k}$  has full row rank, where  $X$  is (any) row space approximator  
 260 of  $A$ .<sup>6</sup>

261 When  $\text{rank}(A) \leq k$  and  $\text{rank}(A(I, J)) = \text{rank}(A)$ , we have  $A = A_{IJ}$  [31]. Since  
 262  $A(I, J)$  is assumed to be non-singular,  $A(:, J)$  and  $A(I, :)$  have full column and row  
 263 rank, respectively. Under Assumption 3.2, by Lemma 3.1,

$$264 \quad \|CU^\dagger\| = \|Q_C(I, :)\dagger\|, \|U^\dagger R\| = \|Q_R(J, :)\dagger\|, \|X(J, :)\dagger X\| = \|Q_X(J, :)\dagger\|$$

265 where  $Q_C$ ,  $Q_R$  and  $Q_X$  are the orthonormal matrices spanning the columns of  $C$ ,  
 266  $R^T$  and  $X^T$ , respectively. We now prove the accuracy of the CURCA and its  $\epsilon$ -  
 267 pseudoinverse variant, the SCURCA.

268 **3.1. Accuracy of CUR and its  $\epsilon$ -pseudoinverse variant.** We prove the  
 269 accuracy of the SCURCA,  $A_{IJ}^\epsilon$  first. The accuracy for the CURCA,  $A_{IJ}$  follows by  
 270 setting  $\epsilon = 0$ . The analysis presented in this section is a key contribution and plays  
 271 an essential role for the stability analysis as we show that the SCURCA satisfies a  
 272 similar bound under roundoff errors, establishing its numerical stability. The bound  
 273 in Theorem 3.3 below somewhat resembles that in [21].

274 **THEOREM 3.3.** *Let  $A \in \mathbb{R}^{m \times n}$  be a matrix,  $I$  and  $J$  be a set of row and column  
 275 indices, respectively, with  $|I| = |J| = k$ ,  $\epsilon > 0$  and  $X \in \mathbb{R}^{k \times n}$  be any row space  
 276 approximator of  $A$ . Then under Assumption 3.2,*

$$277 \quad (3.2) \quad \|A - A_{I \cup I_0, J}^\epsilon\| \leq \|Q_C(I \cup I_0, :)\dagger\|_2 \|Q_X(J, :)^{-1}\|_2 (\|A - AX^\dagger X\| + \|E\|)$$

278 for any unitarily invariant norm  $\|\cdot\|$  where  $I_0$  is a set of extra row indices distinct  
 279 from  $I$  with  $|I_0| = p$ , and  $E \in \mathbb{R}^{k_\epsilon \times k_\epsilon}$  is a matrix satisfying  $\|E\|_2 \leq \epsilon$  where  $k_\epsilon \leq k$  is  
 280 the number of singular values of  $A(I \cup I_0, J)$  smaller than  $\epsilon$ .

281 *Proof.* For shorthand, let  $I_* := I \cup I_0$ . Let the thin SVD of  $A(I_*, J) \in \mathbb{R}^{(k+p) \times k}$   
 282 be  $W\Sigma V^T = [W_1, W_2] \text{diag}(\Sigma_1, \Sigma_2) [V_1, V_2]^T$  where  $\Sigma_2$  contains the singular values of  
 283  $A(I_*, J)$  smaller than  $\epsilon$ . Then

$$284 \quad A_{I_*, J}^\epsilon \Pi_J = A \Pi_J (\Pi_{I_*}^T A \Pi_J)^\dagger_\epsilon \Pi_{I_*}^T A \Pi_J = A \Pi_J V_1 \Sigma_1^{-1} W_1^T W \Sigma V^T \\
 285 \quad = A \Pi_J V_1 V_1^T = A \Pi_J - A \Pi_J V_2 V_2^T.$$

286 Therefore,

$$287 \quad A - A_{I_*, J}^\epsilon = \left( I - A \Pi_J (\Pi_{I_*}^T A \Pi_J)^\dagger_\epsilon \Pi_{I_*}^T \right) A \\
 288 \quad (3.3) \quad = (I - A \Pi_J (\Pi_{I_*}^T A \Pi_J)^\dagger_\epsilon \Pi_{I_*}^T) A (I - \Pi_J (X \Pi_J)^\dagger X) + A \Pi_J V_2 V_2^T (X \Pi_J)^\dagger X.$$

289 Note that  $\mathcal{P}_{A \Pi_J, \Pi_{I_*}}^\epsilon := A \Pi_J (\Pi_{I_*}^T A \Pi_J)^\dagger_\epsilon \Pi_{I_*}^T$  is an oblique projector since

$$290 \quad (\mathcal{P}_{A \Pi_J, \Pi_{I_*}}^\epsilon)^2 = A \Pi_J V_1 \Sigma_1^{-1} W_1^T W \Sigma V^T V_1 \Sigma_1^{-1} W_1^T \Pi_{I_*}^T \\
 291 \quad = A \Pi_J V_1 \Sigma_1^{-1} W_1^T \Pi_{I_*}^T \\
 292 \quad = \mathcal{P}_{A \Pi_J, \Pi_{I_*}}^\epsilon$$

<sup>6</sup>The assumptions and the theorems in this section are stated in terms of row space approximators, but similar assumptions and theorems for column space approximators can be obtained, for example, by considering  $A^T$  instead.

293 and similarly,  $\mathcal{P}_{\Pi_J, X^T} = \Pi_J(X\Pi_J)^\dagger X$  is an oblique projector. Now bounding the first  
294 term of (3.3) gives

$$\begin{aligned}
295 \quad & \left\| (I - \mathcal{P}_{A\Pi_J, \Pi_{I_*}^\epsilon}) A (I - \mathcal{P}_{\Pi_J, X^T}) \right\| \leq \left\| I - \mathcal{P}_{A\Pi_J, \Pi_{I_*}^\epsilon} \right\|_2 \left\| A (I - \mathcal{P}_{\Pi_J, X^T}) \right\| \\
296 \quad & = \left\| \mathcal{P}_{A\Pi_J, \Pi_{I_*}^\epsilon} \right\|_2 \left\| A (I - X^\dagger X) (I - \mathcal{P}_{\Pi_J, X^T}) \right\| \\
297 \quad & \leq \left\| \mathcal{P}_{A\Pi_J, \Pi_{I_*}^\epsilon} \right\|_2 \left\| A (I - X^\dagger X) \right\| \left\| I - \mathcal{P}_{\Pi_J, X^T} \right\|_2 \\
298 \quad & = \left\| \mathcal{P}_{A\Pi_J, \Pi_{I_*}^\epsilon} \right\|_2 \left\| \mathcal{P}_{\Pi_J, X^T} \right\|_2 \left\| A (I - X^\dagger X) \right\|
\end{aligned}$$

299 where in the first equality we used  $X\mathcal{P}_{\Pi_J, X^T} = X\Pi_J(X\Pi_J)^\dagger = I$ , as  $X\Pi_J = X(:$   
300  $, J)$  has full row rank by Assumption 3.2. The first term  $\left\| \mathcal{P}_{A\Pi_J, \Pi_{I_*}^\epsilon} \right\|_2$  in the final  
301 expression can be bounded by letting  $A\Pi_J = Q_C R_C$  be the thin QR decomposition  
302 and noting that  $\Pi_{I_*}^T Q_C$  has full column rank, as

$$\begin{aligned}
303 \quad & \left\| \mathcal{P}_{A\Pi_J, \Pi_{I_*}^\epsilon} \right\|_2 = \left\| Q_C R_C (\Pi_{I_*}^T A\Pi_J)^\dagger \right\|_2 = \left\| R_C (\Pi_{I_*}^T A\Pi_J)^\dagger \right\|_2 \\
304 \quad & = \left\| (\Pi_{I_*}^T Q_C)^\dagger \Pi_{I_*}^T Q_C R_C (\Pi_{I_*}^T A\Pi_J)^\dagger \right\|_2 \\
305 \quad & \leq \left\| (\Pi_{I_*}^T Q_C)^\dagger \right\|_2 \left\| \Pi_{I_*}^T A\Pi_J (\Pi_{I_*}^T A\Pi_J)^\dagger \right\|_2 \\
306 \quad & \leq \left\| (\Pi_{I_*}^T Q_C)^\dagger \right\|_2.
\end{aligned}$$

307 Therefore, using Lemma 3.1 on  $\mathcal{P}_{\Pi_J, X^T}$ , the first term in (3.3) can be bounded as

$$308 \quad \left\| (I - \mathcal{P}_{A\Pi_J, \Pi_{I_*}^\epsilon}) A (I - \mathcal{P}_{\Pi_J, X^T}) \right\| \leq \left\| Q_C(I_*, \cdot)^\dagger \right\|_2 \left\| Q_X(J, \cdot)^{-1} \right\|_2 \left\| A (I - X^\dagger X) \right\|.$$

309 The second term in (3.3) can be bounded using a similar argument as

$$\begin{aligned}
310 \quad & \left\| A\Pi_J V_2 V_2^T (X\Pi_J)^\dagger X \right\| = \left\| Q_C R_C V_2 V_2^T (X\Pi_J)^\dagger X \right\| = \left\| R_C V_2 V_2^T (X\Pi_J)^\dagger X \right\| \\
311 \quad & = \left\| (\Pi_{I_*}^T Q_C)^\dagger \Pi_{I_*}^T Q_C R_C V_2 V_2^T (X\Pi_J)^\dagger X \right\| \\
312 \quad & = \left\| (\Pi_{I_*}^T Q_C)^\dagger W_2 \Sigma_2 V_2^T (X\Pi_J)^\dagger X \right\| \\
313 \quad & \leq \left\| (\Pi_{I_*}^T Q_C)^\dagger \right\|_2 \left\| W_2 \Sigma_2 V_2^T \right\| \left\| (X\Pi_J)^\dagger X \right\|_2 \\
314 \quad & = \left\| Q_C(I_*, \cdot)^\dagger \right\|_2 \left\| Q_X(J, \cdot)^{-1} \right\|_2 \left\| \Sigma_2 \right\|.
\end{aligned}$$

315 Putting everything together and letting  $E = \Sigma_2$ , we get the desired result.  $\square$

316 **COROLLARY 3.4.** *Under the same assumptions as in Theorem 3.3,*

$$317 \quad (3.4) \quad \left\| A - A_{I \cup I_0, J} \right\| \leq \left\| Q_C(I \cup I_0, \cdot)^\dagger \right\|_2 \left\| Q_X(J, \cdot)^{-1} \right\|_2 \left\| A - A X^\dagger X \right\|$$

318 *for any unitarily invariant norm  $\|\cdot\|$ .*

319 *Proof.* Set  $\epsilon = 0$  in Theorem 3.3.  $\square$

320 **Remark 3.5.**

- 321 1. The condition number of the CURCA is  $\kappa = \left\| Q_C(I \cup I_0, \cdot)^\dagger \right\|_2 \left\| Q_X(J, \cdot)^{-1} \right\|_2$ ,  
322 as indicated by (3.2). Theorem 3.3 tells us that the SCURCA,  $A_{I_*^\epsilon, J}$  is worse  
323 than the CURCA by at most a factor  $\kappa\sqrt{k\epsilon}$  in the Frobenius norm.
- 324 2. The bound in Theorem 3.3 and Corollary 3.4 has two factors involving the  
325 (pseudo)inverse, which is in contrast to the CURBA,  $C(C^\dagger A R^\dagger)R$  (see Ap-  
326 pendix A) having only one factor. This makes the CURCA,  $A_{I, J}$  usually worse

327 than the CURBA, which is expected as the CURCA is cheaper to compute.  
 328 However, the CURCA can still be very accurate; see for example [13, 56],  
 329 which establishes the existence of a rank- $r$  CURCA that has error within a  
 330 factor  $r + 1$  of the best rank- $r$  approximation via the truncated SVD. The  
 331 second multiplicative factor in Theorem 3.3 and Corollary 3.4 comes from  
 332 the fact that  $A_{IJ}$  is associated with the oblique projector  $\mathcal{P}_{A\Pi_J, \Pi_I}$  rather  
 333 than the orthogonal projectors  $CC^\dagger$  and  $R^\dagger R$  for the CURBA. Nonetheless,  
 334 the CURCA and the CURBA have comparable accuracy when employed with  
 335 good row and column indices and oversampling, with the CURCA being much  
 336 more computationally efficient.

337 3. The row space approximator  $X \in \mathbb{R}^{k \times n}$  can be chosen in various ways, for  
 338 example, the  $k$ -dominant right singular vectors of  $A$  or the row sketch of  $A$ .  
 339 The bounds in Theorem 3.3 and Corollary 3.4 can both be computed a poste-  
 340 riori when  $X$  can be computed easily. The  $k$ -dominant right singular vectors  
 341 would make the right-most term,  $\|A - AX^\dagger X\|$  in the bound of Theorem  
 342 3.3 and Corollary 3.4, optimal. However, the singular vectors are often too  
 343 expensive to compute.

344 Theorem 3.3 and Corollary 3.4 provide a bound for the SCURCA and the CURCA,  
 345 respectively. It also demonstrates the benefit of oversampling through the extra set  
 346 of row indices  $I_0$ . We obtain a pseudoinverse for  $\|Q_C(I \cup I_0, :)\dagger\|_2$  instead of the ma-  
 347 trix inverse  $\|Q_C(I, :)^{-1}\|_2$  with the former always being smaller. Additionally, since  
 348  $\|Q_C(I \cup I_0, :)\dagger\|_2 = \|A(:, J)A(I_*, J)^\dagger\|_2$  and we want to minimize this quantity, the  
 349 row indices  $I$  and  $I_0$  should be chosen in terms of the already-chosen columns of  
 350  $A$ . This has been suggested and employed in other works such as [14, 20, 54, 55].  
 351 This comes with the benefit that the resulting core matrix  $A(I_*, J)$  will generally be  
 352 better-conditioned, improving the accuracy of the CURCA. For example, consider the  
 353 following  $2 \times 2$  matrix,

$$354 \quad A = \begin{bmatrix} \epsilon & 1 \\ 1 & 0 \end{bmatrix}$$

355 where  $0 < \epsilon < 1$ . For a rank-1 approximation of  $A$ , we need to choose a column and a  
 356 row for the CURCA. If we choose a column and a row separately in the best possible  
 357 way, we would choose the first row and the first column, giving us

$$358 \quad A_{1,1} = \begin{bmatrix} \epsilon \\ 1 \end{bmatrix} \epsilon^{-1} [\epsilon \quad 1] = \begin{bmatrix} \epsilon & 1 \\ 1 & 1/\epsilon \end{bmatrix},$$

359 which is a poor approximation as  $\|A - A_{1,1}\|_F = 1/\epsilon$  can be arbitrary large<sup>7</sup> as  $\epsilon \rightarrow 0$ .  
 360 On the other hand, if we choose a column first and then a row, we choose the first  
 361 column  $[\epsilon, 1]^T$  and the second row as  $\epsilon < 1$ , giving us

$$362 \quad A_{2,1} = \begin{bmatrix} \epsilon \\ 1 \end{bmatrix} 1^{-1} [1 \quad 0] = \begin{bmatrix} \epsilon & 0 \\ 1 & 0 \end{bmatrix},$$

363 which is a reasonable approximation as  $\|A - A_{2,1}\|_F = 1 \approx \sigma_2(A) \approx 1 - \epsilon/2$ .

364 The significance of controlling the  $\|Q_C(I, :)^{-1}\|_2$  term will also be highlighted  
 365 when we analyze the numerical stability of the CUR decomposition in the subsequent  
 366 section.

<sup>7</sup>The accuracy of CURBA, by contrast, is good as long as the chosen columns and rows are good approximators for the range and co-range of  $A$  [20, Remark 1].

367 **3.2. Numerical Stability of the CUR decomposition.** In the absence of  
 368 rounding errors, the error for the CURCA can be bounded by Theorem 3.3 and  
 369 Corollary 3.4. In this section, we derive an error bound that accounts for rounding  
 370 errors.

371 We use the standard model of floating-point arithmetic as in [34, Section 2.2]:

$$372 \quad fl(x \text{ op } y) = (x \text{ op } y)(1 + \delta), \quad |\delta| \leq u$$

373 where  $\text{op} \in \{+, -, *, /\}$  is the basic arithmetic operations and  $u \ll 1$ , the unit round-  
 374 off, is the precision at which the computations are performed. We use  $fl(\cdot)$  and  $\hat{\cdot}$   
 375 to denote the computed value of the expression. We define  $\gamma := p(m, n, k)u$  where  
 376  $p(m, n, k)$  is a low-degree polynomial in  $m, n$  and  $k$ , and use  $\gamma$  to suppress any constant  
 377 factors and terms related to the size of the matrix and the target rank, e.g.,  
 378  $\sqrt{m}, \sqrt{n}$  and  $k$ , but not  $\sigma_i(A)$  or  $1/\epsilon$ . While this may appear as an oversimplification,  
 379 this approach is standard practice in stability analysis; see e.g. [42, 43].

380 In this section, we denote the  $i$ th row of a matrix  $B$  as  $[B]_i$  and we use  $C =$   
 381  $A(:, J)$ ,  $U = A(I_*, J)$  and  $R = A(I, :)$  for shorthand. We assume that the rows are  
 382 oversampled, i.e.,  $I_* = I \cup I_0$  is such that the number of the oversampling indices  
 383  $I_0$  is bounded by a constant times  $k$  where  $k = |I| = |J|$  is the target rank, and  
 384 the truncation parameter  $\epsilon$  satisfies  $\|A\|_2 \gg \epsilon > \gamma \|A\|_2$ . We begin by stating two  
 385 lemmas that will be used in the CURCA stability analysis. Lemma 3.6 proves the  
 386 perturbation bound for the projector  $CU_\epsilon^\dagger \Pi_J$  when only  $C$  and  $U$  get perturbed and  
 387 in Lemma 3.7, we prove that under perturbation on  $C$  and  $U$ ,  $\tilde{C}\tilde{U}_\epsilon^\dagger \Pi_J$  approximately  
 388 projects  $C$  onto itself. The proofs for the two lemmas can be found in Appendix C.

389 LEMMA 3.6. *Under Assumption 3.2, for any  $\Delta C$  and  $\Delta U$ ,*

$$390 \quad (3.5) \quad \|(C + \Delta C)(U + \Delta U)_\epsilon^\dagger\|_2 \leq \|Q_C(I_*, :)\|_2 \left(1 + \frac{1}{\epsilon} \|\Delta U\|_2\right) + \frac{1}{\epsilon} \|\Delta C\|_2.$$

391 LEMMA 3.7. *Under Assumption 3.2, for any  $\Delta C$  and  $\Delta U$ ,*

$$392 \quad (3.6) \quad (C + \Delta C)(U + \Delta U)_\epsilon^\dagger R \Pi_J = C + E_*$$

393 *where*

$$394 \quad \|E_*\|_2 \leq \|(Q_C(I_*, :)\|_2 \left(\epsilon + 2\|\Delta U\|_2 + \frac{1}{\epsilon} \|\Delta U\|_2^2\right) + \|\Delta C\|_2 \left(1 + \frac{\|\Delta U\|_2}{\epsilon}\right).$$

395 We now begin with the stability analysis. Note that the stability depends on the  
 396 specific implementation used. In the forthcoming analysis, we assume the SCURCA  
 397  $A_{I_* J}^\epsilon$  is computed as follows.<sup>8</sup>

- 398 1. Compute the factors  $\hat{C} = fl(A \Pi_J)$ ,  $\hat{U} = fl(\Pi_{I_*}^T A \Pi_J)$  and  $\hat{R} = fl(\Pi_{I_*}^T A)$ .
- 399 2. Solve the (rank-deficient) underdetermined linear systems,

$$400 \quad \hat{s}_i^{(1)} = fl\left((\hat{U}^T)_\epsilon^\dagger [\hat{C}]_i^T\right) \in \mathbb{R}^{k+p}$$

401 for all  $i \in [m]$ , where  $[\hat{C}]_i$  is the  $i$ th row of  $\hat{C}$ .

<sup>8</sup>Note that the stability crucially depends on the implementation. Other implementations are possible, an obvious one being one that computes  $C(U^\dagger R)$  rather than  $(CU^\dagger)R$  as done here. This is seen to work well too, although we do not have a proof.

402 3. Compute matrix-vector multiply  $\hat{s}_i^{(2)} = fl\left(\hat{R}^T \hat{s}_i^{(1)}\right) \in \mathbb{R}^n$ .

403 4. Let  $\left(\hat{s}_i^{(2)}\right)^T$  be the  $i$ th row of the computed CUR decomposition  $fl\left(A_{I_* J}^\epsilon\right)$ .

404 We now analyze each step. The first step is matrix-matrix multiplications with or-  
 405 thormal matrices. Using the forward error bound<sup>9</sup> for matrix-matrix multiplication  
 406 [34, Section 3.5], we obtain the following:

- 407 •  $\hat{C} = C + E_C$  where  $\|E_C\|_2 \leq \gamma \|A\|_2$ ,
- 408 •  $\hat{U} = U + E_U$  where  $\|E_U\|_2 \leq \gamma \|A\|_2$ ,
- 409 •  $\hat{R} = R + E_R$  where  $\|E_R\|_2 \leq \gamma \|A\|_2$ .

410 In many cases, these error matrices  $E_C, E_U$  and  $E_R$  are the zeros matrix as  $C, U$  and  
 411  $R$  are simply submatrices of the original matrix  $A$ .

412 In the second step, we solve the (rank-deficient) underdetermined linear systems  
 413 row by row. The error analysis for the (rank-deficient) underdetermined linear systems  
 414 can be summarized in the following theorem. The proof of Theorem 3.8 can be found  
 415 in Appendix B.

416 THEOREM 3.8. Consider the (rank-deficient) underdetermined linear system,

$$417 \quad (3.7) \quad \min_{x'} \|B_\epsilon x' - b\|_2$$

418 where  $B_\epsilon \in \mathbb{R}^{m \times n}$  ( $m \leq n$ ) is (possibly) rank-deficient ( $\text{rank}(B_\epsilon) \leq m$ ) with singular  
 419 vales larger than  $\epsilon$  and  $b \in \mathbb{R}^m$ . Then assuming  $\epsilon > \gamma \|B_\epsilon\|_2$ , the minimum norm  
 420 solution to (3.7) can be computed in a backward stable manner, i.e., the computed  
 421 solution  $\hat{s}$  satisfies

$$422 \quad (3.8) \quad \hat{s} = (B_\epsilon + E_1)^\dagger (b + E_2)$$

423 where  $\|E_1\|_2 \leq \gamma \|B_\epsilon\|_2$  and  $\|E_2\|_2 \leq \gamma \|b\|_2$ .

424 Theorem 3.8 tells us that the computed solution to a (rank-deficient) underde-  
 425 termined linear system is the exact solution to a slightly perturbed problem. Now,  
 426 using Theorem 3.8, for each  $i \in [m]$  we obtain

$$427 \quad (3.9) \quad \hat{s}_i^{(1)} = \left( (\hat{U}^T)_\epsilon + E_i^{(U)} \right)^\dagger \left( [\hat{C}]_i^T + E_i^{(C)} \right),$$

428 where  $\|E_i^{(U)}\|_2 \leq \gamma \|A\|_2$  and  $\|E_i^{(C)}\|_2 \leq \gamma \|A\|_2$ . It is worth emphasizing that the  
 429 backward errors  $E_i^{(C)}, E_i^{(U)}$  depend on  $i$ . Now since  $\epsilon > \gamma \|A\|_2$  by assumption and  
 430  $\sigma_{\min}(\hat{U}_\epsilon^T + E_i^{(U)}) \geq \epsilon - \gamma \|A\|_2$  by Weyl's inequality, there exists a perturbation  
 431  $E_i \in \mathbb{R}^{k \times (k+p)}$  with  $\|E_i\|_2 \leq \epsilon + \gamma \|A\|_2$ <sup>10</sup> such that

$$432 \quad (3.10) \quad \hat{s}_i^{(1)} = \left( \hat{U}^T + E_i \right)_{\epsilon - \gamma \|A\|_2}^\dagger \left( [\hat{C}]_i^T + E_i^{(C)} \right).$$

433 For shorthand, let  $\hat{S}$  be a matrix with its  $i$ th row equal to  $\left(\hat{s}_i^{(1)}\right)^T$ .

<sup>9</sup>The error is termed the forward error as it indicates how close the computed version  $\hat{C}$  is to the exact version  $C$  [34, Section 1.5].

<sup>10</sup>If  $\hat{U}^T = W_1 \Sigma_1 V_1^T + W_2 \Sigma_2 V_2^T$  is the SVD of  $\hat{U}^T$  where  $\Sigma_2$  contains the singular values of  $\hat{U}^T$  smaller than  $\epsilon$ , then we can take  $E_i = -W_2 \Sigma_2 V_2^T + E_i^{(U)}$  for each  $i$ .

434 In the third step, we compute a matrix-vector product [34, Section 3.5] with  $\hat{R}^T$ ,  
 435 which gives us  $\hat{s}_i^{(2)} = fl\left(\hat{R}^T \hat{s}_i^{(1)}\right) = \hat{R}^T \hat{s}_i^{(1)} + E_{s_i}$  with

$$\begin{aligned}
 436 \quad \|E_{s_i}\|_2 &\leq \gamma \left\| \hat{R}^T \right\|_2 \left\| \hat{s}_i^{(1)} \right\|_2 \\
 437 \quad &\leq \gamma \|A\|_2 \left( \|Q_C(I_*, \cdot)^\dagger\|_2 \left( 1 + \frac{\|E_i\|_2 + \|E_U\|_2}{\epsilon - \gamma \|A\|_2} \right) + \frac{\|E_i^{(C)}\|_2 + \|E_C\|_2}{\epsilon - \gamma \|A\|_2} \right) \\
 438 \quad &\leq \gamma \|A\|_2 \left( \|Q_C(I_*, \cdot)^\dagger\|_2 \left( 1 + \frac{(\epsilon + \gamma \|A\|_2) + \gamma \|A\|_2}{\epsilon - \gamma \|A\|_2} \right) + \frac{\gamma \|A\|_2 + \gamma \|A\|_2}{\epsilon - \gamma \|A\|_2} \right) \\
 439 \quad (3.11) \quad &\leq \gamma \|A\|_2 \|Q_C(I_*, \cdot)^\dagger\|_2,
 \end{aligned}$$

440 where Lemma 3.6 was used in the penultimate line with  $\Delta U = E_U + E_i^T$  and  $\Delta C =$   
 441  $E_C + e_i E_i^{(C)}$  where  $e_i \in \mathbb{R}^m$  is the  $i$ th canonical basis vector. In the last line of (3.11),  
 442 we used the fact that  $\gamma$  suppresses any low-degree polynomial in  $m, n$  and  $k$ , and  
 443  $\epsilon > \gamma \|A\|_2$ .

444 The following shows the expression for  $\hat{s}_i^{(2)}$ ,

$$\begin{aligned}
 445 \quad \hat{s}_i^{(2)} &= \hat{R}^T \hat{s}_i^{(1)} + E_{s_i} = \hat{R}^T \left( \hat{U}^T + E_i \right)_{\epsilon - \gamma \|A\|_2}^\dagger \left( [\hat{C}]_i^T + E_i^{(C)} \right) + E_{s_i} \\
 446 \quad &= (R + E_R)^T \left( U^T + E_U^T + E_i \right)_{\epsilon - \gamma \|A\|_2}^\dagger \left( [C + E_C]_i^T + E_i^{(C)} \right) + E_{s_i}.
 \end{aligned}$$

447 Finally, in the fourth step, we combine  $\hat{s}_i^{(2)} \in \mathbb{R}^n$  into a matrix to form  $\widehat{A}_{I_* J}^\epsilon =$   
 448  $fl\left(A_{I_* J}^\epsilon\right) \in \mathbb{R}^{m \times n}$  by setting the  $i$ th row of  $\widehat{A}_{I_* J}^\epsilon$  to be  $\left(\hat{s}_i^{(2)}\right)^T$ , giving us

$$449 \quad (3.12) \quad fl\left(A_{I_* J}^\epsilon\right) = \hat{S} \hat{R} + E$$

450 where  $E \in \mathbb{R}^{m \times n}$  is a matrix with its  $i$ th row equal to  $E_{s_i}^T$  and satisfies  $\|E\|_2 \leq$   
 451  $\gamma \|A\|_2 \|Q_C(I_*, \cdot)^\dagger\|_2$ , since  $\gamma$  suppresses any low-degree polynomial in  $m$ .

452 We now state the main stability result of the CURCA with the  $\epsilon$ -pseudoinverse,  
 453  $A_{I_* J}^\epsilon$ .

454 **THEOREM 3.9.** *Let  $0 < \epsilon \ll 1$  be a truncation parameter for the pseudoinverse*  
 455 *such that  $\epsilon > \gamma \|A\|_2$ . Suppose that  $A_{I_* J}^\epsilon = A(:, J)A(I_*, J)_\epsilon^\dagger A(I_*, :)$  is computed in the*  
 456 *following order:*

- 457 1. Compute  $C = A(:, J)$ ,  $U = A(I_*, J)$  and  $R = A(I_*, :)$ ,
- 458 2. Compute each row of  $CU_\epsilon^\dagger$  using a backward stable (rank-deficient) underde-
- 459 *termined linear solver,*
- 460 3. Compute  $CU_\epsilon^\dagger$  times  $R$ . Let  $fl\left(A_{I_* J}^\epsilon\right)$  denote the output.

461 Then under Assumption 3.2,

$$\begin{aligned}
 462 \quad (3.13) \quad \|A - fl\left(A_{I_* J}^\epsilon\right)\|_F &\leq 4\sqrt{m} \|Q_C(I_*, \cdot)^\dagger\|_2 \|Q_X(J, \cdot)^\dagger\|_2 \left( \|A(I - X^\dagger X)\|_F + 2\epsilon \right) \\
 463 \quad &\quad + \gamma \|A\|_2 \|Q_C(I_*, \cdot)^\dagger\|_2.
 \end{aligned}$$

464 Theorem 3.9 tells us that the bound for the computed version of the stabilized  
 465 CURCA  $\widehat{A}_{I_* J}^\epsilon$  is at most a factor  $\mathcal{O}(\sqrt{m})$  worse than its exact arithmetic counterpart  
 466  $A_{I_* J}^\epsilon$  plus an error of  $\gamma \|A\|_2 \|Q_C(I_*, \cdot)^\dagger\|_2$ . More specifically, we have

$$467 \quad \|A - fl\left(A_{I_* J}^\epsilon\right)\|_F \leq 4\sqrt{m} (\text{Bound (3.2)}) + \gamma \|A\|_2 \|Q_C(I_*, \cdot)^\dagger\|_2.$$

468 Roughly, this shows that the computed error is in the same order as the error in exact  
 469 arithmetic, up to a factor  $\mathcal{O}(\sqrt{m})$ . The factor  $\mathcal{O}(\sqrt{m})$  is likely an artifact of the proof  
 470 given below; in stability analysis it is common to see bounds with such overestimates,  
 471 and also standard practice to expect to observe much better performance in practice.  
 472 Throughout various parts of the proof, we loosely bound the 2-norm by the Frobenius  
 473 norm, typically because we only have information about the norms of rows or columns  
 474 of a matrix. For example in the proof of Theorem 3.9, we bound  $\|\hat{S}\|_2 \leq \|\hat{S}\|_F$  in  
 475 (3.15), which is likely a pessimistic overestimate.

476 *Proof of Theorem 3.9.* From the above analysis, we have

$$477 \quad (3.12) \quad fl(A_{I_*J}^\epsilon) = \hat{S}\hat{R} + E = \hat{S}R + \hat{S}E_R + E,$$

478 where  $\|E\|_2 \leq \gamma \|A\|_2 \|Q_C(I_*, \cdot)^\dagger\|_2$ . Let us apply Lemma 3.7 to each row of  $\hat{S}R$  to  
 479 obtain

$$480 \quad \hat{S}_i R \Pi_J = \left( [C + E_C]_i + \left( E_i^{(C)} \right)^T \right) (U + E_U + E_i^T)_{\epsilon - \gamma \|A\|_2}^\dagger R \Pi_J = [C]_i + E_i^{(P)},$$

481 where  $\|E_i^{(P)}\|_2 \leq 8\epsilon \|Q_C(I_*, \cdot)^\dagger\|_2$ . Therefore, letting  $E^{(P)}$  be a matrix with  $E_i^{(P)}$  as  
 482 its  $i$ th row, we get

$$483 \quad \hat{S} R \Pi_J = C + E^{(P)},$$

484 where  $\|E^{(P)}\|_F \leq 8\epsilon\sqrt{m} \|Q_C(I_*, \cdot)^\dagger\|_2$ . Now proceeding similarly to the proof of  
 485 Theorem 3.3, we obtain

$$486 \quad A - \hat{S}R = \left( I - \hat{S}\Pi_{I_*}^T \right) A = \left( I - \hat{S}\Pi_{I_*}^T \right) A \left( I - \Pi_J (X\Pi_J)^\dagger X \right) - E^{(P)} (X\Pi_J)^\dagger X \\ 487 \quad (3.14) \quad = \left( I - \hat{S}\Pi_{I_*}^T \right) A \left( I - X^\dagger X \right) \left( I - \Pi_J (X\Pi_J)^\dagger X \right) - E^{(P)} (X\Pi_J)^\dagger X,$$

488 where  $E^{(P)} (X\Pi_J)^\dagger X$  satisfies

$$489 \quad \left\| E^{(P)} (X\Pi_J)^\dagger X \right\|_F \leq \left\| E^{(P)} \right\|_F \left\| (X\Pi_J)^\dagger X \right\|_2 \leq 8\epsilon\sqrt{m} \|Q_C(I_*, \cdot)^\dagger\|_2 \|Q_X(J, \cdot)^\dagger\|_2.$$

490 The first term in (3.14) can be bound by

$$491 \quad \left\| I - \hat{S}\Pi_{I_*}^T \right\|_2 \left\| A \left( I - X^\dagger X \right) \right\|_F \left\| I - \Pi_J (X\Pi_J)^\dagger X \right\|_2 \\ 492 \quad \leq \left( 1 + \left\| \hat{S} \right\|_2 \right) \left\| A \left( I - X^\dagger X \right) \right\|_F \left\| \Pi_J (X\Pi_J)^\dagger X \right\|_2 \\ 493 \quad \leq 4\sqrt{m} \|Q_C(I_*, \cdot)^\dagger\|_2 \|Q_X(J, \cdot)^\dagger\|_2 \left\| A \left( I - X^\dagger X \right) \right\|_F$$

494 where in the final line we used (3.11), noting that

$$495 \quad (3.15) \quad \left\| \hat{S} \right\|_2 \leq \left\| \hat{S} \right\|_F \leq \sqrt{\sum_{i=1}^m \left\| \hat{s}_i^{(1)} \right\|_2^2} \leq 3\sqrt{m} \|Q_C(I_*, \cdot)^\dagger\|_2.$$

496 Using (3.15), we can also bound

$$497 \quad \left\| \hat{S}E_R \right\|_F \leq \left\| \hat{S} \right\|_F \|E_R\|_2 \leq \gamma \|A\|_2 \|Q_C(I_*, \cdot)^\dagger\|_2.$$

498 Finally putting everything together, we obtain

$$\begin{aligned}
499 \quad \|A - fl(A_{I_*J}^\epsilon)\|_F &= \|A - \hat{S}R - \hat{S}E_R - E\|_F \\
500 \quad &\leq 4\sqrt{m} \|Q_C(I_*, \cdot)^\dagger\|_2 \|Q_X(J, \cdot)^\dagger\|_2 (\|A(I - X^\dagger X)\|_F + 2\epsilon) \\
501 \quad &\quad + \gamma \|A\|_2 \|Q_C(I_*, \cdot)^\dagger\|_2. \quad \square
\end{aligned}$$

502 It is natural to wonder what could go wrong without the  $\epsilon$ -pseudoinverse. Two  
503 problems may arise without the  $\epsilon$ -pseudoinverse. First, the matrix  $U$  may not be  
504 numerically full rank, so Theorem 3.8 cannot be used as  $\epsilon > \gamma \|A\|_2$  may no longer  
505 hold. In addition, the two lemmas, Lemma 3.6 and Lemma 3.7 become meaning-  
506 less as both require division by  $\epsilon$ . Without the  $\epsilon$ -pseudoinverse, the error bound can  
507 become uncontrollably large, causing issues in several places in the proof of Theo-  
508 rem 3.9. For example, the bound for the error  $\|E_{s_i}\|_2$  in (3.11) for computing the  
509 matrix-vector product and the bound for  $\|E_i^{(P)}\|_2$  in the proof of Theorem 3.9 may no  
510 longer hold as they can become arbitrarily large without the  $\epsilon$ -pseudoinverse. Never-  
511 theless, in practice we observe stability without the  $\epsilon$ -truncation, so we recommend a  
512 careful implementation of the pseudoinverse (without the  $\epsilon$ -truncation) for practical  
513 purposes.<sup>11</sup> A similar observation, commenting on the role of the  $\epsilon$ -pseudoinverse,  
514 has been mentioned in [42, Section 4.2]. See Section 5.1 for numerical experiments.

515 In Theorem 3.9, one might question how large  $\|Q_C(I_*, \cdot)^\dagger\|_2$  can be, given that it is  
516 part of the added term,  $\gamma \|A\|_2 \|Q_C(I_*, \cdot)^\dagger\|_2$  in Theorem 3.9. This could pose a prob-  
517 lem if  $\|Q_C(I_*, \cdot)^\dagger\|_2$  grows exponentially in  $m$  or  $n$ . However, Theorem 3.9 is enough to  
518 conclude that  $A_{I_*J}^\epsilon$  is numerically stable when the indices are chosen reasonably. This  
519 means that we first choose the column indices  $J$ , and sensibly select the row indices  $I_*$   
520 from the the selected columns  $A(:, J)$  such that  $\|A(:, J)A(I_*, J)^\dagger\|_2 = \|Q_C(I_*, \cdot)^\dagger\|_2$   
521 is bounded by a low-degree polynomial involving  $m, n$  and  $k$ ; see Section 5.2. A simi-  
522 lar approach is also employed in other works such as [14, 20, 54, 55]. For example,  
523 Gu-Eisenstat’s strong rank-revealing QR factorization [29] can be used on the chosen  
524 columns  $A(:, J)$  to obtain  $\|Q_C(I, \cdot)^{-1}\|_2 \leq \sqrt{mk}$ , which can be reduced further with  
525 oversampling. This also highlights the importance of oversampling. In the absence  
526 of oversampling, poorly selected indices can make  $\|Q_C(I, \cdot)^{-1}\|_2$  exponentially large.  
527 Therefore, oversampling can be employed to stabilize the CURCA, which we discuss  
528 further in the following section.

529 **4. Oversampling for the CURCA.** In the previous section, we proved theo-  
530 retical results involving the CURCA (Corollary 3.4), the stabilized CURCA (Theorem  
531 3.3) and the stabilized CURCA in the presence of rounding errors (Theorem 3.9). All  
532 of the results involved an oversampling parameter  $p$  and an extra set of row indices  $I_0$   
533 and bounding  $\|Q_C(I \cup I_0, \cdot)^\dagger\|_2$  was important for the accuracy and stability of the  
534 CURCA. In this section, we discuss oversampling in the context of the CURCA and  
535 devise an algorithm that naturally arises from the discussion.

536 We first describe the setting. Suppose we have obtained the row indices  $I$  and the  
537 column indices  $J$  with  $|I| = |J| = k$  by applying some algorithm, for example, the ones  
538 discussed in the introduction (Section 1), on a row space approximator  $X \in \mathbb{R}^{k \times n}$  of

<sup>11</sup>To be clear, implementing the  $\epsilon$ -truncation does not increase the complexity and can be recom-  
mended for guaranteed stability. We have simply not observed instability without the  $\epsilon$ -truncation.

539  $A \in \mathbb{R}^{m \times n}$ .<sup>12</sup> To have a concrete algorithm in mind for getting the set of indices  $I$   
 540 and  $J$ , we present pivoting on a random sketch [20, 25, 54] below in Algorithm 4.1.

---

**Algorithm 4.1** Pivoting on a random sketch ([20, Algorithm 1])

---

**Require:**  $A \in \mathbb{R}^{m \times n}$  of rank  $r$ , target rank  $k \leq r$  (typically  $k \ll \min\{m, n\}$ )  
**Ensure:** Column indices  $J$  and row indices  $I$  with  $|I| = |J| = k$

```

function  $[I, J] = \text{Rand\_Pivot}(A, k)$ 
1: Draw a random embedding  $\Omega \in \mathbb{R}^{k \times m}$ .
2: Set  $X = \Omega A \in \mathbb{R}^{k \times n}$ , a row sketch of  $A$ 
3: Apply CPQR on  $X$ . Let  $J$  be the  $k$  column pivots.
4: Apply CPQR on  $A(:, J)^T$ . Let  $I$  be the  $k$  row pivots.
    
```

---

541 Algorithm 4.1 is a version of Algorithm 1 from [20], which selects the column  
 542 indices first by applying column pivoted QR (CPQR)<sup>13</sup> on the row sketch  $X = \Omega A$   
 543 and then selects the row indices by applying CPQR on the chosen columns  $A(:, J)^T$ .  
 544 Algorithm 4.1 is an example where we obtain the row indices from the already-chosen  
 545 column indices, which was recommended in the previous section. Here, the row space  
 546 approximator is the row sketch  $X = \Omega A$  and the column space approximator is the  
 547 columns  $A(:, J)$ . A bound for the CURCA (see Corollary 3.4) without oversampling  
 548 is given by

549 (4.1) 
$$\|A - A_{IJ}\|_F \leq \|Q_C(I, :)^{-1}\|_2 \|Q_X(J, :)^{-1}\|_2 \|A - AX^\dagger X\|_F$$

550 where  $A_{IJ} = A(:, J)A(I, J)^\dagger A(I, :)$  is the CURCA and  $Q_C$  and  $Q_X$  are orthonormal  
 551 matrices spanning the columns of  $C$  and  $X^T$  respectively.

552 Many existing algorithms focus on minimizing the first two terms on the right-  
 553 hand side of (4.1) as they control the accuracy of the CURCA. In the CURCA and  
 554 the other CUR decompositions such as the CURBA,  $C(C^\dagger AR^\dagger)R$ , we often take  
 555  $|I| = |J|$ , however, without increasing the overall rank of the approximation, we  
 556 can oversample either  $I$  or  $J$ .<sup>14</sup> Suppose we oversample the rows to  $I_* := I \cup I_0$   
 557 where  $|I_0| = p$  is an extra set of row indices for oversampling. Then the first term  
 558 of the bound changes from  $\|Q_C(I, :)^{-1}\|_2$  to  $\|Q_C(I_*, :)^{-1}\|_2$  (see Corollary 3.3). Now  
 559  $\sigma_{\min}(Q_C(I_*, :)) \geq \sigma_{\min}(Q_C(I, :))$  by the Courant-Fischer min-max theorem, which  
 560 improves the bound in (4.1) as  $\|Q_C(I_*, :)^{-1}\|_2 \leq \|Q_C(I, :)^{-1}\|_2$ . Now, to maximize  
 561 the effect of oversampling, we ought to find unchosen indices that enrich the trailing  
 562 singular subspace of  $Q_C(I, :)$ , which in turn increases the minimum singular value(s)  
 563 of  $Q_C(I, :)$ . It turns out that we can achieve this by projecting  $Q_C([m] - I, :)$  onto  
 564 the trailing singular subspace of  $Q_C(I, :)$  and use a good row selection algorithm to  
 565 choose  $p$  extra rows. Before stating the algorithm, we first motivate our rationale  
 566 behind our approach using the cosine-sine (CS) decomposition.

567 Let  $Q \in \mathbb{R}^{n \times k}$  be any orthonormal matrix with partition

568 (4.2) 
$$Q = \begin{bmatrix} Q_1 \\ Q_2 \end{bmatrix} \begin{matrix} k \\ n_1 \\ n_2 \end{matrix}$$

---

<sup>12</sup>A similar version for column space approximator can also be devised by considering  $A^T$  instead. In the case when  $X$  and  $Y$  are not available, for example, when the initial set of indices were obtained using uniform sampling,  $R = A(I, :)$  and  $C = A(:, J)$  can be used as the row space approximator and the column space approximator, respectively.

<sup>13</sup>LU with partial pivoting is also effective in practice [20].

<sup>14</sup>Oversampling both the row indices  $I$  and the column indices  $J$  independently is not recommended. See Section 5.3 for a further discussion.

569 where  $n_1, n_2 \geq k$  with  $n_1 + n_2 = n$ . Then the CS decomposition [45] gives us

$$570 \quad (4.3) \quad Q_1 = U_1 C V^T \text{ and } Q_2 = U_2 S V^T,$$

571 where  $U_1 \in \mathbb{R}^{n_1 \times k}, U_2 \in \mathbb{R}^{n_2 \times k}, V \in \mathbb{R}^{k \times k}$  are matrices with orthonormal columns  
 572 and  $C = \text{diag}(c_1, c_2, \dots, c_k), S = \text{diag}(s_1, s_2, \dots, s_k)$  are diagonal matrices satisfying  
 573  $c_i^2 + s_i^2 = 1$  for all  $i$ . In our context, we can view  $Q_C(I, :)$  as  $Q_1$  and  $Q_C([m] - I, :)$   
 574 as  $Q_2$ . Now assume, without loss of generality, that the  $c_i$ 's are in non-increasing  
 575 order so the  $s_i$ 's are in non-decreasing order. Then in order to increase the minimum  
 576 singular value of  $Q_1$ , we could add the rows of  $Q_2$  that contribute the most to the  
 577 trailing right singular subspace of  $Q_1$ , i.e.,  $V_{-p} = V(:, k-p+1 : k) \in \mathbb{R}^{k \times p}$ . By (4.3),  
 578  $V_{-p}$  is also the dominant right singular subspace of  $Q_2$ . When we add the rows of  
 579  $Q_2$  that lies in the subspace spanned by  $V_{-p}$  to  $Q_1$ , we can increase the minimum  
 580 singular value(s) of  $Q_1$ , i.e., increase  $c_k$  or the last few  $c_i$ 's. Therefore we can apply  
 581 any algorithm (such as those discussed in Section 1) that finds good row indices on  
 582  $Q_2 V_{-p} \in \mathbb{R}^{n_2 \times p}$  and append them to  $Q_1$  to increase the minimum singular value of  
 583  $Q_1$ .

584 If the algorithm requires the dominant left singular vectors of  $Q_2 V_{-p}$  such as in  
 585 DEIM or leverage scores sampling, we can simply scale the columns of  $Q_2 V_{-p}$  using  
 586 the singular values of  $Q_1, C = \text{diag}(c_1, \dots, c_k)$  because

$$587 \quad Q_2 V_{-p} = U_{2,-p} \text{diag}(s_{k-p+1}, \dots, s_k) = U_{2,-p} \sqrt{1 - \text{diag}(c_{k-p+1}^2, \dots, c_k^2)}$$

588 where  $U_{2,-p} = U_2(:, k-p+1 : k)$ .

589 In light of the observation made above, we propose the following algorithm (Al-  
 590 gorithm 4.2) for oversampling indices. For simplicity, we use column pivoted QR  
 591 (CPQR) on  $Q_2 V_{-p}$  in Algorithm 4.2 to obtain a good set of oversampling indices.  
 592 However, any algorithm that finds a good set of rows can replace line 4 of Algorithm  
 593 4.2.

---

#### Algorithm 4.2 Oversampling indices

---

**Require:** A full column rank matrix  $B \in \mathbb{R}^{n \times k}$  with  $n \geq k$ , an index set  $I$  with  $|I| = k$  and an  
oversampling parameter  $p \leq k$

**Ensure:** Extra indices  $I_0$  with  $|I_0| = p$

```

function  $I_0 = \text{OS}(B, I, p)$ 
1:  $[Q_B, \sim] = \text{qr}(B, 0)$ , ▷ Skip this step if  $B$  is already orthonormal.
2:  $[\sim, \sim, V] = \text{svd}(Q_B(I, :))$ ,
3: Set  $V_{-p} = V(:, k-p+1 : k)$ , the trailing  $p$  right singular vectors of  $Q_B(I, :)$ .
4: Apply CPQR on  $(Q_B([m] - I, :) V_{-p})^T$ . Let  $I_0$  be the extra  $p$  indices for oversampling.

```

---

594 Given a full column rank matrix  $B$ , an index set  $I$  and an oversampling parameter  
 595  $p(\leq k)$ , Algorithm 4.2 finds the extra indices for oversampling by projecting  $Q_B$  onto  
 596 the unchosen indices  $[m] - I$  from the left and the trailing  $p$  right singular vectors  
 597 of  $Q_B(I, :)$  from the right and performing CPQR to obtain the extra indices  $I_0$  with  
 598  $|I_0| = p$ . When  $p > k$ , we can iterate Algorithm 4.2 to obtain at most  $k$  oversampling  
 599 indices at each iteration. We can also devise a version of Algorithm 4.2 with a tolerance  
 600 parameter  $0 < \epsilon < 1$  rather than taking  $p$  as input; for example  $\epsilon = \sqrt{k/n}$  may be a  
 601 reasonable choice (given that singular values of a  $O(k) \times k$  submatrix of an  $n \times n$  Haar  
 602 distributed orthogonal matrix are of this order [41]). This version uses projection onto  
 603 the trailing right singular subspace of  $Q_B(I, :)$  corresponding to the singular values  
 604 that are less than  $\epsilon$ . While this version does not guarantee  $\sigma_{\min}(Q_B(I \cup I_0, :)) \geq \epsilon$ ,

605 it can be applied iteratively to achieve this. To guarantee  $\sigma_{\min}(Q_B(I \cup I_0, :)) \geq \epsilon$ ,  
 606 one can alternatively use the GappyPOD+E oversampling algorithm [47]. However, this  
 607 approach comes with a higher computational cost. The complexity of Algorithm 4.2  
 608 is  $\mathcal{O}(nk^2)$  where the dominant cost comes from taking the QR decomposition of  $B$   
 609 (line 1). If  $B$  were orthonormal to begin with, then the dominant cost comes from  
 610 forming  $Q_B([m] - I, :)$ , which costs  $\mathcal{O}(nkp)$ .

611 In the analysis for the CURCA in Section 3, oversampling played two roles: (i)  
 612 improve the accuracy of the CURCA (see Theorem 3.3, Corollary 3.4, Theorem 3.9),  
 613 and (ii) improve the stability of the CURCA in the presence of roundoff errors (see  
 614 Theorem 3.9). Therefore, it is important to oversample, especially if the original set  
 615 of indices are not good. For example, if the core matrix  $A(I, J)$  is (nearly) singular.  
 616 Oversampling should be done in such a way that the term  $\|Q_C(I, :)^{-1}\|_2$  is reduced  
 617 further by adding an extra set of indices  $I_0$  that lie in the trailing singular subspace  
 618 of  $Q_C(I, :)$ . Algorithm 4.2 achieves this by picking good unchosen indices from the  
 619 trailing singular subspace of  $Q_C$ . We further highlight the significance of oversampling  
 620 through numerical illustrations in Section 5.

621 **5. Numerical Illustration.** In this section, we illustrate the concepts discussed  
 622 in the previous sections through numerical experiments. We first discuss implemen-  
 623 tation details of (S)CURCA in MATLAB. Then we show that, without loss of gen-  
 624 erality, after selecting the rows first, the columns should be chosen based on those  
 625 rows, i.e., the rows and the columns for the CURCA should not be chosen indepen-  
 626 dently. We show that oversampling can improve the quality when the indices are  
 627 poorly selected. Additionally, we also illustrate the effectiveness of the oversampling  
 628 algorithm, Algorithm 4.2, for the CURCA and show its competitiveness against some  
 629 existing methods. In all the experiments, the best rank- $k$  approximation error using  
 630 the truncated SVD (TSVD) is used as reference. The experiments were conducted in  
 631 MATLAB version 2021a using double precision arithmetic.

632 **5.1. Implementation of (S)CURCA.** The main concern in the computation  
 633 of the CURCA is that (i) the representation of the CURCA typically involves three  
 634 (highly) ill-conditioned matrix and (ii) we take the pseudoinverse of a (highly) ill-  
 635 conditioned matrix. When the original matrix  $A$  is low-rank and we have a good  
 636 low-rank approximation of  $A$ , then we expect  $C, U$  and  $R$  to be highly ill-conditioned.  
 637 We test four possible implementations in MATLAB,

- 638 1.  $A_{IJ}^{(1)} = (A(:, J)/A(I, J)) * A(I, :)$ ,
- 639 2.  $A_{IJ}^{(2)} = A(:, J) * (V/S * W') * A(I, :)$  where  $[W, S, V] = \text{svd}(U, 'econ')$ ,
- 640 3.  $A_{IJ}^{(3)} = (A(:, J) * V/S) * (W' * A(I, :))$  where  $[W, S, V] = \text{svd}(U, 'econ')$ ,
- 641 4.  $A_{IJ}^{(\epsilon)} = (A(:, J) * V(:, II)/S(II, II)) * (W(:, II))' * A(I, :)$   
 642 where  $[W, S, V] = \text{svd}(U, 'econ')$  and  $II = (\text{diag}(S) > \epsilon)$  with  $\epsilon = 10^{-15}$ .

643 The third implementation,  $A_{IJ}^{(3)}$  is the suggested implementation of the CURCA.  
 644 For guaranteed stability, we suggest the fourth implementation,  $A_{IJ}^{(\epsilon)}$  with the  $\epsilon$ -  
 645 pseudoinverse. However we have not observed instability without the  $\epsilon$ -pseudoinverse  
 646 using the third implementation. We use two test matrices: (1)  $1000 \times 1000$  test matrix  
 647 generated using the MATLAB command, `randn(1000, 30) * randn(30, 1000)` and (2)  
 648  $1374 \times 1374$  matrix named `nnc1374` from the SuiteSparse Matrix Collection [15]. We  
 649 choose the rows and columns by first choosing the column indices  $J$  using column  
 650 pivoted QR (CPQR) on  $A$  and then using CPQR on the chosen columns  $A(:, J)$   
 651 to get the row indices  $I$ .

652 The results are depicted in Figure 1. First, we notice that the third implemen-

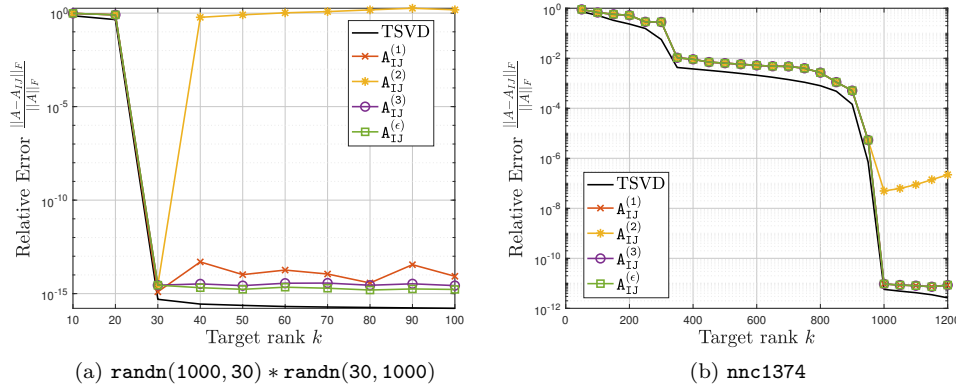


FIG. 1. Tests for different implementations of the CURCA. The third and fourth implementation performs stably. We recommend the third implementation.

653 tation,  $A_{IJ}^{(3)}$ , which is the one we suggest, yields stable approximation throughout the  
 654 experiment. In particular, the third implementation performs similarly to the fourth  
 655 implementation with the  $\epsilon$ -pseudoinverse. On the other hand, the second implementa-  
 656 tion,  $A_{IJ}^{(2)}$  suffers from numerical errors as we lose accuracy when the singular values  
 657 decay extremely rapidly. The reason is because we explicitly form the pseudoinverse  
 658 of the core matrix  $A(I, J)$ , which is highly ill-conditioned, an observation also noted in  
 659 [42]. Therefore, it is important to not form the pseudoinverse of  $U = A(I, J)$  explicitly  
 660 in the CURCA. Note also that the order in which the factors are multiplied is essen-  
 661 tial for numerical stability, as the second and third implementations only differ by the  
 662 order in which the factors are computed. The first implementation,  $A_{IJ}^{(1)}$  may lose 1 or  
 663 2 orders of magnitude in Figure 1a once the target rank becomes larger than the rank  
 664 of the test matrix and appear to be less stable than the third implementation. The  
 665 slash commands ( $/$ ,  $\backslash$ ) in MATLAB should be used with caution for underdetermined  
 666 problems, as its backslash command applied to numerically rank-deficient underde-  
 667 termined problems output a sparse solution based on a pivoting strategy [2, § 2.4],  
 668 which can differ significantly from the minimum-norm solution and may not satisfy  
 669 the assumptions in our analysis. Note that using the `pinv` command with tolerance  
 670 parameter 0 for the CURCA in MATLAB, i.e.,  $A(:, J) * \text{pinv}(A(I, J), 0) * A(I, :)$ , is  
 671 equivalent to the second implementation, which can be unstable. For the rest of the  
 672 numerical experiments, we use the third implementation

$$673 \quad (5.1) \quad (A(:, J) * V/S) * (W' * A(I, :)) \text{ where } [W, S, V] = \text{svd}(U, \text{'econ'}).$$

674 The fourth implementation computes the stabilized CURCA and is implemented  
 675 in such a way that the analysis in Section 3.2 hold. A less expensive alternative  
 676 is to use a rank-revealing QR factorization and truncate the bottom-right corner  
 677 of the upper triangular factor and the relevant columns of the orthonormal factor  
 678 corresponding to the diagonal elements less than  $\epsilon$ . In the rank-revealing QR, the  
 679 diagonal elements give a good approximation to the singular values; see for example  
 680 [10, 29]. A possible workaround without the  $\epsilon$ -pseudoinverse is also discussed in [42]  
 681 where we perturb the core matrix  $A(I, J)$  by a small noise matrix such that the  
 682 singular values of  $A(I, J)$  are all larger than the unit roundoff.

683 **5.2. Importance of not choosing rows and columns independently.** In  
 684 this section, we demonstrate the importance of choosing the rows and columns that  
 685 are dependent on one another in the CURCA. If the rows and columns are chosen  
 686 independently of each other, the matrix  $U = A(I, J)$  can be (nearly) singular. In  
 687 such a scenario, we show that a sufficient amount of oversampling can remedy this  
 688 problem. We use two test matrices: (1) synthetic matrix generated using

689 (5.2) 
$$A = \begin{bmatrix} 10^{-10} \cdot \text{randn}(50, 50) & \text{randn}(50, 950) \\ \text{randn}(950, 50) & 0 \end{bmatrix} \in \mathbb{R}^{1000 \times 1000},$$

690 and (2) `nnc1374` matrix from the SuiteSparse Matrix Collection used in the previous  
 691 section. We test the following four cases:

- 692 1. Choose columns and rows independently by applying CPQR on  $A$  and  $A^T$   
 693 respectively,
- 694 2. Choose columns and rows dependently by applying CPQR on  $A$  and then  
 695 applying CPQR on the chosen columns  $A(:, J)$  to obtain the rows,
- 696 3. Apply Case 1 and additionally do row oversampling using Algorithm 4.2 with  
 697  $p = k$ ,
- 698 4. Apply Case 2 and additionally do row oversampling using Algorithm 4.2 with  
 699  $p = k$ ,

700 where  $k$  is the target rank.

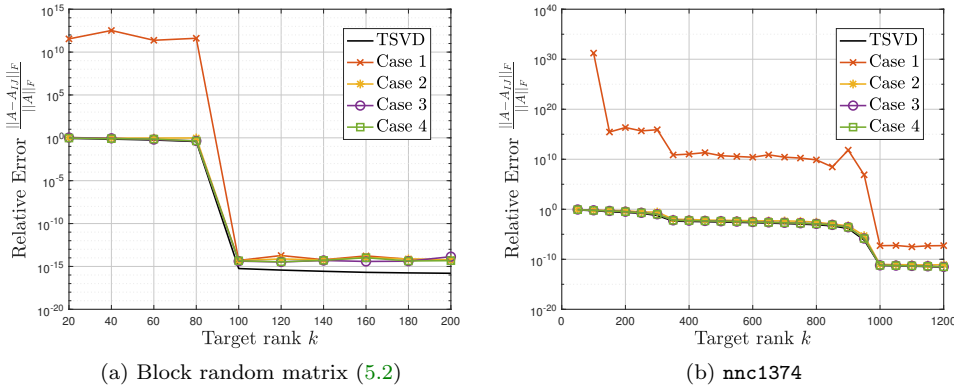


FIG. 2. Relationship between the rows and columns in the CURCA. In Case 1, the rows and columns are chosen independently from one another and in Case 2, we select the columns first and then the rows were computed from the selected columns. Cases 3 and 4 correspond to Cases 1 and 2 with row oversampling ( $p = k$ ), respectively. When the rows and columns are chosen independently (Case 1), the resulting approximation can be catastrophic.

701 In Figure 2, we show the importance of choosing the columns and rows in the  
 702 CURCA. When the columns and rows are chosen independently (Case 1), the resulting  
 703 CURCA can be catastrophically poor. In Figure 2a, poor approximation happens  
 704 because when we choose the columns and rows independently, we choose the first 50  
 705 rows and columns, as they are the most important, implying that the core matrix is  
 706 the small (1, 1)-block of  $A$ . Since the chosen rows and columns are  $\mathcal{O}(1)$  and their

707 intersection is  $\mathcal{O}(10^{-10})$ , the CURCA becomes inaccurate.<sup>15</sup> A similar issue also arises  
 708 in Figure 2b where the approximation becomes noticeably unreliable throughout; see  
 709 also the  $2 \times 2$  example at the end of Section 3.1. However, when we oversample  
 710 sufficiently in Case 3, we obtain a good CUR approximation. When we choose the  
 711 rows and columns dependently in Cases 2 and 4, we have good stable approximations.  
 712 Therefore, it is highly recommended to choose the rows and columns dependently.  
 713 However, if that is not possible, a sufficient amount of oversampling is recommended.

714 **5.3. Oversampling algorithm comparison.** In this section, we illustrate ad-  
 715 vantages of oversampling through numerical experiments. In all experiments, unless  
 716 stated otherwise, we use pivoting on a random sketch (Algorithm 4.1) with column  
 717 pivoted QR [20, 54] and the Gaussian sketch to get the initial set of  $k$  row and column  
 718 indices, where  $k$  is the target rank. We then obtain the oversampling indices using  
 719 the following algorithms:

- 720 1. OS + P: Algorithm 4.2,
- 721 2. OS + L: Choose  $p$  extra indices corresponding to the largest  $p$  leverage scores  
 722 out of the unchosen indices as in [26],
- 723 3. OS + E: GappyPOD + E algorithm in [47].

724 We set the number of oversampling indices to be  $p = 0$ ,  $p = 10$  and  $p = 0.5k$ .

725 We consider three different classes of test matrices, which are summarized below:

- 726 1. CIFAR10: The CIFAR-10 training set [37] consists of 60000 images of size  
 727  $32 \times 32 \times 3$ . We choose 10000 random images, each flattened to a vector, and  
 728 treat it as a  $10000 \times 3072$  data matrix.
- 729 2. YaleFace64x64: Yale face is a full-rank dense matrix of size  $165 \times 4096$   
 730 consisting of 165 face images each of size  $64 \times 64$ . The flattened image vectors  
 731 are centered and normalized such that the mean is zero and the entries lie  
 732 within  $[-1, 1]$ .
- 733 3. SNN: Random sparse non-negative matrices are test matrices used in [49, 54]  
 734 that is given by,

$$735 \quad \text{SNN} = \sum_{j=1}^r s_j x_j y_j^T,$$

736 where  $s_1 \geq \dots \geq s_r > 0$  and  $x_j \in \mathbb{R}^m, y_j \in \mathbb{R}^n$  are random sparse vectors  
 737 with non-negative entries. We take  $m = 100000, n = 300$  with

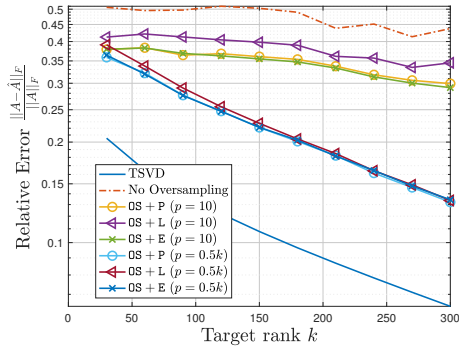
$$738 \quad \sum_{j=1}^{50} \frac{2}{j} x_j y_j^T + \sum_{j=51}^{300} \frac{1}{j} x_j y_j^T,$$

739 where the sparse vectors  $x_j$ 's and  $y_j$ 's are computed in MATLAB using the  
 740 command `sprand(m, 1, 0.025)` and `sprand(n, 1, 0.025)`, respectively.

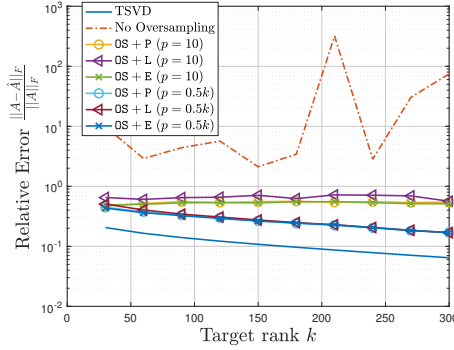
741 The results are depicted in Figure 3. We compare different oversampling algo-  
 742 rithms for the three test matrices. We use pivoting on a random sketch (Algorithm  
 743 4.1) to obtain the initial set of indices except for Figure 3b where we use uniform

---

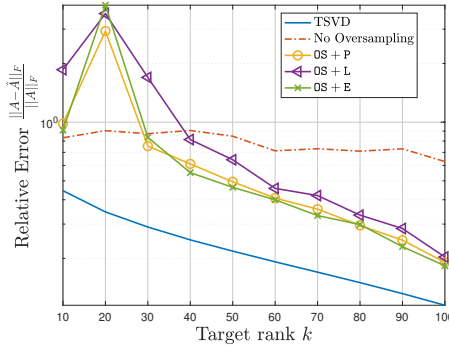
<sup>15</sup>A similar problem may also arise in random sampling where sampling rows and columns indepen-  
 dently may lead to a poor CUR decomposition with high probability. For example, when  $A = \begin{bmatrix} 0 & 1 \\ 1 & 0 \end{bmatrix}$   
 and the target rank is 1, the probability of sampling zero as the core matrix using leverage scores or  
 column norms is 0.5 when the rows and columns are sampled independently. Independent sampling  
 leads to a poor CUR approximation 50% of the time. However, if we sample them dependently, we  
 choose 1 as the core matrix with probability 1, leading to a sensible CUR approximation.



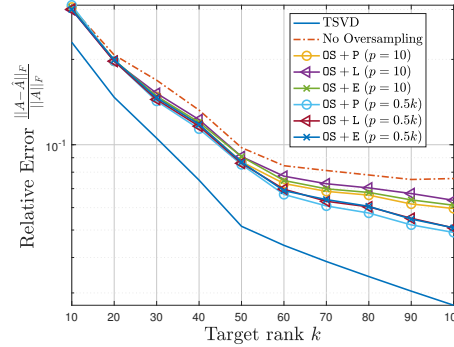
(a) CIFAR10 dataset with pivoting on a random sketch with row oversampling.



(b) CIFAR10 dataset with uniform sampling with row oversampling.



(c) YaleFace64x64 dataset with pivoting on a random sketch with row ( $p = 0.5k$ ) and column ( $p = 10$ ) oversampling.



(d) SNN dataset with pivoting on a random sketch with row oversampling

FIG. 3. Comparison of different oversampling methods for various test matrices. We use pivoting on a random sketch to obtain the initial set of indices except for Figure 3b where we use uniform sampling. Then we oversample the row indices using the three oversampling algorithms, OS+P, OS+L and OS+E. We also oversample column indices in Figure 3c to demonstrate that oversampling both row and column indices can be harmful.

744 sampling; pivoting on a random sketch usually results in a good set of indices, while  
 745 uniform sampling can give a poor set of indices. Then we oversample the row indices  
 746 using the three oversampling algorithms, OS+P, OS+L and OS+E. We also oversample  
 747 column indices in Figure 3c. For Figure 3c,  $p = 0.5k$  for row oversampling and  $p = 10$   
 748 for column oversampling.

749 We begin with the top two plots involving the CIFAR10 dataset, Figures 3a and  
 750 3b. We observe that as the oversampling parameter increases, the accuracy improves.  
 751 The different oversampling techniques yield a similar result with OS+L being usually  
 752 slightly worse. In Figures 3a and 3b, we observe that oversampling plays a big role  
 753 for the accuracy of the CURCA. In Figure 3a,  $p = 10$  improves the accuracy slightly,  
 754 but when the oversampling parameter  $p$  becomes proportional to the target rank,  
 755 we make further progress and we approximately capture the singular value decay  
 756 rate. In Figure 3b, when the initial set of indices are worse, as it can be when we  
 757 perform uniform sampling, we observe that oversampling makes the approximation

758 more robust even when the CURCA without oversampling yields unstable results.  
 759 Therefore, oversampling helps the CURCA yield more accurate and stable results.

760 In Figure 3c, using the `YaleFace64x64` dataset, we show what happens when we  
 761 oversample both row and column indices. In this experiment, we oversample the row  
 762 indices by  $p = 0.5k$  and the column indices by  $p = 10$ . In Figure 3c, the CURCA  
 763 can get worse with oversampling and the approximation may become unstable when  
 764 both column and row indices are oversampled. This is caused by the core matrix  
 765  $A(I \cup I_0, J \cup J_0)$  underestimating the singular values of  $A$  as oversampling in only one  
 766 of row or column indices improve the condition number of the core matrix, but when  
 767 we oversample both row and column indices, condition number of the core matrix  
 768 can become worse. This is similar to the phenomenon happening in Section 5.2,  
 769 where the CURCA can yield poor accuracy when we choose the rows and columns  
 770 independently. Therefore, we advocate oversampling in only one of row or column  
 771 indices for the CURCA. Note that if we do not oversample column indices in Figure  
 772 3c, so that we only oversample rows for the `YaleFace64x64` dataset, the relative error  
 773 decreases steadily, as observed in the other figures in Figure 3. In Figure 3d, we use  
 774 the `SNN` dataset. The effect of oversampling is less immediate in this example, but the  
 775 accuracy still improves and the decay rate is more in line with the spectral decay.

776 Lastly, in all the experiments in Figure 3, we observe that our algorithm for  
 777 oversampling OS + P (Algorithm 4.2) is competitive with OS + E with OS + L usually  
 778 being slightly worse. Therefore, Algorithm 4.2, which runs with complexity  $\mathcal{O}(nk^2 +$   
 779  $nkp)$  is competitive with OS+E, which run with complexity  $\mathcal{O}(nk^2p + k^4)$ .

780 **6. Conclusion.** In this work, we study the accuracy and stability of the CUR  
 781 decomposition with oversampling. We prove the relative norm bounds for the CUR  
 782 decomposition,  $A \approx CU^\dagger R$ , and its stabilized version,  $A \approx CU_\epsilon^\dagger R$ . We further show  
 783 that the stabilized version satisfies a similar relative norm bound in the presence of  
 784 roundoff errors under the assumption that the rows and columns for the CUR decom-  
 785 position are chosen reasonably. This means that the rows should be selected based on  
 786 the chosen columns or vice versa; see Section 5.2. This aims to reduce the quantity  
 787  $\|Q_C(I, \cdot)^{-1}\|_2$  (see Theorem 3.9), which can be further reduced by oversampling the  
 788 row indices. For a stable implementation of the CURCA,  $A \approx CU^\dagger R$ , we recommend  
 789 the MATLAB implementation

$$790 \quad \mathbf{A}_{IJ}^{(3)} = (\mathbf{A}(:, \mathbf{J}) * \mathbf{V}/\mathbf{S}) * (\mathbf{W}' * \mathbf{A}(\mathbf{I}, :)) \text{ where } [\mathbf{W}, \mathbf{S}, \mathbf{V}] = \text{svd}(\mathbf{U}, \text{'econ'})$$

791 or the corresponding version where QR factorization is used instead of the SVD.

792 We also proposed how oversampling should be done. Oversampling should be  
 793 done such that it increase the minimum singular value of a certain square matrix that  
 794 is a submatrix of an orthonormal matrix; see Section 4. We suggest doing so through  
 795 projecting the unchosen rows of an orthonormal matrix onto the trailing singular  
 796 subspace of the square matrix and finding the important unchosen indices that will  
 797 enrich the trailing singular subspace; see Algorithm 4.2. Oversampling improves the  
 798 stability as the core matrix becomes rectangular and rectangular matrices are more  
 799 well-conditioned than square matrices. We recommend oversampling in one only one  
 800 of row or column indices, but not both (see Figure 3c) and choose the oversampling  
 801 parameter to be proportional to the target rank when possible. Experiments show  
 802 that the algorithm is competitive with other existing methods. Therefore, we advo-  
 803 cate oversampling for the accuracy and stability of the CUR decomposition whenever  
 804 possible.

- 806 [1] D. ANDERSON, S. DU, M. MAHONEY, C. MELGAARD, K. WU, AND M. GU, *Spectral Gap Error*  
807 *Bounds for Improving CUR Matrix Decomposition and the Nyström Method*, in Proceedings  
808 of the Eighteenth International Conference on Artificial Intelligence and Statistics,  
809 G. Lebanon and S. V. N. Vishwanathan, eds., vol. 38 of Proceedings of Machine Learning  
810 Research, San Diego, California, USA, 09–12 May 2015, PMLR, pp. 19–27.
- 811 [2] E. ANDERSON, Z. BAI, C. BISCHOF, J. DEMMEL, J. DONGARRA, J. DU CROZ, A. GREENBAUM,  
812 S. HAMMARLING, A. MCKENNEY, S. OSTROUCHOV, ET AL., *LAPACK Users' guide*, SIAM,  
813 1995.
- 814 [3] H. AVRON, P. MAYMOUNKOV, AND S. TOLEDO, *Blendenpik: Supercharging LAPACK's least-*  
815 *squares solver*, SIAM J. Sci. Comput., 32 (2010), pp. 1217–1236, [https://doi.org/10.1137/](https://doi.org/10.1137/090767911)  
816 [090767911](https://doi.org/10.1137/090767911).
- 817 [4] O. BALABANOV AND L. GRIGORI, *Randomized Gram–Schmidt process with application to*  
818 *GMRES*, SIAM J. Sci. Comput., 44 (2022), pp. A1450–A1474, [https://doi.org/10.1137/](https://doi.org/10.1137/20M138870X)  
819 [20M138870X](https://doi.org/10.1137/20M138870X).
- 820 [5] J. BATSON, D. A. SPIELMAN, AND N. SRIVASTAVA, *Twice-ramanujan sparsifiers*, SIAM J. Com-  
821 put., 41 (2012), pp. 1704–1721, <https://doi.org/10.1137/090772873>.
- 822 [6] C. BOUTSIDIS, P. DRINEAS, AND M. MAGDON-ISMAIL, *Near-optimal column-based matrix recon-*  
823 *struction*, SIAM J. Comput., 43 (2014), pp. 687–717, <https://doi.org/10.1137/12086755X>.
- 824 [7] C. BOUTSIDIS AND D. P. WOODRUFF, *Optimal CUR matrix decompositions*, SIAM J. Comput.,  
825 46 (2017), pp. 543–589, <https://doi.org/10.1137/140977898>.
- 826 [8] D. CAI, J. NAGY, AND Y. XI, *Fast deterministic approximation of symmetric indefinite kernel*  
827 *matrices with high dimensional datasets*, SIAM J. Matrix Anal. Appl., 43 (2022), pp. 1003–  
828 1028, <https://doi.org/10.1137/21M1424627>.
- 829 [9] K. CARLBERG, C. FARHAT, J. CORTIAL, AND D. AMSALLEM, *The GNAT method for nonlinear*  
830 *model reduction: Effective implementation and application to computational fluid dynam-*  
831 *ics and turbulent flows*, J. Comput. Phys., 242 (2013), pp. 623–647, [https://doi.org/https://doi.org/https://doi.org/10.1016/j.jcp.2013.02.028](https://doi.org/https://doi.org/10.1016/j.jcp.2013.02.028).
- 832 [10] T. F. CHAN, *Rank revealing QR factorizations*, Linear Algebra Appl., 88–89 (1987), pp. 67–82,  
833 [https://doi.org/https://doi.org/10.1016/0024-3795\(87\)90103-0](https://doi.org/https://doi.org/10.1016/0024-3795(87)90103-0).
- 834 [11] S. CHATURANTABUT AND D. C. SORENSEN, *Nonlinear model reduction via discrete empirical*  
835 *interpolation*, SIAM J. Sci. Comput., 32 (2010), pp. 2737–2764, [https://doi.org/10.1137/](https://doi.org/10.1137/090766498)  
836 [090766498](https://doi.org/10.1137/090766498).
- 837 [12] J. CHIU AND L. DEMANET, *Sublinear randomized algorithms for skeleton decompositions*, SIAM  
838 J. Matrix Anal. Appl., 34 (2013), pp. 1361–1383, <https://doi.org/10.1137/110852310>.
- 839 [13] A. CORTINOVIS AND D. KRESSNER, *Low-rank approximation in the Frobenius norm by column*  
840 *and row subset selection*, SIAM J. Matrix Anal. Appl., 41 (2020), pp. 1651–1673.
- 841 [14] A. CORTINOVIS AND L. YING, *A sublinear-time randomized algorithm for column and*  
842 *row subset selection based on strong rank-revealing QR factorizations*, arXiv preprint  
843 arXiv:2402.13975, (2024).
- 844 [15] T. A. DAVIS AND Y. HU, *The University of Florida sparse matrix collection*, ACM Trans.  
845 Math. Softw., 38 (2011), <https://doi.org/10.1145/2049662.2049663>, [https://doi.org/10.](https://doi.org/10.1145/2049662.2049663)  
846 [1145/2049662.2049663](https://doi.org/10.1145/2049662.2049663).
- 847 [16] M. DEREZIŃSKI AND M. MAHONEY, *Determinantal point processes in randomized numerical*  
848 *linear algebra*, Notices Amer. Math. Soc., 60 (2021), p. 1, <https://doi.org/10.1090/noti2202>.
- 849 [17] A. DESHPANDE, L. RADEMACHER, S. S. VEMPALA, AND G. WANG, *Matrix approximation and*  
850 *projective clustering via volume sampling*, Theory Comput., 2 (2006), pp. 225–247, <https://doi.org/10.4086/toc.2006.v002a012>.
- 851 [18] A. DESHPANDE AND S. VEMPALA, *Adaptive sampling and fast low-rank matrix approximation*, in  
852 Approximation, Randomization, and Combinatorial Optimization. Algorithms and Tech-  
853 niques, J. Díaz, K. Jansen, J. D. P. Rolim, and U. Zwick, eds., Berlin, Heidelberg, 2006,  
854 Springer Berlin Heidelberg, pp. 292–303.
- 855 [19] M. DONELLO, G. PALKAR, M. H. NADERI, D. C. DEL REY FERNÁNDEZ, AND H. BABAEI,  
856 *Oblique projection for scalable rank-adaptive reduced-order modelling of nonlinear sto-*  
857 *chastic partial differential equations with time-dependent bases*, Proceedings of the Royal  
858 Society A: Mathematical, Physical and Engineering Sciences, 479 (2023), p. 20230320,  
859 <https://doi.org/10.1098/rspa.2023.0320>.
- 860 [20] Y. DONG AND P.-G. MARTINSSON, *Simpler is better: a comparative study of randomized pivoting*  
861 *algorithms for CUR and interpolative decompositions*, Adv. Comput. Math., 49 (2023),  
862 <https://doi.org/10.1007/s10444-023-10061-z>.
- 863 [21] P. DRINEAS AND I. C. F. IPSEN, *Low-rank matrix approximations do not need a singular*  
864 *value decomposition*, SIAM J. Matrix Anal. Appl., 63 (2021), pp. 1001–1020, <https://doi.org/10.1137/20M1001001>.

- 866 *value gap*, SIAM J. Matrix Anal. Appl., 40 (2019), pp. 299–319, [https://doi.org/10.1137/](https://doi.org/10.1137/18M1163658)  
867 [18M1163658](https://doi.org/10.1137/18M1163658).
- 868 [22] P. DRINEAS, M. MAGDON-ISMAIL, M. W. MAHONEY, AND D. P. WOODRUFF, *Fast approximation*  
869 *of matrix coherence and statistical leverage*, J. Mach. Learn. Res., 13 (2012), pp. 3475–3506,  
870 <http://jmlr.org/papers/v13/drineas12a.html>.
- 871 [23] P. DRINEAS, M. W. MAHONEY, AND S. MUTHUKRISHNAN, *Relative-error CUR matrix decom-*  
872 *positions*, SIAM J. Matrix Anal. Appl., 30 (2008), pp. 844–881, [https://doi.org/10.1137/](https://doi.org/10.1137/07070471X)  
873 [07070471X](https://doi.org/10.1137/07070471X).
- 874 [24] Z. DRMAČ AND S. GUGERCIN, *A new selection operator for the discrete empirical interpolation*  
875 *method—improved a priori error bound and extensions*, SIAM J. Sci. Comput., 38 (2016),  
876 pp. A631–A648, <https://doi.org/10.1137/15M1019271>.
- 877 [25] J. A. DUERSCH AND M. GU, *Randomized projection for rank-revealing matrix factorizations*  
878 *and low-rank approximations*, SIAM Rev., 62 (2020), pp. 661–682, [https://doi.org/10.](https://doi.org/10.1137/20M1335571)  
879 [1137/20M1335571](https://doi.org/10.1137/20M1335571).
- 880 [26] P. Y. GIDISU AND M. E. HOCHSTENBACH, *A hybrid DEIM and leverage scores based method for*  
881 *CUR index selection*, in Progress in Industrial Mathematics at ECMI 2021, M. Ehrhardt  
882 and M. Günther, eds., Cham, 2022, Springer International Publishing, pp. 147–153.
- 883 [27] G. H. GOLUB AND C. F. VAN LOAN, *Matrix Computations*, Johns Hopkins University Press,  
884 4 ed., 2013.
- 885 [28] S. GOREINOV, E. TYRTYSHNIKOV, AND N. ZAMARASHKIN, *A theory of pseudoskeleton approxi-*  
886 *mations*, Linear Algebra Appl., 261 (1997), pp. 1–21, [https://doi.org/https://doi.org/10.](https://doi.org/https://doi.org/10.1016/S0024-3795(96)00301-1)  
887 [1016/S0024-3795\(96\)00301-1](https://doi.org/https://doi.org/10.1016/S0024-3795(96)00301-1).
- 888 [29] M. GU AND S. C. EISENSTAT, *Efficient algorithms for computing a strong rank-revealing QR*  
889 *factorization*, SIAM J. Sci. Comput., 17 (1996), pp. 848–869, [https://doi.org/10.1137/](https://doi.org/10.1137/0917055)  
890 [0917055](https://doi.org/10.1137/0917055).
- 891 [30] N. HALKO, P.-G. MARTINSSON, AND J. A. TROPP, *Finding structure with randomness: Prob-*  
892 *abilistic algorithms for constructing approximate matrix decompositions*, SIAM Rev., 53  
893 (2011), p. 217–288, <https://doi.org/10.1137/090771806>.
- 894 [31] K. HAMM AND L. HUANG, *Perspectives on CUR decompositions*, Appl. Comput. Harmon. Anal.,  
895 48 (2020), pp. 1088–1099, <https://doi.org/https://doi.org/10.1016/j.acha.2019.08.006>.
- 896 [32] K. HAMM AND L. HUANG, *Stability of sampling for CUR decompositions*, Foundations of Data  
897 Science, 2 (2020), pp. 83–99, <https://doi.org/10.3934/fods.2020006>.
- 898 [33] K. HAMM AND L. HUANG, *Perturbations of CUR decompositions*, SIAM J. Matrix Anal. Appl.,  
899 42 (2021), pp. 351–375, <https://doi.org/10.1137/19M128394X>.
- 900 [34] N. J. HIGHAM, *Accuracy and Stability of Numerical Algorithms*, SIAM, second ed., 2002, <https://doi.org/10.1137/1.9780898718027>.
- 901 [35] R. A. HORN AND C. R. JOHNSON, *Matrix Analysis*, Cambridge University Press, 2 ed., 2012,  
902 <https://doi.org/10.1017/9781139020411>.
- 903 [36] I. C. F. IPSEN AND B. NADLER, *Refined perturbation bounds for eigenvalues of Hermitian*  
904 *and non-Hermitian matrices*, SIAM J. Matrix Anal. Appl., 31 (2009), pp. 40–53, <https://doi.org/10.1137/070682745>.
- 905 [37] A. KRIZHEVSKY, *Learning multiple layers of features from tiny images*, 2009.
- 906 [38] M. W. MAHONEY AND P. DRINEAS, *CUR matrix decompositions for improved data analysis*,  
907 Proc. Natl. Acad. Sci., 106 (2009), pp. 697–702, <https://doi.org/10.1073/pnas.0803205106>.
- 908 [39] P.-G. MARTINSSON, *Randomized methods for matrix computations*, The Mathematics of Data,  
909 25 (2019), pp. 187–231.
- 910 [40] P.-G. MARTINSSON AND J. A. TROPP, *Randomized numerical linear algebra: Founda-*  
911 *tions and algorithms*, Acta Numer., 29 (2020), p. 403–572, [https://doi.org/10.1017/](https://doi.org/10.1017/s0962492920000021)  
912 [s0962492920000021](https://doi.org/10.1017/s0962492920000021).
- 913 [41] E. S. MECKES, *The random matrix theory of the classical compact groups*, vol. 218, Cambridge  
914 University Press, 2019.
- 915 [42] Y. NAKATSUKASA, *Fast and stable randomized low-rank matrix approximation*, arXiv preprint  
916 arXiv:2009.11392, (2020), <https://arxiv.org/abs/2009.11392>.
- 917 [43] Y. NAKATSUKASA AND N. J. HIGHAM, *Backward stability of iterations for computing the polar*  
918 *decomposition*, SIAM J. Matrix Anal. Appl., 33 (2012), pp. 460–479, [https://doi.org/10.](https://doi.org/10.1137/110857544)  
919 [1137/110857544](https://doi.org/10.1137/110857544).
- 920 [44] Y. NAKATSUKASA AND T. PARK, *Randomized low-rank approximation for symmetric indefinite*  
921 *matrices*, SIAM J. Matrix Anal. Appl., 44 (2023), pp. 1370–1392, [https://doi.org/10.1137/](https://doi.org/10.1137/22M1538648)  
922 [22M1538648](https://doi.org/10.1137/22M1538648).
- 923 [45] C. PAIGE AND M. WEI, *History and generality of the CS decomposition*, Linear Al-  
924 gebra Appl., 208–209 (1994), pp. 303–326, [https://doi.org/https://doi.org/10.1016/](https://doi.org/https://doi.org/10.1016/0024-3795(94)90446-4)  
925 [0024-3795\(94\)90446-4](https://doi.org/https://doi.org/10.1016/0024-3795(94)90446-4).

928 [46] D. PAPAILOPOULOS, A. KYRILLIDIS, AND C. BOUTSIDIS, *Provable deterministic leverage score*  
 929 *sampling*, in Proceedings of the 20th ACM SIGKDD International Conference on Knowl-  
 930 edge Discovery and Data Mining, KDD '14, New York, NY, USA, 2014, Association  
 931 for Computing Machinery, p. 997–1006, <https://doi.org/10.1145/2623330.2623698>, <https://doi.org/10.1145/2623330.2623698>.  
 932  
 933 [47] B. PEHERSTORFER, Z. DRMAČ, AND S. GUGERCIN, *Stability of discrete empirical interpola-*  
 934 *tion and gappy proper orthogonal decomposition with randomized and deterministic sam-*  
 935 *pling points*, SIAM J. Sci. Comput., 42 (2020), pp. A2837–A2864, [https://doi.org/10.1137/](https://doi.org/10.1137/19M1307391)  
 936 [19M1307391](https://doi.org/10.1137/19M1307391).  
 937 [48] A. SHIMIZU, X. CHENG, C. MUSCO, AND J. WEARE, *Improved active learning via dependent*  
 938 *leverage score sampling*, arXiv preprint arXiv:2310.04966, (2023).  
 939 [49] D. C. SORENSEN AND M. EMBREE, *A DEIM induced CUR factorization*, SIAM J. Sci. Comp.,  
 940 38 (2016), pp. A1454–A1482, <https://doi.org/10.1137/140978430>.  
 941 [50] D. B. SZYLD, *The many proofs of an identity on the norm of oblique projections*, Numer.  
 942 Algorithms, 42 (2006), pp. 309–323, <https://doi.org/10.1007/s11075-006-9046-2>.  
 943 [51] L. N. TREFETHEN AND D. BAU, *Numerical Linear Algebra*, SIAM, 1997, [https://doi.org/10.](https://doi.org/10.1137/1.9780898719574)  
 944 [1137/1.9780898719574](https://doi.org/10.1137/1.9780898719574).  
 945 [52] J. A. TROPP, A. YURTSEVER, M. UDELL, AND V. CEVHER, *Practical sketching algorithms for*  
 946 *low-rank matrix approximation*, SIAM J. Matrix Anal. Appl., 38 (2017), p. 1454–1485,  
 947 <https://doi.org/10.1137/17m1111590>.  
 948 [53] M. UDELL AND A. TOWNSEND, *Why are big data matrices approximately low rank?*, SIAM J.  
 949 Math. Data Sci., 1 (2019), pp. 144–160, <https://doi.org/10.1137/18M1183480>.  
 950 [54] S. VORONIN AND P.-G. MARTINSSON, *Efficient algorithms for CUR and interpolative matrix*  
 951 *decompositions*, Adv. Comput. Math., 43 (2017), pp. 495–516, [https://doi.org/10.1007/](https://doi.org/10.1007/s10444-016-9494-8)  
 952 [s10444-016-9494-8](https://doi.org/10.1007/s10444-016-9494-8).  
 953 [55] J. XIA, *Making the Nyström method highly accurate for low-rank approximations*, SIAM J. Sci.  
 954 Comput., 46 (2024), pp. A1076–A1101, <https://doi.org/10.1137/23M1585039>.  
 955 [56] N. L. ZAMARASHKIN AND A. I. OSINSKY, *On the existence of a nearly optimal skeleton ap-*  
 956 *proximation of a matrix in the Frobenius norm*, Dokl. Math., 97 (2018), pp. 164–166,  
 957 <https://doi.org/10.1134/S1064562418020205>.  
 958 [57] R. ZIMMERMANN AND K. WILLCOX, *An accelerated greedy missing point estimation procedure*,  
 959 SIAM J. Sci. Comput., 38 (2016), pp. A2827–A2850, <https://doi.org/10.1137/15M1042899>.

960 **Appendix A. Analysis of the CURBA.** In this section, we analyze the  
 961 CURBA,

$$962 \quad (A.1) \quad A_{IJ}^{(BA)} = A(:, J) (A(:, J)^\dagger A A(I, :))^\dagger A(I, :) = C(C^\dagger A R^\dagger) R$$

963 with oversampling. For the CURBA, there exists a numerically stable algorithm given  
 964 by the StableCUR algorithm in [1]. The algorithm computes the CURBA in the  
 965 following way in MATLAB,

$$966 \quad (A.2) \quad [Q_C, \sim] = \text{qr}(A(:, J), 0), [Q_R, \sim] = \text{qr}(A(I, :)', 0), A_{IJ}^{(BA)} = Q_C * (Q_C' * A * Q_R) * Q_R'.$$

967 We use this implementation of the CURBA in the experiments at the end of this  
 968 section.

969 We first prove a relative norm bound for the CURBA with oversampling. We make  
 970 the following standard assumptions, which is analogous to the CURCA counterpart  
 971 in Section 3.

972 *Assumption A.1.*

- 973 1.  $|I| = |J| = k$  where  $k$  is the target rank. The oversampling indices are  $I_0$   
 974 with  $|I_0| = p_1$  for the rows and  $J_0$  with  $|J_0| = p_2$  for the columns and they  
 975 satisfy  $k + \max\{p_1, p_2\} \leq \text{rank}(A)$ .
- 976 2.  $A(:, J \cup J_0)$  has full column rank and  $A(I \cup I_0, :)$  has full row rank.
- 977 3.  $X(:, J) \in \mathbb{R}^{k \times (k+p_2)}$  has full row rank, where  $X \in \mathbb{R}^{k \times n}$  is a row space  
 978 approximator of  $A$ .
- 979 4.  $Y(I, :) \in \mathbb{R}^{k \times (k+p_2)}$  has full row rank, where  $Y \in \mathbb{R}^{m \times k}$  is a column space  
 980 approximator of  $A$ .

981 The assumption that  $X(:, J)$  and  $A(I, :)$  having a full row rank and  $Y(I, :)$  and  $A(:, J)$   
 982 having a full column rank is generic, since most methods that pick good row and  
 983 column indices satisfy this assumption. If one of  $X$  or  $Y$  are not available then  $A(I, :)$   
 984 or  $A(:, J)$  can be used instead, respectively.

985 We now prove the result for the CURBA with oversampling. The results shown  
 986 below is a simple extension of [49, Lemma 4.2] and [20, Theorem 1]. Lemma A.2  
 987 considers one-sided projection with oversampling where we project  $A$  onto the chosen  
 988 columns of  $A$ . In Theorem A.3, we consider the CURBA,  $C(C^\dagger AR^\dagger)R$  where we  
 989 project  $A$  onto the chosen rows and columns of  $A$ .

990 LEMMA A.2. *Under Assumption A.1,*

$$991 \quad (\text{A.3}) \quad \|A - CC^\dagger A\| \leq \|Q_X(J \cup J_0, :)\|_2 \|A - AX^\dagger X\|$$

992 where  $\|\cdot\|$  is any unitarily invariant norm and  $Q_X \in \mathbb{R}^{n \times k}$  is an orthonormal matrix  
 993 spanning the columns of  $X^T$ .

994 *Proof.* For shorthand let  $J_* = J \cup J_0$  and let  $\Pi_{J_*} = I_n(:, J_*) \in \mathbb{R}^{n \times (k+p)}$  such  
 995 that  $C = A(:, J_*) = A\Pi_{J_*}$ . We first define two oblique projectors

$$996 \quad \mathcal{P}_X := \Pi_{J_*}(X\Pi_{J_*})^\dagger X, \quad \mathcal{P}_C := \Pi_{J_*}C^\dagger A \in \mathbb{R}^{n \times n}.$$

997 Note that since  $C$  has full column rank,  $C^\dagger A\Pi_{J_*} = C^\dagger C = I_{k+p}$ , and

$$998 \quad \mathcal{P}_C \mathcal{P}_X = \Pi_{J_*}C^\dagger A\Pi_{J_*}(X\Pi_{J_*})^\dagger X = \mathcal{P}_X.$$

999 Therefore we get

$$1000 \quad A - CC^\dagger A = A(I - \mathcal{P}_C) = A(I - \mathcal{P}_C)(I - \mathcal{P}_X) = (I - CC^\dagger)A(I - \mathcal{P}_X).$$

1001 Since  $X\Pi_{J_*} = X(:, J_*)$  has full row rank,

$$1002 \quad X\mathcal{P}_X = X\Pi_{J_*}(X\Pi_{J_*})^\dagger X = X$$

1003 and we obtain

$$1004 \quad (I - \mathcal{P}_X) = (I - X^\dagger X)(I - \mathcal{P}_X).$$

1005 Now putting these together, we obtain

$$\begin{aligned} 1006 \quad \|A - CC^\dagger A\| &= \|(I - CC^\dagger)A(I - X^\dagger X)(I - \mathcal{P}_X)\| \\ 1007 \quad &\leq \|I - CC^\dagger\|_2 \|A(I - X^\dagger X)\| \|I - \mathcal{P}_X\|_2 \\ 1008 \quad &= \|I - \mathcal{P}_X\|_2 \|A(I - X^\dagger X)\|, \end{aligned}$$

1009 since  $\|I - CC^\dagger\|_2 = 1$  as  $CC^\dagger$  is an orthogonal projector. Now since  $\mathcal{P}_X$  is an oblique  
 1010 projector [50],  $\|I - \mathcal{P}_X\|_2 = \|\mathcal{P}_X\|_2$  so the result follows by noting that  $\square$

$$1011 \quad \|\mathcal{P}_X\|_2 = \|(X\Pi_{J_*})^\dagger X\|_2 = \|Q_X(J, :)\|_2.$$

1012 THEOREM A.3. *Under Assumption A.1,*

$$1013 \quad (\text{A.4}) \quad \left\| A - A_{I \cup J_0, J \cup J_0}^{(BA)} \right\| \leq \|Q_X(J \cup J_0, :)\|_2 \|A - AX^\dagger X\|$$

$$1014 \quad + \|Q_Y(I \cup I_0, :)^{\dagger}\|_2 \|A - YY^{\dagger}A\|$$

1015 where  $\|\cdot\|$  is any unitarily invariant norm and  $Q_X \in \mathbb{R}^{n \times k}$  and  $Q_Y \in \mathbb{R}^{m \times k}$  are the  
 1016 orthonormal matrices spanning the columns of  $X^T$  and  $Y$ , respectively.

1017 *Proof.* The proof follows by Lemma A.2 and noting that

$$\begin{aligned} 1018 \quad \|A - CC^{\dagger}AR^{\dagger}R\| &\leq \|A - CC^{\dagger}A\| + \|CC^{\dagger}A - CC^{\dagger}AR^{\dagger}R\| \\ 1019 \quad &\leq \|A - CC^{\dagger}A\| + \|CC^{\dagger}\|_2 \|A - AR^{\dagger}R\| \\ 1020 \quad &= \|A - CC^{\dagger}A\| + \|A - AR^{\dagger}R\| \\ 1021 \quad &\leq \|Q_X(J, :)^{\dagger}\|_2 \|A - AX^{\dagger}X\| + \|Q_Y(I, :)^{\dagger}\|_2 \|A - YY^{\dagger}A\| \end{aligned}$$

1022 where the inequality for the second term  $\|A - AR^{\dagger}R\|$  in the final line can be shown  
 1023 by considering  $A^T$  in Lemma A.2.  $\square$

1024 *Remark A.4.*

1025 1. The difference between Theorem A.3 and its counterpart without oversam-  
 1026 pling, e.g. [20, Theorem 1], are the terms  $\|Q_X(J, :)^{\dagger}\|_2$  and  $\|Q_Y(I, :)^{\dagger}\|_2$   
 1027 where instead of the matrix inverse, we have the pseudoinverse. This tight-  
 1028 ens the bound (A.4) as

$$1029 \quad \|Q_X(J \cup J_0, :)^{\dagger}\|_2 \leq \|Q_X(J, :)^{-1}\|_2,$$

1030 and

$$1031 \quad \|Q_Y(I \cup I_0, :)^{\dagger}\|_2 \leq \|Q_Y(I, :)^{-1}\|_2,$$

1032 where  $I_0$  and  $J_0$  are the extra indices for oversampling.

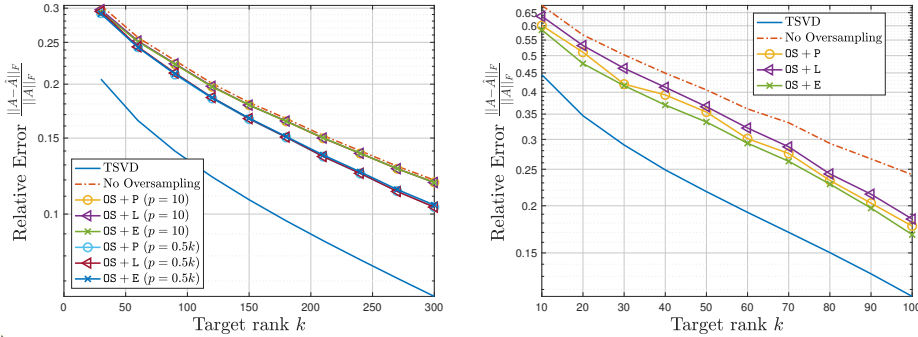
1033 2. There are many possible choices for  $X$  and  $Y$ . For example, if we choose  
 1034 them to be the dominant right and left singular vectors of  $A$ , then we get  
 1035 the bound similar to the one in [49] involving the best rank- $k$  approximation  
 1036 since  $\|A - AX^{\dagger}X\| = \|A - YY^{\dagger}A\| = \|A - \llbracket A \rrbracket_k\|$  where  $\llbracket A \rrbracket_k$  is the best  
 1037 rank- $k$  approximation to  $A$ . If we choose  $X$  and  $Y$  to be the row sketch and  
 1038 the column sketch (see Algorithm 4.1) then we involve the randomized SVD  
 1039 error [30]. Choosing  $X$  and  $Y$  to be the dominant singular vectors in Theorem  
 1040 A.3 can be the optimal choice, but when  $X$  and  $Y$  are approximators, which  
 1041 can be computed easily, the bound (A.4) can be computed a posteriori.

1042 3. When  $k + \min\{p_1, p_2\} \geq \text{rank}(A)$  and  $\text{rank}(C) = \text{rank}(R) = \text{rank}(A)$ , we  
 1043 have  $A = CC^{\dagger}AR^{\dagger}R$  [31].

1044 Theorem A.3 suggests that oversampling helps to improve the accuracy of the  
 1045 CUR decomposition  $A_{IJ}^{(BA)}$ . However, the numerical simulations shown below illus-  
 1046 trate that the improvement is not too significant.

1047 The numerical results are depicted in Figure 4, illustrating the effect of oversam-  
 1048 pling for the CURBA. We use pivoting on a random sketch (Algorithm 4.1) to obtain  
 1049 the initial set of indices. Then we oversample the row indices using the three over-  
 1050 sampling algorithms, OS+P, OS+L and OS+E as in Section 5.3. We also oversample  
 1051 column indices in Figure 4b.

1052 We start with the CIFAR10 dataset in Figure 4a. We first observe that as the over-  
 1053 sampling parameter increases, the accuracy improves and the different oversampling  
 1054 techniques yield a similar result. However, there is no significant improvement when



(a) CIFAR10 dataset with pivoting on a random sketch with row oversampling. (b) YaleFace64x64 dataset with pivoting on a random sketch with row ( $p = 0.5k$ ) and column ( $p = 10$ ) oversampling.

FIG. 4. The effect of oversampling for the CURBA. We use pivoting on a random sketch (Algorithm 4.1) to obtain the initial set of indices. Then we oversample the row indices using the three oversampling algorithms, OS+P, OS+L and OS+E. We also oversample column indices in Figure 4b.

1055 oversampling is used for the CURBA. This shows that the CURBA gives a robust  
 1056 approximation and oversampling only plays a slight role in its accuracy. In Figure  
 1057 4b, we used the YaleFace64x64 dataset to illustrate the effect of oversampling both  
 1058 rows and columns. We observe that when we oversample both rows and columns, we  
 1059 obtain a higher accuracy for the CURBA. This is expected as we are enlarging the  
 1060 subspace that we project onto  $A$ , which increases the accuracy of the CURBA. This  
 1061 is contrary to the CURCA, as the CURCA can become more inaccurate and unstable  
 1062 when we oversample both rows and columns; see Section 5.3.

1063 To conclude, while oversampling improves the accuracy of the CURBA, it has  
 1064 less significant effect on its accuracy and stability.

1065 **Appendix B. Stability of rank-deficient systems.** In this section, we  
 1066 prove the stability result for solving rank-deficient underdetermined linear systems  
 1067 (Theorem 3.8). The result can also be extended to rank-deficient overdetermined  
 1068 problems; see Appendix B.1. Theorem 3.8 was used to derive an expression for the  
 1069 computed solution to

$$1070 \quad (\text{B.1}) \quad \min_x \left\| (\hat{U}^T)_\epsilon x - [\hat{C}]_i^T \right\|_2$$

1071 for each row of  $\hat{C}$ . The statement and the proof is shown below. The proof follows a  
 1072 similar method outlined in [42, Section 4.1].

1073 **THEOREM B.1** (Theorem 3.8). Consider the (rank-deficient) underdetermined lin-  
 1074 ear system,

$$1075 \quad (3.7) \quad \min_{x'} \|B_\epsilon x' - b\|_2$$

1076 where  $B_\epsilon \in \mathbb{R}^{m \times n}$  ( $m \leq n$ ) is (possibly) rank-deficient ( $\text{rank}(B_\epsilon) \leq m$ ) with singular  
 1077 vales larger than  $\epsilon$  and  $b \in \mathbb{R}^m$ . Then assuming  $\epsilon > \gamma \|B_\epsilon\|_2$ , the minimum norm  
 1078 solution to (3.7) can be computed in a backward stable manner, i.e., the computed

1079 solution  $\hat{s}$  satisfies

$$1080 \quad \hat{s} = (B_\epsilon + E_1)^\dagger (b + E_2)$$

1081 where  $\|E_1\|_2 \leq \gamma \|B_\epsilon\|_2$  and  $\|E_2\|_2 \leq \gamma \|b\|_2$ .

1082 *Proof.* First, if  $B_\epsilon$  has full row-rank then the statement follows by the stability  
 1083 result in [34, Theorem 21.4]. If  $B_\epsilon$  is rank-deficient, the stability result in [34, The-  
 1084 orem 21.4] cannot be invoked as it is only applicable for underdetermined full-rank  
 1085 linear systems. So we project  $B_\epsilon$  onto its column space  $Q_\epsilon$  to make the problem an  
 1086 underdetermined full-rank linear system:

$$1087 \quad \min_{x'} \|Q_\epsilon^T B_\epsilon x' - Q_\epsilon^T b\|_2.$$

1088 Let  $\hat{Q}_\epsilon$  be the computed column space of  $B_\epsilon$ . Then  $\hat{Q}_\epsilon$  is the exact column space  
 1089 of  $B_\epsilon + \Delta B$  where  $\|\Delta B\|_2 \leq \gamma \|B_\epsilon\|_2$ ,  $\hat{Q}_\epsilon^T B_\epsilon$  is numerically full-rank with singular  
 1090 values larger than  $\epsilon$  and  $\hat{Q}_\epsilon \hat{Q}_\epsilon^T B_\epsilon = B_\epsilon + E$  where  $\|E\|_2 \leq \gamma \|B_\epsilon\|_2$  [27, Chapter  
 1091 5.4.1]. Note the following from matrix-matrix (or matrix-vector) multiplication [34,  
 1092 Section 3.5],

$$1093 \quad fl(\hat{Q}_\epsilon^T B_\epsilon) = \hat{Q}_\epsilon^T B_\epsilon + E^{(1)}$$

1094 where  $\|E^{(1)}\|_2 \leq \gamma \|B_\epsilon\|_2$  and

$$1095 \quad fl(\hat{Q}_\epsilon^T b) = \hat{Q}_\epsilon^T b + E^{(2)}$$

1096 where  $\|E^{(2)}\|_2 \leq \gamma \|b\|_2$ . Since  $Q_\epsilon^T B_\epsilon \in \mathbb{R}^{\text{rank}(B_\epsilon) \times n}$  is a fat rectangular matrix, we  
 1097 solve the following underdetermined full-rank linear system,

$$1098 \quad (\text{B.2}) \quad \min_{x'} \left\| \left( \hat{Q}_\epsilon^T B_\epsilon \right) x' - \hat{Q}_\epsilon^T b \right\|_2.$$

1099 Assuming  $\gamma \kappa_2(B_\epsilon) \leq \gamma \|B_\epsilon\|_2 / \epsilon < 1$  (where  $\kappa_2$  denotes the 2-norm condition number),  
 1100 we obtain the computed solution of (B.2),  $\hat{s}$ , satisfying [34, Theorem 21.4],

$$\begin{aligned} 1101 \quad \hat{s} &= fl \left( \left( \hat{Q}_\epsilon^T B_\epsilon \right)^\dagger \hat{Q}_\epsilon^T b \right) = \left( fl(\hat{Q}_\epsilon^T B_\epsilon) + E^{(3)} \right)^\dagger fl(\hat{Q}_\epsilon^T b) \\ 1102 &= \left( \hat{Q}_\epsilon^T B_\epsilon + E^{(1)} + E^{(3)} \right)^\dagger \left( \hat{Q}_\epsilon^T b + E^{(2)} \right) \\ 1103 &= \left( \hat{Q}_\epsilon \hat{Q}_\epsilon^T B_\epsilon + \hat{Q}_\epsilon E^{(1)} + \hat{Q}_\epsilon E^{(3)} \right)^\dagger \left( b + \hat{Q}_\epsilon E^{(2)} \right) \\ 1104 &= \left( B_\epsilon + E + \hat{Q}_\epsilon E^{(1)} + \hat{Q}_\epsilon E^{(3)} \right)^\dagger \left( b + \hat{Q}_\epsilon E^{(2)} \right) \end{aligned}$$

1105 where  $\|E^{(3)}\|_2 \leq \gamma \|B_\epsilon\|_2$  and we use  $(QA)^\dagger b = A^\dagger Q^T b$  for a tall-skinny orthonormal  
 1106 matrix  $Q$  in line 3. Therefore, assuming that  $\gamma \|B_\epsilon\|_2 < \epsilon$ , the computed solution of  
 1107 (3.7),  $\hat{s}$  satisfies

$$1108 \quad (\text{B.3}) \quad \hat{s} = fl(B_\epsilon^\dagger b) = (B_\epsilon + E_1)^\dagger (b + E_2),$$

1109 where  $\|E_1\|_2 \leq \gamma \|B_\epsilon\|_2$  and  $\|E_2\|_2 \leq \gamma \|b\|_2$ . □

1110 **B.1. Extension to rank-deficient overdetermined problems.** Theorem 3.8  
 1111 can be extended to include rank-deficient overdetermined least-squares problems.  
 1112 Consider the rank-deficient overdetermined least-squares problem,

$$1113 \quad \min_{x'} \|C_\epsilon x' - b\|_2$$

1114 where  $C_\epsilon \in \mathbb{R}^{m \times n}$  ( $m \geq n$ ) is a rank-deficient ( $\text{rank}(C_\epsilon) < n$ ), tall-skinny matrix with  
 1115 singular values larger than  $\epsilon$  and  $b \in \mathbb{R}^m$ . Then under the same assumption as Theorem  
 1116 3.8, i.e.,  $\epsilon > \gamma \|B_\epsilon\|_2$ , the solution to the overdetermined least-squares problem can  
 1117 be computed in a backward stable manner.

1118 The proof proceeds in the same way as Theorem 3.8, as once the tall-skinny  
 1119 matrix  $C_\epsilon$  is projected onto its column space  $W_\epsilon$ ,  $W_\epsilon^T C_\epsilon \in \mathbb{R}^{\text{rank}(C_\epsilon) \times n}$  becomes a fat  
 1120 full-rank matrix, and we are now solving the underdetermined full-rank linear system:

$$1121 \quad \min_{x'} \left\| (\hat{W}_\epsilon^T C_\epsilon) x' - \hat{W}_\epsilon^T b \right\|_2.$$

1122 **Appendix C. Two lemmas.** We give the proofs for the two lemmas, Lemma  
 1123 3.6 and Lemma 3.7 in this section.

1124 **LEMMA C.1** (Lemma 3.6). *Under Assumption 3.2, for any  $\Delta C$  and  $\Delta U$ ,*

$$1125 \quad (\text{C.1}) \quad \|(C + \Delta C)(U + \Delta U)_\epsilon^\dagger\|_2 \leq \|Q_C(I_*, \cdot)^\dagger\|_2 \left(1 + \frac{1}{\epsilon} \|\Delta U\|_2\right) + \frac{1}{\epsilon} \|\Delta C\|_2.$$

1126 *Proof.* Let  $C = Q_C R_C$  be the thin QR factorization of  $C$ . Then

$$\begin{aligned} 1127 \quad \|(C + \Delta C)(U + \Delta U)_\epsilon^\dagger\|_2 &\leq \|Q_C R_C (U + \Delta U)_\epsilon^\dagger\|_2 + \|\Delta C (U + \Delta U)_\epsilon^\dagger\|_2 \\ 1128 &= \|R_C (U + \Delta U)_\epsilon^\dagger\|_2 + \|\Delta C (U + \Delta U)_\epsilon^\dagger\|_2 \\ 1129 &= \|(\Pi_{I_*}^T Q_C)^\dagger \Pi_{I_*}^T Q_C R_C (U + \Delta U)_\epsilon^\dagger\|_2 + \|\Delta C (U + \Delta U)_\epsilon^\dagger\|_2 \\ 1130 &\leq \|(\Pi_{I_*}^T Q_C)^\dagger\|_2 \|U (U + \Delta U)_\epsilon^\dagger\|_2 + \frac{1}{\epsilon} \|\Delta C\|_2 \end{aligned}$$

1131 where we use  $(\Pi_{I_*}^T Q_C)^\dagger \Pi_{I_*}^T Q_C = I$  for the penultimate line. The result follows by  
 1132 noting that  $\|U (U + \Delta U)_\epsilon^\dagger\|_2$  simplifies to

$$1133 \quad \|U (U + \Delta U)_\epsilon^\dagger\|_2 = \|(U + \Delta U)(U + \Delta U)_\epsilon^\dagger - \Delta U (U + \Delta U)_\epsilon^\dagger\|_2 \leq 1 + \frac{1}{\epsilon} \|\Delta U\|_2. \quad \square$$

1134 **LEMMA C.2** (Lemma 3.7). *Under Assumption 3.2, for any  $\Delta C$  and  $\Delta U$ ,*

$$1135 \quad (\text{C.2}) \quad (C + \Delta C)(U + \Delta U)_\epsilon^\dagger R \Pi_J = C + E_*$$

1136 *where*

$$1137 \quad \|E_*\|_2 \leq \|(Q_C(I_*, \cdot)^\dagger)\|_2 \left( \epsilon + 2 \|\Delta U\|_2 + \frac{1}{\epsilon} \|\Delta U\|_2^2 \right) + \|\Delta C\|_2 \left( 1 + \frac{\|\Delta U\|_2}{\epsilon} \right).$$

1138 *Proof.* We divide the expression into two pieces:

$$1139 \quad (C + \Delta C)(U + \Delta U)_\epsilon^\dagger R \Pi_J = \underbrace{C(U + \Delta U)_\epsilon^\dagger R \Pi_J}_{(i)} + \underbrace{\Delta C(U + \Delta U)_\epsilon^\dagger R \Pi_J}_{(ii)}$$

1140 and treat them separately. Let the thin SVD of  $U + \Delta U$  be

$$1141 \quad U + \Delta U = W\Sigma V^T = [W_1, W_2] \begin{bmatrix} \Sigma_1 & \\ & \Sigma_2 \end{bmatrix} [V_1, V_2]^T$$

1142 where  $\Sigma_2$  contains the singular values of  $U + \Delta U$  smaller than  $\epsilon$  and  $C = Q_C R_C$  be  
 1143 the thin QR decomposition of  $C$ .

1144 We begin by examining the matrix (i).

$$\begin{aligned} 1145 \quad C(U + \Delta U)_\epsilon^\dagger R \Pi_J &= C(U + \Delta U)_\epsilon^\dagger U = C(U + \Delta U)_\epsilon^\dagger (U + \Delta U) - C(U + \Delta U)_\epsilon^\dagger \Delta U \\ 1146 &= C V_1 V_1^T - C(U + \Delta U)_\epsilon^\dagger \Delta U = C - C V_2 V_2^T - C(U + \Delta U)_\epsilon^\dagger \Delta U \\ 1147 &= C + E_1 \end{aligned}$$

1148 where  $E_1 = -C V_2 V_2^T - C(U + \Delta U)_\epsilon^\dagger \Delta U$  satisfies

$$\begin{aligned} 1149 \quad \|E_1\|_2 &\leq \|C V_2 V_2^T\|_2 + \|C(U + \Delta U)_\epsilon^\dagger \Delta U\|_2 \\ 1150 &\leq \|C V_2 V_2^T\|_2 + \|Q_C(I_*, :)\|_2 \left(1 + \frac{1}{\epsilon} \|\Delta U\|\right) \|\Delta U\|_2 \end{aligned}$$

1151 by Lemma 3.6. We bound  $\|C V_2 V_2^T\|_2$  as in the final part of the proof of Theorem  
 1152 3.3.

$$\begin{aligned} 1153 \quad \|C V_2 V_2^T\|_2 &= \|R_C V_2 V_2^T\|_2 = \|(\Pi_{I_*}^T Q_C)^\dagger \Pi_{I_*}^T Q_C R_C V_2 V_2^T\|_2 = \|(\Pi_{I_*}^T Q_C)^\dagger U V_2 V_2^T\|_2 \\ 1154 &\leq \|(\Pi_{I_*}^T Q_C)^\dagger\|_2 \|(U + \Delta U) V_2 V_2^T - \Delta U V_2 V_2^T\|_2 \\ 1155 &\leq \|(\Pi_{I_*}^T Q_C)^\dagger\|_2 (\|W_2 \Sigma_2\|_2 + \|\Delta U\|_2) \\ 1156 &\leq \|(\Pi_{I_*}^T Q_C)^\dagger\|_2 (\epsilon + \|\Delta U\|_2). \end{aligned}$$

1157 Therefore, (i) can be bounded as

$$1158 \quad C(U + \Delta U)_\epsilon^\dagger R \Pi_J = C + E_1$$

1159 where  $\|E_1\|_2 \leq \|Q_C(I_*, :)\|_2 \left(\epsilon + 2\|\Delta U\|_2 + \frac{1}{\epsilon} \|\Delta U\|_2^2\right)$ .

1160 Next, we bound the matrix (ii). Let  $E_2 = \Delta C(U + \Delta U)_\epsilon^\dagger R \Pi_J$ , then

$$\begin{aligned} 1161 \quad \|E_2\|_2 &= \|\Delta C(U + \Delta U)_\epsilon^\dagger R \Pi_J\|_2 \\ 1162 &\leq \|\Delta C\|_2 \|(U + \Delta U)_\epsilon^\dagger U\|_2 \\ 1163 &\leq \|\Delta C\|_2 (\|(U + \Delta U)_\epsilon^\dagger (U + \Delta U)\|_2 + \|(U + \Delta U)_\epsilon^\dagger \Delta U\|_2) \\ 1164 &\leq \|\Delta C\|_2 \left(1 + \frac{\|\Delta U\|_2}{\epsilon}\right). \end{aligned}$$

1165 Putting everything together and setting  $E_* = E_1 + E_2$ , we get the desired result:

$$1166 \quad (C + \Delta C)(U + \Delta U)_\epsilon^\dagger R \Pi_J = C + E_*$$

1167 where

$$1168 \quad \|E_*\|_2 \leq \|Q_C(I_*, :)\|_2 \left(\epsilon + 2\|\Delta U\|_2 + \frac{1}{\epsilon} \|\Delta U\|_2^2\right) + \|\Delta C\|_2 \left(1 + \frac{\|\Delta U\|_2}{\epsilon}\right). \quad \square$$