

RANDOMIZED LOW-RANK APPROXIMATION FOR SYMMETRIC INDEFINITE MATRICES*

YUJI NAKATSUKASA[†] AND TAEJUN PARK[†]

Abstract. The Nyström method is a popular choice for finding a low-rank approximation to a symmetric positive semi-definite matrix. The method can fail when applied to symmetric indefinite matrices, for which the error can be unboundedly large. In this work, we first identify the main challenges in finding a Nyström approximation to symmetric indefinite matrices. We then prove the existence of a variant that overcomes the instability, and establish relative-error nuclear norm bounds of the resulting approximation that hold when the singular values decay rapidly. The analysis naturally leads to a practical algorithm, whose robustness is illustrated with experiments.

Key words. Symmetric matrices, Nyström method, Low-rank approximation, Randomized linear algebra

AMS subject classifications. 15A23, 65F55

1. Introduction. Low-rank structures are ubiquitous in the computational sciences. They appear frequently as matrices having low numerical rank [35]. A low-rank approximation to a matrix provides an efficient way to store and process the matrix when the dimension is large. The Nyström method [14, 24, 38] has been a popular choice for finding low-rank approximations to symmetric positive semi-definite (SPSD) matrices, especially in the machine learning community for kernel-based methods.

Let $A \in \mathbb{R}^{n \times n}$ be a SPSP matrix and let the positive integer r be the target rank. Then the Nyström method is given by $A_{nys}^{(s)} = CW^\dagger C^T$ where $C := AX \in \mathbb{R}^{n \times s}$ and $W := X^T A X \in \mathbb{R}^{s \times s}$ with $r \leq s < n$ and $X \in \mathbb{R}^{n \times s}$ is a sketching matrix. The positive integer s is called the sketch size, and typically $r < s \ll n$. Traditionally, X is chosen to be a column sampling matrix, which has exactly one non-zero entry equal to 1 in each column [14, 38]. In this case, C is a subset of s columns of A and W is an $s \times s$ principal submatrix of A . There are different sampling schemes for column sampling, including uniform sampling, leverage score sampling [14, 19, 38, 39, 21] and k-means++ sampling [25]. In recent years, other choices for X have been shown to be practical, including Gaussian matrices, subsampled randomized trigonometric transforms (SRTTs) and sparse maps [15, 20]. These are *random embeddings*, which are the focus of this paper, and unlike column sampling, they mix up the coordinates of a vector when applied [20].

In this paper, we investigate the effect of using $A_{nys}^{(s)}$ and its rank-restricted variants for symmetric matrices that are possibly *indefinite*. Low-rank approximation of symmetric indefinite matrices arises in many applications, such as learning in reproducing kernel Krein spaces [26], natural language processing [8, 27] and non-metric proximity transformations [12], which has applications in bioinformatics and social networks. The original matrix A does not have to be SPSP for one to form the Nyström approximation $A_{nys}^{(s)}$. However, the theory does not translate directly to symmetric indefinite matrices because it uses the fact that the original matrix is SPSP [13, 14, 36]. Indeed, the Nyström approximation can be very poor for indefinite A , as we illustrate below. In this work, we show that a judiciously constructed

*Date: May 31, 2023

Funding: TP was supported by the Heilbronn Institute for Mathematical Research.

[†]Mathematical Institute, University of Oxford, Oxford, OX2 6GG, UK, (nakatsukasa@maths.ox.ac.uk, park@maths.ox.ac.uk).

rank-restricted variant of the Nyström approximation, when used with random embeddings, is robust even for symmetric indefinite matrices, which often outperforms other existing methods as we show for synthetic datasets (Figure 5) and real datasets (Figure 6) in Section 4. We also show in Section 3 that there exists a projection for the core matrix W such that the Nyström approximation gives a good low-rank approximation to *any* symmetric matrix when the singular values decay sufficiently fast.

1.1. Nyström methods and related work. There are several variants of the Nyström method for SPSD matrices. There are two rank-restricted versions that give a rank- r approximation to $A_{nys}^{(s)}$ where $r < s$. The first version, which is more traditional, is defined by $A_{nys}^{(s,r)} = C \llbracket W \rrbracket_r^\dagger C^T$ [9, 14, 18] where $\llbracket W \rrbracket_r$ denotes the best rank- r approximation to the matrix W using the truncated SVD. The second version is given by $\llbracket A_{nys}^{(s)} \rrbracket_r = \llbracket CW^\dagger C^T \rrbracket_r$ [28, 32, 36], which was suggested more recently. The difference between the two methods is that $A_{nys}^{(s,r)}$ performs rank-truncation in the core matrix, W , which makes this method cheaper to compute, while $\llbracket A_{nys}^{(s)} \rrbracket_r$ performs rank-truncation in the Nyström approximation $A_{nys}^{(s)}$, which makes this method take advantage of the full Nyström approximation, C and W , when performing the rank-truncation. There are also other variants of the Nyström method, including one for rectangular matrices [22, 33] and one that guarantees numerical stability [22]. This paper will mostly focus on $A_{nys}^{(s,r)}$.

It is known that for SPSD matrices, $A_{nys}^{(s)}$ [14] and $\llbracket A_{nys}^{(s)} \rrbracket_r$ [36] satisfy relative-error bounds in the nuclear norm. This means that if \hat{A} is a low-rank approximation to A (in this case, $A_{nys}^{(s)}$ or $\llbracket A_{nys}^{(s)} \rrbracket_r$) and $\epsilon > 0$ then

$$(1.1) \quad \left\| A - \hat{A} \right\|_* \leq (1 + \epsilon) \|A - \llbracket A \rrbracket_r\|_*$$

holds with high probability under some conditions on the sketch X and the sketch size $s > r$ where $\|\cdot\|_*$ is the nuclear norm (the sum of the singular values). The details are in the relevant papers [14, 36]. On the other hand, it is not known whether $A_{nys}^{(s,r)}$ satisfies a relative-error norm bound mentioned above [36]. In [28], an example of a 3×3 SPSD matrix is given, showing the downside of using $A_{nys}^{(s,r)}$ for kernel approximations which commonly uses a column sampling matrix. The authors propose $\llbracket A_{nys}^{(s)} \rrbracket_r$ ¹ as an alternative, for which later Wang, Gittens and Mahoney derived a relative-error norm bound [36]. For this example, the problem persists even if we use random embeddings. However, this is a small example that can yield results with high variability, and random embeddings do give a smaller expected relative-error in the nuclear norm and a smaller variance result than column sampling, especially when the dimension of the matrix is large. This hints that random embeddings can be more robust and reliable than column sampling. This type of phenomena have been discussed before, for example in [20] where the authors point out that column sampling is less reliable than random embeddings due to their relatively high variance results.

For symmetric indefinite matrices, which are the focus of this paper, not much has been shown. It is however known that the problem is rather difficult. We can easily see that the plain Nyström approximation, $A_{nys}^{(s)}$ can behave poorly for sym-

¹As in [28], for SPSD matrices, it should be noted that $\left\| A - \llbracket A_{nys}^{(s)} \rrbracket_r \right\|_* \leq \left\| A - A_{nys}^{(s,r)} \right\|_*$ will hold in the spectral norm and the Frobenius norm.

85 metric indefinite matrices.² In Figure 1, the two plots were generated using 100×100
 86 symmetric indefinite matrices with Haar distributed eigenvectors. In the left plot, the
 87 eigenvalues decay geometrically from 1 to 10^{-8} with random signs, and in the right
 88 plot, the first 20 eigenvalues are equal to ± 1 and the other 80 eigenvalues are equal
 89 to $\pm 10^{-10}$ where the signs were applied randomly with equal probability. We apply
 90 the plain Nyström approximation $A_{nys}^{(r)}$ using the Gaussian sketch to A . We can see
 that the plain Nyström approximation can be unstable. This type of issue has also

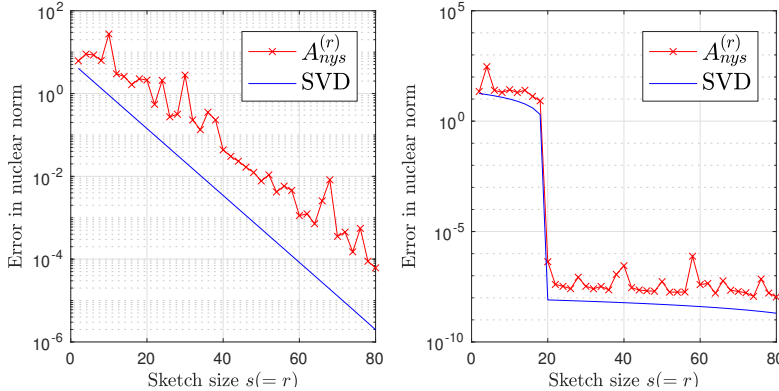


FIG. 1. Plain Nyström approximation $A_{nys}^{(r)}$ using the Gaussian sketch to 100×100 symmetric indefinite matrices. We can see that $A_{nys}^{(r)}$ can be unstable.

91 been observed in a different context for CUR approximations of rectangular matrices
 92 [6]. Essentially, the issue arises from the possible presence of an eigenvalue of X^TAX
 93 much smaller than $\sigma_r(A)$ or even $\sigma_{\min}(A)$ —a phenomenon that is absent when A is
 94 SPSD. This blows up the norm of the core matrix $(X^TAX)^\dagger$, causing instability.

95 *Contributions.* Our first contribution is to identify the main challenges in finding
 96 a good Nyström approximation to symmetric indefinite matrices. We find that the
 97 accuracy of the Nyström method is related to controlling the singular values of the
 98 core matrix $W = X^TAX$, and show that the accuracy can be lost even if the singular
 99 values of W are sufficiently larger than the unit roundoff if W severely underestimates
 100 the leading eigenvalues of A . We then perform an analysis in Section 3 that overcomes
 101 the challenges. The analysis shows that a certain truncation in the core matrix can
 102 give a reliable Nyström approximation that guarantees (1.1) to symmetric indefinite
 103 matrices when the singular values decay sufficiently quickly. To our knowledge, this
 104 is the first relative-error norm bound for the Nyström method concerning general
 105 symmetric matrices that are possibly indefinite.

106 Our second contribution is providing a practical algorithm (Algorithm 2.1) that
 107 gives a Nyström approximation to symmetric indefinite matrices. We show its robust-
 108 ness by comparing the algorithm to some of the existing methods in Section 4 and
 109 show that the algorithm performs robustly for symmetric indefinite matrices even in
 110 the presence of small singular values in the core matrix, whereas the other algorithms
 111 can fail. This algorithm is not new in the context of the Nyström method for SPSD
 112 matrices. However, to our knowledge, it has not been suggested or studied before for
 113 symmetric indefinite matrices.
 114

²For a contrived 2×2 example, where the plain Nyström approximation fails to succeed because the matrix is indefinite, see the arXiv version <https://arxiv.org/abs/2212.01127>.

115 *Existing methods.* We review three existing ideas for using the Nyström method
 116 for indefinite matrices, among others. Cai, Nagy and Xi [3] derive an error bound for
 117 the Nyström method, $A_{nys}^{(s)}$ for symmetric indefinite matrices that arise from a sym-
 118 metric function. This bound depends on how close the function values of the sampled
 119 points are, which is not an attractive dependence and may not be very useful in more
 120 general or practical situations. They suggest the plain Nyström method $A_{nys}^{(r)}$, which
 121 can be unstable. They also suggest $AX(X^TAX)_\epsilon^\dagger(AX)^T$ for the Nyström approxima-
 122 tion motivated by [22] with the aim of improving the stability. This version truncates
 123 the core matrix $W = X^TAX$ so that $\sigma_{\min}((X^TAX)_\epsilon) > \epsilon$ where ϵ is of the order of
 124 the unit roundoff. However, this version can give worse approximations than $A_{nys}^{(s)}$
 125 [3] and does not always improve the stability of the Nyström approximation. Second,
 126 Ray *et al.* [29] suggest submatrix-shifted (SMS) Nyström to provide an efficient al-
 127 gorithm that deals with symmetric matrices that have only few negative eigenvalues.
 128 This method uses an eigenvalue shift based on the minimum eigenvalue of a small
 129 principal submatrix before applying the plain Nyström method $A_{nys}^{(r)}$. The downside
 130 of this method is that the eigenvalue shift can have serious negative impact on the
 131 approximation quality. Lastly, the authors in [12, 26] devise strategies to form the
 132 Nyström approximation to symmetric indefinite matrices. However, these methods
 133 use eigenvalue information of the original matrix, which is expensive to compute. The
 134 three existing methods described above use column sampling matrices for X , which
 135 is different from random embeddings. In the final section (Section 5), we will revisit
 136 their differences in relation to our method and discuss the implications.

137 *Non-Nyström approaches.* In [15], a low-rank approximation for symmetric ma-
 138 trices in the form of the randomized SVD is given. This approximation is given by
 139 $QQ^T AQQ^T$ where $Q \in \mathbb{R}^{n \times s}$ is the orthonormal matrix in the thin QR decomposi-
 140 tion of AX and is known to satisfy a relative-error norm bound. The dominant cost
 141 is $O(n^2s)$ flops for forming $Q^T A$ (assuming A is dense), which becomes prohibitive
 142 when n, s are large. Wang, Luo and Zhang derived in [37] a relative-error norm bound
 143 to any symmetric matrices (possibly indefinite) for the prototype model. This model
 144 computes the low-rank approximation by first forming the sketch $C = AX$ and then
 145 approximating A by CXC^T where $X = C^\dagger A(C^\dagger)^T$. The authors show that if C
 146 contains $s = O(k/\epsilon)$ columns of A chosen by adaptive sampling then the prototype
 147 model has relative-error of at most $(1 + \epsilon)$. The dominant costs for the algorithm in
 148 [37] are $O(n^2r \log r)$ for computing C and $O(n^2r)$ for computing $C^\dagger A$, which becomes
 149 very costly with large n .

150 *Non-symmetric approaches.* We can use non-symmetric low-rank approximation
 151 to symmetric indefinite matrices. Examples are the randomized SVD [15], which is
 152 given by $QQ^T A$ using the notation in the previous paragraph and the generalized
 153 Nyström method [4, 22, 33] given by $AX(Y^TAX)^\dagger Y^T A$ where X and Y are indepen-
 154 dent random embeddings of different dimensions. The details can be found in the
 155 relevant papers. For both methods, since their representation is not symmetric, if we
 156 want to force symmetry in their representations (e.g. by taking the symmetric part
 157 $(M^T + M)/2$), we may risk doubling the rank in the approximation. In addition, as
 158 mentioned in the previous paragraph, the randomized SVD has the cost of comput-
 159 ing $Q^T A$, which becomes prohibitive when n, s are large. For generalized Nyström,
 160 we approximately double the number of matrix-vector multiplications needed as A
 161 needs to be multiplied by two independent random embeddings X and Y and this,
 162 in turn doubles the storage requirement (in fact, more than double because Y (or X)
 163 is recommended to be larger [22]). In this paper, we focus on symmetric low-rank

164 approximations.

165 *Notation.* Throughout, we use $\|\cdot\|_2$ for the spectral norm or the vector- ℓ_2 norm,
 166 $\|\cdot\|_*$ for the nuclear norm (sum of singular values) and $\|\cdot\|_F$ for the Frobenius norm. We
 167 use dagger \dagger to denote the pseudoinverse of a matrix and $\llbracket A \rrbracket_r$ to denote the best rank-
 168 r approximation to A in any unitarily invariant norm, i.e., the approximation derived
 169 from truncated SVD [16]. Unless specified otherwise, $\sigma_i(A)$ denotes the i th largest
 170 singular value of the matrix A and $\lambda_i(A)$ the i th largest eigenvalue in magnitude.
 171 Lastly, we use MATLAB style notation for matrices and vectors. For example, for
 172 the k th to $(k + j)$ th columns of a matrix A we write $A(:, k : k + j)$.

173 **2. Proposed method.** When we use the Nyström method on symmetric indefi-
 174 nite matrices, it can lead to problems. The main concern is in the core matrix
 175 $W = X^TAX$ because the positive and negative eigenvalues of A can ‘cancel’ each
 176 other out when forming W , making the eigenvalues of W much smaller than $\sigma_r(A)$.
 177 This causes inaccuracies and instabilities when computing the pseudo-inverse of W .
 178 More specifically, if we use column sampling then W would be a principal submatrix
 179 of A . By Cauchy’s interlacing theorem, the spectrum of W is contained in the interval
 180 $[\lambda_{\min}(A), \lambda_{\max}(A)]$ which contains both positive and negative values since A is indefi-
 181 nite. Therefore the magnitude of the eigenvalues of W can be significantly smaller in
 182 magnitude from those of A , resulting in the matrix W^\dagger blowing up. In addition, the
 183 computation of the pseudo-inverse of W can be numerically unstable if $\sigma_{\min}(W) < u$
 184 where u is the unit roundoff. Thus, the main challenge is to ensure that W^\dagger does not
 185 ruin the Nyström approximation quality. One approach is to introduce a potentially
 186 large shift to make A SPSD, but this can severely affect the approximation quality
 187 unless A is nearly definite, that is, the negative eigenvalues of A are very small in
 188 magnitude, for example, on the order of machine precision. This idea is used for
 189 SPSD matrices where a small shift is introduced to gain numerical stability, however
 190 the shift here needs to be small enough to ensure that accuracy is still high [17, 32].

191 In light of these observations, we propose

$$192 \quad A_{indef}^{(c,r)} = AX \llbracket X^TAX \rrbracket_r^\dagger (AX)^T$$

193 for symmetric *indefinite* matrices $A \in \mathbb{R}^{n \times n}$ where $X \in \mathbb{R}^{n \times cr}$ is a random embedding,
 194 $c > 1$ is a modest constant, say $c = 1.5$ or $c = 2$, and r is the target rank. When A is
 195 SPSD and the sketch size s is proportional to the target rank, $A_{indef}^{(c,r)}$ is equivalent to
 196 $A_{nys}^{(cr,r)}$. This rank-restricted version truncates the bottom $(c - 1)r$ singular values of
 197 $W \in \mathbb{R}^{cr \times cr}$, which can potentially be harmful even if they are sufficiently larger than
 198 the unit roundoff. This is different to the truncation used in [3] as they use truncation
 199 based on the magnitudes of the singular values of W , whereas for our method, the
 200 number of bottom singular values we truncate is proportional to the target rank. This
 201 intuition is justified by Andoni and Nguyễn [1], who prove that the largest eigenvalues
 202 (whose proportional to the sketch size) of symmetric matrices with rapidly decaying
 203 singular values are approximately preserved under conjugation by a Gaussian sketch
 204 with an appropriate normalization factor.

205 Now, let us define a quantity that will measure how well the singular values are
 206 preserved in the core matrix W of the Nyström method. For a symmetric matrix
 207 $A \in \mathbb{R}^{n \times n}$, a target rank r and a sketch size $s \geq r$, define

$$208 \quad (2.1) \quad \kappa_W(A, r, s) := \frac{\max_{1 \leq i \leq r} \sigma_i(X^TAX) / \sigma_i(A)}{\min_{1 \leq j \leq r} \sigma_j(X^TAX) / \sigma_j(A)} = \max_{1 \leq i \leq r} \max_{1 \leq j \leq r} \frac{\sigma_i(X^TAX) \sigma_j(A)}{\sigma_j(X^TAX) \sigma_i(A)}$$

209 where $X \in \mathbb{R}^{n \times s}$ is a Gaussian embedding matrix. This quantity measures the ratio
 210 between the worst over-approximation and the worst under-approximation of the lead-
 211 ing singular values of A using the singular values in the core matrix W . $\kappa_W(A, r, s)$
 212 will help us see how much the singular values of W have deviated from the leading
 213 singular values of A , which directly affects the Nyström approximation quality as we
 214 illustrate below.

215 In Figure 2, we show how important it is to ensure that the spectrum of W
 216 does not ruin the approximation quality. In this experiment³, $A \in \mathbb{R}^{1000 \times 1000}$ is a
 217 symmetric indefinite matrix constructed as in the left plot of Figure 1. The smallest
 218 singular value in the core matrix was larger than 10^{-7} throughout this experiment.
 219 For the truncated cases, $A_{indef}^{(1.5,r)}$ and $A_{nys}^{(r+5,r)}$, the approximation is robust as seen
 220 in Figure 2a. This robustness we see is illustrated in Figure 2b where the singular
 221 values of $W = X^T A X$ behaves well in the sense that there is no wild fluctuations in
 222 $\kappa_W(A, r, r+5)$ and $\kappa_W(A, r, 1.5r)$. However, when the sketch size is not proportional
 223 to the target rank ($s = r + 5$), the relative approximation error for $A_{nys}^{(r+5,r)}$ (when
 224 compared with the truncated SVD) and $\kappa_W(A, r, r+5)$ grow as we increase the target
 225 rank. This problem can become worse and the approximation can become unstable
 226 when we use SRTT matrices for efficiency with the sketch size $s = r + 5$ (See Figure
 227 3 and Subsection 2.1). When the sketch size is proportional to the target rank,
 228 $\kappa_W(A, r, 1.5r)$ and the relative approximation error for $A_{indef}^{(1.5,r)}$ are approximately a
 229 constant, which motivates us to choose the oversample size to be proportional to
 230 the target rank. On the other hand, without the truncation in the core matrix we
 231 see that $\kappa_W(A, r, r)$ behaves wildly. This indicates that the singular values of W
 232 inaccurately approximates the leading singular values of A . As a result, the Nyström
 233 approximations $A_{nys}^{(1.5r)}$ and $\llbracket A_{nys}^{(1.5r)} \rrbracket_r$ can yield unstable results. Empirically, this
 234 provides a reason to favour $A_{indef}^{(c,r)}$ over other variants of the Nyström method for
 235 symmetric indefinite matrices.

236 **2.1. Random embeddings.** A subspace embedding [30] is a linear map which
 237 preserves the 2-norm of every vector in a given subspace, that is, $S \in \mathbb{R}^{s \times n}$ is a
 238 subspace embedding for the span of $A \in \mathbb{R}^{n \times n}$ with distortion $\epsilon \in (0, 1)$ if

$$239 \quad (2.2) \quad (1 - \epsilon) \|Ax\|_2 \leq \|SAx\|_2 \leq (1 + \epsilon) \|Ax\|_2$$

240 for every $x \in \mathbb{R}^n$. A random embedding is a subspace embedding drawn at random
 241 that satisfy Equation (2.2) with high probability.

242 Random embeddings have more attractive properties than column sampling ma-
 243 trices [11, 20], one of which is that the results obtained using random embeddings
 244 generally have smaller variance than the results obtained using column sampling.
 245 Below are few important examples of random embeddings.

246 **2.1.1. Gaussian matrices.** A Gaussian embedding is a random matrix $G \in$
 247 $\mathbb{R}^{s \times n}$ with i.i.d. entries $G_{ij} \sim N(0, 1/s)$. The scaling ensures that $\mathbb{E}[\|Gx\|_2^2] = \|x\|_2^2$
 248 for every $x \in \mathbb{R}^n$. Gaussian embedding is the most widely used random embedding for
 249 theoretical analysis⁴ and often has optimal guarantees [15, 20]. The cost of applying
 250 a Gaussian embedding to an $n \times n$ matrix is $O(n^2 s)$. This becomes prohibitive for
 251 large n , so a more structured random embeddings are often used in practice.

³All experiments were performed in MATLAB version 2021a using double precision arithmetic.

⁴Other random embeddings often lack strong theoretical guarantees, however they behave similarly to a Gaussian embedding in practice. For this reason, Gaussian theory is often used to provide a rule of thumb for the general behavior [20].

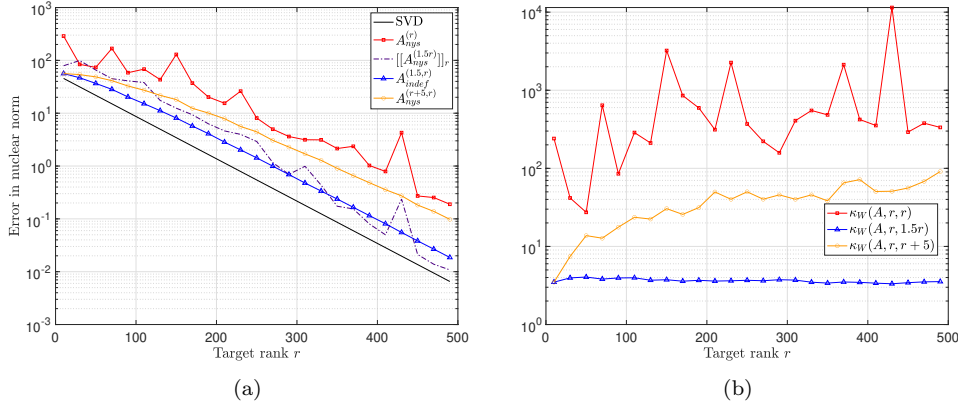


FIG. 2. Accuracy of the Nyström approximations $A_{indef}^{(1.5,r)}$, $A_{nys}^{(r)}$, $A_{nys}^{(r+5,r)}$ and $\llbracket A_{nys}^{(1.5r)} \rrbracket_r$ to a symmetric indefinite matrix $A \in \mathbb{R}^{1000 \times 1000}$. Figure 2a shows the Nyström error in the nuclear norm and Figure 2b shows the accuracy of the singular values of $W = X^T A X$ when compared with the singular values of A . We observe that the truncation in the core matrix W can significantly increase the robustness and the accuracy of the Nyström approximation.

252 **2.1.2. SRTTs.** A subsampled randomized trigonometric transform (SRTT) ma-
 253 trix is an $n \times s$ matrix with $n \geq s$ of the form

$$254 \quad (2.3) \quad S = \sqrt{\frac{n}{s}} D F R^T$$

255 where $D \in \mathbb{R}^{n \times n}$ is a random diagonal matrix whose entries are independent and
 256 take ± 1 with equal probability, $F \in \mathbb{C}^{n \times n}$ is a unitary trigonometric transform and
 257 $R \in \mathbb{R}^{s \times n}$ is a random restriction. In the complex case, F is the unitary discrete
 258 Fourier transform (DFT) and in the real case, F is commonly the discrete cosine
 259 transform (DCT). The sketch size needs to be $s = O(r \log r)$ for theoretical guarantees
 260 [31], but in practice $s = O(r)$ often suffices⁵ [15, 20]. The cost of applying SRTT to
 261 an $n \times n$ matrix is $O(n^2 \log r)$ [2] using the subsampled FFT algorithm [40].

262 **2.1.3. Sparse maps.** Sparse maps are sparse matrices with nonzero entries that
 263 are random signs [4, 20, 23, 39]. They are particularly useful for sparse data and they
 264 take the form

$$265 \quad (2.4) \quad S = \frac{1}{\sqrt{s}} [s_1, \dots, s_n] \in \mathbb{R}^{s \times n}$$

266 where the columns of S , the s_i 's are statistically independent and has exactly ξ
 267 nonzero entries that take ± 1 with equal probability, placed uniformly at random
 268 coordinates. We need the sketch size to be $s = O(r \log r)$ and the sparsity parameter to
 269 be $\xi = O(\log r)$ for theoretical guarantees [5]. In [34], $\xi = \min\{s, 8\}$ was recommended
 270 in practice. The cost of applying sparse maps to a matrix A is $O(\xi \cdot nnz(A))$ where
 271 $nnz(A)$ is the number of nonzero entries of A if sparse data structures and arithmetic
 272 are available.

⁵For difficult examples, say a coherent example, the $\log r$ factor is necessary. (See Figure 3)

273 **2.2. Suggested algorithm.** For a general symmetric matrix $A \in \mathbb{R}^{n \times n}$ with
 274 the target rank r , we suggest

$$275 \quad (2.5) \quad A_{indef}^{(c,r)} = AX[X^TAX]_r^\dagger(AX)^T = C[W]_r^\dagger C^T$$

276 where $X \in \mathbb{R}^{n \times s}$ is a random embedding with the sketch size $s = cr$ where $c > 1$ is a
 277 modest constant. The algorithm is given in Algorithm 2.1. For the choice of random
 278 embeddings, if A is sparse then we suggest sparse maps with sparsity $\xi = \min\{cr, 8\}$
 279 and when A is dense we suggest SRTT matrices. The recommended sketch size is
 280 $s = 1.5r$ for efficiency, but if one wants a better approximation quality guarantee
 281 then the sketch size can be increased to, for example, $s = 2r$ or $s = 4r$. Note that
 282 the truncation is performed irrespectively of the singular values of W (unlike previous
 283 studies, e.g. [3]); our analysis in Section 3 suggests that it is important that the
 284 number of singular values to be truncated $(s - r) = (c - 1)r$ is proportional to r .

Algorithm 2.1 Judiciously truncated Nyström approximation for indefinite matrices

Require: Symmetric matrix $A \in \mathbb{R}^{n \times n}$, target rank $r < n$, sketch size $r < s < n$ (rec. $s = 1.5r$)

Ensure: $C \in \mathbb{R}^{n \times s}$ and $W_r^\dagger \in \mathbb{R}^{s \times s}$ with $\text{rank}(W) \leq r$ as in (2.5)

- 1: Draw a random embedding $X \in \mathbb{R}^{n \times s}$ ▷ Sparsity $\xi = \min\{s, 8\}$ for sparse maps
 - 2: $C \leftarrow AX$
 - 3: $W \leftarrow X^T C$
 - 4: $[V, \Lambda] = \text{eig}(W)$, eigendecomposition of W
 - 5: $W_r^\dagger = V(:, 1:r)\Lambda(1:r, 1:r)^\dagger V(:, 1:r)^T$, pseudoinverse of the best rank- r approximation of W
 - 6: Output $C \in \mathbb{R}^{n \times s}$ and $W_r^\dagger \in \mathbb{R}^{s \times s}$
-

285 *Complexity.* When a sparse map is used, the cost of Algorithm 2.1 is $O(\xi \cdot \text{nnz}(A) +$
 286 $r^3)$ which consists of $O(\xi \cdot \text{nnz}(A))$ flops for forming the sketch and $O(r^3)$ flops for the
 287 eigendecomposition. With an SRTT sketch, the total cost is $O(n^2 \log r + r^3)$, where
 288 $O(n^2 \log r)$ is needed for forming the sketch and $O(r^3)$ for computing the eigendecom-
 289 position.⁶

290 *Eigendecomposition of $A_{indef}^{(c,r)}$.* Algorithm 2.1 as presented does not output the
 291 eigendecomposition of $A_{indef}^{(c,r)}$. To do this, we require an extra $O(nr^2 + r^3)$ flops. We
 292 need $O(nr^2)$ flops to compute the thin QR decomposition of $C = QR$, $O(r^3)$ flops to
 293 form and compute the eigendecomposition of $R[W]_r^\dagger R^T = U\Sigma U^T$ and $O(nr^2)$ flops
 294 to form $U_1 = QU$ giving us the eigendecomposition, $A_{indef}^{(c,r)} = U_1 \Sigma U_1^T$.

295 In Figure 3, we illustrate Algorithm 2.1 for the SRFT sketch and the sparse
 296 map. The experiment was conducted with synthetic 2000×2000 symmetric indefinite
 297 matrices. The top two plots have eigenvalues that decay geometrically from 1 to
 298 10^{-12} each assigned a random sign with equal probability and the eigenvectors are
 299 in a 2×2 block diagonal form, $\text{diag}(I_{200}, U)$ where I_{200} is the 200×200 identity
 300 matrix and $U \in \mathbb{R}^{1800 \times 1800}$ is a Haar distributed orthogonal matrix. This eigenvector
 301 matrix is a more coherent example than our previous examples and is known to be a
 302 difficult example for SRTT matrices [2] (when the eigenvectors are Haar distributed,
 303 SRTT (or essentially any sketch) behaves the same as a Gaussian sketch, giving good
 304 results). The bottom two plots were generated using the same eigenvector matrix, but
 305 with eigenvalues equal to ± 1 for the first 100, $\pm 10^{-4}$ for the next 100, $\pm 10^{-8}$ for the
 306 100 eigenvalues after that and $\pm 10^{-16}$ for the last 1700 eigenvalues each assigned a

⁶Since we are using random embeddings for robustness, Algorithm 2.1 is strictly more expensive than classical Nyström methods (column subsampling) if the columns can be sampled quickly.

307 random sign with equal probability. In the two left plots, we see that the SRFT sketch
 308 can fail if the sketch size is not large enough. This instability in the approximation
 309 can be fixed by enlarging the sketch size. We see that $s = r + 5$ does not do well,
 310 but when $s = 4r$ the approximation becomes more accurate and robust. In the right
 311 plot, we see that the SRFT sketch with the sketch size $s = r \log r$, which comes with
 312 theoretical guarantees has excellent approximation quality. Finally, we see that the
 313 sparse map with sparsity $\xi = 8$ gives a robust approximation throughout, which can
 314 be improved by enlarging the sketch size.

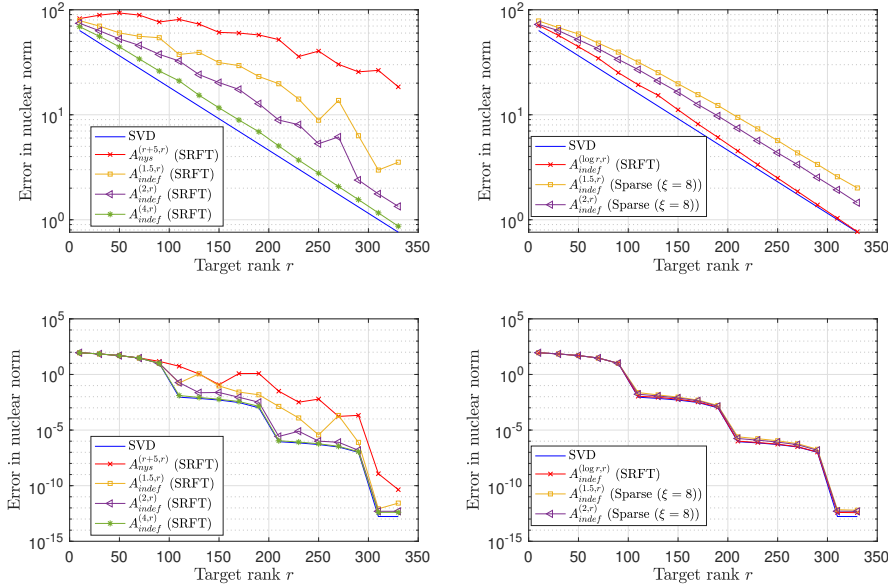


FIG. 3. Algorithm 2.1: A difficult (coherent) example for the SRFT sketch. The approximation can be unstable if the sketch size is too small for the SRFT sketch (left plots). This problem can be fixed by enlarging the sketch size. The right plots show that sparse maps have no issue with this example and the approximation is robust.

315 **3. Analysis.** For a general symmetric matrix $A \in \mathbb{R}^{n \times n}$, there are no known
 316 relative-error norm bounds for the Nyström method. Here we show that for general
 317 symmetric matrices, the Nyström method when used with a Gaussian sketch satisfies
 318 in expectation a relative-error nuclear norm bound under some orthogonal projection
 319 in the core matrix, when the singular values decay sufficiently fast. The analysis that
 320 follows establishes the accuracy not of Algorithm 2.1, but of a closely related variant
 321 of the Nyström method. The last paragraph of this section discusses this in more
 322 detail.

323 Let $A \in \mathbb{R}^{n \times n}$ be a symmetric matrix and let the eigendecomposition of A be

$$324 \quad (3.1) \quad A = V \Lambda V^T = [V_1, V_2, V_3] \begin{bmatrix} \Lambda_1 & 0 & 0 \\ 0 & \Lambda_2 & 0 \\ 0 & 0 & \Lambda_3 \end{bmatrix} [V_1, V_2, V_3]^T$$

325 where $V \in \mathbb{R}^{n \times n}$ is the orthogonal eigenvector matrix of A and $\Lambda \in \mathbb{R}^{n \times n}$ is a diagonal
 326 matrix containing the eigenvalues of A . The matrices with subscript 1 have r columns,

327 those with subscript 2 have $(c_1 - 1)r$ columns and subscript 3 have $(n - c_1r)$ columns
 328 where $r < c_1r < n$ and $c_1 > 1$ is a constant such that c_1r is a positive integer. The
 329 eigenvalues are ordered in non-increasing order with respect to their magnitude, so
 330 we have $\sigma_i(A) = |\lambda_i(A)|$ for all i .

331 Now we state our main theorem, and discuss the three key facts that will accom-
 332 pany our proof before getting to the proof immediately.

333 **THEOREM 3.1.** *Let $A \in \mathbb{R}^{n \times n}$ be a symmetric matrix as in (3.1) and assume that*
 334 *$\lambda_r(A) \neq 0$. Let c_1 and c_2 be constants with $1 < c_1 < c_2 < \frac{n}{r} - 1$ such that c_1r*
 335 *and c_2r are positive integers. Define $X_i := V_i^T X$ for $i = 1, 2, 3$ where $X \in \mathbb{R}^{n \times c_2r}$*
 336 *is a Gaussian matrix, and set $B = X_3 Q_\perp (X_1 Q_\perp)^\dagger \in \mathbb{R}^{(n - c_1r) \times r}$ where $Q_\perp \in$*
 337 *$\mathbb{R}^{c_2r \times (c_2 - c_1 + 1)r}$ is an orthogonal complement of $X_2^T \in \mathbb{R}^{c_2r \times (c_1 - 1)r}$. Let $(X_1 Q_\perp)^\dagger =$*
 338 *$\hat{Q} \hat{R}$ be the thin QR decomposition of $(X_1 Q_\perp)^\dagger$ and set $U := Q_\perp \hat{Q} \in \mathbb{R}^{c_2r \times r}$. Then*
 339 *the orthogonal projector $P = UU^T \in \mathbb{R}^{c_2r \times c_2r}$ satisfies*

$$340 \quad (3.2) \quad \mathbb{E} [\|E\|_* | \Omega_F] \leq (1 + \epsilon_{r,A}) \|A - [A]_r\|_*$$

341 where

$$342 \quad (3.3) \quad E := A - AX(PX^TAXP)^\dagger X^T A$$

343 is the associated Nyström error, Ω_F is an event defined as

$$344 \quad (3.4) \quad \Omega_F := \left\{ \left\| \|\Lambda_3\|^{1/2} B \right\|_F^2 \leq 0.5 |\lambda_r(A)| \right\}$$

345 where $|\Lambda_3|$ is defined element-wise and

$$346 \quad (3.5) \quad \epsilon_{r,A} := 2b\sqrt{r} \left(1 + \frac{|\lambda_{c_1r+1}(A)|}{|\lambda_r(A)|} + \frac{2}{\sqrt{b}} \right) \frac{\|\Lambda_3\|_*}{\|\Lambda_2\|_* + \|\Lambda_3\|_*}$$

347 where $b = \frac{r}{(c_2 - c_1)r - 1}$.

348 In the above theorem, c_1 and c_2 are oversampling factors which are of modest size,
 349 say $c_1 = 1.5$ and $c_2 = 2$. We need two factors because we need $X_1 Q_\perp \in \mathbb{R}^{r \times (c_2 - c_1 + 1)r}$
 350 and $X_3 Q_\perp \in \mathbb{R}^{(n - c_1r) \times (c_2 - c_1 + 1)r}$ to be *rectangular* Gaussian matrices, which makes
 351 them well-conditioned with high probability [7]. We can view c_1 as c in Algorithm 2.1
 352 and c_2 to be the oversampling factor introduced to make the analysis possible. By
 353 making c_1 , c_2 and $(c_2 - c_1)$ larger, we can improve the bound in the above Theorem.
 354 The orthogonal projector $P = UU^T$ truncates the core matrix $W = X^T A X$ by remov-
 355 ing the largest ‘unwanted’ eigenvalues of A , i.e. the eigenvalues in Λ_2 , using X_\perp factor
 356 in U . This helps the core matrix to not be corrupted by the interaction between the
 357 target and the large ‘unwanted’ singular values and singular vectors of A , which can
 358 happen when forming $X^T A X$. Lastly, the $\epsilon_{r,A}$ in the theorem plays a similar role to
 359 the distortion ϵ in Equation (1.1) and Ω_F is roughly the event that the eigenvalues of
 360 A decay rapidly enough. If we assume that A has a low-rank structure, for example,
 361 $|\lambda_r(A)| \gg |\lambda_{c_1r+1}(A)|$, then Ω_F would hold with high probability and $\epsilon_{r,A}$ would be a
 362 moderately-sized constant, which tells us that the relative-error nuclear norm bound
 363 in (3.2) is good.

364 We now introduce three key facts that will be useful for our proof. The first fact
 365 follows closely the analysis in [22]. Let $\mathcal{P} := \Lambda V^T X (P X^T A X P)^\dagger X^T V$ be an oblique
 366 projector. Then we can rewrite the associated Nyström error as

$$367 \quad (3.6) \quad E = V(I - \mathcal{P})\Lambda V^T.$$

368 As shown in [22], it is straightforward to see that we can rewrite the associated
 369 Nyström error as

$$370 \quad (3.7) \quad V^T E V = (I - \mathcal{P})\Lambda = (I - \mathcal{P})\Lambda(I - V^T X U M)$$

371 for any $M \in \mathbb{R}^{r \times n}$. Let $V_r = [I_r, 0]^T \in \mathbb{R}^{n \times r}$ and set $M = (V_r^T V^T X U)^\dagger V_r^T$ then we
 372 get

$$373 \quad (3.8) \quad V^T E V = (I - \mathcal{P})\Lambda(I - V_r V_r^T)(I - V^T X U (V_r^T V^T X U)^\dagger V_r^T).$$

374 This modification of the associated Nyström error will be important for our proof.

375 The second fact is the following. Let $f(x)$ be convex in the interval $[x_1, x_2]$ with
 376 $x_1 < x_2$. Define $g(x)$ on $[x_1, x_2]$ to be the linear function joining the endpoints of f on
 377 $[x_1, x_2]$, that is, $g(x) = \frac{f(x_2) - f(x_1)}{x_2 - x_1}x + \frac{f(x_1)x_2 - f(x_2)x_1}{x_2 - x_1}$. Then $f(x) \leq g(x)$ on $[x_1, x_2]$.
 378 Let Y be a random variable with $Y \in [x_1, x_2]$ almost surely. Then $f(Y) \leq g(Y)$ almost
 379 surely. Furthermore, if $Y \in [x_1, x_2]$ conditional on an event Ω , then conditional on Ω
 380 we get

$$381 \quad (3.9) \quad f(Y) \leq g(Y).$$

382 The last fact is based on expected norm bounds for Gaussian matrices from [15,
 383 App. A]. We can deduce the following lemma.

384 **LEMMA 3.2.** *Let B be the matrix as in Theorem 3.1 and let S be a fixed real*
 385 *matrix such that SB is defined. Then*

$$386 \quad (3.10) \quad \mathbb{E} \|SB\|_F^2 = b \|S\|_F^2$$

387 where $b = \frac{r}{(c_2 - c_1)r - 1}$ as in Theorem 3.1.

388 *Proof.* Since conditional on X_2 , $X_3 Q_\perp$ and $X_1 Q_\perp$ are two independent Gaussian
 389 matrices, we have

$$\begin{aligned} 390 \quad \mathbb{E}_{X_1, Q_\perp, X_3} \|SB\|_F^2 &= \mathbb{E}_{X_1, Q_\perp} \left[\mathbb{E}_{X_3} \left[\|S X_3 Q_\perp (X_1 Q_\perp)^\dagger\|_F^2 \mid X_1, X_2 \right] \right] \\ 391 &= \|S\|_F^2 \mathbb{E}_{X_1, Q_\perp} \|(X_1 Q_\perp)^\dagger\|_F^2 \\ 392 &= \frac{r}{(c_2 - c_1)r - 1} \|S\|_F^2 \\ 393 \end{aligned}$$

394 using the tower property and the propositions in [15, App. A]. □

395 Now using these three key facts we are ready to prove Theorem 3.1.

396 *Proof of Theorem 3.1.* Since U is an orthonormal matrix we have

$$\begin{aligned} 397 \quad AX(PX^TAXP)^\dagger X^T A &= AX(UU^T X^T AXUU^T)^\dagger X^T A \\ 398 &= AXU(U^T X^T AXU)^\dagger U^T X^T A. \end{aligned}$$

400 Now since $X_1 Q_\perp \in \mathbb{R}^{r \times (c_2 - c_1 + 1)r}$ is a fat rectangular Gaussian matrix, hence full
 401 rank with probability 1, we have $X_1 Q_\perp (X_1 Q_\perp)^\dagger = I_r$. Therefore

$$402 \quad (3.11) \quad X_1 Q_\perp \hat{Q} = \hat{R}^{-1}$$

403 and we get

$$404 \quad (3.12) \quad V^T X U = \begin{bmatrix} X_1 \\ X_2 \\ X_3 \end{bmatrix} Q_\perp \hat{Q} = \begin{bmatrix} X_1 Q_\perp \hat{Q} \\ 0 \\ X_3 Q_\perp \hat{Q} \end{bmatrix} = \begin{bmatrix} \hat{R}^{-1} \\ 0 \\ B \hat{R}^{-1} \end{bmatrix}$$

405 where $B = X_3 Q_\perp (X_1 Q_\perp)^\dagger \in \mathbb{R}^{(n-c_1 r) \times r}$.

406 Now we use the first key fact (Equation (3.8)) and get

$$407 \quad (3.13) \quad V^T E V = (I - \mathcal{P}) \Lambda (I - V_r V_r^T) (I - V^T X U (V_r^T V^T X U)^\dagger V_r^T)$$

408 where $\mathcal{P} = \Lambda V^T X (P X^T A X P)^\dagger X^T V$ and V_r is as below. Using

$$409 \quad (3.14) \quad V^T X U = \begin{bmatrix} \hat{R}^{-1} \\ 0 \\ B \hat{R}^{-1} \end{bmatrix}, \Lambda = \begin{bmatrix} \Lambda_1 & 0 & 0 \\ 0 & \Lambda_2 & 0 \\ 0 & 0 & \Lambda_3 \end{bmatrix}, V_r = \begin{bmatrix} I_r \\ 0 \end{bmatrix}$$

410 we get

$$411 \quad V^T E V = (I - \mathcal{P}) \Lambda \begin{bmatrix} 0 \\ I_{n-r} \end{bmatrix} [0, I_{n-r}] \left(I - \begin{bmatrix} \hat{R}^{-1} \\ 0 \\ B \hat{R}^{-1} \end{bmatrix} (\hat{R}^{-1})^\dagger [I_r, 0] \right) \\ 412 \quad = (I - \mathcal{P}) \begin{bmatrix} 0 & 0 & 0 \\ 0 & \Lambda_2 & 0 \\ -\Lambda_3 B & 0 & \Lambda_3 \end{bmatrix}. \\ 413$$

414 We also get

$$415 \quad \mathcal{P} = \begin{bmatrix} \Lambda_1 \hat{R}^{-1} \\ 0 \\ B \hat{R}^{-1} \end{bmatrix} \left(\begin{bmatrix} \hat{R}^{-1} \\ 0 \\ B \hat{R}^{-1} \end{bmatrix}^T \begin{bmatrix} \Lambda_1 \hat{R}^{-1} \\ 0 \\ \Lambda_3 B \hat{R}^{-1} \end{bmatrix} \right)^\dagger \begin{bmatrix} \hat{R}^{-1} \\ 0 \\ B \hat{R}^{-1} \end{bmatrix}^T \\ 416 \quad = \begin{bmatrix} \Lambda_1 \hat{R}^{-1} \\ 0 \\ \Lambda_3 B \hat{R}^{-1} \end{bmatrix} \left(\hat{R}^{-T} (\Lambda_1 + B^T \Lambda_3 B) \hat{R}^{-1} \right)^\dagger \begin{bmatrix} \hat{R}^{-1} \\ 0 \\ B \hat{R}^{-1} \end{bmatrix}^T \\ 417 \quad = \begin{bmatrix} \Lambda_1 \\ 0 \\ \Lambda_3 B \end{bmatrix} (\Lambda_1 + B^T \Lambda_3 B)^\dagger [I_r, 0, B^T] \\ 418$$

419 by taking out a factor of \hat{R}^{-1} and \hat{R}^{-T} from the pseudo-inverse. This is possible

420 because if we condition on Ω_F then $(\Lambda_1 + B^T \Lambda_3 B)$ is a non-singular $r \times r$ matrix.

421 Now for shorthand let $S := \Lambda_1 + B^T \Lambda_3 B$. Then

$$422 \quad I - \mathcal{P} = \begin{bmatrix} I_r - \Lambda_1 S^\dagger & 0 & -\Lambda_1 S^\dagger B^T \\ 0 & I_{(c_1-1)r} & 0 \\ -\Lambda_3 B S^\dagger & 0 & I_{n-c_1 r} - \Lambda_3 B S^\dagger B^T \end{bmatrix}. \\ 423$$

424 Therefore

$$\begin{aligned}
425 \quad V^T E V &= (I - \mathcal{P}) \begin{bmatrix} 0 & 0 & 0 \\ 0 & \Lambda_2 & 0 \\ -\Lambda_3 B & 0 & \Lambda_3 \end{bmatrix} \\
426 \quad &= \begin{bmatrix} \Lambda_1 S^\dagger B^T \Lambda_3 B & 0 & -\Lambda_1 S^\dagger B^T \Lambda_3 \\ 0 & \Lambda_2 & 0 \\ -\Lambda_3 B + \Lambda_3 B S^\dagger B^T \Lambda_3 B & 0 & \Lambda_3 - \Lambda_3 B S^\dagger B^T \Lambda_3 \end{bmatrix} \\
427 \quad &= \begin{bmatrix} \Lambda_1 S^\dagger B^T \Lambda_3 B & 0 & -\Lambda_1 S^\dagger B^T \Lambda_3 \\ 0 & \Lambda_2 & 0 \\ -\Lambda_3 B S^\dagger \Lambda_1 & 0 & \Lambda_3 - \Lambda_3 B S^\dagger B^T \Lambda_3 \end{bmatrix}. \\
428 \quad &
\end{aligned}$$

429 We now bound E in the nuclear norm. For shorthand, define the following

$$\begin{aligned}
430 \quad a_1 &= \Lambda_1 S^\dagger B^T \Lambda_3 B \\
431 \quad a_2 &= \Lambda_1 S^\dagger B^T \Lambda_3 \\
432 \quad a_3 &= \Lambda_3 - \Lambda_3 B S^\dagger B^T \Lambda_3. \\
433 \quad &
\end{aligned}$$

434 We then have

$$435 \quad \|E\|_* \leq \|\Lambda_2\|_* + \|a_1\|_* + 2\|a_2\|_* + \|a_3\|_*.$$

437 Let us note

$$438 \quad (3.15) \quad \Lambda_1 S^\dagger = \Lambda_1 (\Lambda_1 + B^T \Lambda_3 B)^\dagger = (I_r + B^T \Lambda_3 B \Lambda_1^{-1})^\dagger$$

439 conditional on Ω_F since $\lambda_r(A) \neq 0$ and

$$440 \quad (3.16) \quad \|B^T \Lambda_3 B\|_F = \|B^T |\Lambda_3|^{1/2} \operatorname{sgn}(\Lambda_3) |\Lambda_3|^{1/2} B\|_F \leq \| |\Lambda_3|^{1/2} B \|_F^2$$

441 where $|\Lambda_3|$ and $\operatorname{sgn}(\Lambda_3)$ are defined element-wise.

442 We now bound $\mathbb{E}[\|a_1\|_* | \Omega_F]$, $\mathbb{E}[\|a_2\|_* | \Omega_F]$ and $\mathbb{E}[\|a_3\|_* | \Omega_F]$ using the second
443 (Equation (3.9)) and the third (Lemma 3.2) key fact. We start with a_1 . Conditional
444 on Ω_F , we have

$$\begin{aligned}
445 \quad \|a_1\|_* &\leq \sqrt{r} \|\Lambda_1 S^\dagger B^T \Lambda_3 B\|_F \\
446 \quad &\leq \sqrt{r} \left\| (I_r + B^T \Lambda_3 B \Lambda_1^{-1})^\dagger \right\|_2 \|B^T \Lambda_3 B\|_F \\
447 \quad &\leq \sqrt{r} \frac{\|B^T \Lambda_3 B\|_F}{1 - \|B^T \Lambda_3 B\|_F \|\Lambda_1^{-1}\|_2} \\
448 \quad &\leq \frac{\sqrt{r}}{\|\Lambda_1^{-1}\|_2} \frac{\| |\Lambda_3|^{1/2} B \|_F^2 \|\Lambda_1^{-1}\|_2}{1 - \| |\Lambda_3|^{1/2} B \|_F^2 \|\Lambda_1^{-1}\|_2} \\
449 \quad &\leq \frac{\sqrt{r}}{\|\Lambda_1^{-1}\|_2} \left(2 \| |\Lambda_3|^{1/2} B \|_F^2 \|\Lambda_1^{-1}\|_2 \right) \\
450 \quad &
\end{aligned}$$

451 where the last inequality was obtained using the second fact with $Y =$
452 $\| |\Lambda_3|^{1/2} B \|_F^2 \|\Lambda_1^{-1}\|_2$, the event Ω_F , the interval $[0, 0.5]$, $f(x) = \frac{x}{1-x}$ which is con-
453 vex on $[0, 0.5]$ and $g(x) = 2x$. Now taking conditional expectation and using the third

454 fact (Lemma 3.2) we get

$$\begin{aligned}
455 \quad \mathbb{E}[\|a_1\|_* | \Omega_F] &\leq 2\sqrt{r}\mathbb{E}\left[\left\|\left|\Lambda_3\right|^{1/2}B\right\|_F^2 \middle| \Omega_F\right] \\
456 \quad &= 2\sqrt{r}b\left\|\left|\Lambda_3\right|^{1/2}\right\|_F^2 \\
457 \quad &= 2\sqrt{r}b\|\Lambda_3\|_*.
\end{aligned}$$

459 For a_2 , it is similar to a_1 . Conditional on Ω_F we have

$$\begin{aligned}
460 \quad \|a_2\|_* &\leq \sqrt{r}\|\Lambda_1(\Lambda_1 + B^T\Lambda_3B)^\dagger B^T\Lambda_3\|_F \\
461 \quad &\leq \sqrt{r}\|(I_r + B^T\Lambda_3B\Lambda_1^{-1})^\dagger\|_2\left\|\left|\Lambda_3\right|^{1/2}B\right\|_F\left\|\left|\Lambda_3\right|^{1/2}\right\|_F \\
462 \quad &\leq \frac{\sqrt{r}\sqrt{\|\Lambda_3\|_*}\left\|\left|\Lambda_3\right|^{1/2}B\right\|_F\sqrt{\|\Lambda_1^{-1}\|_2}}{\sqrt{\|\Lambda_1^{-1}\|_2} - \left\|\left|\Lambda_3\right|^{1/2}B\right\|_F\|\Lambda_1^{-1}\|_2} \\
463 \quad &\leq \frac{\sqrt{r}\sqrt{\|\Lambda_3\|_*}2\left\|\left|\Lambda_3\right|^{1/2}B\right\|_F\sqrt{\|\Lambda_1^{-1}\|_2}}{\sqrt{\|\Lambda_1^{-1}\|_2}}
\end{aligned}$$

465 where we used the second fact for the last inequality with $Y = \left\|\left|\Lambda_3\right|^{1/2}B\right\|_F\sqrt{\|\Lambda_1^{-1}\|_2}$,
466 the interval $[0, \sqrt{0.5}]$, $f(x) = \frac{x}{1-x^2}$ and $g(x) = 2x$. Therefore we get

$$\begin{aligned}
467 \quad \mathbb{E}[\|a_2\|_* | \Omega_F] &\leq 2\sqrt{r}\sqrt{\|\Lambda_3\|_*}\mathbb{E}\left[\left\|\left|\Lambda_3\right|^{1/2}B\right\|_F \middle| \Omega_F\right] \\
468 \quad &\leq 2\sqrt{r}\sqrt{\|\Lambda_3\|_*}\sqrt{\mathbb{E}\left[\left\|\left|\Lambda_3\right|^{1/2}B\right\|_F^2 \middle| \Omega_F\right]} \\
469 \quad &\leq 2\sqrt{r}b\|\Lambda_3\|_*.
\end{aligned}$$

471 using Lemma 3.2.

472 Finally for a_3 , we get

$$473 \quad \|a_3\|_* \leq \|\Lambda_3\|_* + \sqrt{r}\left\|\left|\Lambda_3\right|^{1/2}\right\|_2^2\left\|\left|\Lambda_3\right|^{1/2}B\right\|_F^2\|(\Lambda_1 + B^T\Lambda_3B)^\dagger\|_2$$

474 in a similar manner, and conditional on Ω_F we have

$$\begin{aligned}
475 \quad \left\|\left|\Lambda_3\right|^{1/2}B\right\|_F^2\|(\Lambda_1 + B^T\Lambda_3B)^\dagger\|_2 &\leq \left\|\left|\Lambda_3\right|^{1/2}B\right\|_F^2\|\Lambda_1^{-1}\|_2\|\Lambda_1(\Lambda_1 + B^T\Lambda_3B)^\dagger\|_2 \\
476 \quad &\leq \frac{\left\|\left|\Lambda_3\right|^{1/2}B\right\|_F^2\|\Lambda_1^{-1}\|_2}{1 - \left\|\left|\Lambda_3\right|^{1/2}B\right\|_F^2\|\Lambda_1^{-1}\|_2} \\
477 \quad &\leq 2\left\|\left|\Lambda_3\right|^{1/2}B\right\|_F^2\|\Lambda_1^{-1}\|_2
\end{aligned}$$

479 using the second fact with the same values as the a_1 case. Therefore

$$\begin{aligned}
480 \quad \mathbb{E}[\|a_3\|_* | \Omega_2] &\leq \|\Lambda_3\|_* + 2\sqrt{r}\|\Lambda_3\|_2\|\Lambda_1^{-1}\|_2\mathbb{E}\left[\left\|\left|\Lambda_3\right|^{1/2}B\right\|_F^2 \middle| \Omega_F\right] \\
481 \quad &\leq \|\Lambda_3\|_* + 2\sqrt{r}b\|\Lambda_3\|_2\|\Lambda_1^{-1}\|_2\|\Lambda_3\|_*.
\end{aligned}$$

483 Finally, combining everything together we get

$$484 \quad (3.17) \quad \mathbb{E} [\|E\|_* | \Omega_F] \leq \|\Lambda_2\|_* + \|\Lambda_3\|_* + 2b\sqrt{r} \left(1 + \frac{|\lambda_{c_1 r+1}(A)|}{|\lambda_r(A)|} + \frac{2}{\sqrt{b}} \right) \|\Lambda_3\|_*.$$

485 Therefore

$$486 \quad (3.18) \quad \mathbb{E} [\|E\|_* | \Omega_F] \leq (1 + \epsilon_{r,A}) (\|\Lambda_2\|_* + \|\Lambda_3\|_*) = (1 + \epsilon_{r,A}) \|A - \llbracket A \rrbracket_r\|_*$$

487 with

$$488 \quad (3.19) \quad \epsilon_{r,A} = 2b\sqrt{r} \left(1 + \frac{|\lambda_{c_1 r+1}(A)|}{|\lambda_r(A)|} + \frac{2}{\sqrt{b}} \right) \frac{\|\Lambda_3\|_*}{\|\Lambda_2\|_* + \|\Lambda_3\|_*}. \quad \square$$

489 *Remark 3.3.*

490 1. The relative-error nuclear norm bound is informative if $\epsilon_{r,A}$ is small. Now
491 since $b \approx (c_2 - c_1)^{-1} = O(1)$, we have

$$492 \quad (3.20) \quad \epsilon_{r,A} = O \left(\frac{\sqrt{r} \|\Lambda_3\|_*}{\|\Lambda_2\|_* + \|\Lambda_3\|_*} \right).$$

493 Therefore the relative-error nuclear norm bound is good if

$$494 \quad (3.21) \quad \sqrt{r} \sum_{j=c_1 r+1}^n |\lambda_j(A)| = \sqrt{r} \|\Lambda_3\|_* \lesssim \|\Lambda_2\|_* = \sum_{j=r+1}^{c_1 r} |\lambda_j(A)|.$$

495 2. Using a similar proof technique we can obtain mixed norm bounds. The
496 2-norm version of Theorem 3.1 would give

$$497 \quad (3.22) \quad \mathbb{E} [\|E\|_2 | \Omega_F] \leq \|A - \llbracket A \rrbracket_r\|_2 + \frac{\epsilon_{r,A}}{\sqrt{r}} \|A - \llbracket A \rrbracket_r\|_*$$

498 and the Frobenius norm version would give

$$499 \quad (3.23) \quad \mathbb{E} [\|E\|_F | \Omega_F] \leq \|A - \llbracket A \rrbracket_r\|_F + \frac{\epsilon_{r,A}}{\sqrt{r}} \|A - \llbracket A \rrbracket_r\|_*$$

500 where $\epsilon_{r,A}$ is as in Theorem 3.1. Therefore the constant in front of the best
501 rank- r nuclear norm error improves to $\epsilon_{r,A}/\sqrt{r} = O(1)$ using the second
502 remark (3.20). This type of mixed norm bounds along with the relative-error
503 nuclear norm bound in Theorem 3.1 are fairly consistent with the SPSD
504 versions in Table 1 of [14].

505 3. We can relax the condition Ω_F to $\Omega_2 := \left\{ \|\Lambda_3^{1/2} B\|_2^2 \leq 0.5 |\lambda_r(A)| \right\}$ at the
506 cost of a slightly worse bound in Equation (3.2). It is easy to show that the
507 bound in Equation (3.2) then changes to

$$508 \quad (3.24) \quad \mathbb{E} [\|E\|_* | \Omega_2] \leq (1 + \sqrt{r} \epsilon_{r,A}) \|A - \llbracket A \rrbracket_r\|_*.$$

509 *Probability of Ω_F .* The probability of the event Ω_F happening can be computed
510 by following the proof of Theorem 10.8 in [15] using $k = r$ and $p = (c_2 - c_1)r$ and
511 Lemma 3.2. We get

$$512 \quad \mathbb{P} \left(\left\| \Lambda_3^{1/2} B \right\|_F \leq \sqrt{\|\Lambda_3\|_*} \sqrt{\frac{3r}{(c_2 - c_1)r + 1}} t + \sqrt{\|\Lambda_3\|_2} \frac{e\sqrt{(c_2 - c_1)r}}{(c_2 - c_1)r + 1} tu \right) \\ 513 \quad \geq 1 - 2t^{-(c_2 - c_1)r} - e^{-u^2/2}$$

515 for $u, t > 0$. Now using $(x + y)^2 \leq 2(x^2 + y^2)$, we get

$$\begin{aligned}
516 & \left(\sqrt{\|\Lambda_3\|_*} \sqrt{\frac{3r}{(c_2 - c_1)r + 1}} t + \sqrt{\|\Lambda_3\|_2} \frac{e\sqrt{(c_2 - c_1)r}}{(c_2 - c_1)r + 1} tu \right)^2 \\
517 & \leq 2t^2 \left(\|\Lambda_3\|_* \frac{3r}{(c_2 - c_1)r + 1} + \|\Lambda_3\|_2 \frac{e^2(c_2 - c_1)r}{((c_2 - c_1)r + 1)^2} u^2 \right). \\
518 &
\end{aligned}$$

519 Therefore

$$520 \quad (3.25) \quad \mathbb{P}(\Omega_F) = \mathbb{P} \left(\left\| \|\Lambda_3\|^{1/2} B \right\|_F^2 \leq 0.5 |\lambda_r(A)| \right) \geq 1 - 2t^{-(c_2 - c_1)r} - e^{-u^2/2}$$

521 if

$$522 \quad (3.26) \quad 0.5 |\lambda_r(A)| \geq 2t^2 \left(\|\Lambda_3\|_* \frac{3r}{(c_2 - c_1)r + 1} + \|\Lambda_3\|_2 \frac{e^2(c_2 - c_1)r}{((c_2 - c_1)r + 1)^2} u^2 \right),$$

523 i.e., Ω_F holds with high probability when the tail singular values of A decay rapidly.

524 A similar result can also be derived for Ω_2 by following the same results in [15].

525 *Mixed norm bounds.* We can obtain mixed norm bounds for Theorem 3.1. The
526 2-norm version of Theorem 3.1 would give

$$527 \quad (3.27) \quad \mathbb{E} [\|E\|_2 | \Omega_F] \leq \|A - \llbracket A \rrbracket_r\|_2 + \frac{\epsilon_{r,A}}{\sqrt{r}} \|A - \llbracket A \rrbracket_r\|_*$$

528 and the Frobenius norm version would give

$$529 \quad (3.28) \quad \mathbb{E} [\|E\|_F | \Omega_F] \leq \|A - \llbracket A \rrbracket_r\|_F + \frac{\epsilon_{r,A}}{\sqrt{r}} \|A - \llbracket A \rrbracket_r\|_*$$

530 where $\epsilon_{r,A}$ is as in Theorem 3.1. This improves the constant in front of the best rank- r

531 nuclear norm error to $\epsilon_{r,A}/\sqrt{r} = O(1)$ using the first remark (3.20) in Remark 3.3.

532 The proof for the two mixed norm bounds above can be obtained by following the

533 proof of Theorem 3.1. More specifically, the proof for the mixed norm bounds stay

534 the same until we bound a_1 , a_2 and a_3 . To get the mixed norm bound, we use the

535 appropriate norms to bound a_1 , a_2 and a_3 . For example, to bound $\|a_1\|_F$, we start

536 similarly as in the nuclear norm case by conditioning on Ω_F to obtain

$$\begin{aligned}
537 & \|a_1\|_F \leq \left\| (I_r + B^T \Lambda_3 B \Lambda_1^{-1})^\dagger \right\|_2 \|B^T \Lambda_3 B\|_F \\
538 & \leq \frac{\|B^T \Lambda_3 B\|_F}{1 - \|B^T \Lambda_3 B\|_F \|\Lambda_1^{-1}\|_2} \\
539 & \leq \frac{1}{\|\Lambda_1^{-1}\|_2} \left(2 \left\| \|\Lambda_3\|^{1/2} B \right\|_F^2 \|\Lambda_1^{-1}\|_2 \right). \\
540 &
\end{aligned}$$

541 We then get

$$542 \quad \mathbb{E} [\|a_1\|_* | \Omega_F] \leq 2 \mathbb{E} \left[\left\| \|\Lambda_3\|^{1/2} B \right\|_F^2 \middle| \Omega_F \right] = 2b \|\Lambda_3\|_*. \\
543$$

544 The bound for $\|a_2\|_F$, $\|a_3\|_F$, $\|a_1\|_2$, $\|a_2\|_2$ and $\|a_3\|_2$ follows similarly. The mixed

545 norm bounds (3.27) and (3.28) along with the relative-error nuclear norm bound in

546 Theorem 3.1 are fairly consistent with the SPSSD versions in Table 1 of [14].

547 Theorem 3.1 and its proof cannot simply be translated into an algorithm because
 548 the proof relies on the eigendecomposition of A , which is too expensive to compute.
 549 However, the proof naturally suggests Algorithm 2.1. From the proof of Theorem
 550 3.1, under the condition that the matrix has a low-rank structure discussed in this
 551 section, for example in the paragraph after the statement of Theorem 3.1 or in the
 552 remark above, we have that a projection is desired in the core matrix. This projection
 553 gets rid of the large ‘unwanted’ eigenvalues of A , i.e. the eigenvalues in Λ_2 . In the
 554 Nyström method, a natural analogue is to truncate the smallest few singular values
 555 in the core matrix $W = X^T A X$ to achieve the target rank r , which is what has been
 556 done in Algorithm 2.1. The theorem also suggests that the sketch size should be
 557 proportional to the target rank r , which is what we suggest in Algorithm 2.1. Despite
 558 Algorithm 2.1 lacking complete theory (even for the SPSD case), we suggest it because
 559 the algorithm does seem to work well in practice as we illustrate below.

560 **4. Numerical illustration.** We first illustrate Theorem 3.1 and Algorithm 2.1
 561 through experiments. In Figure 4, we show a priori and a posteriori error in Theorem
 562 3.1, and Algorithm 2.1 using 1000×1000 symmetric indefinite matrices. In the left
 563 plot, the matrix A has eigenvalues that decay geometrically from 1 to 10^{-12} each
 564 assigned a random sign with equal probability. In the right plot, A has eigenvalues
 565 equal to ± 1 for the first 100 eigenvalues and $\pm 10^{-10}$ for the other 900 eigenvalues
 566 each assigned a random sign with equal probability; this example illustrates the per-
 567 formance when there is a gap in the singular values. The eigenvectors for both plots
 568 are in a 2×2 block diagonal form, $\text{diag}(I_{100}, U)$ where I_{100} is the 100×100 identity
 569 matrix and $U \in \mathbb{R}^{900 \times 900}$ is a Haar distributed orthogonal matrix. Both the algorithm
 570 and the theorem were constructed using the Gaussian sketch with the sketch size $1.5r$
 571 for the algorithm and $c_1 r = 1.5r$ and $c_2 r = 2r$ for the theorem. We see that Ω_F
 572 holds whenever there is a rapid decay of eigenvalues, i.e., when $|\lambda_r| \gg |\lambda_{c_1 r + 1}|$. But
 573 more importantly, we see that the bound holds when the event Ω_F occurs (circles)
 574 and frequently holds even if the event Ω_F did not occur (crosses). The theorem does
 575 extremely well when Ω_F has occurred. We see that the algorithm gives a good robust
 576 approximation that is a modest factor worse than the best approximation given by
 577 the SVD. Although the theorem does better than the algorithm when Ω_F holds, the
 578 theorem can give unstable approximation when Ω_F does not hold. This illustrates
 579 that the algorithm, which arose from the theorem, works well in practice.

580 In experiments not shown here, we compared Algorithm 2.1 with randomized
 581 SVD [15] and the generalized Nyström method [4, 22, 33, 40], which are applicable to
 582 nonsymmetric (and rectangular) matrices and do not preserve symmetry. We observe
 583 that Algorithm 2.1 tends to obtain a slightly better approximant for a fixed rank r .

584 **4.1. Synthetic examples.** We now compare some of the existing algorithms
 585 against Algorithm 2.1 using different kernel functions and synthetic dataset. We
 586 illustrate the following algorithms

- 587 1. Algorithm 2.1 with the SRFT sketch and the sketch size $s = 2r$,
- 588 2. Algorithm 2.1 with uniform column sampling and the sketch size $s = 2r$,
- 589 3. Algorithm 2.1 with leverage score column sampling and the sketch size $s = 2r$,
- 590 4. Submatrix-Shifted (SMS) Nyström [29] with uniform column sampling and
 591 $s_1 = r$, $s_2 = 2r$ and $\alpha = 1.5$,
- 592 5. Submatrix-Shifted (SMS) Nyström [29] with the Gaussian sketch and $s_1 = r$,
 593 $s_2 = 2r$ and $\alpha = 1.5$,
- 594 6. Stabilized Nyström [3] with the SRFT sketch, $s = r$ and $\epsilon = 10^{-14}$

595 where r is the target rank and the parameters for SMS Nyström and Stabilized

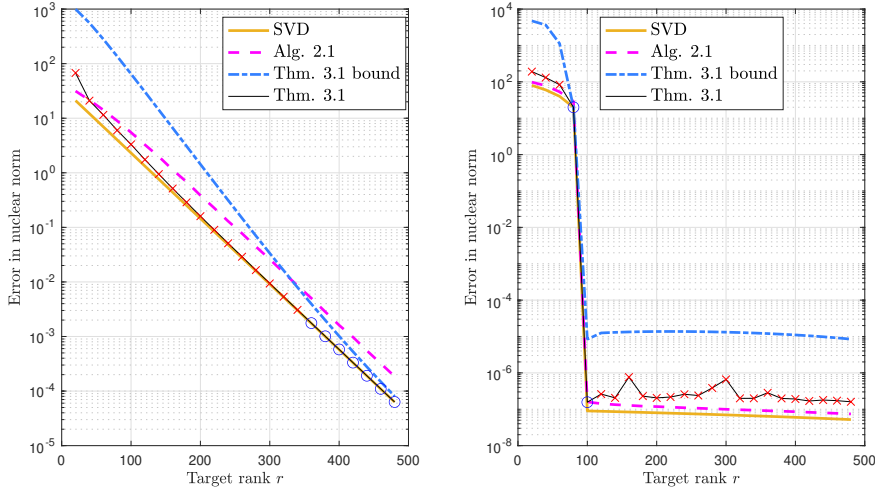


FIG. 4. Two plots showing the empirical results for Theorem 3.1 and Algorithm 2.1. Algorithm 2.1 is robust with the approximation being a modest factor worse than the best approximation. Theorem 3.1 bound holds when Ω_F has occurred (circles on Theorem 3.1) and also frequently holds even if Ω_F has not occurred (crosses on Theorem 3.1). Theorem 3.1 does extremely well when Ω_F has occurred.

596 Nyström are as recommended in their original papers.⁷ For SMS Nyström method,
 597 the Gaussian sketch was not used in the original paper [29]. We use the following
 598 kernel functions

- 599 1. Epanechnikov kernel: $k_1(x, y) = \max\{1 - \|x - y\|^2, 0\}$
 600 2. Multiquadric kernel: $k_2(x, y) = \sqrt{1 + \|x - y\|^2}$
 601 3. Thin plate spline: $k_3(x, y) = \|x - y\|^2 \ln(\|x - y\|^2)$

602 to generate the kernel matrices. The kernel matrices $K^{(1)}, K^{(2)}$ and $K^{(3)}$ correspond-
 603 ing to the kernel functions k_1, k_2 and k_3 were generated by sampling 1000 random
 604 numbers $\{x_i\}_{i=1}^{1000}$ from the standard normal distribution, i.e., $K_{ij}^{(\ell)} = k_\ell(x_i, x_j)$. All
 605 the kernel matrices are symmetric indefinite.

606 In Figure 5, we illustrate the results. The eigenvalue histogram is shown in the left
 607 plots. The right plots show the approximation. We see that SMS Nyström performs
 608 poorly in all 3 examples except the Gaussian case for the multiquadric kernel. This
 609 is possibly because the extreme eigenvalues are large in magnitude so the large shift
 610 is ruining the approximation quality. The stabilized Nyström works well for the
 611 multiquadric kernel and the thin plate spline, but the approximation is very unstable
 612 for the Epanechnikov kernel. This is possibly because the number of positive and the
 613 negative eigenvalues are about the same with similar magnitudes for the Epanechnikov
 614 kernel, which can increase the chance of instability in the core matrix.⁸ This also tells

⁷For stabilized Nyström method, $s = r$ was chosen to ensure that all approximations in the experiment have rank at most r and $\epsilon = 10^{-14}$ as suggested in the original paper was chosen to try diminish the error that might come from taking the pseudo-inverse of the core matrix W .

⁸To our knowledge, the numerical behavior of stabilized Nyström method is an open problem; the stability analysis in [22] applies only to an algorithm where A is sketched from both sides using independent sketches of different dimensions.

615 us that the truncation in the core matrix should not depend on the magnitude of the
 616 singular values of W , but the truncation should always happen proportional to the
 617 target rank. Algorithm 2.1 using uniform column sampling and leverage score column
 618 sampling are both unstable for all 3 examples, which shows the unreliability of using
 619 column sampling matrices. On the other hand, Algorithm 2.1 using the SRFT sketch
 620 works well in all cases.

621 **4.2. Dataset examples.** We now compare the three different methods using
 622 two different high-dimensional datasets, the Coverttype and the Anuran Calls (MFCC)
 623 from the UC Irvine Machine Learning Repository [10]. We illustrate the following
 624 algorithms

- 625 1. Algorithm 2.1 with the SRFT sketch and the sketch size $s = 2r$,
- 626 2. Algorithm 2.1 with k-means++ samples and the sketch size $s = 2r$,
- 627 3. Algorithm 2.1 with uniform column sampling and the sketch size $s = 2r$,
- 628 4. Stabilized Nyström with k-means++ samples, the sketch size $s = r$ and
 629 $\epsilon = 10^{-14}$

630 where r is the target rank. We use the following kernel functions

- 631 1. Thin plate spline kernel: $\|x - y\|^2 \log(\|x - y\|^2)$
- 632 2. Sigmoid kernel: $\tanh\left(\frac{1 + \|x - y\|^2}{\sqrt{1 + \|x - y\|^2}}\right)$
- 633 3. Multiquadric kernel: $\sqrt{1 + \|x - y\|^2}$

634 with the datasets

- 635 1. Coverttype ($n = 581012$) with dimension $d = 54$,
- 636 2. Anuran Calls (MFCC) ($n = 7195$) with dimension $d = 22$.

637 For each dataset, we sample $n = 4000$ data uniformly at random and then center the
 638 mean and normalize all features to have variance 1.

639 The results are illustrated in Figure 6. We observe that the cause of instability in
 640 the Nyström approximation for symmetric indefinite matrices is not necessarily com-
 641 ing from the core matrix W having very small singular values as Stabilized Nyström
 642 can give unstable approximations as seen in Figure 6. Also, although Algorithm 2.1
 643 using k-means++ samples is more accurate than uniform column sampling, they both
 644 do not give robust low-rank approximations. This shows that it is difficult to find
 645 a column sampling scheme that guarantees stable Nyström approximation for sym-
 646 metric indefinite matrices. On the other hand, Algorithm 2.1 using the SRFT sketch
 647 gives robust approximation throughout the experiment and sometimes outperforms
 648 the other methods in this experiment such as in Figure 6a and 6e.

649 **5. Discussion.** Much of the literature on approximating symmetric matrices us-
 650 ing any of the variants of the Nyström method is based on column sampling. In this
 651 work, we used random embeddings for our algorithm (Algorithm 2.1) and a special
 652 class of random embeddings for the analysis, namely Gaussian embeddings. Random
 653 embeddings were used as they are more robust than column sampling, and Gauss-
 654 ian embeddings were used for analysis because we can leverage their rich theoretical
 655 properties. The general behaviour when we use the Nyström method with column
 656 sampling matrices on symmetric indefinite matrices is unknown. In Figure 5, we
 657 see that the two frequently used column sampling schemes, uniform sampling and
 658 leverage score sampling can be unstable. It appears to be difficult to find a column
 659 sampling scheme that guarantees robust Nyström approximation for symmetric indef-
 660 inite matrices and, to our knowledge, is an open problem. We hope that our results
 661 would shed light on the development of a robust indefinite Nyström method based on

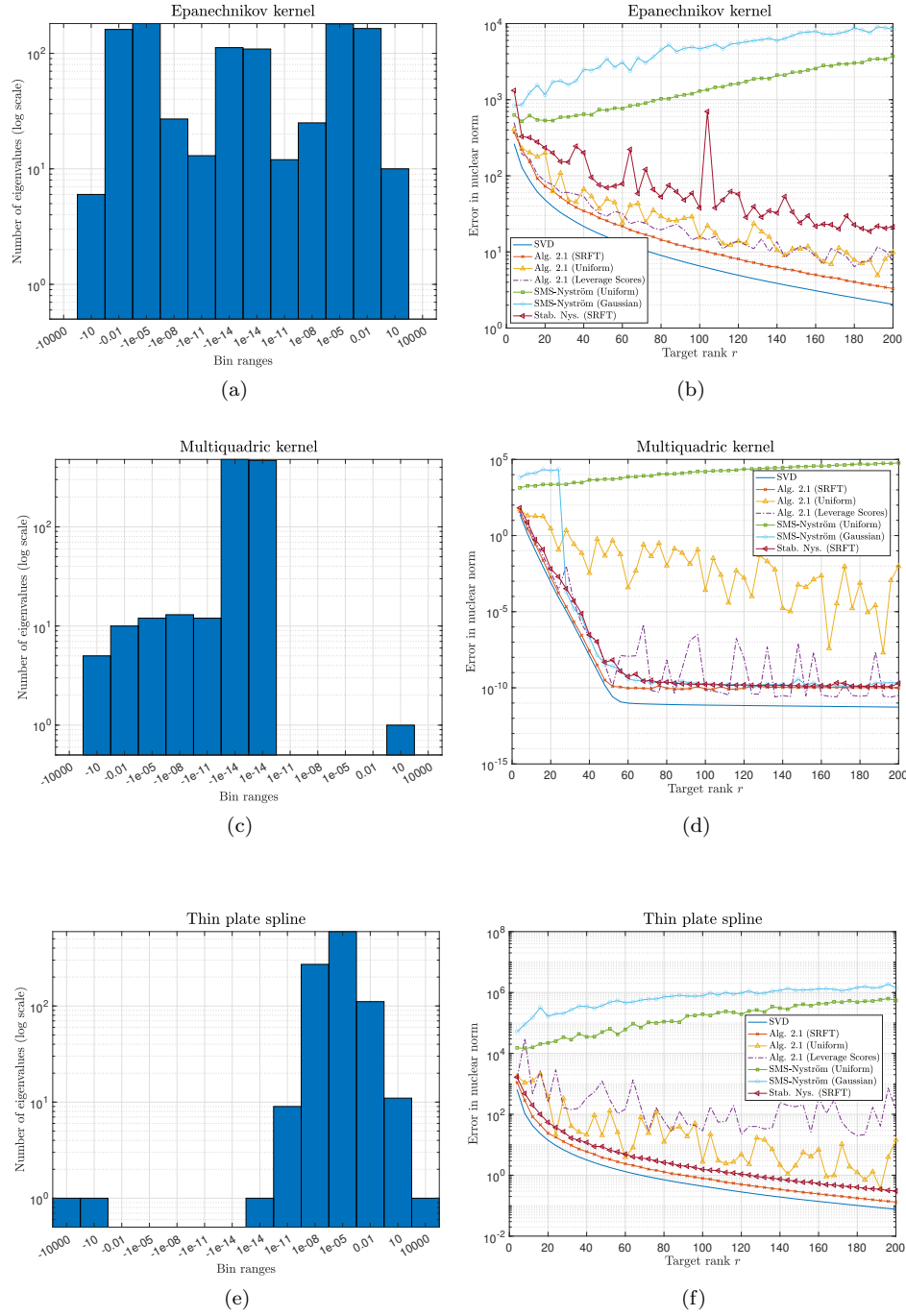


FIG. 5. Comparison of different methods for symmetric indefinite matrices: SMS-Nyström [29], stabilized Nyström [3] and Algorithm 2.1. The first two methods and Algorithm 2.1 using uniform column sampling and leverage score column sampling can fail on some kernels while Algorithm 2.1 using the SRFT sketch (random embedding) works well for all the kernels in the experiment.

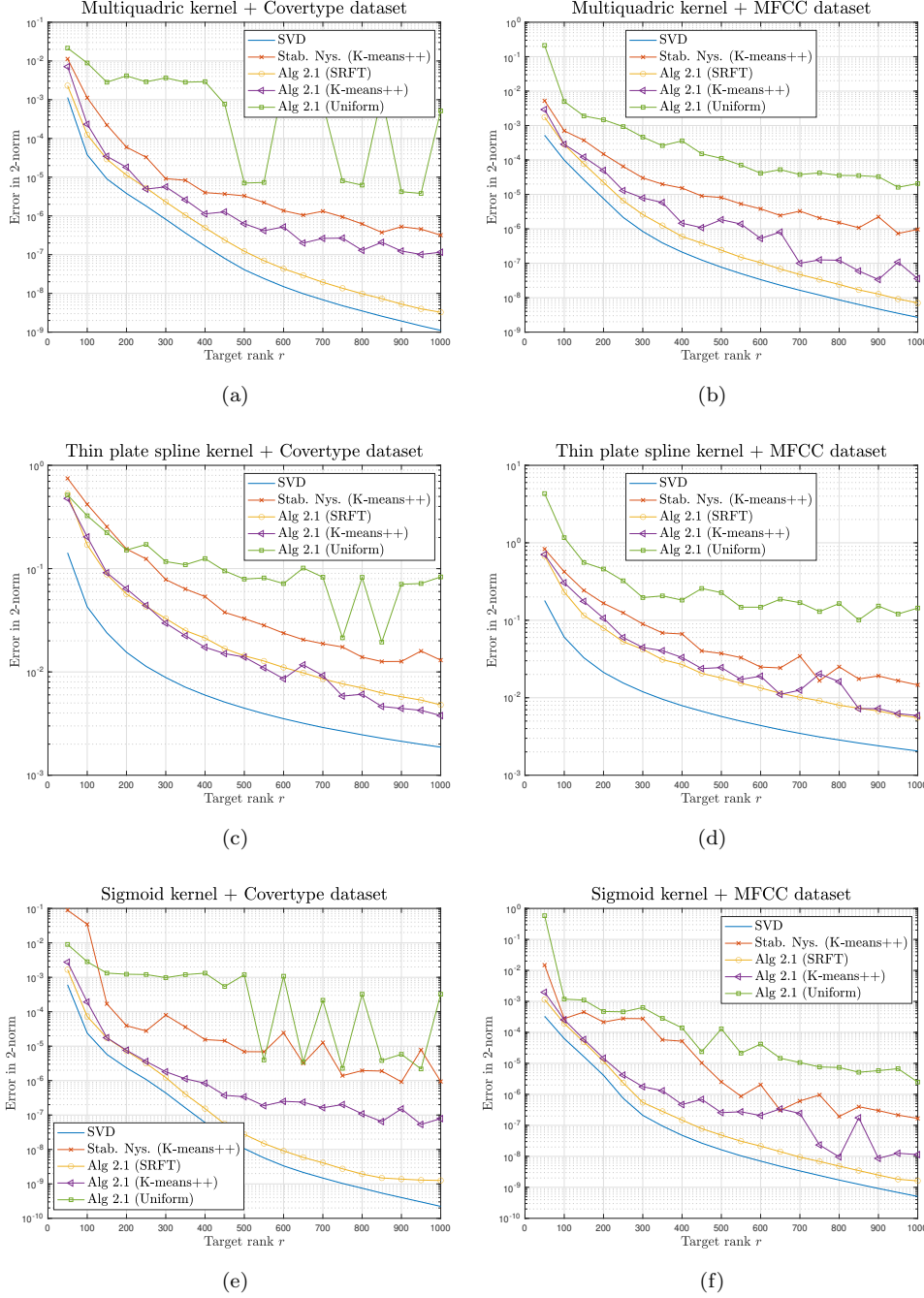


FIG. 6. Comparison of stabilized Nyström [3] and Algorithm 2.1 for symmetric indefinite matrices using three different indefinite kernels and two different datasets. Stabilized Nyström method and Algorithm 2.1 using k -means++ samples and uniform column sampling can give unstable low-rank approximation while Algorithm 2.1 using the SRFT sketch (random embedding) gives robust approximation throughout the experiment.

662 column subsampling.

663 **Acknowledgements.** We thank the anonymous referees and the editor for their
 664 many insightful comments and suggestions, which helped us to improve the quality
 665 of the paper.

666

REFERENCES

- 667 [1] A. ANDONI AND H. L. NGUYÊN, *Eigenvalues of a matrix in the streaming model*, in Proceedings of the Twenty-Fourth Annual ACM-SIAM Symposium on Discrete Algorithms, 2013,
 668 pp. 1729–1737, <https://doi.org/10.1137/1.9781611973105.124>.
- 669 [2] C. BOUTSIDIS AND A. GITTENS, *Improved matrix algorithms via the subsampled randomized*
 670 *Hadamard transform*, SIAM J. Matrix Anal. Appl., 34 (2013), pp. 1301–1340, <https://doi.org/10.1137/120874540>.
- 671 [3] D. CAI, J. NAGY, AND Y. XI, *Fast deterministic approximation of symmetric indefinite kernel*
 672 *matrices with high dimensional datasets*, SIAM J. Matrix Anal. Appl., 43 (2022), pp. 1003–
 673 1028, <https://doi.org/10.1137/21M1424627>.
- 674 [4] K. L. CLARKSON AND D. P. WOODRUFF, *Low-rank approximation and regression in input*
 675 *sparsity time*, J. ACM, 63 (2017), pp. 1–45, <https://doi.org/10.1145/3019134>, <https://doi.org/10.1145/3019134>.
- 676 [5] M. B. COHEN, *Nearly tight oblivious subspace embeddings by trace inequalities*, in Proceedings
 677 of the Twenty-Seventh Annual ACM-SIAM Symposium on Discrete Algorithms, 2016,
 678 pp. 278–287, <https://doi.org/10.1137/1.9781611974331.ch21>.
- 679 [6] A. CORTINOVIS AND D. KRESSNER, *Low-rank approximation in the Frobenius norm by column*
 680 *and row subset selection*, SIAM J. Matrix Anal. Appl., 41 (2020), pp. 1651–1673.
- 681 [7] K. R. DAVIDSON AND S. J. SZAREK, *Local operator theory, random matrices and Banach*
 682 *spaces*, in Handbook of the Geometry of Banach Spaces, W. Johnson and J. Linden-
 683 strauss, eds., vol. 1, Elsevier, 2001, pp. 317–366, [https://doi.org/https://doi.org/10.1016/](https://doi.org/10.1016/S1874-5849(01)80010-3)
 684 [S1874-5849\(01\)80010-3](https://doi.org/10.1016/S1874-5849(01)80010-3).
- 685 [8] J. DEVLIN, M.-W. CHANG, K. LEE, AND K. TOUTANOVA, *BERT: Pre-training of deep bidi-*
 686 *rectional transformers for language understanding*, in Proceedings of the 2019 Confer-
 687 ence of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers), Minneapolis, Minnesota, June 2019, Association for Computational Linguistics, pp. 4171–4186,
 688 <https://doi.org/10.18653/v1/N19-1423>, <https://aclanthology.org/N19-1423>.
- 689 [9] P. DRINEAS AND M. W. MAHONEY, *On the Nystrom method for approximating a Gram matrix*
 690 *for improved kernel-based learning*, J. Mach. Learn. Res., 6 (2005), pp. 2153–2175, <http://jmlr.org/papers/v6/drineas05a.html>.
- 691 [10] D. DUA AND C. GRAFF, *UCI machine learning repository*, 2017, <http://archive.ics.uci.edu/ml>.
- 692 [11] Z. FRANGELLA, J. A. TROPP, AND M. UDELL, *Randomized Nystrom preconditioning*, arXiv
 693 preprint arXiv:2110.02820, (2021), <https://doi.org/10.48550/ARXIV.2110.02820>.
- 694 [12] A. GISBRECHT AND F.-M. SCHLEIF, *Metric and non-metric proximity transformations at linear*
 695 *costs*, Neurocomputing, 167 (2015), pp. 643–657, [https://doi.org/https://doi.org/10.1016/](https://doi.org/https://doi.org/10.1016/j.neucom.2015.04.017)
 696 [j.neucom.2015.04.017](https://doi.org/10.1016/j.neucom.2015.04.017).
- 697 [13] A. GITTENS, *The spectral norm error of the naïve Nystrom extension*, arXiv preprint
 698 arXiv:1110.5305, (2011), <https://doi.org/10.48550/ARXIV.1110.5305>.
- 699 [14] A. GITTENS AND M. W. MAHONEY, *Revisiting the Nystrom method for improved large-scale*
 700 *machine learning*, J. Mach. Learn. Res., 17 (2016), p. 3977–4041.
- 701 [15] N. HALKO, P.-G. MARTINSSON, AND J. A. TROPP, *Finding structure with randomness: Prob-*
 702 *abilistic algorithms for constructing approximate matrix decompositions*, SIAM Rev., 53
 703 (2011), p. 217–288, <https://doi.org/10.1137/090771806>.
- 704 [16] R. A. HORN AND C. R. JOHNSON, *Matrix Analysis*, Cambridge University Press, 2 ed., 2012,
 705 <https://doi.org/10.1017/9781139020411>.
- 706 [17] H. LI, G. C. LINDERMAN, A. SZLAM, K. P. STANTON, Y. KLUGER, AND M. TYGERT, *Algorithm*
 707 *971: An implementation of a randomized algorithm for principal component analysis*,
 708 ACM Trans. Math. Softw., 43 (2017), <https://doi.org/10.1145/3004053>.
- 709 [18] M. LI, W. BI, J. T. KWOK, AND B.-L. LU, *Large-scale Nystrom kernel matrix approximation*
 710 *using randomized SVD*, IEEE Trans. Neural Netw. Learn. Syst., 26 (2015), pp. 152–164,
 711 <https://doi.org/10.1109/TNNLS.2014.2359798>.
- 712 [19] M. W. MAHONEY AND P. DRINEAS, *CUR matrix decompositions for improved data analysis*,
 713 Proceedings of the National Academy of Sciences, 106 (2009), pp. 697–702, <https://doi.org/10.1073/pnas.0810051106>.

- 720 [org/10.1073/pnas.0803205106](https://doi.org/10.1073/pnas.0803205106).
- 721 [20] P.-G. MARTINSSON AND J. A. TROPP, *Randomized numerical linear algebra: Founda-*
 722 *tions and algorithms*, Acta Numer., 29 (2020), p. 403–572, [https://doi.org/10.1017/](https://doi.org/10.1017/s0962492920000021)
 723 [s0962492920000021](https://doi.org/10.1017/s0962492920000021).
- 724 [21] C. MUSCO AND C. MUSCO, *Recursive sampling for the Nyström method*, in Adv. Neural Inf.
 725 Process. Syst., I. Guyon, U. V. Luxburg, S. Bengio, H. Wallach, R. Fergus, S. Vish-
 726 wanathan, and R. Garnett, eds., vol. 30, Curran Associates, Inc., 2017, [https://proceedings.](https://proceedings.neurips.cc/paper_files/paper/2017/file/a03fa30821986dff10fc66647c84c9c3-Paper.pdf)
 727 [neurips.cc/paper_files/paper/2017/file/a03fa30821986dff10fc66647c84c9c3-Paper.pdf](https://proceedings.neurips.cc/paper_files/paper/2017/file/a03fa30821986dff10fc66647c84c9c3-Paper.pdf).
- 728 [22] Y. NAKATSUKASA, *Fast and stable randomized low-rank matrix approximation*, arXiv preprint
 729 arXiv:2009.11392, (2020), <https://arxiv.org/abs/2009.11392>.
- 730 [23] J. NELSON AND H. L. NGUYÊN, *OSNAP: Faster numerical linear algebra algorithms via sparser*
 731 *subspace embeddings*, in Proc. IEEE 54th Annu. Symp. Found. Comput. Sci., 2013, pp. 117–
 732 126, <https://doi.org/10.1109/FOCS.2013.21>.
- 733 [24] E. J. NYSTRÖM, *Über die praktische auflösung von integralgleichungen mit anwendungen*
 734 *auf randwertaufgaben*, Acta Math., 54 (1930), pp. 185 – 204, [https://doi.org/10.1007/](https://doi.org/10.1007/BF02547521)
 735 [BF02547521](https://doi.org/10.1007/BF02547521), <https://doi.org/10.1007/BF02547521>.
- 736 [25] D. OGLIC AND T. GÄRTNER, *Nyström method with kernel k-means++ samples as landmarks*,
 737 in Proceedings of the 34th International Conference on Machine Learning, D. Precup and
 738 Y. W. Teh, eds., vol. 70 of Proceedings of Machine Learning Research, PMLR, 06–11 Aug
 739 2017, pp. 2652–2660, <https://proceedings.mlr.press/v70/oglic17a.html>.
- 740 [26] D. OGLIC AND T. GÄRTNER, *Scalable learning in reproducing kernel Krein spaces*, in Interna-
 741 tional Conference on Machine Learning, PMLR, 2019, pp. 4912–4921.
- 742 [27] B. PICCOLI AND F. ROSSI, *Generalized Wasserstein Distance and its Application to Transport*
 743 *Equations with Source*, Archive for Rational Mechanics and Analysis, 211 (2014), pp. 335–
 744 358, <https://doi.org/10.1007/s00205-013-0669-x>, <https://arxiv.org/abs/1206.3219>.
- 745 [28] F. POURKAMALI-ANARAKI, S. BECKER, AND M. WAKIN, *Randomized clustered Nyström for*
 746 *large-scale kernel machines*, Proceedings of the AAAI Conference on Artificial Intelligence,
 747 32 (2018), pp. 3960–3967, <https://doi.org/10.1609/aaai.v32i1.11614>.
- 748 [29] A. RAY, N. MONATH, A. MCCALLUM, AND C. MUSCO, *Sublinear time approximation of text sim-*
 749 *ilarity matrices*, Proceedings of the AAAI Conference on Artificial Intelligence, 36 (2022),
 750 pp. 8072–8080, <https://doi.org/10.1609/aaai.v36i7.20779>.
- 751 [30] T. SARLOS, *Improved approximation algorithms for large matrices via random projections*, in
 752 Proc. IEEE 47th Annu. Symp. Found. Comput. Sci., 2006, p. 143–152, [https://doi.org/10.](https://doi.org/10.1109/FOCS.2006.37)
 753 [1109/FOCS.2006.37](https://doi.org/10.1109/FOCS.2006.37).
- 754 [31] J. A. TROPP, *Improved analysis of the subsampled randomized Hadamard transform*, Ad-
 755 vances in Adaptive Data Analysis, 03 (2011), pp. 115–126, [https://doi.org/10.1142/](https://doi.org/10.1142/S1793536911000787)
 756 [S1793536911000787](https://doi.org/10.1142/S1793536911000787).
- 757 [32] J. A. TROPP, A. YURTSEVER, M. UDELL, AND V. CEVHER, *Fixed-rank approximation of a*
 758 *positive-semidefinite matrix from streaming data*, in Proceedings of the 31st International
 759 Conference on Neural Information Processing Systems, 2017, p. 1225–1234.
- 760 [33] J. A. TROPP, A. YURTSEVER, M. UDELL, AND V. CEVHER, *Practical sketching algorithms for*
 761 *low-rank matrix approximation*, SIAM J. Matrix Anal. Appl., 38 (2017), p. 1454–1485,
 762 <https://doi.org/10.1137/17m1111590>.
- 763 [34] J. A. TROPP, A. YURTSEVER, M. UDELL, AND V. CEVHER, *Streaming low-rank matrix ap-*
 764 *proximation with an application to scientific simulation*, SIAM J. Sci. Comp., 41 (2019),
 765 pp. A2430–A2463, <https://doi.org/10.1137/18M1201068>.
- 766 [35] M. UDELL AND A. TOWNSEND, *Why are big data matrices approximately low rank?*, SIAM
 767 Journal on Mathematics of Data Science, 1 (2019), pp. 144–160, [https://doi.org/10.1137/](https://doi.org/10.1137/18M1183480)
 768 [18M1183480](https://doi.org/10.1137/18M1183480).
- 769 [36] S. WANG, A. GITTENS, AND M. W. MAHONEY, *Scalable kernel k-means clustering with Nyström*
 770 *approximation: Relative-error bounds*, J. Mach. Learn. Res., 20 (2019), p. 431–479.
- 771 [37] S. WANG, L. LUO, AND Z. ZHANG, *SPSD matrix approximation via column selection: Theories,*
 772 *algorithms, and extensions*, J. Mach. Learn. Res., 17 (2014), pp. 49:1–49:49.
- 773 [38] C. WILLIAMS AND M. SEEGER, *Using the Nyström method to speed up kernel machines*,
 774 in Advances in Neural Information Processing Systems, T. Leen, T. Dietterich, and
 775 V. Tresp, eds., vol. 13, MIT Press, 2000, [https://proceedings.neurips.cc/paper/2000/file/](https://proceedings.neurips.cc/paper/2000/file/19de10adbba1b2ee13f77f679fa1483a-Paper.pdf)
 776 [19de10adbba1b2ee13f77f679fa1483a-Paper.pdf](https://proceedings.neurips.cc/paper/2000/file/19de10adbba1b2ee13f77f679fa1483a-Paper.pdf).
- 777 [39] D. P. WOODRUFF, *Sketching as a tool for numerical linear algebra*, Found. Trends Theor.
 778 Comput. Sci., 10 (2014), p. 1–157, <https://doi.org/10.1561/04000000060>.
- 779 [40] F. WOOLFE, E. LIBERTY, V. ROKHLIN, AND M. TYGERT, *A fast randomized algorithm for*
 780 *the approximation of matrices*, Appl. Comput. Harmon. Anal., 25 (2008), pp. 335–366,
 781 <https://doi.org/https://doi.org/10.1016/j.acha.2007.12.002>.